# Challenges in Machine Learning for Predicting Psychological Attributes from Smartphone Data

Jiew-Quay Au

München 2024

# Challenges in Machine Learning for Predicting Psychological Attributes from Smartphone Data

Jiew-Quay Au

Dissertation
an der Fakultat für Mathematik, Informatik und Statistik
der Ludwig–Maximilians–Universitat München

eingereicht von
Jiew-Quay Au
am 04.09.2023

Erster Berichterstatter: Prof. Dr. Bernd Bischl
Zweiter Berichterstatter: Prof. Dr. Markus Bühner
Dritter Berichterstatter: Prof. Dr. Achim Zeileis

Tag der Disputation: 30.01.2024

# Acknowledgments

*I am deeply grateful to the many individuals whose help, support, guidance, and advice made this thesis possible. In particular, I would like to express my sincere thanks to the following people: . . .*

. . . *Prof. Dr. Bernd Bischl for his unwavering support, encouragement, and valuable guidance throughout the years. Our seamless collaboration and the trust he placed in me have been instrumental in making this thesis a reality.*

. . . *Prof. Dr. Markus Bühner and Prof. Dr. Achim Zeileis for their willingness to serve as the second and third reviewers for my PhD thesis.*

. . . *Prof. Dr. Christian Heumann and Prof. Dr. Helmut Küchenhoff for their willingness to be part of the examination panel at my PhD defense.*

. . . *Prof. Dr. Clemens Stachl, Dr. Sarah Theres Völkel, and Dr. Ramona Schödel for their exceptional support, collaborative spirit, and teamwork during our time working together. I would also like to express my sincere appreciation to Dr. Stefan Hummel of AUDI AG for his consistent support and guidance. His valuable insights and expertise were instrumental in helping us navigate complex challenges and achieve our goals.*

. . . *all my coauthors for their supportive collaboration.*

. . . *my parents, my brother, and my sister for their support and encouragement throughout my life.*

. . . *my partner Vanessa and our children Eliza, Estelle, and Eleanor for their unwavering love and constant encouragement, which have enriched my life beyond measure. Although their lively and charismatic personalities often posed a challenge to maintaining focus during my research, they always stood by me with unrivaled dedication to help me achieve my goal of completing my PhD.*

# Summary

Predicting psychological attributes using psychometric approaches is a complex task that involves estimating latent constructs that cannot be directly measured. Psychometrics focuses on the measurement and assessment of psychological attributes, such as personality traits, behavioral patterns, or psychological disorders. Traditionally, personality assessment relied on self-report questionnaires, but advancements in technology have opened up new possibilities for assessment, particularly through the analysis of digital footprints.

Smartphone sensor data has become particularly valuable in this context. By analyzing data related to movement, conversation patterns, activities, and interests, it is possible to gather insights that can contribute to predicting psychological attributes. Machine learning techniques are commonly employed to develop predictive models in this field. However, it is essential to ensure that the predictions are meaningful, accepted, and interpretable to gain trust from users.

Interpreting machine learning models is crucial in the context of psychometric prediction. Interpreting the models helps identify biases, understand their operations, and determine the variables they rely on. This process enhances the accuracy of the models, establishes trust in their predictions, and promotes fairness in the prediction process. Given the large datasets involved in using smartphone sensor data, the issue of multicollinearity arises, making it challenging to identify which features are truly essential for predicting psychological attributes. To address this challenge, this thesis focuses on grouping similar features and quantifying their importance, aiming to reduce data complexity and highlight the most relevant factors. Additionally, visualizing the impact of these feature groups can provide a deeper understanding in the behavior of the predictive models.

# Zusammenfassung

Psychometrie bezieht sich auf die Messung psychologischer Merkmale wie Persönlichkeitsmerkmale, Verhaltensmuster oder psychischer Störungen. Üblicherweise werden hierfür Selbstauskunftsfragebögen verwendet, da psychologische Merkmale oft nicht direkt messbar sind. Dank technologischer Fortschritte eröffnen sich jedoch moderne Möglichkeiten, psychologische Merkmale vorherzusagen, insbesondere durch die Analyse digitaler Fußspuren.

Besonders relevant sind in diesem Zusammenhang Smartphone-Sensordaten. Durch die Auswertung von Daten zu Bewegungsmustern, Gesprächsverhalten, Aktivitäten und Interessen können Erkenntnisse gewonnen werden, die zur Vorhersage psychologischer Merkmale beitragen können. Hierbei kommen häufig maschinelle Lernverfahren zum Einsatz. Dabei ist es wichtig sicherzustellen, dass die Vorhersagen sinnvoll, akzeptiert und interpretierbar sind.

Die Interpretation maschineller Lernverfahren spielt bei der Vorhersage psychologischer Merkmale eine entscheidende Rolle. Sie hilft dabei, die Funktionsweise der Modelle zu verstehen und wichtige Variablen zu identifizieren. Bei der Verwendung von Smartphone-Daten entstehen große Datensätze, was das Problem der Multikollinearität mit sich bringt. Dies erschwert die Bestimmung, welche Merkmale tatsächlich relevant sind, um psychologische Merkmale vorherzusagen. Um dieser Herausforderung zu begegnen, konzentriert sich diese Arbeit darauf, ähnliche Merkmale zu gruppieren und ihre Bedeutung zu quantifizieren. Dadurch kann die Komplexität der Daten reduziert und die relevantesten Faktoren hervorgehoben werden. Darüber hinaus kann die Visualisierung der Effekte dieser Merkmalsgruppen ein besseres Verständnis für das Verhalten der Vorhersagemodelle liefern.

x

# Contents

# Introduction

## 1.1  Outline

This thesis focuses on the utilization of smartphone sensor data in machine learning to predict psychological traits and discusses the challenges that arise in the field of machine learning.

In Chapter 2, the application of smartphone sensor data in the context of psychological research is discussed. The PhoneStudy (Stachl et al., 2022), a comprehensive research project involving three distinctive phases of data collection, is introduced. This project has led to a number of research articles, underscoring the significant impact of the PhoneStudy.

Chapter 3 introduces fundamental machine learning concepts, with a focus on their practical use in predicting psychological traits using sensor data. Substantial challenges that arise in the field of machine learning are discussed here, especially the importance of interpreting and explaining the inner workings of machine learning models.

The remaining sections of this thesis are divided into two main parts, labeled as Part I and Part II. These parts include the contributing articles, structured as individual chapters (Chapter 4 - 12). Each chapter starts with a comprehensive reference to the original publication, accompanied by a detailed description of the specific contributions made by the doctoral candidate. Moreover, where applicable, supplementary materials, accompanying software, and essential copyright information for the articles are included.

Finally, the thesis concludes with Part III, offering insights into potential future directions and ongoing research.

## 1.2  Motivation and Scope

In recent years, the integration of machine learning techniques with psychological research has significantly advanced our understanding of human behavior and cognitive processes (Montag and Elhai, 2019; Kosinski et al., 2013). This interdisciplinary approach has opened up a wide range of possibilities, with one area of particular interest being the utilization of smartphone sensor data in psychological studies. Smartphones have become widespread in modern society,

and they are equipped with an array of sensors that can capture rich and diverse data about individuals' behavior, activities, and interactions.

The adoption of smartphones in daily life has led to the generation of vast amounts of sensor data, including GPS location, accelerometer readings, call and text logs, application usage, screen touches, and more. This vast amount of data provides researchers with unprecedented opportunities to gain insights into various aspects of human behavior, such as mobility patterns, social interactions, sleep patterns, physical activity, and even mental health indicators. Psychological researchers have embraced the potential of smartphone sensor data to address research questions and explore novel avenues of investigation. For instance, accelerometer data has been used to predict cognitive decline and detect symptoms of depression (Saeb et al., 2015).

The use of machine learning algorithms to analyze smartphone sensor data has proven to be instrumental in uncovering meaningful patterns and relationships within the data. Techniques such as supervised learning, unsupervised learning, and deep learning have been employed to extract valuable information from these rich datasets. Moreover, the integration of data from multiple sensors has enabled researchers to create more comprehensive models of human behavior, yielding a deeper understanding of psychological phenomena (Wan et al., 2020). The integration of machine learning and smartphone sensor data in psychological research shows great promise. Still, certain challenges must be addressed thoughtfully. These challenges include safeguarding data privacy, handling ethical concerns, and acknowledging potential biases in the data collection process. The responsible and ethical use of this technology in research is essential to ensure both its effectiveness and the protection of the individuals involved (Fuller et al., 2017).

The thesis presents significant contributions in distinct areas. Firstly, it focuses on predicting personality traits, exemplified by the utilization of the Big Five Inventory, as detailed in Chapter 8 (Stachl et al., 2020). Furthermore, the thesis explores the prediction of the sensation seeking personality trait, a topic discussed in Chapter 9 (Schoedel et al., 2018). Additionally, the thesis explores the relationship between smartphone usage and various aspects, such as autism (Chapter 12, Schuwerk et al. (2019)), circadian rhythm (Chapter 11, Schoedel et al. (2020)), and personality traits based on the Big Five model (Chapter 10, Stachl et al. (2017)).

The main focus in this work is on supervised machine learning and the challenges that arise when applying these techniques to smartphone sensor data for predicting psychological attributes. These challenges can be complex, especially when predicting personality traits (e.g., Big Five Inventory). One challenge arises when there are multiple target variables that need to be predicted simultaneously, which can make the modeling process more complicated. To tackle this issues, the article in Chapter 5 (Probst et al., 2017) introduces an implementation of Multilabel algorithms within the machine learning software mlr Bischl et al. (2016). An updated version of this machine learning framework for the statistical software R is introduced in Chapter 6 (Lang et al., 2019).

Furthermore, especially in psychological research, there is a need for explaining the statistical models, which are used. There is a trade-off between simple models, which are easy to interpret, and complex, so-called "black box" models, which are typically highly accurate, but they lack transparency, making it challenging to interpret the underlying reasons for their predictions (Rudin, 2019). Conventional techniques for interpreting black box models primarily focus on quantifying the importance and visualizing the effects of individual features. In the specific use case in the PhoneStudy, where smartphone sensor data is used to predict psychological attributes, these features displayed high intercorrelation but could naturally be categorized into distinct groups, such as music listening data, call behavior, GPS sensor data, and more. To address the need to interpret machine learning models that incorporate grouped features, Chapter 4 (Au et al., 2022) provides an in-depth discussion of methods explicitly designed for assessing the importance and visualizing the effects from groups of features.

# Smartphone Sensor Data

The widespread adoption of smartphones has made them an essential aspect of our daily routines. These devices come equipped with a variety of sensors and logging features that enable the collection of extensive data regarding human behavior. They can record diverse aspects such as app usage, media consumption, location tracking, and communication patterns, providing a comprehensive picture of individual's daily activities.

The large amount of data collected on smartphones creates substantial research opportunities across multiple fields, particularly in psychology. Researchers have already started leveraging this data to discover valuable insights. For example, studies have investigated smartphone data regarding mental health (Zbiciak and Markiewicz, 2023; Servia-Rodríguez et al., 2017), movement patterns and social behavior (Harari et al., 2017), as well as sleep patterns and cognitive performance (Wilmer et al., 2017). Moreover, personality traits have also been a subject of investigation using smartphone data (Chittaranjan et al., 2013; de Montjoye et al., 2013; Mønsted et al., 2018; Harari et al., 2020; Montag et al., 2014).

## 2.1 PhoneStudy

The data used in this thesis was collected as part of the PhoneStudy mobile sensing research project (Stachl et al., 2022) conducted at the Ludwig-Maximilians-Universität München. The data collection process spanned from 2014 to 2018 and consisted of three separate data collections, each following distinct study procedures. Further information regarding the procedures of these individual studies can be found in the respective research articles (Stachl et al. (2017), Schuwerk et al. (2019), Schoedel et al. (2020)).

In summary, behavioral data was gathered from 624 volunteers over a period of 30 consecutive days. The majority of participants (91%) had completed A levels, while 20% held a university degree. The sample had a relatively young age distribution, with a mean age of 23.56 years and a standard deviation of 6.63. However, there was an imbalance in terms of gender representation, with 377 women, 243 men, and 4 participants who chose not to disclose their gender. All participants provided informed consent and willingly participated in the study. They retained the right to withdraw their participation and request the deletion of their data, as long as re-identification

remained possible. The research procedures were approved by the Institutional Review Board of the Psychology Department at Ludwig-Maximilians-Universität München, and the study was conducted in accordance with the applicable laws and regulations of the European Union. The successful publications resulting from the PhoneStudy research project have played a central role in shaping this thesis, with contributions that can be summarized as follows:

***Predicting Personality from Patterns of Behavior Collected with Smartphones*** – Chapter 8 (Stachl et al., 2020) investigates the use of machine learning for predicting Big Five personality dimensions using behavioral data collected from smartphones. The research identifies distinct behavioral patterns in communication, social behavior, music consumption, app usage, mobility, overall phone activity, and day-night activity that serve as predictors of personality traits. The study highlights the potential benefits for research and raises concerns about privacy and psychological targeting when collecting and modeling behavioral data from smartphones.

***Digital Footprints of Sensation Seeking*** – Chapter 9 (Schoedel et al., 2018) explores the impact of new technologies on personality research, where personality traits are investigated through digital footprints. Using data collected on smartphones, the study predicts self-reported sensation seeking scores with machine learning methods, demonstrating the potential of mobile sensing techniques in enhancing the understanding of human behavior.

***Personality Traits Predict Smartphone Usage*** – The contribution in Chapter 10 (Stachl et al., 2017) examines how psychological attributes, like personality traits, fluid intelligence, and demographics, can predict mobile application usage on smartphones. It reveals that personality traits, particularly extraversion, conscientiousness, and agreeableness, are better predictors of app usage in specific categories like communication, photography, gaming, transportation, and entertainment. Additionally, the study shows that fluid intelligence and demographics also have stable associations with categorical app usage.

***To Challenge the Morning Lark and the Night Owl: Using Smartphone Sensing Data to Investigate Day–Night Behaviour Patterns*** – Chapter 11 (Schoedel et al., 2020) highlights the potential of smartphone sensing in investigating day-night patterns and related traits. This article explores individual differences in day-night patterns, their association with personality traits, and the impact of traits and work behaviors on day-night patterns during weekends.

***Enter the Wild: Autistic Traits and Their Relationship to Mentalizing and Social Interaction in Everyday Life*** – In Chapter 12 (Schuwerk et al., 2019), smartphone-based experience sampling was used to investigate the relationship between mentalizing, autistic traits, and social behavior in everyday life. The results showed that mentalizing occurred less frequently compared to reasoning about actions, and individuals with higher autistic traits displayed reduced communication via smartphone. However, there was no significant association between autistic traits and social media usage for connecting with others.

## 2.2 Psychometric Assessments in the PhoneStudy

During the PhoneStudy, several psychometric assessments were conducted. Across all three data collections, the Big Five Inventory was assessed to measure participants' personality traits. In the second data collection (Schuwerk et al., 2019), the focus was also on assessing autism-related characteristics. In the third data collection (Schoedel et al., 2020), the psychometric assessment focused on measuring participants' sensation seeking tendencies.

### 2.2.1 Big 5 Inventory

The Big Five Inventory (BFI) is a widely accepted tool for assessing personality, which is based on the Five-Factor Model of personality, commonly known as the Big Five personality traits (John et al., 2008). In the PhoneStudy, participants were given a questionnaire consisting of 300 items. Participants were asked how well each item described them on a four-point Likert scale, ranging from "untypical for me" to "typical for me". Data collection was conducted in two different ways, depending on the respective study's design. Some participants completed the questionnaire on provided computers. Other participants were provided access to the questionnaire on their smartphones. The Five-Factor Model consists of five major personality factors:

**Openness**: This trait reflects an individual's willingness to engage with imagination, creativity, and openness to new experiences. People high in openness tend to be curious, innovative, and willing to explore unconventional ideas. Conversely, those low in openness may exhibit resistance to change, a preference for routine, and a tendency to avoid novel situations.

**Conscientiousness**: Conscientious individuals are characterized by their organization, responsibility, and goal-oriented behavior. They are generally disciplined, reliable, and tend to plan and complete tasks meticulously. On the negative end, low conscientiousness may lead to disorganization, lack of follow-through, and a more relaxed approach to responsibilities.

**Extraversion**: Extraverts are outgoing, social, and enjoy interacting with others. They tend to be energetic, assertive, and seek stimulation from their external environment. Conversely, introverted individuals might be reserved, prefer solitary activities, and may find social interactions draining.

**Agreeableness**: This factor relates to an individual's interpersonal tendencies, such as kindness, cooperation, and empathy. People high in agreeableness are considerate, nurturing, and tend to prioritize harmonious relationships. On the other hand, low agreeableness could manifest as being more competitive, less empathetic, and less inclined to cooperate with others.

**Emotional Stability**: Also known as neuroticism, this trait refers to an individual's emotional resilience and stability. Those with high emotional stability are more composed, calm, and less

prone to experiencing negative emotions. Individuals with low emotional stability, however, might be more susceptible to stress, anxiety, and mood fluctuations.

The questionnaire provides ratings for each of the five factors. These ratings are expressed as numerical values and are typically standardized to allow for a better comparison among different individuals. In addition, each factor also consists of six sub-facets, making a total of 30 facets and 5 factors, that contribute to a comprehensive personality assessment.

### 2.2.2   Sensation Seeking

Sensation seeking is a personality trait that indicates an individual's tendency to pursue new, thrilling, and exciting experiences (Zuckerman, 1994). The sensation seeking personality trait of the participants was measured using the impulsive sensation seeking (ImpSS) subscale of the Zuckerman–Kuhlman personality questionnaire (Zuckerman, 2002). The 19 individual item scores were summed up to derive an ImpSS score, which ranged between 0 and 19. For more detailed information, please consult the research article by Schoedel et al..

### 2.2.3   Autism

The level of autistic traits was assessed using the following three most commonly used and validated self-report questionnaires for adults: the Autism-Spectrum Quotient (AQ), the Empathy Quotient (EQ), and the Broader Autism Phenotype (BAP) questionnaire. See the respective research article (Schuwerk et al., 2019) for more details. These questionnaires are known for their sensitivity in measuring the prevalence of autistic traits in the general population, with each targeting slightly different aspects of autistic personality traits. To analyze the data in this study, individual scores from these three questionnaires were combined into a single compound score of autistic traits, calculated as the mean of z-transformed scores from each questionnaire. All participants completed the questionnaires, including the control questionnaires, using computers in the laboratory.

### 2.2.4   Machine Learning for PhoneStudy Data

In the PhoneStudy project, one objective was to apply machine learning algorithms to predict the big five personality traits based on smartphone data. However, preparing the data for machine learning posed a significant challenge due to the nature of the dataset. The data consisted of longitudinal records for each of the 624 participants, capturing various smartphone events such as screen interactions, app usage, GPS coordinates, call and SMS events, and more. As a result,

the dataset grew quite large, with over 30 million rows and 70 columns, amounting to approximately 18GB of data stored in a MariaDB database (see Table 2.1). This sheer volume of data made it impractical to directly use it for machine learning purposes. To address this issue, extensive feature engineering was undertaken to extract meaningful and relevant information from the raw smartphone data. A dedicated effort was made to derive insightful features that could serve as inputs for the machine learning algorithms. In total, over 100 individual functions were developed to process the data and extract more than 15,000 features for each participant. The high number of features is a result of calculating features for each app used by the participants. For instance, this includes metrics such as the *daily mean number of app uses* for each app, resulting in a large set of features. The feature engineering process aimed to capture relevant patterns, behaviors, and interactions from the smartphone data that could potentially be indicative of the participants' big five personality traits. Through careful selection and extraction of features, the dataset was transformed into a more manageable and informative representation (see Figure 2.1). For further details on the feature engineering process and the specific features derived, interested readers can refer to the research paper (Stachl et al., 2020) on the topic. Additional information can also be found on the project's website[1].

Various algorithms were considered, starting with simpler and interpretable linear models like Lasso regression (Tibshirani, 1996). These models facilitated a better understanding of the relationship between features and the target variable (e.g., Big 5 personality traits, or sensation seeking score). In addition, decision trees were explored as an alternative approach, offering interpretable rules for decision-making (Breiman et al., 2017). However, they had limitations in capturing complex data interactions, which affected their predictive power. For a better predictive performance, more advanced algorithms such as random forests (Breiman, 2001), neural networks (Rumelhart et al., 1986), and boosting algorithms (Friedman, 2001a) were investigated.

The performance of these models is commonly compared in benchmark experiments using nested cross-validation (Bischl et al., 2012). In this approach, hyperparameter tuning is exclusively performed on the inner datasets, while the evaluation of model performance takes place on the outer datasets that were left out during the tuning process. The contribution in Chapter 6 offers a machine learning framework in the statistical programming language R and aids in this process by making it easier to apply different algorithms on various datasets while tuning hyperparameters.

Furthermore, since many factors and facets contribute to an individual's personality assessment, it could be advantageous to use algorithms, that predict the target variables simultaneously, while also addressing potential interdependencies among these. The contribution presented in Chapter 5 provides an implementation of multilabel algorithms, serving as an initial approach to address scenarios with multiple target variables. For a more comprehensive understanding of these learning tasks, please refer to Section 3.1.2.

---

[1] https://compstat-lmu.shinyapps.io/Personality_Prediction/

| ID | Timestamp | Source | Event | Call Length | .... |
|----|-----------|--------|-------|-------------|------|
| 1 | 1597996890 | SCREEN | SCREEN_ON | | |
| 1 | 1597996892 | APP | Homescreen | | |
| 1 | 1597996895 | APP | WhatsApp | | |
| ... | ... | ... | ... | | |
| 624 | 1565265600 | WIFI | WIFI_ON | | |
| 624 | 1565265610 | APP | Homescreen | | |
| 624 | 1565265613 | PHONE | Incoming | 132 | |

Table 2.1: Example for the raw smartphone logging data. This data is timestamped, providing information about events that occur on the smartphone, and it includes additional metadata for certain information (e.g., call length). The raw data consists of more than 30 million rows and 70 columns.



Figure 2.1: Schematic representation of feature extraction from raw data. This process results in one row per participant (ID). Exemplary features include the *daily mean number of incoming calls*, or the *daily mean number of app uses* for each app.

In conclusion, Section 3.2 tackles the crucial task of interpreting machine learning models. This interpretation is vital because it allows to gain a deeper understanding of how these models make predictions. The `iml`-package (Molnar et al., 2018) is a valuable resource, offering practical model interpretation methods within the R environment.

# Methodological Background

In this thesis, the focus lies on the prediction of psychological attributes from smartphone data using machine learning algorithms.

By leveraging various types of smartphone data, including communication behavior, app usage, GPS location, or music usage, valuable insights can be gained into people's behavior. These insights can potentially serve as indicators of their psychological traits. Machine learning techniques play a crucial role in this process. By applying these algorithms, valuable patterns and relationships within the data can be uncovered, enabling accurate predictions of psychometric measures. Furthermore, these approaches contribute to the generation of hypotheses about human behavior and provide insights into individuals' psychological characteristics. This chapter offers introductory insights into supervised machine learning, accompanied by mathematical definitions and essential background information.

## 3.1 Supervised Machine Learning

In supervised machine learning, the goal is to learn from labeled data, where input-output pairs are given. The primary aim is to train the model so that it can make precise predictions on new, unseen data. In contrast, unsupervised machine learning differs from supervised learning, as the model learns from unlabeled data without explicit input-output pairs. The objective in unsupervised learning is to explore the inherent structure within the data, identify patterns, and discover relationships among the variables. Common unsupervised learning techniques include clustering, dimensionality reduction, and density estimation. While unsupervised learning holds significant value with its wide range of applications, this thesis primarily focuses on supervised machine learning.

### 3.1.1 Notation

In the domain of supervised machine learning, the underlying assumption is the existence of an unknown functional relationship between a feature space $\mathcal{X}$ and a target space $\mathcal{Y}$

$$f : \mathcal{X} \longrightarrow \mathcal{Y}, \tag{3.1}$$

where $\mathcal{X}$ is a $p$-dimensional feature space, defined as $\mathcal{X} = \mathcal{X}_1 \times \ldots \times \mathcal{X}_p$, and $\mathcal{Y}$ is a $d$-dimensional target space, defined as $\mathcal{Y} = \mathcal{Y}_1 \times \ldots \times \mathcal{Y}_d$. The corresponding random variables associated with these spaces are denoted as $X = (X_1, \ldots, X_p)$ and $Y = (Y_1, \ldots, Y_d)$.

Machine learning algorithms aim to discover the underlying functional relationship by using a set of $n \in \mathbb{N}$ independently and identically distributed (i.i.d.) observations drawn from the joint space $\mathcal{X} \times \mathcal{Y}$, where the underlying probability distribution $\mathcal{P}$ is unknown. This dataset is represented as $\mathcal{D} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^n$, where the vector $\mathbf{x}^{(i)} = (x_1^{(i)}, \ldots, x_p^{(i)})^\mathsf{T} \in \mathcal{X}$ is the $i$-th observation associated with the target variable $\mathbf{y}^{(i)} = (y_1^{(i)}, \ldots, y_d^{(i)})^\mathsf{T} \in \mathcal{Y}$. The $j$-th feature is denoted by $\mathbf{x}_j = (x_j^{(1)}, \ldots, x_j^{(n)})^\mathsf{T}$, for $j = 1, \ldots, p$.

The performance of a fitted model $\hat{f}$ is assessed by the general error measure

$$\rho(\hat{f}, \mathcal{P}) = \mathbb{E}(L(\hat{f}(X), Y)). \tag{3.2}$$

This measure represents the expected value of a loss function $L$ when applied to test data drawn independently from $\mathcal{P}$. To estimate this general error, the model is evaluated on unseen test data denoted as $\mathcal{D}_\text{test}$, and the estimation is calculated as

$$\hat{\rho}(\hat{f}, \mathcal{D}_\text{test}) = \frac{1}{|\mathcal{D}_\text{test}|} \sum_{(\mathbf{x}, y) \in \mathcal{D}_\text{test}} L(\hat{f}(\mathbf{x}), \mathbf{y}). \tag{3.3}$$

When a machine learning algorithm (or *learner*) $\mathcal{I}$ is applied to a given dataset $\mathcal{D}$, it produces a fitted model $\mathcal{I}(\mathcal{D}) = \hat{f}_\mathcal{D}$. The *expected generalization error* of a learner $\mathcal{I}$ considers the variability by sampling different datasets $\mathcal{D}$ of equal size $n$ from the underlying probability distribution $\mathcal{P}$ and is defined by

$$GE(\mathcal{I}, \mathcal{P}, n) = \mathbb{E}_{|\mathcal{D}|=n}(\rho(\mathcal{I}(\mathcal{D}), \mathcal{P})). \tag{3.4}$$

In practice, to estimate Eq. (3.4), resampling techniques like cross-validation or bootstrapping are applied to the available dataset $\mathcal{D}$. These techniques involve creating multiple subsets of the data, which are then used for training and evaluation. The dataset $\mathcal{D}$ is divided into $k \in \mathbb{N}$ training datasets $\mathcal{D}_\text{train}^i$ and their corresponding test datasets $\mathcal{D}_\text{test}^i$, $i = 1, \ldots, k$. Each training dataset contains roughly the same size $n_\text{train}^i < n$. The expected generalization error can be estimated by

$$\widehat{GE}(\mathcal{I}, \mathcal{D}, n_\text{train}) = \frac{1}{k} \sum_{i=1}^k \hat{\rho}(\hat{f}_{\mathcal{D}_\text{train}^i}, \mathcal{D}_\text{test}^i). \tag{3.5}$$

### 3.1.2 Learning Tasks

The type of machine learning task is determined by the target space $\mathcal{Y} = \mathcal{Y}_1 \times \ldots \times \mathcal{Y}_d$. In the case, where $d = 1$, there are two fundamental machine learning tasks: regression and classification.

***Regression*** – Regression involves predicting a continuous target, i.e., $\mathcal{Y} = \mathbb{R}$. An illustrative example in psychometric prediction is the utilization of machine learning models to accurately predict an individual's IQ based on a variety of cognitive tests, educational background, and other pertinent factors (Schütze et al., 2018).

***Classification*** – On the other hand, classification focuses on predicting discrete classes, which can be either binary classification involving two classes or multiclass classification involving more than two classes. Mathematically, binary classification can be represented as $\mathcal{Y} = \{0, 1\}$, where $0$ and $1$ represent the two distinct classes. Multiclass classification, on the other hand, is denoted as $\mathcal{Y} = \{1, 2, ..., m\}$, where $m$ denotes the total number of classes. In the field of psychometric prediction, classification techniques find diverse applications. For instance, binary classification enables machine learning models to discern whether an individual belongs to a particular category, such as a diagnostic group, by analyzing specific features or test results. Furthermore, multiclass classification plays a crucial role in classifying individuals into multiple distinct groups. One illustrative example is in the context of Attention Deficit Hyperactivity Disorder (ADHD) classification, where individuals can be categorized into different subtypes, including typically developing (TDC), ADHD-inattentive (ADHD-I), and ADHD-combined (ADHD-C) (Qureshi et al., 2016).

In addition to the fundamental machine learning tasks of regression and classification, when encountering a more complex target space $\mathcal{Y} = \mathcal{Y}_1 \times ... \times \mathcal{Y}_d$ where $d > 1$, several other notable machine learning tasks come into play (Waegeman et al., 2018):

***Multilabel Classification*** – When all target variables are binary, this problem is known as multi-label classification. Unlike traditional classification, where an instance belongs to a single class, multilabel classification allows for more complex and nuanced categorization. The target space can be defined as $\mathcal{Y} = \{0, 1\}^d$. This task is commonly encountered in real-world applications, such as image recognition, where an image may contain multiple objects, and the goal is to identify all relevant objects present in the image.

In the contribution in Chapter 5 several methods for multilabel classification were implemented in the machine learning software *mlr* (Bischl et al., 2016) and compared in a benchmark experiment.

***Multivariate Regression*** – Multivariate regression is a powerful technique used to predict multiple continuous target variables simultaneously. Here the target space can be defined as $\mathcal{Y} = \mathbb{R}^d$. Instead of using separate regression models for each target variable, multivariate regression captures the relationships and dependencies between the targets, resulting in more accurate and comprehensive predictions. When applying multivariate regression to predict a person's Big Five personality traits (openness, conscientiousness, extraversion, agreeableness, and neuroticism), it gives us a comprehensive understanding of their personality by considering how these traits influence each other (Omheni et al., 2014).

***Multioutput Prediction*** – Multi-output prediction is considered the most generalized and flexible form of learning, as it allows the target variables to be of mixed types. Here the target space can be defined as $\mathcal{Y} = \mathcal{Y}_1 \times ... \times \mathcal{Y}_d$, where $\mathcal{Y}_i$ can be either $\{0, 1\}$, $\{1, ..., m\}$, or $\mathbb{R}$, for $i = 1, ..., d$.

In psychological research, mixed target variables are relevant, for example, predicting personality and demographic traits based on behavioral data. Predicting personality through regression, gender through classification, and age through ordinal regression simultaneously provides valuable insights compared to predicting them independently. Traits like gender and age have been found to be related to personality, making multi-output prediction highly useful in such contexts (Goldberg et al., 1998).

### 3.1.3  Evaluation Measures

The performance of machine learning models (James et al., 2013) are evaluated by comparing their predictions $\hat{f}(x)$ with the actual values $y$ through the means of a loss function $L$ (see Eq. 3.3). For the one dimensional ($d = 1$) tasks, the dataset consists of observations $\mathbf{x}^{(i)}$ and scalar target variables $y^{(i)}$. For regression tasks, typical performance metrics are computed by using the estimated model's residuals $\hat{\epsilon} = y - \hat{f}(\mathbf{x})$. A frequently used measure is the mean squared error (MSE), which uses the $L_2$ loss function: $L_2(\hat{f}(\mathbf{x}), y) = (y - \hat{f}(\mathbf{x}))^2$. For a given test dataset $\mathcal{D}_{\text{test}} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^{m}$ the MSE can be calculated as

$$\widehat{MSE} = \sum_{i=1}^{m} \left( y^{(i)} - \hat{f}(\mathbf{x}^{(i)}) \right). \tag{3.6}$$

For classification tasks, one commonly used metric for is the accuracy score $ACC$, which calculates the proportion of correctly classified instances out of the total number of instances in the test dataset $\mathcal{D}_{\text{test}}$. The loss function here is the 0/1-loss: $L_{0/1}(\hat{f}(\mathbf{x}), y) = \mathbb{1}_{\hat{f}(\mathbf{x}) \neq y}$. The accuracy score can thus be calculated as

$$\widehat{ACC} = 1 - \frac{1}{m} \sum_{i=1}^{m} \mathbb{1}_{\hat{f}(\mathbf{x}^{(i)}) \neq y^{(i)}}. \tag{3.7}$$

In practice, $ACC$ is often used as a primary metric to evaluate binary classifiers, especially when the classes are balanced (i.e., have approximately equal representation in the dataset). However, it is essential to consider other metrics, such as precision, recall, and F1-score, especially when dealing with imbalanced datasets, where one class might dominate the other in terms of instances. These additional metrics provide a more comprehensive evaluation of the model's performance in handling different classes and help identify potential biases or limitations in the classifier.

In the case $d > 1$, dealing with multiple target variables makes evaluating performance more challenging. In general, the actual target values $\mathbf{y}^{(i)} = (y_1^{(i)}, ..., y_d^{(i)})^\mathsf{T}$ are compared with the predicted target values $\hat{\mathbf{y}}^{(i)} = (\hat{y}_1^{(i)}, ..., \hat{y}_d^{(i)})^\mathsf{T}$. Various performance measures have been developed specifically for multilabel classification and multivariate regression problems (Borchani et al., 2015; Zhang and Zhou, 2013).

For multi-label classification, an example of a performance metric is the Hamming loss $HL$, which compares predicted labels with actual labels. For one observation, this can be calculated as

$$\widehat{HL} = \frac{1}{d} \sum_{i=1}^{d} \mathbb{1}_{\{y_i \neq \hat{y}_i\}}. \tag{3.8}$$

This value is calculated instance-wise, and the performance of a test set is the mean Hamming loss of each instance.

For multivariate regression, an example of a performance metric is the multivariate mean squared error ($MMSE$), which is the mean $MSE$ of every target. For one instance, this can be calculated as

$$MMSE = \frac{1}{d} \sum_{i=1}^{d} (y_i - \hat{y}_i)^2. \tag{3.9}$$

Again, the performance of a test dataset is calculated by the mean of each $MMSE$ score of each instance.

In the more generalized multi-output prediction problem, however, calculating a single performance value from possible mixed target spaces is not straightforward. Since datasets with mixed target spaces can vary significantly and may involve both classification and regression tasks during evaluation, a general definition of a performance metric is impractical and should be left to the user's discretion.

### 3.1.4 $R^2$ as an Evaluation Measure

The coefficient of determination, commonly represented as $R^2$, plays a fundamental role in regression analysis (Montgomery et al., 2021). It quantifies the extent to which the variance in the dependent variable can be anticipated by the independent variables. $R^2$ takes values between 0 and 1, where 0 indicates that the model does not explain any of the variability in the dependent variable, and 1 indicates that the model perfectly explains all the variability. $R^2$ is defined as:

$$R^2 = 1 - \frac{\sum_{i=1}^{n} (y_i - \hat{y}_i)^2}{\sum_{i=1}^{n} (y_i - \bar{y})^2}, \tag{3.10}$$

where $\bar{y}$ is the mean value of the dependent variable for all data points. In the context of machine learning, $R^2$ can be used as a performance measure for regression models, but it can indeed take on negative values in certain scenarios. This can happen when the model's predictions are worse than simply using the mean of the dependent variable as the prediction for all data points (Alexander et al., 2015). When $R^2$ is negative, it indicates that the model's performance is worse than a simple baseline that predicts the mean of the dependent variable. This can occur due to overfitting, underfitting, or dataset shift, where the model doesn't generalize well to new data. It emphasizes the importance of assessing a model's performance on unseen data that shares similar characteristics with the training data.

## 3.2   Interpretable Machine Learning

Interpretable Machine Learning (IML) is a vital aspect of the PhoneStudy project, enabling a deeper understanding of the predictive models and their impact on psychological attribute prediction based on smartphone data. IML methods play a central role in enhancing model transparency and accountability, ensuring the trustworthiness of the predictions in real-world applications (Lundberg and Lee, 2017a; Ribeiro et al., 2016). Researchers often prefer simpler models due to their ease of interpretation and understandability. These models, such as linear regression, decision trees, or logistic regression, have transparent structures that allow users to understand how the input features contribute to the final predictions. This interpretability is especially beneficial when researchers aim to gain insights into the relationships between variables and to communicate the model's findings effectively. Furthermore, in critical domains like medicine or criminology, predictive models are frequently employed to make decisions that directly impact human life (Caruana et al., 2015). In such high-stakes applications, the adoption of interpretable models becomes essential to offer clear explanations for their decisions. Interpretable models play a significant role in improving the understanding of prediction processes and ensuring accountability in these sensitive fields.

Complex black box models, like deep neural networks or ensemble methods, may achieve higher predictive accuracy but are challenging to interpret. These models are often used in certain use cases where interpretability may not be the primary concern. For example, in recommendation systems or image recognition tasks, the focus is often on maximizing predictive accuracy without requiring detailed explanations for each recommendation or classification (Doshi-Velez and Kim, 2017).

The desire to explain predictions and outcomes from any model has led to the development of model-agnostic methods. Model-agnostic techniques are designed to be applicable to any type of predictive model. These methods focus on providing post hoc explanations for individual predictions, enabling users to understand the factors that influenced a specific outcome.

### 3.2.1 Model-Agnostic Methods

In the context of IML, model-agnostic methods focus on making any prediction model interpretable (Molnar, 2022). These methods are not tied to specific types of models and can be applied universally to provide explanations for any machine learning algorithm. The area of model-agnostic methods is wide-ranging, including both global and local approaches. Global model-agnostic methods, such as the partial dependence plot (Friedman, 2001b), accumulated local effects plot (Apley and Zhu, 2020a), and permutation feature importance (Fisher et al., 2019), offer insights into the overall behavior of the model and the significance of features on a global scale. Local methods in the context of model-agnostic techniques focus on providing explanations for individual predictions, offering insights into how specific input features influence the model's output for a particular instance. These methods, such as individual conditional expectation (ICE) plots (Goldstein et al., 2015), local interpretable modelagnostic explanations (LIME) (Ribeiro et al., 2016), or the use of shapley values (SHAP) (Lundberg and Lee, 2017b) help understand the model's behavior at a more granular level, providing valuable insights into individual predictions.

***Performance-Based Feature Importance*** – The technique of random forest permutation feature importance was first introduced in the seminal paper by Breiman (2001). This method assesses the importance of each feature in a random forest model by measuring the extent of predictive performance reduction when the values of a particular feature are randomly shuffled. This idea of calculating the change in performance, when breaking the link between a feature and the target value, can easily be adapted for the use of any machine learning model. Fisher et al. (2019) first implemented a model agnostic permutation based feature importance called *model reliance*.

Another idea to measure a feature's importance stems from game theory. The concept of Shapley values originated from cooperative game theory and was first introduced by Shapley (1953). Shapley values provide a principled and fair method for allocating contributions among individual features within a predictive model. Unlike other feature importance methods that focus solely on the impact of each feature in isolation, Shapley values consider the cooperative interactions between features, taking into account all possible combinations of features and their effects on predictions (Štrumbelj and Kononenko, 2013; Lundberg and Lee, 2017b).

***Feature Effects*** – Interpreting feature effects in linear regression models is straightforward, as the feature effect is represented by the regression coefficient of that specific feature. A positive coefficient indicates that an increase in the feature's value leads to an increase in the target variable, while a negative coefficient implies the opposite effect.

On the local level, Individual Conditional Expectation (ICE) curves provide valuable insights into the effects of a specific feature on individual predictions (Goldstein et al., 2015). ICE curves illustrate how the model's output changes for an instance as the value of the selected feature

varies while keeping other features constant. By visualizing multiple ICE curves, one for each instance in the dataset, users can understand the feature's impact on different predictions individually.

For a global perspective, Partial Dependence (PD) plots, introduced by Friedman (Friedman, 2001b), are commonly used. PD plots visualize the marginal relationship between a feature and the target variable across its range of feature values while averaging out the effects of other features. This provides a comprehensive understanding of the overall behavior of the feature in relation to the target variable on a global scale. Another global method is Accumulated Local Effects (ALE) plots (Apley and Zhu, 2020b), which offer similar insights to PD plots but are more suitable for dealing with interactions between features. ALE plots also provide a smoothed view of the feature's effect on the target variable, enabling a clearer understanding of its impact across various regions of the feature space.

### 3.2.2   Challenges in IML in the PhoneStudy Project

In the PhoneStudy project, the prediction of psychological attributes through machine learning algorithms presents challenges within the realm of Interpretable Machine Learning (IML). One such challenge arises from the nature of smartphone data, where each participant's data represents only one observation, resulting in a relatively low number of observations compared to the high number of features captured by smartphones. Moreover, model-agnostic methods typically focus on the individual feature level, making it difficult to handle high correlations between features. However, a promising way to address this is by categorizing these features into meaningful groups, such as calling behavior, music listening, and movement.

The contribution in Chapter 4 deals with quantifying and visualizing the effect of feature groups. It provides a comprehensive overview of model-agnostic techniques for groups of features and introduces an importance-based sequential procedure to identify well-performing combinations of feature groups. Additionally, it introduces the combined features effect plot, which visualizes the effect of a feature group using sparse, linear combinations of features. These approaches significantly improve the interpretability of machine learning models, especially when handling complex datasets, where features can be grouped.

# Part I - Advancing Machine Learning Approaches

# Grouped Feature Importance and Combined Features Effect Plot

This article focuses on assessing the importance and visualizing the effect of feature groups in machine learning with model-agnostic techniques. It also introduces an importance-based sequential procedure, and a visualization technique, which are validated through simulations and real data.

***Contributing article***

Au, Q., Herbinger, J., Stachl, C., Bischl, B., and Casalicchio, G. (2022). Grouped feature importance and combined features effect plot. *Data Mining and Knowledge Discovery*, 36(4):1401–1450

***Copyright information***

***Declaration of contributions***

The PhD candidate made significant contributions to the manuscript creation by developing and implementing methodologies to quantify the importance of groups of variables, which were determined using various methods, including permutation and retraining. The PhD candidate also implemented the forward search method to identify combinations of groups that had a positive impact on prediction accuracy. Additionally, the combined feature effect plot to better understand and interpret the effects of variables within identified groups was also a contribution made by the candidate. Furthermore, the PhD candidate demonstrated the usefulness of these methodologies by providing simulations and a real data example. Moreover, the candidate provided a mathematical proof for a property of the grouped Shapley method used to determine the importance of groups.

*Contribution of the coauthors*

Julia Herbinger also made a significant contribution to the manuscript and shares the first authorship with Jiew-Quay Au at a 50% split. In addition to her large role in collaborating on the

manuscript and conducting research, she also performed simulations to demonstrate additional properties of the implemented methods. Giuseppe Casalicchio provided significant assistance in selecting and developing the various methods, as well as in conducting simulations. Bernd Bischl was instrumental in offering insightful discussions and valuable guidance throughout the whole project. Clemens Stachl supported the team in research and in specific sections of the manuscript, particularly in the real data example.

# Grouped feature importance and combined features effect plot

**Quay Au[1] · Julia Herbinger[1]  · Clemens Stachl[2] · Bernd Bischl[1] · Giuseppe Casalicchio[1]**

## Abstract

Interpretable machine learning has become a very active area of research due to the rising popularity of machine learning algorithms and their inherently challenging interpretability. Most work in this area has been focused on the interpretation of single features in a model. However, for researchers and practitioners, it is often equally important to quantify the importance or visualize the effect of feature groups. To address this research gap, we provide a comprehensive overview of how existing model-agnostic techniques can be defined for feature groups to assess the grouped feature importance, focusing on permutation-based, refitting, and Shapley-based methods. We also introduce an importance-based sequential procedure that identifies a stable and well-performing combination of features in the grouped feature space. Furthermore, we introduce the combined features effect plot, which is a technique to visualize the effect of a group of features based on a sparse, interpretable linear com-

✉ Julia Herbinger
    julia.herbinger@stat.uni-muenchen.de

    Quay Au
    quayau@gmail.com

    Clemens Stachl
    clemens.stachl@unisg.ch

    Bernd Bischl
    bernd.bischl@stat.uni-muenchen.de

    Giuseppe Casalicchio
    giuseppe.casalicchio@stat.uni-muenchen.de

[1]  Department of Statistics, Ludwig-Maximilians-University Munich, 80539 Munich, Germany

[2]  Institute of Behavioral Science and Technology, University of St. Gallen, 9000 St. Gallen, Switzerland

 Springer

bination of features. We used simulation studies and real data examples to analyze, compare, and discuss these methods.

**Keywords** Grouped feature importance · Combined features effects · Dimension reduction · Interpretable machine learning

## 1 Introduction

Machine learing (ML) algorithms are nowadays used in many diverse fields e.g. in medicine (Shipp et al. 2002), criminology (Berk et al. 2009), and increasingly in the social sciences (Stachl et al. 2020b; Yarkoni and Westfall 2017). Interpretable models are paramount in many high-stakes settings, such as medical and juridical applications (Lipton 2018). However, well-performing ML models often bear a lack of interpretability. In the context of interpretable ML (IML) research, several model-agnostic methods to produce explanations for single features have been developed (Molnar 2019). Examples include the permutation feature importance (PFI; Fisher et al. 2019), leave-one-covariate out (LOCO) importance (Lei et al. 2018), SHAP values (Lundberg and Lee 2017), or partial dependence plots (PDP; Friedman 2001).

In many applications, it can be more informative to produce explanations for the importance or effect of a group of features (which we refer to as grouped interpretations) rather than for single features. It is important to note that the meaning of grouped interpretations, in general, differs from single feature interpretations, and resulting interpretations are usually not directly comparable (e.g., as Gregorutti et al. (2015) shows for the permutation feature importance). Hence, our aim is not to challenge single feature interpretations as both single and grouped feature interpretation methods measure different things and are useful on their own.

Grouped interpretations might be especially interesting for high-dimensional settings with hundreds or thousands of features. In particular, when analyzing the influence of these features visually (e.g., by plotting the marginal effect of a feature on the target) on a single feature level, this might result in an information overload which might not provide a comprehensive understanding of the learned effects (Molnar et al. 2020b). Furthermore, the runtime of some interpretation methods—such as Shapley values—does not scale linearly in the number of features. Hence, calculating them on a single feature level might not be computationally feasible for high-dimensional settings, making grouped computations a feasible remedy (Lundberg and Lee 2017; Covert et al. 2020; Molnar et al. 2020b).

From a use case perspective, the concept of grouped interpretations is particularly useful when the feature grouping is available *a priori* based on the application context. In that sense, features that either belong to the same semantic area (e.g., behaviors in psychology or biomarkers in medicine) or are generated by the same mechanism or device (e.g., fMRI, EEG, smartphones) can be grouped together to assess their joint effect or importance. For example, in our application in Sect. 7, we use a real-world use case from psychology that studies how the human behavior on smartphone app usage is associated to different personality traits (Stachl et al. 2020a). Features were extracted from longitudinal data collected from smartphones of 624 participants,

**Fig. 1** A possible process from group definition to grouped interpretations. First, the feature groups must be defined. A model is then fitted, typically on the feature space where the information of the pre-defined grouping might be used (e.g., if the fitting process is combined with a feature selection procedure) or ignored. When the best model is found, model-agnostic grouped interpretation methods are applied on the previously defined feature groups. A commonly used approach is to first obtain an overview of which groups are most important for achieving a good model performance (grouped feature importance) to subsequently analyze how the most important feature groups influence the model's prediction (grouped feature effect) (Color figure online)

and can be grouped into different behavioral classes (i.e., communication and social activity, app-usage, music consumption, overall phone activity, mobility). Another example is applications with sensor data (Chakraborty and Pal 2008), where multiple features measured by a single sensor naturally belong together, and hence grouped interpretations on sensor-level might be more informative.

There are also situations where the interpretation of single features might be misleading and where grouped interpretations can provide a remedy. Examples include datasets with time-lagged or categorical features (e.g., dummy or one-hot encoded categories) and the presence of feature interactions (Gregorutti et al. 2015). A concrete example for dummy encoded categorical features is shown in Appendix A.

Even in situations where feature groups are not naturally given in advance, it still might be beneficial to define groups in a data-driven manner and apply interpretation methods on groups of features (for examples, see Sect. 1.2).

Hence, compared to single feature interpretation methods, the grouping structure must be defined beforehand. A possible process—from group membership definition to modeling up to post-hoc interpretations—is illustrated in Fig. 1. Since defining the underlying group structure is a relevant step in this process, we discuss some applied techniques on how to find groups of features in Sect. 1.2. However, in this paper, we focus on the interpretation component once the groups are known (the green part in Fig. 1).

Although the grouped feature perspective is relevant in many applications, most IML research has focused on methods that attempt to provide explanations on a single-feature level. Model-agnostic methods for feature groups are rare and not well-studied.

## 1.1 Real data use cases with grouped features

In the following we summarize further exemplary predictive tasks with pre-specified feature groupings. These tasks will also be used in Sect. 3.4 for further empirical analysis. For more details on features and associated groups see Table 1.

*Heat value of fossil fuels* In this small scale regression task ($n = 129$), the objective is to predict the heat value of fossil fuels from spectral data (Fuchs et al. 2015). In addition to one scalar feature (humidity), the dataset contains two groups of curve data, the first from the ultraviolet-visible spectrum (UVVIS) and the second from the near infrared spectrum (NIR).

**Table 1** Real world datasets with grouped features and their pre-specified group memberships

| Dataset | Single features | Group membership | Description |
|---|---|---|---|
| *Birthweight* | age1, age2, age3 | Age | Mother's age represented by 3 orthogonal polynomials |
| | lwt1, lwt2, lwt3 | lwt | Mother's weight represented by 3 orthogonal polynomials |
| | White, black | Race | Mother's race (indicator functions) |
| | Smoke | Smoke | Smoking status (indicator function) |
| | ptl1, ptl2m | ptl | One, or two or more previous premature labors |
| | ht | ht | History of hypertension (indicator function) |
| | ui | ui | Presence of uterine irritability (indicator function) |
| | ftv1, ftv2, ftv3m | ftv | One, two, or three or more physician visits during first trimester |
| *Colon* | x1, ..., x5 | Gene1 | Gene expression data for gene 1 |
| | ⋮ | ⋮ | ⋮ |
| | x96, ..., x100 | Gene20 | Gene expression data for gene 20 |
| *Fuelsubset* | H20 | H20 | Humidity in percent |
| | UVVIS1, ..., UVVIS134 | UVVIS | Data from the ultraviolet-visible spectrum (134 wavelength points) |
| | NIR1, ..., NIR231 | NIR | Data from the near infrared spectrum (231 wavelength points) |

*Birthweight* The *birthweight* dataset has data on 189 births at the Baystate Medical Centre in Massachusetts during 1986 (Venables and Ripley 2002). The objective is to predict the birth weight in kilograms from a set of 16 features, some of which are grouped (e.g., dummy encoded categorical features).

*Colon cancer* The *colon* dataset contains gene expression data of 20 genes (5 basis B-Splines each) for 62 samples from microarray experiments of colon tissue (Alon et al. 1999). The task is to predict cancerous tissue from the resulting 100 predictors.

## 1.2 Grouping procedures

Following the definitions of He and Yu (2010), we provide a brief overview of different procedures to define feature groups in a knowledge-driven and data-driven manner. In data-driven grouping, an algorithmic approach such as clustering or density estimation is used to define groups of features. Knowledge-driven grouping, on the other hand, uses domain knowledge to define the grouping structure of features. Throughout our

paper, we mainly assume a user-defined grouping structure. However, all methods introduced in this paper should also be compatible with an appropriate data-driven method if the defined groups have a meaningful interpretation.

### *Data-driven grouping*

One method to group features in a data-driven manner is to use clustering approaches such as hierarchical clustering (Park et al. 2006; Toloşi and Lengauer 2011; Rapaport et al. 2008) or fuzzy clustering (Jaeger et al. 2003). These approaches often work well in highly correlated feature spaces, such as in genomics or medicine, where correlated features are grouped together so that no relevant information is discarded (Toloşi and Lengauer 2011). For instance, Jaeger et al. (2003) tackles a feature selection problem for a high-dimensional and intercorrelated feature space when working with microarray data. To simultaneously select informative and distinct genes, they first apply fuzzy clustering to obtain groups of similar genes from microarray data. Next, the informative representatives of each group are selected based on a suitable test statistic. The disadvantage of data-driven grouping is that groups depend only on the statistical similarity between features, which might not coincide with domain-specific interpretations (Chakraborty and Pal 2008).

### *Knowledge-driven grouping*

Knowledge-driven group formation has the advantage that the dimensionality reduction might lead to better interpretability than the data-driven path. Gregorutti et al. (2015) apply a knowledge-driven approach in the context of multiple functional data analysis, where they then select groups for subsequent modeling based on their group importance values. Chakraborty and Pal (2008) also select groups of features, where data from one sensor (e.g., to capture satellite images in different spectral bands) represents a group. Hence, features are grouped based on their topical character (e.g., measurement device) rather than their shared statistical properties. Another use case of knowledge-driven grouping is described in Lozano et al. (2009), who group time-lagged features of the same time series for gene expression data. They use the given grouping structure in a group feature selection procedure and apply group LASSO as well as a boosting method.

### 1.3 Related work

A well-known model that handles groups of features is the *group LASSO* (Yuan and Lin 2006), which extends the LASSO (Tibshirani 1996) for feature selection based on groups. Moreover, other extensions—e.g., to obtain sparse groups of features (Friedman et al. 2010), to support classification tasks (Meier et al. 2008) or non-linear effects (Gregorova et al. 2018)—also exist. However, group LASSO is a modeling technique that focuses on selecting groups in the feature space rather than quantifying their importance.

A large body of research already exists regarding the importance of individual features (see, e.g., Fisher et al. 2019; Hooker and Mentch 2019; Scholbeck et al. 2020). Hooker and Mentch (2019) distinguish between two loss-based feature importance approaches, namely permutation methods and refitting methods. Permutation meth-

ods measure the increase in expected loss (or error) after permuting a feature while the model remains untouched. Refitting methods measure the increase in expected loss after leaving out the feature of interest completely and hence require refitting the model (Lei et al. 2018). Since the model remains untouched in the former approach, interpretations refer to a specific fitted model, while interpretations for refitting methods refer to the underlying ML algorithm. Gregorutti et al. (2015) introduced a model-specific, grouped PFI score for random forests and applied this approach to functional data analysis. Valentin et al. (2020) introduced a model-agnostic grouped version of the model reliance score (Fisher et al. 2019). However, they focus more on the application and omit a detailed theoretical foundation. Recently, a general refitting framework to measure the importance of (groups of) features was introduced by Williamson et al. (2020). In their approach, the feature importance measurement is detached from the model level and defined by an algorithm-agnostic version to measure the intrinsic importance of features. The importance score is defined as the difference between the performance of the full model and the performance based on all features *except* the group of interest.

Permutation methods can be computed much faster than refitting methods. However, the PFI, for example, has issues when features are correlated and interact in the model due to extrapolation in regions without any or just a few observations (Hooker and Mentch 2019). Hence, interpretations in these regions might be misleading. To avoid this problem, alternatives based on conditional distributions or refitting have been suggested (e.g., Strobl et al. 2008; Nicodemus et al. 2010; Hooker and Mentch 2019; Watson and Wright 2019; Molnar et al. 2020a). Although the so-called conditional PFI provides a solution to this problem, its interpretation is different and "must be interpreted as the additional, unique contribution of a feature given all remaining features we condition on were known" (Molnar et al. 2020a). This property complicates the comparison with non-conditional interpretation methods. Therefore, we do not consider any conditional variants in this paper.

A third class of importance measures is based on Shapley values (Shapley 1953), a theoretical concept of game theory. The SHAP (Lundberg and Lee 2017) approach quantifies the contribution of each feature to the predicted outcome and is a permutation-based method. It has the advantage that contributions of interactions are distributed fairly between features. Besides being computationally more expensive, SHAP itself is based on the model's predicted outcome rather than the model's performance (e.g., measured by the model's expected loss). Casalicchio et al. (2019) extended the concept of Shapley values to fairly distribute the model's performance among features and called it Shapley Feature IMPortance (SFIMP). A similar approach called SAGE has also been proposed by Covert et al. (2020), who showed the benefits of the method on various simulation studies. One approach that uses Shapley values to explain grouped features was introduced by de Mijolla et al. (2020). However, instead of directly computing Shapley importance on the original feature space, they first apply a semantically-meaningful latent representation (e.g. by projecting the original feature space into a lower dimensional latent variable space using disentangled representations) and compute the Shapley importance on the resulting latent variables. Williamson and Feng (2020) mention that their feature importance method based on Shapley values can also be extended to groups of features. Additionally,

Amoukou et al. (2021) investigated grouping approaches for Shapley values in the case of encoded categorical features and subset selection of important features for tree-based methods. The calculation of Shapley values on groups of features based on performance values has only been applied with regard to feature subset selection methods and not for interpretation purposes (Cohen et al. 2005; Tripathi et al. 2020).[1]

After identifying which groups of features are important, the user is often interested in how they (especially the important groups) influence the model's prediction. Several techniques to visualize single-feature effects exist. These include partial dependence plots (PDP) (Friedman 2001), individual conditional expectation (ICE) curves (Goldstein et al. 2013), SHAP dependence plots (Lundberg et al. 2018), and accumulated local effects (ALE) plots (Apley and Zhu 2019). However, in the case of high-dimensional feature spaces, it is often not feasible to compute, visualize, and interpret single-feature plots for all (important) features. If features are grouped, visualization techniques become computationally more complex, and it may become even harder to visualize the results in an easily interpretable way. In the case of low-dimensional feature spaces, this might still be feasible, for example by using two-feature PDPs or ALE plots. Recently, effect plots that visualize the combined effect of multiple features have been introduced by Seedorff and Brown (2021) and Brenning (2021). They use principal component analysis (PCA) to reduce the dimension of the feature space and calculate marginal effect curves for the principal components. However, the employed dimension reduction method does not include information about the target variable and lacks sparsity (and hence, interpretability).

### 1.4 Contribution

Our contributions can be summarized as follows: We extend the permutation-based and refitting-based grouped feature importance methods introduced by Valentin et al. (2020) and Williamson et al. (2020) by comparing these methods to not only the full model (i.e., taking into account all features), but also to a null model (i.e., ignoring all features). Hence, we can quantify to what extent a group itself contributes to the prediction of a model without the presence of other groups. Furthermore, we introduce Shapley importance for feature groups and describe how these scores can be decomposed into single-feature importance scores of the respective groups. Our main contributions are: (1) We define a new algorithm to sequentially add groups of features depending on their importance, thereby enabling identification of well-performing combinations of groups. (2) We compare all grouped feature importance methods with respect to the main challenges that arise when applying these methods by creating small simulation examples. Subsequently, we provide recommendations for using and interpreting the respective methods correctly. (3) We introduce a model-agnostic method to visualize the joint effect of a group of features. To that end, we use a suitable dimension reduction technique and the conceptual idea of PDPs to calculate and plot the mean prediction of a sparse group of features with regard to their linear

---

[1] Feature subset selection methods usually aim to find sparse, well-performing feature combinations. Hence, the intended purpose of employing these methods is not to produce interpretability, but rather to generate a sufficient performance with fewer features.

combination. This novel method finally enables the user to visualize effects for groups of features. Finally, we showcase the usefulness of all these methods in real data examples.

The structure of this paper is as follows: First, we provide some general notation and definitions in Sect. 2. We formally define the grouped feature importance methods and introduce the sequential grouped feature importance procedure in Sects. 3 and 4, respectively. We compare these methods for different scenarios in Sect. 5. In Sect. 6, we introduce the combined features effect plot (CFEP) to visualize the effects of feature groups based on a supervised dimension reduction technique. Moreover, we also show the suitability of this technique compared to its unsupervised counterpart in a simulation study. Finally, in Sect. 7, all methods are applied to a real data example before summarizing and offering an outlook for future research in Sect. 8.

## 2 Background and notation

Analogous to Casalicchio et al. (2019), we use the term *feature importance* to refer to the influence of features on a model's predictive performance, which we measure by the expected loss when we perturb these features in a permutation approach or remove these features in a refitting approach.

### 2.1 General notation

Consider a $p$-dimensional feature space $\mathcal{X} = (\mathcal{X}_1 \times \cdots \times \mathcal{X}_p)$ and a one-dimensional target space $\mathcal{Y}$. The corresponding random variables that are generated from these spaces are denoted by $X = (X_1, \ldots, X_p)$ and $Y$. We denote a ML prediction function that maps the $p$-dimensional feature space to a one-dimensional target space by $\hat{f} : \mathcal{X} \to \mathbb{R}$ for regression tasks.[2] ML algorithms try to learn this functional relationship using $n \in \mathbb{N}$ i.i.d. observations drawn from the joint space $\mathcal{X} \times \mathcal{Y}$ with unknown probability distribution $\mathcal{P}$. The resulting dataset is denoted by $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$, where the vector $\mathbf{x}^{(i)} = (x_1^{(i)}, \ldots, x_p^{(i)})^\intercal \in \mathcal{X}$ is the $i$-th observation associated with the target variable $y^{(i)} \in \mathcal{Y}$. The $j$-th feature is denoted by $\mathbf{x}_j = (x_j^{(1)}, \ldots, x_j^{(n)})^\intercal$, for $j = 1, \ldots, p$. The dataset $\mathcal{D}$ can also be written in matrix form:

$$\begin{pmatrix} x_1^{(1)} & \ldots & x_p^{(1)} & y^{(1)} \\ \vdots & \ddots & \vdots & \vdots \\ x_1^{(n)} & \ldots & x_p^{(n)} & y^{(n)} \end{pmatrix} = (\mathbf{X}, \mathbf{Y}), \text{ with } \mathbf{X} = \begin{pmatrix} x_1^{(1)} & \ldots & x_p^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(n)} & \ldots & x_p^{(n)} \end{pmatrix}, \mathbf{Y} = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{pmatrix}. \quad (1)$$

The general error measure $\rho(\hat{f}, \mathcal{P}) = \mathbb{E}(L(\hat{f}(X), Y))$ of a learned model $\hat{f}$ is measured by a loss function $L$ on test data drawn independently from $\mathcal{P}$ and can be

---

[2] The target space is defined by $\mathbb{R}^g$ in the case of scoring classifiers with $g$ classes.

estimated using unseen test data $\mathcal{D}_{\text{test}}$ by

$$\hat{\rho}(\hat{f}, \mathcal{D}_{\text{test}}) = \frac{1}{|\mathcal{D}_{\text{test}}|} \sum_{(\mathbf{x}, y) \in \mathcal{D}_{\text{test}}} L(\hat{f}(\mathbf{x}), y). \tag{2}$$

The application of an ML algorithm (or *learner*) $\mathcal{I}$ to a given dataset $\mathcal{D}$ results in a fitted model $\mathcal{I}(\mathcal{D}) = \hat{f}_{\mathcal{D}}$. The *expected generalization error* of a learner $\mathcal{I}$ takes into account the variability introduced by sampling different datasets $\mathcal{D}$ of equal size $n$ from $\mathcal{P}$ and is defined by

$$GE(\mathcal{I}, \mathcal{P}, n) = \mathbb{E}_{|\mathcal{D}|=n}(\rho(\mathcal{I}(\mathcal{D}), \mathcal{P})). \tag{3}$$

In practice, resampling techniques such as cross-validation or bootstrapping on the available dataset $\mathcal{D}$ are used to estimate Eq. (3). Resampling techniques usually split the dataset $\mathcal{D}$ into $k \in \mathbb{N}$ training datasets $\mathcal{D}_{\text{train}}^i$, $i = 1, \ldots, k$, of roughly the same size $n_{\text{train}} < n$. Eq. (3) can be estimated by

$$\widehat{GE}(\mathcal{I}, \mathcal{D}, n_{\text{train}}) = \frac{1}{k} \sum_{i=1}^{k} \hat{\rho}(\hat{f}_{\mathcal{D}_{\text{train}}^i}, \mathcal{D}_{\text{test}}^i). \tag{4}$$

In the following, we often associate the set of numbers $\{1, \ldots, p\}$ in a one-to-one manner with the features $\mathbf{x}_1, \ldots, \mathbf{x}_p$ by referring a number $j \in \{1, \ldots, p\}$ as feature $x_j$. We call $G \subset \{1, \ldots, p\}$ a *group of features*.

## 2.2 Permutation feature importance (PFI)

Fisher et al. (2019) proposed a model-agnostic version of the PFI measure used in random forests (Breiman 2001). The PFI score of the $j-$th feature of a fitted model $\hat{f}$ is defined as the increase in expected loss after permuting feature $X_j$:

$$\text{PFI}_j(\hat{f}) = \mathbb{E}(L(\hat{f}(X_{[j]}), Y)) - \mathbb{E}(L(\hat{f}(X), Y)). \tag{5}$$

Here, $X_{[j]} = (X_1, \ldots, X_{j-1}, \tilde{X}_j, X_{j+1}, \ldots, X_p)$ is the $p$-dimensional random variable vector of features, where $\tilde{X}_j$ is an independent replication of $X_j$ following the same distribution. The idea behind this method is to break the association between the $j-$th feature and the target variable by permuting its feature values. If a feature is not useful for predicting an outcome, changing its values by permutation will not increase the expected loss.[3] For an accurate estimation of Eq. (5), we would need to calculate all possible permutation vectors over the index set $\{1, \ldots, n\}$ (see Casalicchio et al. (2019) for an in-depth discussion on this topic). However, Eq. (5) can be approximated on a dataset $\mathcal{D}$ with $n$ observations by Monte Carlo integration using $m$

---

[3] We consider the case of loss functions that are to be minimized. Hence, the larger $\text{PFI}_j$, the more substantial the increase in expected loss and the more important the $j-$th feature.

random permutations:

$$\widehat{\text{PFI}}_j(\hat{f}, \mathcal{D}) = \frac{1}{nm} \sum_{i=1}^{n} \sum_{k=1}^{m} \left( L\left( \hat{f}((x_1^{(i)}, \ldots, x_j^{(\tau_k^{(i)})}, \ldots, x_p^{(i)}), y^{(i)}) \right) - L\left( \hat{f}(\mathbf{x}^{(i)}, y^{(i)}) \right) \right), \tag{6}$$

where $\tau_k$ is a random permutation vector of the index set $\{1, \ldots, n\}$ for $k = 1, \ldots, m$ permutations.[4]

Equation (6) could also be embedded into a resampling technique, where the permutation is always applied on the held-out test set of each resampling iteration (Fisher et al. 2019). However, this leads to refits and is computationally more expensive. The resulting resampling-based PFI of a learner $\mathcal{I}$ is estimated by

$$\widehat{\text{PFI}}_j^{\text{res}}(\mathcal{I}, \mathcal{D}, n_{\text{train}}) = \frac{1}{k} \sum_{i=1}^{k} \widehat{\text{PFI}}_j(\hat{f}_{\mathcal{D}_{\text{train}}^i}, \mathcal{D}_{\text{test}}^i), \tag{7}$$

where the permutation strategy is applied on the test sets $\mathcal{D}_{\text{test}}^i$.

## 3 Feature importance for groups

In our first minor contribution, we provide a general notation and formal definitions for grouped permutation and refitting methods and explain them by answering the following questions:

(a) To what extent does a group of features contribute to the model's performance in the presence of other groups?
(b) To what extent does a group itself increase the expected loss if it is added to a null model like the mean prediction of the target for refitting methods?
(c) How can we fairly distribute the expected loss among all groups and all features within a group?

The definitions of all grouped feature importance scores are based on loss functions. They are defined in such a way that important groups will yield positive grouped feature importance scores. The question of how to interpret the differing results of these methods is addressed in Sect. 5.

### 3.1 Permutation methods

Here, we extend the existing definition of PFI to groups of features and introduce the GPFI (Grouped Permutation Feature Importance) and GOPFI (Group Only Permutation Feature Importance) scores. For ease of notation, we will only define these scores for a fitted model $\hat{f}$ (see Eq. 5).

---

[4] An example for $n = 3$ would be $\tau_1 = (1, 3, 2)^\mathsf{T}$ with $\tau_1^{(i)}$ being the $i-$th entry of that vector.

### 3.1.1 Grouped permutation feature importance (GPFI)

For the definition of GPFI—which is based on the definitions of Gregorutti et al. (2015) and Valentin et al. (2020)—let $G \subset \{1, \ldots, p\}$ be a group of features. Let $\tilde{X}_G = (\tilde{X}_j)_{j \in G}$ be a $|G|$-dimensional random vector of features, which is an independent replication of $X_G = (X_j)_{j \in G}$ following the same joint distribution. This random vector is independent of both the target variable and the random vector of the remaining features, which we define by $X_{-G} := (X_j)_{j \in \{1, \ldots, p\} \setminus G}$. With slight abuse of notation to index the feature groups included in $G$, we define the grouped permutation feature importance of $G$ as

$$\text{GPFI}_G = \mathbb{E}(L(\hat{f}(\tilde{X}_G, X_{-G}), Y)) - \mathbb{E}(L(\hat{f}(X), Y)). \tag{8}$$

Equation (8) extends Eq. (5) to groups of features so that the interpretation of GPFI scores always refers to the importance when the feature values of the group defined by $G$ are permuted jointly (i.e., without destroying the dependencies of the features within the group). Similar to Eq. (7), the grouped permutation feature importance can be estimated by Monte Carlo integration:

$$\widehat{\text{GPFI}}_G = \frac{1}{nm} \sum_{i=1}^{n} \sum_{k=1}^{m} \left( L(\hat{f}(\mathbf{x}_G^{(\tau_k^{(i)})}, \mathbf{x}_{-G}^{(i)}), y^{(i)}) - L(\hat{f}(\mathbf{x}^{(i)}, y^{(i)})) \right). \tag{9}$$

The GPFI measures the contribution of one group to the model's performance if all other groups are present in the model (see (a) from Sect. 3).

### 3.1.2 Group only permutation feature importance (GOPFI)

To evaluate the extent to which a group itself contributes to a model's performance (see (b) from Sect. 3), one can also use a slightly different measure. As an alternative to Eq. 9, we can compare the expected loss after permuting all features jointly with the expected loss after permuting all features except the considered group. We define this GOPFI for a group $G \subset \{1, ..., p\}$ as

$$\text{GOPFI}_G = \mathbb{E}(L(\hat{f}(\tilde{X}), Y)) - \mathbb{E}(L(\hat{f}(X_G, \tilde{X}_{-G}), Y)), \tag{10}$$

which can be approximated by

$$\widehat{\text{GOPFI}}_G = \frac{1}{nm} \sum_{j=1}^{n} \sum_{k=1}^{m} \left( L(\hat{f}(\mathbf{x}^{(\tau_k^{(j)})}, y^{(j)})) - L(\hat{f}(\mathbf{x}_G^{(j)}, \mathbf{x}_{-G}^{(\tau_k^{(j)})}), y^{(j)}) \right). \tag{11}$$

While the relevance of GOPFI as an importance measure might be limited, it is technically useful for the grouped Shapley importance (see Eq. 14).

## 3.2 Refitting methods

Here, we introduce two refitting-based methods for groups of features. The first definition is similar to the one introduced in Williamson et al. (2020).

### 3.2.1 Leave-one-group-out importance (LOGO)

For a subset $G \subset \{1, \ldots, p\}$, we define the reduced dataset $\tilde{\mathcal{D}} := \{(\mathbf{x}_{-G}^{(i)}, y^{(i)})\}_{i=1}^n$. Given a learner $\mathcal{I}$, which generates models $\mathcal{I}(\mathcal{D}) = \hat{f}_{\mathcal{D}}$ and $\mathcal{I}(\tilde{\mathcal{D}}) = \hat{f}_{\tilde{\mathcal{D}}}$, we define the Leave-One-Group-Out Importance (LOGO) as

$$LOGO(G) = \mathbb{E}(L(\hat{f}_{\tilde{\mathcal{D}}}(X_{-G}), Y)) - \mathbb{E}(L(\hat{f}_{\mathcal{D}}(X), Y)). \tag{12}$$

The LOGO can be estimated by using a learner $\mathcal{I}$ on $\tilde{\mathcal{D}}$ and should be embedded in a resampling technique:

$$\widehat{LOGO}(G) = \widehat{GE}(\mathcal{I}, \tilde{\mathcal{D}}, n_{\text{train}}) - \widehat{GE}(\mathcal{I}, \mathcal{D}, n_{\text{train}})$$
$$= \frac{1}{k} \sum_{i=1}^k \hat{\rho}(\hat{f}_{\tilde{\mathcal{D}}_{\text{train}}^i}, \tilde{\mathcal{D}}_{\text{test}}^i) - \frac{1}{k} \sum_{i=1}^k \hat{\rho}(\hat{f}_{\mathcal{D}_{\text{train}}^i}, \mathcal{D}_{\text{test}}^i).$$

Consequently, we compare the increase in expected loss compared to the full model's expected loss when leaving out a group of features and performing a refit (see (a) from Sect. 3).

While GPFI can be calculated with a resampling-based strategy by using refits to receive the algorithm-based instead of model-based GPFI, the meaning still varies from LOGO. For the algorithm-based GPFI, we calculate for each fitted model the importance score by permuting the regarded group and predicting with the same model. Then we average over all models from our resampling strategy and receive an importance score, which tells us how important a group of features is for some learner $\mathcal{I}$ when we break the association between this group and all other groups and the target. LOGO, on the other hand, leaves the group out and then performs the refit to calculate the importance of the group, and hence, it addresses the question: Can we remove this group from our dataset without reducing our model's performance? This is not answered by permutation-based methods.

### 3.2.2 Leave-one-group-in importance (LOGI)

While it may be too limiting to estimate the performance of a model based on one feature only, it can be informative to determine the extent to which a group of features (e.g., all measurements from a specific medical device) can reduce the expected loss in contrast to a null model (see (b) from Sect. 3). The Leave-One-Group-In (LOGI) method could be particularly helpful in settings where information on additional groups of measures will induce significant costs (e.g., adding functional imaging

data for a diagnosis) and/or limited resources are available (e.g., in order to be cost-covering, only one group of measures can be acquired). The LOGI method can also be useful for theory development in the natural and social sciences (e.g., which group of behaviors is most predictive by itself).

Let $\mathcal{I}_{\text{null}}$ be a null algorithm, which results in a null model $\hat{f}_{\text{null}}$ that only guesses the mean (or majority class for classification) of the target variable for any dataset. We additionally define a learner $\mathcal{I}$, which generates a model $\mathcal{I}(\mathring{\mathcal{D}}) = \hat{f}_{\mathring{\mathcal{D}}}$ for a dataset $\mathring{\mathcal{D}} := \{(\mathbf{x}_G^{(i)}, y^{(i)})\}_{i=1}^n$, which only contains features defined by $G \subset \{1, \ldots, p\}$. We define the $LOGI$ of a group $G$ as

$$LOGI(G) = \mathbb{E}(L(\hat{f}_{\text{null}}, Y)) - \mathbb{E}(L(\hat{f}_{\mathring{\mathcal{D}}}(X_G), Y)). \tag{13}$$

The LOGI can be estimated by using a learner $\mathcal{I}$ on $\mathring{\mathcal{D}} = \{(\mathbf{x}_G^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^n$ and should be embedded in a resampling technique:

$$\widehat{LOGI}(G) = \widehat{GE}(\mathcal{I}_{\text{null}}, \mathcal{D}, n_{\text{train}}) - \widehat{GE}(\mathcal{I}, \mathring{\mathcal{D}}, n_{\text{train}})$$
$$= \frac{1}{k}\sum_{i=1}^k \hat{\rho}(\hat{f}_{\text{null}}, \mathcal{D}_{\text{test}}^i) - \frac{1}{k}\sum_{i=1}^k \hat{\rho}(\hat{f}_{\mathring{\mathcal{D}}_{\text{train}}^i}, \mathring{\mathcal{D}}_{\text{test}}^i).$$

### 3.3 Grouped Shapley importance (GSI)

The importance measures defined above either exclude (or permute) individual groups of features from the total set of features or consider only the importance of groups by omitting (or permuting) all other features. The grouped importance scores are usually not affected if interactions within the groups are present. However, they can be affected if features from different groups interact, since permuting a group of features jointly destroys any interactions with other features outside the considered group. Therefore, we define the grouped Shapley importance (GSI) based on Shapley values (Shapley 1953). GSI scores account for feature interactions, as they measure the average contribution of a given group to all possible combinations of groups and fairly distribute the importance value caused by interactions among all groups (see (c) from Sect. 3).

We assume a set of distinct groups $\mathcal{G} = \{G_1, \ldots, G_l\}$, with $G_i \subset \{1, \ldots, p\}$, for $i = 1, \ldots, l$. In our grouped feature context, the value function $v : \mathcal{P}(\mathcal{G}) \longrightarrow \mathbb{R}$ assigns a "payout" to each possible group or combination of groups included in $\mathcal{G}$. With slight abuse of notation, we define the value function for a subset $S \subset \mathcal{G}$ as

$$v(S) := v\left(\cup_{G_i \in S} G_i\right).$$

We define the value function for a group $G \in \mathcal{G}$ calculated by a refitting or a permutation method by

$$v_{\text{refit}}(G) = LOGI(G) \quad \text{or} \quad v_{\text{perm}}(G) = GOPFI(G), \tag{14}$$

respectively. The marginal contribution of a group $G \in \mathcal{G}$, with $S \subset \mathcal{G}$ is

$$\Delta_G(S) = v(S \cup G) - v(S).$$

The GSI of the feature group $G$ is then defined as

$$\phi(G) = \sum_{S \subset \mathcal{G} \setminus G} \frac{(|\mathcal{G}| - 1 - |S|)! \cdot |S|!}{|\mathcal{G}|!} \Delta_G(S), \tag{15}$$

which is a weighted average of marginal contributions to all possible combinations of groups.

The GSI cannot always be calculated in a time-efficient way, because the number of coalitions $S \subset \mathcal{G} \setminus G$ can become large very quickly. In practice, the Shapley value is often approximated (Casalicchio et al. 2019; Covert et al. 2020) by drawing $M \leq |\mathcal{G}|!$ different coalitions $S \subset \mathcal{G} \setminus G$ and averaging the marginal, weighted contributions:

$$\hat{\phi}_M(G) = \frac{1}{M} \sum_{m=1}^{M} (|\mathcal{G}| - 1 - |S_m|)! \cdot |S_m|! \cdot \Delta_G(S_m), \tag{16}$$

with $S_m \subset \mathcal{G} \setminus G$, for all $m = 1, \ldots, M$.

The GSI can in general not be exactly decomposed into the sum of the Shapley importances for single features of the regarded group. In Appendix B, we show that the remainder term $R = \phi(G) - \sum_{i \in G} \phi(x_i)$ depends only on higher-order interaction effects between features of the regarded group and features of other groups. Hence, if one is interested in which features contributed most within a group, the Shapley importances for single features can be calculated, which provide a fair distribution of feature interactions within the group but not necessarily of feature interactions across groups. However, the remainder term can be used as a quantification of learned higher-order interaction effects between features of different groups.

While the GSI can be calculated with permutation- as well as refitting-based approaches, we will only apply the permutation-based approach in the upcoming simulation studies and the real-world example.

### 3.4 Real world use cases

For each dataset from Sect. 1.1, we fitted a random forest and summarized the three most important groups according to different grouped feature importance methods. For the importance scores of LOGI and LOGO, we used a 10-fold cross-validation (Table 2).

For the *birthweight* task, the feature **lwt** (mother's weight) was the most important group to predict the birthweight for all grouped feature importance methods except for LOGI. While all methods except LOGI also agree on the second most important group **ui** (presence of uterine irritability), feature groups differ for the third rank. However, this may also be due to statistical variability, as the importance values become very

🖄 Springer

**Table 2** Best 3 groups for each grouped feature importance score

| Dataset | GPFI | GOPFI | GSI | LOGI | LOGO |
|---|---|---|---|---|---|
| *Birthweight* | lwt (0.067) | lwt (0.056) | lwt (0.062) | ui (0.041) | lwt (0.036) |
| | ui (0.056) | ui (0.047) | ui (0.046) | Race (0.017) | ui (0.029) |
| | Smoke (0.009) | Race (0.045) | ptl (0.019) | ptl (0.015) | Race (0.005) |
| *Colon* | Gene14 (0.143) | Gene14 (0.174) | Gene14 (0.125) | Gene14 (0.128) | Gene14 (0.131) |
| | Gene10 (0.007) | Gene16 (0.087) | Gene16 (0.042) | Gene20 (0.045) | Gene17 (0.036) |
| | Gene7 (0.001) | Gene12 (0.057) | Gene13 (0.019) | Gene13 (0.028) | Gene18 (0.033) |
| *Fuelsubset* | NIR (30.51) | NIR (42.20) | NIR (36.21) | NIR (27.35) | NIR (8.34) |
| | UVVIS (2.85) | UVVIS (14.38) | UVVIS (7.99) | UVVIS (15.74) | $H_2O$ (0.14) |
| | $H_2O$ (0.01) | $H_2O$ (1.26) | $H_2O$ (0.24) | $H_2O$ ($-12.17$) | UVVIS ($-2.14$) |

For the classification task (*colon*) the scores were calculated as differences in classification accuracy. For the other two regression tasks the scores result from differences in MSE

small. It is interesting that **lwt**, despite being the most important group for all other scores, is not very important in terms of LOGI. Thus, **lwt** is less important as a stand-alone group, but appears important if the other feature groups are included in the model.

In the *colon* task, the feature group *gene14* is by far the most important group to predict cancerous tissue for all grouped feature important methods. However, there are variations in the second and third most important groups.

For the *fuelsubset* task, the permutation-based grouped importance methods (GPFI, GOPFI and GSI) show the same importance ranking for the three most important feature groups. However, for the refitting-based grouped importance methods (LOGI and LOGO), we can observe interesting differences. The features from the *UVVIS* group are important as a stand-alone group as can be seen by their positive LOGI score. However, the negative LOGO score of the *UVVIS* group indicates that the algorithm seems to perform better with only the *NIR* and *H2O* groups.

GPFI, GOPFI and GSI provide importance scores for feature groups of a given trained model without the necessity to refit the model. In contrast, LOGI and LOGO provide grouped importance scores based on the underlying algorithm and should always be considered together.

## 4 Sequential grouped feature importance

In general, feature groups do not necessarily have to be distinct or independent of each other. When groups partly contain the same or highly correlated features, we may obtain high grouped feature importance scores for similar groups. This can lead to misleading conclusions regarding the importance of groups. Quantifying the importance of different combinations of groups is especially relevant in applications where extra costs are associated with using additional features from other data sources. In this case, one might be interested in the sparsest, yet most important combination of groups or in understanding the interplay of different combinations of groups. Hence,

in practical settings, it is often important to decide which additional group of features to make available (e.g., buy or implement) for modeling and how groups should be prioritized under economic considerations.

Gregorutti et al. (2015) introduced a method called *grouped variable selection*, which is an adaptation of the recursive feature elimination algorithm from Guyon et al. (2002) and uses permutation-based grouped feature importance scores for the selection of feature groups. In Algorithm 1, we introduce a sequential procedure that is based on the idea of stability selection (Meinshausen and Bühlmann 2010). The procedure primarily aims at understanding the interplay of different combinations of groups by analyzing how the importance scores change after including other groups in a sequential manner. The feature groups must be pre-specified by the user. We prefer a refitting-based over a permutation-based grouped feature importance score when the secondary goal is to find well-performing combinations of groups. Here, the fundamental idea is to start with an empty set of features and to sequentially add the next best group in terms of LOGI until no further substantial improvement can be achieved. Our sequential procedure is based on a greedy forward search and creates an implicit ranking by showing the order in which feature groups are added to the model. To account for the variability introduced by the model, we propose to use repeated subsampling or bootstrap with sufficient repetitions (e.g., 100 repetitions).

To better understand Algorithm 1, we will demonstrate it with a small example with four groups $\mathcal{G} = \{G_1, G_2, G_3, G_4\}$ here. As a reminder, each group is a subset of $\{1, \ldots, p\}$, and we want to find a subset $B \subset \{1, \ldots, p\}$, which consists of the union of groups in $\mathcal{G}$. The subset $B$ is found by our sequential grouped feature importance procedure. To account for variability, the whole dataset is split into two sets (training and test set) repeatedly so that the train-test splits are different in each repetition of the resampling strategy (bootstrap or subsampling). For each training set, Algorithm 1 starts with an empty set $B = \emptyset$ (line 2, Algorithm 1). In line 5 of Algorithm 1, the candidate set $\mathcal{B} \subset \mathcal{P}(\mathcal{G})$ is defined as all subsets of the power set with cardinality 1. These are all individual groups $\mathcal{B} = \{\{G_1\}, \{G_2\}, \{G_3\}, \{G_4\}\}$. The LOGI score of each single group is then calculated. In our example, let $G_1$ have the highest LOGI score, which also exceeds the threshold $\delta$. The desired combination $B$ is preliminarily defined as $G_1$ (line 8), and for the comparison in the next step, the LOGI score of $G_1$ is defined as $L_0$ (line 9). Then, a new candidate set $\mathcal{B}$ is defined (line 11), which consists of all subsets of the power set of $\mathcal{G}$ of size $i$ (at this step, we have $i = 2$), where $B = G_1$ is also a subset of $\mathcal{B}$. Hence, $\mathcal{B} := \{\{G_1, G_2\}, \{G_1, G_3\}, \{G_1, G_4\}\}$. The LOGI score of elements of $\mathcal{B}$ is calculated as the LOGI score of the union of all subsets. Now, let $\widehat{LOGI}(G_1 \cup G3)$ have the highest score. This score is compared to the LOGI score of the previous iteration $L_0$ (line 13). Let the difference exceed the threshold $\delta$ for our example. In line 14 and 15, the desired combination $B$ is now defined as $G_1 \cup G_3$ and the LOGI score is again defined as $L_1$. Algorithm 1 now jumps to line 10 again with $i = 3$. The candidate set is now $\mathcal{B} = \{\{G_1, G_3, G_2\}, \{G_1, G_3, G_4\}\}$ (line 11). The LOGI scores are now calculated again for each element of $\mathcal{B}$. Let no LOGI score exceed $L_0$ by the threshold $\delta$ (line 13). Algorithm 1 now ends for this dataset split and returns $B = G_1 \cup G_3$ as the best combination. This procedure is repeated for each train-test split in each repetition.

---

**Algorithm 1:** Sequential Grouped Feature Importance

---

**input** : Set of groups $\mathcal{G} = \{G_1, ..., G_k\}$.
      Improvement threshold $\delta > 0$.
      Number of repetitions for the data splitting.
**output**: For every data split: a combination $B \subset \{1, ..., p\}$ and the order in which feature groups
      were added.

1  **for** *Every outer data split* **do**
2      Let $B = \emptyset$ **for** $i = 1, ..., k$ **do**
3         **if** $i = 1$ **then**
4            Define candidate set $\tilde{B} := \left\{ \tilde{G} \in \mathcal{P}(\mathcal{G}) \,\middle|\, |\tilde{G}| = 1 \right\}$
5            Find best single group $G^* = \underset{\tilde{G} \in \mathcal{B}}{\arg\max} \left( \widehat{LOGI}(\tilde{G}) \right)$
6            **if** $\widehat{LOGI}(G^*) > \delta$ **then**
7               $B = G^*$
8               $L_{i-1} = \widehat{LOGI}(B)$
9
10        **if** $i > 1$ *and* $B \neq \emptyset$ **then**
11           Define candidate set $\tilde{B} := \left\{ \tilde{G} \in \mathcal{P}(\mathcal{G}) \,\middle|\, |\tilde{G}| = i \text{ and } B \subset \tilde{G} \right\}$
12           Find best combination $G^* = \underset{\tilde{G} \in \mathcal{B}}{\arg\max} \left( \widehat{LOGI} \left( \bigcup_{G' \in \tilde{G}} G' \right) \right)$
13           **if** $\widehat{LOGI} \left( \bigcup_{G' \in G^*} G' \right) - L_{i-1} > \delta$ **then**
14              $B = \bigcup_{G' \in G^*} G'$
15              $L_{i-1} = \widehat{LOGI}(B)$
16           **else**
17              **break** for loop
18

---

Since the order in which feature groups are added is also known, alluvial charts (Allaire et al. 2017) can be created for visualization purposes (see Figs. 2 and 10). In these charts, we included the number of times feature groups were added as well as the performance on the test datasets. These charts show how frequently a group was selected given that another group was already included and thereby highlight robust combinations of groups.

## 5 Comparison of grouped feature importance methods

After introducing the methodological background of the different loss-based grouped feature importance measures in Sect. 3, we will now compare them in different simulation settings. We analyze the impact on all methods for settings where (1) groups are dependent, (2) correlations within groups vary, and (3) group sizes differ.

### 5.1 Dependencies between groups and sparsity

In this section, we compare refitting- and permutation-based grouped feature importance methods and show how different dependencies between groups can influence the importance scores. We demonstrate the benefits of the sequential grouped feature importance procedure and conclude with a recommendation of when to use refitting or permutation-based methods depending on the use-case.

We simulate a data matrix $\mathbf{X}$ with $n = 1000$ instances and 3 groups $G_1, G_2, G_3$, with each of them containing 10 normally distributed features. Features are simulated in such a way that features within each group are highly correlated. However, features in $G_3$ are independent of features in $G_1$ and $G_2$, while features in $G_1$ and $G_2$ are also highly correlated with each other. To generate normally distributed features with such correlation patterns, we follow the approach of Toloşi and Lengauer (2011) and use prototype vectors in the following way: (1) We draw $n$ instances of the prototype vector $\mathbf{U} \sim \mathcal{N}(0, 1)$. (2) We generate features in $G_1$ by adding a normally distributed error term $\epsilon \sim \mathcal{N}(0, 0.5)$ to 10% of the instances of the prototype vector $\mathbf{U}$. (3) Features in $G_2$ are generated by copying features of $G_1$ and adding a small normally distributed error term $\epsilon \sim \mathcal{N}(0, 0.01)$ to the copied features. It follows that features within $G_1$ and $G_2$ as well as features between the two groups are highly correlated. (4) We generate a new prototype vector $\mathbf{V}$, which is independent of $\mathbf{U}$. (5) We generate features for $G_3$ in the same way as done for $G_1$ in step (2) but with the prototype vector $\mathbf{V}$.

The target vector $\mathbf{Y}$ is generated by $\mathbf{Y} = 2\mathbf{U} + 1\mathbf{V} + \epsilon$, with $\epsilon \sim \mathcal{N}(0, 0.1)$. We fitted a support vector machine with a radial basis function kernel[5], as an example of a black-box algorithm.
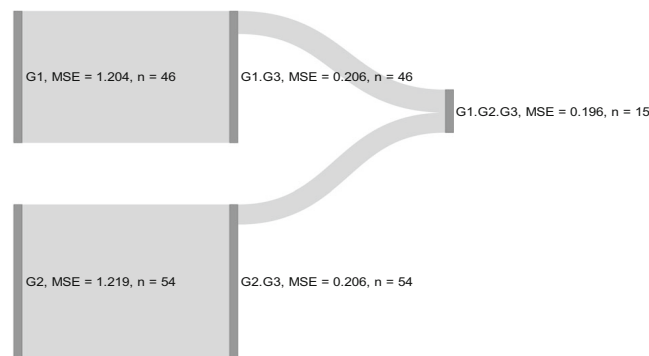
The results in Table 3 show that there can be major differences depending on how the grouped feature importance is calculated. Permutation methods (GOPFI & GPFI & GSI) reflect the importance of the groups based on a model trained on a fixed dataset. In contrast, refitting methods (LOGI & LOGO) retrain the model on a reduced dataset and can therefore learn new relationships. Looking at the results from the permutation methods, we can see that the groups $G_1$ and $G_2$ are approximately equally important while both being more important than $G_3$. However, the results from the refitting methods can reveal some interesting relationships between the groups. The refitting methods highlight that $G_1$ and $G_2$ are more or less interchangeable if we only consider a performance-based interpretation (which might not coincide with a domain-specific

---

[5] Epsilon regression, $\epsilon = 0.1$, $C = 1$ with heuristically chosen kernel width according to (Caputo et al. 2002) (here: $\sigma = 0.079$).

**Table 3** Results of different feature importance calculations of the simulation

| Group | GOPFI | GPFI | GSI | LOGI | LOGO |
|-------|-------|------|-----|------|------|
| $G_1$ | 6.04 ($\pm$ 0.37) | 2.64 ($\pm$ 0.07) | 4.12 ($\pm$ 0.45) | 3.93 ($\pm$ 0.75) | $-$0.01 ($\pm$ 0.02) |
| $G_2$ | 5.90 ($\pm$ 0.35) | 2.57 ($\pm$ 0.09) | 4.01 ($\pm$ 0.47) | 3.93 ($\pm$ 0.76) | $-$0.00 ($\pm$ 0.02) |
| $G_3$ | 1.76 ($\pm$ 0.39) | 1.75 ($\pm$ 0.05) | 1.54 ($\pm$ 0.39) | 0.58 ($\pm$ 1.01) | 1.01 ($\pm$ 0.22) |

GSI scores were calculated without approximation, with $v_{\text{perm}}$ as value function (see Eq. 14). All results were averaged by a 10-fold cross-validation scheme, with standard deviations reported in parentheses



**Fig. 2** Sequential grouped feature importance for the simulation in Sect. 5.1. 100 times repeated subsampling. Improvement threshold $\delta = 0.001$. Vertical bars show one step of the sequential procedure (left to right). Height of the vertical bars represent the number of subsampling iterations that a combination of groups was chosen. $MSE$ scores show predictive performance. Streams represent the addition of a group

perspective)[6]. Hence, the two groups do not complement each other. This is reflected by the near-zero LOGO scores, which indicate that leaving each group out of the full model does not considerably change the model's expected loss.

Figure 2 illustrates the results of the sequential procedure introduced in Algorithm 1. We see that across 100 subsampling iterations, $G_1$ was chosen 46 times as the most important first group, and $G_2$ was chosen 54 times with similar predictive performance for both groups, while $G_3$ was never chosen as the first most important group. Hence, similar to LOGI, we can see that if only one group can be chosen, it would either be $G_1$ or $G_2$ with approximately the same probability. In the second step, the group $G_3$ was added in all cases to either $G_1$ or $G_2$ (depending on which group had been chosen in the first step). This step resulted in an on-average drop in the MSE score from 1.2 to 0.2. In only a few cases (15 out of 100), the final addition of either $G_1$ or $G_2$ to a full model in step 3 exceeded the very low chosen threshold of $\delta = 0.001$. This rather unlikely improvement is represented by the proportionally narrower band that connects the second and the third step (dark gray bars) in the chart in Fig. 2. This reveals that these two groups are—from a performance or loss perspective—rather interchangeable and do not benefit from one another.

---

[6] It is possible that adding a group of features to the model might not lead to a better model performance, but the group may still be relevant due to the domain-specific context. However, this depends on the regarded use case. All our interpretations here are purely statistical.

The choice between using permutation-based or refitting-based grouped feature importance methods might depend on the number of groups and correlation strength between the different groups. If feature groups are distinct and features between the groups are almost uncorrelated, we might prefer permutation over refitting methods due to lower computation time. In cases where groups are correlated with each other (e.g., because some features belong to multiple groups), refitting methods might be preferable, as they are not misleading in correlated settings. Since the number of groups is usually smaller than the number of features in a dataset, refitting methods for groups of features could become a viable choice. Furthermore, with the sequential grouped feature importance procedure, it is possible to find sparse and well performing combinations of groups in an interpretable manner. Thus, this approach helps to better understand which groups of features were important (e.g., as they were more frequently selected) given that certain groups were already selected.

### 5.2 Varying correlations within groups

In many use cases, it is quite common to group similar (and therefore, often correlated) features together, while groups of features may be almost independent of each other. However, compared to Sect. 5.1, correlations of features within groups might differ. We created a data matrix $\mathbf{X}$ with $n = 1000$ instances and 4 groups $G_1$, $G_2$, $G_3$, and $G_4$, with each of these groups containing 10 normally distributed features. Using fivefold cross-validation, we fitted a random forest with 2000 trees and a support vector regression with a radial basis function kernel.[7] The univariate target vector $\mathbf{Y}$ is defined as follows:

$$\mathbf{Z}_j = 3\mathbf{X}_{G_j,3}^2 - 4\mathbf{X}_{G_j,5} - 6\mathbf{X}_{G_j,7} + 5\mathbf{X}_{G_j,9} \cdot d_j, \quad j \in \{1, 2, 3\}$$

$$\mathbf{Y} = \sum_{j=1}^{3} \mathbf{Z}_j + \epsilon$$

with

$$d_j = \begin{cases} 1, & \text{if mean}(\mathbf{X}_{G_j,8}) > 0 \\ 0, & \text{otherwise} \end{cases}$$

and $\epsilon \overset{iid}{\sim} N(0, 1)$. The $i$−th feature of the $j$-th group is denoted by $\mathbf{X}_{G_j,i}$. We repeated the simulation 500 times.

It follows that $G_1$, $G_2$, and $G_3$ have the same influence on the target variable, while $G_4$ has no influence on $\mathbf{Y}$. We generate the feature space $\mathbf{X}$—similar to the approach in Sect. 5.1—as follows: (1) For each feature group $j$, we generate a prototype vector $\mathbf{U}_j \sim \mathcal{N}(0, 1)$ with $n$ instances. (2) We generate the features of a group $G_j$ by altering a proportion $\alpha$ with $0 \leq \alpha \leq 1$ of the n instances of $\mathbf{U}_j$. We alter these instances by taking

---

[7] We used a cost parameter of $C = 1$ and estimate the kernel width based on the heuristic introduced by Caputo et al. (2002)

a weighted average between the respective values of $\mathbf{U}_j$ (20%) and a standard normally distributed random variable $\mathbf{W}_i$ (80%). For the results shown in Fig. 3, we set $\alpha$ to 0.1 for all features within the same group. Hence, correlations within groups are the same (around 90%) for all groups, while groups themselves are independent of each other. The plots show that all methods correctly attribute the same importance to the first three groups, while the fourth group is not important for predicting $\mathbf{Y}$. The lower plots in Fig. 3, on the other hand, correlations within groups vary across groups. The altering proportion parameter $\alpha$ is set to 0.1 for features of $G_1$ and $G_4$, to 0.3 for features of $G_2$, and to 0.6 for features of $G_3$. Hence, features in $G_1$ and $G_4$ are highly correlated within the respective group, while features within $G_2$ and $G_3$ show a medium and small correlation, respectively. While $G_4$ is still recognized to be unimportant, the relative importance of groups 1 to 3 drops with decreasing within-group correlation. This artifact seems—at least, in this simulation setting—to be even more severe for the random forest compared to the support vector machine. For example, $G_3$ is on average less than half as important as $G_1$ for permutation-based methods. Thus, none of the methods reflect the true importance of the different groups of the underlying data generating process. A possible reason for this artifact is that the regarded model learned effects different from those given by the underlying true relationship. Especially for the random forest, this has already been studied extensively in the presence of different correlation patterns in the feature space (Strobl et al. 2008; Nicodemus et al. 2010). Additionally, Hooker and Mentch (2019) showed that permutation-based methods are more sensitive in this case than refitting methods, which is also visible for both models in Fig. 3. Since the model is learned on the original feature space and group structures are not considered in the modelling process, we can also observe this effect when applying grouped feature importance methods. This is due to the fact that we can only quantify which groups are important for the model or algorithm performance but not for the underlying data generating process, which is usually unknown. Another approach to quantify feature importance when using random forests is to extract the information on how often a feature has been used as a splitting variable for the different trees. The feature chosen for the first split has the most influence within each tree. Hence, we calculated for each repetition the percentage of how often a feature is chosen as the first splitting feature. The distribution over all repetitions is displayed in Fig. 4. Each of the features of $G_1$ is on average chosen more often as the first splitting feature than all features of the other groups, no matter if it has an influence on the target or not. The influential features of $G_3$ (which has the lowest within-group correlation) are rarely chosen as the first splitting feature. This observation confirms the results of the grouped importance methods in Fig. 3, since all of them rank $G_3$ as least important from the influential feature groups.

Note that while GPFI and LOGO are calculated with reference to the full model's performance—which on average leads to higher absolute values than the two counter-methods based on the null model's performance—GOPFI and LOGI might lead to less robust results, as the newly learned effects as well as the approximation of the permutation effect underlie a higher uncertainty. This effect might increase when relative values instead of absolute values are considered due to smaller absolute importance scores of GOPFI and LOGI. However, the methods are only comparable on a relative scale. This effect is also visible in the boxplots of Fig. 3. Furthermore, LOGI can also
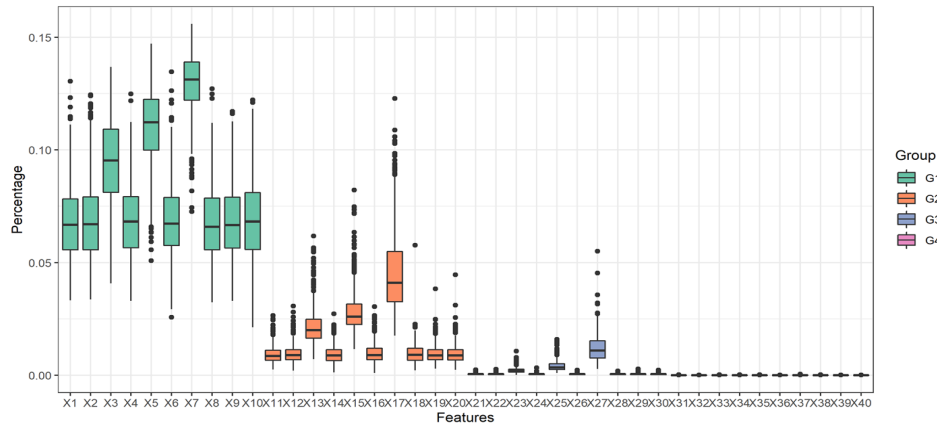
**Fig. 3** Upper (lower) plots: Grouped relative importance scores in the case of *equally sized (varying sizes of)* within-group correlations for random forest (left) and SVM (right). Relative importance is calculated by dividing each of the absolute group importance scores by the importance score of $G_2$. Hence, the relative importance of $G_1$ is 1. Boxplots illustrate the variation between different repetitions

take negative values in the case of $G_4$, as the feature group does not affect the target in the underlying data generating process, and hence it might be counterproductive to only include $G_4$ compared to the null model.

### 5.3 Varying sizes of groups

Another factor to consider when calculating grouped rather than individual feature importance scores is that differing group sizes might influence the ranking of the scores. Groups with more features might often have higher grouped importance scores and

🖉 Springer

**Fig. 4** Percentage of how often each feature is chosen as the first splitting feature within the trained random forests. Results have been averaged over the cross-validation folds for each repetition. Boxplots show the distribution over all 500 repetitions

might contain more noise features than smaller groups. Therefore, Gregorutti et al. (2015) argue that in case one must decide between two groups that have an equal importance score, one would prefer the group with fewer features. Following from that, they normalize the grouped feature importance scores regarding the group size with the factor $|G|^{-1}$. This is also used in the default definition of the grouped model reliance score in Valentin et al. (2020). However, the usefulness of normalization highly depends on the question the user would like to answer. This is illustrated in a simulation example in Fig. 5. We created a data matrix $\mathbf{X}$ with $n = 2000$ instances and 2 groups, with $G_1$ containing $\{x_1, \ldots x_6\}$ and $G_2$ containing $\{x_7, x_8\}$ i.i.d. uniformly distributed features on the interval [0, 1]. The univariate target variable $\mathbf{Y}$ is defined as follows:

$$\mathbf{Y} = 2\mathbf{X}_1 + 2\mathbf{X}_3 + 2\mathbf{X}_7 + \epsilon, \quad \text{with} \quad \epsilon \overset{iid}{\sim} N(0, 1).$$

We used 1000 observations for fitting a random forest with 2000 trees and 1000 observations for prediction and calculating the GSI as defined in Sect. 3.3 with a permutation-based value function. This was repeated 500 times. Figure 5 shows that $G_1$ is about twice as important as $G_2$. As shown in Sect. 3.3 and Appendix B, we can compare the GSI with the Shapley importance on feature level. In case there are no higher-order interaction terms between groups modeled by the random forest, the single feature importance scores will approximately sum up to the grouped importance score, as shown in this example. This provides a more detailed view of how many and which features are important within each group. In this case, there are two equally important features in $G_1$ and one equally important feature in $G_2$. If we use the normalization constant in this example, we would divide the grouped importance score of $G_1$ (which is on average approximately 1.1) by 6 and the one of $G_2$ (which is on average approximately 0.55) by 2. Consequently, $G_2$ with a normalized score of

**Fig. 5** Shapley importance on group (left) and on feature level (right). Boxplots show the variation between the 500 repetitions of the experiment

approximately 0.27 would be regarded as more important than $G_1$ with a normalized score of approximately 0.18. It follows that if we must decide between two groups, we would choose $G_2$ when we follow the approach of Gregorutti et al. (2015). However, since $G_1$ contains two features with the same importance as the one important feature of $G_2$, and hence $G_1$ contains more information from a statistical perspective, the user might prefer $G_1$. Furthermore, breaking down the GSI to the single-feature Shapley importance scores puts the user in the position of defining sparser groups by excluding non-influential features.

Finally, Table 4 presents a summary of the key takeaways regarding all discussed grouped feature importance methods.

## 6 Feature effects for groups

Feature effect methods quantify or visualize the influence of features on the model's prediction. For a linear regression model, we can easily summarize the feature effect in one number, thus making interpretation very simple: If we change feature $x_1$ by one unit, our prediction will change by the corresponding coefficient estimate $\hat{\beta}_1$ (positively or negatively depending on the sign of the coefficient). For more complex non-linear models like generalized additive models, such a simplified summary of the feature effect is not adequate, as the magnitude and sign of the effect might change over the feature's value range. Hence, it is more common to visualize the marginal effect of the feature of interest on the predicted outcome. Since ML models are often complex non-linear models, different visualization techniques for the feature effect have been introduced in recent years. Common methods are PDP, ICE curves or ALE (Friedman 2001; Goldstein et al. 2013; Apley and Zhu 2019), which show how changes in the feature values affect the predictions of the model. However, these are

**Table 4** Overview of pros and cons of the grouped feature importance methods

| Criteria | GOPFI & GPFI | GSI | LOGI & LOGO |
|---|---|---|---|
| Time efficient | Yes (in comparison to alternatives) | Depends on number of groups | Depends on number of groups |
| Dependencies between groups (Sect. 5.1) | No full picture | No full picture | More insights than permutation-based if regarded together |
| Identify well performing combinations of groups (Sect. 5.1) | Not in general | Not in general | Only LOGI wihin Algorithm 1 |
| Correlations within groups but independence between groups (Sect. 5.2) | Depends on learned effects of the model, less problematic if within group correlations do not differ strongly between groups | Depends on learned effects of the model, less problematic if within group correlations do not differ strongly between groups | More robust than permutation-based methods but still dependent on learned effects |
| Drilldown of grouped importance score on feature level (Sect. 5.3) | No | Yes (approximately depending on the influence of higher-order interactions) | No |

While GOPFI is less relevant on its own, LOGI can provide insightful interpretations, e.g., if feature groups are correlated with each other or when used within the sequential procedure introduced in Sect. 4. The sequential procedure is the only method that can identify well performing and sparse combination of groups. Note that GSI is only evaluated w.r.t. a permutation-based calculation

usually only defined for a maximum of two features. For larger groups of features, this becomes more challenging, since it is difficult to visualize the influence of several features simultaneously. The approach described in this section aims to create effect plots for a predefined group of features that have an interpretation similar to that of the single-feature PDP. To achieve this, we transform the high-dimensional space of the feature group into a low-dimensional space by using a supervised dimension reduction method, which is discussed in Sect. 6.1. We want to find a few underlying factors that are attributed to a sparse and interpretable combination of features that explain the effect of the regarded group on the model's expected loss. We provide a detailed description of this method in Sect. 6.3 and introduce the resulting combined features effect plot (CFEP). In Sect. 6.4, we illustrate the advantages of applying a supervised rather than an unsupervised dimension reduction method and compare our method to the main competitor, which is the totalvis effect plot introduced in Seedorff and Brown (2021).

## 6.1 Choice of dimension reduction method

The most prominent dimension reduction technique is arguably PCA (Jolliffe 1986). PCA is restricted to explaining most of the variance of the feature space, and the identified projections are not related to the target variable (for more details see Appendix

C.1). Because we want to visualize the mean prediction of combined features as a result of the dimension reduction process, we prefer supervised procedures that maximize dependencies between the projected data $\mathbf{XV}$—with $\mathbf{V}$ being a projection $\mathbf{V} \in \mathbb{R}^{p \times p}$—and the target vector $\mathbf{Y}$ (as we show in Sect. 6.4). Many methods for supervised PCA have been established. For example, see Bair et al. (2006), who used a subset of features that were selected based on their linear correlation with the target variable. Another very popular method that maximizes the covariance between features and the target variable is partial least squares (PLS) (Wold et al. 1984). The main difference between these methods and the supervised PCA (SPCA) introduced by Barshan et al. (2011) is that the SPCA is based on a more general measure of dependence, called the Hilbert-Schmidt Independence Criterion (HSIC). This independence measure is constructed to be zero, if and only if any bounded continuous function between the feature and target space is uncorrelated. In practice, an empirical version of the HSIC criterion is calculated with kernel matrices. It follows that while this SPCA technique can cover a variety of linear and non-linear dependencies between $\mathbf{X}$ and $\mathbf{Y}$ by choosing an appropriate kernel, the other suggested methods are only able to model linear dependencies between the features and the target variable. The approach that is probably best suited for our application of finding *interpretable* sets of features in a high-dimensional dataset is the method called sparse SPCA, described in Sharifzadeh et al. (2017). Similar to the SPCA method from Barshan et al. (2011), sparse SPCA not only uses the HSIC criterion to maximize the dependency between projected data $\mathbf{XV}$ and the target $\mathbf{Y}$, but also incorporates an $L_1$ penalty of the projection $\mathbf{V}$ for sparsity. The sparse SPCA problem can be solved with a *penalized matrix decomposition* (Witten et al. 2009). More theoretical details on the sparse SPCA, including the HSIC criterion and how it can be calculated empirically, and the choice of kernels and hyperparameters can be found in Appendix C.

### 6.2 Totalvis effect plot

Seedorff and Brown (2021) recently introduced a method that aims to plot the combined effect of multiple features by using PCA. Their approach can be described as follows: First, they apply PCA on the regarded feature space to receive the principal components matrix after rotation. For the principal component of interest, they create an equidistant grid. Second, for each grid value, they replace all values of the selected principal component with this grid value and transform the matrix back to the original feature space. Third, The ML model is applied on these feature values and a mean prediction for the grid point of the regarded principal component is calculated. Steps 2 and 3 are repeated for all grid points.

Hence, with this method, combined effect plots for up to $p$ principal components can be created. Thus, Seedorff and Brown (2021) do not focus on explaining groups of features explicitly. Furthermore, they use PCA for unsupervised dimension reduction, and thus, projections might not be related to the target. Due to using PCA and not sparse PCA, the results might be difficult to interpret, as many or all features might have an influence on the principal component. Lastly, with the back-transformation from the principal component matrix to the original feature space, all feature values

change and might not be meaningful anymore. For example, in the case of integer features, the back-transformation might lead to real feature values. We illustrate the drawbacks of the method compared to the CFEP in Sect. 6.4.

### 6.3 Combined features effect plot (CFEP)

The CFEP picks up the idea of PDPs (Friedman 2001) and extends it to groups of features. The partial dependence function is defined as

$$f_S^{PD}(\mathbf{x}_S) = \mathbb{E}_{X_C}[\hat{f}(\mathbf{x}_S, X_C)] \tag{17}$$

with $S \subset \{1, \ldots, p\}$ and $C = \{1, \ldots, p\} \backslash S$. Since the joint distribution of $X_C$ is usually unknown, the Monte Carlo method is used to estimate $f_S^{PD}(\mathbf{x}_S)$:

$$\hat{f}_S^{PD}(\mathbf{x}_S) = \frac{1}{n} \sum_{i=1}^{n} \hat{f}(\mathbf{x}_S, \mathbf{x}_C^{(i)}) \tag{18}$$

Hence, we marginalize over all features in $C$ and with that we obtain the average marginal effect for the feature subset in $S$. The PDP usually visualizes this average marginal effect for $|S| \leq 2$ by plotting $\left(\mathbf{x}_S^{(k)}, \hat{f}_S^{PD}(\mathbf{x}_S^{(k)})\right)$ for some pre-specified grid points $k = \{1, \ldots, m\}$.[8] However, this is usually only possible for $|S| \leq 2$ and thus not directly applicable to visualize the combined effect of feature groups. To obtain a visualization in the case of $|S| > 2$, we need to reduce the dimensions and therefore define the CFEP of a certain group of features $G$ as follows:

(1) We first apply a suitable (preferably supervised) dimension reduction method (e.g., here we use the sparse SPCA, however, the CFEP follows a modular approach and hence the dimension reduction method is exchangeable) on features in $G \subset \{1, \ldots, p\}$ to obtain a low dimensional representation of the feature group $G$. We denote these principle component functions—which are ordered according to relevance[9] and which possibly depend on a reduced set of features[10] $S_j \subseteq G$ with $j \in \{1, \ldots, |G|\}$—by $g_j : \mathcal{X}_{S_j} \longrightarrow \mathbb{R}$.

(2) For visualization purposes, we choose from all possible $g_j$ with $j \in \{1, \ldots, |G|\}$ a principle component function

$$g : \mathcal{X}_S \longrightarrow \mathbb{R} \tag{19}$$

(with $S$ being its reduced set of features) which serves as a proxy for the feature group $G$. We usually only consider the first few principle components.

---

[8] For example, by using an equidistant grid or a random sample of values of $\mathbf{x}_S$.

[9] The relevance is defined by the objective that is optimized by the dimension reduction method. For sparse SPCA this is the HSIC criterion (see also Appendix C) and for PCA it is the explained variance.

[10] If a dimension reduction method which results in a sparse solution (e.g., sparse SPCA) is applied, then $S_j$ is only a subset of $G$ and might differ for different principal components.

(3) We calculate the average marginal effect $\hat{f}_S^{PD}(\mathbf{x}_S)$ of the feature set S exactly as in Eq. (18).

(4) We visualize the CFEP by plotting $\left(g(\mathbf{x}_S^{(i)}), \hat{f}_S^{PD}(\mathbf{x}_S^{(i)})\right)$ for each observation in the dataset.
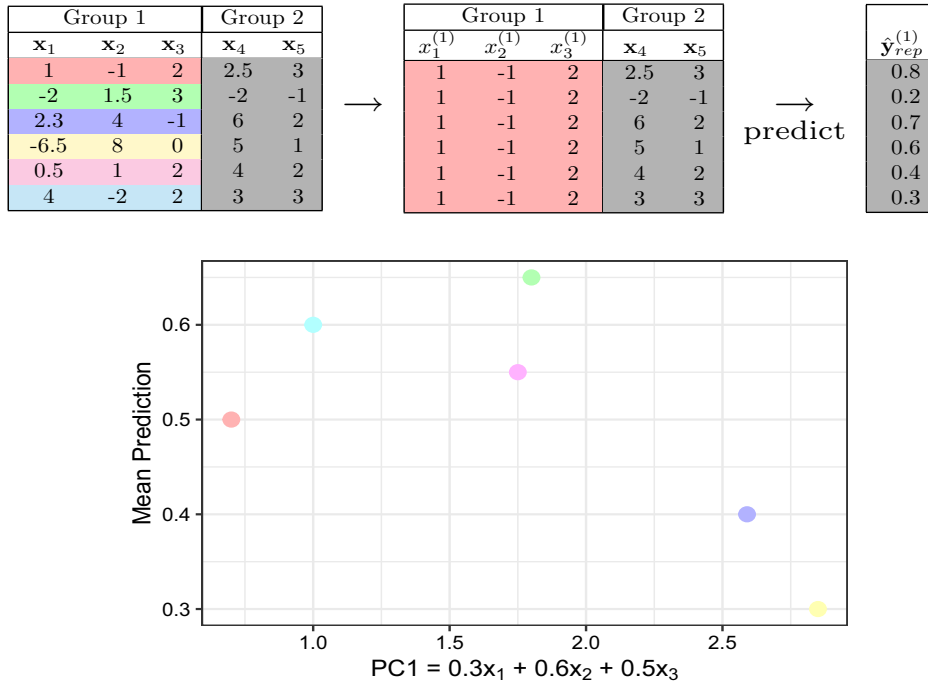
Hence, the CFEP visualizes the average marginal effect of features in $S$ against the combinations of features received by the dimension reduction method (e.g., a linear combination of a principal component in the case of sparse SPCA) and thus shows how different values of $g(\mathbf{x}_S)$ affect the predictions of a given model. For a feature group, several principle components $g_j$ and hence several CFEPs may be of interest.

The CFEP is defined in Algorithm 2, but we will demonstrate the procedure of constructing a CFEP with the illustrative example in Fig. 6. In this example, we have two predefined groups of features, where the first group contains $x_1$, $x_2$, and $x_3$, and the second group contains features $x_4$ and $x_5$. The sparse SPCA on the first group yields a first principal component ($g_1$) with the loadings 0.3 for $x_1$, 0.6 for $x_2$ and 0.5 for $x_3$ (step 1 to 3 of Algorithm 2). It follows that $S = \{1, 2, 3\}$ and that the low dimensional representation of interest is $g_1$. For the construction of a CFEP for $g_1$, mean predictions for the principal component are calculated for each observation. To calculate the mean prediction of the first observation (shown in red), we replace the values of features with non-zero loadings of $g_1$ of each instance in the dataset by the feature values of the first observation (step 6 in Algorithm 2). A prediction vector $\hat{\mathbf{y}}_{rep}^{(1)}$ is then calculated with the previously trained model (step 7 in Algorithm 2). The value on the y-axis for the red point in the graph below corresponds to the mean over all predictions for the first observation: $\bar{\hat{y}}_{rep}^{(1)} = (0.8 + 0.2 + 0.7 + 0.6 + 0.4 + 0.3)/6 = 0.5$. The value on the x-axis is the linear projection of the first observation for the regarded principal component (step 8 and 9 in Algorithm 2). Hence, it is calculated by the weighted sum of feature values $x_1^{(i)}$ to $x_3^{(i)}$, where the weights are defined by the loadings of the respective principal component that we receive with sparse SPCA.

In contrast to PDP or totalvis effect plots, CFEP produces a point cloud instead of a curve. The CFEP is, mathematically speaking, not a function, since points on the x-axis correspond to linear projections of features within a group. A point $z$ on the x-axis can have multiple combinations of features, which lead to $z$ and have different mean predictions on the y-axis. However, we now have the possibility to interpret the shape of the point cloud and can draw conclusions about the behavior of the mean prediction of the model regarding a linear combination of features of interest.

### 6.4 Experiments on supervised versus unsupervised dimension reduction

As discussed in Sect. 6.1, PCA might be the most popular dimension reduction method. However, since PCA is unsupervised, it does not account for the dependency between the feature space and the target variable. To evaluate the degree to which this drawback influences CFEP, we examine two regression problems on simulated data. The first is defined by a single underlying factor depending on a sparse set of features, which can be represented by a single principal component. The linear combination of this feature set is also linearly correlated with the target variable. The second regression problem

| Group 1 | | | Group 2 | |
|---|---|---|---|---|
| $\mathbf{x}_1$ | $\mathbf{x}_2$ | $\mathbf{x}_3$ | $\mathbf{x}_4$ | $\mathbf{x}_5$ |
| 1 | -1 | 2 | 2.5 | 3 |
| -2 | 1.5 | 3 | -2 | -1 |
| 2.3 | 4 | -1 | 6 | 2 |
| -6.5 | 8 | 0 | 5 | 1 |
| 0.5 | 1 | 2 | 4 | 2 |
| 4 | -2 | 2 | 3 | 3 |

$\rightarrow$

| Group 1 | | | Group 2 | |
|---|---|---|---|---|
| $x_1^{(1)}$ | $x_2^{(1)}$ | $x_3^{(1)}$ | $\mathbf{x}_4$ | $\mathbf{x}_5$ |
| 1 | -1 | 2 | 2.5 | 3 |
| 1 | -1 | 2 | -2 | -1 |
| 1 | -1 | 2 | 6 | 2 |
| 1 | -1 | 2 | 5 | 1 |
| 1 | -1 | 2 | 4 | 2 |
| 1 | -1 | 2 | 3 | 3 |

$\xrightarrow{\text{predict}}$

| $\hat{\mathbf{y}}_{rep}^{(1)}$ |
|---|
| 0.8 |
| 0.2 |
| 0.7 |
| 0.6 |
| 0.4 |
| 0.3 |



**Fig. 6** Explanation of estimating and visualizing CFEP; the x-coordinate reflects the linear combination of features with non-zero loadings for $g_1$, and the y-coordinate reflects the mean predictions $\bar{\hat{y}}_{rep}^{(i)}$ for each observation $i$. The substitution of values for each observation is only done for features with non-zero loadings

contains two underlying factors that depend on two sparse sets of features. While the linear combination of the first feature set is also linearly correlated with the target, the second factor has a quadratic effect on $\mathbf{Y}$. In both cases, we compare the usage of sparse supervised and unsupervised PCA (sparse SPCA and sparse PCA) as dimension reduction methods within CFEP and compare them to the totalvis effect plot. Here, we investigate if the respective dimension reduction method does correctly identify the sparse set of features for each group. Additionally, we determine how accurately we can predict the true underlying relationship between the linear combination of these features and the target variable. Since we simulated the data, we know the number of underlying factors (principal components).

### 6.4.1 One factor

In this example, we created a data matrix $\mathbf{X}$ with 500 instances of 50 standard normally distributed features with decreasing correlations. Therefore, all features are generated as done in Sect. 5.2. The altering proportion $\alpha$ is set to 0.2 for the first 10 features, to 0.4 for the next 10 features, and to 1 for the last 10 features. Thus, while the first 10 features are highly correlated with each other, the last 10 features are approximately

---

**Algorithm 2:** Combined Features Effect Plot

---

   **input**  : Dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^{n}$,
              group $G \subset \{1, ..., p\}$,
              model $\hat{f}$ trained on $\mathcal{D}$.
   **output**: Combined Features Effect Plot

**1** Perform sparse SPCA on $\mathring{\mathcal{D}} := \{(\mathbf{x}_G^{(i)}, y^{(i)})\}_{i=1}^{n}$;
**2** Choose a principle component function of interest $g$;
**3** Let $S \subseteq G$ be the sparse set of features of $g$;
**4** **for** $i \in \{1, ..., n\}$ **do**
**5**     get feature values $\mathbf{x}_S^{(i)}$;
**6**     create $\mathcal{D}_{rep}^{(i)}$ by replacing feature values from $S$ of every observation with $\mathbf{x}_S^{(i)}$;
**7**     predict vector $\hat{\mathbf{y}}_{rep}^{(i)}$ by applying $\hat{f}$ on $\mathcal{D}_{rep}^{(i)}$ row-wise;
**8**     calculate the mean prediction $\bar{\hat{y}}_{rep}^{(i)}$ of $\hat{\mathbf{y}}_{rep}^{(i)}$;
**9**     save $g(\mathbf{x}_S^{(i)})$ as x-coordinate and $\bar{\hat{y}}_{rep}^{(i)}$ as y-coordinate of observation $i$ for the CFEP (see Eq. (19));

---

The CFEP can be used as a descriptive method to better understand the effect of a group of features on the target variable. The dimension reduction method in step 1 is exchangeable.

uncorrelated with each other. The sparse subgroup defined by the variable $\mathbf{Z}$ is a linear combination of 5 features from $\mathbf{X}$ and has itself a linear effect on the target variable $\mathbf{Y}$:

$$\mathbf{Z} = \mathbf{X}_5 - 2\mathbf{X}_8 - 4\mathbf{X}_{25} + 8\mathbf{X}_{47} + 4\mathbf{X}_{49}$$

$$\mathbf{Y} = \mathbf{Z} + \epsilon, \quad \text{with} \quad \epsilon \overset{iid}{\sim} N(0, 1).$$

Hence, according to our notation, $G_{\mathbf{Z}}$ is defined by $G_{\mathbf{Z}} = \{5, 8, 25, 47, 49\}$, and thus, $X_{G_Z}$ is the related subset of all features. We drew 100 samples and fitted a random forest with 2000 trees with each sample drawing. We used the 10-fold cross-validated results to perform sparse SPCA. For each dimension reduction method, we estimate $\hat{\mathbf{Z}}$ by summing up the (sparse) loading vector (estimated by the dimension reduction method) multiplied by the feature matrix $\mathbf{X}$. Therefore, $\mathbf{X}_{G_{\hat{Z}}}$ is defined by the received sparse feature set. The mean prediction $\bar{\hat{\mathbf{Y}}}_{rep}$ for the CFEP is calculated as described in Sect. 6.3.

    The impact of choosing a supervised over an unsupervised sparse PCA approach is shown in Fig. 7, which also shows the average linear trend and 95% confidence bands of CFEP for the simulation results. To evaluate how well the estimated mean prediction $\hat{\mathbf{Y}}_{rep}$ approximates the underlying trend, we assume that we know that $\mathbf{Z}$ has a linear influence on the target. Thus, we fit a linear model on each simulation result. To compare the received regression lines, we evaluate each of them on a predefined grid and average over all 100 samples (represented by the red line). The confidence bands are then calculated by taking the standard deviation over all estimated regression lines on grid level and calculating the 2.5% and 97.5% quantiles using the standard normal approximation. The associated calculation steps for each of the 100 samples can be summarized as follows:

**Fig. 7** Average linear trend and confidence bands of CFEP over all samples using sparse SPCA (left) and sparse PCA (middle) compared to estimated totalvis effect curves over all 100 samples for the first principal component (black) and the average linear trend (red) (right) (Color figure online)

(1) Estimate a linear model $\hat{f}(\mathbf{X}_{G_{\hat{\mathbf{Z}}}}) \sim \mathbf{Z}$.
(2) Define an equidistant grid of length 50 within the range of $\mathbf{Z}$.
(3) Apply the linear model estimated in 1) on the grid defined in 2).
(4) Repeat steps 1 to 3 for $\hat{f}(\mathbf{X}_{G_{\mathbf{Z}}})$ by using the true underlying features of $\mathbf{Z}$ to calculate the combined features dependencies that we call the ground truth.

The left plot in Fig. 7 shows a similar linear trend of the estimated CFEP compared to the average ground truth (represented by the blue line), while the red line in the right plot varies around 0. By using sparse SPCA, the underlying feature set $\mathbf{X}_{G_{\hat{\mathbf{Z}}}}$ is better approximated than with sparse PCA, which is reflected in the MSE between $\mathbf{Z}$ and $\hat{\mathbf{Z}}$ of 0.7 for sparse SPCA and 1.9 for sparse PCA. Figure 8 provides an explanation for those differences. While sparse SPCA (on average) more strongly weights features that have a large influence on the target, impactful loading weights for sparse PCA are solely distributed over highly correlated features in $\mathbf{X}$ that explain the most variance in the feature space. Thus, including the relationship between the target and $\mathbf{X}$ in the dimension reduction method may have a huge influence on correctly approximating the underlying factor and, hence, also on the CFEP.

Similar to using sparse PCA as a dimension reduction method within CFEP, on average, the totalvis effect curves based on PCA do not show a clear positive linear trend (see Fig. 7). For almost half of the samples, we even receive a negative instead of a positive trend for the underlying factor. The interpretation is opposite to the actual effect and, hence, is misleading.

### 6.4.2 Two factors

In real-world data settings are often more complex by containing non-linear relationships and the target variable is described by more than one underlying factor. Hence,

**Fig. 8** Distribution of feature loadings in sparse SPCA (top) and sparse PCA (bottom) over all samples. Rhombuses denote the mean values, with the blue rhombuses indicating the features that have an influence on the target in the underlying model formula (Color figure online)

we now examine a more complex simulation setting to assess if we can observe the same behavior that we observed for the simple case. To that end, we simulated a data matrix $\mathbf{X}$ with 500 instances for two feature sets, each containing 20 standard normally distributed features. The data for each feature set is generated as described in Sect. 5.2 but with an altering proportion of 0.15 and 0.35 for the features in the first set and 0.55 and 0.85 in the second set. Hence, within each set, the first ten features show a higher correlation among each other than the last ten features. Additionally, all features of the first set are on average more highly correlated than all features of the second set. Features between the two sets are uncorrelated. The first factor $\mathbf{Z}_1$ is a linear combination of four features from the first set and $\mathbf{Z}_2$ of two features from the second set. $\mathbf{Z}_1$ has a linear and $\mathbf{Z}_2$ a quadratic effect on $\mathbf{Y}$.

$$\mathbf{Z}_1 = 3\mathbf{X}_3 - 2\mathbf{X}_8 - 4\mathbf{X}_{13} + 8\mathbf{X}_{18}$$
$$\mathbf{Z}_2 = 2\mathbf{X}_{25} + 4\mathbf{X}_{35}$$
$$\mathbf{Y} = \mathbf{Z}_1 + \mathbf{Z}_2^2 + \epsilon, \quad \text{with} \quad \epsilon \overset{iid}{\sim} N(0, 1).$$

Again, we drew 100 samples and fitted a random forest with 2000 trees with each sample drawing. The approach is almost the same as described for one factor, with the difference being that we use the first two principal components (as we want to find two sparse feature sets instead of one).

In Fig. 9, the average linear and quadratic trend of the underlying CFEPs of $\mathbf{Z}_1$ and $\mathbf{Z}_2$ are depicted for both dimension reduction methods. While the average linear regression line of sparse SPCA matches the average ground truth almost perfectly for $\mathbf{Z}_1$, the associated line of sparse PCA shows only a slightly positive trend and differs substantially from the ground truth. Regarding $\mathbf{Z}_2$, a similar propensity can be observed for the quadratic shape. Again, this behavior results from sparse SPCA (on

**Fig. 9** Top ($\mathbf{Z}_1$): Average linear trend and confidence bands of CFEP over all samples using sparse SPCA (left) and sparse PCA (middle) compared to estimated totalvis effect curves over all 100 samples for first principal component (black) and the average linear trend (red) (right). Bottom ($\mathbf{Z}_2$): Same structure as for $\mathbf{Z}_1$, but showing the quadratic trend of $\mathbf{Z}_2$ (Color figure online)

average) more strongly weighting features that have a large effect on the target, while the unsupervised version focuses on features that explain the most variance in $\mathbf{X}$.

The estimated linear trend of the totalvis effect curves for the first principal component is negative instead of positive. Thus, for most of the samples and on average, these results are completely misleading (see Fig. 9). The quadratic shape of the second component is (on average and for almost all samples) steeper than the average ground truth. Additionally, the deviation is higher here than for CFEP with sparse SPCA.

## 7 Real data example: smartphone sensor data

Smartphones and other consumer electronics have increasingly been used to collect data for research (Miller 2012; Raento et al. 2009). The emerging popularity of these devices for data collection is grounded in their connectivity, the number of built-in sensors, and their widespread use. Moreover, smartphones enable users to perform a wide variety of activities (e.g., communication, shopping, dating, banking, navigation, listening to music) and thus provide an ideal means to study human behavior in naturalistic contexts, over extended periods of time, and at fine granularity (Harari et al. 2015, 2016, 2017). In this regard, smartphone data has been used to investigate

<span style="float:right">🍎 Springer</span>

individual differences in personality traits (Stachl et al. 2017; Harari et al. 2019), in human emotion and well-being (Servia-Rodríguez et al. 2017; Rachuri et al. 2010; Saeb et al. 2016; Thomée 2018; Onnela and Rauch 2016; Kolenik and Gams 2021), and in daytime and nighttime activity patterns (Schoedel et al. 2020).

We use a dataset on human behavior, collected with smartphones, to illustrate methods for group-based feature importance. The PhoneStudy dataset was consolidated from three separate datasets (Stachl et al. 2017; Schuwerk et al. 2019; Schoedel et al. 2018). It consists of 1821 features on smartphone-sensed behavior and 35 target variables on self-reported Big Five personality trait dimensions (domains) and subdimensions (facets). The dataset has been published online and is openly available.[11] The Big Five personality trait taxonomy is the most widely used conceptualization of stable individual differences in human patterns of thoughts, feelings, and behavior (Goldberg 1990). In their original study, Stachl et al. (2020a) used the behavioral variables to predict self-reported Big Five personality trait scores (five dimensions and 30 subdimensions) and used grouped feature importance measures to explore which classes of behaviors were most predictive for each personality trait dimension. The groups in this study were created based on theoretical considerations from past work.

The personality prediction task is challenging because (1) the dataset contains many variables on similar behaviors, (2) these variables are often correlated, and (3) effects with the targets are interactive, very small, and partially non-linear. Many variables in the dataset can be manually grouped into classes of behavior (e.g., communication and social activity, app-usage, music consumption, overall phone activity, mobility).

We use this dataset to illustrate the idea of grouped feature importance with regard to the prediction of personality trait scores for the dimension of conscientiousness (Table 5). Conscientiousness is a personality trait dimension that globally describes people's propensity to be reliable, dutiful, orderly, ambitious, and cautious (Jackson et al. 2010). We chose this personality trait because it has high practical relevance due to its ability to predict important life outcomes and behaviors (Ozer and Benet-Martínez 2006). Here, we (1) fit a random forest model to predict the personality dimension of conscientiousness, (2) compute the introduced methods for grouped feature importance (GOPFI, GPFI, GSI, LOGI, LOGO), (3) use the proposed sequential grouped feature importance procedure to investigate which groups are most important in combination, and (4) visualize the effect of different groups with CFEPs. Once the importance of individual groups has been quantified, CFEPs can be helpful to further explore the variables in these groups with regard to the criterion variable of interest (i.e., conscientiousness) to generate new hypotheses for future research.

In Fig. 10, we show a sequential procedure for our personality prediction example. The figure shows that the groups *overall phone usage* and *app usage* lead to the best model performance if used alone and, in many cases, lead to even better performances if combined. The results also suggests that if only one group can be selected, the initial selection of the feature group app usage more often leads to the smallest expected loss (mean MSE = 0.519). For a practical application, this would indicate that if only one type of feature may be collected from smartphones to predict the personality trait conscientiousness, features on app usage should be used. If two groups of data can

---

[11] https://osf.io/kqjhr/.

**Table 5** Grouped feature importance values for predicting the personality trait conscientiousness based on MSE

| Group | GOPFI | GPFI | GSI | LOGI | LOGO |
|---|---|---|---|---|---|
| Mobility (Mo) | − 0.002 (± 0.011) | − 0.002 (± 0.001) | 0.000 (± 0.003) | − 0.011 (± 0.075) | 0.000 (± 0.006) |
| Music (Mu) | − 0.001 (± 0.011) | 0.002 (± 0.002) | 0.001 (± 0.006) | − 0.019 (± 0.074) | 0.001 (± 0.012) |
| Communication and social (C) | 0.000 (± 0.008) | 0.001 (± 0.003) | 0.004 (± 0.006) | 0.008 (± 0.070) | 0.001 (± 0.010) |
| Overall phone usage (O) | 0.007 (± 0.011) | 0.009 (± 0.003) | 0.012 (± 0.008) | 0.032 (± 0.080) | 0.009 (± 0.014) |
| App usage (A) | 0.032 (± 0.009) | 0.028 (± 0.005) | 0.031 (± 0.012) | 0.041 (± 0.069) | 0.011 (± 0.019) |

All values were calculated using a resampling method (10-times cross-validation)

be collected, overall phone usage should also be added (mean MSE = 0.513). Finally, the plot indicates that in some cases (n = 9), the additional consideration of music listening behaviors in the model could lead to additional, small improvements of the expected loss (mean MSE= 0.508). If a feature group is not added, this means that it did not make a significant contribution in this iteration of the data split. Interestingly, the feature group *music* alone shows very low (or even negative) grouped feature importance scores. This would mean that music features are only predictive in the presence of other features.

To additionally explore meaningful and predictive directions in the feature space of the app usage group, we use CFEPs for the visualization. Subplot (a) in Fig. 11 shows that combinations of higher values in features on weather app usage on average lead to higher mean values in the personality trait conscientiousness. The increased frequency in weather app usage could signify the propensity of conscientious people to be prepared for future eventualities (e.g., bad weather; Jackson et al. 2010). Subplot (b) shows an interesting non-monotonic relationship between the number of different apps used each day and the mean value in conscientiousness. Subplot (c) shows that the combinations of higher values in overall phone activities lead to lower mean values in conscientiousness. Finally, plot (d) shows a similar, negative effect pattern with regard to music listening behaviors.

## 8 Conclusion

We introduced various techniques to analyze the importance and effect of user-defined feature groups on predictions of ML models. We provided formal definitions and dis-

**Fig. 10** Sequential grouped feature importance procedure for smartphone sensor data predicting *conscientiousness*. 100 times repeated subsampling. Inner resampling strategy: 10-fold cross-validation. Improvement threshold $\delta = 0.01$. Abbreviations: app-usage (A), communication & social (C), music (Mu), overall phone activity (O), mobility (Mo). Vertical bars show one step in the greedy forward search algorithm. Height of the vertical bars represent the number of subsampling iterations in which a combination of groups was chosen (for example, out of 100 subsampling iterations the group app-usage (A) was chosen 82 times as the best first group. Streams indicate the proportion of iterations that additionally benefited from a consequent step. Only streams containing at least 5 iterations and better mean performance at the end are displayed

tinction criteria for grouped feature importance methods and distinguished between permutation- and refitting-based methods. For both approaches, we defined two calculation strategies that either start with a null model or with the full model. Based on these two definitions, we introduced Shapley importance scores for groups, which we defined for permutation as well as refitting methods. Moreover, we introduced as our first main contribution a sequential grouped feature importance procedure to find good and stable combinations of feature groups. To contrast the newly proposed methods with existing ones, we compared them for different scenarios. The key recommendations for the user can be summarized for four scenarios: (1) If high correlations between groups are present, refitting methods should be preferred over permutation methods, since they often deliver more meaningful results in these scenarios. Moreover, if the number of groups is reasonably small, refitting methods become computationally feasible. (2) If a sparse set of feature groups is of interest (e.g., due to data availability), the introduced sequential procedure can be useful. It provides insights regarding the most important groups: which sparse group combinations are stable in the sense that they are frequently selected and achieve a good performance. These criteria can be critically informative in situations where feature groups must be obtained from different data sources that are associated with further costs. (3) If the correlation strengths of features within groups are very diverse, all of the introduced methods might fail to reflect the true underlying importance of the feature groups. The size of this effect

**(a)** $g_1$ *app usage*



**(b)** $g_2$ *app usage*



**(c)** $g_1$ *overall phone usage*



**(d)** $g_1$ *music*

**Fig. 11** CFEPs for the prediction of the personality trait conscientiousness. $g_1$ describes the first principal component of the respective group, and $g_2$ describes the second. More details about the features can be found in Appendix D and on the supplemental website https://compstat-lmu.shinyapps.io/Personality_Prediction/ for Stachl et al. (2020a)

depends heavily on how well the fitted model captures the true underlying relationship between features. Especially when using random forests, we showed that all of the methods lead to misleading results. (4) Groups with many features might tend to have a higher grouped importance score than groups with fewer features. Normalizing the grouped importance score leads to an average score per feature. However, this might result in choosing groups where grouped scores are smaller than those of other groups and, hence, contain less (performance-based) information than others. When using GSI, users can extract additional feature-level information to gain more insights into the group scores. Specifically, we showed that single feature Shapley importance scores add up to GSI when no higher-order interactions between groups are present. As third main contribution we proposed the CFEP, which is another global interpretation method that allows visualizations of the combined effect of multiple features on the prediction of an ML model. By applying a sparse SPCA, we received more meaningful and interpretable results for the final CFEPs compared to its unsupervised counterpart. We also demonstrated the suitability of the method in our real data example from computational psychology. Although, we only considered a numeric feature space here, all methods are in general also applicable to mixed feature spaces. However, in the presence of categorical features, a suitable dimension reduction method for CFEP must be chosen.

Here, we have focused on knowledge-driven feature groupings. However, the introduced methods could also be applied to data-driven groups (e.g., via shared variance).

<span style="float:right">⌂ Springer</span>

Notably, their interpretation is only meaningful if groups can be described by some underlying factor. This might be a good application for interpretable latent variables to find causal relationships between feature groups and predictions of ML models. Additionally, with regard to highly correlated feature groups that cannot be grouped naturally, a data-driven approach might be more suitable.

It is our goal that this article not only provides a helpful reference for researchers in selecting appropriate interpretation methods when features can be grouped, but also that it inspires future research in this area.

## Declarations

## Appendix A Motivational example for grouped importance methods

In some settings, permuting single features individually might not be meaningful, for example, when categorical features are dummy-encoded. Table 6 shows for two

**Table 6** We draw 1000 samples of two independent categorical random variables $X_1, X_2 \in \{1, 2, 3, 4\}$ where the categories 1 and 2 occur four times more frequently than 3 and 4

| Method | $X_1$ | $X_{2,2}, X_{2,3}, X_{2,4}$ | $X_{2,2}$ | $X_{2,3}$ | $X_{2,4}$ |
|---|---|---|---|---|---|
| Individually permuted | 2.63 | – | 2.45 | 1.00 | 1.71 |
| Group-wise permuted | 2.63 | 2.65 | | – | |

Consider the target $y = 5 \cdot \mathbb{1}_{X_1 \neq 1} + 5 \cdot \mathbb{1}_{X_2 \neq 1} + \epsilon$ with $\epsilon \sim \mathcal{N}(0, 1)$. Both categorical features have the same influence on the target. We explicitly dummy encode $X_2$ using $X_2 = 1$ as the reference category to obtain 3 binary features $X_{2,k} = \mathbb{1}_{X_2=k}, k \in \{2, 3, 4\}$. We fit a linear model using the categorical feature $X_1$ and the binary features $X_{2,2}, X_{2,3}, X_{2,4}$. Here, we want to illustrate why it makes more sense to permute the 3 binary features jointly rather than individually, since they naturally belong together. As expected, permuting the binary features $X_{2,2}, X_{2,3}, X_{2,4}$ jointly as a group yields a comparable importance to $X_1$. However, permuting each binary feature individually gives different importance scores making it unclear how important $X_2$ is compared to $X_1$

equally important categorical features that if one feature is dummy-encoded (here: $X_2$), then all resulting binary features must be permuted as a group to obtain a comparable importance score to $X_1$. Hence, settings like in Table 6 or as described in Sects. 1 or 1.2 point out the need of grouped importance methods.

## Appendix B Shapley importance

### B.1 Properties of the grouped Shapley importance

For single features[12] $x_i \in \{1, \ldots, p\}$, which are divided into $l$ groups, we define the marginal contribution for $x_i$ as

$$\Delta_{\{x_i\}}(S) = v(S \cup \{x_i\}) - v(S),$$

for $S \subset \{1, \ldots, p\} \backslash \{x_i\}$. The Shapley importance for single features $\phi(x_i)$ can also be defined analogously to (15). One interesting question is, does the GSI for a group $G \subset \{1, \ldots, p\}$ decompose into the sum of Shapley importances of features in $G$? In the following, we want to analyze the remainder

$$R = \phi(G) - \sum_{i \in G} \phi(x_i). \tag{B1}$$

Similar to the functional ANOVA decomposition (Hooker 2004), we assume, that the value function for a coalition $S \subset \{1, \ldots, p\}$ can be broken down into main and interaction effects

$$v(S) = v_0 + \sum_{x_i \in S} v(x_i) + \sum_{i \neq j} \epsilon_{ij} + \sum_{i \neq j \neq k} \epsilon_{ijk} + \cdots, \tag{B2}$$

---

[12] Remember the one-to-one association of the numbers $1, \ldots, p$ and the features $\mathbf{x}_1, \ldots, \mathbf{x}_p$

where $\epsilon_{i...m}$ is the effect of the interaction between the features $x_i, \ldots, x_m \in S$. A needed requirement to apply this decomposition is that each of the functional terms has zero means, hence they need to be centralized. The considered intercept shift is stored in $v_0$. To receive a unique decomposition, the orthogonality between the functional terms needs to be fulfilled which is not the case in the presence of correlated features. Hooker (2007) therefore suggests the generalized functional ANOVA which replaces the orthogonality property with a hierarchical orthogonality condition and which is a weighted version of the standard functional ANOVA (Hooker 2004). However, we do not try to estimate or calculate the decomposed function terms, we only use the (valid) assumption that a function can be decomposed as in Eq. (B2) to show how GSI relates to Shapley importance for individual features. Hence, we are not directly interested in a unique solution of the decomposition.

With the assumption in Eq. (B2), it follows that the Shapley importance of a single feature $x_1$ (without loss of generality) can be written as

$$\phi(x_1) = v(x_1) + \frac{1}{2}\left(\sum_{i \neq 1}^{p} \epsilon_{1i}\right) + \frac{1}{3}\left(\sum_{i \neq j \neq 1}^{p} \epsilon_{1ij}\right) + \cdots + \frac{1}{p}\epsilon_{1...p}. \quad \text{(B3)}$$

The value function of the feature $x_1$ contributes to the Shapley importance with the weight 1 and all possible interaction effects with feature $x_1$ contribute with the reciprocal length of the interaction effect. We proved this assertion in Appendix B.2. Similar to (B3), the GSI of a group $G_1$ (w.l.o.g.) can be written as

$$\phi(G_1) = v(G_1) + \frac{1}{2}\left(\sum_{i \neq 1}^{k} \epsilon_{G_1 G_i}\right) + \frac{1}{3}\left(\sum_{i \neq j \neq 1}^{k} \epsilon_{G_1 G_i G_j}\right) + \cdots + \frac{1}{k}\epsilon_{G_1...G_k} \quad \text{(B4)}$$

where $\epsilon_{G_1...G_k}$ is the (non-computable) interaction effect between features of groups $G_1, \ldots, G_k$, where each group provides at least one feature. By using Eq. (B2) on $v(G_1)$, we get:

$$v(G_1) = \sum_{i \in G_1} v(x_i) + \sum_{i \neq j \in G_1} \epsilon_{ij} + \sum_{i \neq j \neq k \in G_1} \epsilon_{ijk} + \cdots \quad \text{(B5)}$$

Looking back at Eq. (B1), a lot of terms cancel out by using Eqs. (B3) and (B5). The term $v(G_1)$, meaning all main effects $v(x_i), i \in G_1$, and all interaction effects $\epsilon_{i,...,k}, 1 \leq k \leq |G_1|$ between features within $G_1$, cancels out entirely.[13] Furthermore, at least all two-way interaction effects between groups $\epsilon_{G_1 G_i}, i = 2, \ldots, k$ cancel out. A combination of higher-order interaction terms between features of $G_1$ and $\{1, \ldots, p\}\backslash G_1$ remain.[14] This means that the remainder $R$ is (usually) not equal to zero in case the applied algorithm learned a higher-order interaction between features

---

[13] Note, $v(G_1)$ cancels out, meaning that these interaction terms cannot be computed directly but are assumed to affect the "payout" of the value function.

[14] They mostly only partly cancel out, depending on the number of features within the groups $G_1, \ldots, G_k$.

of the regarded group and other groups. The higher the remainder, the larger the higher-order interaction effect. Thus, the remainder can be used as a quantification of learned higher-order interaction effects between features of different groups.

## B.2 Proof of Properties

Assume, that the value function for a coalition $S \subset \{x_1, \ldots, x_p\}$ can be broken down into main and interaction effects:

$$v(S) = \sum_{x_i \in S} v(x_i) + \sum_{i_1 \neq i_2} \epsilon_{i_1 i_2} + \sum_{i_1 \neq i_2 \neq i_3} \epsilon_{i_1 i_2 i_3} + \cdots,$$

the Shapley importance of a single feature $x_1$ can be written as

$$\phi(x_1) = v(x_1) + \frac{1}{2}\left(\sum_{i \neq 1}^{p} \epsilon_{1i}\right) + \frac{1}{3}\left(\sum_{i \neq j \neq 1}^{p} \epsilon_{1ij}\right) + \cdots + \frac{1}{p}\epsilon_{1\ldots p}.$$

**Proof** Let $N = \{x_2, \ldots, x_p\}$. The general formula for the Shapley importance is given by:

$$\phi_p(x_1) = \sum_{S \subset N \setminus \{x_1\}} \frac{(p-1-|S|)! \cdot |S|!}{p!} \left(v(S \cup \{x_1\}) - v(S)\right) \tag{B6}$$

With assumption (B2) the term $v(S \cup \{x_1\}) - v(S)$ will reduce to:

$$v(S \cup \{x_1\}) - v(S) = v(x_1) + \sum_{i_1 \neq 1}^{p} \epsilon_{1i_1} + \cdots + \sum_{i_1 \neq \cdots \neq i_{|S|} \neq 1}^{p} \epsilon_{1i_1 \ldots i_{|S|}} \tag{B7}$$

It is the sum of $v(x_1)$ and all interactions with feature $x_1$ of sizes $2, \ldots, |S|+1$. All other terms without feature $x_1$ cancel out.

Equation (B6) consists of many summands of the form (B7). The term $v(x_1)$ appears once for every subset $S \subset N \setminus \{x_1\}$. There are $\binom{p-1}{|S|}$ different subsets of size $|S|$. Only looking at the summands with the term $v(x_1)$, Eq. (B6) reduces to

$$\sum_{|S|=0}^{p-1} \frac{(p-1-|S|)! \cdot |S|!}{p!} \binom{p-1}{|S|} v(x_1) = v(x_1). \tag{B8}$$

For the interaction terms, we first start counting the interaction term $\epsilon_{12}$ of size 2, as an example. For $|S| = 0$, there are zero terms of $\epsilon_{12}$. For $|S| = 1$, the term $\epsilon_{12}$ only appears once, when $S = \{x_2\}$. For $|S| = 2$, the term $\epsilon_{12}$ appears $p - 2$ times, once for each subset $S = \{x_2, x_j\}$, for $3 \leq j \leq p$. For $|S| = 3$, we have $\binom{p-2}{2}$ times the term $\epsilon_{12}$, again, once for each subset $S = \{x_2, x_j, x_k\}$, for $3 \leq j \neq k \leq p$. This pattern goes on until there are $\binom{p-2}{p-2}$ terms of $\epsilon_{12}$ for $|S| = p - 1$. Now, we look at

the interaction terms $\epsilon_{1i_1\ldots i_{k-1}}$ of size $k$. Following the pattern, which we just derived, there are zero terms of $\epsilon_{1i_1\ldots i_{k-1}}$ for $|S| \leq k - 2$ and $\binom{p-k}{|S|-k+1}$ terms of $\epsilon_{1i_1\ldots i_{k-1}}$ for $k \leq |S| \leq p - 1$. If we only look at the interaction terms $\epsilon_{1i_1\ldots i_{k-1}}$ of size $k$ and following the Eq. (B6), we get

$$\sum_{|S|=k-1}^{p-1} \frac{(p-1-|S|)! \cdot |S|!}{p!} \binom{p-k}{|S|-k+1} \epsilon_{1i_1\ldots i_{k-1}} = \frac{1}{k} \epsilon_{1i_1\ldots i_{k-1}},$$

which was left to show the assertion. □

## Appendix C More details on dimension reduction techniques

### C.1 Principal component analysis

PCA only considers the data matrix $\mathbf{X}$ and does not take the target vector $\mathbf{Y}$ into account. This procedure is thus unsupervised.

Given a centering Matrix

$$\mathbf{H} = \mathbf{I} - n^{-1}ee^T, \tag{C9}$$

where $e$ is an $n$-dimensional vector of all ones. The centered matrix is $\mathbf{X}_C = \mathbf{H}\mathbf{X}$. The sample covariance matrix of $\mathbf{X}$ can be written as:

$$\mathbf{S_X} := \frac{1}{n}\mathbf{X}_C^{\mathsf{T}}\mathbf{X}_C = \frac{1}{n}\mathbf{X}^{\mathsf{T}}\mathbf{HHX} \tag{C10}$$

The goal is to maximize the total variance of projected data, which is equivalent to maximizing trace of the sample covariance matrix. Equation (C10) can also be written as $\mathbf{S_X} = \frac{1}{n}\sum_{i=1}^{n}\mathbf{x}_C^{(i)}\mathbf{x}_C^{(i)\mathsf{T}}$, where $\mathbf{x}_C^{(i)}$ corresponds to the $i-$th row of $\mathbf{X}_C$. By projecting each data point by some unknown vectors $\mathbf{v}_j$, $j = 1, \ldots, p$, we get the projected variance for each $j = 1, \ldots, p$, which is:

$$\frac{1}{n}\sum_{i=1}^{n}\mathbf{v}_j^{\mathsf{T}}\mathbf{x}_C^{(i)}\mathbf{x}_C^{(i)\mathsf{T}}\mathbf{v}_j = \mathbf{v}_j^{\mathsf{T}}\left(\frac{1}{n}\sum_{i=1}^{n}\mathbf{x}_C^{(i)}\mathbf{x}_C^{(i)\mathsf{T}}\right)\mathbf{v}_j = \mathbf{v}_j^{\mathsf{T}}\mathbf{S_X}\mathbf{v}_j.$$

Let $\mathbf{V} \in \mathbb{R}^{p \times p}$ be the full projection matrix. The projected total variance is $tr(\mathbf{V}^{\mathsf{T}}\mathbf{S_X}\mathbf{V})$, and by ignoring constant terms, PCA finds a solution to the problem

$$\underset{\mathbf{V}}{\text{argmax}}\, tr(\mathbf{V}^{\mathsf{T}}\mathbf{S_X}\mathbf{V}) = \underset{\mathbf{V}}{\text{argmax}}\, tr(\mathbf{V}^{\mathsf{T}}\mathbf{X}^{\mathsf{T}}\mathbf{HHX}\mathbf{V}) \tag{C11}$$

with an Eigen decomposition of the covariance matrix $\mathbf{S_X}$. The resulting Eigen vectors thus maximize the variation of projected data.

## C.2 Measuring statistical dependence with Hilbert Schmidt norms

In Gretton et al. (2005) a more generalized measure of dependence between variables X and Y was introduced:

Two random variables $X$ and $Y$ are independent if and only if any bounded continuous function of them are uncorrelated.

In more detail, this means that any pairs $(X, Y), (X, Y^2), (X^2, Y), (cos(X), log(Y)), ...$ have to be uncorrelated. The resulting independence measure is called the Hilbert-Schmidt Independence Criterion (HSIC). For the analysis of this independence measure, it is necessary to analyze functions on random variables. Therefore theory of Hilbert spaces and concepts of functional analysis are necessary for a thorough analysis, but they are not part of this paper. For an extensive discussion of Hilbert spaces, especially reproducing kernel hilbert spaces (RKHS) we refer to Hein and Bousquet (2004).

Let $\mathcal{F}$ be a separable RKHS containing all bounded continuous functions from $\mathcal{X}$ to $\mathbb{R}$. The associated kernel shall be denoted by $\mathbf{K} \in \mathbb{R}^{n \times n}$, with $\mathbf{K}_{ij} = k(x_i, x_j)$. Concurrently, let $\mathcal{G}$ be a separable RKHS with bounded continuous functions from $\mathcal{Y}$ to $\mathbb{R}$ and associated kernel $\mathbf{L} \in \mathbb{R}^{n \times n}$, with $\mathbf{L}_{ij} = l(y_i, y_j)$.

We are particularly interested in the cross variance between $f$ and $g$:

$$Cov(f(x), g(y)) = \mathbb{E}_{x,y}[f(x)g(y)] - \mathbb{E}_x[f(x)]\mathbb{E}_y[g(y)] \tag{C12}$$

A function, which maps one element from one hilbert space to another hilbert space is called *operator*. A theorem (see e.g. Fukumizu et al. 2004) states, that there exists a unique operator $C_{X,Y} : \mathcal{G} \longrightarrow \mathcal{F}$ with

$$\langle f, C_{x,y}(g) \rangle_{\mathcal{F}} = Cov(f(x), g(y)). \tag{C13}$$

The Hilbert-Schmidt Independence Criterion (HSIC) is defined as the squared Hilbert-Schmidt norm of the cross-covariance operator C:

$$\text{HSIC}(P_{\mathcal{X}, \mathcal{Y}}, \mathcal{F}, \mathcal{G}) = \|C_{x,y}\|_{HS}^2 \tag{C14}$$

$\|C_{x,y}\|_{HS}^2 = 0$ if and only if the random variables $\mathcal{X}$ and $\mathcal{Y}$ are independent. For a detailed discussion and derivation of the HSIC independence measure, we refer to Gretton et al. (2005). The HSIC measure was used for feature selection in Song et al. (2007) or for supervised principal components in Barshan et al. (2011).

### C.2.1 Empirical HSIC

For a dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^{n}$ the empirical HSIC is:

$$HSIC(\mathcal{D}, F, G) = (n-1)^{-2} tr(\mathbf{KHLH}) = (n-1)^{-2} tr(\mathbf{HKHL}), \tag{C15}$$

where $\mathbf{H}$ is the centering matrix from (C9). A high level of dependency between two kernels yields a high HSIC value.

<span style="float:right">&#x2A02; Springer</span>

## C.3 Supervised sparse principal components

In the process of finding interpretable latent variables, which also incorporate dependencies to a target variable, the Sparse Supervised Principal Components (SPCA), which was introduced in Sharifzadeh et al. (2017), is a suitable method for our application.

For sparse SPCA the kernel matrix $K$ ist defined as $K = XVV^\mathsf{T}X^\mathsf{T}$ with a constraint for unit length and an $L_1$ penalty for sparsity. By ignoring constant terms, we get the optimization problem:

$$\underset{\mathbf{V}}{\operatorname{argmax}}\, tr(\mathbf{HKHL}) = \underset{\mathbf{V}}{\operatorname{argmax}}\, tr(\mathbf{HXVV^\mathsf{T}X^\mathsf{T}HL}) \tag{C16}$$

$$= \underset{\mathbf{V}}{\operatorname{argmax}}\, tr(\mathbf{V^\mathsf{T}X^\mathsf{T}HLHXV}) \tag{C17}$$

$$s.t. \quad \mathbf{V^\mathsf{T}V = I}, \quad |\mathbf{V}| \le c. \tag{C18}$$

Note, that without the sparsity constraint, (C17) reduces to (C11), when choosing $\mathbf{L} = \mathbf{I}$. Already explained in Barshan et al. (2011), PCA is a special form of their Supervised PCA, where setting $\mathbf{L} = \mathbf{I}$ is a kernel, which only captures similarity between a point and itself. Maximizing dependency between $\mathbf{K}$ and the identiy matrix corresponds to retaining maximal diversity between observations.

Now, an arbitrary $\mathbf{L}$ can be decomposed as $\mathbf{L} = \Delta\Delta^\mathsf{T}$, since $\mathbf{L}$, as a kernel matrix, is positive definite and symmetric. Defining $\Psi := \Delta^\mathsf{T}\mathbf{HX} \in \mathbb{R}^{n \times p}$, the objective function (C17) can be rewritten as:

$$\underset{\mathbf{V}}{\operatorname{argmax}}\, tr(\mathbf{V^\mathsf{T}\Psi^\mathsf{T}\Psi V})\, s.t.\, \mathbf{V^\mathsf{T}V = I},\, |\mathbf{V}| \le c. \tag{C19}$$

Using the singular value decomposition (SVD), the matrix $\Psi$ with $\text{rank}(\Psi) = m \le n$ can be written as a product of matrices:

$$\Psi = \mathbf{U}\Lambda\mathbf{V^\mathsf{T}} \quad s.t. \quad \mathbf{U^\mathsf{T}U} = I_n, \mathbf{VV^\mathsf{T}} = I_p, \Lambda = I(\lambda_1, \dots, \lambda_m, 0, \dots, 0), \tag{C20}$$

where $\mathbf{U} \in \mathbb{R}^{n \times n}$ and $\mathbf{V} \in \mathbb{R}^{p \times p}$ are orthogonal matrices, and $\Lambda \in \mathbb{R}^{n \times p}$ is a diagonal matrix, with descending diagonal entries $\lambda_1 \ge \lambda_2 \ge \cdots \ge \lambda_m \ge 0$. It is easy to see that the columns of $\mathbf{V}$ are Eigen vectors of the matrix $\Psi^\mathsf{T}\Psi$, since the following Eigen value decomposition holds:

$$\Psi^\mathsf{T}\Psi = \mathbf{V}\Lambda\mathbf{U^\mathsf{T}U}\Lambda\mathbf{V^\mathsf{T}} = \mathbf{V}(\Lambda^2)\mathbf{V^\mathsf{T}}. \tag{C21}$$

The sparse SPCA problem (C19) now becomes a matrix decomposition problem of the matrix $\Psi$, when adding an $L_1$ penalty on the matrix $\mathbf{V}$, since the columns of $\mathbf{V}$, being Eigen vectors of $\Psi^\mathsf{T}\Psi$, maximize $tr(\mathbf{V^\mathsf{T}\Psi^\mathsf{T}\Psi V})$.

With an $L_1$ penalty on $\mathbf{V}$, this problem is a *penalized matrix decomposition* problem (PMD, Witten et al. (2009)).

⁂ Springer

Recalling our original problem of finding interpretable latent variables that also depend on a target variable, the rank $m$ matrix decomposition of $\Psi$ may not be desirable. It can be shown (e.g. Eckart and Young 1936) that the best low rank ($r \leq m$) approximation of $\Psi$ is calculated by the first $r$ singular values of $\Lambda$ and the first $r$ singular vectors of $\mathbf{U}$ and $\mathbf{V}$. With $\mathbf{u}_i$ being the $i-$th column of $\mathbf{U}$ and $\mathbf{v}_i$ being the $i-$th column of $\mathbf{V}$, the best low rank approximation can thus be written as:

$$\sum_{i=1}^{r} \lambda_i \mathbf{u}_i \mathbf{v}_i^{\mathsf{T}} = \underset{\hat{\Psi}}{\mathrm{argmin}} \|\Psi - \hat{\Psi}\|_F^2, \tag{C22}$$

subject to the squared Frobenius-norm ($A \in \mathbb{R}^{m \times n}$: $\|A\|_F^2 = \sum_{i=1}^{n} \sum_{j=1}^{m} |a_{ij}|^2$). The following equality was demonstrated in Witten et al. (2009):

$$\frac{1}{2}\|\Psi - \mathbf{U}\Lambda\mathbf{V}^{\mathsf{T}}\|_F^2 = \frac{1}{2}\|\Psi\|_F^2 - \sum_{i=1}^{r} \mathbf{u}_i^{\mathsf{T}} \Psi \mathbf{v}_i \lambda_i + \frac{1}{2}\sum_{i=1}^{r} \lambda_i^2. \tag{C23}$$

The minimization problem (C22) thus becomes a maximization problem, by ignoring the constant terms. Sharifzadeh et al. (2017) added additional $L_2$ constraints on $\mathbf{u}_i$ and $\mathbf{v}_i$, an $L_1$ constaint on $v_i$ for sparsity and an orthogonality constraint for $u_i$:

$$\underset{\mathbf{u}_i \mathbf{v}_i}{\mathrm{argmax}} \; \mathbf{u}_i^{\mathsf{T}} \Psi \mathbf{v}_i \; s.t. \|\mathbf{u}_i\|_2 \leq 1, \|\mathbf{v}_i\|_2 \leq 1, \|\mathbf{v}_i\|_1 \leq c, \mathbf{u}_i \perp \mathbf{u}_1, \ldots, \mathbf{u}_{i-1} \tag{C24}$$

The $L_2$ constraints do not force unit length to avoid non convex optimization problems. Witten et al. (2009) discuss how to solve many penalized matrix decomposition problems of this kind. Without the orthogonality constraint, they call this particular problem PMD(., $L_1$). The solution to this problem is discussed in detail in Sharifzadeh et al. (2017). A software implementation is available with the R-package PMA by Witten and Tibshirani (2020), which we will use for our demonstrations. Problem (C24) does not yield orthogonal sparse vectors $\mathbf{v}_i$, Witten et al. (2009) state that these vectors are unlikely to be very correlated, since the vectors $\mathbf{v}_i$ are associated with orthogonal vectors $\mathbf{u}_i$, $i = 1, \ldots, r$.

### C.3.1 Choice of the Kernel

For sparse SPCA the kernel $\mathbf{K}$ has been predefined as. The choice of the kernel $\mathbf{L}$, however, has a decisive impact on how the dependencies are modeled. Song et al. (2012) discuss the kernel choice for different situations. For binary classification, one may simply choose

$$l(y_i, y_j) = y_i y_j, \; \text{where } y_i, y_j \in \{\pm 1\}, \tag{C25}$$

or a weighted version, giving different weights on positive and negative labels. For multiclass classification a possible kernel is

$$l(y_i, y_j) = c_y \delta_{y_i, y_j}, \; \text{where } c_y > 0. \tag{C26}$$

For regression one can also use a linear kernel $l(y_i, y_j) = y_i, y_j$, but then only simple linear correlations between features and the target variable can be detected. A more universal choice is the radial basis function (RBF) kernel:

$$l(y_i, y_j) = exp\left(-\frac{\|y_i - y_j\|^2}{2\sigma^2}\right). \tag{C27}$$

The choice of the bandwidth $2\sigma^2$ is extremely important. For example, if $2\sigma^2 \to 0$, the matrix L becomes the identity matrix. Or if $2\sigma^2 \to \infty$, all entries of $\mathbf{L}$ are 1. In both cases, all relevant information of the dependency between features and the target variable is lost. Besides the bandwidth $2\sigma$, the kernel matrix $L$ depends only on the pairwise distances $\|y_i - y_j\|^2$. A reasonable, and heuristically well performing (Pfister et al. 2017) choice is $2\sigma^2 = \text{median}(\|y_i - y_j\|^2 : i > j)$. However, it might also be possible and advantageous to use other kernels that are selected to be particularly efficient in detecting certain kinds of dependencies.

### C.3.2 Choice of c

Witten et al. (2009) explained how PMD can be used to impute missing data. The main idea is simply to exclude missing entries from the maximization problem (C24) and impute missing values by the low rank approximation matrix $\mathbf{U}\Lambda\mathbf{V}^\mathsf{T}$. This procedure can also be used for finding optimal values for $c$ by a cross-validation approach. The test data consists of leaving out some entries of the matrix $\Psi$ (not entire rows or columns, but individual elements of the matrix), yielding a matrix with missing entries $\tilde{\Psi}$. For candidate values $c_i, i = 1, \ldots, k$, calculate the PMD$(., L_1)$ and record the mean squared error over the missing elements of $\tilde{\Psi}$ and the estimate $\mathbf{U}\Lambda\mathbf{V}^\mathsf{T}$. The true values of the missing values of $\tilde{\Psi}$ are available in the original data $\Psi$. The optimal value $c^*$ corresponds to the best candidate value $c_j$, which minimizes the mean squared error.

However, such a cross-validation approach for the search for $c$ is not always necessary. If the method is used as a descriptive method to better understand the underlying structure of the data, a small value of $c$ can be chosen to achieve a desired sparsity.

### Appendix D Feature description for smartphone sensor data

See Table 7.

**Table 7** Description of features used for CFEPs in Sect. 7

| Feature | Description |
| --- | --- |
| daily_mean_num_unique_Weather_weekend | Mean number of different weather apps used each day on weekends |
| daily_mean_num_Weather | Mean number of weather apps used each day |
| daily_mean_num_unique_Weather_week | Mean number of different weather apps used each day on weekdays |
| daily_mean_num_unique_Weather | Mean number of different weather apps used each day |
| daily_mean_num_unique_apps | Mean number of different apps used each day |
| daily_mean_num_unique_apps_week | Mean number of different apps used each day on weekdays |
| daily_mean_num_unique_apps_weekend | Mean number of different apps used each day on weekends |
| daily_mean_sum_events_night | Number of all events during the night averaged for each day |
| daily_mean_dur_all | Duration of all events averaged for each day |
| daily_sd_sum_intereventall | Sd of the sum of all inter-event time intervals for each day |
| daily_mean_num_uniq_song | Mean number of different songs listened to each day |
| daily_mean_num_song | Mean number of songs listened to each day |
| daily_mean_duration_music | Mean duration of music apps used each day |

# References

Allaire J, Gandrud C, Russell K, et al (2017) networkD3: D3 JavaScript network graphs from R. https://CRAN.R-project.org/package=networkD3, R package version 0.4

Alon U, Barkai N, Notterman DA et al (1999) Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. Proc Natl Acad Sci 96(12):6745–6750

Amoukou SI, Brunel NJB, Salaün T (2021) The shapley value of coalition of variables provides better explanations. arXiv:2103.13342

Apley DW, Zhu J (2019) Visualizing the effects of predictor variables in black box supervised learning models. arXiv:1612.08468

Bair E, Hastie T, Paul D et al (2006) Prediction by supervised principal components. J Am Stat Assoc 101(473):119–137

Barshan E, Ghodsi A, Azimifar Z et al (2011) Supervised principal component analysis: visualization, classification and regression on subspaces and submanifolds. Pattern Recogn 44(7):1357–1371

Berk R, Sherman L, Barnes G et al (2009) Forecasting murder within a population of probationers and parolees: a high stakes application of statistical learning. J R Stat Soc A Stat Soc 172(1):191–211

Breiman L (2001) Random forests. Mach Learn 45(1):5–32

Brenning A (2021) Transforming feature space to interpret machine learning models. arXiv:2104.04295

Caputo B, Sim K, Furesjö F, et al (2002) Appearance-based object recognition using SVMS: Which kernel should I use. In: Proceedings of the NIPS workshop on statistical methods for computational experiments in visual processing and computer vision, Red Hook, NY, USA

Casalicchio G, Molnar C, Bischl B (2019) Visualizing the feature importance for black box models. Springer International Publishing. Machine Learning and Knowledge Discovery in Databases, pp 655–670

Chakraborty D, Pal NR (2008) Selecting useful groups of features in a connectionist framework. IEEE Trans Neural Netw 19(3):381–396

Cohen SB, Ruppin E, Dror G (2005) Feature selection based on the Shapley value. In: Kaelbling LP, Saffiotti A (eds) IJCAI-05, Proceedings of the nineteenth international joint conference on artificial intelligence, Edinburgh, Scotland, UK, July 30–August 5, 2005. Professional Book Center, pp 665–670

Covert I, Lundberg SM, Lee SI (2020) Understanding global feature contributions with additive importance measures. Adv Neural Inf Process Syst 33:17212–17223

de Mijolla D, Frye C, Kunesch M, et al (2020) Human-interpretable model explainability on high-dimensional data. CoRR arXiv:2010.07384

Eckart C, Young G (1936) The approximation of one matrix by another of lower rank. Psychometrika 1(3):211–218

Fisher A, Rudin C, Dominici F (2019) All models are wrong, but many are useful: learning a variable's importance by studying an entire class of prediction models simultaneously. J Mach Learn Res 20(177):1–81

Friedman JH (2001) Greedy function approximation: a gradient boosting machine. Ann Stat, 1189–1232

Friedman J, Hastie T, Tibshirani R (2010) A note on the group lasso and a sparse group lasso. arXiv:1001.0736

Fuchs K, Scheipl F, Greven S (2015) Penalized scalar-on-functions regression with interaction term. Comput Stat Data Anal 81:38–51

Fukumizu K, Bach FR, Jordan MI (2004) Dimensionality reduction for supervised learning with reproducing kernel Hilbert spaces. J Mach Learn Res 5:73–99

Goldberg LR (1990) An alternative "description of personality": the big-five factor structure. J Person Soc Psychol 59:1216–1229

Goldstein A, Kapelner A, Bleich J et al (2013) Peeking inside the black box: visualizing statistical learning with plots of individual conditional expectation. J Comput Gr Stat 24:44–65

Gregorova M, Kalousis A, Marchand-Maillet S (2018) Structured nonlinear variable selection. In: Globerson A, Silva R (eds) Proceedings of the thirty-fourth conference on uncertainty in artificial intelligence, UAI 2018, Monterey, California, USA, August 6–10, 2018. AUAI Press, pp 23–32

Gregorutti B, Michel B, Saint-Pierre P (2015) Grouped variable importance with random forests and application to multiple functional data analysis. Comput Stat Data Anal 90:15–35

Gretton A, Bousquet O, Smola A, et al (2005) Measuring statistical dependence with Hilbert-Schmidt norms. In: International conference on algorithmic learning theory. Springer, pp 63–77

Guyon I, Weston J, Barnhill S et al (2002) Gene selection for cancer classification using support vector machines. Mach Learn 46(1–3):389–422

Harari GM, Gosling SD, Wang R et al (2015) Capturing situational information with smartphones and mobile sensing methods. Eur J Pers 29(5):509–511

Harari GM, Lane ND, Wang R et al (2016) Using smartphones to collect behavioral data in psychological science: opportunities, practical considerations, and challenges. Perspect Psychol Sci 11(6):838–854

Harari GM, Müller SR, Aung MS et al (2017) Smartphone sensing methods for studying behavior in everyday life. Curr Opin Behav Sci 18:83–90

Harari GM, Müller SR, Stachl C et al (2019) Sensing sociability: individual differences in young adults' conversation, calling, texting, and app use behaviors in daily life. J Person Soc Psychol 119:204

He Z, Yu W (2010) Stable feature selection for biomarker discovery. Comput Biol Chem 34:215–225

Hein M, Bousquet O (2004) Kernels, Associated structures and generalizations, Max Planck Institute for Biological Cybernetics

Shipp MA, Ross KN, Tamayo P et al (2002) Diffuse large b-cell lymphoma outcome prediction by gene-expression profiling and supervised machine learning. Nat Med 8(1):68–74

Hooker G (2004) Discovering additive structure in black box functions. In: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining. ACM, pp 575–580

Hooker G (2007) Generalized functional anova diagnostics for high-dimensional functions of dependent variables. J Comput Graph Stat 16(3):709–732

Hooker G, Mentch L (2019) Please stop permuting features: an explanation and alternatives. arXiv:1905.03151

Jackson JJ, Wood D, Bogg T et al (2010) What do conscientious people do? Development and validation of the behavioral indicators of conscientiousness (bic). J Res Pers 44(4):501–511

Jaeger J, Sengupta R, Ruzzo W (2003) Improved gene selection for classification of microarrays. Pac Symp Biocomput Pac Symp Biocomput 8:53–64

Jolliffe IT (1986) Principal component analysis. Springer, New York

Kolenik T, Gams M (2021) Intelligent cognitive assistants for attitude and behavior change support in mental health: state-of-the-art technical review. Electronics 10(11):1250

Lei J, G'Sell M, Rinaldo A et al (2018) Distribution-free predictive inference for regression. J Am Stat Assoc 113(523):1094–1111

Lipton ZC (2018) The mythos of model interpretability: in machine learning, the concept of interpretability is both important and slippery. Queue 16(3):31–57

Lozano AC, Abe N, Liu Y et al (2009) Grouped graphical granger modeling for gene expression regulatory networks discovery. Bioinformatics 25(12):i110–i118

Lundberg SM, Erion GG, Lee S (2018) Consistent individualized feature attribution for tree ensembles. CoRR arXiv:1802.03888

Lundberg SM, Lee SI (2017) A unified approach to interpreting model predictions. In: Proceedings of the 31st international conference on neural information processing systems. Curran Associates Inc., Red Hook, NY, USA, NIPS'17, pp 4768–4777

Meier L, Van De Geer S, Bühlmann P (2008) The group lasso for logistic regression. J R Stat Soc Ser B (Stat Methodol) 70(1):53–71

Meinshausen N, Bühlmann P (2010) Stability selection. J R Stat Soc Ser B (Stat Methodol) 72(4):417–473

Miller G (2012) The smartphone psychology manifesto. Perspect Psychol Sci 7(3):221–237

Molnar C (2019) Interpretable machine learning. https://christophm.github.io/interpretable-ml-book/

Molnar C, König G, Bischl B, et al (2020a) Model-agnostic feature importance and effects with dependent features—a conditional subgroup approach. arXiv:2006.04628

Molnar C, König G, Herbinger J, et al (2020b) General pitfalls of model-agnostic interpretation methods for machine learning models. arXiv preprint arXiv:2007.04131

Nicodemus K, Malley J, Strobl C, et al (2010) The behaviour of random forest permutation-based variable importance measures under predictor correlation. BMC Bioinform 11–110

Onnela JP, Rauch SL (2016) Harnessing smartphone-based digital phenotyping to enhance behavioral and mental health. Neuropsychopharmacology 41(7):1691–1696

Ozer DJ, Benet-Martínez V (2006) Personality and the prediction of consequential outcomes. Annu Rev Psychol 57:401–421

Park MY, Hastie T, Tibshirani R (2006) Averaged gene expressions for regression. Biostatistics 8(2):212–227

Pfister N, Bühlmann P, Schölkopf B et al (2017) Kernel-based tests for joint independence. J R Stat Soc Ser B (Stat Methodol) 80(1):5–31

Rachuri KK, Musolesi M, Mascolo C, et al (2010) Emotionsense: a mobile phones based adaptive platform for experimental social psychology research. In: UbiComp'10—Proceedings of the 2010 ACM conference on ubiquitous computing

Raento M, Oulasvirta A, Eagle N (2009) Smartphones: an emerging tool for social scientists. Sociol Methods Res 37(3):426–454

Rapaport F, Barillot E, Vert JP (2008) Classification of Arraycgh data using fused SVM. Bioinformatics 24(13):i375–i382

Saeb S, Lattie EG, Schueller SM et al (2016) The relationship between mobile phone location sensor data and depressive symptom severity. PeerJ 4:e2537

Schoedel R, Au Q, Völkel ST et al (2018) Digital footprints of sensation seeking. Zeitschrift für Psychologie 226(4):232–245

Schoedel R, Pargent F, Au Q et al (2020) To challenge the morning lark and the night owl: using smartphone sensing data to investigate day-night behaviour patterns. Eur J Personal 34:733–752

Scholbeck CA, Molnar C, Heumann C et al (2020) Sampling, intervention, prediction, aggregation: a generalized framework for model-agnostic interpretations. In: Cellier P, Driessens K (eds) Machine learning and knowledge discovery in databases. Springer, Cham, pp 205–216

Schuwerk T, Kaltefleiter LJ, Au JQ et al (2019) Enter the wild: autistic traits and their relationship to mentalizing and social interaction in everyday life. J Autism Dev Disorders 49:4193–4208

Seedorff N, Brown G (2021) totalvis: a principal components approach to visualizing total effects in black box models. SN Comput Sci 2(3):1–12

Servia-Rodríguez S, Rachuri KK, Mascolo C, et al (2017) Mobile sensing at the service of mental well-being: A large-scale longitudinal study. In: 26th international world wide web conference, WWW 2017. International World Wide Web Conferences Steering Committee, pp 103–112

Shapley LS (1953) A value for n-person games. Contrib Theory Games 2(28):307–317

Sharifzadeh S, Ghodsi A, Clemmensen LH et al (2017) Sparse supervised principal component analysis (sspca) for dimension reduction and variable selection. Eng Appl Artif Intell 65:168–177

Song L, Smola A, Gretton A et al (2012) Feature selection via dependence maximization. J Mach Learn Res 13:1393–1434

Song L, Smola A, Gretton A, et al (2007) Supervised feature selection via dependence estimation. In: Proceedings of the 24th international conference on Machine learning, pp 823–830

Stachl C, Hilbert S, Au JQ et al (2017) Personality traits predict smartphone usage. Eur J Pers 31(6):701–722

Stachl C, Au Q, Schoedel R et al (2020a) Predicting personality from patterns of behavior collected with smartphones. Proc Natl Acad Sci 117:17680–17687

Stachl C, Pargent F, Hilbert S et al (2020b) Personality research and assessment in the era of machine learning. Eur J Personal 34:613–631

Strobl C, Boulesteix AL, Kneib T et al (2008) Conditional variable importance for random forests. BMC Bioinform 9:307

Thomée S (2018) Mobile phone use and mental health; A review of the research that takes a psychological perspective on exposure. Int J Environ Res Public Health 15(12):2692

Tibshirani R (1996) Regression shrinkage and selection via the lasso. J Roy Stat Soc Ser B (Methodol) 58(1):267–288

Toloşi L, Lengauer T (2011) Classification with correlated features: unreliability of feature ranking and solutions. Bioinformatics 27(14):1986–1994

Tripathi S, Hemachandra N, Trivedi P (2020) Interpretable feature subset selection: a Shapley value based approach. In: Proceedings of 2020 IEEE international conference on big data, special session on explainable artificial intelligence in safety critical systems

Valentin S, Harkotte M, Popov T (2020) Interpreting neural decoding models using grouped model reliance. PLOS Comput Biol 16(1):e1007148

Venables B, Ripley B (2002) Modern applied statistics with S

Watson DS, Wright MN (2019) Testing conditional independence in supervised learning algorithms. arXiv:1901.09917

Williamson BD, Gilbert PB, Simon NR, et al (2020) A unified approach for inference on algorithm-agnostic variable importance. arXiv:2004.03683

Williamson B, Feng J (2020) Efficient nonparametric statistical inference on population feature importance using Shapley values. In: International conference on machine learning, PMLR, pp 10282–10291

Witten D, Tibshirani R (2020) PMA: penalized multivariate analysis. R Package Vers 1(2):1

Witten DM, Tibshirani R, Hastie T (2009) A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. Biostatistics 10(3):515–534

Wold S, Albano C, Dunn WJ et al (1984) Multivariate data analysis in chemistry. Springer, Dordrecht, pp 17–95

Yarkoni T, Westfall J (2017) Choosing prediction over explanation in psychology: lessons from machine learning. Perspect Psychol Sci 12(6):1100–1122

Yuan M, Lin Y (2006) Model selection and estimation in regression with grouped variables. J R Stat Soc Ser B (Stat Methodol) 68(1):49–67

# Multilabel Classification with `R` Package `mlr`

This article introduces a range of multilabel classification algorithms implemented in the `mlr` package.

### Contributing article[1]

Probst, P., Au, Q., Casalicchio, G., Stachl, C., and Bischl, B. (2017). Multilabel Classification with R Package mlr. *The R Journal*, 9(1):352–369

### Copyright information

This article is licensed under the Creative Commons Attribution 4.0 International License (CC BY 4.0). See `https://creativecommons.org/licenses/by/4.0/`. The license allows to copy and redistribute the material in any medium or format.

Jiew-Quay Au was primarily responsible for the content of the manuscript, although Philipp Probst took the role of the corresponding author. Therefore, his contribution can be indicated as follows:

### Declaration of contributions

The PhD candidate was responsible for drafting the manuscript. During this process, Jiew-Quay Au integrated various algorithms for multi-label prediction into the open-source software mlr and analyzing them in a benchmark experiment. He implemented problem transformation methods, enabling the adaptation of classical univariate algorithms for multi-label applications. Additionally, he was in charge of creating visualizations and tables and interpreting the results.

*Contribution of the coauthors*

Philipp Probst took on the role of the corresponding author and provided support in implementing alternative methods for multi-label prediction and conducting the comparative experiment. Giuseppe Casalicchio assisted in implementing the procedures into the software and

---

[1]Note: This publication was based on Jiew-Quay Au's master's thesis (`https://msnat.statistik.tu-dortmund.de/studium/abschlussarbeiten/`). The theoretical part of the master's thesis, including the implementation of the algorithms, are part of this publication.

offered valuable guidance in the experiment's elaboration. The remaining co-authors support-
ed the publication with valuable advice and in revising the manuscript.

### *Software:*

- `https://cran.r-project.org/web/packages/mlr/index.html`

- `https://mlr.mlr-org.com/`

- `https://mlr.mlr-org.com/articles/tutorial/multilabel.html`

# Multilabel Classification with R Package mlr

*by Philipp Probst, Quay Au, Giuseppe Casalicchio, Clemens Stachl and Bernd Bischl*

**Abstract** We implemented several multilabel classification algorithms in the machine learning package **mlr**. The implemented methods are binary relevance, classifier chains, nested stacking, dependent binary relevance and stacking, which can be used with any base learner that is accessible in **mlr**. Moreover, there is access to the multilabel classification versions of **randomForestSRC** and **rFerns**. All these methods can be easily compared by different implemented multilabel performance measures and resampling methods in the standardized **mlr** framework. In a benchmark experiment with several multilabel datasets, the performance of the different methods is evaluated.

## Introduction

Multilabel classification is a classification problem where multiple target labels can be assigned to each observation instead of only one, like in multiclass classification. It can be regarded as a special case of multivariate classification or multi-target prediction problems, for which the scale of each response variable can be of any kind, for example nominal, ordinal or interval.

Originally, multilabel classification was used for text classification (McCallum, 1999; Schapire and Singer, 2000) and is now used in several applications in different research fields. For example, in image classification, a photo can belong to the classes *mountain* and *sunset* simultaneously. Zhang and Zhou (2008) and others (Boutell et al., 2004) used multilabel algorithms to classify scenes on images of natural environments. Furthermore, gene functional classifications is a popular application of multilabel learning in the field of biostatistics (Elisseeff and Weston, 2002; Zhang and Zhou, 2008). Additionally, multilabel classification is useful to categorize audio files. Music genres (Sanden and Zhang, 2011), instruments (Kursa and Wieczorkowska, 2014), bird sounds (Briggs et al., 2013) or even emotions evoked by a song (Trohidis et al., 2008) can be labeled with several categories. A song could, for example, be classified both as a *rock song* and a *ballad*.

An overview of multilabel classification was given by Tsoumakas and Katakis (2007). Two different approaches exist for multilabel classification. On the one hand, there are algorithm adaptation methods that try to adapt multiclass algorithms so they can be applied directly to the problem. On the other hand, there are problem transformation methods, which try to transform the multilabel classification into binary or multiclass classification problems.

Regarding multilabel classification software, there is the **mldr** (Charte and Charte, 2015) R package that contains some functions to get basic characteristics of specific multilabel datasets. The package is also useful for transforming multilabel datasets that are typically saved as ARFF-files (Attribute-Relation File Format) to data frames and vice versa. This is especially helpful because until now only the software packages MEKA (Read and Reutemann, 2012) and Mulan (Tsoumakas et al., 2011) were available for multilabel classification and both require multilabel datasets saved as ARFF-files to be executed. Additionally, the **mldr** package provides a function that applies the binary relevance or label powerset transformation method which transforms a multilabel dataset into several binary datasets (one for each label) or into a multiclass dataset using the set of labels for each observation as a single target label, respectively. However, there is no R package that provides a standardized interface for executing different multilabel classification algorithms. With the extension of the **mlr** package described in this paper, it will be possible to execute several multilabel classification algorithms in R with many different base learners.

In the following section of this paper, we will describe the implemented multilabel classification methods and then give a practical instruction of how to execute these algorithms in **mlr**. Finally, we present a benchmark experiment that compares the performance of all implemented methods on several datasets.

## Multilabel classification methods implemented in mlr

In this section, we present multilabel classification algorithms that are implemented in the **mlr** package (Bischl et al., 2016), which is a powerful and modularized toolbox for machine learning in R. The package offers a unified interface to more than a hundred learners from the areas classification, regression, cluster analysis and survival analysis. Furthermore, the package provides functions and tools that facilitate complex workflows such as hyperparameter tuning (see, e.g., Lang et al., 2015) and

feature selection that can now also be applied to the multilabel classification methods presented in this paper. In the following, we list the algorithm adaptation methods and problem transformation methods that are currently available in **mlr**.

### Algorithm adaptation methods

The **rFerns** (Kursa and Wieczorkowska, 2014) package contains an extension of the random ferns algorithm for multilabel classification. In the **randomForestSRC** (Ishwaran and Kogalur, 2016) package, multivariate classification and regression random forests can be created. In the classification case, the difference to standard random forests is that a composite normalized Gini index splitting rule is used. Multilabel classification can be achieved by using binary encoding for the labels.

### Problem transformation methods

Problem transformation methods try to transform the multilabel classification problem so that a simple binary classification algorithm, the so-called base learner, can be applied.

Let $n$ be the number of observations, let $p$ be the number of predictor variables and let $Z = \{z_1, \ldots, z_m\}$ be the set of all labels. Observations follow an unknown probability distribution $\mathcal{P}$ on $\mathcal{X} \times \mathcal{Y}$, where $\mathcal{X}$ is a $p-$dimensional input space of arbitrary measurement scales and $\mathcal{Y} = \{0, 1\}^m$ is the target space. In our notation, $\mathbf{x}^{(i)} = \left( x_1^{(i)}, \ldots, x_p^{(i)} \right)^\top \in \mathcal{X}$ refers to the $i$-th observation and $\mathbf{x}_j = \left( x_j^{(1)}, \ldots, x_j^{(n)} \right)^\top$ refers to the $j$-th predictor variable, for all $i = 1, \ldots, n$ and $j = 1, \ldots, p$. The observations $\mathbf{x}^{(i)}$ are associated with their multilabel outcomes $\mathbf{y}^{(i)} = \left( y_1^{(i)}, \ldots, y_m^{(i)} \right)^\top \in \mathcal{Y}$, for all $i = 1, \ldots, n$. For all $k = 1, \ldots, m$, setting $y_k^{(i)} = 1$ indicates the relevance, i.e., the occurrence, of label $z_k$ for observation $\mathbf{x}^{(i)}$ and setting $y_k^{(i)} = 0$ indicates the irrelevance of label $z_k$ for observation $\mathbf{x}^{(i)}$. The set of all instances thus becomes $D = \left\{ \left( \mathbf{x}^{(1)}, \mathbf{y}^{(1)} \right), \left( \mathbf{x}^{(2)}, \mathbf{y}^{(2)} \right), \ldots, \left( \mathbf{x}^{(n)}, \mathbf{y}^{(n)} \right) \right\}$. Furthermore, $\mathbf{y}_k = \left( y_k^{(1)}, \ldots, y_k^{(n)} \right)^\top$ refers to the $k$-th target vector, for all $k = 1, \ldots, m$. Throughout this paper, we visualize multilabel classification problems in the form of tables ($n = 6$, $p = 3$, $m = 3$):

$$D \triangleq \quad \begin{array}{|ccc|ccc|} \hline x_1 & x_2 & x_3 & y_1 & y_2 & y_3 \\ \hline & & & 0 & 0 & 1 \\ & & & 1 & 0 & 1 \\ & & & 1 & 1 & 0 \\ & & & 1 & 1 & 1 \\ & & & 1 & 1 & 0 \\ & & & 1 & 1 & 0 \\ \hline \end{array} \tag{1}$$

The entries of $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3$ can be of any (valid) kind, like continuous, binary, or categorical. The table in (1) visualizes this as an empty gray background. The target variables are indicated by a red background and can only take the binary values 0 or 1.

### Binary relevance

The binary relevance method (BR) is the simplest problem transformation method. BR learns a binary classifier for each label. Each classifier $C_1, \ldots, C_m$ is responsible for predicting the *relevance* of their corresponding label by a 0/1 prediction:

$$C_k : \mathcal{X} \longrightarrow \{0, 1\}, \quad k = 1, \ldots, m$$

These binary prediction are then combined to a multilabel target. An unlabeled observation $\mathbf{x}^{(l)}$ is assigned the prediction $\left( C_1 \left( \mathbf{x}^{(l)} \right), C_2 \left( \mathbf{x}^{(l)} \right), \ldots, C_m \left( \mathbf{x}^{(l)} \right) \right)^\top$. Hence, labels are predicted independently of each other and label dependencies are not taken into account. BR has linear computational complexity with respect to the number of labels and can easily be parallelized.

### Modeling label dependence

In the problem transformation setting, the arguably simplest way (Montañés et al., 2014) to model label dependence is to condition classifier models not only on $\mathcal{X}$, but also on other label information. The idea is to augment the input space $\mathcal{X}$ with information of the output space $\mathcal{Y}$, which is available in the training step. There are different ways to realize this idea of augmenting the input space. In essence, they can be distinguished in the following way:

- Should the true label information be used? (True vs. predicted label information)

- For predicting one label $z_k$, should all other labels augment the input space, or only a subset of labels? (Full vs. partial conditioning)

### True vs. predicted label information

During the training of a classifier $C_k$ for the label $z_k$, the label information of other labels are available in the training data. Consequently, these true labels can directly be used as predictors to train the classifier. Alternatively, the predictions that are produced by some classifier can be used instead of the true labels.

A classifier, which is trained on additional labels as predictors, needs those additional labels as input variables. Since these labels are not available at prediction time, they need to be predicted first. When the true label information is used to augment the feature space in the training of a classifier, the assumption that the training data and the test data should be identically distributed is violated (Senge et al., 2013). If the true label information is used in the training data and the predicted label information is used in the test data, the training data is not representative for the test data. However, experiments (Montañés et al., 2014; Senge et al., 2013) show that none of these methods should be dismissed immediately. Note that we use the superscript "true" or "pred" to emphasize that a classifier $C_k^{\text{true}}$ or $C_k^{\text{pred}}$ used true labels or predicted labels as additional predictors during training, respectively.

Suppose there are $n = 6$ observations with $p = 3$ predictors and $m = 3$ labels. The true label $\mathbf{y}_3$ shall be used to augment the feature space of a binary classifier $C_1^{\text{true}}$ for label $\mathbf{y}_1$. $C_1^{\text{true}}$ is thus trained on all predictors and the true label $\mathbf{y}_3$. The binary classification task for label $\mathbf{y}_1$ is therefore:

$$\text{Train } C_1^{\text{true}} \text{ on } \boxed{\phantom{XXXX}} \text{ to predict } \mathbf{y}_1 \tag{2}$$

For an unlabeled observation $\mathbf{x}^{(l)}$, only the three predictor variables $x_1^{(l)}, \dots, x_3^{(l)}$ are available at prediction time. However, the classifier $C_1^{\text{true}}$ needs a 4-dimensional observation $\left(\mathbf{x}^{(l)}, y_3^{(l)}\right)$ as input. The input $y_3^{(l)}$ therefore needs to be predicted first. A new *level-1* classifier $C_3^{\text{lvl1}}$, which is trained on the set $D' = \cup_{i=1}^6 \left\{ \left(\mathbf{x}^{(i)}, y_3^{(i)}\right) \right\}$, will make those predictions for $y_3^{(l)}$. The training task is:

$$\text{Train } C_3^{\text{lvl1}} \text{ on } D' \triangleq \boxed{\phantom{XXXX}} \text{ to predict } \mathbf{y}_3 \tag{3}$$

Therefore, for a new observation $\mathbf{x}^{(l)}$, the predicted label $\hat{y}_3^{(l)}$ is obtained by using $C_3^{\text{lvl1}}$ on $\mathbf{x}^{(l)}$. The final prediction for $y_1^{(l)}$ is then obtained by using $C_1^{\text{true}}$ on $\left(\mathbf{x}^{(l)}, \hat{y}_3^{(l)}\right)$.

The alternative to (2) would be to use predicted labels $\hat{\mathbf{y}}_3$ instead of true labels $\mathbf{y}_3$. These labels should be produced by means of an out-of-sample prediction procedure (Senge et al., 2013). This can be done by an internal leave-one-out cross-validation procedure, which can of course be computationally intensive. Because of this, coarser resampling strategies can be used. As an example, an internal 2-fold cross-validation will be shown here. Again, let $D' = \cup_{i=1}^6 \left\{ \left(\mathbf{x}^{(i)}, y_3^{(i)}\right) \right\}$ be the set of all predictor variables with $\mathbf{y}_3$ as target variable. Using 2-fold cross-validation, the dataset $D'$ is split into two parts $D_1' = \cup_{i=1}^3 \left\{ \left(\mathbf{x}^{(i)}, y_3^{(i)}\right) \right\}$ and $D_2' = \cup_{i=4}^6 \left\{ \left(\mathbf{x}^{(i)}, y_3^{(i)}\right) \right\}$:

$$\boxed{\phantom{XXXX}} \tag{4}$$

Two classifiers $C_{D_1'}$ and $C_{D_2'}$ are then trained on $D_1'$ and $D_2'$, respectively, for the prediction of $\mathbf{y}_3$:

$$\text{Train } C_{D_1'} \text{ on } \boxed{D_1'} \text{ to predict } \mathbf{y}_3, \quad \text{Train } C_{D_2'} \text{ on } \boxed{D_2'} \text{ to predict } \mathbf{y}_3$$

Following the cross-validation paradigm, $D_1'$ is used as test set for the classifier $C_{D_2'}$, and $D_2'$ is used as a test set for $C_{D_1'}$:

$$C_{D_2'}: \boxed{\begin{matrix} x_1 & x_2 & x_3 \\ & D_1' & \end{matrix}} \mapsto \boxed{\begin{matrix} \hat{y}_3 \\ 1 \\ 0 \\ 0 \end{matrix}}, \quad C_{D_1'}: \boxed{\begin{matrix} x_1 & x_2 & x_3 \\ & D_2' & \end{matrix}} \mapsto \boxed{\begin{matrix} \hat{y}_3 \\ 0 \\ 0 \\ 1 \end{matrix}}$$

These predictions are merged for the final predicted label $\hat{\mathbf{y}}_3$, which is used to augment the feature space. The classifier $C_1^{\text{pred}}$ is then trained on that augmented feature space:

$$\text{Train } C_1^{\text{pred}} \text{ on } \boxed{\begin{matrix} x_1 & x_2 & x_3 & \hat{y}_3 & y_1 \\ & & & 1 & 0 \\ & & & 0 & 1 \\ & & & 0 & 1 \\ & & & 0 & 1 \\ & & & 1 & 1 \end{matrix}} \text{ to predict } \mathbf{y}_1 \tag{5}$$

The prediction phase is completely analogous to (3). It is worthwhile to mention that the level-1 classifier $C_3^{\text{lvl1}}$, which will be used to obtain predictions $\hat{\mathbf{y}}_3$ at prediction time, is trained on the whole set $D' = D_1' \cup D_2'$, following Simon (2007).

### Full vs. partial conditioning

Recall the set of all labels $Z = \{z_1, \ldots, z_m\}$. The prediction of a label $z_k$ can either be conditioned on all remaining labels $\{z_1, \ldots, z_{k-1}, z_{k+1}, \ldots, z_m\}$ (*full conditioning*) or just on a subset of labels (*partial conditioning*). The only method for partial conditioning, which is examined in this paper, is the chaining method. Here, labels $z_k$ are conditioned on all previous labels $\{z_1, \ldots, z_{k-1}\}$ for all $k = 1, \ldots, m$. This sequential structure is motivated by the product rule of probability (Montañés et al., 2014):

$$P\left(\mathbf{y}^{(i)} \middle| \mathbf{x}^{(i)}\right) = \prod_{k=1}^{m} P\left(y_k^{(i)} \middle| \mathbf{x}^{(i)}, y_1^{(i)}, \ldots, y_{k-1}^{(i)}\right) \tag{6}$$

Methods that make use of this chaining structure are e.g., *classifier chains* or *nested stacking* (these methods will be discussed further below).

To sum up the discussions above: there are four ways in modeling label dependencies through conditioning labels $z_k$ on other labels $z_\ell$, $k \neq \ell$. They can be distinguished by the subset of labels, which are used for conditioning, and by the use of predicted or real labels in the training step. In Table 1 we show the four methods, which implement these ideas and describe them consequently.

| | True labels | Pred. labels |
|---|---|---|
| Partial cond. | Classifier chains | Nested stacking |
| Full cond. | Dependent binary relevance | Stacking |

**Table 1:** Distinctions in modeling label dependence and models

### Classifier chains

The classifier chains (CC) method implements the idea of using partial conditioning together with the true label information. It was first introduced by Read et al. (2011). CC selects an order on the set of labels $\{z_1, \ldots, z_m\}$, which can be formally written as a bijective function (permutation):

$$\tau : \{1, \ldots, m\} \longrightarrow \{1, \ldots, m\} \tag{7}$$

Labels will be chained along this order $\tau$:

$$z_{\tau(1)} \rightarrow z_{\tau(2)} \rightarrow \ldots \rightarrow z_{\tau(m)} \tag{8}$$

However, for this paper the permutation shall be $\tau = id$ (only for simplicity reasons). The labels therefore follow the order $z_1 \rightarrow z_2 \rightarrow \ldots \rightarrow z_m$. In a similar fashion to the binary relevance (BR) method, CC trains $m$ binary classifiers $C_k$, which are responsible for predicting their corresponding label $z_k$, $k = 1, \ldots, m$. The classifiers $C_k$ are of the form

$$C_k : \mathcal{X} \times \{0,1\}^{k-1} \longrightarrow \{0,1\}, \tag{9}$$

where $\{0,1\}^0 := \varnothing$. For a classifier $C_k$ the feature space is augmented by the true label information of all previous labels $z_1, z_2, \ldots, z_{k-1}$. Hence, the training data of $C_k$ consists of all observations $\left( \left( \mathbf{x}^{(i)}, y_1^{(i)}, y_2^{(i)}, \ldots, y_{k-1}^{(i)} \right), y_k^{(i)} \right)$, $i = 1, \ldots, n$, with the target $y_k^{(i)}$. In the example from above, this would look like:

$$\text{Train } C_1 \text{ on} \qquad \text{Train } C_2 \text{ on} \qquad \text{Train } C_3 \text{ on} \qquad (10)$$

At prediction time, when an unlabeled observation $\mathbf{x}^{(l)}$ is labeled, a prediction $\left( \hat{y}_1^{(l)}, \ldots, \hat{y}_m^{(l)} \right)$ is obtained by successively predicting the labels along the chaining order:

$$
\begin{aligned}
\hat{y}_1^{(l)} &= C_1 \left( \mathbf{x}^{(l)} \right) \\
\hat{y}_2^{(l)} &= C_2 \left( \mathbf{x}^{(l)}, \hat{y}_1^{(l)} \right) \\
&\;\;\vdots \\
\hat{y}_m^{(l)} &= C_m \left( \mathbf{x}^{(l)}, \hat{y}_1^{(l)}, \hat{y}_2^{(l)}, \ldots, \hat{y}_{m-1}^{(l)} \right)
\end{aligned}
\qquad (11)
$$

The authors of Senge et al. (2013) summarize several factors, which have an impact on the performance of CC:

- *The length of the chain.* A high number $(k-1)$ of preceding classifiers in the chain comes with a high potential level of feature noise for the classifier $C_k$. One may assume that the probability of a mistake will increase with the level of feature noise in the input space. Then the probability of a mistake will be reinforced along the chain, due to the recursive structure of CC.

- *The order of the chain.* Some labels may be more difficult to predict than others. The order of a chain can therefore be important for the performance. It can be advantageous to put simple to predict labels in the beginning and harder to predict labels more towards the end of the chain. Some heuristics for finding an optimal chain ordering have been proposed in da Silva et al. (2014); Read et al. (2013). Alternatively Read et al. (2011) developed an ensemble of classifier chains, which builds many randomly ordered CC-classifiers and put them on a voting scheme for a prediction. However, these methods are not subject of this article.

- *The dependency among labels.* For an improvement of performance through chaining, there should be a dependence among labels, CC cannot gain in case of label independence. However, CC is also only likely to lose if the binary classifiers $C_k$ cannot ignore the added features $\mathbf{y}_1, \ldots, \mathbf{y}_{k-1}$.

**Nested stacking**

The nested stacking method (NST), first proposed in Senge et al. (2013), implements the idea of using partial conditioning together with predicted label information. NST mimics the chaining structure of CC, but does not use real label information during training. Like in CC the chaining order shall be $\tau = id$, again for simplicity reasons. CC uses real label information $\mathbf{y}_k$ during training and predicted labels $\hat{\mathbf{y}}_k$ at prediction time. However, unless the binary classifiers are perfect, it is likely that $\mathbf{y}_k$ and $\hat{\mathbf{y}}_k$ do not follow the same distribution. Hence, the key assumption of supervised learning, namely that the training data should be representative for the test data, is violated by CC. Nested stacking tries to overcome this issue by using predicted labels $\hat{\mathbf{y}}_k$ instead of true labels $\mathbf{y}_k$.

NST trains $m$ binary classifiers $C_k$ on $D_k := \cup_{i=1}^n \left\{ \left( \left( \mathbf{x}^{(i)}, \hat{y}_1^{(i)}, \ldots, \hat{y}_{k-1}^{(i)} \right), y_k^{(i)} \right) \right\}$, for all $k = 1, \ldots, m$. The predicted labels should be obtained by an internal out-of-sample method (Senge et al., 2013). How these predictions are obtained was already explained in the **True vs. Predicted Label Information** chapter. The prediction phase is completely analogous to (11).

The training procedure is visualized in the following with 2-fold cross-validation as an internal out-of-sample method:

$$\text{Train } C_1 \text{ on} \qquad \text{Use 2-fold CV on} \qquad \text{to obtain} \qquad (12)$$

$$\text{Train } C_2 \text{ on } \qquad \text{Use 2-fold CV on } \qquad \text{to obtain} \qquad (13)$$

$$\text{Train } C_3 \text{ on } \qquad (14)$$

The factors which impact the performance of CC (i.e., length and order of the chain, and the dependency among labels), also impact NST, since NST mimics the chaining method of CC.

### Dependent binary relevance

The dependent binary relevance method (DBR) implements the idea of using full conditioning together with the true label information. DBR is built on two main hypotheses (Montañés et al., 2014):

(i) Taking conditional label dependencies into account is important for performing well in multil-abel classification tasks.

(ii) Modeling and learning these label dependencies in an overcomplete way (take all other labels for modeling) may further improve model performance.

The first assumption is the main prerequisite for research in multilabel classification. It has been shown theoretically that simple binary relevance classifiers cannot achieve optimal performance for specific multilabel loss functions (Montañés et al., 2014). The second assumption, however, is harder to justify theoretically. Nonetheless, the practical usefulness of learning in an overcomplete way has been shown in many branches of (classical) single-label classification (e.g., ensemble methods (Dietterich, 2000)).

Formally, DBR trains $m$ binary classifiers $C_1, \ldots, C_m$ (as many classifiers as labels) on the corresponding training data

$$D_k = \cup_{i=1}^n \left\{ \left( \left( \mathbf{x}^{(i)}, y_1^{(i)}, \ldots, y_{k-1}^{(i)}, y_{k+1}^{(i)}, \ldots, y_m^{(i)} \right), y_k^{(i)} \right) \right\}, \qquad (15)$$

$k = 1, \ldots, m$. Thus, each classifier $C_k$ is of the form

$$C_k : \mathcal{X} \times \{0,1\}^{m-1} \longrightarrow \{0,1\}.$$

Hence, for each classifier $C_k$ the true label information of all labels except $\mathbf{y}_k$ is used as augmented features. Again, here is a visualization with the example from above:

$$\text{Train } C_1 \text{ on } \qquad \text{Train } C_2 \text{ on } \qquad \text{Train } C_3 \text{ on }$$

$$(16)$$

To make these classifiers applicable, when an unlabeled instance $\mathbf{x}^{(l)}$ needs to be labeled, the help of other multilabel classifiers is needed to produce predicted labels $\hat{y}_1^{(l)}, \ldots, \hat{y}_m^{(l)}$ as additional features. The classifiers, which produce predicted labels as additional features, are called *base learners* (Montañés et al., 2014). Theoretically any multilabel classifier can be used as base learner. However, in this paper, the analysis is focused on BR as base learner only. The prediction of an unlabeled instance $\mathbf{x}^{(l)}$ formally works as follows:

(i) First level: Produce predicted labels by using the BR base learner:

$$C_{BR} \left( \mathbf{x}^{(l)} \right) = \left( \hat{y}_1^{(l)}, \ldots, \hat{y}_m^{(l)} \right)$$

(ii) Second level, which is also called meta level (Montañés et al., 2014): Produce final prediction $\hat{\mathbf{y}}_k = \left( \hat{y}_1^{(l)}, \ldots, \hat{y}_m^{(l)} \right)$ by applying DBR classifiers $C_1, \ldots, C_m$:

$$C_1 \left( \mathbf{x}^{(l)}, \hat{y}_2^{(l)}, \ldots, \hat{y}_m^{(l)} \right) = \hat{y}_1^{(l)}$$
$$C_2 \left( \mathbf{x}^{(l)}, \hat{y}_1^{(l)}, \hat{y}_3^{(l)}, \ldots, \hat{y}_m^{(l)} \right) = \hat{y}_2^{(l)}$$
$$\vdots$$
$$C_m \left( \mathbf{x}^{(l)}, \hat{y}_1^{(l)}, \ldots, \hat{y}_{m-1}^{(l)} \right) = \hat{y}_m^{(l)}$$

### Stacking

Stacking (STA) implements the last variant of Table 1, namely the use of full conditioning together with predicted label information. Stacking is short for *stacked generalization* (Wolpert, 1992) and was first proposed in the multilabel context by Godbole and Sarawagi (2004). Like in classical stacking, for each label it takes predictions of several other learners that were trained in a first step to get a new learner to make predictions for the corresponding label. Both hypotheses on which DBR is built on also apply to STA, of course.

STA trains $m$ classifiers $C_1, \ldots, C_m$ on the corresponding training data

$$D_k = \cup_{i=1}^n \left\{ \left( \left( \mathbf{x}^{(i)}, \hat{y}_1^{(i)}, \ldots, \hat{y}_m^{(i)} \right), y_k^{(i)} \right) \right\}, k = 1, \ldots, m. \tag{17}$$

The classifiers $C_k$, $k = 1, \ldots, m$, are therefore of the following form:

$$C_k : \mathcal{X} \times \{0,1\}^m \longrightarrow \{0,1\}$$

Like in NST, the predicted labels should be obtained by an internal out-of-sample method (Sill et al., 2009). STA can be seen as the alternative to DBR using predicted labels (like NST is for CC). However, the classifiers $C_k$, $k = 1, \ldots, m$, are trained on all predicted labels $\hat{\mathbf{y}}_1, \ldots, \hat{\mathbf{y}}_m$ for the STA approach (in DBR the label $\mathbf{y}_k$ is left out of the augmented training set).

The training procedure is outlined in the following:

For i=1,2,3 use 2-fold CV on  to obtain  $\tag{18}$

For i=1,2,3 train $C_k$ on  $\tag{19}$

Like in DBR, STA depends on a BR base learner, to produce predicted labels as additional features. Again, the use of BR as a base learner is not mandatory, but it is the proposed method in Godbole and Sarawagi (2004).

The prediction of an unlabeled instance $\mathbf{x}^{(l)}$ works almost identically to the DBR case and is illustrated here:

(i) First level. Produce predicted labels by using the BR base learner:

$$C_{BR} \left( \mathbf{x}^{(l)} \right) = \left( \hat{y}_1^{(l)}, \ldots, \hat{y}_m^{(l)} \right)$$

(ii) Meta level. Apply STA classifiers $C_1, \ldots, C_m$:

$$C_1 \left( \mathbf{x}^{(l)}, \hat{y}_1^{(l)}, \ldots, \hat{y}_m^{(l)} \right) = \hat{y}_1^{(l)}$$
$$\vdots$$
$$C_m \left( \mathbf{x}^{(l)}, \hat{y}_1^{(l)}, \ldots, \hat{y}_m^{(l)} \right) = \hat{y}_m^{(l)}$$

**Multilabel performance measures**

Analogously to multiclass classification there exist multilabel classification performance measures. Six multilabel performance measures can be evaluated in **mlr**. These are: *Subset 0/1 loss*, *hamming loss*, *accuracy*, *precision*, *recall* and $F_1$-*index*. Multilabel performance measures are defined on a per instance basis. The performance on a test set is the average over all instances.

Let $D_{\text{test}} = \left\{ \left( \mathbf{x}^{(1)}, \mathbf{y}^{(1)} \right), \ldots, \left( \mathbf{x}^{(n)}, \mathbf{y}^{(n)} \right) \right\}$ be a test set with $\mathbf{y}^{(i)} = \left( y_1^{(i)}, \ldots, y_m^{(i)} \right) \in \{0,1\}^m$ for all $i = 1, \ldots, n$. Performance measures quantify how good a classifier $C$ predicts the labels $z_1, \ldots, z_n$.

(i) The subset 0/1 loss is used to see if the predicted labels $C(\mathbf{x}^{(i)}) = \left( \hat{y}_1^{(i)}, \ldots, \hat{y}_m^{(i)} \right)$ are equal to the actual labels $\left( y_1^{(i)}, \ldots, y_m^{(i)} \right)$:

$$\text{subset}_{0/1} \left( C, \left( \mathbf{x}^{(i)}, \mathbf{y}^{(i)} \right) \right) = \mathbb{1}_{\left( \mathbf{y}^{(i)} \neq C(\mathbf{x}^{(i)}) \right)} := \begin{cases} 1 & \text{if } \mathbf{y}^{(i)} \neq C \left( \mathbf{x}^{(i)} \right) \\ 0 & \text{if } \mathbf{y}^{(i)} = C \left( \mathbf{x}^{(i)} \right) \end{cases}$$

The subset 0/1 loss of a classifier $C$ on a test set $D_{\text{test}}$ thus becomes:

$$\text{subset}_{0/1} \left( C, D_{\text{test}} \right) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{\mathbf{y}^{(i)} \neq C(\mathbf{x}^{(i)})}$$

The subset 0/1 loss can be interpreted as the analogon of the mean misclassification error in multiclass classifications. In the multilabel case it is a rather drastic measure because it treats a mistake on a single label as a complete failure (Senge et al., 2013).

(ii) The hamming loss also takes into account observations where only some labels have been predicted correctly. It corresponds to the proportion of labels whose relevance is incorrectly predicted. For an instance $\left( \mathbf{x}^{(i)}, \mathbf{y}^{(i)} \right) = \left( \mathbf{x}^{(i)}, \left( y_1^{(i)}, \ldots, y_m^{(i)} \right) \right)$ and a classifier $C \left( \mathbf{x}^{(i)} \right) = \left( \hat{y}_1^{(i)}, \ldots, \hat{y}_m^{(i)} \right)$ this is defined as:

$$\text{HammingLoss} \left( C, \left( \mathbf{x}^{(i)}, \mathbf{y}^{(i)} \right) \right) = \frac{1}{m} \sum_{k=1}^{m} \mathbb{1}_{\left( y_k^{(i)} \neq \hat{y}_k^{(i)} \right)}$$

If one label is predicted incorrectly, this accounts for an error of $\frac{1}{m}$. For a test set $D_{\text{test}}$ the hamming loss becomes:

$$\text{HammingLoss}(C, D_{\text{test}}) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{m} \sum_{k=1}^{m} \mathbb{1}_{\left( y_k^{(i)} \neq \hat{y}_k^{(i)} \right)}$$

The following measures are scores instead of loss function like the two previous ones.

(iii) The accuracy, also called Jaccard-Index, for a test set $D_{\text{test}}$ is defined as:

$$\text{accuracy}(C, D_{\text{test}}) = \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{k=1}^{m} \mathbb{1}_{\left( y_k^{(i)} = 1 \text{ and } \hat{y}_k^{(i)} = 1 \right)}}{\sum_{k=1}^{m} \mathbb{1}_{\left( y_k^{(i)} = 1 \text{ or } \hat{y}_k^{(i)} = 1 \right)}}$$

(iv) The precision for a test set $D_{\text{test}}$ is defined as:

$$\text{precision}(C, D_{\text{test}}) = \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{k=1}^{m} \mathbb{1}_{\left( y_k^{(i)} = 1 \text{ and } \hat{y}_k^{(i)} = 1 \right)}}{\sum_{k=1}^{m} \mathbb{1}_{\left( \hat{y}_k^{(i)} = 1 \right)}}$$

(v) The recall for a test set $D_{\text{test}}$ is defined as:

$$\text{recall}(C, D_{\text{test}}) = \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{k=1}^{m} \mathbb{1}_{\left( y_k^{(i)} = 1 \text{ and } \hat{y}_k^{(i)} = 1 \right)}}{\sum_{k=1}^{m} \mathbb{1}_{\left( y_k^{(i)} = 1 \right)}}$$

(vi) For a test set $D_{\text{test}}$ the $F_1$-index is defined as follows:

$$F_1(C, D_{\text{test}}) = \frac{1}{n} \sum_{i=1}^{n} \frac{2 \sum_{k=1}^{m} \mathbb{1}_{\left( y_k^{(i)} = 1 \text{ and } \hat{y}_k^{(i)} = 1 \right)}}{\sum_{k=1}^{m} \left( \mathbb{1}_{\left( y_k^{(i)} = 1 \right)} + \mathbb{1}_{\left( \hat{y}_k^{(i)} = 1 \right)} \right)}$$

The F$_1$-index is the harmonic mean of recall and precision on a per instance basis.

All these measures lie between 0 and 1. In the case of the subset 0/1 loss and the hamming loss the values should be low, in all other cases the scores should be high. Demonstrative definitions with sets instead of vectors can be seen in Charte and Charte (2015).

## Implementation

In this section, we briefly describe how to perform multilabel classifications in **mlr**. We provide small code examples for better illustration. A short tutorial is also available at http://mlr-org. github.io/mlr-tutorial/release/html/multilabel/index.html. The first step is to transform the multilabel dataset into a 'data.frame' in R. The columns must consist of vectors of features and one logical vector for each label that indicates if the label is present for the observation or not. To fit a multilabel classification algorithm in **mlr**, a multilabel task has to be created, where a vector of targets corresponding to the column names of the labels has to be specified. This task is an S3 object that contains the data, the target labels and further descriptive information. In the following example, the yeast data frame is extracted from the yeast.task, which is provided by the **mlr** package. Then the 14 label names of the targets are extracted and the multilabel task is created.

```
yeast = getTaskData(yeast.task)
labels = colnames(yeast)[1:14]
yeast.task = makeMultilabelTask(id = "multi", data = yeast, target = labels)
```

### Problem transformation methods

To generate a problem transformation method learner, a binary classification base learner has to be created with 'makeLearner'. A list of available learners for classifications in **mlr** can be seen at http://mlr-org.github.io/mlr-tutorial/release/html/integrated_learners/. Specific hyperparameter settings of the base learner can be set in this step through the 'par.vals' argument in 'makeLearner'. Afterwards, a learner for any problem transformation method can be created by applying the function 'makeMultilabel[...]Wrapper', where [...] has to be substituted by the desired problem transformation method. In the following example, two multilabel variants with rpart as base learner are created. The base learner is configured to output probabilities instead of discrete labels during prediction.

```
lrn = makeLearner("classif.rpart", predict.type = "prob")
multilabel.lrn1 = makeMultilabelBinaryRelevanceWrapper(lrn)
multilabel.lrn2 = makeMultilabelNestedStackingWrapper(lrn)
```

### Algorithm adaptation methods

Algorithm adaptation method learners can be created directly with 'makeLearner'. The names of the specific learner can be looked up at http://mlr-org.github.io/mlr-tutorial/release/html/integrated_learners/ in the multilabel section.

```
multilabel.lrn3 = makeLearner("multilabel.rFerns")
multilabel.lrn4 = makeLearner("multilabel.randomForestSRC")
```

### Train, predict and evaluate

Training and predicting on data can be done as usual in **mlr** with the functions 'train' and 'predict'. Learner and task have to be specified in 'train'; trained model and task or new data have to be specified in 'predict'.

```
mod = train(multilabel.lrn1, yeast.task, subset = 1:1500)
pred = predict(mod, task = yeast.task, subset = 1501:1600)
```

The performance of the prediction can be assessed via the function 'performance'. Measures are represented as S3 objects and multiple objects can be passed in as a list. The default measure for multilabel classification is the hamming loss (*multilabel.hamloss*). All available measures for multilabel classification can be shown by 'listMeasures' or looked up in the appendix of the tutorial page[1] (http://mlr-org.github.io/mlr-tutorial/release/html/measures/index.html).

---

[1]In the **mlr** package *precision* is named *positive predictive value* and *recall* is named *true positive rate*.

```
performance(pred, measures = list(multilabel.hamloss, timepredict))
multilabel.hamloss      timepredict
0.230            0.174
listMeasures("multilabel")
# [1] "multilabel.ppv" "timepredict"       "multilabel.hamloss" "multilabel.f1"
# [5] "featperc"       "multilabel.subset01" "timeboth"          "timetrain"
# [9] "multilabel.tpr" "multilabel.acc"
```

### Resampling

To properly evaluate the model, a resampling strategy, for example k-fold cross-validation, should be applied. This can be done in **mlr** by using the function 'resample'. First, a description of the subsequent resampling strategy, in this case three-fold cross-validation, is defined with 'makeResampleDesc'. The resample is executed by a call to the 'resample' function. The hamming loss is calculated for the binary relevance method.

```
rdesc = makeResampleDesc(method = "CV", stratify = FALSE, iters = 3)
r = resample(learner = multilabel.lrn1, task = yeast.task, resampling = rdesc,
measures = list(multilabel.hamloss), show.info = FALSE)
r
# Resample Result
# Task: multi
# Learner: multilabel.classif.rpart
# multilabel.hamloss.aggr: 0.23
# multilabel.hamloss.mean: 0.23
# multilabel.hamloss.sd: 0.00
# Runtime: 6.36688
```

### Binary performance

To calculate a binary performance measure like, e.g., the accuracy, the mean misclassification error (mmce) or the AUC for each individual label, the function 'getMultilabelBinaryPerformances' can be used. This function can be applied to a single multilabel test set prediction and also on a resampled multilabel prediction. To calculate the AUC, predicted probabilities are needed. These can be obtained by setting the argument 'predict.type = "prob"' in the 'makeLearner' function.

```
head(getMultilabelBinaryPerformances(r$pred, measures = list(acc, mmce, auc)))
#        acc.test.mean mmce.test.mean auc.test.mean
# label1    0.7389326     0.2610674     0.6801810
# label2    0.5908151     0.4091849     0.5935160
# label3    0.6512205     0.3487795     0.6631469
# label4    0.6921804     0.3078196     0.6965552
# label5    0.7517584     0.2482416     0.6748458
# label6    0.7343815     0.2656185     0.6054968
```

### Parallelization

In the case of a high number of labels and larger datasets, parallelization in the training and prediction process of the multilabel methods can reduce computation time. This can be achieved by using the package parallelMap in mlr (see also the tutorial section of parallelization: http://mlr-org.github.io/mlr-tutorial/release/html/multilabel/index.html). Currently, only the binary relevance method is parallelizable, the classifier for each label is trained in parallel, as they are independent of each other. The other problem transformation methods will also be parallelizable (as far as possible) soon.

```
library(parallelMap)
parallelStartSocket(2)
lrn = makeMultilabelBinaryRelevanceWrapper("classif.rpart")
mod = train(lrn, yeast.task)
pred = predict(mod, yeast.task)
```

## Benchmark experiment

In a similar fashion to Wang et al. (2014), we performed a benchmark experiment on several datasets in order to compare the performances of the different multilabel algorithms.

**Datasets**: In Table 2 we provide an overview of the used datasets. We retrieved most datasets from the Mulan Java library for multilabel learning[2] as well as from other benchmark experiments of multilabel classification methods. See Table 2 for article references. We uploaded all datasets to the open data platform OpenML (Casalicchio et al., 2017; Vanschoren et al., 2013), so they now can be downloaded directly from there. In some of the used datasets, sparse labels had to be removed in order to avoid problems during cross-validation. Several binary classification methods have difficulties when labels are sparse, i.e., a strongly imbalanced binary target class can lead to constant predictions for that target. That can sometimes lead to direct problems in the base learners (when training on constant class labels is simply not allowed) or, e.g., in classifier chains, when the base learner cannot handle constant features. Furthermore, one can reasonably argue that not much is to be learned for such a label. Hence, labels that appeared in less than 2% of the observations were removed. We computed *cardinality* scores (based on the remaining labels) indicating the mean number of labels assigned to each case in the respective dataset. The following description of the datasets refers to the final versions after removal of sparse labels.

- The first dataset (*birds*) consists of 645 audio recordings of 15 different vocalizing bird species (Briggs et al., 2013). Each sound can be assigned to various bird species.

- Another audio dataset (*emotions*) consists of 593 musical files with 6 clustered emotional labels (Trohidis et al., 2008) and 72 predictors. Each song can be labeled with one or more of the labels {*amazed-surprised, happy-pleased, relaxing-calm, quiet-still, sad-lonely, angry-fearful*}.

- The *genbase* dataset contains protein sequences that can be assigned to several classes of protein families (Diplaris et al., 2005). The entire dataset contains 1186 binary predictors.

- The *langLog*[3] dataset includes 998 textual predictors and was originally compiled in the doctorial thesis of Read (2010). It consists of 1460 text samples that can be assigned to one or more topics such as *language, politics, errors, humor* and *computational linguistics*.

- The UC Berkeley *enron*[4] dataset represents a subset of the original *enron*[5] dataset and consists of 1702 cases of emails with 24 labels and 1001 predictor variables (Klimt and Yang, 2004).

- A subset of the *reuters*[6] dataset includes 2000 observations for text classification (Zhang and Zhou, 2008).

- The *image*[7] benchmark dataset consists of 2000 natural scene images. Zhou and ling Zhang (2007) extracted 135 features for each image and made it publicly available as *processed* image dataset. Each observation can be associated with different label sets, where all possible labels are {*desert, mountains, sea, sunset, trees*}. About 22% of the images belong to more than one class. However, images belonging to three classes or more are very rare.

- The *scene* dataset is an image classification task where labels like *Beach, Mountain, Field, Urban* are assigned to each image (Boutell et al., 2004).

- The *yeast* dataset (Elisseeff and Weston, 2002) consists of micro-array expression data, as well as phylogenetic profiles of yeast, and includes 2417 genes and 103 predictors. In total, 14 different labels can be assigned to a gene, but only 13 labels were used due to label sparsity.

- Another dataset for text-classification is the *slashdot*[8] dataset (Read et al., 2011). It consists of article titles and partial blurbs. Blurbs can be assigned to several categories (e.g., *Science, News, Games*) based on word predictors.

**Algorithms**: We used all multilabel classification methods currently implemented in **mlr**: binary relevance (BR), classifier chains (CC), nested stacking (NST), dependent binary relevance (DBR) and stacking (STA) as well as algorithm adaption methods of the **rFerns** (RFERN) and **randomForestSRC** (RFSRC) packages. For DBR and STA the first level and meta level classifiers were equal. For CC and NST we chose random chain orders for each resample iteration.

---

[2] http://mulan.sourceforge.net/datasets-mlc.html
[3] http://languagelog.ldc.upenn.edu/nll/
[4] http://bailando.sims.berkeley.edu/enron_email.html
[5] http://www.cs.cmu.edu/~enron/
[6] http://lamda.nju.edu.cn/data_MIMLtext.ashx
[7] http://lamda.nju.edu.cn/data_MIMLimage.ashx
[8] http://slashdot.org

| Dataset | Reference | # Inst. | # Pred. | # Labels | Cardinality |
|---------|-----------|--------:|--------:|---------:|------------:|
| birds* | Briggs et al. (2013) | 645 | 260 | 15 | 0.96 |
| emotions | Trohidis et al. (2008) | 593 | 72 | 6 | 1.87 |
| genbase* | Diplaris et al. (2005) | 662 | 112 | 16 | 1.20 |
| langLog* | Read (2010) | 1460 | 998 | 18 | 0.85 |
| enron* | Klimt and Yang (2004) | 1702 | 1001 | 24 | 3.12 |
| reuters | Zhang and Zhou (2008) | 2000 | 243 | 7 | 1.15 |
| image | Zhou and ling Zhang (2007) | 2000 | 135 | 5 | 1.24 |
| scene | Boutell et al. (2004) | 2407 | 294 | 6 | 1.07 |
| yeast* | Elisseeff and Weston (2002) | 2417 | 103 | 13 | 4.22 |
| slashdot* | Read et al. (2011) | 3782 | 1079 | 14 | 1.13 |

**Table 2:** Used benchmark datasets including number of instances, number of predictor, number of label and label cardinality. Datasets with an asterisk differ from the original dataset as sparse labels have been removed. The genbase dataset contained many constant factor variables, which were automatically removed by mlr.

**Base Learners**: We employed two different binary classification base learner for each problem transformation algorithm: random forest (rf) of the **randomForest** package (Liaw and Wiener, 2002) with ntree = 100 and adaboost (ad) from the **ada** package (Culp et al., 2012), each with standard hyperparameter settings.

**Performance Measures:** We used the six previously proposed performance measures. Furthermore, we calculated the reported values by means of a 10-fold cross-validation.

**Code:** For reproducibility, the complete code and results can be downloaded from Probst (2017). The R package **batchtools** (Bischl et al., 2015) was used for parallelization.

The results for hamming loss and $F_1$-index are illustrated in Figure 1. Tables 3 and 4 contain performance values with the best performing algorithms highlighted in blue. For all remaining measures one may refer to the Appendix. We did not perform any threshold tuning that would potentially improve some of the performance of the methods.

The results of the problem transformation methods in this benchmark experiment concur with the general conclusions and results in Montañés et al. (2014). The authors ran a similar benchmark study with penalized logistic regression as base learner. They concluded that, on average, DBR performs well in $F_1$ and accuracy. Also, CC outperform the other methods regarding the subset 0/1 loss most of the time. For the hamming loss measure they got mixed results, with no clear winner concordant to our benchmark results. As base learner, on average, adaboost performs better than random forest in our benchmark study.

Considering the measure $F_1$, the problem transformation methods DBR, CC, STA and NST outperform RFERN and RFSRC on most of the datasets and also almost always perform better than BR, which does not consider dependencies among the labels. RFSRC and RFERN only perform well on either precision or recall, but in order to be considered as good classifiers they should perform well on both. The generally poor performances of RFERN can be explained by the working mechanism of the algorithm which randomly chooses variables and split points at each split of a fern. Hence, it cannot deal with too many features that are useless for the prediction of the target labels.

## Summary

In this paper, we describe the implementation of multilabel classification algorithms in the R package **mlr**. The problem transformation methods binary relevance, classifier chains, nested stacking, dependent binary relevance and stacking are implemented and can be used with any base learner that is accessible in **mlr**. Moreover, there is access to the multilabel classification versions of **randomForestSRC** and **RFerns**. We compare all of these methods in a benchmark experiment with several datasets and different implemented multilabel performance measures. The dependent binary relevance method performs well regarding the measures $F_1$ and *accuracy*. Classifier chains outperform the other methods in terms of the subset 0/1 loss most of the time. Parallelization is available for the binary relevance method and will be available soon for the other problem transformation methods. Algorithm adaptation methods and problem transformation methods that are currently not available can be incorporated in the current **mlr** framework easily. In our benchmark experiment we had to remove labels which occured too sparsely, because some algorithms crashed due to one class problems, which appeared during cross-validation. A solution to this problem and an implementation into the **mlr** framework is of great interest.
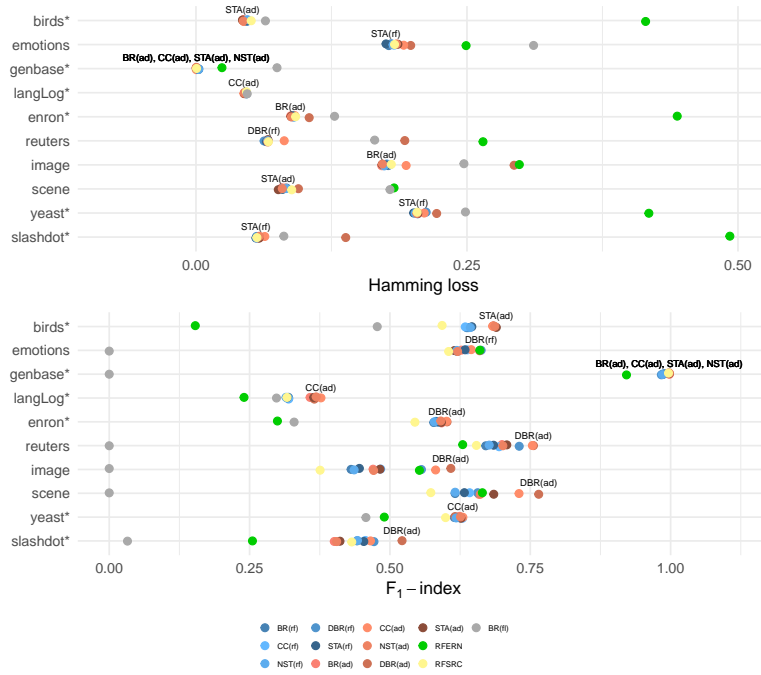
**Figure 1:** Results for hamming loss and $F_1$-index. The best performing algorithms are highlighted on the plot.

|           | BR(rf) | CC(rf) | NST(rf) | DBR(rf) | STA(rf) | BR(ad) | CC(ad) | NST(ad) | DBR(ad) | STA(ad) | RFERN  | RFSRC  | BR(fl) |
|-----------|--------|--------|---------|---------|---------|--------|--------|---------|---------|---------|--------|--------|--------|
| birds*    | 0.0477 | 0.0479 | 0.0475  | 0.0472  | 0.0468  | 0.0442 | 0.0441 | 0.0436  | 0.0431  | 0.0429  | 0.4148 | 0.0510 | 0.0641 |
| emotions  | 0.1779 | 0.1832 | 0.1818  | 0.1801  | 0.1753  | 0.181  | 0.1916 | 0.1849  | 0.1981  | 0.1863  | 0.2492 | 0.1832 | 0.3114 |
| genbase*  | 0.0021 | 0.0023 | 0.0025  | 0.0027  | 0.0023  | 0.0003 | 0.0003 | 0.0003  | 0.0004  | 0.0003  | 0.0240 | 0.0006 | 0.0748 |
| langLog*  | 0.0464 | 0.0465 | 0.0467  | 0.0464  | 0.0466  | 0.0451 | 0.0442 | 0.0446  | 0.0447  | 0.0448  | 0.4440 | 0.0466 | 0.0473 |
| enron*    | 0.0903 | 0.0904 | 0.0902  | 0.0909  | 0.0891  | 0.0874 | 0.0913 | 0.0881  | 0.1045  | 0.0877  | 0.2648 | 0.0668 | 0.1649 |
| reuters   | 0.0663 | 0.0654 | 0.0661  | 0.0629  | 0.065   | 0.0666 | 0.0814 | 0.0664  | 0.1926  | 0.0664  | 0.2983 | 0.1802 | 0.2472 |
| image     | 0.1774 | 0.1791 | 0.1737  | 0.1761  | 0.1754  | 0.1714 | 0.1939 | 0.1721  | 0.2935  | 0.1717  | 0.2983 | 0.1802 | 0.2472 |
| scene     | 0.0836 | 0.0809 | 0.0832  | 0.0796  | 0.0799  | 0.0791 | 0.0821 | 0.0796  | 0.0945  | 0.076   | 0.1827 | 0.0884 | 0.1790 |
| yeast*    | 0.2038 | 0.2044 | 0.2023  | 0.2123  | 0.2008  | 0.2048 | 0.2105 | 0.2038  | 0.2221  | 0.2046  | 0.4178 | 0.2040 | 0.2486 |
| slashdot* | 0.0558 | 0.0560 | 0.0559  | 0.0559  | 0.0554  | 0.059  | 0.0635 | 0.0586  | 0.1382  | 0.0582  | 0.4925 | 0.0562 | 0.0811 |

**Table 3:** Hamming loss

|           | BR(rf) | CC(rf) | NST(rf) | DBR(rf) | STA(rf) | BR(ad) | CC(ad) | NST(ad) | DBR(ad) | STA(ad) | RFERN  | RFSRC  | BR(fl) |
|-----------|--------|--------|---------|---------|---------|--------|--------|---------|---------|---------|--------|--------|--------|
| birds*    | 0.6369 | 0.6342 | 0.6433  | 0.64    | 0.6459  | 0.6835 | 0.683  | 0.6867  | 0.6846  | 0.6895  | 0.1533 | 0.5929 | 0.4774 |
| emotions  | 0.6199 | 0.6380 | 0.6192  | 0.6625  | 0.6337  | 0.6274 | 0.6449 | 0.6206  | 0.6598  | 0.615   | 0.6603 | 0.6046 | 0.0000 |
| genbase*  | 0.9885 | 0.9861 | 0.9855  | 0.9835  | 0.9861  | 0.9977 | 0.9977 | 0.9977  | 0.9962  | 0.9977  | 0.9214 | 0.9962 | 0.0000 |
| langLog*  | 0.3192 | 0.3194 | 0.3148  | 0.3199  | 0.3167  | 0.3578 | 0.3772 | 0.3686  | 0.3653  | 0.3643  | 0.2401 | 0.3167 | 0.2979 |
| enron*    | 0.5781 | 0.5822 | 0.5791  | 0.5866  | 0.5826  | 0.592  | 0.6009 | 0.5906  | 0.6017  | 0.5917  | 0.2996 | 0.5446 | 0.3293 |
| reuters   | 0.6708 | 0.6944 | 0.6769  | 0.7303  | 0.6846  | 0.6997 | 0.7537 | 0.7012  | 0.7556  | 0.7082  | 0.6296 | 0.6541 | 0.0000 |
| image     | 0.4308 | 0.4835 | 0.4362  | 0.5561  | 0.4456  | 0.47   | 0.5814 | 0.4709  | 0.6085  | 0.4824  | 0.5525 | 0.3757 | 0.0000 |
| scene     | 0.6161 | 0.6420 | 0.6161  | 0.6563  | 0.6326  | 0.6585 | 0.73   | 0.661   | 0.765   | 0.685   | 0.6647 | 0.5729 | 0.0000 |
| yeast*    | 0.6148 | 0.6294 | 0.6180  | 0.6195  | 0.6244  | 0.6238 | 0.63   | 0.6257  | 0.616   | 0.6266  | 0.4900 | 0.5991 | 0.4572 |
| slashdot* | 0.4415 | 0.4562 | 0.4422  | 0.4716  | 0.4535  | 0.4009 | 0.4654 | 0.4052  | 0.5216  | 0.411   | 0.2551 | 0.4320 | 0.0325 |

**Table 4:** $F_1$-index

## Bibliography

B. Bischl, M. Lang, O. Mersmann, J. Rahnenführer, and C. Weihs. BatchJobs and BatchExperiments: Abstraction mechanisms for using R in batch environments. *Journal of Statistical Software*, 64(11): 1–25, 2015. URL https://doi.org/10.18637/jss.v064.i11. [p363]

B. Bischl, M. Lang, L. Kotthoff, J. Schiffner, J. Richter, E. Studerus, G. Casalicchio, and Z. M. Jones. Mlr: Machine learning in R. *Journal of Machine Learning Research*, 17(170):1–5, 2016. [p352]

M. R. Boutell, J. Luo, X. Shen, and C. M. Brown. Learning multi-label scene classification. *Pattern Recognition*, 37(9):1757–1771, 2004. URL https://doi.org/10.1016/j.patcog.2004.03.009. [p352, 362, 363]

F. Briggs, H. Yonghong, R. Raich, and others. New methods for acoustic classification of multiple simultaneous bird species in a noisy environment. In *IEEE International Workshop on Machine Learning for Signal Processing*, pages 1–8, 2013. URL https://doi.org/10.1109/mlsp.2013.6661934. [p352, 362, 363]

G. Casalicchio, J. Bossek, M. Lang, D. Kirchhoff, P. Kerschke, B. Hofner, H. Seibold, J. Vanschoren, and B. Bischl. OpenML: An R package to connect to the networked machine learning platform OpenML. *ArXiv e-prints*, 2017. [p362]

F. Charte and D. Charte. Working with multilabel datasets in R: The mldr package. *The R Journal*, 7(2): 149–162, 2015. [p352, 360]

M. Culp, K. Johnson, and G. Michailidis. *ada: An R Package for Stochastic Boosting*, 2012. URL https://cran.r-project.org/package=ada. [p363]

P. N. da Silva, E. C. Gonçalves, A. Plastino, and A. A. Freitas. Distinct chains for different instances: An effective strategy for multi-label classifier chains. In *European Conference, ECML PKDD 2014*, pages 453–468, 2014. URL https://doi.org/10.1007/978-3-662-44851-9_29. [p356]

T. G. Dietterich. Ensemble methods in machine learning. *Lecture Notes in Computer Science*, 1857:1–15, 2000. URL https://doi.org/10.1007/3-540-45014-9_1. [p357]

S. Diplaris, G. Tsoumakas, P. A. Mitkas, and I. Vlahavas. Protein classification with multiple algorithms. In *Advances in Informatics*, pages 448–456. Springer-Verlag, 2005. URL https://doi.org/10.1007/11573036_42. [p362, 363]

A. Elisseeff and J. Weston. A kernel method for multi-labelled classification. In T. G. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems 14*, pages 681–687. MIT Press, 2002. [p352, 362, 363]

S. Godbole and S. Sarawagi. Discriminative methods for multi-labeled classification. In *Advances in Knowledge Discovery and Data*, volume LNCS3056, pages 22–30, 2004. URL https://doi.org/10.1007/978-3-540-24775-3_5. [p358]

H. Ishwaran and U. B. Kogalur. *Random Forests for Survival, Regression and Classification (RF-SRC)*, 2016. URL http://cran.r-project.org/package=randomForestSRC. [p353]

B. Klimt and Y. Yang. The enron corpus: A new dataset for email classification research. *Machine Learning: ECML 2004*, pages 217–226, 2004. URL https://doi.org/10.1007/978-3-540-30115-8_22. [p362, 363]

M. B. Kursa and A. A. Wieczorkowska. Multi-label ferns for efficient recognition of musical instruments in recordings. In *International Symposium on Methodologies for Intelligent Systems*, pages 214–223. Springer, 2014. URL https://doi.org/10.1007/978-3-319-08326-1_22. [p352, 353]

M. Lang, H. Kotthaus, P. Marwedel, C. Weihs, J. Rahnenführer, and B. Bischl. Automatic model selection for high-dimensional survival analysis. *Journal of Statistical Computation and Simulation*, 85 (1):62–76, 2015. URL https://doi.org/10.1080/00949655.2014.929131. [p352]

A. Liaw and M. Wiener. Classification and regression by randomForest. *R News: The Newsletter of the R Project*, 2(3):18–22, 2002. [p363]

A. McCallum. Multi-label text classification with a mixture model trained by EM. *AAAI'99 Workshop on Text Learning*, pages 1–7, 1999. [p352]

E. Montañés, R. Senge, J. Barranquero, J. R. Quevedo, J. J. del Coz, and E. Hüllermeier. Dependent binary relevance models for multi-label classification. *Pattern Recognition*, 47(3):1494–1508, 2014. URL https://doi.org/10.1016/j.patcog.2013.09.029. [p353, 354, 355, 357, 358, 363]

P. Probst. Multilabel classification with R package mlr. figshare. Code may be downloaded here, 2017. URL https://doi.org/10.6084/m9.figshare.3384802.v5. [p363]

J. Read. Scalable multi-label classification. *Hamilton, New Zealand: University of Waikato*, 2010. [p362, 363]

J. Read and P. Reutemann. Meka: A multi-label extension to WEKA, 2012. URL http://meka.sourceforge.net/. [p352]

J. Read, B. Pfahringer, G. Holmes, and E. Frank. Classifier chains for multi-label classification. *Machine Learning*, 85:333–359, 2011. URL https://doi.org/10.1007/s10994-011-5256-5. [p355, 356, 362, 363]

J. Read, L. Martino, and D. Luengo. Efficient Monte Carlo methods for multi-dimensional learning with classifier chains. *Pattern Recognition*, (Mdc):1–36, 2013. URL https://doi.org/10.1016/j.patcog.2013.10.006. [p356]

C. Sanden and J. Z. Zhang. Enhancing multi-label music genre classification through ensemble techniques. In *Proceedings of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 705–714, 2011. URL https://doi.org/10.1145/2009916.2010011. [p352]

R. E. Schapire and Y. Singer. BoosTexter: A boosting-based system for text categorization. *Machine Learning*, 39:135–168, 2000. [p352]

R. Senge, J. J. del Coz Velasco, and E. Hüllermeier. Rectifying classifier chains for multi-label classification. *Space, 2 (8)*, 2013. [p354, 356, 359]

J. Sill, G. Takacs, L. Mackey, and D. Lin. Feature-Weighted Linear Stacking. *ArXiv e-prints*, 2009. [p358]

R. Simon. Resampling strategies for model assessment and selection. In W. Dubitzky, M. Granzow, and D. Berrar, editors, *Fundamentals of Data Mining in Genomics and Proteomics SE - 8*, pages 173–186. Springer-Verlag, 2007. URL https://doi.org/10.1007/978-0-387-47509-7_8. [p355]

K. Trohidis, G. Tsoumakas, G. Kalliris, and I. P. Vlahavas. Multi-label classification of music into emotions. *ISMIR*, 8:325–330, 2008. URL https://doi.org/10.1186/1687-4722-2011-426793. [p352, 362, 363]

G. Tsoumakas and I. Katakis. Multi label classification: An overview. *International Journal of Data Warehousing and Mining*, 3(3):1–13, 2007. URL https://doi.org/10.4018/jdwm.2007070101. [p352]

G. Tsoumakas, E. Spyromitros-Xioufis, J. Vilcek, and I. P. Vlahavas. Mulan: A java library for multi-label learning. *Journal of Machine Learning Research*, 12:2411–2414, 2011. [p352]

J. Vanschoren, J. N. van Rijn, B. Bischl, and L. Torgo. OpenML: Networked science in machine learning. *SIGKDD Explorations*, 15(2):49–60, 2013. URL https://doi.org/10.1145/2641190.2641198. [p362]

H. Wang, X. Liu, B. Lv, F. Yang, and Y. Hong. Reliable multi-label learning via conformal predictor and random forest for syndrome differentiation of chronic fatigue in traditional chinese medicine. *PLoS ONE*, 9(6), 2014. URL https://doi.org/10.1371/journal.pone.0099565. [p362]

D. H. Wolpert. Stacked generalization. *Neural Networks*, 5(2):241–259, 1992. URL https://doi.org/10.1016/s0893-6080(05)80023-1. [p358]

M. L. Zhang and Z. H. Zhou. M3MIML: A maximum margin method for multi-instance multi-label learning. In *Proceedings - IEEE International Conference on Data Mining, ICDM*, pages 688–697, 2008. URL https://doi.org/10.1109/icdm.2008.27. [p352, 362, 363]

Z.-H. Zhou and M. ling Zhang. Multi-instance multilabel learning with application to scene classification. *Neural Information Processing Systems*, 40(7):2038–2048, 2007. [p362, 363]

*Philipp Probst*
*Department of Medical Informatics, Biometry and Epidemiology*
*LMU Munich*
*81377 Munich*

*Germany*
probst@ibe.med.uni-muenchen.de

*Quay Au*
*Department of Statistics*
*LMU Munich*
*80539 Munich*
*Germany*
quay.au@stat.uni-muenchen.de

*Giuseppe Casalicchio*
*Department of Statistics*
*LMU Munich*
*80539 Munich*
*Germany*
giuseppe.casalicchio@stat.uni-muenchen.de

*Clemens Stachl*
*Department of Psychology*
*LMU Munich*
*80802 Munich*
*Germany*
clemens.stachl@psy.lmu.de

*Bernd Bischl*
*Department of Statistics*
*LMU Munich*
*80539 Munich*
*Germany*
bernd.bischl@stat.uni-muenchen.de

92

## Appendices

|  | BR(rf) | CC(rf) | NST(rf) | DBR(rf) | STA(rf) | BR(ad) | CC(ad) | NST(ad) | DBR(ad) | STA(ad) | RFERN | RFSRC | BR(fl) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| birds* | 0.4481 | 0.4481 | 0.4466 | 0.4497 | 0.4451 | **0.4156** | 0.4218 | 0.4171 | 0.4233 | 0.4202 | 0.9830 | 0.4777 | 0.5226 |
| emotions | 0.6846 | 0.6575 | 0.6728 | **0.6457** | 0.6626 | 0.6777 | 0.6643 | 0.7031 | 0.6845 | 0.6828 | 0.7992 | 0.6829 | 1.0000 |
| genbase* | 0.0333 | 0.0363 | 0.0393 | 0.0423 | 0.0363 | **0.0045** | **0.0045** | **0.0045** | 0.0060 | **0.0045** | 0.2115 | 0.0091 | 1.0000 |
| langLog* | 0.6836 | 0.6829 | 0.6884 | 0.6842 | 0.6856 | 0.6521 | **0.6349** | 0.6418 | 0.6438 | 0.6466 | 0.8589 | 0.6856 | 0.7021 |
| enron* | 0.8531 | 0.8413 | 0.8560 | 0.8408 | 0.8484 | 0.8496 | **0.819** | 0.8484 | 0.8320 | 0.8408 | 1.0000 | 0.8619 | 0.9982 |
| reuters | 0.3620 | 0.3405 | 0.3575 | 0.311 | 0.3515 | 0.349 | **0.2945** | 0.338 | 0.3495 | 0.3385 | 0.5830 | 0.3695 | 1.0000 |
| image | 0.6635 | 0.6150 | 0.6505 | 0.575 | 0.6445 | 0.63 | **0.539** | 0.6275 | 0.6225 | 0.619 | 0.8365 | 0.6955 | 1.0000 |
| scene | 0.4225 | 0.3926 | 0.4217 | 0.3835 | 0.4046 | 0.3913 | **0.3095** | 0.3805 | 0.3610 | 0.3648 | 0.7540 | 0.4570 | 1.0000 |
| yeast* | 0.8316 | 0.7600 | 0.8201 | 0.8167 | 0.8155 | 0.8304 | **0.7563** | 0.8134 | 0.8217 | 0.806 | 0.9338 | 0.8337 | 0.9855 |
| slashdot* | 0.6140 | 0.5994 | 0.6116 | **0.5859** | 0.6052 | 0.6489 | 0.5923 | 0.6449 | 0.6658 | 0.6396 | 0.9966 | 0.6142 | 0.9675 |

**Table 5:** Subset 0/1 loss

|  | BR(rf) | CC(rf) | NST(rf) | DBR(rf) | STA(rf) | BR(ad) | CC(ad) | NST(ad) | DBR(ad) | STA(ad) | RFERN | RFSRC | BR(fl) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| birds* | 0.6153 | 0.6126 | 0.6197 | 0.6169 | 0.6232 | 0.6589 | 0.657 | 0.6604 | 0.6581 | **0.6621** | 0.0999 | 0.5753 | 0.4774 |
| emotions | 0.5453 | 0.5649 | 0.5464 | **0.5849** | 0.5609 | 0.5519 | 0.5676 | 0.5408 | 0.5727 | 0.5427 | 0.5503 | 0.5332 | 0.0000 |
| genbase* | 0.9834 | 0.9806 | 0.9796 | 0.9773 | 0.9806 | **0.9972** | **0.9972** | **0.9972** | 0.9957 | **0.9972** | 0.8884 | 0.9950 | 0.0000 |
| langLog* | 0.3185 | 0.3188 | 0.3140 | 0.3188 | 0.3161 | 0.3553 | **0.3741** | 0.366 | 0.363 | 0.3615 | 0.1953 | 0.4394 | 0.2241 |
| enron* | 0.4693 | 0.4757 | 0.4694 | 0.4804 | 0.4742 | 0.483 | **0.4987** | 0.4824 | 0.4919 | 0.4847 | 0.1859 | 0.4394 | 0.2241 |
| reuters | 0.6625 | 0.6856 | 0.6682 | 0.7199 | 0.6754 | 0.6873 | **0.7414** | 0.6912 | 0.7197 | 0.6964 | 0.5620 | 0.6482 | 0.0000 |
| image | 0.4068 | 0.4585 | 0.4142 | 0.5225 | 0.4228 | 0.4446 | **0.5508** | 0.4458 | 0.5366 | 0.4564 | 0.4467 | 0.3578 | 0.0000 |
| scene | 0.6064 | 0.6333 | 0.6067 | 0.6463 | 0.6233 | 0.646 | 0.7201 | 0.6505 | **0.7313** | 0.6725 | 0.5513 | 0.5654 | 0.0000 |
| yeast* | 0.5091 | 0.5320 | 0.5138 | 0.514 | 0.5205 | 0.5182 | **0.5345** | 0.522 | 0.5068 | 0.5239 | 0.3674 | 0.4945 | 0.3361 |
| slashdot* | 0.4274 | 0.4421 | 0.4285 | 0.4569 | 0.4385 | 0.3883 | 0.4507 | 0.3925 | **0.4613** | 0.3982 | 0.1651 | 0.4202 | 0.0325 |

**Table 6:** Accuracy

|  | BR(rf) | CC(rf) | NST(rf) | DBR(rf) | STA(rf) | BR(ad) | CC(ad) | NST(ad) | DBR(ad) | STA(ad) | RFERN | RFSRC | BR(fl) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| birds* | 0.2763 | 0.2752 | 0.2897 | 0.2859 | 0.2936 | 0.3755 | 0.3865 | 0.3772 | 0.3687 | 0.3784 | **0.8352** | 0.1949 | 0.0000 |
| emotions | 0.6197 | 0.6474 | 0.6187 | 0.6847 | 0.6358 | 0.6335 | 0.6708 | 0.6293 | 0.7189 | 0.6237 | **0.8276** | 0.6001 | 0.0000 |
| genbase* | 0.9846 | 0.9819 | 0.9809 | 0.9786 | 0.9819 | **0.9977** | **0.9977** | **0.9977** | 0.9962 | **0.9977** | 0.9962 | 0.9955 | 0.0000 |
| langLog* | 0.0334 | 0.0330 | 0.0270 | 0.0331 | 0.0308 | 0.0971 | 0.1191 | 0.1056 | 0.0995 | 0.102 | **0.9264** | 0.0301 | 0.0000 |
| enron* | 0.5426 | 0.5487 | 0.5421 | 0.5580 | 0.5466 | 0.5611 | 0.5902 | 0.5619 | 0.6314 | 0.5633 | **0.771** | 0.4959 | 0.2613 |
| reuters | 0.6733 | 0.6959 | 0.6801 | 0.7338 | 0.6875 | 0.7038 | 0.754 | 0.7046 | **0.9032** | 0.7123 | 0.8598 | 0.6559 | 0.0000 |
| image | 0.4192 | 0.4696 | 0.4228 | 0.5562 | 0.4335 | 0.4581 | 0.5691 | 0.4603 | **0.7787** | 0.4724 | 0.7374 | 0.3598 | 0.0000 |
| scene | 0.6148 | 0.6373 | 0.6134 | 0.6555 | 0.6306 | 0.6614 | 0.7243 | 0.6613 | 0.8174 | 0.6879 | **0.9173** | 0.5662 | 0.0000 |
| yeast* | 0.5722 | 0.6097 | 0.5788 | 0.6035 | 0.5874 | 0.5951 | 0.6229 | 0.5978 | 0.6104 | 0.6013 | **0.6296** | 0.5442 | 0.3365 |
| slashdot* | 0.4267 | 0.4412 | 0.4270 | 0.4574 | 0.4391 | 0.3834 | 0.4526 | 0.3868 | 0.6984 | 0.3931 | **0.8065** | 0.4094 | 0.0000 |

**Table 7:** Recall

|  | BR(rf) | CC(rf) | NST(rf) | DBR(rf) | STA(rf) | BR(ad) | CC(ad) | NST(ad) | DBR(ad) | STA(ad) | RFERN | RFSRC | BR(fl) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| birds* | 0.8812 | 0.8889 | 0.8764 | **0.9056** | 0.8874 | 0.8461 | 0.8349 | 0.8401 | 0.8648 | 0.8605 | 0.0859 | 0.8996 |  |
| emotions | 0.7627 | 0.7242 | 0.7499 | 0.7265 | **0.7644** | 0.7537 | 0.7014 | 0.739 | 0.6783 | 0.7347 | 0.5869 | 0.7577 |  |
| genbase* | 0.9987 | 0.9987 | 0.9987 | 0.9987 | 0.9987 | **0.9995** | **0.9995** | **0.9995** | **0.9995** | **0.9995** | 0.8917 | **0.9995** |  |
| langLog* | 0.7267 | **0.7356** | 0.7058 | 0.7207 | 0.6882 | 0.6874 | 0.7228 | 0.7133 | 0.7014 | 0.6965 | 0.0632 | 0.7233 |  |
| enron* | 0.7283 | 0.7188 | 0.7305 | 0.7092 | 0.7331 | 0.7235 | 0.6807 | 0.7198 | 0.6371 | 0.7233 | 0.1973 | **0.7448** | 0.5135 |
| reuters | 0.9411 | 0.9168 | 0.9346 | 0.8995 | 0.9298 | 0.9014 | 0.7689 | 0.8983 | 0.7465 | 0.8931 | 0.5715 | **0.9562** |  |
| image | 0.7899 | 0.7333 | 0.8029 | 0.7086 | 0.7865 | 0.7841 | 0.6281 | 0.7814 | 0.6036 | 0.7737 | 0.4813 | **0.83** |  |
| scene | 0.9071 | 0.8956 | 0.9112 | 0.8917 | 0.9143 | 0.8936 | 0.81 | 0.8856 | 0.7879 | 0.8872 | 0.5662 | **0.9233** |  |
| yeast* | 0.7372 | 0.7218 | 0.7351 | 0.7055 | 0.7389 | 0.7225 | 0.6947 | 0.7233 | 0.6827 | 0.7159 | 0.4361 | **0.7508** | 0.7478 |
| slashdot* | 0.8365 | 0.8127 | 0.8298 | 0.7927 | 0.8277 | 0.8119 | 0.6804 | 0.8161 | 0.5025 | 0.8196 | 0.1679 | **0.8366** |  |

**Table 8:** Precision (For the featureless learner we have no precision results for several datasets. The reason is that the featureless learner does not predict any value in all observations in these datasets. Hence, the denominator in the precision formula is always zero. *mlr* predicts NA in this case.)
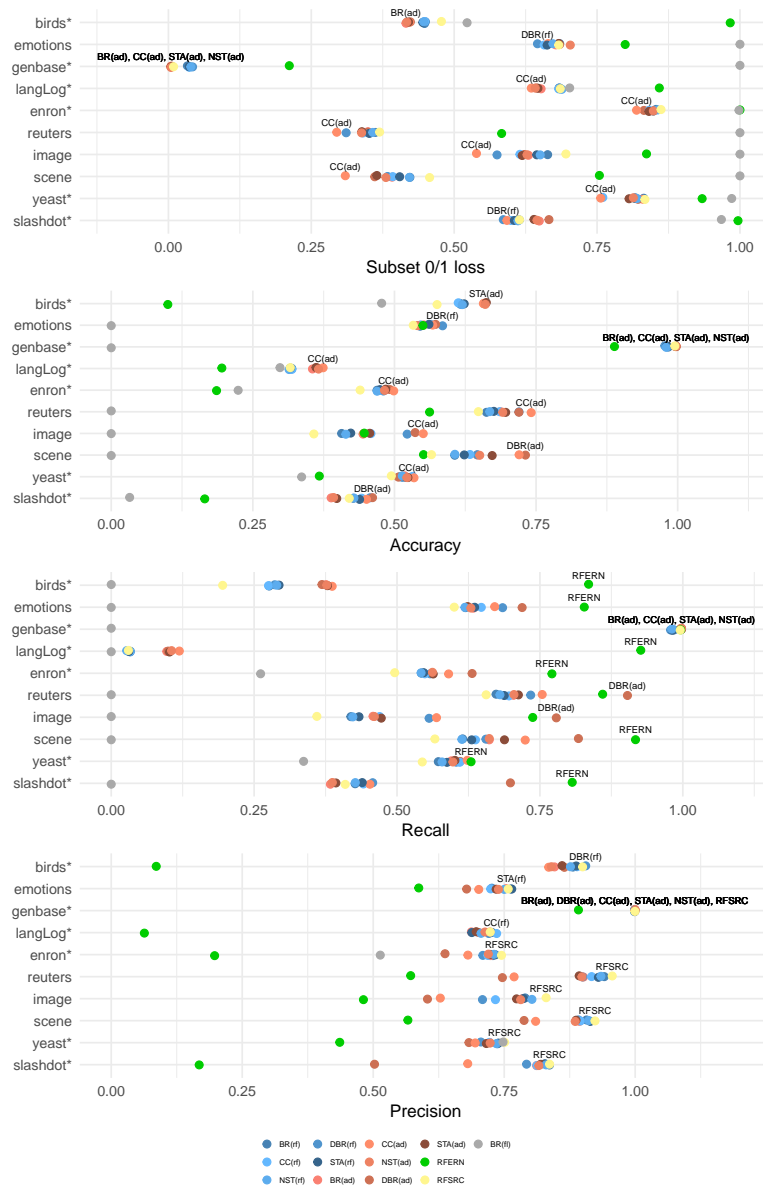
**Figure 2:** Results for the remaining measures.

94

# mlr3: A modern object-oriented machine learning framework in R

The mlr3 framework offers a versatile and expandable platform for machine learning in R. With a modular design and reliance on modern R packages, this ecosystem supports a wide range of ML tasks and serves both practitioner and researchers in model development, experimentation, and evaluation.

***Contributing article***

Lang, M., Binder, M., Richter, J., Schratz, P., Pfisterer, F., Coors, S., Au, Q., Casalicchio, G., Kotthoff, L., and Bischl, B. (2019). mlr3: A modern object-oriented machine learning framework in R. *Journal of Open Source Software*

***Copyright information***

This article is licensed under the Creative Commons Attribution 4.0 International License (CC BY 4.0). See https://creativecommons.org/licenses/by/4.0/. The license allows to copy and redistribute the material in any medium or format.

***Declaration of contributions***

The doctoral candidate assisted the software development team of mlr3 in implementing algorithms for regression and classification.

*Contribution of the coauthors:*

Michel Lang authored the manuscript as the primary author, while the other co-authors provided substantial support in software programming.

***Software:***

- https://cran.r-project.org/web/packages/mlr3/index.html

- https://mlr3.mlr-org.com/

# mlr3: A modern object-oriented machine learning framework in R

**Michel Lang**[1,2]**, Martin Binder**[2]**, Jakob Richter**[1]**, Patrick Schratz**[2]**, Florian Pfisterer**[2]**, Stefan Coors**[2]**, Quay Au**[2]**, Giuseppe Casalicchio**[2]**, Lars Kotthoff**[3]**, and Bernd Bischl**[2]

**1** TU Dortmund University **2** LMU Munich **3** University of Wyoming

## Summary

The R (R Core Team, 2019) package `mlr3` and its associated ecosystem of extension packages implements a powerful, object-oriented and extensible framework for machine learning (ML) in R. It provides a unified interface to many learning algorithms available on CRAN, augmenting them with model-agnostic general-purpose functionality that is needed in every ML project, for example train-test-evaluation, resampling, preprocessing, hyperparameter tuning, nested resampling, and visualization of results from ML experiments. The package is a complete reimplementation of the `mlr` (Bischl et al., 2016) package that leverages many years of experience and learned best practices to provide a state-of-the-art system that is powerful, flexible, extensible, and maintainable. We target both **practitioners** who want to quickly apply ML algorithms to their problems and **researchers** who want to implement, benchmark, and compare their new methods in a structured environment. `mlr3` is suitable for short scripts that test an idea, for complex multi-stage experiments with advanced functionality that use a broad range of ML functionality, as a foundation to implement new ML (meta-)algorithms (for example AutoML systems), and everything in between. Functional correctness is ensured through extensive unit and integration tests.

Several other general-purpose ML toolboxes exist for different programing languages. The most widely used ones are scikit-learn (Pedregosa et al., 2011) for Python , Weka (Hall et al., 2009) for Java, and mlj (Blaom, Kiraly, Lienart, & Vollmer, 2019) for Julia. The most important toolboxes for R are mlr, caret (Kuhn, 2008) and tidymodels (Kuhn & Wickham, 2019).

## Lessons Learned from 6 Years of Machine Learning in R

The predecessor package `mlr` was first released to CRAN in 2013, with the core design and architecture dating back much further. As with most software, more code was added over time to integrate more ML algorithms, more approaches for feature selection or hyperparameter tuning, more methods to analyze trained models, and many other things. With each addition, the code base became larger and more difficult to test and maintain, in particular as changes in the dozens of packages that we integrated with `mlr` would break our code and prevent releases. Installing the package with all dependencies and a complete build with all tests would take hours – we had arrived at a point where adding **any** new functionality became a major undertaking. Further, some of the architectural and design decisions made it essentially impossible to support new cross-cutting functionality, for example ML pipelines, or using new R packages for better performance.

`mlr3` takes these lessons learned to heart and now follows these design principles:

- Be modular and light on dependencies. The core `mlr3` package provides only the basic building blocks of ML: tasks, a few learners, resampling methods, and performance measures. Everything else can be installed and loaded separately through additional packages in the `mlr3` ecosystem, for example support for other kinds of data, methods for tuning hyperparameters, or integrations for additional ML packages.
- Leverage modern R packages, especially `data.table` for fast and efficient computations on rectangular data.
- Embrace `R6` for a clean object-oriented design, object state changes, and reference semantics.
- Defensive programming and type safety. All user input is checked with `checkmate` (Lang, 2017). Return types are documented and automatic type casting for "simplification" is avoided.

In addition, we simplified the API considerably by unifying container and result classes. Many result objects are now tabular by mixing `data.table`'s list-column feature with R6 objects, which also allows for easy and efficient selection and "split-apply-combine" type operations.

## Ecosystem

In addition to the main `mlr3` package, `mlr3learners` provides integrations to a careful selection of the most important ML algorithms and packages in R. Complex ML workflows (using directed acyclic graphs) that can incorporate preprocessing, (stacking) ensembles, alternative-branch execution, and much more can be built with the `mlr3pipelines` package. Funtionality for hyperparameter tuning and nested resampling of learners and complex pipelines is provided by the `mlr3tuning` package. `mlr3filters` integrates many feature filtering techniques and `mlr3db` allows direct use of databases as data sources for out-of-memory data. We are planning and working on many more packages; for example for Bayesian optimization, Hyperband, probabilistic regression, survival analysis, and spatial and temporal data. A complete list of existing and planned extension packages can be found on the mlr3 wiki.

`mlr3` and its ecosystem are documented in numerous manual pages and a comprehensive book (work in progress). All packages are licensed under GNU Lesser General Public License (LGPL-3).

## Acknowledgments

## References

Bischl, B., Lang, M., Kotthoff, L., Schiffner, J., Richter, J., Studerus, E., Casalicchio, G., et al. (2016). Mlr: Machine learning in r. *Journal of Machine Learning Research*, *17*(170), 1–5. Retrieved from http://jmlr.org/papers/v17/15-066.html

Blaom, A., Kiraly, F., Lienart, T., & Vollmer, S. (2019). *Alan-turing-institute/mlj.jl: V0.5.3*. Zenodo. doi:10.5281/zenodo.3541506

2

Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., & Witten, I. H. (2009). The WEKA Data Mining Software: An Update. *ACM SIGKDD explorations newsletter*, *11*(1), 10–18. doi:10.1145/1656274.1656278

Kuhn, M. (2008). Building predictive models in r using the caret package. *Journal of Statistical Software, Articles*, *28*(5), 1–26. doi:10.18637/jss.v028.i05

Kuhn, M., & Wickham, H. (2019). *Tidymodels: Easily install and load the 'tidymodels' packages*. Retrieved from https://CRAN.R-project.org/package=tidymodels

Lang, M. (2017). checkmate: Fast Argument Checks for Defensive R Programming. *The R Journal*, *9*(1), 437–445. doi:10.32614/RJ-2017-028

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., et al. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, *12*, 2825–2830. Retrieved from http://jmlr.org/papers/v12/pedregosa11a.html

R Core Team. (2019). *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing. Retrieved from https://www.R-project.org/

# PART II - Uncovering Patterns in Smartphone Sensor Data

# The PhoneStudy Project

The PhoneStudy project[1] aims to develop smartphone applications, a server infrastructure, and adjacent data-tools to efficiently collect and analyze behavioral data in psychology and other sciences.

***Contributing project***

Stachl, C., Schoedel, R., Au, Q., Völkel, S., Buschek, D., Hussmann, H., Bischl, B., and Bühner, M. (2022). The phonestudy project

***Copyright information***

Since this is a research project, there are no applicable copyright restrictions.

***Declaration of contributions***

The PhD candidate played a crucial role in maintaining the data quality and preprocessing of the project. He was responsible for extracting relevant variables from the smartphone data and developing an internal documented R package that was available for the entire project team, making it accessible for other members. In addition, the PhD candidate conducted a thorough data cleaning process to ensure that the data were reliable and suitable for analysis.

*Contribution of the coauthors*

The PhoneStudy[2] project is an ongoing development with changes in responsibilities of the team members over time. During the time of the PhD candidate's involvement, Clemens Stachl was the project lead responsible for conducting studies and defining the necessary features for data analyses. Ramona Schoedel took over as project lead and also conducted studies while defining features for analyses. Sarah Theres Völkel was in charge of leading the app development team and ensured the continuous improvement of the smartphone application and the server infrastructure. Heinrich Hussmann, Bernd Bischl, and Markus Bühner provided valuable support to the team with their expertise and advice.

---

[1] https://osf.io/ut42y/
[2] https://phonestudy.org/

# Predicting Personality from Patterns of Behavior Collected with Smartphones

Big Five personality traits were predicted from smartphone-based behavioral data using machine learning. These findings highlight the potential benefits and risks of smartphone data collection and modeling.

### *Contributing article*

Stachl, C., Au, Q., Schoedel, R., Gosling, S. D., Harari, G. M., Buschek, D., Völkel, S. T., Schuwerk, T., Oldemeier, M., Ullmann, T., Hussmann, H., Bischl, B., and Bühner, M. (2020). Predicting personality from patterns of behavior collected with smartphones. *Proceedings of the National Academy of Sciences*, 117(30):17680–17687

This publication was part of the PhoneStudy project (see Chapter 7).

### *Copyright information*

### *Declaration of contributions*

Besides the contributions to the PhoneStudy project, the PhD candidate employed novel methods to interpret and visualize the results of the data analysis. Additionally, he wrote parts of the chapter about the data analysis and the interpretation of the results. He also created corresponding graphics to make the analysis results more accessible and understandable. In addition, the PhD candidate also developed an interactive website[1] that allows users to explore the results of the analysis interactively.

*Contribution of the coauthors*

As the lead author, Clemens Stachl was responsible for writing the manuscript. Samuel Gosling and Gabriella Harari provided valuable support by engaging in discussions and contributing to

---

[1] `https://compstat-lmu.shinyapps.io/Personality_Prediction/`

the writing of the manuscript. The remaining co-authors played important roles in the Phone-Study project and provided support with final revisions to the manuscript.

# Predicting personality from patterns of behavior collected with smartphones

Clemens Stachl[a,1], Quay Au[b], Ramona Schoedel[c], Samuel D. Gosling[d,e], Gabriella M. Harari[a], Daniel Buschek[f], Sarah Theres Völkel[g], Tobias Schuwerk[h], Michelle Oldemeier[c], Theresa Ullmann[i], Heinrich Hussmann[g], Bernd Bischl[b], and Markus Bühner[c]

[a]Department of Communication, Media and Personality Laboratory, Stanford University, Stanford, CA 94305; [b]Department of Statistics, Computational Statistics, Ludwig-Maximilians-Universität München, 80539 Munich, Germany; [c]Department of Psychology, Psychological Methods and Assessment, Ludwig-Maximilians-Universität München, 80802 Munich, Germany; [d]Department of Psychology, University of Texas at Austin, Austin, TX 78712; [e]School of Psychological Sciences, University of Melbourne, Parkville, VIC 3010, Australia; [f]Research Group Human Computer Interaction and Artificial Intelligence, Department of Computer Science, University of Bayreuth, 95447 Bayreuth, Germany; [g]Media Informatics Group, Ludwig-Maximilians-Universität München, 80337 Munich, Germany; [h]Department of Psychology, Developmental Psychology, Ludwig-Maximilians-Universität München, 80802 Munich, Germany; and [i]Institute for Medical Information Processing, Biometry, and Epidemiology, Ludwig-Maximilians-Universität München, 81377 Munich, Germany

**Smartphones enjoy high adoption rates around the globe. Rarely more than an arm's length away, these sensor-rich devices can easily be repurposed to collect rich and extensive records of their users' behaviors (e.g., location, communication, media consumption), posing serious threats to individual privacy. Here we examine the extent to which individuals' Big Five personality dimensions can be predicted on the basis of six different classes of behavioral information collected via sensor and log data harvested from smartphones. Taking a machine-learning approach, we predict personality at broad domain ($r_{median}$ = 0.37) and narrow facet levels ($r_{median}$ = 0.40) based on behavioral data collected from 624 volunteers over 30 consecutive days (25,347,089 logging events). Our cross-validated results reveal that specific patterns in behaviors in the domains of 1) communication and social behavior, 2) music consumption, 3) app usage, 4) mobility, 5) overall phone activity, and 6) day- and night-time activity are distinctively predictive of the Big Five personality traits. The accuracy of these predictions is similar to that found for predictions based on digital footprints from social media platforms and demonstrates the possibility of obtaining information about individuals' private traits from behavioral patterns passively collected from their smartphones. Overall, our results point to both the benefits (e.g., in research settings) and dangers (e.g., privacy implications, psychological targeting) presented by the widespread collection and modeling of behavioral data obtained from smartphones.**

personality | behavior | machine learning | mobile sensing | privacy

It has been well documented that "digital footprints" derived from social network platforms (e.g., Facebook likes) can reveal individuals' psychological characteristics, such as their personality traits (1). This is consequential because the Big Five personality traits have been shown to predict a broad range of life outcomes in the domains of health, political participation, personal and romantic relationships, purchasing behaviors, and academic and job performance (2–4). Data-driven inferences about individuals' personality traits present great opportunities for research; but they also have major implications for individual privacy because they allow for personality-based targeting and manipulation (5, 6).

Even greater threats to privacy are posed by smartphones, which can collect a far broader, fine-grained array of daily behaviors than can be scraped from social media platforms and which are pervasive in most societies around the globe (7). The on-board sensors of a smartphone and the device's logging capabilities (e.g., app-usage logs, media and website consumption, location, communications, screen activity) can be harnessed by apps to record daily behaviors performed both on the devices themselves and in close proximity to them (8–10). These data

have great potential for psychological research and have already begun to yield valuable findings, including studies relating physical activity and communication data to human emotion and mental wellbeing (11–14). However, behavioral data from smartphones can contain private information and should therefore be collected and processed only when informed consent is given (15). In theory, users must give permission for apps to access certain types of data on their phones (e.g., to record location or audio data). However, people are often unaware of the data they are providing, are tricked into giving access to more data (16), and struggle to understand current permission systems that are unspecific and ineffective in preventing the collection of personal data from smartphones (17–19). Finally, many apps find creative side channels to routinely extract data from people's phones (20, 21)—regardless of whether permission has been provided.

Here we evaluate whether individuals' Big Five personality trait levels can be predicted on the basis of six different classes of

behavioral information collected via smartphones. Moreover, we examine which behaviors reveal most about each personality trait and how predictive each behavioral class is on average. Using sensor and log data from volunteers' smartphones, we extracted thousands of variables, categorized into six classes of daily behavior derived from previous research: 1) app usage (e.g., mean duration of gaming app usage), 2) music consumption (e.g., mean valence of played songs), 3) communication and social behavior (e.g., number of outgoing calls per day), 4) mobility behaviors (e.g., mean radius of gyration), 5) overall phone activity (e.g., number of unlock events per day), and 6) a higher-level behavioral class that captured the extent of daytime versus nighttime activity (e.g., outgoing calls at night). Together these six classes of behavior provided a broad sampling of the data that can easily be derived from smartphones and which may provide clues to individuals' personalities and allow for a robust investigation of our research question.

We assessed personality in terms of the Big Five dimensions, the most widely used and well-established system in psychological science for organizing personality traits (22–24). This taxonomy describes human personality in terms of five broad and relatively stable dimensions: openness, conscientiousness, extraversion, agreeableness, and emotional stability (22, 23), with each dimension subsuming a larger number of more specific facets. The Big Five have been found to have a strong genetic basis and to replicate across cultures and contexts (25–27).

Past studies have highlighted the promise of using smartphones to associate behavioral data with personality traits and other private attributes (28–39). A subset of these studies has used machine learning in analyses with the goal of predicting personality traits from behavioral measures (28–30, 38, 39). However, this subset of studies was subject to a number of key limitations, including the following: 1) focusing on just a single class of behavior or a small number of similar behaviors (e.g., communication behavior, refs. 31 and 39); 2) using small samples (28–30, 39); 3) being confined to the broad personality trait domains, not their more specific facets (28–30, 38, 39); 4) using classification instead of regression for the prediction of continuous personality scores (28, 29, 31); 5) likely overestimating model performance (28–30) (see ref. 31, for a discussion of the problem); 6) not providing enough information to reproduce findings (e.g., open data and materials, refs. 28–30, 38, and 39); and 7) not determining the relative effects of variables in the prediction models (28–30, 38).

To address these issues, we use smartphone sensing to gather behaviors from a wide variety of behavioral classes from a large sample, measure personality at both the domain and facet levels, train linear and nonlinear regression models (elastic net, random forest), properly evaluate our models out of sample using a (nested) cross-validated approach, and explore which behaviors are most predictive of personality overall and with respect to the individual personality domains and facets using interpretable machine learning and corrected significance tests. As a benchmark for the performance of our models, we compare the predictive performance with that of previous research using digital footprints from social media platforms (e.g., ref. 1).

## Results

**Personality Trait Prediction with Behavioral Patterns.** Descriptive statistics can be found in *SI Appendix*, Tables S1 and S2, and in extensive detail on the project's website, accessible via the project repository (40). The results show that we successfully predicted levels of Big Five personality traits from behavioral patterns, derived from smartphone data, for more than half of the domains and facets (57% of all personality dimensions). In multiple instances both model types performed well above

the baseline model (i.e., a model that constantly predicts the mean in the respective training set). Furthermore, our results suggest differences in how well the trait dimensions were predicted, as can be seen in Fig. 1 and in *SI Appendix*, Table S4 (e.g., sociableness most accurately and agreeableness not at all). The results also show that the nonlinear random forest models on average outperformed the linear elastic net models in both prediction performance and the number of successfully predicted criteria, hinting at the presence of nonlinear correlational structures in the data. Table 1 shows the top five most-important predictor variables per criterion. In Fig. 2 we provide a comprehensive visualization of all model results and effects of the behavioral classes. Fig. 2, *Top* shows the median prediction performance in $R^2$, and Fig. 2, *Upper Middle* shows the contribution and significance of a behavioral class by itself for the respective model (unique class importance). Fig. 2, *Lower Middle* shows the contribution of a behavioral class in the context of all other classes (combined class importance). Red circles indicate significant effects. In Fig. 2, *Bottom*, color-coded behavioral patterns ranked by variable importance are displayed across all models.

Here we report median prediction performances for all personality trait models, aggregated across the outer cross-validation folds. We report all metrics for both model types in *SI Appendix*, Table S4. In *SI Appendix*, Fig. S1 we also show exploratory predictor effects in accumulated local effect plots (ALEs). Additionally, we provide $P$ values for the behavioral class effects, in *SI Appendix*, Table S5. For clarity and due to the model's superiority in prediction, we report performance metrics only for the random forest models in the text. However, results for both types of models, including plots, variable importance measures, and all exploratory single-predictor effects, are available on the project's website, accessible via the project's repository (40). In addition to results from predictive modeling, we also summarize findings from the interpretable machine-learning analyses. Below we describe which classes of behavior were significantly predictive for the respective personality dimension and provide some illustrative examples of single-variable effects, which should not be generalized beyond our sample. Finally, by refitting models on all combinations of the behavioral classes, we evaluate the average effect of each class for the prediction of personality trait dimensions. Data and code to reproduce all analyses are available in the project's repository (40).

Except for openness to imagination ($r_{md} = 0.19$, $r_{sd} = 0.13$), openness ($r_{md} = 0.29$, $r_{sd} = 0.11$) and its facets were successfully predicted in our dataset. With regard to facets, openness to aesthetics showed the highest median prediction performance ($r_{md} = 0.29$, $r_{sd} = 0.12$) and openness to actions ($r_{md} = 0.23$, $r_{sd} = 0.11$) the lowest, with openness to feelings ($r_{md} = 0.24$, $r_{sd} = 0.09$) and openness to ideas falling in between ($r_{md} = 0.24$, $r_{sd} = 0.11$). The top predictors in Table 1 and behavioral patterns in Fig. 2 suggest that music consumption also played a role in the prediction models for openness (e.g., quieter music), but this could not be confirmed by the unique and combined class-based variable importance scores in Fig. 2. Those scores suggest that overall patterns in app-usage behavior (e.g., increased camera usage, more photos, less usage of sports news apps) and for openness to actions communication and social behavior (e.g., ringing events, calls at night) were most important for the prediction of openness and its facets.

Conscientiousness ($r_{md} = 0.31$, $r_{sd} = 0.13$) was also successfully predicted above baseline, as were its facets, except for competence ($r_{md} = 0.19$, $r_{sd} = 0.11$). In terms of prediction performance, the facet love of order ranked first ($r_{md} = 0.31$, $r_{sd} = 0.13$), followed by sense of duty ($r_{md} = 0.29$, $r_{sd} = 0.10$), ambition ($r_{md} = 0.26$, $r_{sd} = 0.12$), discipline ($r_{md} = 0.22$, $r_{sd} = 0.12$),

**Fig. 1.** Box and whisker plot of prediction performance measures from repeated cross-validation for each personality domain and facet. The middle symbol represents the median, boxes include values between the 25 and 75% quantiles, and whiskers extend to the 2.5 and 97.5% quantiles. Outliers are depicted by single points. Names of significant models are in boldface type. Figure is available at https://osf.io/kqjhr/, under a CC-BY4.0 license.

and caution ($r_{md} = 0.20$, $r_{sd} = 0.12$). Inspection of behavioral patterns and class importance indicators in Fig. 2 suggests that in the context of all other variables, predominantly variables related to overall phone activity (e.g., earlier first and last phone use per day), day and nighttime activity (e.g., less variable nightly duration of inactivity), and most unique app usage (e.g., increased usage of weather apps, timers, and checkup-monitoring apps) were especially important for the prediction of higher scores in the models of conscientiousness and its facets. Additionally, for the facets love of order and sense of duty, a very specific behavior was found to be important—the mean charge of the phone when it was disconnected from a charging cable. ALEs in *SI Appendix*, Fig. S1 suggest that in the context of all predictors higher average scores in love of order were predicted for charges above 60%.

With the exception of the cheerfulness facet ($r_{md} = 0.16$, $r_{sd} = 0.12$), the personality trait of extraversion ($r_{md} = 0.37$, $r_{sd} = 0.09$) and its facets were successfully predicted above baseline. Most notably, the facet of sociableness was predicted with the highest performance of all criteria ($r_{md} = 0.40$, $r_{sd} = 0.10$). Besides sociableness, the facets friendliness ($r_{md} = 0.24$, $r_{sd} = 0.09$), assertiveness ($r_{md} = 0.29$, $r_{sd} = 0.11$), dynamism ($r_{md} = 0.29$, $r_{sd} = 0.10$), and adventurousness ($r_{md} = 0.29$, $r_{sd} = 0.11$) were predicted above baseline. Behavioral patterns and class importance (unique and combined) in Fig. 2 suggest that variables related to communication and social behavior (e.g., higher mean number of outgoing calls per day, higher irregularity of all calls,

higher mean number of WhatsApp uses per day) were important in the prediction of higher scores in the models of extraversion and its facets.

In the present analyses, the personality dimension of agreeableness could not be successfully predicted from the data, either on domain or on facet levels ($r_{md} = 0.05$, $r_{sd} = 0.11$).

For the personality dimension of emotional stability, only the facets of carefreeness ($r_{md} = 0.22$, $r_{sd} = 0.10$) and self-consciousness ($r_{md} = 0.32$, $r_{sd} = 0.09$) were predicted significantly. Behavioral patterns in Fig. 2 are rather distinct for the individual facets of emotional stability. Whereas communication and social behavior were significantly predictive for the facet self-consciousness (e.g., higher number of calls), the model of carefreeness did not show any significant effects at the class level.

In summary, all behavioral classes had some impact on the prediction of personality trait scores (as seen in Fig. 2). However, behaviors related to communication and social behavior and app usage showed as most significant in the models. This pattern can be discerned in Fig. 2. To estimate the average effect of each behavioral class on the prediction of personality trait dimensions overall (successfully and unsuccessfully predicted in the main analyses), we used a linear mixed model (details of the analysis are described in *Materials and Methods*). Results of the model show that communication and social behavior had the biggest impact on model performance on domain ($\beta = 0.027$, $CI_{95\%} = [.026, .028]$) and facet levels ($\beta = 0.019$,

107

**Table 1. Top five predictors per prediction model**

| Personality dimension | Top five predictors |
| --- | --- |
| O, openness | Daily mean length text messages \| robust mean dur sports news apps \| daily robust variation dur phone ringing \| daily robust mean no. photos \| robust mean dur sports news apps night |
| O2, openness to aesthetics | Robust mean dur sports news apps \| daily mean no. photos \| daily mean no. unique sports news apps \| robust mean dur nightly sports news app \| daily mean no. sports news apps |
| O3, openness to feelings | Excess music acousticness \| daily mean no. unique sports news apps per week \| robust variation dur shared transportation apps \| daily robust variation in dur phone ringing \| daily mean no. unique sports news apps |
| O4, openness to actions | Mean no. of phone ringing night \| daily mean no. of ringing events \| daily mean no. Google Maps \| mean no. calls night \| irregularity of phone ringing |
| O5, openness to ideas | Loudness fourth most listened song \| robust mean dur sports news apps \| daily SD no. of photos \| robust mean dur *Süddeutsche Zeitung* (newspaper) \| robust mean dur Samsung Notes |
| O6, openness to value and norm | Daily mean no. unique sports news week \| daily mean no. Facebook \| daily mean no. sports news \| daily mean no. unique sports news weekend \| daily mean no. Kicker (soccer news) |
| C, conscientiousness | Robust mean dur weather app night \| daily SD sum interevent time \| robust mean time last event \| robust variation dur checkup monitoring apps \| robust variation first event weekdays |
| C2, love of order | Daily SD sum interevent time \| robust mean dur news-magazine apps \| daily mean no. unique email apps \| mean mean charge disconnection \| robust variation dur TV-filmguide apps |
| C3, sense of duty | SD dur nightly downtime \| robust mean time first event weekdays \| robust variation time last event weekdays \| robust mean dur Stadtwerke München Fahrinfo München (public transportation) |
| C4, ambition | Robust mean time first event \| robust variation first event weekdays \| robust mean time last event \| robust variation time first event weekends \| daily mean no. Google Playstore |
| C5, discipline | Robust variation time first event weekdays \| robust mean time first event weekdays \| robust mean dur weather apps night \| robust variation time first event weekends \| daily SD sum interevent time |
| C6, caution | Robust variation time last event weekdays \| SD dur nightly downtime Sunday til Thursday \| similarity contacts phone and messaging \| robust variation time last event \| mean music valence weekends |
| E, extraversion | Nightly mean no. phone ringing \| nightly mean no. calls \| daily mean no. outgoing calls \| daily mean no. phone ringing \| nightly mean no. outgoing calls |
| E1, friendliness | Daily mean no. phone ringing \| irregularity of phone ringing weekend \| daily SD no. incoming calls \| daily robust variation sum dur phone ringing \| daily SD sum dur incoming calls |
| E2, sociableness | Mean no. calls night \| daily mean no. outgoing calls \| mean no. phone ringing night \| mean no. outgoing calls night \| irregularity of phone ringing weekend |
| E3, assertiveness | Daily mean no. outgoing calls \| daily mean no. contacts per week \| daily mean no. contacts outgoing calls \| daily mean no. contacts calls \| mean no. calls night |
| E4, dynamism | Daily mean no. outgoing calls \| mean no. phone ringing night \| daily mean no. contacts outgoing calls \| mean no. calls night \| daily mean no. phone ringing |
| E5, adventurousness | Mean no. phone ringing night \| mean no. calls night \| irregularity of phone ringing \| mean no. outgoing calls night \| irregularity of calls |
| ES1, carefreeness | Daily mean no. Android-Email (app) \| daily mean no. screen unlocks \| robust variation dur system apps \| robust variation dur strategy games \| daily mean no. phone ringing |
| ES4, self-consciousness | Nightly mean no. calls \| daily mean no. phone ringing \| daily mean no. contacts calls \| daily mean no. outgoing calls \| daily mean no. contacts incoming calls |

The top five most predictive features are shown for each successfully predicted personality dimension in the random forest models. The ranking is based on permutation feature importance and goes from left (high) to right (low). dur = duration.
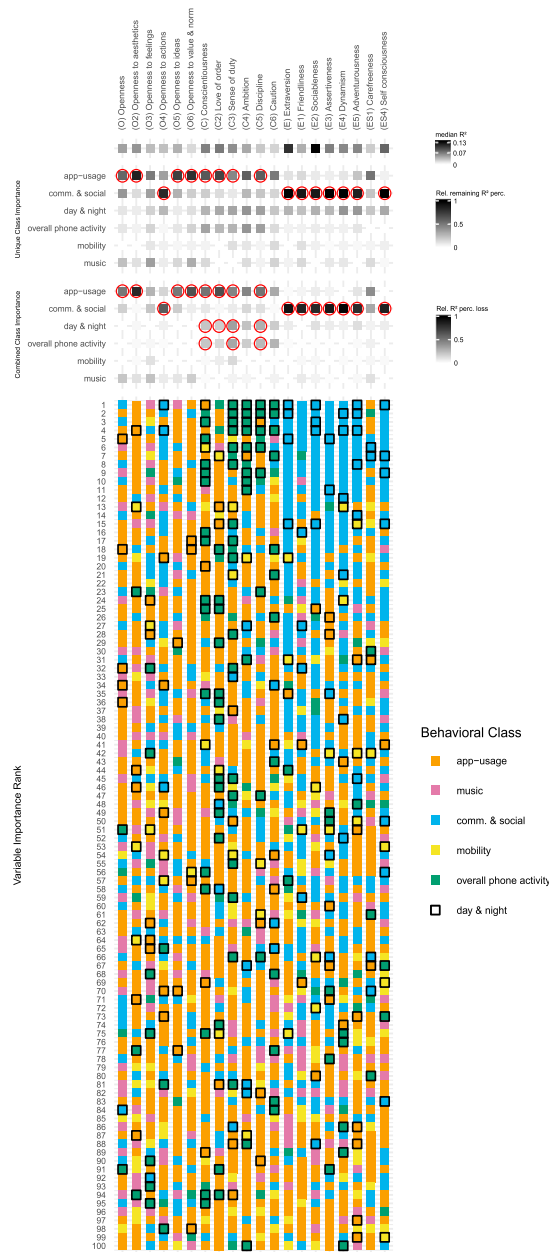
$\text{CI}_{95\%} = [.019, .020]$). App usage was second ($\beta_{\text{domains}} = 0.014$, $\text{CI}_{95\%} = [.013, .015]$, $\beta_{\text{facets}} = 0.014$, $\text{CI}_{95\%} = [.014, .015]$) followed by day and nighttime activity ($\beta_{\text{domains}} = 0.013$, $\text{CI}_{95\%} = [.012, .014]$, $\beta_{\text{facets}} = 0.011$, $\text{CI}_{95\%} = [.011, .012]$), overall phone activity ($\beta_{\text{domains}} = 0.006$, $\text{CI}_{95\%} = [.005, .007]$, $\beta_{\text{facets}} = 0.004$, $\text{CI}_{95\%} = [.004, .005]$), and music ($\beta_{\text{domains}} = 0.001$, $\text{CI}_{95\%} = [.000, .002]$, $\beta_{\text{facets}} = 0.001$, $\text{CI}_{95\%} = [.001, .002]$). The behavioral class of mobility was least important for the prediction of Big Five personality trait dimensions ($\beta_{\text{domains}} = -0.001$, $\text{CI}_{95\%} = [-.002, -.001]$, $\beta_{\text{facets}} = -0.001$, $\text{CI}_{95\%} = [-.001, .000]$). In *SI Appendix*, Fig. S2, we provide additional, exploratory results of a resampled greedy forward search analysis, indicating which combinations of behavioral classes were most predictive overall, in our dataset.

### Discussion

The results presented here demonstrate that information about individuals' everyday behaviors detected from smartphone sensors and logs can be used to infer their Big Five personality trait dimensions. Specific classes of behavior (app usage, music consumption, communication and social behavior, mobility behavior, overall phone activity, daytime vs. nighttime activity) were distinctively informative about the different Big Five trait dimensions. Our models were able to predict personality on the broad domain level and the narrow facet level for openness, conscientiousness, and extraversion. For emotional stability, only single facets could be predicted above baseline. Finally, scores for agreeableness could not be predicted at all. The behavioral class of communication and social behavior was most important for the prediction of personality trait dimensions on average, but app usage and day and nighttime activity were also important[*]. We found performance levels across all significant models ($r_{range} = [0.20, 0.40]$) to be on average similar to those identified in a metaanalysis of previous studies predicting personality from digital footprints, which reported a mean effect size of $r = 0.34$ (1). As benchmarks for gauging these effect sizes,

---

[*]As can be seen in Fig. 2, in roughly half of the models the behavioral class communication and social was most important and, for the other half, app usage was most important.

108

**Fig. 2.** (*Top*) Median prediction performance in $R^2$. (*Upper Middle*) Relative remaining performance when keeping variables of the respective class intact and permuting variables of other classes (unique class importance). (*Lower Middle*) Relative drop in performance when permuting variables of respective groups (combined class importance). Red circles indicate significant effects tested with the PIMP algorithm (41). (*Bottom*) Behavioral patterns of ranked permutation-based variable importance, color coded by class of behavior. Black frames indicate additional day–night dependency. Figure is available at https://osf.io/kqjhr/, under a CC-BY4.0 license.

consider that the highest median effect size ($r_{md} = 0.40$) is comparable to the tendency of people in a bad mood to be more aggressive than people in a good mood, and the smallest significant effect size is equal to the average reported effect size in personality psychology (42). These performance levels highlight the practical relevance of our results beyond significance.

The results here point to the breadth of behavior that can easily be obtained from the sensors and logs of smartphones and, more importantly, the breadth and specificity of personality predictions that can be made from the behavioral data so obtained. However, it is important to note that these findings are, if anything, a conservative estimate of what can be learned about people's personalities using information obtainable from their smartphones. Greater prediction accuracies would almost certainly be obtained when using more sensors (e.g., accelerometers, microphones, cameras; ref. 11); more log data (e.g., over longer time periods); content-level data (e.g., the content of texts, calls, emails, photos, videos, or all visible information on the screen; ref. 43); bigger, more diverse, and more representative samples (e.g., iPhone operating system [iOS] and Android users, nonwestern, educated, industrialized, rich, and democratic [WEIRD] samples; ref. 44); and by combining these data with other information about the user, derived from other sources (e.g., purchase histories, digital footprints from social media). Furthermore, models in this paper are still limited by the sparsity in the data (e.g., app usage), because some apps were used by only very few participants. Larger samples (e.g., as used in studies on personality social media use; ref. 6) could also allow for more accurate predictions.

As such, the present work serves as a harbinger of both the benefits and the dangers presented by the widespread use of behavioral data obtained from smartphones. On the positive side, obtaining behavior-based estimates of personality stands to open additional avenues of research on the causes and consequences of personality traits, as well as permitting consequential decisions (e.g., in personnel selection) to draw on behavioral data rather than estimates derived from self-report questionnaires, which are subject to a range of biases (e.g., responses biases, social desirability, different reference standards, memory limitations; refs. 45 and 46).

At the same time, we should not underestimate the potential negative consequences of the routine collection, modeling, and uncontrolled trade of personal smartphone data (20, 21, 47). For example, organizations and companies can obtain information about individuals' private traits (e.g., the Big Five personality traits), without the personality information ever being deliberately provided or explicitly requested (48). Mounting evidence suggests that these data can and are being used for psychological targeting to influence people's actions, including purchasing decisions (5, 47) and potentially voting behaviors, which are related to personality traits (49, 50).

Many commercial actors already collect a subset of the behavioral data that we have used in this work using publicly available applications (20). In academic settings, such data collection requires institutional review board (IRB) approval of the research study. However, current data protection laws in many nations do not adequately regulate data collection practices in the private sector. For example, in online real-time bidding on advertisements multiple actors exchange cross-device data to win bids to cater personalized ads to single users; this process is complex, happens within milliseconds, and is poorly understood outside of the industry (47). In such cases, once the data are collected from people's smartphones, the data's distribution seems to largely escape legislative oversight and legal enforcement (21, 47). This is the case even though legal frameworks against the routine collection of these data exist (e.g., the General Data Protection Regulation [GDPR] in the European Union; ref. 51) and reflects the growing asymmetry between one-click privacy

109

permissions and the untraceable ways behavioral data from peoples' phones can wander.

Hence, a more differentiated choice with regard to the types of data and their intended usage should be given to users. For example, users should be made aware that behavioral data from phones are required for the completion of a specific task (e.g., finding a café); could be reused or sold to third parties, combined with other data; or used to create user models to make indirect predictions (e.g., personality, financial, credit scoring). In other words, it must be more obvious to consumers whether they are consenting to the measurement of their app use or to the automatic prediction of their private traits (e.g., personality).

Under most legislation, all of these actions are currently possible after initially providing the permission to access data on phones. One idea is for user data to have an automatic expiration date, after which data attributable to a unique identity must be deleted. Finally, the manifold techniques that online marketing companies use to link datasets of individuals to facilitate personalized ads (i.e., unique identifiers; ref. 47) could also be used to opt out of all advertisements and data-processing activities. Some variations of these suggestions are already implemented in the European Union's GDPR (51). We hope our findings stimulate further debate on the sensitivity of behavioral data from smartphones and how privacy rights can be protected at the individual (15) and aggregate levels (52).

A large portion of current economic and scientific progress depends on the availability of data about individuals' behaviors. The smartphone represents an ideal instrument to gather such information. Therefore, our results should not be taken as a blanket argument against the collection and use of behavioral data from phones. Instead, the present work points to the need for increased research at the intersection of machine learning, human computer interaction, and psychology that should inform policy makers. We believe that to understand complex social systems, while at the same time protecting the privacy of smartphone users, more sophisticated technical and methodological approaches combined with more dynamic and more transparent approaches to informed consent will be necessary (e.g., distributed privacy, federated learning, privacy nudges; refs. 53–56). These approaches could help balance the tradeoff between the collection of behavioral smartphone data and the protection of individual privacy rights, resulting in higher standards for consumers and industry alike.

## Materials and Methods

**Participants and Dataset.** The dataset was collected in three separate studies as part of the PhoneStudy mobile sensing research project at the Ludwig-Maximilians-Universität München (LMU) (57). Parts of the data have been used in other publications (32, 33, 58, 59), but the joint dataset of common parameters has not been analyzed before. A total of 743 volunteers were recruited via forums, social media, blackboards, flyers, and direct recruitment, between September 2014 and January 2018 (33, 58, 59). All subjects participated willingly and provided informed consent prior to their participation in the study. Volunteers could withdraw from participation and demand the deletion of their data as long as their reidentification was possible. Dependent on the respective study (33, 58, 59), we provided different rewards for participation. Procedures for all studies were approved by the IRB of the Psychology Department at Ludwig-Maximilians-Universität München and have been conducted according to European Union laws. In *SI Appendix*, Table S3 we provide an overview of the datasets. We excluded data from volunteers with less than 15 d of logging data (29), no app usage (39), and missing questionnaire data (52). The final sample ($n = 624$) was skewed in favor of more educated (91% completed A levels, 20% had a university degree), younger participants (M = 23.56, SD = 6.63) and was not equally balanced with regard to gender (377 women, 243 men, and 4 with undisclosed gender).

**Procedures.** Study procedures were somewhat different across the three studies (33, 58, 59). However, in all three studies, Big Five personality trait levels were measured with the German version of the Big Five Structure

Inventory (BFSI) (60) and naturalistic smartphone usage in the field was automatically recorded over a period of 30 d. The data were regularly transferred to our encrypted server using Secure Sockets Layer (SSL) encryption, when phones were connected to WiFi. In study 2, volunteers had to answer experience sampling questionnaires during the data collection period on their smartphones (59). Volunteers in studies 2 and 3 completed the demographic and BFSI personality questionnaires via smartphone at a convenient time (58). In cases where volunteers turned off location services, they were reminded to reactivate them. At the end of mobile data collection, volunteers were instructed to contact the research staff to receive compensation (studies 1 to 3) and to schedule a final laboratory session (study 2). More details about the procedures of the individual studies are available in the respective research articles (33, 58, 59).

**Self-Reported Personality Measures and Demographics.** Big Five personality dimensions were assessed with the German version of the BFSI (60). The test consists of 300 items and measures the Big Five personality dimensions (openness to experience, conscientiousness, extraversion, agreeableness, and emotional stability) on five domains and 30 facets. Participants indicated their agreement with items using a four-point Likert scale ranging from untypical for me to typical for me.

Additionally, we collected age, gender, highest completed education, and a number of other questionnaires that were used in other research projects. More information can be found in the respective online repositories and articles (33, 58, 59). Questionnaires were administered either via desktop computer (studies 1 and 2) or via smartphone (studies 2 and 3). We used the laboratory version scores from study 2 in this study. Descriptive statistics including confidence intervals of internal consistencies ($\alpha$) are provided in *SI Appendix*, Table S1.

**Behavioral Data from Smartphone Sensing.** We used the PhoneStudy smartphone research app for Android to collect behavioral data from the volunteers' privately owned smartphones. This app has been continuously developed at the Ludwig-Maximilians-Universität München since September 2013.

Initially, activities were recorded in the form of time-stamped logs of events. Those events included calls, contact entries, texting, global positioning system (GPS) locations, app starts/installations, screen de/activations, flight mode de/activations, Bluetooth connections, booting events, played music, battery charging status, photo and video events, and connections to wireless networks (WiFi). Additionally, the character length of text messages and technical device characteristics were collected. Irreversibly hash-encoded versions of contacts and phone numbers were collected to enable us to measure the number of distinct contacts while preventing the possibility of reidentification. Information such as names, phone numbers, and contents of messages, calls, etc., was not recorded at any time.

**Data Analysis.** The final dataset consisted of 1,821 behavioral predictors and 35 personality criteria (five domains and 30 facets). Gender, age, and education were used solely for descriptive statistics and were not included as predictors in the models.

*Variable extraction.* In a first step, we extracted 15,692 variables from the raw dataset. The extracted variables roughly correspond to the aforementioned behavioral classes of app usage, music consumption, communication and social behavior, mobility, overall phone activity, and day- and night-time dependency. Variables with regard to day and night dependency were not computed for music consumption behaviors. Besides common estimators (e.g., arithmetic mean, SD sum, etc.), we computed more complex variables containing information about the irregularity, the entropy, the similarity, and the temporal correlation of behaviors. These variables provided information about specific data types (e.g., mobility data) and were used for the quantification of behavioral structures within person and across time while avoiding more complex time-series models. The large amounts of data meant it was unfeasible to check for outliers manually, so we used robust estimators (e.g., Huber M Estimator; ref. 61) for most variables (except for call and messaging variables that were checked manually). Details about the calculation of variables and the full set of extracted variables and a detailed overview of all sensed data are provided in the project repository (40).

*Machine learning.* We fitted machine-learning models with an inner cross-validation loop (5-fold cross-validation [CV]) for preprocessing and hyperparameter tuning and an outer cross-validation loop ($10 \times 10$-fold CV) for unbiased model evaluation. We compared the predictive performance of elastic net regularized linear regression models (62) with those of nonlinear tree-based random forest models (63) and a baseline model. The baseline

110

model predicted the mean of the respective training set for all cases in a test set. We chose these standard models due to their ability to cope with $P \gg N$ problems (i.e., few cases, many predictors). Furthermore, the usage of random forest models allowed us to include nonlinear predictor effects and high-dimensional interactions in the models.

We evaluated the predictive performance of the models based on the Pearson correlation ($r$) and the coefficient of determination ($R^2$). Specifically, we compared the predicted values from our models with the latent person-parameter trait estimates from the self-reported values of the personality trait measures. Because the personality scores in our analyses already represent latent trait scores, correlation measures were not adjusted for the reliability of the personality trait scales (all attenuated). Thus, the absolute size of the correlations is limited by the reliability of the personality trait measures. Disattenuated correlation coefficients are provided in *SI Appendix*, Table S5. We computed performance measures within each fold of the cross-validation procedure and averaged across all outer resampling folds within a single prediction model (e.g., for extraversion). To determine whether a model was predictive at all, we carried out $t$ tests by comparing the $R^2$ measures of the random forest model with those of the baseline model. The $t$ tests were based on 10-times repeated 10-fold cross-validation and used a variance correction to specifically address the dependence structure of cross-validation experiments (64). All comparisons were adjusted for multiple comparisons ($n = 35$) via Holm correction. Significant prediction models ($\alpha = 0.05$) are marked in boldface type in Fig. 1.

In addition to measures of predictive performance, we used interpretable machine-learning techniques with significant models to gain insights into our models' inner workings. Specifically, we used permutation strategies to determine the unique contribution of the respective behavioral class and the importance of a class within the context of all other classes. These effects were also tested for significance (41) and adjusted for multiple comparisons.

To determine which of the behavioral classes was the most important overall for the prediction of Big Five personality traits, we performed an additional resampling analysis: 1) We created predictor sets with all possible combinations of subsets of the six behavioral classes ($2^6 = 64$); 2) we created 100 resampling folds of the complete dataset (10-times repeated 10-fold cross-validation); train and test data splits remained the same across all combinations); 3) for each of these combinations in all folds ($64 \times 100 = 6,400$), we fitted (on training data) and evaluated (on test data) models to predict each personality criterion (30 facets or 5 domains = 30 or 5 $R^2$ coefficients);

4) we averaged $R^2$ across all personality criteria, within each fold of a combination (100 mean $R^2$ values); and 5) we used two maximum-likelihood linear mixed models (domains vs. facets) with the mean $R^2$ as the outcome variable, the resampling iteration as the random factor (fold 1 to 100), and the behavioral classes (dummy encoded) as fixed factors. This procedure allowed us to determine the effects of each behavioral class on the average prediction performance across all personality trait dimensions. $P$ values in the linear mixed models were adjusted for multiple testing with the Holm method. All procedures were performed on domain and facet levels, separately. Further details about preprocessing, the modeling procedures, and the performance metrics are available in *SI Appendix* and in the project's repository (40).

***Software.*** Due to the high computational load of the machine-learning analyses, we parallelized the computations on the Linux Cluster of the LRZ-Supercomputing Center, in Garching, near Munich, Germany. For computations on the cluster, R-version 3.5.0 was used (65). We used R 3.5.2 for all other analyses (65). We used the fxtract package (66) for variable extraction from the raw data. Furthermore, we used the mlrCPO (67) and caret (68) packages for preprocessing. For machine learning we used the mlr (69), glmnet (70), iml (71), and ranger (72) packages.

***Open data and materials and additional resources.*** We provide the dataset and the code for variable extraction, preprocessing, and modeling in the project's repository (40). Raw data files cannot be provided (due to unsolved privacy implications); full reproducibility is possible for the analyses but not for preprocessing and variable extraction. In the repository, we link to the interactive project website where readers can find an exhaustive data dictionary, additional methodological descriptions, references, and results for all models in much greater detail. This paper is based on a preprint (73).

1. M. Settanni, D. Azucar, D. Marengo, Predicting individual characteristics from digital traces on social media: A meta-analysis. *Cyberpsychol. Behav. Soc. Netw.* **21**, 217–228 (2018).
2. D. J. Ozer, V. Benet-Martínez, Personality and the prediction of consequential outcomes. *Annu. Rev. Psychol.* **57**, 401–421 (2006).
3. B. W. Roberts, N. R. Kuncel, R. Shiner, A. Caspi, L. R. Goldberg, The power of personality: The comparative validity of personality traits, socioeconomic status, and cognitive ability for predicting important life outcomes. *Perspect. Psychol. Sci.* **2**, 313–45 (2007).
4. C. J. Soto, How replicable are links between personality traits and consequential life outcomes? The life outcomes of personality replication project. *Psychol. Sci.* **30**, 711–727 (2019).
5. S. C. Matz, M. Kosinski, G. Nave, D. J. Stillwell, Psychological targeting as an effective approach to digital mass persuasion. *Proc. Natl. Acad. Sci. U.S.A.* **114**, 12714–12719 (2017).
6. W. Youyou, M. Kosinski, D. Stillwell, Computer-based personality judgments are more accurate than those made by humans. *Proc. Natl. Acad. Sci. U.S.A.* **112**, 1036–1040 (2015).
7. International Telecommunication Union, Measuring the information society report 2018. *ITU Publ.* **1**, 2–18 (2018).
8. G. M. Harari *et al.*, Using smartphones to collect behavioral data in psychological science: Opportunities, practical considerations, and challenges. *Perspect. Psychol. Sci.* **11**, 838–854 (2016).
9. G. M. Harari, S. D. Gosling, R. Wang, A. T. Campbell, Capturing situational information with smartphones and mobile sensing methods. *Eur. J. Pers.* **29**, 509–511 (2015).
10. G. Miller, The smartphone psychology manifesto. *Perspect. Psychol. Sci.* **7**, 221–237 (2012).
11. S. Servia-Rodríguez *et al.*, "Mobile sensing at the service of mental well-being: A large-scale longitudinal study" in *26th International World Wide Web Conference, WWW 2017* (International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, Switzerland, 2017), pp. 103–112.
12. K. K. Rachuri *et al.*, "EmotionSense: A mobile phones based adaptive platform for experimental social psychology research" in *UbiComp'10 - Proceedings of the 2010 ACM Conference on Ubiquitous Computing* (Association for Computing Machinery, New York, NY, 2010), pp. 281–290.
13. S. Saeb, E. G. Lattie, S. M. Schueller, K. P. Kording, D. C. Mohr, The relationship between mobile phone location sensor data and depressive symptom severity. *PeerJ* **4**, e2537 (2016).
14. S. Thomée, Mobile phone use and mental health. A review of the research that takes a psychological perspective on exposure. *Int. J. Environ. Res. Publ. Health* **15**, 2692 (2018).
15. G. M. Harari, A process-oriented approach to respecting privacy in the context of mobile phone tracking. *Curr. Opin. Psychol.* **31**, 141–147 (2019).
16. C. M. Gray, Y. Kou, B. Battles, J. Hoggatt, A. L. Toombs, "The dark (patterns) side of ux design" in *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems, CHI '18* (Association for Computing Machinery, New York, NY, 2018), pp. 1–14.
17. S. Barocas, H. Nissenbaum, "On notice: The trouble with notice and consent" in *Proceedings of the Engaging Data Forum: The First International Forum on the Application and Management of Personal Electronic Information* (2009). https://ssrn.com/abstract=2567409. Accessed 7 July 2020.
18. A. P. Felt *et al.*, "Android permissions: User attention, comprehension, and behavior" in *Proceedings of the Eighth Symposium on Usable Privacy and Security, SOUPS '12* (Association for Computing Machinery, New York, NY, 2012).
19. P. Wijesekera *et al.*, "The feasibility of dynamically granted permissions: Aligning mobile privacy with user preferences" in *Proceedings - IEEE Symposium on Security and Privacy* (Institute of Electrical and Electronics Engineers Inc., New York, NY, 2017), pp. 1077–1093.
20. J. Reardon *et al.*, "50 ways to leak your data: An exploration of apps' circumvention of the android permissions system" in *28th USENIX Security Symposium (USENIX Security 19)* (USENIX Association, Santa Clara, CA, 2019), pp. 603–620.
21. J. Valentino-DeVries, N. Singer, M. Keller, A. Krolik, Your apps know where you were last night, and they're not keeping it secret. *New York Times*, 10 December 2018.
22. L. R. Goldberg, An alternative "description of personality": The big-five factor structure. *J. Pers. Soc. Psychol.* **59**, 1216–1229 (1990).
23. R. R. McCrae, O. P. John, An introduction to the five-factor model and its applications. *J. Pers.* **60**, 175–215 (1992).
24. B. De Raad, The big five personality factors: The psycholexical approach to personality (Hogrefe & Huber Publishers, Seattle, WA, 2000).
25. J. C. Loehlin, R. R. McCrae, P. T. Costa, O. P. John, Heritabilities of common and measure-specific components of the big five personality factors. *J. Res. Pers.* **32**, 431–453 (1998).
26. P. Costa Jr, A. Terracciano, R. R. McCrae, Gender differences in personality traits across cultures: Robust and surprising findings. *J. Pers. Soc. Psychol.* **81**, 322–331 (2001).
27. C. M. Ching *et al.*, The manifestation of traits in everyday behavior and affect: A five-culture study. *J. Res. Pers.* **48**, 1–16 (2014).

Stachl et al.

111

28. G. Chittaranjan, J. Blom, D. Gatica-Perez, Mining large-scale smartphone data for personality studies. *Personal Ubiquitous Comput.* **17**, 433–450 (2013).

29. Y. A. De Montjoye, J. Quoidbach, F. Robic, A. Pentland, "Predicting personality using novel mobile phone-based metrics" in *Proceedings of the 6th International Conference on Social Computing, Behavioral-Cultural Modeling and Prediction, SBP'13*, A. M. Greenberg, W. G. Kennedy, N. D. Bos, Eds. (Springer-Verlag, Berlin/Heidelberg, Germany, 2013), pp. 48–55.

30. W. Wang *et al.*, Sensing behavioral change over time: Using within-person variability features from mobile sensing to predict personality traits. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* **2**, 1–21 (2018).

31. B. Mønsted, A. Mollgaard, J. Mathiesen, Phone-based metric as a predictor for basic personality traits. *J. Res. Pers.* **74**, 16–22 (2018).

32. G. M. Harari *et al.*, Sensing sociability: Individual differences in young adults' conversation, calling, texting, and app use behaviors in daily life. *J. Pers. Soc. Psychol.*, 10.1037/pspp0000245 (2019).

33. C. Stachl *et al.*, Personality traits predict smartphone usage. *Eur. J. Pers.* **31**, 701–722 (2017).

34. C. Montag *et al.*, Smartphone usage in the 21st century: Who is active on WhatsApp? *BMC Res. Notes* **8**, 331 (2015).

35. C. Montag *et al.*, Correlating personality and actual phone usage: Evidence from psychoinformatics. *J. Indiv. Differ.* **35**, 158–165 (2014).

36. P. Ai, Y. Liu, X. Zhao, Big Five personality traits predict daily spatial behavior: Evidence from smartphone data. *Pers. Indiv. Differ.* **147**, 285–291 (2019).

37. L. Alessandretti, S. Lehmann, A. Baronchelli, Understanding the interplay between social and spatial behaviour. *EPJ Data Sci.* **7**, 36 (2018).

38. N. K. Kambham, K. G. Stanley, S. Bell, "Predicting personality traits using smartphone sensor data and app usage data" in *2018 IEEE 9th Annual Information Technology, Electronics and Mobile Communication Conference, IEMCON 2018* (IEEE, New York, NY, 2019), pp. 125–132.

39. N. Gao, W. Shao, F. D. Salim, Predicting personality traits from physical activity intensity. *Computer* **52**, 47–56 (2019).

40. C. Stachl *et al.* Repository: Predicting personality from patterns of behavior collected with smartphones. Open Science Framework. https://osf.io/kqjhr/. Deposited 18 June 2019.

41. A. Altmann, L. Toloşi, O. Sander, T. Lengauer, Permutation importance: A corrected feature importance measure. *Bioinformatics* **26**, 1340–1347 (2010).

42. D. C. Funder, D. J. Ozer, Evaluating effect size in psychological research: Sense and nonsense. *Adv. Methods Pract. Psychol. Sci.* **2**, 156–168 (2019).

43. B. Reeves, T. Robinson, N. Ram, Time for the human screenome project. *Nature* **577**, 314–317 (2020).

44. J. Henrich, S. J. Heine, A. Norenzayan, The weirdest people in the world? *Behav. Brain Sci.* **33**, 61–83 (2010).

45. P. M. Podsakoff, S. B. MacKenzie, N. P. Podsakoff, Sources of method bias in social science research and recommendations on how to control it. *Annu. Rev. Psychol.* **63**, 539–569 (2012).

46. Y. V. Vaerenbergh, T. D. Thomas, Response styles in survey research: A literature review of antecedents, consequences, and remedies. *Int. J. Publ. Opin. Res.* **25**, 195–217 (2013).

47. ICO, "Update report into adtech and real time bidding" (Tech. Rep., Information Commissioner's Office, UK, 2019). https://ico.org.uk/media/about-the-ico/documents/2615156/adtech-real-time-bidding-report-201906.pdf. Accessed 7 July 2020.

48. M. Kosinski, D. Stillwell, T. Graepel, Private traits and attributes are predictable from digital records of human behavior. *Proc. Natl. Acad. Sci. U.S.A.* **110**, 5802–5805 (2013).

49. A. Roets, I. Cornelis, A. Van Hiel, Openness as a predictor of political orientation and conventional and unconventional political activism in western and eastern Europe. *J. Pers. Assess.* **96**, 53–63 (2014).

50. R. L. Bach *et al.*, Predicting voting behavior using digital trace data. *Soc. Sci. Comput. Rev.*, 10.1177/0894439319882896 (2019).

51. European Parliament, REGULATION (EU) 2016/679 OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation, 2016).

52. Y. A. de Montjoye *et al.*, On the privacy-conscientious use of mobile phone data. *Sci. Data* **5**, 180286 (2018).

53. S. C. Matz, R. E. Appel, M. Kosinski, Privacy in the age of psychological targeting. *Curr. Opin. Psychol.* **31**, 116–121 (2020).

54. B. McMahan, E. Moore, D. Ramage, S. Hampson, B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data" in *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, Proceedings of Machine Learning Research*, A Singh, J Zhu, Eds. (PMLR, Fort Lauderdale, FL, 2017), **vol. 54**, pp. 1273–1282.

55. J. Hong, The privacy landscape of pervasive computing. *IEEE Pervasive Comput.* **16**, 40–48 (2017).

56. H. Almuhimedi *et al.*, "Your location has been shared 5,398 times! A field study on mobile app privacy nudging" in *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems, CHI '15* (Association for Computing Machinery, New York, NY, 2015), pp. 787–796.

57. C. Stachl *et al.*, Data from "The PhoneStudy project." Open Science Framework. https://osf.io/ut42y/. Accessed 7 July 2020.

58. R. Schoedel *et al.*, Digital footprints of sensation seeking. *Z. Psychol.* **226**, 232–245 (2018).

59. T. Schuwerk, L. J. Kaltefleiter, J. Q. Au, A. Hoesl, C. Stachl, Enter the wild: Autistic traits and their relationship to mentalizing and social interaction in everyday life. *J. Autism Dev. Disord.* **49**, 4193–4208 (2019).

60. M. Arendasy, *BFSI: Big-Five Struktur-Inventar (Test & Manual)* (Schuhfried GmbH, Mödling, Austria, 2009).

61. P. J. Huber, "Robust statistics" in *Wiley Series in Probability and Statistics* (John Wiley & Sons, Inc., 1981).

62. H. Zou, T. Hastie, Regularization and variable selection via the elastic net. *J. R. Stat. Soc. B Stat. Methodol.* **67**, 301–320 (2005).

63. L. Breiman, Random forests. *Mach. Learn.* **45**, 5–32 (2001).

64. R. R. Bouckaert, E. Frank, "Evaluating the replicability of significance tests for comparing learning algorithms" in *Advances in Knowledge Discovery and Data Mining*, H Dai, R Srikant, C Zhang, Eds. (Springer Berlin Heidelberg, Berlin/Heidelberg, Germany, 2004), pp. 3–12.

65. R Core Team, R Development Core Team, R: A Language and Environment for Statistical Computing (R Foundation for Statistical Computing, Vienna, 2018).

66. Q. Au, C. Stachl, R. Schoedel, T. Ullmann, A. Hofheinz, *fxtract: Feature Extraction from Grouped Data, R Package Version 0.9.2* (The Comprehensive R Archive Network, 2019).

67. M. Binder, *mlrCPO: Composable Preprocessing Operators and Pipelines for Machine Learning, R Package Version 0.3.4* (The Comprehensive R Archive Network, 2018).

68. M. Kuhn *et al.*, *caret: Classification and Regression Training, R package version 6.0-79* (The Comprehensive R Archive Network, 2018).

69. B. Bischl *et al.*, mlr: Machine learning in R. *J. Mach. Learn. Res.* **17**, 1–5 (2016).

70. J. Friedman, T. Hastie, R. Tibshirani, Regularization paths for generalized linear models via coordinate descent. *J. Stat. Softw.* **33**, 1–22 (2010).

71. C. Molnar, B. Bischl, G. Casalicchio, iml: An R package for interpretable machine learning. *JOSS* **3**, 786 (2018).

72. M. N. Wright, A. Ziegler, ranger: A fast implementation of random forests for high dimensional data in C++ and R. *J. Stat. Software* **77**, 1–17 (2017).

73. C. Stachl *et al.*, Behavioral patterns in smartphone usage predict big five personality traits. https://doi.org/10.31234/osf.io/ks4vd (12 June 2019).

PSYCHOLOGICAL AND COGNITIVE SCIENCES

112

# SI Correction

**PSYCHOLOGICAL AND COGNITIVE SCIENCES**

Correction to Supporting Information for "Predicting personality from patterns of behavior collected with smartphones," by Clemens Stachl, Quay Au, Ramona Schoedel, Samuel D. Gosling, Gabriella M. Harari, Daniel Buschek, Sarah Theres Völkel, Tobias Schuwerk, Michelle Oldemeier, Theresa Ullmann, Heinrich Hussmann, Bernd Bischl, and Markus Bühner, which was first published July 14, 2020; 10.1073/pnas.1920484117 (*Proc. Natl. Acad. Sci. U.S.A.* **117**, 17680–17687).

The authors note that, in the *SI Appendix*, page 3, first full paragraph, "After we enriched the music listening records with additional parameters from the Spotify API, we manually checked whether the retrieved music parameters were correctly matched to the listened artist-title-album triples" should instead appear as "After we enriched the music listening records with additional parameters from the Spotify API, we automatically checked whether the retrieved music parameters were correctly matched to the listened artist-title-album triples." The *SI Appendix* has been corrected online.

**SI CORRECTION**

# Digital Footprints of Sensation Seeking

This article explores the impact of new technologies on personality research, particularly focusing on the sensation seeking trait.

### Contributing article

Schoedel, R., Au, Q., Völkel, S. T., Lehmann, F., Becker, D., Bühner, M., Bischl, B., Hussmann, H., and Stachl, C. (2018). Digital footprints of sensation seeking. *Zeitschrift für Psychologie*, 226(4):232–245

This publication was part of the PhoneStudy project (see Chapter 7).

### Copyright information

This version of the article may not completely replicate the final authoritative version published in Hogrefe Publishing at `https://doi.org/10.1027/2151-2604/a000342`. It is not the version of record and is therefore not suitable for citation. The *accepted manuscript version* is allowed to be shared and posted at any time[1].

### Declaration of contributions

This manuscript greatly benefited from preliminary work carried out by the doctoral candidate on other projects within the PhoneStudy project. Jiew-Quay Au contributed by conducting the data preparation and calculation of important variables for the data analysis, where machine learning techniques were employed for predicting the personality trait Sensation-Seeking. The PhD candidate also performed a benchmark experiment to compare different learning methods and conducted a detailed analysis of the selected model.

#### Contribution of the coauthors

Ramona Schödel, as the lead author, developed the idea for the study, which aimed to predict the stable personality trait sensation-seeking using sensor data. She conducted the study, managed data collection, and wrote the manuscript. Sarah Theres Völkel led the app development team responsible for data acquisition. The other co-authors also contributed equally to the manuscript revision.

---

[1] `https://www.hogrefe.com/us/resources/publishing-with-hogrefe/for-journals/usage-guidelines-for-journal-articles`

Running head: DIGITAL FOOTPRINTS OF SENSATION SEEKING

Digital Footprints of Sensation Seeking: A Traditional Concept in the Big Data Era

Ramona Schoedel[1], Quay Au[2], Sarah Theres Völkel[3], Florian Lehmann[3], Daniela Becker[3],

Markus Bühner[1], Bernd Bischl[2], Heinrich Hussmann[3], & Clemens Stachl[1]

[1] Department of Psychology, Psychological Methods and Assessment,

Ludwig-Maximilians-Universität München

[2] Department for Statistics, Computational Statistics, Ludwig-Maximilians-Universität

München

[3] Institute of Informatics, Ludwig-Maximilians-Universität München

**Author note**

Correspondence concerning this article should be addressed to Ramona Schoedel,

Department of Psychology, Psychological Methods and Assessment,

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Ludwig-Maximilians-Universität München, 80802 Munich, Germany. E-mail:

Ramona.Schoedel.psy.lmu.de

DIGITAL FOOTPRINTS OF SENSATION SEEKING

## Abstract

The increasing usage of new technologies implies changes for personality research. First, human behavior gets measurable by digital data and second, replaces conventional behavior in the analog world. This offers the opportunity to investigate personality traits by means of digital footprints. In this context, the investigation of the personality trait sensation seeking attracted our attention as objective behavioral correlates have been missing so far. By collecting behavioral markers (e.g. communication or app usage) via Android smartphones, we examined whether self-reported sensation seeking scores can be reliably predicted. Overall 260 subjects participated in our 30-days data logging real-life study. Using a machine learning approach, we evaluated cross-validated model fit based on how accurate sensation seeking scores in unseen samples can be predicted. Our findings highlight the potential of mobile sensing techniques in personality research and show exemplarily how prediction approaches can help to foster understanding human behavior.

*Keywords:* sensation seeking, machine learning, big data, behaviour, smartphone sensing

Word count: 7666

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Digital Footprints of Sensation Seeking: A Traditional Concept in the Big Data Era

Only recently researchers have started to discover the potential of big data for research in Psychology. E. E. Chen and Wojcik (2016) pointed out that the rather theory-driven field of Psychology could benefit from an additional focus on big data methods such as prediction modeling (Yarkoni & Westfall, 2017). Inversely, Cheung and Jak (2016) highlighted that the discipline of Psychology traditionally aims to explain complex issues and consequently could help to develop a understanding of big data. Although, only a small number of studies has so far combined both approaches, an increasing potential for such studies exists. Nowadays, people produce vast amounts of user data throughout their daily life by the means of increased technology usage (E. E. Chen & Wojcik, 2016). Thereby, human behavior becomes more and more quantifiable in terms of data (e.g. mobility can be measured via GPS data (Harari et al., 2016)). Furthermore, according to Mayer-Schönberger and Cukier (2013), digital behavior even replaces formerly "analog" behavior (e.g. using gaming apps on a smartphone instead of playing a card game). Such digital footprints can be used for personality research as they offer the opportunity for traits to manifest in a new context and to investigate those manifestation in terms of daily usage behavior.

**The personality trait of sensation seeking**

Why do some people go skydiving, while others read detective stories to feel aroused? Systematic, individual differences in the need for external stimulation have been described as the personality trait sensation seeking (Zuckerman, 1994). Initially proposed by Zuckerman, it refers to "seeking of varied, novel, complex, and intense sensations and experiences, and the willingness to take physical, social, legal, and financial risks for the sake of such experience"

DIGITAL FOOTPRINTS OF SENSATION SEEKING

(Zuckerman, 1994, p. 27). The construct of sensation seeking has been defined from a

biopsychological personality perspective and is explained by genetic, biological,

psychophysiological, but also social factors (Roberti, 2004; Zuckerman, 1994). Accordingly, age

and sex were found to be related to sensation seeking, namely younger and male individuals

showed higher trait scores (Roberti, 2004). After reviewing the vast amout of existing studies on

sensation seeking, we have identified three key issues that provide room for new research.

First, the majority of studies has dealt with an unsocialized form of sensation seeking.

This term refers to actions like criminal behaviors, alcohol and substance usage, excessive

gambling, risky sexual activities, or reckless driving (Roberti, 2004). However, Zuckerman

(1994) also postulated the existence of a non-impulsive, socialized type of sensation seeking.

This type was described by individual characteristics such as being against conventionalism,

lacking planning skills (Glicksohn & Abulafia, 1998), and by an affinity for unfamiliar

international travel destinations (Lepp & Gibson, 2008).

Second, most previous studies have been focused on high risk activities including taking

financial risks (Zabel, Christopher, Marek, Wieth, & Carlson, 2009) or doing extreme sports

(Jack & Ronan, 1998). Thus, Guszkowska and Bołdak (2010) found that individual levels of

sensation seeking are positiveley related to practicing sports like parachuting, snowboarding, or

alpinism. However, according to Roberti (2004), sensation seeking is not limited to the seeking of

risks per se. Rather, a certain amount of risk is accepted to obtain an ideal level of arousal. In

contrast to research focusing on high risk activities, studies about everyday-expressions of

sensation seeking have been rare and have investigated, for example, the association between

sensation seeking and the need for social stimulation (Weisskirch & Murphy, 2004).

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Third, traditionally the collection of data about actual behavior has been very difficult and costly to achieve. Behavioral correlates of sensation seeking like reckless driving (Dahlen, Martin, Ragan, & Kuhlman, 2005) or smartphone usage (Leung, 2008) have almost exclusively been measured via retrospective self-reports. However, it is commonly known that self-report questionnaires are subject to a series of biases, such as memory and social desirability (M. Ziegler & Buehner, 2009). Accordingly, Baumeister, Vohs, and Funder (2007) argued that self-reported behavior can greatly differ from actual behavior and highlighted the necessity to investigate behavior directly. To summarize our three key points, previous studies have mainly focussed on self-reports of unsocialized, and high-risk related types of sensation seeking. This motivates our research effort to re-investigate the socialized and every-day expression of sensation seeking by using objective behavioral data collected via smartphone sensing.

**Smartphone sensing and automated trait recognition**

Within the last years, smartphone-sensing has established itself as an active area of research within the field of Psychology (Harari et al., 2016). An increasing number of consumer electronics are equipped with sensors capable of logging data about its user's natural mobility and everyday-activities, and habits. These developments enable researchers to develop applications (apps) to collect extensive records of individual behavior in an efficient and unobtrusive manner (Harari et al., 2016). Smartphone sensing seems especially promising for personality psychology, as more and more behaviors (e.g. shopping, listening to music, playing games) can be exerted via smartphones, reflecting potential dimensions of individual difference. Accordingly, a growing body of research has investigated associations of smartphone usage and individual traits. So far,there has been consensus that individual traits are related to smartphone

DIGITAL FOOTPRINTS OF SENSATION SEEKING

usage behavior in some way. Andone et al. (2016) reported that age and gender was systematically related to individual smartphone usage. Montag et al. (2015) reported associations of extraversion and conscientiousness with daily WhatsApp usage. Smartphone usage in a broader sense was examined by Stachl et al. (2017). They evaluated the predictive performance of personality traits, fluid intelligence and demographic variables for the frequency and duration of categorial app usage.

Beyond mere association, patterns in sensing data could also be used to directly predict individual trait levels. The idea of inferring states and traits from the everyday digital technology usage has recently gained importance in the field of Psychology. So far, studies have focused on the investigation of social network data (e.g. Kosinski, Stillwell, & Graepel, 2013 Youyou, Kosinski, and Stillwell (2015)).

Researchers from other fields have started to investigate the automatic inference of traits based on data collected via smartphones. Chittaranjan, Blom, and Gatica-Perez (2013) and Montjoye, Quoidbach, Robic, and Pentland (2013) used machine learning algorithms to predict Big Five traits based on smartphone logging data. Whereas Chittaranjan et al. (2013) focused on features derived from app, text message, and call logs, Montjoye et al. (2013) additionally included features based on location data. Despite their slightly different approaches, both Chittaranjan et al. (2013) and Montjoye et al. (2013) reported that their machine learning algorithms could predict personality traits above chance.

If successful, the automated recognition of trait variables from usage data could have impact on both the academic and industrial sector. First, predicted traits could be used in recommender systems to develop personalized services or interfaces (Brinkman & Fine, 2005; Tkalcic & Chen, 2015). Second, the recognition of pathological traits like depression could help

DIGITAL FOOTPRINTS OF SENSATION SEEKING

to develop smartphone-based prevention programs (Saeb et al., 2015). Third, Yarkoni and Westfall (2017) argued that prediction approaches could also help to understand and consequently explain systematic variations in human behavior. It might be promising to revisit theory-based findings with objective data within a machine learning framework to detect possible underlying mechanisms of individual differences in human behavior.

**Rationale**

The aim of this study was to investigate the traditional concept of sensation seeking as reflected in natural smartphone usage. We think that for observing objective behavioral manifestations of sensation seeking in everyday-contexts, appropriate investigation methods have been missing so far. We therefore combined smartphone sensing data with traditional self-report measures, to gain new insights into the behavioral manifestations of sensation seeking. Using a large number of literature-derived predictor variables, we evaluated whether individual sensation seeking scores can be reliably predicted from the data. Additionally, we compared the prediction performance of different machine learning algorithms and investigated the importance of single variables for the models. Moreover, we want to replicate the often reported finding that sensation seeking is related to age as well as gender.

**Method**

This study was pre-registered prior to analysing the data. The pre-registration form is available in our open science framework project (OSF; Schoedel, Au, Völkel, Bühner, & Stachl, 2018). Our data was collected within the framework of the larger, ongoing "PhoneStudy" project - an interdisciplinary research project between the chair of psychological assessment and the

DIGITAL FOOTPRINTS OF SENSATION SEEKING

working groups computational statistics as well as media informatics at

Ludwig-Maximilians-Universität München (LMU), Germany (see Stachl et al., 2018). The

present study obtained approval from the responsible institutional review board and data

protection office.

**Participants**

All participants were recruited by student researchers during a seminar. Participation

requirements included speaking German fluently as well as a minimum age of 18. For technical

reasons, only participants with smartphones running Android 4.4 or higher could participate in

the study. Initially, our dataset contained data entries from 361 participants. However, as defined

in our pre-registration, we only included participants with completed questionnaire data and at

least 15 days of logging data in our analyses. This resulted in a final sample size of $N = 260$

participants (68% women). Participants age ranged from 18 to 72 with an average age of 24 ($SD$

$= 8.84$). The sample was skewed towards younger and highly educated participants as

recruitment took mainly place in the university context. Accordingly, 73% of all participants had

a high school degree, 16% had a university degree.

**Data collection procedure**

After being informed about the study, the participants provided informed consent via an

online form. In the consequent 30-days data collection period, rich behaviorally-focused log-data

were collected on the participants` smartphones. Participants were instructed to answer a series of

self-report questionnaires integrated in the app at a time convenient for them during the study

period. The PhoneStudy research app enables unobstrusive data logging utilizing background

services to monitor smartphone usage and location tracking. For this study we focused on the

DIGITAL FOOTPRINTS OF SENSATION SEEKING

logging of app usage, phone calls, and GPS data. For privacy reasons we did not collect

content-related data (e.g. text or notification contents). App usage and phone calls were recorded

event based, location data time based every 15 minutes. Data was synchronized hourly, if users

were connected to WiFi. In the case of missing WiFi connectivity, synchronization was forced

using any available network connection after one week. The data was synchronized with a

backend server using SSL encryption. Data were stored in encypted form on the backend server

and secured via two-factor authentication. The entire data collection for this study took place

between October 2017 and January 2018.

**Measures**

 **Self-report measures.** In previous studies, a series of sensation seeking questionnaires

had been used. Although the 40-item Sensation Seeking Scale Form V (SSS-V; Zuckerman,

Eysenck, & Eysenck, 1978) was used in most studies, this scale shows weakness in terms of

psychometric properties and its factorial structure (Beauducel, Strobel, & Brocke, 2003). Thus,

we employed the Impulsive Sensation Seeking (ImpSS) subscale of the Zuckerman-Kuhlman

Personality Questionnaire (ZKPQ-III-R; Zuckerman, 2002) which represents a more reliable and

valid alternative (Roberti, 2004). Zuckerman (2002) reports good internal consistency

(Cronbach's $\alpha$ = 0.83 for a German subsample). The ImpSS consists of 19 items (e.g. "I am an

impulsive person") and participants are instructed to indicate if statements are either true or false.

The ImpSS is defined by two facets: impulsivity (eight items) and sensation seeking (eleven

items). According to Zuckerman and Aluja (2015), facets can be cumulated to one score due to

their joint biological basis. Therefore, we summed up the 19 individual item scores to one ImpSS

score (ranging between 0 and 19). For our sample we found Cronbach's $\alpha$ = 0.80, $CI_{95\%}$ [0.77,

DIGITAL FOOTPRINTS OF SENSATION SEEKING

0.84]. Moreover, participants were asked to indicate their demographics. In addition, participants completed the German version of the Big-Five-Structure-Inventory (BFSI; Arendasy, Sommer, & Feldhammer, 2011), the newer German version of the Big-Five-Inventory 2 (BFI-2; Danner et al., 2016), and the Smartphone-Addiction Scale (SAS; Kwon et al., 2013). As those questionnaires were used for additional research, not covered in this study, we will not continue to elaborate on it in this article.

### Behavioral measures and extracted features

. Originally, the data existed as timestamped event-data. Each row represented a registered event (e.g. call, app usage), each column an event characteristic (e.g. outgoing, timestamp, duration, contact-hash). Thus, before modeling, we pre-processed our dataset in order to create meaningful predictors (also called features in machine learning) for our models. The feature extraction was carried out with specifically created aggregation functions from an R-package, currently under development by the working group of computational statistics at LMU.

### *Identification and quantification of behavioral categories*

. Initially, we performed an extensive literature review to identify behaviors characteristic for sensation seeking. As we could not find research about sensation seeking and smartphone usage, we identified behavioral manifestations of sensation seeking from "traditional" literature and matched those to measures of possibly equivalent smartphone usage. For example, sensation seeking was commonly associated with gambling in previous studies (McDaniel, 2002). Consequently, we "translated" gambling into gaming app usage behavior. Afterwards, we quantified the literature-derived categories (e.g. gaming app usage) by following previous research investigating the relationship of smartphone usage and user characteristics (e.g.

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Chittaranjan, Blom, & Gatica-Perez, 2011; Montjoye et al., 2013; Stachl et al., 2017). Used

quantification measures were for example mean/variation of frequency and duration, entropy,

irregularity, ratio, or radius of gyration. For their detailed explanation see table A1 in the

Appendix. The complete feature list was pre-registered prior to data analyses and is available in

our OSF project.

### *Categorization of apps*

. In order to effectively analyze app usage data, we chose to categorize all used apps into a finite

number of categories. The Google Play store offers a categorization of apps (Google, 2018).

However, this categorization is based on the subjective labeling by app developers and might be

influenced by reasons like popularity of certain app categories. We therefore pre-defined our own

app categories relevant for our research question: gaming, dating, communication, social media,

listening to music/audio clips, watching video clips, planning and organising, traveling, trading,

browsing, shopping, reading news, personalizing the own smartphone, informing about risky

driving behavior, and apps related to running as well as to outdoor sporting activities. In order to

increase transparency of our categorization approach, we provide the full list of apps, assigned

labels and the definition of all categories in our OSF project.

In the course of data preprocessing, all apps were categorized manually by one coder who

read the descriptions provided in the Google Play store. A second coder checked the reliability of

these codings and ambigious cases were discussed with a third coder. Only apps available in the

Google Play store at the time of re-categorization (18.01.2018) were included. Background and

launcher apps were excluded, as they do not reflect intentional app usage behavior.

DIGITAL FOOTPRINTS OF SENSATION SEEKING

**Data pre-processing**

In order to prepare the data set for prediction modeling, we applied a series of pre-processing steps according to Kuhn and Johnson (2013) and Schiffner et al. (2016). We removed predictors with more than 90% missing values and predictors with zero or near-zero variance (10% cutoff). To avoid overfitting and to get a reliable estimate of the predictive performance on new data, the pre-processing steps transformation (scaling and centering) and imputation of missing values were performed within the respective inner resampling iterations. In our pre-registration, we planned to use a k-nearest neighbours algorithm for imputation. Due to software-related bugs, we had to use the median for imputation.

**Data analyses**

First, we aimed to replicate the often reported finding that impulsive sensation seeking is related to both age and gender (Roberti, 2004). To do so we calculated Bonferroni corrected pairwise Spearman correlations. In addition, we calculated simple pairwise correlations between impulsive sensation seeking scores and the self-reported Big Five personality scores. As suggested in previous literature (Yarkoni, 2010), we consistently used Spearman's correlation coefficients due to non-normally distributed data. Second, we computed descriptive statistics related to smartphone usage and app usage in particular. Third, we used a machine learning approach to predict self-reported sensation seeking scores from the features described in the method section.

**Machine learning algorithms**

. Within algorithmic modeling culture, it is assumed that there is no single best model (Wolpert & Macready, 1997). Rather, various models perform differently well, dependent on the unknown

DIGITAL FOOTPRINTS OF SENSATION SEEKING

true relationship between predictors and outcome. Therefore, we carried out a benchmark

experiment in which we compared the generalized predictive performance of different algorithms

(also called learners) against a common guessing baseline. This baseline is also called

"featureless learner" and constantly predicts the mean value of the training data's outcome value.

The learners we chose for the benchmark experiment represent various trade-off levels between

interpretability and expected prediction performance. First, we used an elastic net model (J.

Friedman, Hastie, & Tibshirani, 2010; Zou & Hastie, 2005). It is a linear regression method

applying a mixed L1-L2-regularization which allows to model linear relationships on

high-dimensional spaces. Furthermore, the L1-penalty drives irrelevant predictor variables out of

the model for model sparsity and therefore better interpretability. We chose the elastic net model

because it has often been proven to be competitive in contrast to non-linear methods and provides

well interpretable coefficients (Zou & Hastie, 2005). Second, we included a random forest

(Breiman, 2001; Wright & Ziegler, 2015) which is an ensemble technique of multiple

bootstrapped, decorrelated decision trees. The random forest as non-linear model is an

all-rounder, which can handle high-dimensional features spaces and small sample sizes usually

very well. Third, a support vector machine (SVM) with RBF kernel was used (Karatzoglou,

Smola, Hornik, & Zeileis, 2004; Vapnik, 1999). Through its kernel function, the SVM implicitly

maps the training observations into a high-dimensional feature space, where a linear decision

boundary is learnt. This results in a non-linear decision boundary in the original feature space.

We included the SVM because it is the most prevalent one used for personality prediction in

psychological research (Chittaranjan et al., 2013; Montjoye et al., 2013). Forth, we used extreme

gradient boosting (T. Chen, He, & Benesty, 2015; xgboost; J. H. Friedman, 2001). This method is

again an ensemble technique based on trees, which are combined via sequential gradient

DIGITAL FOOTPRINTS OF SENSATION SEEKING

boosting. Currently, xgboost is considered one of the most powerful prediction algorithm in the machine learning community.

**Evaluation metrics**

. We consider metrics that are typically used to measure the predictive performance of regression models: mean squared error ($MSE$), root mean squared error ($RMSE$), mean absolute deviation ($MAE$), and the coefficent of determination ($R^2$) (e.g. James, Witten, Hastie, & Tibshirani, 2013; Kuhn & Johnson, 2013). For the three metrics $MSE$, $RMSE$ and $MAE$ it is valid that lower values (approaching zero) indicate better model performance. The measure of $R^2$ is also referred to as the coefficient of determination. According to the conventional, in psychological research prevalent definition of $R^2$, its values range between 0 and 1, whereby the closer $R^2$ is to 1, the better the model explains the data. However, if model training and model evaluation happens on different datasets (e.g. in cross-validation), the mean of the response values between the training and the validation dataset can vary greatly and therefore, $R^2$ can become negative (Alexander, Tropsha, & Winkler, 2015). According to Alexander et al. (2015), negative values indicate that model fit is poor and that the number of observations is too small. As there is no consensus in literature which metric is superior to others, we follow Chai and Draxler (2014) and consider a combination of all metrics for model evaluation. In addition, we will present correlation coefficients between actual and predicted sensation seeking scores.

**Resampling procedure**

. For each learner the optimal choice of hyperparameters is data-dependent (Schiffner et al., 2016). To avoid overfitting, we applied a nested resampling strategy selecting optimal hyperparamters within inner resampling loops. The predictive performance of the tuned learners

DIGITAL FOOTPRINTS OF SENSATION SEEKING

is then evaluated within separate outer resampling loops. This ensures a strict separation of training and test data while allowing for the tuning of hyperparamters as well as pre-processing. More information about the detailed tuning procedure are included in the R code of the benchmark experiment which can be found as supplemental file in our OSF project.

For the inner resampling loops we used simple holdout validation for all learners. In the outer resampling loop, 10 times repeated 10-fold cross-validation was performed. Tuning for all learners was optimized on the *MSE* performance metric.

### Variable importance

. We selected the best prediction model with regard to the presented performance measures. To achieve a better understanding of which variables were important for prediction success, methods-inherent variable importance measures and partial dependence plots are presented (Schiffner et al., 2016). The plots help to explore the partial dependence of the trained function by selecting a subset of the predictor space. That means, the curves show how a trained function takes the values of features into consideration in order to predict sensation seeking scores.

### Statistical software

. Data processing and statistical analyses were performed with the statistical software R 3.4.3 (R Core Team, 2017). For pre-processing we used the *car* and *dplyr* packages (Fox & Weisberg, 2011; Wickham, Francois, Henry, & Müller, 2017). For modeling we used functions from the *mlr*, *caret*, and *psych* packages (Bischl et al., 2016; Revelle, 2017; Wing et al., 2017). See the Supplemental Information section for a link to a complete overview about all used R packages including version information.

DIGITAL FOOTPRINTS OF SENSATION SEEKING

## Results

**Descriptive statistics**

Our dataset contained 222 predictors before and 178 variables after pre-processing. An average, 1263 daily events were recorded for each participant. The participants used a total number of 2205 different apps during the course of the study. Table 1 shows information about the overall usage frequencies of app categories. The most frequently used app categories were related to communication, social media usage and browser usage. The number of different apps within one category was highest for gaming apps. Due to the scope of this article, summary statistics for all included variables are provided as a supplemental in our OSF project.

**Impulsive sensation seeking and demographics**

On average, participants reported an ImpSS score of $M = 7.91$ ($SD = 4.22$) which is in line with previous literature (e.g. Aluja, Garcia, & Garcia, 2003). Contrary to our assumptions, neither age ($r_s = -0.04$, $CI_{95\%}$ [-0.15, 0.07]), nor gender ($r_s = -0.02$, $CI_{95\%}$ [-0.15, 0.14]) were significantly related to sensation seeking.

**Prediction of individual Sensation Seeking scores**

**Benchmark results.** Table 2 presents the results of our benchmark experiment. The mean performance measures $MSE$, $RMSE$ and $MAE$ were lowest and $R^2$ was highest for the random forest compared to extreme gradient boosting, the support vector machine and the elastic net. The mean $MSE$ of the random forest was 10% lower, the mean $RMSE$ and $MAE$ were 5% lower than the guessing baseline. However, the dispersion of the $MSE$ and $R^2$ across all 100 iterations (see figure 1) shows that in some iterations, $R^2$ was negative for the random forest, indicating poor fit

DIGITAL FOOTPRINTS OF SENSATION SEEKING

(Alexander et al., 2015). Despite the relatively low mean $R^2$ (0.06) of the random forest model, we assume that the model grasped systematic variance in sensation seeking related behaviors. Both the constantly better performance measures, and a Pearson correlation of $r = 0.31$ between true and predicted test data, we consider the random forest provided predictions even if only slightly better than predicting by chance.

**Variable importance.**

***Permutation-based feature importance.*** To gain a better understanding of how the random forest model predicted new cases, we investigated the permutation-based feature importance measures for the top ten predictors. According to Breiman (2001), the idea behind is that first, the initial relation of one feature with the criterion variable is dissolved by randomly permutating the respective feature. Second, the permutated feature and all other remaining (unchanged) features are used to predict the criterion. The variable importance measure is the result of taking into account the difference in the prediction performance before and after permuting the respective feature. The larger the reduction in the prediction performance is, the stronger is the initial relation between the particular feature and the criterion variable and consequently, the more important is this respective feature in the model (Breiman, 2001). Table 3 displays the top ten predictors with the highest permutation-based feature importance (Wright & Ziegler, 2017). To illustrate the prediction direction of features, we added pairwise Spearman correlations between predictors and sensation seeking to the table.

The list suggests that the top ten features for predicting sensation seeking belonged to two primary categories: calling and day/night time activity. Calling activity included outgoing and missed calls, represented via different quantification metrics. For example, the random forest

DIGITAL FOOTPRINTS OF SENSATION SEEKING

judged participants as higher sensation seekers if they initialized or missed calls more often. In addition, the entropy of calling turned out to be important. Spearman coefficients suggest positive relationships between entropy of contacts-related variables and sensation seeking scores. Another set of important features for the random forest was related to individual day and night time activity. Spearman correlations suggest positive associations between night time activity indicators and sensation seeking levels. As an illustration, predicted sensation scores were higher, if the average point of time of the last smartphone usage on Friday/Saturday or on Sunday was late, and if the mean range of motion was high during night at weekends.

Some of our features were highly inter-correlated. Multicollinearity has been proven not to be an issue for the predictive performance of the random forest (and other machine learning algorithms), but to be likely to bias variable importance measures (James et al., 2013; Strobl, Boulesteix, Kneib, Augustin, & Zeileis, 2008). We therefore conducted an additional analysis using the conditional forest which is a learner taking into account the correlated structure of features (Strobl et al., 2008). As neither prediction performance, nor variable importance measures considerably differed between the conditional versus the random forest, and not to go beyond the scope of this article, we only report the results for the random forest here. However, corresponding additional analysis including detailed background information can be found as supplemental file in our OSF project.

### *Partial dependence plots*

. In addition to feature importance values, partial dependence plots can help to better understand how values of individual features on average influenced the prediction model (see figure 2). The curves show how predicted sensation seeking scores (y-axis) changed with regard to values of the respective predictor variable (x-axis).

DIGITAL FOOTPRINTS OF SENSATION SEEKING

In the top left of figure 2 the *mean frequency of missed calls per day* is plotted against sensation seeking scores. The plot shows that the average frequency of missed calls per day led to an increase in predicted sensation seeking scores for very low frequency values, but did not change noticeably if a mean value of about 0.4 missed calls per day was exceeded.

At the top right of figure 2 a partial dependence plot for *entropy of contacts for outgoing calls* is visible. Increasing values in contact-entropy on average resulted in higher predicted sensation seeking scores. This increase got sharper with rising entropy values.

As shown in the bottom left of 2, with a rising *mean number of intended events on Fridays/Saturdays nights* predicted sensation seeking scores first slowly and from a value of about 14 intended events sharply increased. Events were counted as "*intended*" when they were carried out intentionally by the participant.

The curve in the bottom right of 2 displays, that the *mean time of the last event on Sunday* on average led to higher predicted sensation seeking scores, when they occured after around 11pm.

Running head: DIGITAL FOOTPRINTS OF SENSATION SEEKING

## Discussion

The results of the present study indicate that individual scores of self-reported sensation seeking can be predicted from digital records of smartphone behavior above chance. Precisely, our results suggest that variables related to calling behavior as well as day and night time activity were particularly predictive for individual sensation seeking scores. In the following subsections we will critically discuss the results within the context of the used machine learning approach and will try to give some post-hoc explanations for important variables in the model. Please note that those interpretations are partially drawn post-hoc and should therefore not be easily generalized.

### A timely approach to a traditional concept

In contrast to all previous studies on sensation seeking, we used data about actual behavior to predict sensation seeking with a machine learning approach. Thus, we compared different statistical models based on their ability to accurately predict sensation seeking scores from unseen data.

Despite a relatively low overall acurracy, our results suggest that the flexible, non-linear random forest model outperformed all the other models. This suggests that resesarch in the field of individual differences might benefit from the additional usage of flexible models for the investigation of behavior-trait relationships (Benson & Campbell, 2007).

Previous studies investigating the automatic inference of personality traits based on features extracted from smartphone logs also reported prediction performances, ranging in size from six to 25 percent points depending on respective traits and used pre-processing procedures (Chittaranjan et al., 2011, 2013; Montjoye et al., 2013). Contrarily, those studies used

DIGITAL FOOTPRINTS OF SENSATION SEEKING

classification approaches and binned participants in low, average, and high on the Big Five

personality traits. Although this hinders a direct comparison with those studies, we argue that the

prediction of continous trait-scores can be more adequately modelled with a regression approach.

Related to this, the findings of Kosinski et al. (2013) offer a possibility to put the present results

into perspective. Kosinski et al. (2013) predicted personality traits from digital footprints

(Facebook likes) and reported correlations between predicted and actually observed trait values in

the range of 0.29 for conscientiousness and 0.43 for openness. Despite very different sample

sizes, our analyses produced coefficients in a similar range, suggesting comparable prediction

performance. Please note that the prediction accuracies of Kosinski et al. (2013) were exceeded

in a later study (Youyou et al., 2015).

Although the obtained prediction performance is comparable to previous studies, one

question still remains: what is the meaning of being 10% better than guessing? First, we want to

point out that in psychological research it is often investigated how well a model fitted a given

dataset (e.g. by considering in-sample $R^2$) and therefore, how trustworthy it is. In contrast, in the

context of prediction modeling, "good" and "trustworthy" are independent criteria. Good model

fit refers to how well new, unseen cases can be predicted with a trained model and "trustworthy"

indicates the correct application of methodological procedures.

Consequently, with regard to the question how "good" our model fit is, it has to be

considered that the prediction performance in the large majority of our folds was above the

guessing baseline, indicating that something more than randomness was going on. However, the

overall prediction performance was low. Reasons for this could be the relatively small sample

size or that self-reported sensation seeking scores cannot be perfectly predicted from the

behavioral indicators used in our study. We carefully selected these indicators by identifying

DIGITAL FOOTPRINTS OF SENSATION SEEKING

manifestations of sensation seeking assessed via self-reports in previous literature. Though, they were not reflected in objective behavioral data to the extent we would have expected from previous research. This could in turn suggest that the theoretical conceptualization of sensation seeking might benefit from additional research efforts in future studies. To sum up, we argue that our results are very well trustworthy in the sense that they indicate that sensation seeking cannot be predicted very accurately from the used behavioral predictors.

In addition, we think that how much better the prediction performance of a learner compared to a guessing baseline has to be, is a context-related question. With regard to practical applications (e.g. mobile computing) our model is certainly far away from being good and therefore applicable. However, in the context of psychological research effect sizes are usually very small and therefore, we would argue that our obtained mean (out of sample) $R^2$ is not unusually small.

**The trait sensation seeking and its correlates**

Beyond the evaluation of prediction performance, our analyses provided more detailed insights into the behavioral correlates of sensation seeking. Following previous studies, we hypothesized that both age and gender are related to sensation seeking (Roberti, 2004). However, associations of demographics with sensation seeking were not present in our data. We suspect that the absence of those effects could be related to our sample characteristics (predominately young females). Although previous studies reported similar gender ratios (Roberti, 2004), age ranges were larger. Possibly, but it can only by suspected, gender and age differences in socialized forms of sensation seeking might also not be as pronounced as in unsocialized forms (e.g. risky driving). However, those post-hoc explanations should be tested in future studies.

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Although machine learning algorithms are often labelled as "black-box models", they can provide additional information beyond prediction performance. In our study, the inspection of variable importance measures suggested that variables related to calling as well as day and night time activity were particularly important for predicting sensation seeking scores.

The variables of our prediction model were in advance derived from a literature review. For example, we "translated" the finding that sensation seeking is positiveley associated with a self-reported preference for social contexts (Roberti, 2004) into predictors related to calling activity. Variables regarding day and night time activity were based on findings about the relation between self-reported preferences for later bedtimes and sensation seeking (Tonetti, Fabbri, & Natale, 2009). The reviewed literature was exclusively based on self-reported behavioral correlates. But our model showed that these variables also turned out to be important for the prediction of sensation seeking when they are operationalized by means of behavioral data collected via smartphones. We think that our results can partially help to underpin questionnaire-based research with objective behavioral data.

Additionally, our prediction modeling approach provides new insights into behavioral manifestations of sensation seeking. Although, our study cannot raise any claims of causality or explanation, it can foster the postulation of new hypotheses for future studies. For example, two of the three most important features for the random forest's sensation seeking prediction were related to missed calls. As a mental game, one could deduce the hypothesis that people scoring high in sensation seeking are very active and busy in their everyday-lifes and therefore miss incoming calls. Such hypotheses can be tested in futures studies and aid the understanding of behavioral expressions of sensation seeking. To take up the current debate whether novel prediction-focused approaches are contradictory to the explanatory goal of Psychology (Yarkoni

DIGITAL FOOTPRINTS OF SENSATION SEEKING

& Westfall, 2017), we argue that our prediction framework suggests otherwise. Following

Yarkoni and Westfall (2017), the present study highlights that prediction approaches can help to

better understand the structure of objective behavioral data and can help to generate new

hypotheses for confirmatory research.

As discussed in the previous two subsections, our study illustrates how Psychology and

big data can work together. Psychological theories and findings can help to understand and

interpret what machine learning algorithms do (Cheung & Jak, 2016). But conversely, prediction

models could help to understand basic structures in complex behaviors (E. E. Chen & Wojcik,

2016). At this point we want to emphasize that our analyses cannot be considered big data due to

the relatively small sample size. However, our data collection tool, the "PhoneStudy app", with

its vast variety of collected variables as well as the methods used in this study hopefully highlight

some potential of the big data-approach in psychological research and inspire future work.

**Limitations & Outlook**

The present study has some limitations which we will discuss below. The categorization of

apps holds the problem that they can be used ambiguously. Hence, an app can be used to fulfill

different purposes and needs. For example, browser apps can be used to do online shopping, to

read news, to visit social media channels and so on. As the PhoneStudy app only provides

meta-data, we only know that participants used certain apps belonging to pre-defined categories.

However, it remains unclear what participants used the app for. Accordingly, we think that for

improving the prediction performance of machine learning algorithms, the inclusion of

content-related logging data such as user preferences (e.g. genres of listened music), browsing

histories, or notification texts might be a promising strategy. Although more fine-grained data

DIGITAL FOOTPRINTS OF SENSATION SEEKING

will likely improve trait prediction accuracy, the protection of individual privacy rights must be prioritized.

Our sample was primarily collected in the university context. Thus, younger and higher educated participants were overrepresented in our sample. Accordingly, some of our literature-derived features (usage of counteracting risky driving apps) were automatically excluded in the pre-processing as only single participants used respective apps. As our sample mainly consisted of younger people, car ownership might be systematically underrepresented. A more representative sample (including elders) could therefore provide more variance in behaviors related to sensation seeking.

Furthermore, machine learning algorithms only perform really well with large samples. Relatedly, the negative range of $R^2$ values of the featureless learner could indicate that our sample size was too small. This study should be replicated with a larger sample size (maybe ten times), to fully benefit from the predictive capabilities of those methods.

Finally, as already stated by Chittaranjan et al. (2013) and Kosinski et al. (2013), personality trait prediction is a challenging task. Traits are defined as latent constructs and can only be measured roughly, via self-report questionnaires. Therefore, prediction efforts using self-reported trait scores as ground truth, can only achieve accuracies that mimic those of self-report questionnaires. As we know that self-report questionnaires are also affected by a series of biases, this problem needs to be adressed eventually. Nevertheless, trait prediction can be improved in many ways. As the biopsychological trait sensation seeking was found to be related to individually as optimally considered levels of arousal (Roberti, 2004), physiological thresholds might be meaningful indicators. It might be helpful to include measures reflecting physiological processes in prediction models of sensation seeking. Measures of heart rate and electro-thermal

DIGITAL FOOTPRINTS OF SENSATION SEEKING

activity could be provided by wearables. Even though we are aware that the performance of our prediction model has to be higher to reach practical importance (e.g. for mobile computing), we think that our study can be a starting point for future research. Accordingly, it is one of the first studies working with such a broad variety of data collected via smartphone sensing in the field of trait prediction, and especially in the context of sensation seeking.

## Conclusion

The present study combined smartphone sensing data with traditional self-report measures, to gain new insights into behavioral manifestations of sensation seeking. The present study shows that self-reported sensation seeking scores can be predicted by smartphone logging data above the level of chance. Despite, limited predicition accuracies our results highlight novel behavioral indicators of sensation seeking and the potential of big data for psychological research.

## Acknowledgements

We want to thank all the student researchers for supporting us with recruiting participants, *Henrike Haase* for her persistence in categorizing apps, *Florian Pargent* for his insightful modeling advices, and the PhoneStudy team (*Daniel Buschek*, *Marius Herget*, *Ferdinand Hof*, *Peter Ehrich*, *Miriam Metz*, *Theresa Ullmann*) for making this kind of research possible.

DIGITAL FOOTPRINTS OF SENSATION SEEKING

**Supplemental Information**

All supplemental files are now accessible via an open science framework project link (and openly after publication): https://osf.io/v4xrf/?view_only=e868921365364fd88786b742edd6b459

*data.csv*: contains the data set with aggregated features used for prediction modeling.

*benchmark.R*: contains the R code for reproducing the reported results.

*features.pdf*: lists our features derived from a literature review.

*app_categories.csv*: contains our categorization of apps and their definition.

*summary_descriptives.pdf*: contains descriptive statistics for all variables.

*packages.pdf*: lists all used R packages including version information.

*cforest_analyses.pdf*: contains additional analyses regarding the conditional forest learner.

DIGITAL FOOTPRINTS OF SENSATION SEEKING

## References

Alexander, D., Tropsha, A., & Winkler, D. A. (2015). Beware of r 2: Simple, unambiguous assessment of the prediction accuracy of qsar and qspr models. *Journal of Chemical Information and Modeling*, *55*(7), 1316–1322. doi:10.1021/acs.jcim.5b00206

Aluja, A., Garcia, Ó., & Garcia, L. F. (2003). Psychometric properties of the zuckerman–Kuhlman personality questionnaire (zkpq-iii-r): A study of a shortened form. *Personality and Individual Differences*, *34*(7), 1083–1097. doi:10.1016/S0191-8869(02)00097-1

Andone, I., Błaszkiewicz, K., Eibes, M., Trendafilov, B., Montag, C., & Markowetz, A. (2016). How age and gender affect smartphone usage. In *Proceedings of the 2016 acm international joint conference on pervasive and ubiquitous computing: Adjunct* (pp. 9–12). ACM. doi:10.1145/2968219.2971451

Arendasy, M., Sommer, M., & Feldhammer, M. (2011). Manual big-five structure inventory bfsi. *Schuhfried Gmbh, Mödling*.

Baumeister, R. F., Vohs, K. D., & Funder, D. C. (2007). Psychology as the science of self-reports and finger movements: Whatever happened to actual behavior? *Perspectives on Psychological Science*, *2*(4), 396–403. doi:10.1111/j.1745-6916.2007.00051.x

Beauducel, A., Strobel, A., & Brocke, B. (2003). Psychometrische eigenschaften und normen einer deutschsprachigen fassung der sensation seeking-skalen, form v. *Diagnostica*. doi:10.1026//0012-1924.49.2.61

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Benson, M. J., & Campbell, J. P. (2007). To be, or not to be, linear: An expanded representation of personality and its relationship to leadership performance. *International Journal of Selection and Assessment*, *15*(2), 232–249. doi:10.1111/j.1468-2389.2007.00384.x

Bischl, B., Lang, M., Kotthoff, L., Schiffner, J., Richter, J., Studerus, E., … Jones, Z. M. (2016). mlr: Machine learning in r. *Journal of Machine Learning Research*, *17*(170), 1–5. Retrieved from http://jmlr.org/papers/v17/15-066.html

Breiman, L. (2001). Random forests. *Machine Learning*, *45*(1), 5–32. doi:10.1023/A:1010933404324

Brinkman, W.-P., & Fine, N. (2005). Towards customized emotional design: An explorative study of user personality and user interface skin preferences. In *Proceedings of the 2005 annual conference on european association of cognitive ergonomics* (pp. 107–114). University of Athens.

Canzian, L., & Musolesi, M. (2015). Trajectories of depression: Unobtrusive monitoring of depressive states by means of smartphone mobility traces analysis. In *Proceedings of the 2015 acm international joint conference on pervasive and ubiquitous computing* (pp. 1293–1304). ACM. doi:10.1145/2750858.2805845

Chai, T., & Draxler, R. R. (2014). Root mean square error (rmse) or mean absolute error (mae)?–Arguments against avoiding rmse in the literature. *Geoscientific Model Development*, *7*(3), 1247–1250. doi:10.5194/gmd-7-1247-2014

Chen, E. E., & Wojcik, S. P. (2016). A practical guide to big data research in psychology. *Psychological Methods*, *21*(4), 458–474. doi:10.1037/met0000111

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Chen, T., He, T., & Benesty, M. (2015). Xgboost: Extreme gradient boosting. *R Package Version 0.4-2*, 1–4.

Cheung, M. W.-L., & Jak, S. (2016). Analyzing big data in psychology: A split/analyze/meta-analyze approach. *Frontiers in Psychology*, *7*. doi:10.3389/fpsyg.2016.00738

Chittaranjan, G., Blom, J., & Gatica-Perez, D. (2011). Who's who with big-five: Analyzing and classifying personality traits with smartphones. In *Wearable computers (iswc), 2011 15th annual international symposium on* (pp. 29–36). IEEE. doi:10.1109/ISWC.2011.29

Chittaranjan, G., Blom, J., & Gatica-Perez, D. (2013). Mining large-scale smartphone data for personality studies. *Personal and Ubiquitous Computing*, *17*(3), 433–450. doi:10.1007/s00779-011-0490-1

Dahlen, E. R., Martin, R. C., Ragan, K., & Kuhlman, M. M. (2005). Driving anger, sensation seeking, impulsiveness, and boredom proneness in the prediction of unsafe driving. *Accident Analysis & Prevention*, *37*(2), 341–348. doi:10.1016/j.aap.2004.10.006

Danner, D., Rammstedt, B., Bluemke, M., Treiber, L., Berres, S., Soto, C., & John, O. (2016, September). Die deutsche version des big five inventory 2 (bfi-2). doi:10.6102/zis247

Fox, J., & Weisberg, S. (2011). *An R companion to applied regression* (Second.). Thousand Oaks CA: Sage. Retrieved from http://socserv.socsci.mcmaster.ca/jfox/Books/Companion

Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 1189–1232. doi:10.1214/aos/1013203451

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Friedman, J., Hastie, T., & Tibshirani, R. (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, *33*(1), 1.

Glicksohn, J., & Abulafia, J. (1998). Embedding sensation seeking within the big three. *Personality and Individual Differences*, *25*(6), 1085–1099. doi:10.1016/S0191-8869(98)00096-8

Guszkowska, M., & Bołdak, A. (2010). Sensation seeking in males involved in recreational high risk sports. *Biology of Sport*, *27*(3). doi:10.5604/20831862.919331

Harari, G. M., Lane, N. D., Wang, R., Crosier, B. S., Campbell, A. T., & Gosling, S. D. (2016). Using smartphones to collect behavioral data in psychological science: Opportunities, practical considerations, and challenges. *Perspectives on Psychological Science*, *11*(6), 838–854. doi:10.1177/1745691616650285

Jack, S., & Ronan, K. R. (1998). Sensation seeking among high-and low-risk sports participants. *Personality and Individual Differences*, *25*(6), 1063–1083. doi:10.1016/S0191-8869(98)00081-6

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning* (Vol. 112). Springer. doi:10.1007/978-1-4614-7138-7

Kafadar, K. (2003). John tukey and robustness. *Statistical Science*, *18*(3), 319–331.

Karatzoglou, A., Smola, A., Hornik, K., & Zeileis, A. (2004). Kernlab-an s4 package for kernel methods in r. *Journal of Statistical Software*, *11*(9), 1–20. doi: 10.18637/jss.v011.i09

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Kosinski, M., Stillwell, D., & Graepel, T. (2013). Private traits and attributes are predictable from digital records of human behavior. *Proceedings of the National Academy of Sciences*, *110*(15), 5802–5805. doi:10.1073/pnas.1218772110

Kuhn, M., & Johnson, K. (2013). *Applied predictive modeling* (Vol. 26). Springer. doi:10.1007/978-1-4614-6849-3

Kwon, M., Lee, J.-Y., Won, W.-Y., Park, J.-W., Min, J.-A., Hahn, C., … Kim, D.-J. (2013). Development and validation of a smartphone addiction scale (sas). *PloS One*, *8*(2), e56936. doi:10.1371/journal.pone.0056936

Lepp, A., & Gibson, H. (2008). Sensation seeking and tourism: Tourist role, perception of risk and destination choice. *Tourism Management*, *29*(4), 740–750. doi:10.1016/j.tourman.2007.08.002

Leung, L. (2008). Leisure boredom, sensation seeking, self-esteem, and addiction. *Mediated Interpersonal Communication*, *1*, 359–381.

Mayer-Schönberger, V., & Cukier, K. (2013). *Big data: A revolution that will transform how we live, work, and think*. Houghton Mifflin Harcourt.

McDaniel, S. R. (2002). Investigating the roles of gambling interest and impulsive sensation seeking on consumer enjoyment of promotional games. *Social Behavior and Personality: An International Journal*, *30*(1), 53–64. doi:10.2224/sbp.2002.30.1.53

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Montag, C., Błaszkiewicz, K., Sariyska, R., Lachmann, B., Andone, I., Trendafilov, B., …
Markowetz, A. (2015). Smartphone usage in the 21st century: Who is active on whatsapp?
*BMC Research Notes*, *8*(1), 331. doi:10.1186/s13104-015-1280-z

Montjoye, Y.-A. de, Quoidbach, J., Robic, F., & Pentland, A. S. (2013). Predicting personality
using novel mobile phone-based metrics. In *International conference on social computing,
behavioral-cultural modeling, and prediction* (pp. 48–55). Springer.
doi:10.1007/978-3-642-37210-0_6

R Core Team. (2017). *R: A language and environment for statistical computing*. Vienna, Austria:
R Foundation for Statistical Computing. Retrieved from https://www.R-project.org/

Revelle, W. (2017). *Psych: Procedures for psychological, psychometric, and personality
research*. Evanston, Illinois: Northwestern University. Retrieved from
https://CRAN.R-project.org/package=psych

Roberti, J. W. (2004). A review of behavioral and biological correlates of sensation seeking.
*Journal of Research in Personality*, *38*(3), 256–279. doi:10.1016/S0092-6566(03)00067-9

Rousseeuw, P. J., & Croux, C. (1993). Alternatives to the median absolute deviation. *Journal of
the American Statistical Association*, *88*(424), 1273–1283. doi:10.2307/2291267

Saeb, S., Zhang, M., Karr, C. J., Schueller, S. M., Corden, M. E., Kording, K. P., & Mohr, D. C.
(2015). Mobile Phone Sensor Correlates of Depressive Symptom Severity in Daily-Life
Behavior: An Exploratory Study. *Journal of Medical Internet Research*, *17*(7), e175.
doi:10.2196/jmir.4273

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Schiffner, J., Bischl, B., Lang, M., Richter, J., M Jones, Z., Probst, P., … Kotthoff, L. (2016). Mlr tutorial.

Schoedel, R., Au, Q., Völkel, S., Bühner, M., & Stachl, C. (2018, March). Digital footprints of sensation seeking: A traditional concept in the big data era. Open Science Framework. doi:10.17605/OSF.IO/V4XRF

Stachl, C., Hilbert, S., Au, J.-Q., Buschek, D., De Luca, A., Bischl, B., … Bühner, M. (2017). Personality traits predict smartphone usage. *European Journal of Personality*, *31*(6), 701–722. doi:10.1002/per.2113

Stachl, C., Schoedel, R., Au, Q., Völkel, S., Buschek, D., Hussmann, H., … Bühner, M. (2018, March). The phonestudy project. Open Science Framework. doi:10.17605/OSF.IO/UT42Y

Strobl, C., Boulesteix, A.-L., Kneib, T., Augustin, T., & Zeileis, A. (2008). Conditional variable importance for random forests. *BMC Bioinformatics*, *9*(1), 307. doi:10.1186/1471-2105-9-307

Tkalcic, M., & Chen, L. (2015). Personality and recommender systems. In *Recommender systems handbook* (pp. 715–739). Springer. doi:10.1007/978-1-4899-7637-6_21

Tonetti, L., Fabbri, M., & Natale, V. (2009). Relationship between circadian typology and big five personality domains. *Chronobiology International*, *26*(2), 337–347. doi:10.1080/07420520902750995

Vapnik, V. N. (1999). An overview of statistical learning theory. *IEEE Transactions on Neural Networks*, *10*(5), 988–999. doi:10.1109/72.788640

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Weisskirch, R. S., & Murphy, L. C. (2004). Friends, porn, and punk: Sensation seeking in

personal relationships, internet activities, and music preference among college students.

*Adolescence*, *39*(154), 189.

Wickham, H., Francois, R., Henry, L., & Müller, K. (2017). *Dplyr: A grammar of data

manipulation*. Retrieved from https://CRAN.R-project.org/package=dplyr

Williams, M. J., Whitaker, R. M., & Allen, S. M. (2012). Measuring individual regularity in

human visiting patterns. In *Privacy, security, risk and trust (passat), 2012 international

conference on and 2012 international confernece on social computing (socialcom)* (pp.

117–122). IEEE. doi:10.1109/SocialCom-PASSAT.2012.93

Wing, M. K. C. from, Weston, S., Williams, A., Keefer, C., Engelhardt, A., Cooper, T., … Hunt.,

T. (2017). *Caret: Classification and regression training*. Retrieved from

https://CRAN.R-project.org/package=caret

Wolpert, D. H., & Macready, W. G. (1997). No free lunch theorems for optimization. *IEEE

Transactions on Evolutionary Computation*, *1*(1), 67–82. doi:1089-778X(97)03422-X.

Wright, M. N., & Ziegler, A. (2015). Ranger: A fast implementation of random forests for high

dimensional data in c++ and r. *arXiv Preprint arXiv:1508.04409*.

Wright, M. N., & Ziegler, A. (2017). ranger: A fast implementation of random forests for high

dimensional data in C++ and R. *Journal of Statistical Software*, *77*(1), 1–17.

doi:10.18637/jss.v077.i01

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Yarkoni, T. (2010). Personality in 100,000 words: A large-scale analysis of personality and word use among bloggers. *Journal of Research in Personality*, *44*(3), 363–373. doi:10.1016/j.jrp.2010.04.001

Yarkoni, T., & Westfall, J. (2017). Choosing prediction over explanation in psychology: Lessons from machine learning. *Perspectives on Psychological Science*, *12*(6), 1100–1122. doi:10.1177/1745691617693393

Youyou, W., Kosinski, M., & Stillwell, D. (2015). Computer-based personality judgments are more accurate than those made by humans. *Proceedings of the National Academy of Sciences*, *112*(4), 1036–1040. doi:10.1073/pnas.1418680112

Zabel, K. L., Christopher, A. N., Marek, P., Wieth, M. B., & Carlson, J. J. (2009). Mediational effects of sensation seeking on the age and financial risk-taking relationship. *Personality and Individual Differences*, *47*(8), 917–921. doi:10.1016/j.paid.2009.07.016

Ziegler, M., & Buehner, M. (2009). Modeling socially desirable responding and its effects. *Educational and Psychological Measurement*, *69*(4), 548–565. doi:10.1177/0013164408324469

Zou, H., & Hastie, T. (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, *67*(2), 301–320. doi:1369–7412/05/67301

Zuckerman, M. (1994). *Behavioral expressions and biosocial bases of sensation seeking*. Cambridge university press.

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Zuckerman, M. (2002). Zuckerman-kuhlman personality questionnaire (zkpq): An alternative five-factorial model. *Big Five Assessment*, 377–396.

Zuckerman, M., & Aluja, A. (2015). Measures of sensation seeking. In *Measures of personality and social psychological constructs* (pp. 352–380). Elsevier. doi:10.1016/B978-0-12-386915-9.00013-9

Zuckerman, M., Eysenck, S. B., & Eysenck, H. J. (1978). Sensation seeking in england and america: Cross-cultural, age, and sex comparisons. *Journal of Consulting and Clinical Psychology*, *46*(1), 139. doi:10.1037//0022-006X.46.1.139

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Table 1

*Summary of mean performance measures of the 10x10 CV benchmark experiment*

| Measures | FL | RF | XG | SVM | EN |
|----------|------|------|------|------|------|
| *M S E* | 17.83 | 16.03 | 16.71 | 17.35 | 17.43 |
| *M AE* | 3.52 | 3.34 | 3.37 | 3.45 | 3.44 |
| *RMSE* | 4.22 | 4.00 | 4.09 | 4.17 | 4.18 |
| $R^2$ | -0.04 | 0.06 | 0.02 | -0.02 | -0.01 |
| *r* | NA | 0.31 | 0.27 | 0.18 | 0.17 |

*Note.* FL = featureless learner; RF = random forest; XG = extreme gradient boosting; SVM = support vector machine; EN = elastic net. *MSE* = mean squared error, *RMSE* = root mean squared error, *MAE* = mean absolute deviation, $R^2$ = coefficient of determination, *r* = Pearson correlation.

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Table 2

*Variable importance and Spearman correlations for the top 10 predictors*

| Feature | $I$ | $r_s$ |
| --- | --- | --- |
| mean frequency of missed calls per day | 0.62 | 0.32 |
| entropy of contacts for outgoing calls | 0.51 | 0.33 |
| entropy of contacts for missed calls | 0.41 | 0.29 |
| variation of frequency of outgoing calls per day | 0.32 | 0.26 |
| mean time of the last event on Friday/Saturday | 0.21 | 0.18 |
| variation of the time of the first event from Monday to Friday | 0.17 | 0.12 |
| mean number of intended events during night on Friday/Saturday | 0.14 | 0.16 |
| mean radius of gyration during night on Friday/Saturday | 0.14 | 0.31 |
| mean time of the last event on Sunday | 0.14 | 0.21 |
| mean frequency of outgoing calls per day | 0.13 | 0.24 |

*Note. I* = permutation-based variable importance. Variables are in descending order of importance scores.

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Table 3

*Descriptive statistics of app category usage*

| App category | $M_{Freq.total}$ | $SD_{Freq.total}$ | $Num_{Users}$ | $Num_{Apps}$ | $M_{ImpSS}$ | $SD_{ImpSS}$ |
|---|---|---|---|---|---|---|
| Communication apps | 1,522.88 | 1,576.08 | 234 | 59 | 7.97 | 4.24 |
| Social media apps | 485.34 | 715.83 | 203 | 70 | 8.18 | 4.20 |
| Browser apps | 210.35 | 316.01 | 231 | 22 | 8.02 | 4.24 |
| Music and audio apps | 183.97 | 522.02 | 190 | 85 | 8.13 | 4.20 |
| Planningtool apps | 92.85 | 185.26 | 209 | 70 | 7.97 | 4.22 |
| Gaming apps | 87.18 | 199.74 | 140 | 415 | 7.71 | 4.28 |
| Video watching apps | 80.43 | 155.22 | 210 | 48 | 8.00 | 4.25 |
| Trip planning apps | 38.60 | 61.20 | 208 | 52 | 8.04 | 4.17 |
| News apps | 31.63 | 138.31 | 125 | 48 | 7.80 | 4.22 |
| Shopping apps | 22.73 | 60.87 | 107 | 60 | 7.89 | 4.03 |
| Dating apps | 22.38 | 132.15 | 24 | 13 | 9.21 | 4.86 |
| Trading apps | 14.27 | 177.11 | 8 | 27 | 11.38 | 6.61 |
| Personalization apps | 12.88 | 147.34 | 47 | 34 | 8.06 | 4.72 |
| Running sports apps | 9.03 | 78.59 | 37 | 8 | 7.08 | 3.74 |
| Risky driving apps | 0.00 | 0.06 | 1 | 2 | 9.00 | NA |
| Outdoorsports apps | 0.05 | 0.74 | 1 | 8 | 12.00 | NA |

*Note.* $M_{Freq.total}$ = average total usage count within 30 days across all participants; $SD_{Freq.total}$ = standard deviation of average total usage count within 30 days across all participants; $Num_{Users}$ = number of all participants that have ever used any apps of the respective app category; $Num_{Apps}$ = number of different apps within one category. $M_{ImpSS}$ = average ImpSS score of all participants that have ever used any apps of the respective app category; $SD_{ImpSS}$ = standard deviation of the ImpSS score of all participants that have

DIGITAL FOOTPRINTS OF SENSATION SEEKING

ever used any apps of the respective app category. App categories are sorted in descending order of $M_{Freq.total}$.
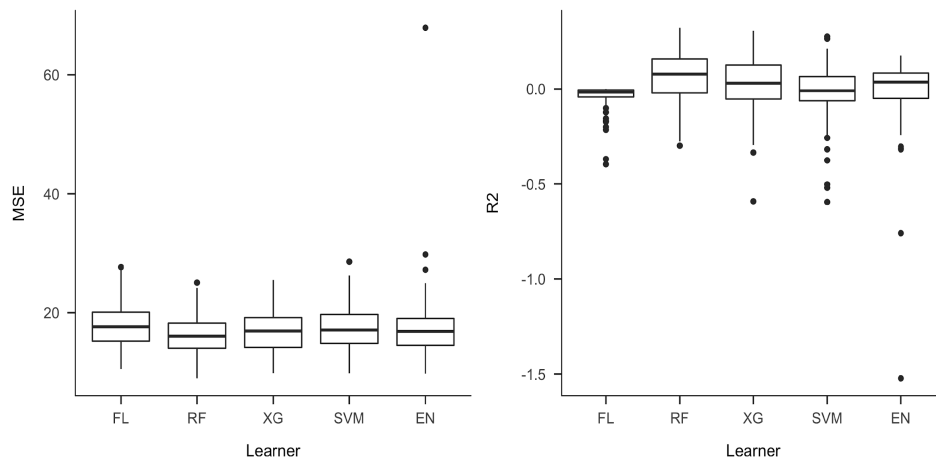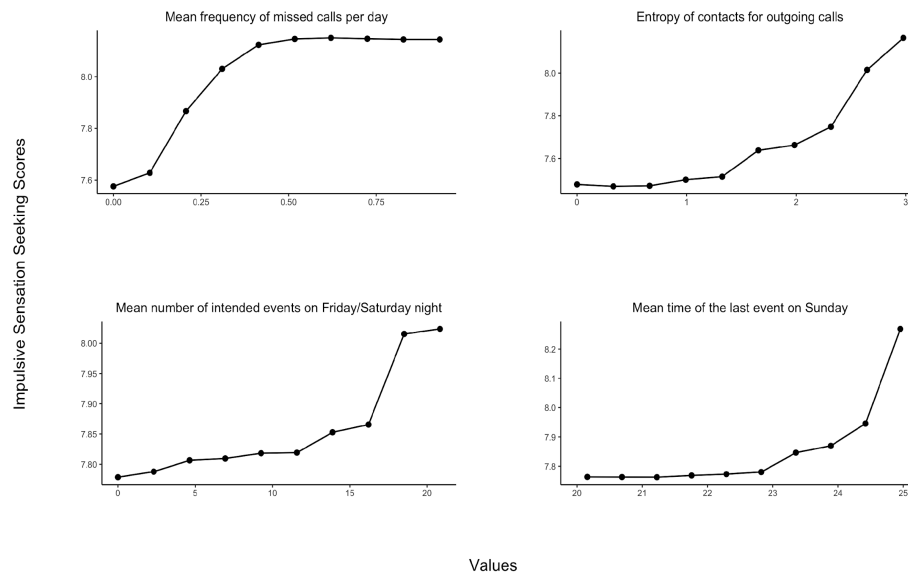
DIGITAL FOOTPRINTS OF SENSATION SEEKING



*Figure 1.* Distribution of the mean squared error (*MSE*) and the coefficient of determination

($R^2$) across all resampling iterations. FL = featureless learner; RF = random forest; XG =

extreme gradient boosting; SVM = support vector machine; EN = elastic net.

DIGITAL FOOTPRINTS OF SENSATION SEEKING



*Figure 2*. Plots displaying the partial dependence of the random forest learner function for four exemplary selected features. The curves show how predicted sensation seeking scores (on the vertical axis) change with increasing values (on the horizontal axis) of the respective displayed features. A last event value of 25 in the bottom right means, that the smartphone was used at 1am in the night between Sunday and Monday. The values of the horizontal axis range between the 10% and the 90% percentile of the respective features in order to clearly show how predicted sensation seeking scores change with increasing values of the respective features.

DIGITAL FOOTPRINTS OF SENSATION SEEKING

Appendix

Table A1

*Description of quantifications of behavioral data collected via smartphones.*

| Quantification | Description |
| --- | --- |
| | **app usage** |
| frequency/duration | Usage frequencies and durations of app usage (Stachl et al., 2017) were aggregated as daily mean and variation per day. As the logging of app usage is generally prone to logging errors, robust estimators were used: the huber mean as measure of central tendency (Kafadar, 2003) and the robust location-free scale estimate Qn as a measure of dispersion (Rousseeuw & Croux, 1993). Robust estimators are less sensitive to outliers which are possibly caused by faulty logs. |
| | **phone usage** |
| frequency/duration | Frequencies and durations of incoming/outgoing/missed calls/text messages were aggregated as daily mean and variation per day. We pre-registered to use robust measures for phone logging data, too. However, as their inspection did not reveal the same potential logging errors as for app usage data, we used the arithmetic mean and variance for feature calculation, because they are more precise estimators if outliers are not an issue. |
| response rate | The response rate was defined as percentage of missed calls and text messages people responded to by calling back within 24 hours (Montjoye et al, 2013). |
| | **app and phone usage** |
| entropy | The entropy describes how many categories one variable has (e.g. total number of contacts), while regarding how equally events (e.g. calls) are distributed across these categories. Therefore, entropy of contacts is high if a person called a broad range of contacts equally |

DIGITAL FOOTPRINTS OF SENSATION SEEKING

| | |
|---|---|
| | often within the study period. We also considered entropy of used apps, measuring how equally often participants used their individual spectrum of installed apps (Montjoye et al, 2013). |
| ratio | The ratio indicates the extent of certain behavioral categories in relation to the overall smartphone usage. For example, we considered the ratio of duration of dating app usage and overall smartphone usage (Montjoye et al, 2013). |
| irregularity | We computed irregularity of the point of time first and last events happen per day. As defined by Williams, M.J., Whitaker, and Allen (2012), this measure represents the dissimilarity of events in a time course. If events happen to very similar points across time (e.g. every day at 10am), dissimilarity and consequently irregularity are very low. |

| mobility | |
|---|---|
| radius of gyration | The radius of gyration was used for the quantification of mobility behavior (Canzian & Musolesi, 2015). It quantifies a person's range of mobility by considering the deviation from the center of all GPS positions, visited within one day. |
| total distance covered | The total distance covered was defined as summed distance between sequent GPS points per day (Canzian & Musolesi, 2015). |
| maximum distance covered | The maximum distance covered was defined as maximum stretch of way per day (Canzian & Musolesi, 2015). |

*Note.* Unlike stated in our pre-registration, we had to exclude the predictors "number of contacts at the beginning of the study" and "number of contacts added within 30 days" as our contact logging data turned out to be corrupted for the majority of our participants due to logging errors. Only after completion of our pre-registration, we conceived the "total usage frequency of app categories within 30 days" as additional important predictors and therefore decided to add them to our feature list.

162

# Personality Traits Predict Smartphone Usage

This article explores how individual traits, fluid intelligence, and demographics predict smartphone app usage across different categories.

### *Contributing article*

Stachl, C., Hilbert, S., Au, J., Buschek, D., Luca, A. D., Bischl, B., Hussmann, H., and Bühner, M. (2017). Personality traits predict smartphone usage. *European Journal of Personality*, 31(6):701–722

This publication was part of the PhoneStudy project (see Chapter 7).

### *Copyright information*

According to *Sage's Author Archiving and Re-Use Guidelines*[1] the final published PDF can be used for publication in a dissertation or thesis, including where the dissertation or thesis will be posted in any electronic Institutional Repository or database.

### *Declaration of contributions*

In this article, the doctoral candidate made substantial contributions to the statistical analysis of the collected sensor and questionnaire data. Jiew-Quay Au handled a range of tasks, including data cleaning and computation of necessary variables for further analysis, as well as the selection and application of statistical models. Moreover, the PhD candidate created all the tables in the manuscript and contributed to writing parts of the chapter on variable selection and data analysis.

#### *Contribution of the coauthors*

Clemens Stachl was the main author and had the responsibility for data collection and manuscript writing. Sven Hilbert contributed to the study design and provided valuable input for the manuscript with creative ideas. Bernd Bischl aided in selecting the most appropriate statistical models to confirm the hypotheses. The remaining co-authors supported the article by contributing to the revision of all sections.

---

[1] https://us.sagepub.com/en-us/nam/journal-author-archiving-policies-and-re-use

# Personality Traits Predict Smartphone Usage🏠🔬

CLEMENS STACHL[1]*, SVEN HILBERT[1,2], JIEW-QUAY AU[3], DANIEL BUSCHEK[4], ALEXANDER DE LUCA[4], BERND BISCHL[3], HEINRICH HUSSMANN[4] and MARKUS BÜHNER[1]

[1]*Department of Psychology, Psychological Methods and Assessment, Ludwig-Maximilians-Universität München, Munich, Germany*
[2]*Faculty of Psychology, Educational Science, and Sport Science, University of Regensburg, Munich, Germany*
[3]*Department of Statistics, Computational Statistics, Ludwig-Maximilians-Universität München, Munich, Germany*
[4]*Media Informatics Group, Ludwig-Maximilians-Universität München, Munich, Germany*

*Abstract: The present study investigates to what degree individual differences can predict frequency and duration of actual behaviour, manifested in mobile application (app) usage on smartphones. In particular, this work focuses on the identification of stable associations between personality on the factor and facet level, fluid intelligence, demography and app usage in 16 distinct categories. A total of 137 subjects (87 women and 50 men), with an average age of 24 (SD = 4.72), participated in a 90-min psychometric lab session as well as in a subsequent 60-day data logging study in the field. Our data suggest that personality traits predict mobile application usage in several specific categories such as communication, photography, gaming, transportation and entertainment. Extraversion, conscientiousness and agreeableness are better predictors of mobile application usage than basic demographic variables in several distinct categories. Furthermore, predictive performance is slightly higher for single factor—in comparison with facet-level personality scores. Fluid intelligence and demographics additionally show stable associations with categorical app usage. In sum, this study demonstrates how individual differences can be effectively related to actual behaviour and how this can assist in understanding the behavioural underpinnings of personality. Copyright © 2017 European Association of Personality Psychology*

Key words: Big Five; factor and facets; behaviour; smartphones; app usage

A core goal of personality psychology is the prediction and explanation of behaviour with regard to individual differences (Back, Schmukle, & Egloff, 2009; Furr, 2009). Although there is a broad consensus in the research community concerning this topic, most studies pursuing personality research heavily rely on self-reported behaviour and neglect observable acts. This evident lack of observable behaviour in personality psychology has been repeatedly criticised (Baumeister, Vohs, & Funder, 2007; Funder, 2001; Furr, 2009). Even though notable exceptions exist (e.g. Back et al., 2009; Mehl, Gosling, & Pennebaker, 2006), most studies still focus on people's reports about how they think they behave rather than their actual behaviour. Moreover, in order to fully understand personality and its impact on real-life behaviour, one must ideally consider behaviour outside the laboratory (Back et al., 2009). Fortunately, the technological development of the last two decades has not only had a strong impact on our everyday life but has also changed the game for psychological research.

*Correspondence to: Clemens Stachl, Psychological Methods and Assessment, Department of Psychology, Ludwig-Maximilians-Universitat München, 80802 Munich, Germany.
E-mail: clemens.stachl@psy.lmu.de

**Digital frontiers of personality research**

The continuous expansion of Internet use and even more so the broad availability of cheap mobile sensor technology makes smartphones a central part of peoples' digital life (e.g. Lane et al., 2010). These developments offer an extraordinary opportunity for personality studies to remove various hindering factors for measuring actual behaviour in personality research (Gosling & Mason, 2015). Smartphones in particular are capable of gathering large and diverse data samples of individual behaviour (Miller, 2012; Yarkoni, 2012) across a broad variety of situations (Harari, Gosling, Wang, & Campbell, 2015). The accumulation of this immense amount of information has led to the stepwise inclusion of these devices in psychological science (see Harari et al., 2016; Wrzus & Mehl, 2015, for reviews) and paved the way for a field of research termed *psychoinformatics* (Yarkoni, 2012).

Besides the opportunities mobile sensing technologies provide for psychological research in general, they yield benefits for personality psychology in particular: in addition to communication data, such as calls (Montag et al., 2014) and messaging (Montag et al., 2015), an ever-growing amount of behaviour-related information is provided through the use of apps. Currently, more than two million apps exist in the leading app stores providing a wide range of functionality to the user, such as playing music, dating, searching the

C. Stachl et al.

web, accessing social networks or messaging friends (Statista, 2016). A wide range of behaviours, all of which are typically of interest to researchers investigating individual differences, are therefore concentrated in a single device that may be tapped for information. This, in turn, opens the door for a convenient form of data collection: many samples of everyday-life behaviour can be continuously collected over a large set of different situations in a most unobtrusive manner. Consequently, behavioural aggregation is significantly facilitated (Epstein, 1983; Harari et al., 2016; Vazire, 2010), and compared with conventional behavioural observations, samples obtained with mobile sensing approaches show high ecological validity (Schmid Mast, Gatica-Perez, Frauendorfer, Nguyen, & Choudhury, 2015).

Naturally, these new opportunities for effective collection of ecologically valid behavioural data have already resulted in several studies on possible associations between smartphone usage and individual differences. The following sections describe the various approaches undertaken to make use of smartphone data in personality research, and Table 1 provides an overview of the most relevant studies.

**Personality and smartphone usage**

Several studies on the prediction of smartphone use through personality collected self-reports of behaviour and related them to scores in personality inventories (Butt & Phillips, 2008; Kim, Briley, & Ocepek, 2015; Lane & Manner, 2011). While Lane and Manner (2011) asked participants to state the importance of various smartphone functions, Butt and Phillips (2008) and Kim et al. (2015) specifically asked participants to quantify how often and how long they used apps of a certain type on their phones. Making use of the aforementioned opportunities of modern sensor-imaging data, several recent investigations directly logged user behaviour to examine associations with self-reported personality scores (Chittaranjan, Blom, & Gatica-Perez, 2013; De Montjoye, Quoidbach, Robic, & Pentland, 2013; Montag et al., 2014, 2015; Xu, Frey, Fleisch, & Ilic, 2016). While the studies of Montag et al. (2014, 2015) exclusively focused on communication behaviour, others collected a broad range of app usage (Chittaranjan et al., 2013) as well listings of installed apps in general (Xu et al., 2016). Notably, Xu et al. (2016) did not record app usage over a period of time but focused on the snapshot of currently installed applications on a smartphone. In contrast to most other studies, De Montjoye et al. (2013) used a classification-based approach to predict personality from smartphone usage. As outlined in Table 1, a wide range of personality traits has been related to differential smartphone usage of various forms, such as frequencies and durations of calls and several categories of apps.

As could be expected based on previous research (e.g. Eaton & Funder, 2003), extraversion was mostly found to be associated with communication app usage and social behaviours, but also with apps for educational purposes. In a study conducted by Butt and Phillips (2008), extraverted participants reported to make more calls for longer durations and used communication-related applications (messaging) intensively. Additionally, it was found that extraverts reported

higher use of personalisation activities, such as changing the wallpaper or the ringtone of the phone. Lane and Manner (2011) observed higher reported importance of communication apps for extraverted participants. Similarly, Kim et al. (2015) found that extraversion positively predicted the reported use of communication and social apps. Furthermore, Kim et al. (2015) also found that extraversion acted as a negative predictor for app usage frequencies related to app categories like books and references and education. In previous studies on observed smartphone behaviour, extraversion was positively associated with both the frequency and the duration of calls (Chittaranjan et al., 2013; Montag et al., 2014). Furthermore, the frequency of communication app usage was positively related to levels of extraversion (Chittaranjan et al., 2013; Montag et al., 2015). Chittaranjan et al. (2013) additionally reported lower frequencies of browser usage for extraverts as well as increased usage of apps for entertainment purposes. Xu et al. (2016) observed a smaller number of game apps, installed on phones of extraverts. Thus, as in direct interpersonal behaviour, extraversion shows a solid link to communication and stimulating activities on smartphones. It is, however, not the sole personality trait that is reflected in a somewhat straightforward manner through smartphone usage: the same holds for agreeableness, as described in the next paragraph.

Traditionally, agreeable people express more pro-social behaviours and positive language to others (e.g. Graziano & Tobin, 2009; Mehl et al., 2006). However, agreeableness was also associated with reported communication and personalisation app usage (Butt & Phillips, 2008; Lane & Manner, 2011). Participants high in agreeableness reported shorter durations of use for communication apps and calls as well as lower frequencies of calls in general. Furthermore, they observed shorter usage durations of personalisation functions for people scoring high in agreeableness (Butt & Phillips, 2008). Agreeable participants also reported high importance of the call function and lower importance for communication apps (messaging; Lane & Manner, 2011). Contrarily, Kim et al. (2015) did not find any association of agreeableness and self-reported app usage. Studies utilising smartphone usage logs solely observed negative associations of app usage frequencies with the trait of agreeableness (Chittaranjan et al., 2013; Xu et al., 2016). Specifically, Chittaranjan et al. (2013) reported negative associations between participants scoring high in agreeableness and usage frequencies of apps related to communication, browser usage, productivity and gaming. In terms of installed applications, Xu et al. (2016) reported lower numbers of apps related to personalisation on agreeable peoples' phones has

In addition to extraversion and agreeableness, the personality trait of conscientiousness was associated with reports of smartphone usage (Kim et al., 2015; Lane & Manner, 2011). So far, high levels of conscientiousness have mostly been related to formal behaviour, conformity to rules and self-organisation (e.g. Jackson et al., 2010). Regarding usage behaviour on smartphones, Lane and Manner (2011) reported shortened durations of communication app usage for participants high in conscientiousness. In a large sample, conscientiousness was also reported as a negative predictor for

Table 1. Literature review—personality and smartphone usage

| Study (year) | Sample size (N) | Personality measure | Emotional stability | Extraversion | Openness | Conscientiousness | Agreeableness |
|---|---|---|---|---|---|---|---|
| **Self-reported smartphone usage** | | | | | | | |
| Butt and Phillips (2008) | 112 | NEO-FFI | – Communication (dur.) | + Calls (freq., dur.)<br>+ Communication (dur.)<br>+ Personalisation (dur.) | | | – Communication (dur.)<br>– Calls (freq., dur.)<br>– Personalisation (dur.)<br>+ Calls (imp.) |
| Lane and Manner (2011) | 312 | BFPI | | + Communication (imp.)<br>+ Communication (freq)<br>+ Social (freq)<br>– Books & References (freq)<br>– Education (freq) | – Shopping (freq.)<br>– Finance (freq.) | | – Communication (imp.) |
| Kim et al. (2015) | 9482 | TIPI | | | | – Communication (dur.) | |
| **Observed smartphone usage** | | | | | | | |
| Chittaranjan et al. (2013) | 117 | TIPI | + Communication (freq.)<br>+ Games (freq.) | – Browser (freq.)<br>+ Calls (freq., dur.)<br>+ Entertainment (freq)<br>+ Communication (freq.)<br>+ Calls (freq., dur.)<br>+ Communication (dur.) | + Entertainment (freq.)<br>– Business (freq.)<br>– Communication (freq.) | + Communication (freq.)<br>+ Business (freq.)<br>– Entertainment (freq.)<br>– Calls (freq.)<br>+ Calls (dur.) | – Communication (freq.)<br>– Browser (freq.)<br>– Productivity (freq.)<br>– Games (freq.) |
| Montag et al. (2014) | 49 | NEO-FFI | | | | –Communication (dur.)<br>–Music and video (num.)<br>–Photography (num.)<br>–Personalisation (num.) | |
| Montag et al. (2015) | 2418 | BFI | | | | | |
| Xu et al. (2016) | 2043 | BF-44Q | + Photography (num.)<br>– Personalisation (num.) | – Games (num.) | | | – Personalisation (num.) |

*Note.* Summary of previous studies reporting on associations of personality with smartphone and app usage. freq, frequencies; dur., durations; imp., importance; num., number of installed apps; NEO-FFI, NEO-five factor inventory; BFPI, Big Five personality inventory; TIPI, ten-item personality inventory; BFI, Big Five inventory; BF-44Q, Big Five 44 questionnaire.
⁺Positive relationship.
⁻Negative relationship.

C. Stachl et al.

frequencies of app usage for shopping and finance matters (Kim et al., 2015). Unlike in previous studies using self-reported smartphone usage, conscientiousness was associated with a wide range of app usage behaviours in studies using app logs. Chittaranjan et al. (2013) reported positive and negative associations of conscientiousness with various app usage behaviours. Communication and business app frequencies on the one hand were observed to be higher for participants with high scores in conscientiousness. The usage of entertainment apps and calls on the other hand was found to be negatively associated with high levels of the trait (Chittaranjan et al., 2013). In a study conducted by Montag et al. (2014), conscientiousness was positively related to the duration of calls, whereas communication app usage was found not to be related to any personality trait. In contrast, durations of communication app usage (WhatsApp) were negatively correlated with a persons' level of conscientiousness in a later study (Montag et al., 2015). Additionally, Xu et al. (2016) reported fewer installed apps of the categories music and video, photography and personalisation for conscientious participants. Taken together, it seems that the self-organisational aspect of conscientious behaviour is especially reflected by smartphone use with fewer installations and less usage of applications related to procrastination behaviour.

The personality traits of emotional stability, on the other hand, have not been frequently associated with app usage in previous studies. In an investigation conducted by Butt and Phillips (2008), shorter durations of communication app usage were reported by emotionally stable participants. Additionally, studies using smartphone logs found emotional stability to be associated with actual app usage in two previous studies (Chittaranjan et al., 2013; Xu et al., 2016). Chittaranjan et al. (2013) reported higher frequencies of communication and gaming app usage for emotionally stable participants. Furthermore, Xu et al. (2016) found that emotionally stable participants had more apps in relation to photography and fewer apps related to personalisation installed on their phones.

Openness was not associated with any app usage behaviours in studies using self-reports. Only Chittaranjan et al. (2013) found that openness was positively associated with the usage of entertainment apps. They additionally reported lower frequencies of business and communication apps for participants with high scores in openness.

Considering these results, it is interesting to note that most associations of app usage behaviour and personality traits could be established for the dimensions of extraversion, agreeableness, and conscientiousness (Table 1). Furthermore, studies using observed behaviour were able to relate personality traits to more distinct types of app usage as well as to openness and emotional stability, for some categories (Chittaranjan et al., 2013; Xu et al., 2016). All in all, these studies show that behaviour associated with each of the Big Five personality traits tends to be (at least partially) reflected by the use of smartphones. Notably, the positive associations of extraversion along with the frequency and duration of communication-related activities provide the most solidly established link. Unfortunately, as the research relating personality and smartphone usage is still young and sparse, the conclusiveness of the existing results is still limited because of methodological differences. Most studies have used different categories of app usage and different definitions of behaviour (frequencies, durations, installed apps, importance and self-reports vs. actual behaviour). Furthermore, it is difficult to know which of the reported associations are reliable, as some results were based on an immense amount of significance tests without type I error correction (e.g. Chittaranjan et al., 2013).

**Personality factors and facets**

Most existing research utilised short-scale personality questionnaires (e.g. Chittaranjan et al., 2013; Kim et al., 2015; Lane & Manner, 2011; Montag et al., 2015; Xu et al., 2016) to measure individual personality scores. From the perspective of personality psychology, an investigation at factor and facet level could be beneficial to increase understanding of observed effects. In order to predict behavioural criteria from individual personality scores, either broad factor values or sub-facet scores can be used. Some previous research suggests that personality facet measures provide independent prediction value in relation to behavioural criteria in addition to factor-level scores (Ziegler et al., 2014). However, disagreement is prevalent in current research concerning this topic (Ashton, Paunonen, & Lee, 2014; Salgado, Moscoso, & Berges, 2013). In particular, uncertainty remains with regard to whether factor-level or facet-level scores are better for predicting behavioural categories. This knowledge would, in turn, also help to understand whether particular types of smartphone usage represent narrow or rather broad categories of behaviour, as it is known that measurement symmetry (or Brunswik Symmetry) enhances the strength of relationships between constructs (and their and related observable behaviour; Wittmann, 2012). In the present study, we therefore relate both factor and facet personality scores to behavioural criteria and compare the obtained results. Although some previous studies have reported on associations between personality and app usage (e.g. Chittaranjan et al., 2013; Kim et al., 2015; Montag et al., 2015; Xu et al., 2016), more guidance on which associations are the most stable and promising ones for prospective confirmatory investigations is necessary. Hence, in the final analysis of the present study, we use directly logged app usage behaviour from 16 distinct categories, combined with fine-grained self-reported personality measures. Furthermore, we utilise novel statistical resampling procedures in combination with regularised regression models to identify the most reliable associations between personality traits and app usage behaviour. We adopt this rather careful statistical approach to acknowledge the exploratory nature of this investigation, which is based on the sometimes diverging findings from only a few previous studies.

**Intelligence**

Moreover, we include a measure of fluid intelligence in the analysis, as it is also likely to also influence the way a smartphone is operated (Zaval, Li, Johnson, & Weber,

2015). The adequate use of cognitive abilities has been repeatedly linked to the acceptance of new technology (e.g. Morris, Venkatesh, & Ackerman, 2005; Venkatesh, Thong, & Xu, 2012). Moreover, as intelligence has been long known to profoundly influence interpersonal behaviour (e.g. Zander & Van Egmond, 1958), which could easily be reflected in the use of particular apps for interpersonal communication. This is also supported by the association between intelligence and several (sub-clinically pathological) personality traits (Rauthmann, 2012), which have additionally been linked to smartphone addiction (Lee, Chang, Lin, & Cheng, 2014).

### Rationale

The current work has three main goals: first, it explores which personality traits, fluid intelligence and demographics predict observed behaviour, manifested in several categories of mobile app usage. Second, this study originally compares personality on factor and facet level in terms of predictive capabilities for smartphone usage. Finally, this work is thought to act as an example of how smartphones can be used to collect large amounts of behavioural data and how it can be analysed for personality research in particular.

## METHOD

Data used in this work constitute a part of a larger research project at Ludwig-Maximilians-Universität München (LMU), investigating relationships between psychological variables and a wide range of behaviour, logged via smartphones (e.g. app usage, music consumption and geolocation). However, this paper only focuses on app usage behaviour, including calls and messages, and explores its relationship with personality, fluid intelligence and demographic variables. Therefore, further descriptions will only include data dimensions related to the present analyses. Data collection took place between September 2014 and August 2015 in Munich, Germany, EU.

### Participants

We recruited 137 participants, 87 women and 50 men, via social media, forums, blackboards, flyers and on campus. The obtained sample was rather young with a mean age of 24 years ($SD = 4.71$) and 75% of the participants being 26 or younger. The reported age ranged from 18 to 50 years. The majority of the sample had at least completed high school (96%), and 31% of all participants had completed education at the university level. The obtained sample consisted primarily of LMU students and employees. Subjects were informed about the data collection procedures as well as anonymization procedures and gave written consent prior to participation. Participants could withdraw participation in the study as well as demand deletion of not yet anonymized data at any time during the data collection period. Because of technical requirements, only people using the Android operation system participated in the study. The study was

approved by the responsible ethics committee and data protection officer. See Table A1 for all descriptive statistics.

### Procedure

The study was conducted in two stages. During a lab session, participants gave written consent and completed a personality inventory, subscales of the Intelligence-Structure Battery (INSBAT; Arendasy et al., 2009), and a demographic questionnaire. The testing took place at the psychological laboratory at the LMU's psychology department. Subsequently, the logging app was installed on the participants' private smartphones and tested for functionality.

The app logged frequencies and durations of mobile application usage for a total of 60 days. The data were regularly uploaded to our servers, when the respective Android smartphone was connected to a wireless Internet connection (WiFi). The parameters were initially logged as timestamp events. The app stayed in the background: participants did not have to complete any tasks or actions to avoid altering their natural smartphone use. As a sole exception, they were reminded (via a pop-up message) to re-enable location sensor and app history access (Android 5.0 and higher) in case they had turned off these settings. After 60 days of logging, the participants were invited to receive their compensation (an individual personality profile and 30 EUR or course credit for students). During this meeting, an additional manual backup of the collected data was retrieved from the device.

### Materials

*Psychometric measures and demographics*
Big Five personality scores were measured with the German version of the Big Five Structure Inventory (BFSI; Arendasy, 2009) in a laboratory setting. The BFSI was selected for personality assessment because of its unambiguous items as well as its favourable psychometric properties. In contrast to more common personality scales such as the new NEO-PI-3 (McCrae, Costa Jr., & Martin, 2005), the authors of the BFSI (Arendasy, 2009) report on the tests conformity to the partial credit model (PCM) (Masters, 1982). The PCM is the most used uni-dimensional item-response model for Likert scale ordinal data. The model assumes that the probability for a single-item response is solely a function of a person's value on an underlying latent variable (person parameter) and the item thresholds. Within the PCM, specific objective comparisons between person and items scores are possible (Masters, 1982).

We therefore used the person parameter of the PCM instead of sum scores for all analyses. The BFSI consists of 300 items and measures the Big Five personality dimensions (*Openness to Experience*, *Conscientiousness*, *Extraversion*, *Emotional Stability/Absence of Neuroticism*, and *Agreeableness*) on the factor and the facet level with a four-point Likert scale ranging from '*untypical for me*' to '*typical for me*'. With regards to the lexical derivation of the Big Five, the BFSI uses adjectives and short phrases as items for personality assessment. This could also help to circumvent previously reported problems regarding the comprehensibility of longer

C. Stachl et al.

sentence-based items, such as in the NEO-FFI (Costa & McCrae, 1992; McCrae, Costa, & Martin, 2005).

Fluid intelligence was assessed with the German version of the INSBAT using aggregated scores of subscales for *Numerical Inductive Reasoning*, *Figural Inductive Reasoning*, and *Verbal Deductive Reasoning*. The INSBAT is an adaptive intelligence test and is based on the hierarchical intelligence theories of Horn and Noll (1997) and Carroll (2003). The sub-factor *Fluid Intelligence* measures a person's ability to identify patterns in stimuli and to draw logical conclusions from given sequences of figures, numbers and statements.

In addition to personality and fluid intelligence, age, gender and the current level of completed education were collected. Gender was recorded dichotomously with '1' representing male and '2' representing female participants. The level of education was split into five categories from 'not completed obligatory schooling' to 'university degree'. This should be taken into consideration when interpreting correlations in Tables A2 and 3 as well as regression coefficients in Tables 4 and 5. Because of a fatal hard drive error of one of our testing laptops, internal consistency scores were calculated based on 120 instead of 137 total participants.[1] We calculated 95% confidence intervals for internal consistencies of the extraversion (95%CI [0.94, 0.97]), agreeableness (95%CI [0.92, 0.95]), conscientiousness (95%CI [0.95, 0.97]), emotional stability (95%CI [0.91, 0.95]) and openness (95%CI [0.91, 0.95]) scales. Situated in the framework of the PCM (Masters, 1982), the INSBAT provides reliability coefficients for every single participant (Arendasy et al., 2009). Therefore, we calculated mean scores across all participants for the *Numerical Inductive Reasoning* scale ($M = 0.71, SD = 0.06$), the *Verbal Deductive Reasoning* scale ($M = 0.74, SD = 0.06$) and the scale for *Figural Inductive Reasoning* ($M = 0.71, SD = 0.06$). Table A1 provides Cronbach's alpha coefficients and descriptive statistics for all psychometric scales included in this study.

*Behavioural measures*

User behaviour was recorded via an Android logging app (available for Android 4.0 or higher), specifically designed for this purpose. The app uses a background service to monitor app usage. The service assesses the currently running app every two seconds, creating a log entry if it had changed. Devices operating on newer versions of the Android operating system support reading the app usage history directly. On capable devices, our app thus automatically switches to this method, retrieving the latest history every 15 min. The data were regularly transferred to our server, once participants were connected to WiFi, using SSL encryption. Afterwards, the logged data were further enriched with information about the respective app categories, retrieved from the Google Play Store (Google, 2016) via web scraping.[2] Originally, the raw

data existed as a timestamp sorted list of events—every event that happened was logged in the file. Participants' demographic and psychometric data were recorded and stored separately and were only combined with logging data for the purpose of statistical analysis.

**Data processing**

Prior to modelling, we had to pre-process and clean the data. In the great majority of cases, we used the app categorisation as provided by the Google Play Store (Google, 2016) for the categorisation of the recorded applications. However, a number of apps were manually identified as mislabeled and had to be re-labelled in order to perform meaningful data analysis. We only re-labelled apps that were clearly assigned to a different category (e.g. an SMS/MMS app was first labelled as personalisation and re-labelled to communication) in order to create a more meaningful data set. Additionally, bloatware[3] and background apps were manually identified and consequently not used for the calculation of app usage durations and frequencies. We share the notion that the labelling of information is always somehow subjective and therefore provide the full list of relabeled apps and bloatware apps as supplemental files to this article.

With regard to the analysis described later, usage logs of apps in a certain category (e.g. communication) were aggregated. Please note that we treated calls in the same fashion as any other category of app usage and analysed total frequency and average duration. For the regression analyses, we used the total number of all app launches in a respective category for each participant over the study period of 60 days as well as the average duration of an app use of a certain category across all events from that category. A single-usage duration was calculated as the time from the start of an app-event until the start of the next event that is not labelled as bloatware.

In order to handle univariate outliers in the data, we first identified robust $z$-transformed values with values larger than 3. Robust $z$-transformation was performed by subtracting values by the median and dividing the result by the median absolute deviation. The median absolute deviation is a robust measure of variability in a univariate data sample. The values were then adjusted to the maximum value of the remaining data points (winsorizing; Ghosh & Vogt, 2012). This procedure allowed us to not waste data while limiting the influence of very extreme, possibly spurious data points. Furthermore, we only included criteria variables with a median absolute deviation larger than zero (0.0001 to be precise), excluding categories with no or almost no variation in the data prior to the analysis (e.g. comics).

**Variable selection and data analysis**

Prior to regression modelling, we investigated descriptive statistics as well as correlations between the Big Five factors and the demographic variables. In order to investigate the effect of personality and demography on app usage behaviour,

---

[1]The hard drive crashed after the factor scores were extracted yet before single-item scores were read out, as this procedure was undertaken at a later time. Because our analyses were run with the factor and facet scores, this did normally not affect the sample size—except for the estimation of the internal consistencies, which is based on single-item values.

[2]Web scraping refers to the practice of extracting information from a website.

[3]Bloatware refers to pre-installed, mostly unwanted, software that often negatively affects system performance of devices.

we performed a two-step analysis for factor and facet scores, respectively. Because of the high number of predictors (facet analysis) and because of the expected multi-collinearity between the personality and demographic variables (visible in Table A2), we used a stability selection procedure (Hofner, Boccuto, & Göker, 2015; Meinshausen & Bühlmann, 2010) in combination with the popular Least Angular Shrinkage Selection Operator (LASSO) penalised regression (Friedman, Hastie, & Tibshirani, 2010). This procedure was chosen in order to only select the most reliable predictors while penalising for correlations between them. The LASSO regressions were modelled under the assumption of a Poisson distribution with each app usage variable as the respective criteria.

Stability selection refers to a relatively new concept that adds resampling procedures to variable selection, such as the LASSO, and therefore makes the selection procedure more reliable (Meinshausen & Bühlmann, 2010). This procedure avoids to fit only one model on the complete data, but instead fits many different ones on subsets. Therefore, variables that repeatedly (above a certain threshold) add predictive value to different models are selected. In other words, stability selection can be considered as another tool of variable selection for the subsequent regression models. Other approaches to variable selection exist but have been heavily criticised for overfitting the data (e.g. stepwise regression or predictor selection based on significance of univariate correlation with the criterion; Rencher & Pun, 1980). We chose this approach as it is supposedly less likely to overfit the data (because of repeated modelling of subsets).

Furthermore, this procedure allows for the assessment of the selection stability of variables while controlling for sample error. Hofner et al. (2015) suggests to set the upper limit of the pairwise family error rate (PFER) to be set at $\alpha < PFER_{max} < m\alpha$, where $m$ represents the number of predictors and $\alpha$ represents the respective significance level ($m_{factor}\alpha = 9 \times 0.05 = 0.45$ and $m_{facet}\alpha = 34 \times 0.05 = 1.7$ in our case). Based on this recommendation, we used an even lower PFER of 0.20. We chose this parameter value because PFER represents the tolerable number of falsely selected noise variables. Therefore, we kept this value well below 1, tolerating less than one noise variable. Furthermore, the stability selection procedure can be described as a parameter tradeoff between the number of to be selected predictors (q), a probability cut-off and the PFER where only two parameters can be specified simultaneously. In the present analysis, we decided to limit the PFER to 0.2 and to only accept variables as stably selected if they appeared in at least 70% of the subsampled LASSO models, which equals to a probability cut-off of 0.7 in the analysis. We did not limit the number of predictors that could be selected.

In a second step, the predictors selected through stability selection were again used as predictors in separate quasi-Poisson regressions, with app usage categories as the respective criteria. This additional step was performed because regression coefficients of a penalised model are hard to interpret. We chose generalised linear regression over linear

regression analysis as count data (and durations) usually follows a Poisson distribution (O'Hara & Kotze, 2010). In order to account for overdispersion in our data set, we assumed quasi-Poisson distributions instead of Poisson distributions for the dependent variables.

Big Five personality scores measured on factor and facet level, as well as demographics, were used as predictors in the regression models. This procedure was repeated for each app usage category respectively. In order to compare factor with facet models, we used the *Dawid–Sebastian* score as a measure of model fit. This measure is similar to the mean squared error but additionally accounts for overdispersion in count data (Czado, Gneiting, & Held, 2009). For a quasi-Poisson-distributed random variable $X$, with $E(X) = \mu$ and $Var(X) = \theta\mu$, the *Dawid–Sebastian* score for an observed value $x$ is calculated as follows:

$$DSS(x) = \frac{(x-\mu)^2}{\theta\mu} + 2\log(\theta\mu).$$

For example, $X$ could be the usage frequency of communication apps, $\mu$ would represent the predicted app usage frequency (needs to be estimated) and the variance of app usage is $\theta\mu$, where $\theta$ is the overdispersion parameter of assumed quasi-Poisson distribution. In order to obtain an unbiased estimation of model fit, we used a Monte Carlo resampling procedure. In particular, we created a test (10%) and a training set (90%), fitted a generalised linear regression model with quasi-Poisson distribution on the training set and calculated the mean *DSS* across all observations. In order to calculate the *DSS*, $\mu$ and $\theta$ were estimated from the training set. This procedure was repeated 100 times for each criterion, and *DSS* scores were averaged across all observations in each test set. In comparison analyses, we made sure that equal test and training set splits were used in the process. Please note that for some app usage categories (visible in Table 2), no modelling was performed as not enough data were available. These categories are therefore not reported in the results section (e.g. comics); see Tables 4 and 5 for all predicted categories.

As we solely ran regression analyses for the predictors identified via stability selection, our sample size of $N = 137$ allowed us to detect the influence of $\delta = 0.24$ regarding the size of regression weights, corresponding to a rather small effect size, with a statistical power of $1\text{-}\beta = 0.80$.

All data processing as well as statistical analyses in this study were performed with statistical software R 3.3.1 (R Core Team, 2016). Additionally, we used the *glmnet* package for statistical modelling and the *stabs* package for stability selection (Friedman et al., 2010; Hofner & Hothorn, 2015). We also used the *mada* and *usdm, psych, haven* and *GPArotation* packages (Bernaards & Jennrich, 2005; Doebler, 2015; Naimi, 2015; Revelle, 2016; Wickham & Miller, 2016). See the Supporting Information section for a link to the full data set, analysis scripts and supplemental files.

Table 2. Descriptive statistics—categories of app usage

| Category | $M_{Num.\ uses}$ | $SD_{Num.\ uses}$ | $M_{usage}$ (s) | $SD_{usage}$ (s) | Num. Apps | Num. Users | Most frequently used apps |
|---|---|---|---|---|---|---|---|
| Communication | 38.48 | 25.87 | 31.14 | 13.60 | 62 | 137 | WhatsApp, Mail, Contacts, Dialer, SMS/MMS |
| Social | 7.26 | 10.87 | 50.52 | 47.26 | 80 | 125 | Facebook, Instagram, Snapchat, Twitter, Weibo |
| Tools | 6.80 | 8.07 | 16.20 | 9.00 | 225 | 137 | Google Search, Clock, Google Play Store, Calculator, S Voice |
| Browser | 6.44 | 6.17 | 71.56 | 29.93 | 6 | 136 | Internet, Firefox, Opera, Dolphin Browser, UC Browser |
| Calls | 4.14 | 4.07 | 103.73 | 162.45 | 1 | 137 | Phone |
| Productivity | 3.45 | 4.21 | 23.22 | 10.21 | 134 | 137 | Settings, S Planner, Calendar, ColorNote, Google Drive |
| Photography | 2.69 | 3.23 | 23.57 | 11.46 | 47 | 137 | Gallery, Camera, SnapApp, Album, PicsArt |
| Games | 2.52 | 5.39 | 153.33 | 196.64 | 229 | 100 | Clash of Clans, Quizduell, Candy Crush Saga, Farm Heroes Saga, Trials Frontier |
| Music and audio | 1.40 | 1.92 | 15.37 | 17.53 | 78 | 134 | Spotify, Music Player, Google Play Music, MP3-Player, SoundCloud |
| Entertainment | 1.08 | 1.46 | 94.01 | 108.69 | 98 | 131 | YouTube, 9GAG, PlayerPro, appinio, PS4-Magazin |
| Travel and local | 1.01 | 1.01 | 54.61 | 24.91 | 85 | 134 | Maps, MVV Companion, TripAdvisor, BlaBlaCar, Airbnb |
| Transportation | 0.90 | 1.02 | 39.47 | 25.77 | 40 | 110 | MVG Fahrinfo, DB Navigator, Oeffi, MeinFernbus, Uber |
| News and magazines | 0.77 | 1.40 | 36.66 | 51.08 | 52 | 118 | FOCUS Online, reddit sync, SPIEGEL ONLINE, Flipboard, SZ.de |
| Lifestyle | 0.54 | 1.17 | 24.79 | 40.00 | 72 | 72 | Tinder, Sleep, Chefkoch, eBay Kleinanzeigen, PAYBACK |
| Sports* | 0.52 | 3.52 | 18.18 | 77.43 | 38 | 33 | kicker, Comunio, Kicktipp, Score!, Sportschau |
| Books and reference* | 0.45 | 0.82 | 31.24 | 63.90 | 71 | 123 | Munpia, dict.cc plus, dict.cc, Wikipedia, LEO |
| Health and fitness* | 0.42 | 1.25 | 20.52 | 54.88 | 59 | 60 | SleepBot, Strava, Fitbit, Freeletics, MyFitnessPal |
| Media and video* | 0.32 | 0.46 | 38.65 | 125.42 | 47 | 118 | Video-Player, Google Play Movies, VLC, Video anzeigen, ZDF |
| Shopping* | 0.30 | 0.92 | 33.20 | 91.04 | 45 | 68 | eBay, mydealz, Amazon, brands4friends, Shpock |
| Business | 0.28 | 0.37 | 36.41 | 44.99 | 40 | 108 | Eigene Dateien, AnyConnect, POLARIS Office Viewer 5, Polaris Viewer 4.1, OfficeSuite |
| Education* | 0.22 | 0.67 | 32.85 | 90.41 | 67 | 54 | UnlockYourBrain, AnkiDroid, TUM Campus App, Duolingo, Web Opac |
| Finance* | 0.22 | 0.75 | 11.37 | 22.43 | 34 | 39 | Sparkasse, Banking 4A, Wuestenrot, YNAB, Banking |
| Weather | 0.18 | 0.24 | 10.65 | 16.66 | 17 | 74 | Weather, wetter.com, WetterOnline, WetterApp, Wetter-Widget |
| Medical* | 0.10 | 0.47 | 2.15 | 10.35 | 10 | 17 | Lady Pill Reminder, PillReminder, Pillsikum, iPhysikum, Remember Your Pill |
| Personalisation* | 0.03 | 0.12 | 4.35 | 19.08 | 14 | 22 | Dokumente, Backgrounds, Zedge, Flatastico, HD Widgets |
| Comics* | 0.00 | 0.03 | 5.33 | 38.01 | 6 | 6 | xkcd Browser, NICHTLUSTIG, Marvel Unlimited, xkcdViewer, xkcd—Now |

*Note.* $M$, mean; $SD$, standard deviation; $M_{Num.\ uses}$, average usage count across all participants and days; $SD_{Num.\ uses}$, standard deviation of average usage count across all participants and days; $M_{usage}$, average single-usage duration in seconds; $SD_{usage}$, standarddeviation of average single-usage duration in seconds; $M_{usage}$, average single-usage duration across all usages in seconds; Num. Apps, total number of apps in the category across all participants in our dataset; Num. Users, respective number of users that ever used an app from the respective category during data collection, top five apps for each category Categories not used in the final analysis are marked with an asterisk. Additionally, the 'unknown' category was removed because of its ambiguous content. The table is sorted indescending order by $M_{Num.\ uses}$.

## RESULTS

### Descriptive statistics

*Personality and demographics*
Several substantial correlations between demographic variables as well as the Big Five personality factors were present in the data. Because of deviations from Gaussian distributions in all app usage categories, we used Spearman correlations for all variables in our analysis (Yarkoni, 2010). As the highest correlation was observed between extraversion and openness ($\rho = 0.58$, $p < 0.001$), we calculated variance inflation factors (*VIF*) for both extraversion (*VIF* = 1.88) and openness (*VIF* = 1.70). Because the *VIF* was smaller than four in both cases, we proceeded with the analysis (Dormann et al., 2013; Fox & Monette, 1992) (see also A2 in Appendix A). Table 3 shows pairwise Spearman correlations of psychometrics and demographics with usage of app categories. Several associations are visible, but as the large number of calculated correlations would induce an immense multiple testing problem for statistical inference, the correlations are reported in a strictly descriptive manner. We will elaborate on the relationships, after variable selection through regression analysis.

*App usage*
In total, 2835 different apps were used by the 137 participants in our study with an average of 12.42 different apps used per day. On average, phones were turned off 27.00 times (*SD* = 26.32) with an average duration of 1.82h (*SD* = 2.06). Apps of the communication category were on average used most frequently (on average 38.48 times a day), whereas apps of the comics category were used most infrequently (0.00 times a day). Game apps show the longest average usage duration, on average 153.33s. More information about app categories as well as the top apps of each category is provided in Table 2. Please note the table is sorted by the average number of app uses per category. Also note that the pre-processing process as well as the stability selection procedure led to the exclusion of categories for all analyses. Categories not used in the final analysis are marked with an asterisk. Additionally, the 'unknown' category was removed because of its ambiguous content. The regression analyses therefore comprises solely up to 16 distinct categories, compared with 25 in the Google Play Store.

### Prediction of app usage

In this section, we report on results of the variable selection procedure as well as regression modelling. This section is divided into two parts. In the first part, results with regard to the frequency of app usage are reported. In the second part, predicted app usage durations are reported.

### Prediction of app usage frequencies

*Factor-level personality, fluid intelligence and demographics*
The variable selection procedure reported stable personality and demography predictors for a total of 13 app usage categories (Table 4). Besides gender, age and fluid intelligence, the three Big Five factors extraversion, conscientiousness and agreeableness were chosen as meaningful behavioural predictors by variable selection. Emotional stability and openness did not provide enough unique predictive value for the app usage criteria. The highest stabilities in feature selection could be observed for gender as a predictor for the use of tools (93% selected), productivity (95%), news and magazines (90%) and music and audio (99%) as well as extraversion as a predictor for the use of communication applications (94%).

The psychometric and demographic variables chosen by stability selection were modelled as predictors in generalised linear regression models using a quasi-Poisson link. In Table 4, positive as well as negative relationships with categorical app usage frequencies can be observed. Female gender, age, fluid intelligence and conscientiousness seem to be mostly negatively associated with app usage frequencies. Only older age was slightly positively associated with higher call frequencies (+4%). Women seem to use less apps related to tools (−49%), productivity (−52%), news and magazines (−60%) and music and audio (−59%). Besides the slight positive association with call frequencies, age showed rather small and mostly negative relationships with app usage frequency. Hence, one unit increase in age was negatively associated with app use in the categories business (−8%), browser (−5%) and social (−10%). Fluid intelligence was negatively associated with lifestyle app usage (−53%). Furthermore, extraversion was generally positively associated with app usage frequency. An increase of one unit in extraversion was associated with app usage increase in calls (+35%), photography (+42%) and communication (+30%). One unit increase in conscientiousness decreased the app usage frequency for games (−46%) apps. Finally, the factor agreeableness was positively associated with the use of transportation apps (+36%).

*Facet-level personality and demographics*
On the facet level, the analysis procedure was performed in the same way as on the factor level, the results are depicted in Table 4. For a more intuitive understanding of the presented relationships, results from stability selection and regression modelling are described in a combined form in this section. Additionally, we elaborate on differences between the factor-level and facet-level analyses.

Although the results of the stability selection procedure show similarities with the factor-level analysis, differences are apparent: calls as well as the games, lifestyle, browser and social application usage could not be reliably predicted with a single facet-level variable. Other relationships are in general very similar to the associations found at the Big Five factor level, as can be seen in Table 4. Further comparisons between factor-level and facet-level predictors show that associations ($\exp(\beta)$ coefficients) with app usage categories are generally higher for factor-level predictors in comparison with facet predictors. This is true for extraversion and agreeableness, compared with their respective comparison with the respective facets (sociableness and willingness to trust). Comparing facet-level models to factor models, the model fit is mostly higher (lower *DSS* values) on the factor level

Table 3. Pairwise spearman correlations between demographics, personality and app usage

| Category | Gender | Age | Education | Fluid Int. | ES | E | O | C | A |
|---|---|---|---|---|---|---|---|---|---|
| Freq. calls | 0.00 [−0.17, 0.17] | 0.16 [−0.01, 0.32] | 0.02 [−0.15, 0.19] | −0.05 [−0.22, 0.12] | 0.06 [−0.11, 0.23] | **0.33** **[0.17, 0.47]** | 0.11 [−0.06, 0.27] | −0.07 [−0.23, 0.10] | 0.07 [−0.10, 0.23] |
| Avg. dur. calls | −0.11 [−0.27, 0.06] | 0.16 [−0.01, 0.32] | 0.09 [−0.08, 0.25] | 0.09 [−0.08, 0.25] | 0.06 [−0.11, 0.23] | −0.06 [−0.23, 0.11] | −0.06 [−0.23, 0.11] | 0.01 [−0.16, 0.18] | 0.04 [−0.13, 0.21] |
| Freq. communication | −0.01 [−0.18, 0.16] | −0.09 [−0.25, 0.08] | −0.12 [−0.28, 0.05] | −0.06 [−0.23, 0.11] | −0.05 [−0.22, 0.12] | **0.27** **[0.11, 0.42]** | −0.01 [−0.18, 0.16] | −0.07 [−0.23, 0.10] | 0.04 [−0.14, 0.20] |
| Avg. dur. communication | **0.21** **[0.04, 0.36]** | −0.14 [−0.31, 0.02] | −0.14 [−0.30, 0.03] | −0.10 [−0.26, 0.07] | −0.02 [−0.19, 0.15] | −0.03 [−0.20, 0.14] | 0.07 [−0.10, 0.23] | 0.01 [−0.16, 0.18] | 0.03 [−0.14, 0.20] |
| Freq. social | 0.07 [−0.10, 0.23] | **−0.29** **[−0.44, −0.13]** | **−0.20** **[−0.36, −0.03]** | −0.14 [−0.30, 0.03] | −0.04 [−0.21, 0.13] | 0.09 [−0.08, 0.25] | −0.10 [−0.26, 0.07] | 0.01 [−0.16, 0.18] | 0.12 [−0.14, 0.20] |
| Avg. dur. social | **0.20** **[0.03, 0.36]** | **−0.25** **[−0.40, −0.09]** | **−0.26** **[−0.41, −0.10]** | −0.11 [−0.27, 0.06] | 0.00 [−0.17, 0.17] | 0.08 [−0.09, 0.25] | 0.07 [−0.10, 0.23] | 0.05 [−0.12, 0.22] | 0.12 [−0.05, 0.28] |
| Freq. photography | 0.05 [−0.12, 0.22] | −0.12 [−0.28, 0.05] | −0.17 [−0.33, 0.00] | −0.03 [−0.20, 0.14] | −0.15 [−0.31, 0.02] | 0.14 [−0.03, 0.30] | 0.03 [−0.14, 0.20] | −0.03 [−0.20, 0.14] | 0.08 [−0.09, 0.24] |
| Avg. dur. photography | **0.33** **[0.17, 0.47]** | 0.05 [−0.12, 0.22] | 0.04 [−0.13, 0.21] | **−0.18** **[−0.34, −0.01]** | −0.08 [−0.24, 0.09] | −0.08 [−0.24, 0.09] | −0.12 [−0.28, 0.05] | 0.07 [−0.17, 0.17] | 0.10 [−0.07, 0.26] |
| Freq. weather | 0.01 [−0.16, 0.18] | 0.09 [−0.12, 0.22] | 0.13 [−0.04, 0.29] | −0.06 [−0.23, 0.11] | 0.01 [−0.16, 0.18] | 0.11 [−0.06, 0.27] | −0.05 [−0.22, 0.12] | 0.11 [−0.10, 0.23] | −0.05 [−0.21, 0.13] |
| Avg. dur. weather | 0.02 [−0.15, 0.19] | 0.09 [−0.08, 0.25] | 0.17 [0.00, 0.33] | −0.02 [−0.19, 0.15] | 0.05 [−0.12, 0.22] | 0.04 [−0.13, 0.21] | −0.01 [−0.18, 0.16] | −0.12 [−0.06, 0.27] | 0.02 [−0.22, 0.12] |
| Freq. browser | −0.09 [−0.25, 0.08] | **−0.25** **[−0.40, −0.09]** | −0.14 [−0.30, 0.03] | −0.01 [−0.18, 0.16] | −0.10 [−0.26, 0.07] | 0.06 [−0.11, 0.23] | −0.03 [−0.20, 0.14] | 0.07 [−0.28, 0.05] | −0.03 [−0.15, 0.19] |
| Avg. dur. browser | **0.20** **[0.03, 0.36]** | −0.01 [−0.18, 0.16] | −0.03 [−0.30, 0.03] | −0.03 [−0.20, 0.14] | 0.03 [−0.14, 0.20] | −0.06 [−0.23, 0.11] | 0.07 [−0.10, 0.23] | 0.07 [−0.10, 0.23] | 0.00 [−0.17, 0.17] |
| Freq. transportation | 0.04 [−0.13, 0.21] | −0.17 [−0.33, 0.00] | 0.04 [−0.13, 0.21] | 0.03 [−0.14, 0.20] | **0.18** **[0.01, 0.34]** | 0.17 [0.00, 0.33] | 0.08 [−0.09, 0.24] | 0.14 [−0.10, 0.23] | **0.20** **[0.03, 0.36]** |
| Avg. dur. transportation | 0.15 [−0.02, 0.31] | −0.09 [−0.26, 0.07] | −0.09 [−0.25, 0.08] | 0.10 [−0.14, 0.20] | 0.01 [−0.03, 0.30] | 0.06 [−0.18, 0.16] | 0.06 [−0.11, 0.23] | −0.07 [−0.03, 0.30] | 0.00 [−0.04, 0.29] |
| Freq. productivity | **−0.23** **[−0.38, −0.06]** | −0.03 [−0.20, 0.14] | 0.04 [−0.13, 0.21] | −0.09 [−0.07, 0.26] | 0.01 [−0.16, 0.18] | −0.01 [−0.11, 0.23] | −0.05 [−0.22, 0.12] | −0.05 [−0.23, 0.10] | −0.04 [−0.17, 0.17] |
| Avg. dur. productivity | 0.16 [−0.01, 0.32] | 0.05 [−0.12, 0.22] | **−0.21** **[−0.36, −0.04]** | −0.09 [−0.25, 0.08] | 0.01 [−0.16, 0.18] | −0.01 [−0.18, 0.16] | 0.06 [−0.11, 0.23] | −0.05 [−0.22, 0.12] | −0.04 [−0.21, 0.13] |
| Freq. business | −0.04 [−0.21, 0.13] | **−0.21** **[−0.36, −0.04]** | −0.15 [−0.31, 0.02] | 0.02 [−0.15, 0.19] | −0.14 [−0.30, 0.03] | 0.00 [−0.17, 0.17] | −0.02 [−0.19, 0.15] | −0.11 [−0.27, 0.06] | 0.04 [−0.24, 0.09] |
| Avg. dur. business | 0.14 [−0.03, 0.30] | 0.01 [−0.16, 0.18] | −0.02 [−0.19, 0.15] | −0.02 [−0.19, 0.15] | −0.10 [−0.26, 0.07] | 0.02 [−0.14, 0.20] | 0.06 [−0.11, 0.23] | −0.01 [−0.18, 0.16] | 0.04 [−0.13, 0.21] |
| Freq. finance* | **−0.26** **[−0.41, −0.10]** | −0.12 [−0.28, 0.05] | 0.11 [−0.06, 0.27] | 0.12 [−0.05, 0.28] | 0.09 [−0.08, 0.25] | 0.01 [−0.15, 0.19] | −0.06 [−0.23, 0.11] | −0.03 [−0.20, 0.14] | −0.13 [−0.29, 0.04] |
| Avg. dur. finance* | **−0.24** **[−0.39, −0.08]** | 0.09 [−0.08, 0.25] | 0.05 [−0.28, 0.05] | −0.12 [−0.07, 0.26] | 0.13 [−0.04, 0.29] | −0.01 [−0.16, 0.18] | −0.06 [−0.23, 0.11] | 0.02 [−0.17, 0.17] | −0.14 [−0.30, 0.03] |
| Freq. news and magazines | −0.09 [−0.25, 0.08] | −0.05 [−0.08, 0.25] | 0.06 [−0.12, 0.22] | −0.12 [−0.28, 0.05] | 0.11 [−0.16, 0.18] | 0.08 [−0.18, 0.16] | −0.03 [−0.23, 0.11] | −0.07 [−0.15, 0.19] | −0.14 [−0.16, 0.18] |
| Avg. dur. news and magazines | −0.17 [−0.33, 0.00] | −0.14 [−0.22, 0.12] | −0.09 [−0.11, 0.23] | 0.03 [−0.22, 0.12] | −0.06 [−0.06, 0.27] | 0.08 [−0.09, 0.24] | −0.04 [−0.20, 0.14] | −0.14 [−0.23, 0.10] | −0.01 [−0.30, 0.03] |
| Freq. tools | **−0.28** **[−0.43, −0.12]** | −0.14 [−0.30, 0.03] | −0.09 [−0.25, 0.08] | 0.03 [−0.14, 0.20] | −0.06 [−0.23, 0.11] | 0.04 [−0.13, 0.21] | −0.04 [−0.21, 0.13] | −0.07 [−0.30, 0.03] | −0.01 [−0.18, 0.16] |

(Continues)

Table 3. (Continued)

| Category | Gender | Age | Education | Fluid Int. | ES | E | O | C | A |
|---|---|---|---|---|---|---|---|---|---|
| Avg. dur. tools | 0.14 [-0.03, 0.30] | 0.06 [-0.11, 0.23] | 0.04 [-0.13, 0.21] | -0.02 [-0.19, 0.15] | 0.03 [-0.14, 0.20] | 0.07 [-0.10, 0.23] | 0.09 [-0.08, 0.25] | -0.03 [-0.20, 0.14] | -0.02 [-0.19, 0.15] |
| Freq. games | -0.08 [-0.24, 0.09] | **-0.29** [**-0.44, -0.13**] | **-0.19** [**-0.35, -0.02**] | **0.20** [**0.03, 0.36**] | 0.02 [-0.15, 0.19] | -0.05 [-0.22, 0.12] | -0.13 [-0.29, 0.04] | -0.15 [-0.31, 0.02] | 0.02 [-0.15, 0.19] |
| Avg. dur. games | 0.10 [-0.07, 0.26] | -0.14 [-0.30, 0.03] | -0.06 [-0.23, 0.11] | 0.12 [-0.05, 0.28] | -0.02 [-0.19, 0.15] | -0.06 [-0.22, 0.12] | -0.13 [-0.29, 0.04] | -0.02 [-0.19, 0.15] | 0.11 [-0.06, 0.27] |
| Freq. entertainment | -0.15 [-0.31, 0.02] | **-0.21** [**-0.36, -0.04**] | **-0.20** [**-0.36, -0.03**] | 0.05 [-0.17, 0.17] | -0.02 [-0.19, 0.15] | 0.02 [-0.15, 0.19] | -0.01 [-0.23, 0.11] | 0.02 [-0.23, 0.11] | 0.00 [-0.20, 0.14] |
| Avg. dur. entertainment | 0.04 [-0.13, 0.21] | -0.09 [-0.25, 0.08] | **-0.20** [**-0.36, -0.03**] | 0.22 [-0.12, 0.22] | -0.06 [-0.19, 0.15] | -0.04 [-0.21, 0.13] | -0.06 [-0.18, 0.16] | 0.02 [-0.15, 0.19] | -0.17 [-0.17, 0.17] |
| Freq. education* | **-0.23** [**-0.38, -0.06**] | 0.06 [-0.19, 0.15] | 0.14 [-0.03, 0.30] | **0.22** [**0.05, 0.37**] | -0.06 [-0.23, 0.11] | -0.03 [-0.20, 0.14] | -0.06 [-0.23, 0.11] | -0.02 [-0.19, 0.15] | 0.00 [-0.17, 0.17] |
| Avg. dur. education* | -0.16 [-0.32, 0.01] | **-0.20** [**-0.36, -0.03**] | **0.22** [**0.05, 0.37**] | **0.19** [**0.02, 0.35**] | 0.01 [-0.16, 0.18] | -0.01 [-0.18, 0.16] | -0.07 [-0.23, 0.10] | 0.06 [-0.11, 0.23] | 0.04 [-0.13, 0.21] |
| Freq. books and reference | -0.06 [-0.23, 0.11] | -0.09 [-0.25, 0.08] | -0.09 [-0.25, 0.08] | 0.07 [-0.10, 0.23] | -0.12 [-0.28, 0.05] | 0.03 [-0.14, 0.20] | 0.03 [-0.14, 0.20] | -0.11 [-0.27, 0.06] | -0.06 [-0.23, 0.11] |
| Avg. dur. books and reference | -0.05 [-0.22, 0.12] | -0.07 [-0.23, 0.10] | 0.02 [-0.24, 0.09] | 0.25 [-0.11, 0.23] | -0.06 [-0.23, 0.11] | -0.09 [-0.14, 0.20] | 0.03 [-0.14, 0.20] | -0.02 [-0.19, 0.15] | -0.09 [-0.25, 0.08] |
| Freq. comics* | -0.14 [-0.30, 0.03] | -0.04 [-0.23, 0.10] | 0.07 [-0.15, 0.19] | **0.25** [**0.09, 0.40**] | 0.07 [-0.10, 0.23] | -0.09 [-0.25, 0.08] | **-0.23** [**-0.38, -0.06**] | -0.16 [-0.32, 0.01] | 0.02 [-0.15, 0.19] |
| Avg. dur. comics* | -0.14 [-0.30, 0.03] | 0.01 [-0.21, 0.13] | 0.04 [-0.10, 0.23] | **0.23** [**0.06, 0.38**] | 0.03 [-0.10, 0.23] | -0.10 [-0.25, 0.08] | **-0.19** [**-0.35, -0.02**] | -0.09 [-0.25, 0.08] | 0.06 [-0.11, 0.23] |
| Freq. travel and local | -0.10 [-0.26, 0.07] | 0.00 [-0.17, 0.17] | 0.07 [-0.13, 0.21] | 0.03 [-0.14, 0.20] | 0.10 [-0.26, 0.07] | -0.04 [-0.04, 0.29] | -0.06 [-0.23, 0.11] | -0.21 [-0.36, -0.04]? | -0.13 [-0.29, 0.04] |
| Avg. dur. travel and local | **0.24** [**0.08, 0.39**] | **-0.21** [**-0.36, -0.04**] | -0.16 [-0.32, 0.01] | 0.08 [-0.09, 0.24] | 0.03 [-0.07, 0.26] | 0.15 [-0.21, 0.13] | 0.10 [-0.07, 0.26] | 0.01 [-0.16, 0.18] | 0.14 [-0.03, 0.30] |
| Freq. music and audio | **-0.33** [**-0.47, -0.17**] | -0.02 [-0.19, 0.15] | -0.09 [-0.32, 0.01] | -0.07 [-0.23, 0.10] | 0.02 [-0.14, 0.20] | **0.20** [**0.03, 0.36**] | 0.06 [-0.19, 0.15] | -0.05 [-0.22, 0.12] | -0.04 [-0.18, 0.16] |
| Avg. dur. music and audio | -0.02 [-0.19, 0.15] | -0.02 [-0.19, 0.15] | -0.09 [-0.25, 0.08] | 0.00 [-0.23, 0.10] | -0.21 [-0.15, 0.19] | 0.05 [-0.02, 0.31] | 0.06 [-0.11, 0.23] | -0.02 [-0.22, 0.12] | -0.09 [-0.21, 0.13] |
| Freq. media and video | 0.07 [-0.10, 0.23] | 0.05 [-0.19, 0.15] | -0.08 [-0.32, 0.01] | -0.11 [-0.27, 0.06] | **-0.36** [**-0.36, -0.04**] | 0.20 [-0.14, 0.20] | -0.09 [-0.25, 0.08] | -0.10 [-0.26, 0.07] | -0.09 [-0.25, 0.08] |
| Avg. dur. media and video | 0.04 [-0.13, 0.21] | 0.05 [-0.20, 0.14] | -0.09 [-0.24, 0.09] | -0.02 [-0.19, 0.15] | -0.07 [-0.20, 0.14] | 0.05 [-0.12, 0.22] | -0.06 [-0.25, 0.08] | -0.05 [-0.22, 0.12] | -0.02 [-0.19, 0.15] |
| Freq. lifestyle | -0.09 [-0.25, 0.08] | 0.02 [-0.20, 0.14] | -0.03 [-0.25, 0.08] | -0.04 [-0.22, 0.12] | 0.04 [-0.20, 0.14] | 0.12 [-0.05, 0.28] | 0.05 [-0.23, 0.10] | 0.01 [-0.16, 0.18] | 0.10 [-0.15, 0.19] |
| Avg. dur. lifestyle* | -0.17 [-0.33, 0.00] | 0.02 [-0.15, 0.19] | 0.03 [-0.20, 0.14] | 0.06 [-0.21, 0.13] | -0.03 [-0.13, 0.21] | **0.18** [**0.01, 0.34**] | 0.05 [-0.12, 0.22] | 0.02 [-0.15, 0.19] | 0.10 [-0.07, 0.26] |
| Freq. personalisation* | -0.17 [-0.33, 0.00] | 0.02 [-0.15, 0.19] | 0.03 [-0.14, 0.20] | 0.06 [-0.11, 0.23] | -0.05 [-0.20, 0.14] | 0.01 [-0.16, 0.18] | 0.00 [-0.17, 0.17] | -0.05 [-0.22, 0.12] | 0.00 [-0.17, 0.17] |
| Avg. dur. personalisation* | -0.16 [-0.32, 0.01] | -0.11 [-0.15, 0.19] | -0.12 [-0.14, 0.20] | -0.13 [-0.12, 0.22] | -0.05 [-0.22, 0.12] | -0.05 [-0.22, 0.12] | -0.02 [-0.19, 0.15] | -0.07 [-0.23, 0.10] | -0.02 [-0.19, 0.15] |
| Freq. shopping* | **-0.18** [**-0.34, -0.01**] | -0.14 [-0.27, 0.06] | -0.15 [-0.28, 0.05] | -0.14 [-0.29, 0.04] | 0.06 [-0.11, 0.23] | 0.11 [-0.06, 0.27] | 0.04 [-0.22, 0.12] | 0.08 [-0.09, 0.24] | 0.09 [-0.13, 0.21] |
| Avg. dur. shopping* | -0.07 [-0.23, 0.10] | -0.14 [-0.30, 0.03] | -0.15 [-0.31, 0.02] | -0.14 [-0.30, 0.03] | 0.10 [-0.07, 0.26] | **0.19** [**0.02, 0.35**] | 0.04 [-0.13, 0.21] | 0.10 [-0.07, 0.26] | -0.08, 0.25] |

(Continues)

Table 3. (Continued)

| Category | Gender | Age | Education | Fluid Int. | ES | E | O | C | A |
|---|---|---|---|---|---|---|---|---|---|
| Freq. sports* | −0.16 [−0.32, 0.01] | −0.10 [−0.26, 0.07] | −0.06 [−0.23, 0.11] | 0.07 [−0.10, 0.23] | −0.06 [−0.23, 0.11] | −0.13 [−0.29, 0.04] | −0.16 [−0.32, 0.01] | −0.04 [−0.21, 0.13] | 0.01 [−0.16, 0.18] |
| Avg. dur. sports* | −0.20 [−0.32, 0.01] | −0.13 [−0.29, 0.04] | −0.08 [−0.24, 0.09] | 0.09 [−0.08, 0.25] | −0.08 [−0.24, 0.09] | −0.16 [−0.32, 0.01] | **−0.20** [**−0.36, −0.03**] | −0.06 [−0.23, 0.11] | 0.02 [−0.15, 0.19] |
| Freq. health and fitness* | **−0.20** [**−0.36, −0.03**] | 0.06 [−0.11, 0.23] | 0.12 [−0.05, 0.28] | 0.09 [−0.08, 0.25] | −0.04 [−0.21, 0.13] | 0.08 [−0.09, 0.24] | −0.03 [−0.20, 0.14] | −0.10 [−0.26, 0.07] | 0.04 [−0.13, 0.21] |
| Avg. dur. health and fitness* | −0.14 [−0.30, 0.03] | −0.21 [−0.37, −0.04] | 0.13 [−0.04, 0.29] | 0.00 [−0.17, 0.17] | −0.01 [−0.18, 0.16] | 0.08 [−0.09, 0.24] | 0.09 [−0.08, 0.25] | −0.04 [−0.21, 0.13] | −0.01 [−0.13, 0.21] |
| Freq. medical* | **0.20** [**0.03, 0.36**] | **−0.20** [**−0.36, −0.04**] | −0.17 [−0.33, 0.00] | 0.04 [−0.13, 0.21] | 0.01 [−0.16, 0.18] | 0.13 [−0.04, 0.29] | −0.15 [−0.31, 0.02] | −0.01 [−0.18, 0.16] | −0.01 [−0.18, 0.16] |
| Avg. dur. medical* | **0.18** [**0.01, 0.34**] | **−0.20** [**−0.36, −0.03**] | **−0.18** [**−0.34, −0.01**] | 0.03 [−0.14, 0.20] | 0.00 [−0.17, 0.17] | 0.11 [−0.06, 0.27] | −0.16 [−0.32, 0.01] | 0.00 [−0.17, 0.17] | −0.03 [−0.20, 0.14] |

*Note.* Fluid Int., fluid intelligence; ES, emotional stability; E, extraversion; O, openness; C, conscientiousness; A, agreeableness. Pairwise Spearman correlations between Big Five measures, fluid intelligence, demographic variables and app usage categories. Variables, which were not used in the analysis because of low variability, are marked with an asterisk. Square brackets contain 95% confidence intervals. Confidence intervals not including zero are bold. Abbreviations left to right;

in five out of eight comparable models. The *Mean DSS* rows in Table 4 provide a comparison of factor and facet level.

*Prediction of app usage durations*

In Table 5, results of the app-duration analysis are depicted. In comparison with the analysis regarding app frequencies, it becomes clear that less variables could be reliably predicted by psychometric and demographic measures. Hence, we will present the results of the app-duration analysis in a combined form later.

Variable selection identified six categories of app usage that could be reliably predicted by a total of four predictor variables. Gender, the level of education and extraversion were identified as reliable predictors in at least one category of app usage duration. Female gender was positively associated with the average duration of photography app usage (+34%), travel and local app usage (+27%) and browser usage (+21%). The level of completed education was identified as a predictor for the usage duration of productivity apps (−15%) and weather apps (+74%). Extraversion proved to be the only personality trait successful in the prediction of app usage durations on factor level. Most interestingly, extraversion was negatively associated with the average duration of calls (−29%).

In comparison with the analyses on the factor level, less categories of app usage duration could be predicted by facet-level predictors. The average durations for weather and browser applications could not be predicted by any variable. Furthermore, with the use of facet-level personality scores, the range of successful predictors was limited to gender and competence—a facet of conscientiousness. A one unit increase in competence was associated with a 23% increase in the average duration of entertainment app usage.

The comparison of models with factor-level and facet-level predictors shows that less app usage durations could be predicted when facet-level predictors were included in the variable selection procedure. Similar to the app frequency analysis, model fits are marginally higher for factor-level models in comparison with facet-level models. The *Mean DSS* rows in Table 5 depict a comparison of factor and facet-level.

**DISCUSSION**

The present results show that individual personality scores, fluid intelligence and demographic variables can predict actual behaviour, manifested as mobile app usage frequencies and durations. Our findings suggest that variations in extraversion, conscientiousness, agreeableness, fluid intelligence and gender and age are associated with increased and decreased application usage on smartphones in general. Both factor and facet models of personality were effective in the prediction of categorical app usage. However, factor-level models could predict more categories of application usage with slightly better model fit. In the following, we will discuss the various effects discovered in our data and suggest possible explanations as a vantage point for prospective research. However, note that these interpretations are drawn

Table 4. Frequencies of app usage—variable selection | prediction

| Predictors | Calls | Tools | Games | Product. | Magazines | Photo. | Comm. | Music | Business | Lifestyle | Transp. | Browser | Social |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Factor level** | | | | | | | | | | | | | |
| Gender | 0.03 | **0.93 \| 0.51** | 0.35 | **0.95 \| 0.48** | **0.9 \| 0.4** | 0.10 | 0.07 | **0.99 \| 0.41** | 0.09 | 0.11 | 0.04 | 0.33 | 0.07 |
| Age | **0.78 \| 1.04** | 0.33 | 0.34 | 0.13 | 0.10 | 0.54 | 0.23 | 0.06 | **0.88 \| 0.92** | 0.07 | 0.51 | **0.79 \| 0.95** | **0.77 \| 0.9** |
| Fludi Int. | 0.08 | 0.09 | 0.34 | 0.10 | 0.15 | 0.19 | 0.09 | 0.28 | 0.05 | 0.12 | 0.12 | 0.25 | 0.31 |
| Extraversion (E) | **0.71 \| 1.35** | 0.19 | 0.05 | 0.15 | 0.06 | **0.81 \| 1.42** | **0.94 \| 1.3** | 0.10 | 0.32 | **0.79 \| 0.47** | 0.38 | 0.30 | 0.34 |
| Conscientiousness (C) | 0.32 | 0.03 | **0.77 \| 0.54** | 0.02 | 0.03 | 0.02 | 0.11 | 0.01 | 0.13 | 0.21 | 0.09 | 0.11 | 0.09 |
| Agreeableness (A) | 0.00 | 0.04 | 0.17 | 0.02 | 0.15 | 0.07 | 0.07 | 0.17 | 0.09 | 0.20 | **0.75 \| 1.36** | 0.07 | 0.07 |
| Mean DSS | 22.60 | 25.21 | 23.73 | 21.94 | 17.92 | 21.35 | 30.15 | 20.15 | 13.33 | 17.81 | 17.54 | 24.62 | 26.35 |
| **Facet level** | | | | | | | | | | | | | |
| Gender | — | **0.86 \| 0.51** | — | **0.94 \| 0.48** | **0.87 \| 0.4** | 0.01 | 0.02 | **0.93 \| 0.41** | 0.04 | — | 0.00 | — | — |
| Age | — | 0.34 | — | 0.10 | 0.09 | 0.26 | 0.12 | 0.02 | **0.75 \| 0.92** | — | 0.38 | — | — |
| Sociableness (E2) | — | 0.08 | — | 0.12 | 0.07 | **0.77 \| 1.2** | **0.94 \| 1.15** | 0.00 | 0.08 | — | 0.13 | — | — |
| Willingness to trust (A1) | — | 0.01 | — | 0.03 | 0.00 | 0.04 | 0.01 | 0.10 | 0.11 | — | **0.7 \| 1.25** | — | — |
| Mean DSS | — | 26.27 | — | 23.39 | 23.82 | 22.58 | 29.79 | 19.27 | 14.86 | — | 17.76 | — | — |

*Note.* Abbreviations of categories stand for productivity, news and magazines, photography, communication, music and audio, transportation, from left to right; DSS, Dawid Sebastian Score. The left values are the respective probabilities of variable selection at factor or facet level, obtained with stability selection. Right values represent exp($\hat{\beta}$) coefficients from quasi-Poisson regression models between Big Five factor scores, demographics and app usage variables. Cells with only one number include selection probabilities. Criteria with empty cells (—) could not be predicted by any variable. Numbers greater than 1 represent a positive relationship, while numbers smaller than 1 represent a negative relationship. Empty cells refer to not-selected variables. Interpretation: Coefficients greater than 1 describe the percentage of increase in app usage that go along with an increase of 1 unit in the personality score (e.g. Comm. ~ extraversion: 1.3 ≙ 30% increase). Scores below 1 indicate a negative relationship (e.g. Games ~ conscientiousness: 0.54 ≙ 100 − 54 = 46% decrease) and indicate the percentage of decrease in app usage per one unit increase in the respective predictor score). Note that '1' represents male and '2' represents female labelling for the variable gender when interpreting coefficients.

post-hoc and should therefore not be readily generalised without additional cross-validation.

### Personality and app usage

*Extraversion*

The extraversion–introversion dimension is often associated with outgoing behaviour, for example, communication (Butt & Phillips, 2008; Montag et al., 2014, 2015). Our data suggest that extraversion and its facet sociableness are related to increased application usage in categories related to calls, photography and communication. A higher frequency of communication app usage and calls is in line with previous literature reporting higher numbers of communication activities (Chittaranjan et al., 2013; Kim et al., 2015; Montag et al., 2014), as well as a higher usage frequencies of the WhatsApp messenger (Montag et al., 2015) for people with higher scores in extraversion. In contrast to the studies of Montag et al. (2014) and Butt and Phillips (2008), our results suggest that calls tend to be shorter for more extraverted people. This finding might reflect extraverts sensitivity for rewards—possibly gained through quantity, not necessarily through quality of social interaction (Ashton, Lee, & Paunonen, 2002).

Furthermore, our results show a positive relationship between extraversion/sociableness and the usage frequency of photography apps. A similar result was found in previous studies reporting increased photo uploads and photo sharing associated with higher values in extraversion (Eftekhar, Fullwood, & Morris, 2014; Hunt & Langstedt, 2014). Xu et al. (2016) found no such association but a relationship between photography apps and conscientiousness as well as emotional stability. This hints towards the difference between our approach and the investigation by Xu et al. They related only the installation of apps to personality traits, and it is possible that conscientious people tend to install camera apps in order to be able to take pictures in the event they want to extraverts, on the other hand, may be more likely to actually use the camera very often. Thus, even though the findings may seem contradictory at first, they point towards the different kinds of implications derived from different methodological approaches.

In general, extraversion was also the personality trait that shows the highest number of positive associations with various categories of app usage. This might also reflect an aspect of the personality trait that is described in the literature as the need for external stimulation (Butt & Phillips, 2008; Eysenck, 1967). People might particularly satisfy this need through communication, or other channels such as the use of photography apps. In line with this argumentation are the results of Chittaranjan et al. (2013), who reported higher usage frequencies of entertainment apps for extraverts.

*Agreeableness*

The personality trait agreeableness describes how cooperative and socially harmonic people tend to act (Graziano & Tobin, 2009). Our data indicate that agreeableness and the respective facet willingness to trust are related to the use of transportation apps. Although it is difficult to draw

C. Stachl et al.

Table 5. Durations of app usage—variable selection | prediction

| Predictors | Calls | Productivity | Photography | Travel and local | Weather | Browser | Entertainment |
|---|---|---|---|---|---|---|---|
| **Factor level** | | | | | | | |
| Gender | 0.12 | 0.53 | **0.91 \| 1.34** | **0.85 \| 1.27** | 0.12 | **0.88 \| 1.21** | — |
| Education | 0.17 | **0.91 \| 0.85** | 0.14 | 0.19 | **0.81 \| 1.74** | 0.12 | — |
| Extraversion (E) | **0.72 \| 0.71** | 0.09 | 0.09 | 0.03 | 0.06 | 0.17 | — |
| Mean DSS | 22.11 | 10.36 | 10.67 | 13.99 | 12.92 | 14.73 | — |
| **Facet level** | | | | | | | |
| Gender | — | — | **0.85 \| 1.34** | **0.74 \| 1.27** | — | — | 0.05 |
| Competence (C1) | — | — | 0.04 | 0.03 | — | — | **0.73 \| 1.23** |
| Mean DSS | — | — | 10.75 | 14.39 | — | — | 19.61 |

*Note.* The left values are the respective probabilities of variable selection at factor or facet level, obtained with stability selection. Right values represent $\exp(\hat{\beta})$ coefficients from quasi-Poisson regression models between Big Five factor scores, demographics and app usage variables. Criteria with empty cells (—) could not be predicted by any variable. Numbers greater than 1 represent a positive relationship, while numbers smaller than 1 represent a negative relationship. Empty cells refer to not-selected variables. Interpretation: Coefficients greater than 1 describe the percentage of increase in app usage that go along with an increase of 1 unit in the personality score (e.g. Entertainment ~ competence: $1.23 \doteq 23\%$ increase). Scores below 1 indicate a negative relationship (e.g. Calls ~ extraversion: $0.71 \doteq 100 - 71 \doteq 29\%$ decrease) and indicate the percentage of decrease in app usage per one unit increase in the respective predictor score). Note that '1' represents male and '2' represents female labelling for the variable gender when interpreting coefficients. DSS, Dawid Sebastian Score.

conclusions about this particular association, it could be interesting to investigate whether higher transportation app usage in agreeable people is related to them being more trusting of others for means of transportation or to help others finding adequate connections. This is backed by the facet willingness to trust being the strongest predictor and the notion that agreeable people tend to be more pro-social in the sense of being tolerable of others, preferring cooperation over competition ([35]). Other research suggests that agreeable people tend to spend more time at public places (e.g. cafés) and less time at home (Mehl et al., 2006) and it may be speculated that they visit people rather than having them come to their homes. This in return could result in more time spent in public transportation.

*Conscientiousness*
Conscientiousness seems to predict the amount of gaming as well as the use of entertainment apps on smartphones. The lower usage of gaming apps in highly conscientious people supports the notion that conscientious people are more focused on their tasks and less likely to engage in procrastination activities (Lee, Kelly, & Edwards, 2006). Previous studies only reported on associations of low agreeableness and extraversion with gaming app usage (Phillips, Butt, & Blaszczynski, 2006) and installations (Xu et al., 2016).

The positive association of conscientiousness/competence and the entertainment app usage durations (e.g. YouTube) could not necessarily be expected, as previous findings reported conscientious individuals to make less frequent use of entertainment apps (Chittaranjan et al., 2013). Possibly, but it can only be speculated, people scoring higher in competence tend to use entertainment apps for longer durations without interruptions because they reserved specific time slots for leisure activities.

*Openness and emotional stability*
The personality factors emotional stability and openness were not predictive for any app usage categories, neither on the factor nor the facet level. Emotional stability or neuroticism is a personality trait that is defined through feelings and emotions rather than actions (John & Robins, 1993; Vazire, 2010) and has been negatively associated with behavioural restraint (Hirsh, Deyoung, & Peterson, 2009). Emotional stability, therefore, is a dimension that is not easily observable and evaluable and has, thus far, mostly linked to text messages (Butt & Phillips, 2008), particularly incoming ones (Chittaranjan et al., 2013). Symptoms of depression (positively associated with very low scores of emotional stability) are hard to detect for that reason (Mehl, 2006). For the same reason, it is not surprising that no consistent associations with emotional stability could be observed with our behavioural logging technique. However, considering the association of neuroticism and its link to depression (Hodgins & Ellenbogen, 2003; Ormel et al., 2013) as well as the link between reduced activity and social contact associated with depression (Burton et al., 2013), variables in relation to these dimensions could be retrieved from data logs related to movement (e.g. GPS) in prospective studies (Saeb et al., 2015).

Openness is considered to be the most heterogeneous Big Five dimension (DeYoung, 2015), related to both intellectual abilities and exploratory behaviour. The missing predictiveness of openness for app usage is somewhat surprising, as it has been linked to smartphone usage before (Chittaranjan et al., 2013). As depicted in Table A2 in Appendix A, openness shares many variance with all other personality factors, highlighting the heterogeneity of the construct. This point also hints towards the problem of inter-correlations of the personality factors (a finding also reported by Chittaranjan et al., 2013), which is discussed in detail further later. Some authors even argue that extraversion and openness could be combined to a single personality dimension related to the engagement in behaviour and the incorporation of new environmental information (Hirsh et al., 2009). However, Spearman correlations in Table 3 do not show any significant correlations between openness and app usage, suggesting other reasons in our case. Possibly, openness can only be related to more specific behaviours

unlike the broad app usage categories used in this study, because of its heterogeneity.

Furthermore, the missing predictiveness of both openness and emotional stability with app usage reflects the picture painted by previous literature with only a few associations between these traits and particular smartphone usage (see Table 1 for an overview).

### Personality on factor and facet level

In order to predict behavioural criteria from individual personality scores, either broad factor values (e.g. extraversion) or facet scores (e.g. sociableness) can be used. Some previous research suggests that personality facet measures provide independent prediction value in relation to behavioural criteria in addition to factor-level scores (Ziegler et al., 2014). However, disagreement is prevalent in current research concerning this topic (Ashton et al., 2014; Salgado et al., 2013). In particular, uncertainty remains with regard to whether factor-level or facet-level scores are better for the prediction of behavioural categories. Even though the results are not completely uniform, the present investigation is able to shed some light on this issue.

On the one hand, model fits as well as directions of effects are similar between factor and facet-level models. On the other hand, some categories of app usage could only be predicted with personality scores on factor level. Furthermore, the direct comparison between both levels shows marginally better (i.e. lower) *DSS* scores for factor-level models. Previous literature partially reports higher predictive performance of personality facets over Big Five factor scores on behavioural criteria (Anglim & Grant, 2014; Paunonen & Ashton, 2001). Although our results technically contradict this notion, it cannot be concluded that factor-level scores are generally better in the prediction of behaviour. As suggested by the analysis of app usage durations, facets can contribute uniquely to the prediction of categorical behaviour (Table 5).

However, it is also not certain that personality scores on factor level are superior to scores on the facet-level variable in prediction contexts, as previously reported (Hogan & Roberts, 1996). A likely candidate for the superiority of the factor scores regarding predictive performance is measurement symmetry, which states that associations of measured constructs should be strongest on the same level of aggregation (see, e.g. Wittmann, 2012). Strong relationships with facets would therefore be expected for rather specific behaviour while an association with the broader personality factors should show for rather aggregated measures of smartphone behaviour. Because, in the present investigation, specific app usages were summed up in categories, the stronger predictive performance of personality factors compared with facets seems not that surprising as they contain more variance in comparison with single facets. Notably, the two predictive facets, sociableness and willingness to trust were related to the rather narrow categories photography, communication and transportation, respectively. It therefore remains to be investigated whether more narrow real behavioural criteria, such as single app usage or even isolated behaviours

performed within apps are better predictable by facets representing narrow traits (Hogan & Roberts, 1996).

Finally, in our study, we only used the most stable predictors, according to the stability selection procedure, for each category of app usage. With a regularisation approach, one could use all predictors simultaneously—possibly resulting in a superior combined predictive performance for facet-level scores. Nevertheless, our results do suggest that if one had to choose the most promising personality trait for the prediction of categorical app usage, a factor-level variable might work best.

### Intelligence

As visible in Tables 4 and 5, fluid intelligence was only selected as a stable negative predictor for the usage frequency of lifestyle apps. Therefore, based on the collected data, it seems that fluid intelligence is not widely associated with variations in app usage behaviour. Yet it has to be kept in mind that the sample of the present investigation was relatively homogeneous in that it was young and well educated, which may have resulted in a restricted range of the variation in fluid intelligence. This, in turn, may have lowered the statistical associations of fluid intelligence, especially considering the rather conservative and competitive stability selection in this investigation. Nevertheless, fluid intelligence still turned out to be related to lifestyle app usage. A possible reason may be that this category was topped by the popular dating app *Tinder* while fluid intelligence has been (loosely) negatively related to the number of sexual partners (e.g. Fergusson, John Horwood, & Ridder, 2005). However, this is a very sensitive post-hoc explanation, which should be interpreted with care. Of course, it may seem surprising that fluid intelligence was not associated more widely with the use of smartphone apps, but several aspects should be considered when interpreting this finding: the association between technology acceptance and intelligence has been predominantly established based on the cognitive decline during age, which has been held responsible for the lesser acceptance of new technologies in the elderly (Venkatesh et al., 2012). Considering the narrow age range of our sample, possible associations between smartphone use and intelligence may have been obscured by even stronger variance restriction because of this characteristic of the sample.

Also, as the link between intelligence and smartphone use may be mediated by personality traits (Lee et al., 2014), the rather conservative stability selection procedure applied in this study may have had a strong influence on the picture of results: even though a link between app usage and intelligence may exist, it was just never as strong as the link between personality and app usage in this study, so that intelligence was never selected as a predictor. Finally, fluid intelligence is a construct expected to be related to how successful or efficient somebody would perform complex tasks. From this perspective, it seems rather unsurprising that only one association with frequencies and durations of broad categories of app usage was found. It is possible a measure of crystallised intelligence may have performed better in predicting apps related to knowledge, such as in the news

and magazines or the books and reference categories. Nevertheless, prospective studies should investigate possible associations between fluid intelligence and performance of app usage at a more fine-grained level.

### Demographics and app usage

Besides personality traits and fluid intelligence, gender, age and level of education predicted app usage behaviour in several categories. Most notably, our results suggest that female gender is associated with less frequent app usage in several categories but longer average usage durations in others.

Our data suggest that on average women used apps related to tools, productivity, news and magazines and music and audio, roughly half as often as men did. The finding regarding music and audio is in accordance with results of Kim et al. (2015), who also reported lower frequencies of entertainment application use (including music) for women. It is unclear why we observed this effect with such stability. As our logging method does not include music consumption on secondary devices (such as mp3 players or iPods), gender-specific differences in listening behaviours could therefore be missed. Furthermore, together with lower frequencies of news and magazines app usage, this could also be related to technology acceptance (Sherman et al., 2000), as many apps in the music and audio category were related to novel music streaming services (Table 2). According to these statistics (based on an American and UK sample) by Globalwebindex (2015), differences between men and women in the adoption of paid music streaming services exist, however, not at the magnitude observed in our sample.

While the finding regarding tools seems not so surprising, considering the role of gender in the emergence of technology (Lerman, Mohun, & Oldenziel, 1997), the finding that women less frequently used productivity-related apps cannot easily be interpreted. Although gender differences in productivity have been investigated in previous research (Leahey, 2006; Reed, Enders, Lindor, McClees, & Lindor, 2011), this result does not allow for causality.

In addition to lower usage frequencies in some app categories for women, our data also suggests that the female gender is associated with longer usage durations of photography, travel and local and browser apps. As we have no conclusive explanation for these findings, we refrain from reading the tea leaves here.

Several small negative effects of age on business, browser and social app usage were present in our data. Additionally, a slight positive effect on the number of calls was observed in our data set. This is in accordance with results of Kim et al. (2015) who also reported lower app usage frequencies for relation apps (e.g. messaging and social) and the results of Montag et al. (2014), who reported positive correlations between call variables and age. In the study of Kim et al. (2015) as well as in the present study, age does not predict app usage to a large degree. However, it is important to point out that it is very likely that these small effects in the present study as well in the study of Kim et al. (2015) are at least partially caused by sample selection effects. Kim et al. (2015) reported a large negative correlation between

age and smartphone ownership (Kim et al., 2015). Furthermore, in our study, only participants with a compatible Android smartphone could participate. As smartphone ownership declines with age and the variation in the data only describes effects of mostly younger participants, the current results cannot rule out different app usage behaviour of older people in general. Age was mostly selected as an important predictor when in competition with factor-level personality scores in predicting app usage frequency. Regarding app usage durations, age was not stably predictive for any of the categories. Because of the limited age range in our sample, the reader is advised not to over-interpret the observed effects of age on app usage.

In contrast to results reported by Kim et al. (2015), no big effects of education level on app usage were found in the present study. The only two effects of education show a negative association with the usage duration of productivity apps and a strong positive association with the usage duration of weather apps. It has to be kept in mind here that Kim et al. (2015) used self-reports, which may differ from the actual log data, due to memory distortion (Lin et al., 2015) and that the present analysis included fluid intelligence, which is known to be strongly related to education (e.g. Mayer, 2000) and may have restricted the unique variance that education shared with the criteria.

### Limitations

There are important limitations to be noted. Our sample was collected purely from the German population in Munich with age and education not perfectly representative of the general population. However, as smartphone usage is less prevalent with older people (Kim et al., 2015), our sample might not be too different from the normal population of smartphone users in this regard. Moreover, usage patterns might differ when compared with, for example, samples from other cities and countries. For instance, availability and popularity of public transportation impacts the use of apps in the related category. Differences in the samples' cultural backgrounds and countries can also be expected to be reflected in app usage. However, although some variation in app usage is to be expected, many popular apps for common tasks are globally available or have popular regional equivalents. Furthermore, associations between app usage and age were similar to previous results even though smaller statistical associations were to be expected in a more homogeneous sample.

Our results also suggest that app usage durations were generally harder to predict in comparison with frequencies of app usage. It is most likely duration measures were not sensitive enough to pick up on psychometric and demographic variance. For example, average durations would not distinguish between a person using an app once for 10 min and a person using the same app 100 times for an average duration of 10 min, as the frequency of app use is not included in the duration analyses. Furthermore, the intention of an app use (e.g. gaming for procrastination purposes) is not adequately grasped by duration-based measures, an aspect likely informative about psychological variables. Usage durations might not sufficiently distinguish between active behaviours

(e.g. having a conversation via phone) and passive behaviours (e.g. listening to music via phone). This is also in accordance with previous research showing predictive performance for active behaviours, such as calling and messaging (e.g. Butt & Phillips, 2008; Montag et al., 2014, 2015).

Furthermore, we observed substantial inter-correlations between the Big Five scores in the present data set. Although based on *VIF* as well as previous research (Van der Linden, te Nijenhuis, & Bakker, 2010), we conclude that the magnitude of the observed correlations is not worrisome, it should be noted that this might have affected our results. Openness for example shares many variance with emotional stability and extraversion. Because of the nature of our variable selection procedure, this could have resulted in some variables (e.g. openness and emotional stability) no being selected as the top predictor for app usage when competing for example against extraversion.

Also, it has to be noted that the present study investigated categorical app usage—a fraction of activities traceable on smartphones. It is likely that the inclusion of additional parameters (e.g. GPS, music consumption and word use; Yarkoni, 2010) and single app usage will make it possible to establish more relationships with personality traits. Furthermore, we want to highlight again that one has to be careful with drawing post-hoc conclusions based on the observed relationships. While our results indicate avenues for both personality research in academics as well personalisation research in industrial settings, we understand the reported relationships as promising starting points for closer investigation.

### Outlook and implications

The present study shows how everyday app usage on smartphones is predicted by individual personality traits on facet and factor level, fluid intelligence and demographic variables. In accordance with previous research, our data suggest that three personality traits (extraversion, conscientiousness and agreeableness) and demographic variables are particularly predictive for specific behavioural categories of app usage.

In addition to further investigating the differences in results between the present study and previous investigations (e.g. extraversion and call duration), future studies could aim at identifying specific usage patterns such as single application usage, as the specific content of apps is likely to be descriptive of the user's personality. Sending images via instant messaging services has for example been related to personality traits (Hunt & Langstedt, 2014). However, the actual content of a photography might be even more interesting, as it should indicate which information that person cares to share with others (Amiel & Sargent, 2004). This approach could, for example, help to also associate openness and emotional stability with traceable information. Furthermore, as different categories of app usage cannot ensure that discovered associations can be generalised, it would be desirable to establish a catalogue of apps as well as their corresponding categorizations. Newer research methodologies such as

large-scale crowd-sourcing could very well be used for this endeavour.

The present results will hopefully stimulate further research, involving actual behaviour and help to relate psychological research better to everyday problems and behaviours. First promising studies involve the prediction of depression (BinDhim et al., 2014; Canzian & Musolesi, 2015; Saeb et al., 2015), bipolar disorder (Grunerbl et al., 2015) and smartphone addiction (Lin et al., 2015).

However, beyond specific implications, this work also provides a general example of how data on behavioural acts can be efficiently collected for personality studies from data logs. App-log collection could be applied to many areas of psychological science and effectively address the current lack of real behaviour in the field of personality research (Back et al., 2009; Baumeister et al., 2007; Furr, 2009). Validation studies of self-report measures for example could be improved with additional data about actual behaviour. Additionally, studies could easily combine behavioural data with, for instance, experience sampling methods in order to enrich subjective with objective data (as suggested by Lin et al., 2015). In general, prospective studies should aim at both replicating the present findings and extending our understanding of the relationship between personality and smartphone use beyond the currently sparse literature with both larger exploratory mobile sensing studies and pre-registered confirmatory experiments.

### CONCLUSION

The present study demonstrates how actual behaviour associated with personality traits is aggregated over different situations and times reflected by different patterns of app usage on smartphones. Extraverted people seemed to satisfy their associated need for stimulation (Butt & Phillips, 2008; Eysenck, 1967) through a high number of calls and intensive use of photography apps. Agreeableness could be related to more use of transportation apps, possibly reflecting cooperative behaviour of agreeable people (Graziano & Tobin, 2009) in that they visit others rather than expecting to be visited. High conscientiousness was associated with low usage of gaming apps, probably reflecting their lack of procrastination behaviour (Lee et al., 2006). Openness and emotional stability, on the other hand, were not associated with any particular kind of app usage, possibly reflecting their heterogeneity and subtleness, respectively. Furthermore, differences in demographics and fluid intelligence also seemed to manifest themselves in different smartphone usage, such as the less extensive use of lifestyle apps by people with higher fluid intelligence.

We therefore conclude that self-reported personality traits on factor level and facet level, as well as fluid intelligence and demographic variables are manifested in a range of different categories of app usage behaviour on smartphones. Further utilisation of traceable user behaviour in psychological research practices could support the incorporation of actual behaviour in personality research. This is likely to help further improving the relevance of psychological and in

C. Stachl et al.

particular personality research, as personality traits and acts of everyday life could be associated in a straightforward manner.

## ACKNOWLEDGEMENTS

## SUPPORTING INFORMATION

Additional Supporting Information may be found online in the supporting information tab for this article.

## REFERENCES

Amiel, T., & Sargent, S. L. (2004). Individual differences in Internet usage motives. *Computers in Human Behavior*, *20*, 711–726. https://doi.org/10.1016/j.chb.2004.09.002.

Anglim, J., & Grant, S. L. (2014). Incremental criterion prediction of personality facets over factors: Obtaining unbiased estimates and confidence intervals. *Journal of Research in Personality*, *53*, 148–157. https://doi.org/10.1016/j.jrp.2014.10.005.

Arendasy, M. (2009). *BFSI: Big-Five Struktur-Inventar (test & manual)*. Mödling: SCHUHFRIED GmbH.

Arendasy, M., Hornke, L., Sommer, M., Häusler, J., Wagner-Menghin, M., Gittler, G., & Wenzl, M. (2009). *Manual intelligence structure battery (INSBAT)*. Mödling: Schuhfried Gmbh.

Ashton, M. C., Lee, K., & Paunonen, S. V. (2002). What is the central feature of extraversion? Social attention versus reward sensitivity. *Journal of Personality and Social Psychology*, *83*, 245–252. https://doi.org/10.1037/0022-3514.83.1.245.

Ashton, M. C., Paunonen, S. V., & Lee, K. (2014). On the validity of narrow and broad personality traits: A response to Salgado, Moscoso, and Berges (2013). *Personality and Individual Differences*, *56*, 24–28. https://doi.org/10.1016/j.paid.2013.08.019.

Back, M. D., Schmukle, S. C., & Egloff, B. (2009). Predicting actual behavior from the explicit and implicit self-concept of personality. *Journal of Personality and Social Psychology*, *97*, 533–548. https://doi.org/10.1037/a0016229.

Baumeister, R. F., Vohs, K. D., & Funder, D. C. (2007). Psychology as the science of self-reports and finger movements: Whatever happened to actual behavior? *Perspectives on Psychological Science*, *2*, 396–403. https://doi.org/10.1111/j.1745-6916.2007.00051.x.

Bernaards, C. A., & Jennrich, R. I. (2005). Gradient projection algorithms and software for arbitrary rotation criteria in factor analysis. *Educational and Psychological Measurement*, *65*, 676–696. https://doi.org/10.1177/0013164404272507.

BinDhim, N. F., Shaman, A. M., Trevena, L., Basyouni, M. H., Pont, L. G., & Alhawassi, T. M. (2014). Depression screening via a smartphone app: Cross-country user characteristics and feasibility. *Journal of the American Medical Informatics Association*, *22*, 29–34. https://doi.org/10.1136/amiajnl-2014-002840.

Burton, C., McKinstry, B., Szentagotai TÄƒtar, A., Serrano-Blanco, A., Pagliari, C., & Wolters, M. (2013). Activity monitoring in patients with depression: A systematic review. *Journal of Affective Disorders*, *145*, 21–28. https://doi.org/10.1016/j.jad.2012.07.001.

Butt, S., & Phillips, J. G. (2008). Personality and self reported mobile phone use. *Computers in Human Behavior*, *24*, 346–360. https://doi.org/10.1016/j.chb.2007.01.019.

Canzian, L., & Musolesi, M. (2015). Trajectories of depression. In *Proceedings of the 2015 ACM International Joint Conference on pervasive and ubiquitous computing—ubiComp '15* (pp. 1293–1304). New York, USA: ACM Press.

Carroll, J. B. (2003). The higher-stratum structure of cognitive abilities. Current evidence supports g and about ten broad factors. In *The scientific study of general intelligence: Tribute to Arthur Jensen*, (5–21). doi:https://doi.org/10.1016/B978-008043793-4/50036-2

Chittaranjan, G., Blom, J., & Gatica-Perez, D. (2013). Mining large-scale smartphone data for personality studies. *Personal Ubiquitous Computing*, *17*, 433–450. https://doi.org/10.1007/s00779-011-0490-1.

Core Team, R. (2016). *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing Retrieved from: https://www.R-project.org/.

Costa, P. T. and R. R. McCrae (1992). *Revised NEO personality inventory (NEO PI-R) and NEO five-factor inventory (NEO-FFI)*. Psychological Assessment Resources Inc. doi:https://doi.org/10.1037//1040-3590.4.1.5

Czado, C., Gneiting, T., & Held, L. (2009). Predictive model assessment for count data. *Biometrics*, *65*, 1254–1261. https://doi.org/10.1111/j.1541-0420.2009.01191.x.

De Montjoye, Y.-A., J. Quoidbach, F. Robic, and A. Pentland (2013). Predicting personality using novel mobile phone-based metrics. In *Proceedings of the 6th International Conference on social computing, behavioral-cultural modeling and prediction* (48–55). SBP'13, Washington, DC: Springer-Verlag. doi:https://doi.org/10.1007/978-3-642-37210-0_6

DeYoung, C. G. (2015). Openness/intellect: A dimension of personality reflecting cognitive exploration. In M. Mikulincer, P. R. Shaver, M. L. Cooper, & R. J. Larsen (Eds.), *APA handbook of personality and social psychology, volume 4: Personality processes and individual differences* (pp. 369–399), APA handbooks in psychology. Washington, DC, US: American Psychological Association.

Doebler, P. (2015). *Mada: Meta-analysis of diagnostic accuracy*. R package version 0.5.7. Retrieved from https://CRAN.R-project.org/package=mada

Dormann, C. F., Elith, J., Bacher, S., Buchmann, C., Carl, G., Carre, G., Marquez, J. R. G., … Lautenbach, S. (2013). Collinearity: A review of methods to deal with it and a simulation study evaluating their performance. *Ecography*, *36*, 27–46. https://doi.org/10.1111/j.1600-0587.2012.07348.x.

Eaton, L. G., & Funder, D. C. (2003). The creation and consequences of the social world: An interactional analysis of extraversion. *European Journal of Personality*, *17*, 375–395. https://doi.org/10.1002/per.477.

Eftekhar, A., Fullwood, C., & Morris, N. (2014). Capturing personality from Facebook photos and photo-related activities: How much exposure do you need? *Computers in Human Behavior*, *37*, 162–170. https://doi.org/10.1016/j.chb.2014.04.048.

Epstein, S. (1983). Aggregation and beyond: Some basic issues on the prediction of behavior. *Journal of Personality*, *51*, 360–392. https://doi.org/10.1111/j.1467-6494.1983.tb00338.x.

Eysenck, H. J. (1967). *The biological basis of personality*. Springfield, Illinois: Thomas Publishing.

Fergusson, D. M., John Horwood, L., & Ridder, E. M. (2005). Show me the child at seven II: Childhood intelligence and later outcomes in adolescence and young adulthood. *Journal of Child*

*Psychology and Psychiatry*, *46*, 850–858. https://doi.org/10.1111/j.1469-7610.2005.01472.x.

Fox, J., & Monette, G. (1992). Generalized collinearity diagnostics. *Journal of the American Statistical Association*, *87*, 178–183. https://doi.org/10.1080/01621459.1992.10475190.

Friedman, J., Hastie, T., & Tibshirani, R. (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, *33*, 1–22. https://doi.org/10.1359/JBMR.0301229.

Funder, D. C. (2001). Personality. *Annual Review of Psychology*, *52*, 197–221. https://doi.org/10.1146/annurev.psych.52.1.197.

Furr, R. M. (2009). Personality psychology as a truly behavioural science. *European Journal of Personality*, *23*, 369–401. https://doi.org/10.1002/per.724.

Ghosh, D., & Vogt, A. (2012). Outliers: An evaluation of methodologies. *Joint Statistical Metings*, 3455–3460 Retrieved from: http://ww2.amstat.org/sections/srms/Proceedings/y2012/files/304068_72402.pdf.

Globalwebindex (2015). *1 in 4 Spotify users pay for the service*. Retrieved March 18, 2015, from: http://www.globalwebindex.net/blog/1-in-4-spotify-users-pay-for-the-service

Google, I. (2016). *Android apps on Google Play*. Retrieved June 16, 2015, from: https://play.google.com/store/apps?hl=en

Gosling, S. D., & Mason, W. (2015). Internet research in psychology. *Annual Review of Psychology*, *66*, 877–902. https://doi.org/10.1146/annurev-psych-010814-015321.

Graziano, W. G., & Tobin, R. M. (2009). Agreeableness. In M. R. Leary, & R. H. Hoyle (Eds.), *Handbook of individual differences in social behavior* (pp. 46–61). New York, NY, US: Guilford Press.

Grunerbl, A., Muaremi, A., Osmani, V., Bahle, G., Ohler, S., Troster, G., Mayora, O., … Lukowicz, P. (2015). Smartphone-based recognition of states and state changes in bipolar disorder patients. *IEEE Journal of Biomedical and Health Informatics*, *19*, 140–148. https://doi.org/10.1109/JBHI.2014.2343154.

Harari, G. M., Gosling, S. D., Wang, R., & Campbell, A. T. (2015). Capturing situational information with smartphones and mobile sensing methods. *European Journal of Personality*, *29*, 509–511.

Harari, G. M., Lane, N. D., Wang, R., Crosier, B. S., Campbell, A. T., & Gosling, S. D. (2016). Using smartphones to collect behavioral data in psychological science: Opportunities, practical considerations, and challenges. *Perspectives on Psychological Science*, *11*, 838–854. https://doi.org/10.1177/1745691616650285.

Hirsh, J. B., Deyoung, C. G., & Peterson, J. B. (2009). Metatraits of the big five differentially predict engagement and restraint of behavior. *Journal of Personality*, *77*, 1085–1102. https://doi.org/10.1111/j.1467-6494.2009.00575.x.

Hodgins, S., & Ellenbogen, M. (2003). Neuroticism and depression. *The British Journal of Psychiatry*, *182*, 79–80. https://doi.org/10.1192/bjp.182.1.79.

Hofner, B., Boccuto, L., & Göker, M. (2015). Controlling false discoveries in high-dimensional situations: Boosting with stability selection. *BMC Bioinformatics*, *16*, 144. https://doi.org/10.1186/s12859-015-0575-3.eprint:arXiv:1411.1285v1.

Hofner, B. and T. Hothorn (2015). *Stabs: Stability selection with error control*. R package version R package version 0.5–1. Retrieved from: http://CRAN.R-project.org/package=stabs

Hogan, J., & Roberts, B. W. (1996). Issues and non-issues in the fidelity-bandwidth trade-off. *Journal of Organizational Behavior*, *17*, 627–637. https://doi.org/10.1002/(SICI)1099-1379(199611)17:6<627::AID-JOB2828>3.0.CO;2-F.

Horn, J. L., & Noll, J. (1997). Human cognitive capabilities: Gf-Gc theory. In D. P. Flanagan, J. L. Genshaft, & P. L. Harrison (Eds.), *Contemporary intellectual assessment: Theories, tests, and issues* (pp. 53–91). New York, NY, US: Guilford Press.

Hunt, D. S., & Langstedt, E. (2014). The influence of personality factors and motives on photographic communication. *The Journal of Social Media in Society*, *3*, Retrieved from http://www.thejsms.org/tsmri/index.php/TSMRI/article/view/68.

Jackson, J. J., Wood, D., Bogg, T., Walton, K. E., Harms, P. D., & Roberts, B. W. (2010). What do conscientious people do? Development and validation of the behavioral indicators of conscientiousness (BIC). *Journal of Research in Personality*, *44*, 501–511. https://doi.org/10.1016/j.jrp.2010.06.005.

John, O. P., & Robins, R. W. (1993). Determinants of interjudge agreement on personality traits: The big five domains, observability, evaluativeness, and the unique perspective of the self. *Journal of Personality*, *61*, 521–551. https://doi.org/10.1111/j.1467-6494.1993.tb00781.x.

Kim, Y., Briley, D. A., & Ocepek, M. G. (2015). Differential innovation of smartphone and application use by sociodemographics and personality. *Computers in Human Behavior*, *44*, 141–147. https://doi.org/10.1016/j.chb.2014.11.059.

Lane, N. D., Miluzzo, E., Lu, H., Peebles, D., Choudhury, T., & Campbell, A. T. (2010). A survey of mobile phone sensing. *IEEE Communications Magazine*, *48*, 140–150. https://doi.org/10.1109/MCOM.2010.5560598.

Lane, W., & Manner, C. (2011). The impact of personality traits on smartphone ownership and use. *International Journal of Business & Social Science*, *2*, 22–28.

Leahey, E. (2006). Gender differences in productivity: Research specialization as a missing link. *Gender & Society*, *20*, 754–780. https://doi.org/10.1177/0891243206293030.

Lee, D. G., Kelly, K. R., & Edwards, J. K. (2006). A closer look at the relationships among trait procrastination, neuroticism, and conscientiousness. *Personality and Individual Differences*, *40*, 27–37. https://doi.org/10.1016/j.paid.2005.05.010.

Lee, Y.-K., Chang, C.-T., Lin, Y., & Cheng, Z.-H. (2014). The dark side of smartphone usage: Psychological traits, compulsive behavior and technostress. *Computers in Human Behavior*, *31*, 373–383. https://doi.org/10.1016/j.chb.2013.10.047.

Lerman, N. E., Mohun, A. P., & Oldenziel, R. (1997). Versatile tools: Gender analysis and the history of technology. *Technology and Culture*, *38*, 1–8. https://doi.org/10.2307/3106781.

Lin, Y.-H., Lin, Y.-C., Lee, Y.-H., Lin, P.-H., Lin, S.-H., Chang, L.-R., Tseng, H.-W., … Kuo, T. B. (2015). Time distortion associated with smartphone addiction: Identifying smartphone addiction via a mobile application (app). *Journal of Psychiatric Research*, *65*, 139–145. https://doi.org/10.1016/j.jpsychires.2015.04.003.

Masters, G. N. (1982). A rasch model for partial credit scoring. *Psychometrika*, *47*, 149–174. https://doi.org/10.1007/BF02296272.

Mayer, R. E. (2000). Intelligence and education. In R. Sternberg (Ed.), *Handbook of intelligence* (pp. 519–533). New York, NY, US: Cambridge University Press.

McCrae, R. R., Costa, P. T. Jr., & Martin, T. A. (2005). The NEO-PI-3: A more readable revised NEO personality inventory. *Journal of Personality Assessment*, *84*, 261–270. https://doi.org/10.1207/s15327752jpa8403_05.

Mehl, M. R. (2006). The lay assessment of subclinical depression in daily life. *Psychological Assessment*, *18*, 340–345. https://doi.org/10.1037/1040-3590.18.3.340.

Mehl, M. R., Gosling, S. D., & Pennebaker, J. W. (2006). Personality in its natural habitat: Manifestations and implicit folk theories of personality in daily life. *Journal of Personality and Social Psychology*, *90*, 862–877. https://doi.org/10.1037/0022-3514.90.5.862.

Meinshausen, N., & Bühlmann, P. (2010). Stability selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, *72*, 417–473. https://doi.org/10.1111/j.1467-9868.2010.00740.x.

Miller, G. (2012). The smartphone psychology manifesto. *Perspectives on Psychological Science*, *7*, 221–237. https://doi.org/10.1177/1745691612441215.

Montag, C., Blaszkiewicz, K., Lachmann, B., Andone, I, Sariyska, R., Trendafilov, B., Reuter, M., et al. (2014). Correlating personality and actual phone usage: Evidence from psychoinformatics. *Journal of Individual Differences*, *35*, 158–165. https://doi.org/10.1027/1614-0001/a000139.

C. Stachl et al.

Montag, C., Blaszkiewicz, K., Sariyska, R., Lachmann, B., Andone, I., Trendafilov, B., Eibes, M., et al. (2015). Smartphone usage in the 21st century: Who is active on WhatsApp? *BMC Research Notes*, *8*, 331. https://doi.org/10.1186/s13104-015-1280-z.

Morris, M. G., Venkatesh, V., & Ackerman, P. L. (2005). Gender and age differences in employee decisions about new technology: An extension to the theory of planned behavior. *IEEE Transactions on Engineering Management*, *52*, 69–84.

Naimi, B. (2015). *Usdm: Uncertainty analysis for species distribution models*. R package version 1.1–15. Retrieved from https://CRAN.R-project.org/package=usdm

O'Hara, R. B., & Kotze, D. J. (2010). Do not log-transform count data. *Methods in Ecology and Evolution*, *1*, 118–122. https://doi.org/10.1111/j.2041-210X.2010.00021.x.

Ormel, J., Jeronimus, B. F., Kotov, R., Riese, H., Bos, E. H., Hankin, B., Rosmalen, J. G. M., et al. (2013). Neuroticism and common mental disorders: Meaning and utility of a complex relationship. *Clinical Psychology Review*, *33*, 686–697. https://doi.org/10.1016/j.cpr.2013.04.003.

Paunonen, S. V., & Ashton, M. C. (2001). Big five factors and facets and the prediction of behavior. *Journal of Personality and Social Psychology*, *81*, 524–539. https://doi.org/10.1037/0022-3514.81.3.524.

Phillips, J. G., Butt, S., & Blaszczynski, A. (2006). Personality and self-reported use of mobile phones for games. *Cyberpsychology & behavior : the impact of the Internet, multimedia and virtual reality on behavior and society*, *9*, 753–758. https://doi.org/10.1089/cpb.2006.9.753.

Rauthmann, J. F. (2012). The dark triad and interpersonal perception: Similarities and differences in the social consequences of narcissism, machiavellianism, and psychopathy. *Social Psychological and Personality Science*, *3*, 487–496. https://doi.org/10.1177/1948550611427608.

Reed, D. a., Enders, F., Lindor, R., McClees, M., & Lindor, K. D. (2011). Gender differences in academic productivity and leadership appointments of physicians throughout academic careers. *Academic Medicine*, *86*, 43–47. https://doi.org/10.1097/ACM.0b013e3181ff9ff2.

Rencher, A., & Pun, F. (1980). Inflation of R 2 in best subset regression. *Technometrics*, *22*, 49–53. https://doi.org/10.2307/1268382.

Revelle, W. (2016). *psych: Procedures for psychological, psychometric, and personality research*. R package version 1.6.4. Evanston, Illinois: Northwestern University Retrieved from http://CRAN.R-project.org/package=psych.

Saeb, S., Zhang, M., Karr, C. J., Schueller, S. M., Corden, M. E., Kording, K. P., & Mohr, D. C. (2015). Mobile phone sensor correlates of depressive symptom severity in daily-life behavior: An exploratory study. *Journal of Medical Internet Research*, *17*. https://doi.org/10.2196/jmir.4273.

Salgado, J. F., Moscoso, S., & Berges, A. (2013). Conscientiousness, its facets, and the prediction of job performance ratings: Evidence against the narrow measures. *International Journal of Selection and Assessment*, *21*, 74–84. https://doi.org/10.1111/ijsa.12018.

Schmid Mast, M., Gatica-Perez, D., Frauendorfer, D., Nguyen, L., & Choudhury, T. (2015). Social sensing for psychology. *Current Directions in Psychological Science*, *24*, 154–160. https://doi.org/10.1177/0963721414560811.

Sherman, R. C., End, C., Kraan, E., Cole, A., Campbell, J., Birchmeier, Z., & Klausner, J. (2000). The internet gender gap among college students: Forgotten but not gone? *Cyberpsychology & Behavior*, *3*, 885–894. https://doi.org/10.1089/10949310050191854.

Statista (2016). *Number of apps available in leading app stores as of June 2016*. Retrieved July 1, 2016, from http://www.statista.com/statistics/276623/numberof-apps-available-in-leading-app-stores/

Van der Linden, D., te Nijenhuis, J., & Bakker, A. B. (2010). The general factor of personality: A meta-analysis of Big Five inter-correlations and a criterion-related validity study. *Journal of Research in Personality*, *44*, 315–327. https://doi.org/10.1016/j.jrp.2010.03.003.

Vazire, S. (2010). Who knows what about a person? The self-other knowledge asymmetry (SOKA) model. *Journal of Personality and Social Psychology*, *98*, 281–300. https://doi.org/10.1037/a0017908.

Venkatesh, V., Thong, J. Y., & Xu, X. (2012). Consumer acceptance and use of information technology: Extending the unified theory of acceptance and use of technology. *MIS Quarterly*, *36*, 157–178 Retrieved from: https://papers.ssrn.com/sol3/papers.cfm?abstract%7B%5C_%7Did=2002388.

Wickham, H. and E. Miller (2016). *Haven: Import and export 'SPSS', 'Stata' and 'SAS' files*. R package version 1.0.0. Retrieved from: https://CRAN.R-project.org/package=haven

Wittmann, W. W. (2012). *Principles of symmetry in evaluation research with implications for offender treatment*. In T. Bliesener, A. Beelmann, & M. Stemmler (Eds.), *Antisocial behavior and crime. Contributions of developmental and evaluation research to prevention and intervention* 2011 (pp. 357–368). Cambridge: Hogrefe.

Wrzus, C., & Mehl, M. R. (2015). Lab and/or field? Measuring personality processes and their social consequences. *European Journal of Personality*, *29*, 250–271. https://doi.org/10.1002/per.1986.

Xu, R., Frey, R. M., Fleisch, E., & Ilic, A. (2016). Understanding the impact of personality traits on mobile app adoption—Insights from a large-scale field study. *Computers in Human Behavior*, *62*, 244–256. https://doi.org/10.1016/j.chb.2016.04.011.

Yarkoni, T. (2010). Personality in 100,000 words: A large-scale analysis of personality and word use among bloggers. *Journal of Research in Personality*, *44*, 363–373. https://doi.org/10.1016/j.jrp.2010.04.001.

Yarkoni, T. (2012). Psychoinformatics: New horizons at the interface of the psychological and computing sciences. *Current Directions in Psychological Science*, *21*, 391–397. https://doi.org/10.1177/0963721412457362.

Zander, A., & Van Egmond, E. (1958). Relationship of intelligence and social power to the interpersonal behavior of children. *Journal of Educational Psychology*, *49*, 257. https://doi.org/10.1037/h0049364.

Zaval, L., Y. Li, E. J. Johnson, and E. U. Weber (2015). Complementary contributions of fluid and crystallized intelligence to decision making across the life span. In T. M. Hess, J. Strough, and L. Corinna (Eds.), *Aging and decision making: empirical and applied perspectives*, Chapter 8, pp. 149–168. Elsevier. doi:https://doi.org/10.1016/B978-0-12-417148-0.00008-X

Ziegler, M., Bensch, D., Maaß, U., Schult, V., Vogel, M., & Bühner, M. (2014). Big Five facets as predictor of job training performance: The role of specific job demands. *Learning and Individual Differences*, *29*, 1–7. https://doi.org/10.1016/j.lindif.2013.10.008.

**APPENDIX**

Table A1. Descriptive statistics—predictors

| Predictor | $M$ | $SD$ | MIN | MAX | 95% CI$_{alpha}$ |
|---|---|---|---|---|---|
| Age | 23.58 | 4.71 | 18.00 | 50.00 | |
| Education | 4.26 | 0.57 | 2.00 | 5.00 | |
| FluidIQ | 0.63 | 0.64 | −0.97 | 3.09 | (see note) |
| Emotional stability (ES) | −0.04 | 0.70 | −2.00 | 2.52 | [0.91, 0.95] |
| Extraversion (E) | 0.03 | 0.74 | −1.98 | 1.88 | [0.94, 0.97] |
| Openness (O) | 0.01 | 0.72 | −1.84 | 2.12 | [0.91, 0.95] |
| Conscientiousness (C) | 0.08 | 0.77 | −1.63 | 1.81 | [0.95, 0.97] |
| Agreeableness (A) | −0.16 | 0.75 | −2.11 | 1.80 | [0.92, 0.95] |
| Carefreeness (ES1) | 0.03 | 1.18 | −2.58 | 3.24 | [0.73, 0.84] |
| Equanimity (ES2) | 0.48 | 1.03 | −2.30 | 3.27 | [0.74, 0.85] |
| Positive mood (ES3) | 0.92 | 1.44 | −4.55 | 5.59 | [0.85, 0.91] |
| Self-consciousness (ES4) | 0.72 | 1.11 | −2.42 | 3.90 | [0.77, 0.87] |
| Self-control (ES5) | 0.70 | 1.01 | −2.10 | 3.36 | [0.65, 0.79] |
| Emotional robustness (ES6) | 0.68 | 1.27 | −1.75 | 5.53 | [0.72, 0.84] |
| Friendliness (E1) | 1.42 | 1.33 | −1.70 | 5.41 | [0.72, 0.83] |
| Sociableness (E2) | 1.35 | 1.72 | −3.41 | 5.64 | [0.89, 0.93] |
| Assertiveness (E3) | 0.79 | 1.42 | −2.30 | 5.61 | [0.79, 0.88] |
| Dynamism (E4) | 1.37 | 1.52 | −2.02 | 5.94 | [0.81, 0.89] |
| Adventurousness (E5) | 0.44 | 1.56 | −3.25 | 5.27 | [0.86, 0.92] |
| Cheerfulness (E6) | 1.82 | 1.66 | −3.23 | 6.09 | [0.83, 0.90] |
| Openness to imagination (O1) | 1.30 | 1.45 | −2.04 | 5.33 | [0.80, 0.88] |
| Openness to aesthetics (O2) | 0.34 | 1.21 | −2.38 | 4.61 | [0.76, 0.86] |
| Openness to feelings (O3) | 2.10 | 2.23 | −5.65 | 6.04 | [0.88, 0.93] |
| Openness to actions (O4) | 1.51 | 1.41 | −2.75 | 5.42 | [0.76, 0.86] |
| Openness to ideas (O5) | 1.88 | 1.44 | −0.85 | 5.51 | [0.78, 0.87] |
| Openness to the value and norm system (O6) | 0.93 | 1.04 | −1.61 | 4.86 | [0.61, 0.77] |
| Competence (C1) | 1.05 | 1.30 | −1.87 | 4.43 | [0.75, 0.85] |
| Love of order (C2) | 1.21 | 1.63 | −4.34 | 5.67 | [0.88, 0.93] |
| Sense of duty (C3) | 2.20 | 1.46 | −1.59 | 5.50 | [0.75, 0.85] |
| Ambition (C4) | 2.20 | 1.62 | −1.40 | 5.86 | [0.81, 0.89] |
| Discipline (C5) | 1.77 | 1.53 | −1.13 | 5.75 | [0.79, 0.87] |
| Caution (C6) | 1.78 | 1.42 | −1.33 | 5.75 | [0.82, 0.89] |
| Willingness to trust (A1) | 0.23 | 1.32 | −3.09 | 4.21 | [0.80, 0.88] |
| Genuineness (A2) | 1.00 | 0.91 | −1.20 | 4.25 | [0.54, 0.73] |
| Helpfulness (A3) | 1.60 | 1.46 | −2.47 | 6.04 | [0.75, 0.85] |
| Obligingness (A4) | 0.89 | 1.15 | −1.86 | 3.70 | [0.72, 0.84] |
| Modesty (A5) | 0.58 | 1.18 | −2.68 | 3.91 | [0.75, 0.85] |
| Good-naturedness (A6) | 1.91 | 1.73 | −2.99 | 6.40 | [0.77, 0.87] |

*Note.* Descriptive statistics for all predictor variables. Internal consistencies are provided as 95% confidence intervals. The minimum and maximum values represent the lowest and highest person parameters, estimated by the fitted partial credit model, respectively. Cronbach alpha scores of the three subscales of fluid intelligence are provided here as mean scores across all participants: Num.: $M = 0.71/SD = 0.06$, Verb.: $M = 0.74/SD = 0.06$, Fig.: $M = 0.71/SD = 0.06$; $M$, mean; $SD$, standard deviation; MIN, minimum value; MAX, maximum value.

Table A2. Pairwise Spearman correlations between psychometrics and demographics

| | Predictor | 1 | 2 | 3 | 4 | ES | E | O | C |
|---|---|---|---|---|---|---|---|---|---|
| 1 | Gender | 1 | | | | | | | |
| 2 | Age | −0.10 [−0.26, 0.07] | 1 | | | | | | |
| 3 | Education | −0.04 [−0.21, 0.13] | **0.42 [0.27, 0.55]** | 1 | | | | | |
| 4 | FluidIQ | −0.14 [−0.30, 0.03] | −0.12 [−0.28, 0.05] | 0.02 [−0.15, 0.19] | 1 | | | | |
| ES | Emotional stability | −0.06 [−0.23, 0.11] | 0.02 [−0.15, 0.19] | 0.06 [−0.11, 0.23] | 0.08 [−0.09, 0.24] | 1 | | | |
| E | Extraversion | **0.20 [0.03, 0.36]** | 0.00 [−0.17, 0.17] | 0.01 [−0.16, 0.18] | −0.13 [−0.29, 0.04] | **0.42 [0.27, 0.55]** | 1 | | |
| O | Openness | 0.14 [−0.03, 0.30] | 0.04 [−0.13, 0.21] | 0.04 [−0.13, 0.21] | −0.05 [−0.22, 0.12] | **0.31 [0.15, 0.45]** | **0.58 [0.46, 0.68]** | 1 | |
| C | Conscientiousness | 0.14 [−0.03, 0.30] | −0.06 [−0.23, 0.11] | 0.03 [−0.14, 0.20] | −0.03 [−0.20, 0.14] | **0.29 [0.13, 0.44]** | **0.19 [0.02, 0.35]** | **0.29 [0.13, 0.44]** | 1 |
| A | Agreeableness | **0.26 [0.10, 0.41]** | 0.07 [−0.10, 0.23] | 0.07 [−0.10, 0.23] | −0.08 [−0.24, 0.09] | **0.23 [0.06, 0.38]** | **0.34 [0.18, 0.48]** | **0.42 [0.27, 0.55]** | **0.17 [0.00, 0.33]** |

*Note.* Pairwise Spearman correlations between Big Five measures, fluid intelligence measures, demographic variables and app usage measures. Square brackets contain 95% confidence intervals. Confidence intervals not including zero are in bold.

# To Challenge the Morning Lark and the Night Owl: Using Smartphone Sensing Data to Investigate Day–Night Behaviour Patterns

This article demonstrates that the growing digitalization of lifestyles allows for investigating day–night patterns and related traits using behavioral data collected through smartphone sensors.

***Contributing article***

Schoedel, R., Pargent, F., Au, Q., Völkel, S. T., Schuwerk, T., Bühner, M., and Stachl, C. (2020). To challenge the morning lark and the night owl: Using smartphone sensing data to investigate day–night behaviour patterns. *European Journal of Personality*, 34(5):733–752

This publication was part of the PhoneStudy project (see Chapter 7).

***Copyright information***

***Declaration of contributions***

Preliminary work of the doctoral candidate on other projects within the PhoneStudy project was crucial for the data preprocessing and extraction of relevant variables in this manuscript. Moreover, the PhD candidate contributed to the manuscript's review and editing process.

*Contribution of the coauthors*

Ramona Schoedel was the main author of this manuscript and was responsible for conducting the study, including data collection. Florian Pargent made significant contributions by providing insightful guidance and interpretations on the results of the study. The remaining coauthors played crucial roles in the PhoneStudy project and helped revise the final version of this publication.

# To Challenge the Morning Lark and the Night Owl: Using Smartphone Sensing Data to Investigate Day–Night Behaviour Patterns

RAMONA SCHOEDEL[1]*, FLORIAN PARGENT[1], QUAY AU[2], SARAH THERES VÖLKEL[3], TOBIAS SCHUWERK[4], MARKUS BÜHNER[1] and CLEMENS STACHL[5]

[1]*Psychological Methods and Assessment, Department of Psychology, LMU Munich, Munich, Germany*
[2]*Computational Statistics, Department of Statistics, LMU Munich, Munich, Germany*
[3]*Media Informatics Group, Institute of Informatics, LMU Munich, Munich, Germany*
[4]*Developmental Psychology, Department of Psychology, LMU Munich, Munich, Germany*
[5]*Department of Communication, Stanford University, USA*

*Abstract: For decades, day–night patterns in behaviour have been investigated by asking people about their sleep–wake timing, their diurnal activity patterns, and their sleep duration. We demonstrate that the increasing digitalization of lifestyle offers new possibilities for research to investigate day–night patterns and related traits with the help of behavioural data. Using smartphone sensing, we collected in vivo data from 597 participants across several weeks and extracted behavioural day–night pattern indicators. Using this data, we explored three popular research topics. First, we focused on individual differences in day–night patterns by investigating whether 'morning larks' and 'night owls' manifest in smartphone-sensed behavioural indicators. Second, we examined whether personality traits are related to day–night patterns. Finally, exploring social jetlag, we investigated whether traits and work weekly day–night behaviours influence day–night patterns on weekends. Our findings highlight that behavioural data play an essential role in understanding daily routines and their relations to personality traits. We discuss how psychological research can integrate new behavioural approaches to study personality. © 2020 The Authors. European Journal of Personality published by John Wiley & Sons Ltd on behalf of European Association of Personality Psychology*

Key words: chronotype; day–night behaviour patterns; diurnal activity; personality; smartphone sensing data

## INTRODUCTION

Are there times of day when you do not use your smartphone at all? Most likely at night. As our everyday companions, smartphones can provide much information about people's day–night patterns (Harari et al., 2016). So far, behavioural manifestations of the underlying circadian system like sleep–wake timing, diurnal activity, or sleep duration have mainly been assessed via self-reports (Adan et al., 2012). However, self-reports about behaviour are known to differ from actual records of behaviour (Baumeister et al., 2007; Gosling et al., 1998). Emphasizing this dilemma, Lauderdale et al. (2008) correlated behaviourally assessed sleep duration with self-reports and concluded that people systematically misjudge it. An alternative approach is to collect actigraphy-based data to study sleep behaviour: movements and environmental factors like ambient brightness are recorded with wristbands and are jointly converted to indicators for sleep–wake timing by special algorithms (e.g.

Križan & Hisler, 2019; Tonetti et al., 2016; Vitale et al., 2015). Regarding the trade-off between measurement accuracy and ecological validity, another interesting complement for studying sleep behaviour could be the use of smartphone sensing data. These data cannot provide a direct measurement of sleep–wake phases, but only periods of nightly inactivity of smartphone use in which physiological sleep occurs. In contrast to actigraphy, these measurements do not take body signals such as movements or pulse into account. However, first studies have indicated that smartphone data provide useful information about sleep–wake timing as smartphones are meanwhile considered to be part of new sleeping habits (Chen et al., 2013; Min et al., 2014; Lin et al., 2019; Borger et al., 2019). Borger et al. (2019) have shown that indicators for sleep onset and offset derived via actigraphy and smartphone touch interactions are highly correlated. In addition, independence from sensors worn on the body also offers advantages in terms of ecological validity. With the help of commercially available smartphones, behavioural indicators for sleep–wake timing can be collected efficiently and unobtrusively in everyday life over a more extended period, even for large samples. To illustrate this, we use smartphone-sensed indicators for sleep–wake timing to investigate traits related to day–night patterns. For this purpose, we chose to study three frequently researched questions, which we will introduce in the following sections.

*Correspondence Ramona Schoedel, Leopoldstraße 13, 80802 Munich, Germany.
Email: ramona.schoedel@psy.lmu.de

### Individual differences in behavioural day–night patterns

The human circadian system has been studied for decades by interdisciplinary research teams. The most prominent finding across all research disciplines is that individuals show stable differences in day–night patterns, a stable trait that is often referred to as the *chronotype* (e.g. Adan et al., 2012; Cavallera & Giudici, 2008; Roenneberg et al., 2003). Literature frequently describes two extremes: the *morning type* ('morning lark') wakes up and goes to bed early, feels fit after getting up, and performs best early in the day. The *evening type* ('night owl') wakes up and goes to bed later, feels tired after waking up, and performs best towards the end of the day (for extensive reviews, see Adan et al., 2012; Cavallera & Giudici, 2008; Takano et al., 2014). The chronotype has been argued to be a genetically predisposed trait with various biological manifestations like body temperature or hormone levels (Bailey & Heitkemper, 2001; Horne & Östberg, 1976; Roenneberg et al., 2003; Katzenberg et al., 1998). In addition, chronotype should be distinguished from sleep duration, which has been argued to be an independent trait (Roenneberg et al., 2007).

Based on the distinction between variable-centred and person-centred personality assessment (Asendorpf, 2003), one might assume that chrono 'types' refer to distinct groups of individuals with similar manifestations in chronotype-related behaviours. However, Putilov (2017) points out in his review that researchers have not yet reached an agreement on the number and content of underlying dimensions, the resultant number of types, and whether the conceptualization as types makes sense at all (Roenneberg et al., 2003). Two different operationalizations of chronotype are most prominent in the literature (see Table 1).

Dating back to Horne and Östberg (1976), chronotype is described as *circadian* or *morningness–eveningness preferences*. The term 'circadian typology' is often used synonymously and shows the emphasis on the categorization of chronotypes in this research tradition (e.g. Adan et al., 2012; Lipnevich et al., 2017). In comparison, Roenneberg et al. (2015) accentuate the chronotype as a continuous variable and describe it as a trait reflecting the *phase of entrainment*, which represents individual differences in the synchronization of the internal circadian rhythm to environmental factors (e.g. light/dark cycle, diurnal temperature curve, social interaction). Despite their different understanding of the underlying construct of chronotype (Roenneberg, 2015), both operationalizations have been found to be strongly correlated (Zavada et al., 2005). In the present study, we take the structural ambiguity of chronotype as our starting point to investigate how smartphone sensing data reflecting day–night activity patterns could help to inform chronotype research, as operationalized both in the Horne–Östberg and in the Roenneberg tradition.

In the Horne–Östberg tradition, the Morningness–Eveningness Questionnaire (MEQ; Horne & Östberg, 1976) still represents the gold standard for chronotype assessment (Putilov, 2017). The MEQ asks for circadian preferences and categorizes people according to *ad hoc* specified cut-off values (Horne & Östberg, 1976). In the development of the MEQ, neither the grouping nor the factorial structure was investigated. Cut-off values were determined using a small but not representative sample (Caci et al., 2009). Meanwhile, various derivates and short scales of the MEQ have been published (Adan et al., 2012; Putilov, 2017). Assumptions on the underlying structure of circadian preferences range from a continuum with two extremes (Natale & Cicogna, 2002; Tonetti et al., 2016), over two dimensions (morningness and eveningness as separate dimensions Lipnevich et al., 2017) to a multidimensional construct with up to four factors (Adan et al., 2012; Randler et al., 2016; Caci et al., 2009). Recently, Preckel et al. (2019) have published pioneering work on a typology of circadian preferences providing empirical evidence on the possible number of types. In an adolescent sample, they found evidence for four types resulting from the combination of the two independent dimensions of morningness and eveningness preference. Joining this search for structure, we translate the questionnaire items typically used to determine the Horne–Östberg chronotype into behavioural smartphone sensing equivalents. Smartphone usage variables can

Table 1. Description of the two most popular approaches to chronotype

| Feature | Horne–Östberg chronotype | Roenneberg chronotype |
|---|---|---|
| Assessment | Morningness–Eveningness Questionnaire (MEQ) by Horne & Östberg (1976) | Munich Chronotype Questionnaire (MCTQ) by Roenneberg et al. (2003) |
| Chronotype as | Time of day preferences | Phase of entrainment |
| Items | Ask for imagined free days: preferred sleeping times, preferred times for mental/physical activity, subjective feeling in the morning/evening, self-reported chronotype | Ask for both free and work days habitual sleeping times |
| Determination of chronotype | Cut-off values classify participants according to their 19-item sum score | Midpoint of sleep for free days without alarm clock usage |
| Emphasized structure | Four dimensions (peak time, morning affect, retiring, rising) according to Caci et al. (2009) | Continuous variable |

The structure for the Horne–Östberg chronotype refers to the original chronotype assessment with the MEQ. However, several derivates of the MEQ have been developed and there is no consensus in research about the factorial structure of the chronotype approximated by the assessment of circadian preferences. Solutions range from one to four dimensions.

approximate many of them. Following Putilov's (2017) recommendation to consider behavioural markers for circadian preferences, we investigate whether we can find types of individuals with similar smartphone usage patterns indicating circadian preferences. Finally, we explore the factorial structure of the behavioural indicators.

In the Roenneberg tradition, freely chosen sleep–wake timing is considered the best approximation of the internal circadian rhythm. Therefore, sleep–wake habits for both work and free days are assessed while controlling for alarm clock usage (Roenneberg et al., 2003; Roenneberg et al., 2015). In this taxonomy, the midpoint between sleep onset and offset determines the chronotype. This reference point for sleep has proven to coincide with nocturnal melatonin production, which in turn controls sleep–wake timing (Terman et al., 2001; Roenneberg et al., 2003; Roenneberg et al., 2007; Roenneberg et al., 2015). In this context, the Munich Chronotype Questionnaire (MCTQ), which has been repeatedly validated by behavioural (actigraphy) and biological (melatonin, cortisol) circadian system markers, is primarily used (Roenneberg et al., 2003; Roenneberg et al., 2007). Only recently, Lin et al. (2019) took up the idea to determine the Roenneberg chronotype by using smartphone sensing data and provided first indications that there is a considerable overlap between sleeping times assessed via smartphones and self-reports. However, their algorithm for characterizing a *digital chronotype* does not explicitly correspond to Roenneberg's chronotype criteria, as they did not differentiate between work and free days and were restricted to the use of a very limited range of data (screen and notification events Lin et al., 2019). We propose a more fine-grained algorithm for determining a smartphone sensing-based proxy by using only free days without alarm clock usage. To explore our smartphone chronotype, we look at descriptives and correlational analyses that were presented by Roenneberg's group to describe the MCTQ-based chronotype. For example, Roenneberg et al. (2007) found that sleep duration depends on chronotype if analysed separately for work and free days and that chronotype is related to age and gender.

**Behavioural day–night patterns and personality traits**

Important research questions are associations between day–night patterns, personality, and demographics. Different aspects of day–night behaviour have been addressed in this context. For example, the *morningness preference* has been linked to personality. Higher values in this dimension indicate a preference for getting up and going to bed early, feeling fit in the morning, and achieving peak performance earlier in the day (Lipnevich et al., 2017). The most established findings in meta-analyses are that conscientiousness and agreeableness are positively related to morningness (Tsaousis, 2010; Lipnevich et al., 2017). No or only small relationships in a specified direction can be found for neuroticism and openness (Adan et al., 2012; Tsaousis, 2010; Lipnevich et al., 2017). Negative relationships between morningness and extraversion were found, but only if the trait extraversion was described with Eysenck's three-factor model (Adan et al., 2012; Tsaousis, 2010). Using the five-factor model, this

association is almost zero (Tsaousis, 2010). For the sake of completeness, please note that morningness has also been found to be related to personality styles or, more precisely, with thinking and behaving styles (Díaz-Morales, 2007). Furthermore, age has been robustly related to morningness. Shifts towards eveningness in adolescence and towards morningness with increasing age (at around 50) have been reported (e.g. Adan et al., 2012; Cavallera & Giudici, 2008). Regarding gender, a meta-analysis has found that the preference for morningness is slightly higher for women compared with men (Randler, 2007). However, complex interactions between age and gender have been reported in previous literature. For example, girls at the age of 13 and 14 have a lower tendency towards morningness than their male counterparts (Mateo et al., 2012), and their peak towards eveningness is earlier (e.g. Adan et al., 2012). In addition, Randler and Engelke (2019) have shown a complex interaction between age and gender with regard to morningness preferences: young women were more and older women less morning oriented than young or older men.

In addition, associations between *sleep duration* and personality traits have been investigated, but findings have been ambiguous so far. For example, there is some evidence that individuals with higher values in neuroticism report to sleep longer (Duggan et al., 2014). According to Križan and Hisler (2019), neuroticism is not related to the mean sleep duration but positively related to the intraindividual variation in sleep duration. Some studies reported correlations between sleep duration and conscientiousness, agreeableness, or openness but not extraversion (Randler, 2008; Križan & Hisler, 2019). In contrast, other researchers did not find any evidence that sleep duration and big five personality traits are associated (Gray & Watson, 2002; Randler et al., 2017; Sutin et al., 2019). Sleep duration decreases with age (Randler, 2008) but was not found to be related to gender (Randler et al., 2017).

In summary, past research provides some evidence for associations between personality traits and day–night behaviour, but past findings are inconsistent. One possible reason for this could be that the majority of studies (except Križan & Hisler, 2019; Sutin et al., 2019) asked participants about their habits but did not include any behavioural measures of sleep. Not only might people differ in their ability to estimate their sleep duration, personality traits themselves might play a role in the evaluation of their day–night behaviours. To circumvent this issue here, we use data from smartphone sensing to derive indicators for sleep–wake behaviour and to consequently investigate their relationship with big five personality traits on factor and facet level. Additionally, we explore *sleep continuity*, which has been defined as a measure of how well people fall asleep and sleep through (Ohayon et al., 2017). Recent actigraphy-based research has found, for example, that conscientiousness and extraversion were negatively related to behavioural indicators of sleep continuity, such as wake after sleep onset. In contrast, higher scores in neuroticism were associated with more wakening (Sutin et al., 2019). As a rough smartphone-based approximation measure, we look at two aspects of sleep continuity: how often and for

how long people check their smartphones during the night. Additionally, we analyse smartphone activity logs to explore how *alarm clock usage*—particularly 'snoozing'—is related to personality.

### Intraindividual and interindividual differences in day–night patterns: The social jetlag

Finally, we explore the so-called *social jetlag* hypothesis (e.g. Adan et al., 2012; Wittmann et al., 2006). Roenneberg et al. (2007) surveyed the sleep habits of more than 55 000 people using the MCTQ and found that sleep behaviour differs for work-free days versus workdays. Specifically, their findings suggest that people, on average, go to bed and awake earlier on work than on free days. Furthermore, the proportion of sleep onset and offset is smaller for workdays than for free days. It has been suggested that this effect is induced by social obligations (Wittmann et al., 2006). Thus, the pairing of late bedtimes with consistent wake-up times leads to a sleep deficit for a week. As a consequence, sleep is compensated on weekends (Roenneberg et al., 2015). This misalignment of the internal biological and the external social clock is associated with health risk behaviours (e.g. increased body mass index and smoking Roenneberg et al., 2012; Wittmann et al., 2006). According to Wittmann et al. (2006) and Roepke and Duffy (2010), late chronotypes are particularly affected by the social jetlag as they stay up until late at night but have to get up early to go to work or to pursue other social obligations on the following day. The assessment of individuals' daily routines through the analysis of smartphone activity logs for several weeks allows us to investigate compensatory nightly rest by considering intraindividual and interindividual factors. Using these indicators, we want to explore whether the smartphone-sensed proxies for sleep duration on weekends and respective weeks are related and whether interindividual factors like the Roenneberg chronotype, demographics, and personality traits have an impact.

### Rationale

Our study aims to reinvestigate selected topics regarding day–night pattern-related traits by using smartphone sensing data. Because we use a new type of data in this field of research, this is exploratory work. A handful of studies have started to use smartphone data in this context (e.g. Chen et al., 2013; Min et al., 2014; Lin et al., 2019). However, these studies have mostly been limited in terms of sample size and types of sensing data.

Here, we show how behavioural records from smartphones can be used to investigate individual differences in day–night patterns, how they relate to personality traits, and how they are influenced by intraindividual and interindividual factors. Besides the examination of whether 'morning larks' and 'night owls' manifest in indicators of sleep–wake timing and diurnal activity patterns, we explore the smartphone-based operationalization of the Roenneberg chronotype. We investigate the associations of day–night behaviour patterns and personality traits. Finally, we illustrate

how continuously logged behavioural data can be used to investigate the contribution of both intraindividual and interindividual factors to predict indicators for sleep behaviour on weekends, using the social jetlag hypothesis as an example.

### METHOD

Our analyses are based on data collected within the long-time project *PhoneStudy* (Stachl et al., 2018). This ongoing interdisciplinary research project at LMU Munich uses the continuously developed smartphone sensing application *PhoneStudy* for Android smartphones for collecting natural smartphone usage behaviours in the field. Data about app usage, calling activity, general phone usage (e.g. calendar, music, power supply), and connectivity (e.g. Bluetooth, WiFi) are logged whenever the respective events occur. GPS data are usually recorded once every 15 minutes. Data are synchronized hourly to the back end server via Secure Sockets Layer (SSL) encryption, whenever a WiFi connection is available. The responsible institutional review board and data protection office approved the project and all associated studies. All materials and aggregated data can be found in our open science framework project (OSF; Schoedel et al., 2020).[1] To protect the data privacy rights of our participants, the raw sensing data cannot be made available due to their granularity.

### Description of data set

We combined data resulting from three studies conducted between 2014 and 2018. In Table 2, we show some basic information about the included studies. Despite some marginal differences, data collection procedures of all studies followed the same principle: after giving informed consent, participants were asked to install the *PhoneStudy* app for at least 30 days on their private smartphones and to complete several questionnaires before, during, or after the smartphone logging period. Participants were mostly recruited in the university context via flyers, mailings lists, social media, and personal contact in Munich, Germany. For more detailed information about study procedures, see also Stachl et al. (2017); Harari et al. (2019); Schuwerk et al. (2019); Schoedel et al. (2018); and Stachl et al. (2019).

We applied several exclusion criteria to our initial data set of 743 participants. We excluded participants with fewer than 21 days of sensing data, more than 50% missing values across all variables, and if questionnaire data were not available. We included data from a maximum of 32 days of continuous logging. This resulted in a final sample size of 597 (61% women). As recruitment took place in the university context, participants were, on average well educated (71% with a high school and 20% with a university degree). With a mean age of 23.56 years ($SD$ = 6.55; $Min$ = 18, $Max$ = 72), the sample was skewed towards younger participants (18–21: 39%; 22–25: 34%; 26–30: 12%; 31–40: 5%; 41 and older: 3%). For a more detailed description of the sample, according to studies, see Table 3.

[1]https://osf.io/a4h3b/

Table 2. Description of data sets used in the study

| Data set | References | N | Study period | Compensation |
|---|---|---|---|---|
| 1 | Stachl et al. (2017), Harari et al. (2019) | 132 (137) | 09/2014–08/2015 | Individualized personality profile and 30 € or course credits |
| 2 | Schuwerk et al. (2019) | 240 (245) | 08/2016–08/2017 | Up to 35 € and lottery (smartphone or tablet worth 400 €) |
| 3 | Schoedel et al. (2018) | 225 (361) | 10/2017–01/2018 | Individualized personality profile and user activity feedback, course credits, and lottery (10 × 50 €) |

N indicates the size of the sample of the respective study after application of our inclusion criteria. The total number of subjects per study is given in parentheses.

Table 3. Description of the sample according to studies

| Data set | N | Age | Education | Students | Employment status |
|---|---|---|---|---|---|
| 1 | 132 | 23.61 (4.73) | No qualification: 0.00%<br>Secondary school: 3.79%<br>High school: 65.15%<br>University: 31.06% | No data available | No data available |
| 2 | 240 | 22.94 (4.57) | No qualification: 0.00%<br>Secondary school: 9.58%<br>High school: 72.50%<br>University: 17.92% | 73.50% | No data available |
| 3 | 225 | 24.20 (8.86) | No qualification: 0.44%<br>Secondary school: 8.88%<br>High school: 72.44%<br>University: 16.44% | 77.33% | Unemployed: 4.89%<br>In training: 24.89%<br>Minor employm.: 41.33%<br>Part-time: 10.67%<br>Full-time: 15.56%<br>Other: 0.88% |

N indicates the size of the samples according to studies. The column Age presents the mean value, and standard deviations are given in parentheses. As procedures slightly varied across studies, not all demographic variables are available for all data sets. The category *other* in the column Employment Status comprises retraining and pension.

## Measures

### Self-report measures

We administered various self-report questionnaires. However, we limit our report to the ones used in our statistical analyses. Besides demographics, personality traits were assessed with the Big Five Structure Inventory (BFSI Arendasy, 2009). Each of the big five factors—openness, conscientiousness, extraversion, agreeableness, and emotional stability—was measured on respectively six subscales (Table 8). Participants were asked to rate 300 personality describing adjectives and short phrases on a 4-point Likert scale with the labels *untypical for me*, *rather untypical for me*, *rather typical for me*, and *typical for me*. Compared with the widely used structure inventory NEO-PI-R (Costa & McCrae, 2008), the BFSI is supposed to have better psychometric properties: Cronbach $\alpha$ values (ranging between 0.72 and 0.92) are partly higher, and subscales are unidimensional in the original paper (Arendasy, 2009). In addition, the BFSI should be less dependent on the participant's reading comprehension ability as it uses short and simple items (Arendasy, 2009). The construction of the BFSI does not follow the classical test theory, but the item response theory framework. Accordingly, the BFSI has been developed in conformity with the partial credit model (Masters, 1982),

which is a probabilistic model describing an individual's observable score on a single item as the result of the functional relationship between the individual's latent trait value (person parameter) and latent item thresholds, which indirectly determine item difficulty (item parameter Arendasy, 2009). Correspondingly, we used the person parameter estimates as personality scores in all our analyses.

### Day–night behavioural measures

Raw smartphone sensing data are sequences of timestamped event data. Whenever a usage event happens, a data entry specified by several event characteristics (e.g. date, study day, details about the event like app package name or type of call) is created. To get an idea of the raw data structure, see also the supplemental codebook (Schoedel et al., 2020). To investigate the research questions specified above, we created variables by reviewing the literature and translating behavioural sleep indicators into smartphone sensing behaviours. Based on our smartphone sensing data, we computed proxy variables to estimate sleep-related behaviours. Please note that our variables are likely to overestimate actual sleep as the last smartphone usage event in the evening has to be before the physiological onset of sleep, and the first smartphone usage event in the morning occurs with delay after waking up. As smartphone sensing data are prone to

logging errors, we extracted robust behavioural estimators when appropriate for the respective variable (Kafadar, 2003; Rousseeuw & Croux, 1993). To stay within the scope of this article, we only summarize our procedure and the engineered variables in the following sections. However, note that variable extraction is usually the most complex and time-consuming task in analyses of smartphone sensing data, and the process includes many researchers' degrees of freedom. For transparency, we provide all code in our OSF project, and the variable extraction procedure is described in detail in the supplemental codebook (Schoedel et al., 2020).

*General indicators for sleep-related behaviours.* We computed the following variables daily while distinguishing between days during the week versus the weekend (Roenneberg et al., 2007). Based on the algorithm specified in Table 7, we determined the *first and last events* according to individual study days and calculated mean and intraindividual variation variables. We defined the smartphone proxy for sleep duration, *nightly inactivity*, as the period between the last event of the day and the first event of the following day. To explore social jetlag, we calculated the average daily inactivity during the night for weekdays and weekends for all study weeks individually.

In addition, we translated two aspects of sleep continuity, sleep fragmentation and waking up after bed, into smartphone usage behaviour by calculating the average number and duration of *checking events* at night. At this point, we would like to point out that our measures do not fully meet the definition of sleep fragmentation and wake after sleep onset by Ohayon et al. (2017). Hence, our measurements only give a rough estimate, taking into account the occurrence of very short smartphone checking events during the nightly inactivity period of smartphone use, which was not part of a more extended usage period in the evening and the next morning. Accordingly, we defined nightly checking events as short periods of less than 2 minutes of smartphone usage during otherwise nightly inactivity. Due to the lack of empirical data in the literature, we have set this threshold value considering that smartphone usage of fewer than 2 minutes might be caused by less significant actions such as checking the clock during the night.

Finally, we calculated some variables related to using the smartphone as an alarm clock: the *mean point of time of alarm app ringing*, the mean daily number, and duration of *snoozing events* (snoozing was defined as the repetition of alarm app events in the morning).

*Horne–Östberg chronotype variables.* To operationalize circadian preferences in terms of smartphone usage behaviour, we computed the following items of the MEQ (Table 1). We translated preferred sleeping times as *mean points of time of the first and the last smartphone usage event on weekends*, as weekends are likely to be organized freely. Following this assumption, we also specified preferred times for activity as diurnal smartphone activity patterns. In this context, we distinguished between different behavioural categories: social communication (social media/communication app usage, calls, and texting), entertainment (browser, gaming, music/video, and news app usage), and general smartphone usage (all active

smartphone usage events). To take into account the distribution of usage events throughout the day, we computed the first quartile, the median, and the third quartile of usage events according to the behavioural categories for each day. In other words, we extracted timestamps that indicate when *25%, 50%, and 75% of the daily events of the respective usage category* took place. Then we computed the mean across all study days for each of the three quantiles. Finally, to depict the subjective feeling of sleepiness in the morning, we considered the *mean number and duration of snoozing events during the week* to indicate how readily people get up in the morning.

*Roenneberg chronotype variables.* Similar to the assessment of the chronotype using the MCTQ, we calculated the *midpoint of sleep* (*MSF*), which is the mean halfway point in time between the last event of a day and the first event of the next day for free (weekend) days without alarm app usage. In addition, we determined the *corrected midpoint of sleep* ($MSF_{corr}$), which has been proposed by Roenneberg et al. (2007) to correct for the sleep debt collected during the week. According to them the $MSF_{corr}$ is better suited for estimating the true underlying chronotype.

### Data analysis

#### Clustering

In the following, we give a short overview of the applied methods. More detailed information can be found in Appendix A. To investigate whether participants can be assigned to groups of similar smartphone usage behaviours indicating circadian preferences, we used clustering as an unsupervised machine learning method. We applied the commonly used *k-means* clustering algorithm with the Euclidean distance as proximity measure. Clustering aims to reduce complexity by finding meaningful structures within the data. According to their similarity in a predefined set of variables, participants are clustered in within-homogeneous groups that are well separated from participants of other clusters (Tan et al., 2006). However, one disadvantage of clustering algorithms is that they sometimes identify random and, therefore, nonreplicable structures (Tan et al., 2006). In line with the literature, we address this problem by using a data-driven approach to determine the number of clusters (Tibshirani & Walther, 2005) and by evaluating the stability and validity of the identified clusters based on bootstrapped metrics (Hennig, 2007; 2008; Tan et al., 2006). We followed the recommendations of Hennig (2018) and used 100 bootstrap iterations. For evaluating cluster stability, we considered the Jaccard coefficient (*JC*, indicates stability if values exceed 0.85) and the criteria of *recovery* and *dissolution*, which count how often each cluster has been successfully recovered and dissolved across all bootstrap iterations (Hennig, 2007; 2008). For evaluating the internal validity of clusters we looked at metrics indicating how similar participants within each cluster are (within-compact) and how different participants from different clusters are (between-separated): the ratio of average within- and between-cluster distances (*wb.ratio* Tan et al., 2006), the *silhouette* coefficient

(Rousseeuw, 1987), and the *Dunn* index (Dunn, 1974; Halkidi et al., 2001). Clusters are within-compact and between-separated if the ratio of distances is small, the silhouette index is close to 1, and the Dunn index is high (Tan et al., 2006; Hennig, 2018). As the *k-means* algorithm cannot handle missing values, we used the multivariate imputation by chained equations technique and specified a random forest imputation model (MICE, van Buuren & Groothuis-Oudshoorn, 2011).

*Exploratory factor analysis*
To explore the factorial structure of our smartphone-based proxy for the Horne–Östberg chronotype, we conducted an exploratory factor analysis based on the averaged correlation matrix of the imputed data sets. We determined the number of factors using the *empirical* Kaiser criterion, which has been shown to perform well for short scales (Braeken & Van Assen, 2017).

*Multilevel modelling*
Measures for nightly inactivity of smartphone usage were repeatedly measured across several study weeks. Considering the intraindividual data dependency, we used multilevel regression modelling with behavioural measures on a weekly basis reflecting level 1 variables that were nested within individuals (level 2). Therefore, we specified a *random-intercept-random-slope model* predicting the mean nightly inactivity duration on weekends based on the mean nightly inactivity duration of the respective preceding workweek (level 1). The averaged nightly inactivity duration, the Roenneberg chronotype, the big five traits, age, and gender, were included as predictors on level 2.

Regarding data preprocessing, we were faced with the challenge of selecting one path from a series of plausible steps. To do justice to these many researcher degrees of freedom and to increase research transparency, we follow the suggestion of Steegen et al. (2016) and present a *multiverse analysis*: for each possible combination of plausible preprocessing steps, a 'new' data set is constructed, and the same multilevel model is estimated for each of those data sets. The multiverse analysis illustrates how much the results depend on the choice of specific preprocessing steps or vice versa, which results are robust across all preprocessing options (Steegen et al., 2016; Simonsohn et al., 2015). Our preprocessing choices include the *coding of the weekend* (Friday to Sunday versus Friday to Monday), the selection of the *number of repeated measurements* (3 versus 4 weeks), the handling of *outliers* (median versus winsorization), and the handling of *missing values* (listwise deletion versus multiple imputation). A detailed description of the alternatives for each decision can be found in supplemental method section in Appendix A. Combining all described decisions resulted in $2 \times 2 \times 2 \times 2 = 16$ choice combinations (see left side in Figure 4).

We used the uncorrected version of the Roenneberg chronotype as a predictor, as we explicitly control for a nightly inactivity deficit in the multilevel model. Gender was dummy coded (0 = male, 1 = female), and all continuous predictor variables were z-standardized based on the grand mean. The level 1 predictor duration of nightly inactivity during the week was centred around the individual mean, which in turn was entered as level 2 predictor (Curran & Bauer, 2011). For a more detailed description of the equation of the multilevel model, we refer the interested reader to the supplemental method section in Appendix A.

*Statistical software*
All data preprocessing and analyses were conducted using R 3.5.0 (R Core Team, 2018). We used *packrat* (Ushey et al., 2018) for package management. For extracting behavioural variables, we mainly used the R packages *dplyr* (Wickham et al., 2019) and *fxtract* (Au, 2019). Multiple imputation was done by using the package *mice* (van Buuren & Groothuis-Oudshoorn, 2011). In addition, we used the following packages to conduct our main analyses: *fpc* for clustering (Hennig, 2018), *psych* for exploratory factor analysis (Revelle, 2018), and *lme4* and *lmerTest* for multilevel modelling (Bates et al., 2015; Kuznetsova et al., 2017). For data visualization, we applied *ggplot2* (Wickham, 2016) and *corrplot* (Wei & Simko, 2017) and created *raincloud plots* (Allen et al., 2019). The complete list of used R packages can be found in our OSF project (Schoedel et al., 2020).

## RESULTS

### Descriptives

We recorded a mean of 22 547 events ($SD = 24 368$) for each participant across the whole study period. Participants had on average smartphone records for 21 ($SD = 1.57$) weekdays and 8 ($SD = 0.92$) weekend days. The mean number of logs per study day was 765 ($SD = 804.70$). As can be seen in Table 4, the average time of first and last smartphone usage was later for weekends than weekdays, and the duration of nightly inactivity was about 20 minutes longer on weekends than on weekdays. However, the mean number and duration of checking events during the night were similar for weekends and weekdays. A total of 91% of our participants used alarm clock apps in the morning, at 7.19 AM on average during the week and about 30 minutes later on weekends. Note that 38% of participants did not use alarm clock apps on any weekend during the entire study period. The number and duration of snoozing events were similar for weekdays and weekends. Descriptive statistics for big five personality traits can be found in Table 8 in the Appendix.

### Individual differences in behavioural day–night patterns

*Person-centred and variable-centred structure of the Horne–Östberg chronotype*
In the first step, we determined the number of clusters. Following the suggestions of Tibshirani and Walther (2005), we looked for solutions resulting in a prediction strength above 0.80. Doing so, in 49 out of 50 imputed data sets, the data-driven proposed number for clustering based on smartphone proxies for circadian preferences was 1. However, decreasing the prediction strength criterion to a value

Table 4. Descriptive statistics for day–night behaviour patterns

| Variable | Week | | Weekend | | |
| --- | --- | --- | --- | --- | --- |
| | Mean | SD | Mean | SD | Cohens' d [CI₉₅%] |
| Mean first event week | 7.89 | 1.31 | 8.96 | 1.30 | 0.82 [0.70, 0.93] |
| Mean last event week | 23.15 | 1.23 | 23.79 | 1.42 | 0.49 [0.37, 0.60] |
| Mean duration nightly inactivity week (h) | 8.68 | 1.20 | 9.02 | 1.45 | 0.26 [0.14, 0.37] |
| Mean number checking events week | 5.59 | 3.97 | 5.61 | 5.37 | 0.00 [-0.11, 0.12] |
| Mean duration checking events week (s) | 26.07 | 26.12 | 25.29 | 40.72 | -0.02 [-0.14, 0.09] |
| Mean first alarm event week | 7.19 | 1.29 | 7.47 | 1.68 | 0.19 [0.06, 0.33] |
| Mean number snoozing week | 1.33 | 1.76 | 1.33 | 2.04 | 0.00 [-0.13, 0.13] |
| Mean duration snoozing week (min) | 23.26 | 23.61 | 23.89 | 34.26 | 0.02 [-0.11, 0.16] |

The coefficients for first and last events represent times of the day. The decimal places indicate the percentage of a full hour. For example, 7.89 means 7:53 AM or 23.15 means 11:09 PM.

of 0.75 yielded a 2-cluster solution for all imputed data sets. Although the recommended predictive power was slightly missed, we further investigated $k-means$ clustering with $k = 2$. The averaged bootstrapped performance measures for the cluster-wise stability assessment show that each component of the 2-cluster solution turned out to be highly stable (cluster 1 $_{n = 296}$: $JC = 0.94$, $dissolved = 0$, $recovered = 100$; cluster 2 $_{n = 301}$: $JC = 0.93$, $dissolved = 0$, $recovered = 100$). However, the internal cluster validation coefficients indicated that the two clusters were poorly separable from each other and were not compact in themselves ($wb.ratio = 0.73$, $silhouette = 0.25$, $Dunn = 0.06$). To get a better understanding of the identified structure in the daily smartphone usage timing, descriptive statistics of the variables that were considered for clustering are displayed in Table 5. On average, participants assigned to cluster 2 had later first and last smartphone usage events on weekends and the daily 25%, 50%, and 75% timestamps for general, social interaction, and entertainment usage events on weekends were on average about 2 hours later. The mean number of snoozing events

was similar in both groups, but participants of cluster 2 on average snoozed approximately 3.5 minutes longer. As an external criterion, we considered the smartphone-based Roenneberg chronotype. The mean midpoint of sleep was $M = 3.90$ ($SD = 1.15$) for cluster 1 and $M = 5.19$ ($SD = 1.38$) for cluster 2.

To return to the question of whether we found different groups of individuals with similar smartphone usage patterns indicating circadian preferences, we refer to Table 5. Effect sizes for variables indicating sleep–wake timing are large, suggesting that participants assigned to cluster 2 have noticeable back-shifted diurnal smartphone usage patterns in comparison with participants assigned to cluster 1. Figure 1 shows, however, that the distributions of the two cluster groups overlap. A considerable proportion of participants could not be clearly assigned to one of the two clusters. Accordingly, the distribution based on the entire sample was not bimodal but only unimodal.

In the second step, we also explored the factorial structure of the smartphone-based proxies for the Horne–Östberg

Table 5. Descriptive statistics for smartphone usage indicating circadian preferences by clusters

| Variable | Cluster 1 | | Cluster 2 | | |
| --- | --- | --- | --- | --- | --- |
| | Mean | SD | Mean | SD | Cohens' d [CI₉₅%] |
| First/last events on weekends | | | | | |
| Mean time of the first event | 8.35 | 1.16 | 9.58 | 1.13 | 1.07 [0.90, 1.25] |
| Mean time of the last event | 23.09 | 1.18 | 24.51 | 1.29 | 1.15 [0.97, 1.32] |
| Mean on weekends daily timestamp of | | | | | |
| 25% general usage | 12.28 | 1.29 | 14.38 | 1.37 | 1.57 [1.39, 1.76] |
| 50% general usage | 15.34 | 1.30 | 17.62 | 1.20 | 1.82 [1.63, 2.01] |
| 75% general usage | 18.37 | 1.34 | 20.62 | 1.17 | 1.79 [1.60, 1.98] |
| 25% social interaction usage | 12.51 | 1.28 | 14.47 | 1.22 | 1.57 [1.38, 1.76] |
| 50% social interaction usage | 15.38 | 1.34 | 17.53 | 1.16 | 1.72 [1.53, 1.91] |
| 75% social interaction usage | 18.18 | 1.53 | 20.29 | 1.16 | 1.56 [1.37, 1.74] |
| 25% entertainment usage | 12.91 | 1.76 | 15.22 | 2.02 | 1.21 [1.03, 1.39] |
| 50% entertainment usage | 15.07 | 1.86 | 17.70 | 1.75 | 1.46 [1.27, 1.64] |
| 75% entertainment usage | 17.25 | 2.05 | 20.03 | 1.74 | 1.47 [1.28, 1.65] |
| Snoozing events on weekdays | | | | | |
| Mean number of snoozing events | 1.31 | 1.88 | 1.35 | 1.65 | 0.02 [-0.15, 0.20] |
| Mean duration of snoozing events | 21.53 | 22.01 | 24.91 | 24.97 | 0.14 [-0.03, 0.32] |

Except the snoozing variables, the coefficients represent times of the day and the corresponding standard deviations are given in hours. The decimal places indicate the percentage of a full hour. The mean daily timestamp of 25% general usage indicates that 25% of all activities on a given day had happened at this point in time. The mean number of snoozing events means the daily mean absolute frequency and the snoozing duration is in minutes.

Figure 1. Plots displaying the distribution of mean daily first and last events on weekdays versus weekends by cluster. The black line shows the distribution based on the total sample. The ordinate axis goes beyond midnight, because last events after midnight were added to 24. An event at 26 therefore means it happened at 2.00 AM.

Table 6. Exploratory factor analysis of the smartphone-sensed circadian preferences

| Variable | F1 | F2 | F3 | U |
|---|---|---|---|---|
| Mean time of the first event on weekends | 0.09 | 0.06 | **0.49** | 0.67 |
| Mean time of the last event on weekends | **0.54** | 0.03 | 0.04 | 0.66 |
| Mean daily timestamp of 25% general usage on weekends | 0.05 | 0.10 | **0.84** | 0.16 |
| Mean daily timestamp of 50% general usage on weekends | **0.44** | 0.17 | **0.47** | 0.20 |
| Mean daily timestamp of 75% general usage on weekends | **0.74** | 0.14 | 0.11 | 0.23 |
| Mean daily timestamp of 25% social interaction usage on weekends | 0.26 | -0.01 | **0.68** | 0.30 |
| Mean daily timestamp of 50% social interaction usage on weekends | **0.63** | -0.01 | **0.38** | 0.22 |
| Mean daily timestamp of 75% social interaction usage on weekends | **0.88** | 0.03 | 0.03 | 0.19 |
| Mean daily timestamp of 25% entertainment usage on weekends | -0.20 | **0.77** | **0.33** | 0.24 |
| Mean daily timestamp of 50% entertainment usage on weekends | 0.02 | **0.99** | 0.00 | 0.01 |
| Mean daily timestamp of 75% entertainment usage on weekends | **0.34** | **0.77** | -0.17 | 0.21 |
| Mean daily number of snoozing events on weekdays | 0.26 | 0.01 | -0.24 | 0.94 |
| Mean daily duration of snoozing events on weekdays | 0.27 | 0.00 | -0.19 | 0.94 |
| F2 | 0.46 | 1.00 | | |
| F3 | 0.52 | 0.47 | 1.00 | |

Maximum likelihood factor analysis, obliquely rotated (oblimin) with three factors. Loadings greater than the amount of 0.30 are in bold. The correlations between the factors are displayed at the bottom of the table. F1 = Factor 1; F2 = Factor 2, F3 = Factor 3; U = Uniqueness.

chronotype. The empirical Kaiser criterion suggested a 3-factorial solution accounting for 62% of the variance. The obliquely (oblimin) rotated factor matrix is displayed in Table 6. Factor 1 explained 23% of the variance and comprised behavioural indicators describing markers for later diurnal smartphone usage. In contrast, the behavioural variables loading high on factor 3 (19% variance explanation) described markers characteristic for early diurnal smartphone usage. The 50% timestamps for daily (general and social interaction) smartphone usage considerably loaded on both, factors 1 and 3. Finally, factor 2 explained 20% of the variance and reflected behavioural indicators of

smartphone usage for entertainment purposes independent of the time of the day. The two snoozing items did not load considerably on any factor. All factors were correlated (see Table 6).

### The Roenneberg chronotype and its correlates

The smartphone-based midpoint of sleep (MSF) and the sleep debt corrected version $MSF_{corr}$, which both indicate the Roenneberg chronotype, were approximately unimodally symmetrically distributed (see Figure 2). As no weekends without alarm clock usage were available for some participants, their $MSF$ could not be computed. Therefore, the

Figure 2.    Plots displaying the distributions of the local time of the midpoint of sleep (*MSF*) and its sleep debt corrected version (*MSF$_{corr}$*) and its relationship with age divided by gender.

following results are based on a subsample of $n = 497$ participants. On average, the mean *MSF* was at 4.52 AM (*SD* = 1.42) and the *MSF$_{corr}$* slightly earlier at 4.26 AM (*SD* = 1.47). The *MSF* and *MSF$_{corr}$* ranged between 0.75 PM and 9.91 AM. The *MSF* was weakly negatively related to nightly inactivity duration during the weeks ($r = -0.13$, $CI_{95\%}$ [-0.21, -0.04]) as well as the weekends ($r = -0.11$, $CI_{95\%}$ [-0.20, -0.03]). As suggested by Roenneberg et al. (2007), we used the *MSF$_{corr}$* for investigating the relationship of chronotype and demographics. Age ($r = -0.16$, $CI_{95\%}$ [-0.24, -0.07]) and gender ($r = -0.15$, $CI_{95\%}$ [-0.23, -0.06]) were both negatively related to the corrected midpoint of sleep, indicating that older and female participants had on average earlier chronotype values. However, the age correlation should be interpreted with caution, as the plot on the right side of Figure 2 indicates that it was probably caused by data points of older participants of whom we only had few in the sample ($Q_3 = 25$). The correlation disappears ($r_s = -0.03$, $CI_{95\%}$ [-0.12, 0.06]) when computing the Spearman correlation, which is only based on ranks.

**Day–night behaviours and personality traits**

Because our analysis of relationships between behavioural day–night patterns and personality is exploratory, we do not perform any hypothesis tests, nor do we speculate about correlations on a variable-by-variable basis. Instead, based on the correlation plot displayed in Figure 3, we want to show the general result pattern and address some conspicuities. Overall, Spearman correlations ranged between $r_s = -0.24$ (mean time of last events during the week and sense of duty) and $r_s = 0.15$ (mean time of the first event on weekends and carefreeness). As can be seen in Figure 3, the most striking aspect is that conscientiousness and its facets (except competence) were related to various day–night behaviours. First, more conscientious people on average had earlier mean and less varying daily points of time of first and last

smartphone usage events both during weeks and on weekends. Furthermore, their duration of nightly inactivity varied less on weekdays and they had lower values on the Roenneberg chronotype. Finally, individuals with higher values on the facet sense of duty snoozed on average less often and shorter on weekdays.

Further but less coherent patterns in Figure 3 can be seen for openness, extraversion, and emotional stability. For example, openness to imagination showed some positive relations to day–night behavioural indicators. Openness to value and norm system was associated positively with the mean number and duration of snoozing events, especially on weekdays. Higher extraversion was related to longer smartphone checking events during nights on weekdays. Furthermore, carefreeness as a facet of emotional stability was associated positively with later day–night activity patterns. Regarding demographics, female participants' first use on weekends and general last use was on average earlier. Accordingly, they also had lower Roenneberg chronotype values. However, no correlations of considerable size were found for age.

**Using multilevel modelling to explore social jetlag**

To investigate social jetlag, we explored compensatory sleep on weekends approximated as nightly inactivity duration by multilevel modelling. The duration of nightly inactivity on weekends was predicted by the duration of nightly inactivity during the week and the interindividual variables Roenneberg chronotype, big five personality traits, age, gender, and the averaged individual mean duration of nightly inactivity. The results are presented in the 12 panels in Figure 4, which show the estimates and their 95% confidence intervals across all multiverse data sets for each predictor in the model. Some aspects were evident across all data sets. There were no relationships between the nightly inactivity duration on weekends and the variables Roenneberg

Figure 3.   Pairwise complete Spearman correlations between smartphone-sensed day–night activities for weekdays versus weekends and personality traits. Male participants were coded as 0. As not all participants used alarm clock apps, the sample size for respective correlations was reduced ($n_{week}$ = 506, $n_{weekend}$ = 371). The colour of the squares indicates the direction and the strength of the respective correlations. For better readability, correlations are presented as percentage (e. g. a value of 3 means $r_s$ = 0.03). Additionally, only correlations with greater absolute values than 0.10 are highlighted in colour.

chronotype, openness, extraversion, agreeableness, emotional stability, and the interaction between the Roenneberg chronotype and the nightly weekday inactivity. Second, the averaged nightly inactivity duration across the study weeks (level 2) was positively associated with the nightly inactivity period on weekends. Nevertheless, estimates for the individual nightly inactivity duration on weekdays (level 1) and conscientiousness, age, and gender (all level 2) varied across the multiverse data sets. Depending on the preprocessing steps, individuals with longer nightly inactivity duration on weekdays in the corresponding week, higher conscientiousness, higher age, and male gender had, on average, longer nightly inactivity periods on weekends.

As can be seen in Figure 4, some patterns can be identified in the multiverse results across different variables: the coding of the weekend seemed to have an influence. In conditions in which the weekend was coded as nights between Friday and Monday, the mean duration of nightly inactivity on weekends was, on average, lower compared with the conditions in which weekends were coded as nights between Friday and Sunday. Also, for gender, a pattern can be determined depending on the coding of the weekend. For conscientiousness, estimates in conditions including 3 weeks were, on average, higher than conditions comprising 4 weeks. Regarding the average duration of nightly inactivity during the week (level 2), estimates were higher when winsorized and imputed.

To get a better understanding of the results concerning social jetlag, we calculated an additional multiverse analysis. For this purpose, we considered a variant of the multilevel model without personality traits and demographics as covariates. As results did not considerably differ and not to go beyond the scope of this paper, they can be found as a supplementary analysis in our OSF project.

## DISCUSSION

We investigated three prominent research questions related to common behavioural day–night patterns by using smartphone sensing data. First, we focused on individual differences in day–night activity patterns. Based on behavioural indicators of circadian preferences, we explored the structure underlying our smartphone proxy for the Horne–Östberg chronotype. Regarding the search for a smartphone chronotype, we found nondiscrete groups of individuals with similar diurnal smartphone usage patterns. In addition, our smartphone-based proxy for the Horne–Östberg chronotype turned out to be a multidimensional construct. In addition, we presented an algorithm for computing the chronotype as defined by Roenneberg et al. (2003). We used smartphone-based indicators for the midpoint of sleep and

Figure 4.    The decision tree on the left side shows how the multiverse of 16 data sets was created. The 12 panels on the right display the estimates and their 95% confidence intervals for the intercept and each predictor, resulting from multilevel modelling across the multiverse of 16 data sets. L1 = level 1 predictors (z-standardized and person-mean-centred); L2 = level 2 predictors (z-standardized, except gender). Male participants were coded as 0. The individual mean of the level 1 predictor was additionally entered as level 2 predictor. Each data set and the corresponding model are coded the same colour.

found associations with age, gender, and duration of nightly inactivity. Regarding personality traits, we found associations of conscientiousness with smartphone-sensed indicators for day–night behaviour. Finally, we explored social jetlag by examining whether people were inactive longer during weekend nights if they accumulated a deficit of nightly inactivity during the preceding workweek while controlling for individual differences. Our findings suggest that nightly inactivity duration on weekends was mainly related to individuals' general level of nightly inactivity across all study weeks. We will critically discuss our results in the following sections. Because our research was explorative, explanations drawn *post hoc* should not be easily generalized but be confirmed by preregistered hypotheses testing in future studies.

### Smartphone sensing in the context of behavioural day–night patterns

*Individual differences in day–night activity patterns*
In contrast to previous research based on self-reports, we used smartphone-sensed behavioural data to investigate the

structure of chronotype and to inform both the variable-centred and the person-centred approach to chronotype. Emphasizing chronotype as a continuous dimension reflecting circadian habits, Roenneberg et al. (2003) have suggested computing the midpoint of sleep. Instead of assessing these habits by questionnaires (e.g. Roenneberg et al., 2003; Roenneberg et al., 2007), we followed Lin et al. (2019) and used smartphone sensing data to determine a smartphone equivalent for the Roenneberg chronotype. We compared our resulting measure with the findings reported by Roenneberg et al. (2007) and found similar descriptive parameters (distribution, mean) and associations with external criteria like gender and sleep duration during the week. In accordance with our assumption that smartphone-based sleep–wake timing indicators overestimate sleep times, the range of values was slightly larger for our measure. Regarding age and chronotype, we found a negative correlation, which was caused by a few older participants with lower chronotype values. However, because the age composition of our sample was highly skewed towards younger participants, we do not want to over-interpret this finding. A nonmatching result was that whereas Roenneberg et al. (2007) found a positive correlation between chronotype and

sleep duration on weekends, we found a negative association. Roenneberg et al. (2007) argued that later chronotypes sleep longer on weekends because they collect a sleep debt during the workweek. In contrast to the large representative sample of their epidemiological study, our sample consisted mainly of students who are more likely to have fewer social obligations during the week than people who have a 9 AM to 5 PM job. Accordingly, compared with nights during the week, our participants' nightly inactivity (indicating sleep) did not differ considerably on weekends. Therefore, one interpretation of our results could be that students have the opportunity to be more flexible in their daily routines during the week following their chronotype. Therefore, late chronotypes do not need disproportionately more sleep on weekends. Accordingly, previous studies have shown that many students report napping after lunch during the week (Vela-Bueno et al., 2008). These naps could serve to use both the weekend and the week for sleep compensation (Gradisar et al., 2008). In line with our interpretation, students with late chronotypes have been found to nap more extensively than students with early chronotypes (Zimmermann, 2011). Please note that this is only our *post hoc* interpretation and further confirmatory research using behavioural data to study the interplay of sleep duration, chronotype, and work schedules is needed..

Keeping the focus on variable-centred trait assessment (Asendorpf, 2003), but following the Horne–Östberg tradition, we operationalized circadian preferences as diurnal smartphone usage behaviours and explored the underlying factorial structure. We found three correlated dimensions reflecting early use of the smartphone during the day, late use of the smartphone during the day, and entertainment usage. In comparison, findings of previous studies investigating the structure of self-reported chronotype have resulted in one to four factors (e.g. Caci et al., 2009; Lipnevich et al., 2017; Natale & Cicogna, 2002). In their recent meta-analysis, Lipnevich et al. (2017) concluded that the preferences for morningness versus eveningness are not the extreme poles of one dimension but two interdependent dimensions. Accordingly, our two correlated dimensions reflecting early and late diurnal smartphone usage activity align with their findings. Regarding our factor entertainment usage, we think that this could be regarded as a methodological artefact, as the content entertainment might have overlaid the diurnal character of the respective behavioural circadian indicators.

Dimensional approaches to personality, such as the two described above, offer the advantage to focus on individual differences. However, in contrast to person-centred approaches, they are not able to describe the structure of traits within persons (e.g. Asendorpf & van Aken, 1999; Asendorpf, 2003). In addition, types might have an advantage for applied purposes as the classification as 'morning larks' or 'night owls' is widely anchored in the popular science literature and scientific research. Therefore, besides examining dimensionality, we also explored the existence of types of individuals with similar diurnal smartphone usage patterns by using unsupervised machine learning. We found two groups that showed earlier versus later smartphone usage over the day. As the effect sizes show, these two groups considerably differed in indicators of diurnal smartphone usage

patterns. However, our results also indicate that despite the high average group differences, a large number of participants could not easily be assigned to one of these two groups, which overlapped considerably in the behavioural indicators used. Therefore, we asked ourselves whether we should call the structure we found types. In previous chronotype literature, types had often been considered as empirically validated, if the resulting groups subsequently proved to be different concerning external criteria (e.g. body temperature, electroencephalography (EEG) recordings Horne & Östberg, 1976; Putilov et al., 2015). In contrast, we did not determine any cut-off values but searched for nonrandom structures in the data. Only recently, Preckel et al. (2019) followed a similar approach identifying four chronotypes in an adolescent sample. However, because circadian preferences change with age (Roenneberg et al., 2007), and our sample was older, and we focused on smartphone-sensed rather than self-reported circadian habits, we argue that the results are not fully comparable.

From a statistical point of view, the existence of types is only justified if underlying variables are multimodally distributed (Hicks, 1984; Fleiss et al., 1971), which was not the case for our behavioural day–night indicators. However, previous research in the social sciences has revealed that nonoverlapping types hardly exist for human behaviours (Meehl, 2004; Costa Jr et al., 2002). Accordingly, Asendorpf and van Aken (1999) distinguish between discrete and nondiscrete types in the context of personality research. Thus, the criteria for defining types are not uniformly defined and applied in the literature. Our results are in line with this argument. Even if there were discrete underlying chronotype groups, it is unlikely that they would appear so clearly in everyday behavioural indicators due to social obligations and societal demands. Nevertheless, the identified nondiscrete groups in our study can be a good starting point towards a smartphone-based behavioural proxy of chronotype operationalized as circadian preferences. Future research should replicate the structure in diurnal smartphone usage indicators across different samples and use external validity criteria.

*Conscientiousness and differences in behavioural day–night patterns*
In contrast to the majority of previous studies, we used behavioural markers for day–night activity patterns to investigate associations with personality traits and demographics. In line with past studies showing women's preference for morningness (Randler, 2007), women in our study were earlier in the day, and their day–night activity varied less. Besides, our results were consistent with previous research showing a majorly coherent pattern of day–night activity and conscientiousness (Adan et al., 2012; Lipnevich et al., 2017) but less clear relations for other big five personality traits (e.g. Gray & Watson, 2002; Randler et al., 2017). Precisely, highly conscientious participants on average showed lower and less varying sleep–wake timing indicators and lower Roenneberg chronotype values. Following questionnaire-based research (Adan et al., 2012; Lipnevich et al., 2017; Tsaousis, 2010; Križan & Hisler, 2019), our

results indicate that more conscientious people on averagen are active earlier during the day and have longer nightly rest periods on weekends. Compared with findings from a meta-analysis ($r = .33$ according to Tsaousis, 2010), our correlations were smaller. However, our findings show that more conscientious people, who describe themselves as dutiful, ambitious, and disciplined (Arendasy, 2009), also act accordingly in everyday life (e.g. getting up early in the morning, longer nightly rest on weekends). Accordingly, Spears et al. (2019) found in a recent longitudinal study that conscientiousness was associated with mortality risk after 10 years and that this association was mediated by sleep duration as an everyday expression of behaviour.

In contrast to previous findings, conscientiousness and emotional stability were not related to indicators for sleep continuity, but extraversion was (Križan & Hisler, 2019; Sutin et al., 2019; Sella et al., 2020). These recent studies measured sleep continuity using actigraphy and therefore used completely different operationalizations of the related indicators sleep fragmentation and wake up after bed (Križan & Hisler, 2019; Sutin et al., 2019; Sella et al., 2020). For example, Sella et al. (2020) defined sleep fragmentation as the number of awakenings exceeding a certain duration. In contrast to actigraphy, smartphone sensing does not provide continuous measurement of wakefulness but approximates this measure via active smartphone usage. This requires the determination of a specific threshold value to classify smartphone usage either as part of a continuous usage phase belonging to the last or first event of the day or as a short usage event during the period of otherwise nightly inactivity. Determining a threshold value according to this principle, our approach has two significant drawbacks. First, using 2 minutes as a threshold was a subjective decision due to the lack of empirical data from previous literature. Second, the derived variable checking duration is restricted in its variance by a maximum value of 2 minutes. Consequently, individual differences in the actual wake after sleep onset might be masked by our smartphone-based operationalization, which in turn could explain the differences in findings compared with actigraphy.

n addition, we did not find some of the relationships that have previously been reported. For example, in our data, we did not find associations between a preference for morningness and agreeableness (Adan et al., 2012; Tsaousis, 2010) or age (Adan et al., 2012). As already discussed in the previous section, our results regarding age should be interpreted with caution due to the restricted variability of age in our sample. Overall, the differing findings could result from the usage of actual behavioural variables in contrast to self-reported preferences in most previous studies. Additionally, differences with past studies might not be surprising considering that previous questionnaire-based research is not clear either (e.g. Duggan et al., 2014; Gray & Watson, 2002). Besides, to the best of our knowledge, we have been the first to explore differences in alarm clock app usage. Our results provide first indications about the relation of snoozing behaviour and personality facets (sense of duty

and openness to value and norm system). They should be further investigated in future research.

*Individual differences in compensatory nightly inactivity on weekends*

To explore social jetlag, we investigated which intraindividual and interindividual factors predict the duration of nightly inactivity of smartphone usage (assumed to indicate sleep duration) on weekends. To explore this research question and to get an impression of the robustness of our estimates, we created a multiverse of 16 data sets resulting from combining different choices of plausible preprocessing steps. In the following, we focus only on those aspects that have been demonstrated across all data sets. Individuals who had higher overall levels of smartphone inactivity during nights on weekdays were also inactive longer on weekend nights. Even though our inactivity measure is not identical to sleep, our results indicate that individuals differ in their nightly rest duration. These findings support the notion that sleep duration is an independent trait (Ferrara & De Gennaro, 2001; Roenneberg et al., 2007). In contrast to the assumptions of social jetlag (Roenneberg et al., 2015; Wittmann et al., 2006), we neither found compensatory nightly inactivity on weekends nor any impact of the Roenneberg chronotype. As already discussed in the section above, our sample was highly skewed towards students. Thus, maybe their social obligations during the week are less pronounced, and therefore, we could not find their need for compensatory sleep on weekends. In addition, previous studies often used self-reports to investigate social jetlag (e.g. Wittmann et al., 2006; Roenneberg et al., 2012). Even though participants are instructed to indicate their habits for the last 4 weeks (Roenneberg et al., 2003), their answers might be biased towards a more general judgment of sleep–wake timing or influenced by short-term experiences like the sleep behaviour of the previous night. In contrast, we looked at behavioural snippets of 3 or 4 concrete weeks.

Finally, our multiverse analysis showed that the results depend on the selected preprocessing steps. Especially for the predictors age, gender, and conscientiousness, the size of the estimates differed depending on the constructed data sets. Our study therefore points to two problems. First, for behavioural indicators extracted from smartphone sensing data, the definition of the weekend and the number of weeks included made a difference to the results. Future research in the field of smartphone sensing should, therefore, carefully explore and report whether decisions made in the preprocessing have an impact on the results. Second, our study highlights the issue of selective reporting in research articles (Simonsohn et al., 2015; Steegen et al., 2016). We could just as well have reported only one of the paths and the results of the corresponding model, and the choice of each path would have been equally plausible. However, depending on the preprocessing decisions, we might or might not have emphasized the effect of conscientiousness or gender or age at this point. In line with Simonsohn et al. (2015) and Steegen et al. (2016), we argue that

decisions that might affect the results should be made transparent.

**Limitations and outlook**

Our study exemplifies the usage of smartphone sensing data in the research field of behavioural day–night patterns. Strictly speaking, the assessment of day–night structures in everyday life and, therefore, sleep-wake phases would require the collection of EEG data (Shambroom et al., 2012). For reasons of efficiency, self-report questionnaires have so far been used to approximate sleep-related behaviours. We propose smartphone sensing as an alternative to collect proxies for these behaviours. However, our approach has some limitations.

First, similar to questionnaires (Lauderdale et al., 2008), our behavioural markers are only proxies for actual sleep–wake timing. In our data set, only app, phone, screen, and notification events were available to determine the nightly inactivity period. Thus, actual sleep times were estimated based on active smartphone usage behaviours. However, for improving the accuracy of smartphone-based sleep–wake indicators, it would be helpful to include sensor data that do not require active usage, for example, brightness and ambient noise (Min et al., 2014). An even better estimate of sleep could be obtained by integrating the idea of actigraphy into the smartphone sensing approach. Meanwhile, many commercial wearables, which can also be used conveniently during bedtime, offer an open interface to integrate motion and physiological data like heart rate variability or galvanic skin response into research apps used for smartphone sensing.

Second, we defined new behavioural variables, which we extracted from smartphone sensing data. Although we derived our variables from previous literature, we had many degrees of freedom. Which period is defined as a weekend? What does active smartphone usage mean? How can daily values be aggregated? These questions are only a few examples for the vast amount of decisions we had to make during data preprocessing. To make this process as transparent as possible, we provide an extensive codebook and analyse a multiverse of data sets where appropriate. However, the researcher community should develop a common standard for sensing data so that the results obtained do not depend on the respective data preprocessing decisions in individual studies.

One further limitation of our study was the skewed sample. In comparison with previous epidemiological studies, it was skewed in terms of age and occupation. As age and work schedules are related to sleep–wake timings (Adan et al., 2012), future studies using smartphone sensing data should use more representative samples.

Finally, in our study, we only focused on smartphone sensing data. Although resulting indicators cannot be equated one-to-one with physiological sleep, smartphone sensing can nevertheless unobtrusively collect data in the field over a long period. This is very beneficial as far as day–night habits are investigated. However, in research focusing on constructs like sleep quality (Križan & Hisler, 2019), it is essential to measure a possible mismatch between behavioural sleep indicators in contrast to individual perceptions and feelings about sleep–wake timings. Consequently, the integration of the experience sampling method (e.g. Takano et al., 2014) could help to gain further interesting insights in individual differences into behavioural day–night patterns. Future studies could additionally benefit from combining actigraphy and smartphone sensing. Both methods assess actual behaviour but highlight different aspects of day–night activity patterns (Borger et al., 2019). In summary, we do not want to discuss whether self-reports, smartphone sensing, or actigraphy are better suitable for depicting actual behavioural day–night patterns. We think that all data collection approaches have their place and could be very fruitfully combined to gain better insights into human day–night behaviour patterns.

**CONCLUSION**

We used smartphone sensing data to extract behavioural variables usually assessed by self-reports in the context of day–night behaviours. Our study contributes to gain new insights into traits related to day–night behaviour patterns. First, we investigated two prominent operationalizations of chronotype: based on indicators for sleep–wake timing and diurnal activity, we found two overlapping groups of smartphone-based 'morning larks' and 'night owls' and two correlated dimensions that were similar to previously reported questionnaire-based factors. By computing a smartphone-based proxy, we presented a smartphone-sensed measure for the Roenneberg chronotype. Second, conscientiousness was related to earlier day schedules. In addition, we found individuals to differ in their overall level of nightly rest. We argue that it is important to understand individual differences in behavioural day–night patterns, as they previously have been found to be related to individuals' well-being and health. This work demonstrates that smartphone sensing provides an efficient and ecologically valid tool that can help to foster this understanding.

**SUPPORTING INFORMATION**

Additional supporting information may be found online in the Supporting Information section at the end of the article.

## REFERENCES

Adan, A., Archer, S. N., Hidalgo, M. P., Di Milia, L., Natale, V., & Randler, C. (2012). Circadian typology: A comprehensive review. *Chronobiology International*, *29*(9), 1153–1175. https://doi.org/10.3109/07420528.2012.719971

Allen, M., Poggiali, D., Whitaker, K., Marshall, T., & Kievit, R. (2019). Raincloud plots: A multi-platform tool for robust data visualization. *Wellcome Open Res*, *4*, 63. https://doi.org/10.12688/wellcomeopenres.15191.1

Arendasy, M. (2009). *BFSI: Big-Five Struktur-Inventar (Test & Manual)*. Mödling: Austria: SCHUHFRIED GmbH.

Asendorpf, J. B. (2003). Head-to-head comparison of the predictive validity of personality types and dimensions. *European Journal of Personality*, *17*(5), 327–346. https://doi.org/10.1002/per.492

Asendorpf, J. B., & van Aken, Marcel A. G. (1999). Resilient, overcontrolled, and undercontrolleed personality prototypes in childhood: Replicability, predictive power, and the trait-type issue. *Journal of Personality and Social Psychology*, *77*(4), 815–832. https://doi.org/10.1037/0022-3514.77.4.815

Au, Q. (2019). *fxtract: Feature extraction from grouped data*. https://github.com/QuayAu/fxtract

Bailey, S. L., & Heitkemper, M. M. (2001). Circadian rhythmicity of cortisol and body temperature: Morningness–eveningness effects. *Chronobiology International*, *18*(2), 249–261. https://doi.org/10.1081/CBI-100103189

Basagaña, X., Barrera-Gómez, J., Benet, M., Antó, J. M., & Garcia-Aymerich, J. (2013). A framework for multiple imputation in cluster analysis. *American Journal of Epidemiology*, *177*(7), 718–725. https://doi.org/10.1093/aje/kws289

Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, *67*(1), 1–48. https://doi.org/10.18637/jss.v067.i01

Baumeister, R. F., Vohs, K. D., & Funder, D. C. (2007). Psychology as the science of self-reports and finger movements: Whatever happened to actual behavior? *Perspectives on Psychological Science*, *2*(4), 396–403. https://doi.org/10.1111/j.1745-6916.2007.00051.x

Borger, J. N., Huber, R., & Ghosh, A. (2019). Capturing sleep–wake cycles by using day-to-day smartphone touchscreen interactions. *NPJ Digital Medicine*, *2*(1), 1–8.

Braeken, J., & Van Assen, M. A. (2017). An empirical Kaiser criterion. *Psychological Methods*, *22*(3), 450–466. https://doi.org/10.1037/met0000074

Caci, H., Deschaux, O., Adan, A., & Natale, V. (2009). Comparing three morningness scales: Age and gender effects, structure and cut-off criteria. *Sleep Medicine*, *10*(2), 240–245. https://doi.org/10.1016/j.sleep.2008.01.007

Cavallera, G., & Giudici, S. (2008). Morningness and eveningness personality: A survey in literature from 1995 up till 2006. *Personality and Individual Differences*, *44*(1), 3–21. https://doi.org/10.1016/j.paid.2007.07.009

Chen, Z., Lin, M., Chen, F., Lane, N. D., Cardone, G., Wang, R., …, & Campbell, A. T.(2013). *Unobtrusive sleep monitoring using smartphones*. In *Proceedings of the 7th International Conference on Pervasive Computing Technologies for Healthcare*, pp. 145–152. https://doi.org/10.4108/icst.pervasivehealth.2013.252148

Costa Jr, P. T., Herbst, J. H., McCrae, R. R., Samuels, J., & Ozer, D. J. (2002). The replicability and utility of three personality types. *European Journal of Personality*, *16*(S1), S73–S87. https://doi.org/10.1002/per.448

Costa, P. T., & McCrae, R. R. (2008). The Revised Neo Personality Inventory (NEO-PI-R), *The SAGE handbook of personality theory and assessment*, pp. 179–198. https://doi.org/10.4135/9781849200479.n9

Curran, P. J., & Bauer, D. J. (2011). The disaggregation of within-person and between-person effects in longitudinal models of change. *Annual Review of Psychology*, *62*, 583–619. https://doi.org/10.1146/annurev.psych.093008.100356

Díaz-Morales, J. F. (2007). Morning and evening-types: Exploring their personality styles. *Personality and Individual Differences*, *43*(4), 769–778. https://doi.org/10.1016/j.paid.2007.02.002

Duggan, K. A., Friedman, H. S., McDevitt, E. A., & Mednick, S. C. (2014). Personality and healthy sleep: The importance of conscientiousness and neuroticism. *PloS ONE*, *9*(3), e90628. https://doi.org/10.1371/journal.pone.0090628

Dunn, J. C. (1974). Well-separated clusters and optimal fuzzy partitions. *Journal of Cybernetics*, *4*(1), 95–104. https://doi.org/10.1080/01969727408546059

Ferrara, M., & De Gennaro, L. (2001). How much sleep do we need? *Sleep Medicine Reviews*, *5*(2), 155–179. https://doi.org/10.1053/smrv.2000.0138

Fleiss, J. L., Lawlor, W., Platman, S. R., & Fieve, R. R. (1971). On the use of inverted factor analysis for generating typologies. *Journal of Abnormal Psychology*, *77*(2), 127.

Gosling, S. D., John, O. P., Craik, K. H., & Robins, R. W. (1998). Do people know how they behave? Self-reported act frequencies compared with on-line codings by observers. *Journal of Personality and Social Psychology*, *74*(5), 1337–1349. https://doi.org/10.1037/0022-3514.74.5.1337

Gradisar, M., Wright, H., Robinson, J., Pain, S., & Gamble, A. (2008). Adolescent napping behavior: Comparisons of school week versus weekend sleep patterns. *Sleep and Biological Rhythms*, *6*(3), 183–186. https://doi.org/10.1111/j.1479-8425.2008.00351.x

Gray, E. K., & Watson, D. (2002). General and specific traits of personality and their relation to sleep and academic performance. *Journal of Personality*, *70*(2), 177–206. https://doi.org/10.1111/1467-6494.05002

Grund, S., Lüdtke, O., & Robitzsch, A. (2018). Multiple imputation of missing data for multilevel models: Simulations and recommendations. *Organizational Research Methods*, *21*(1), 111–149. https://doi.org/10.1177/1094428117703686

Halkidi, M., Batistakis, Y., & Vazirgiannis, M. (2001). On clustering validation techniques. *Journal of Intelligent Information Systems*, *17*(2-3), 107–145. https://doi.org/10.1023/A:1012801612483

Harari, G. M., Müller, S. R., Stachl, C., Wang, R., Wang, W., Bühner, M., Rentfrow, P. J., Campbell, A. T., & Gosling, S. D. (2019). Sensing sociability: Individual differences in young adults' conversation, calling, texting, and app use behaviors in daily life. *Journal of Personality and Social Psychology*. Advance online publication, https://doi.org/10.1037/pspp0000245

Hennig, C. (2007). Cluster-wise assessment of cluster stability. *Computational Statistics & Data Analysis*, *52*(1), 258–271. https://doi.org/10.1016/j.csda.2006.11.025

Hennig, C. (2008). Dissolution point and isolation robustness: Robustness criteria for general cluster analysis methods. *Journal of Multivariate Analysis*, *99*(6), 1154–1176. https://doi.org/10.1016/j.jmva.2007.07.002

Hennig, C. (2018). *fpc: Flexible procedures for clustering*. R package version 2.1-11.1, https://CRAN.R-project.org/package=fpc

Hicks, L. E. (1984). Conceptual and empirical analysis of some assumptions of an explicitly typological theory. *Journal of Personality and Social Psychology*, *46*(5), 1118–1131. https://doi.org/10.1037/0022-3514.46.5.1118

Horne, J. A., & Östberg, O. (1976). A self-assessment questionnaire to determine morningness–eveningness in human circadian rhythms. *International Journal of Chronobiology*, *4*(2), 97–110.

Katzenberg, D., Young, T., Finn, L., Lin, L., King, D. P., Takahashi, J. S., & Mignot, E. (1998). A clock polymorphism associated with human diurnal preference. *Sleep*, *21*(6), 569–576. https://doi.org/10.1093/sleep/21.6.569

Križan, Z., & Hisler, G. (2019). Personality and sleep: Neuroticism and conscientiousness predict behaviourally recorded sleep years

later. *European Journal of Personality*, *33*, 133–153. https://doi.org/10.1002/per.2191

Kuznetsova, A., Brockhoff, P. B., & Christensen, RuneH. B. (2017). lmerTest package: Tests in linear mixed effects models. *Journal of Statistical Software*, *82*(13), 1–26. https://doi.org/10.18637/jss.v082.i13

Lipnevich, A. A., Credé, M., Hahn, E., Spinath, F. M., Roberts, R. D., & Preckel, F. (2017). How distinctive are morningness and eveningness from the big five factors of personality? A meta-analytic investigation. *Journal of Personality and Social Psychology*, *112*(3), 491. https://doi.org/10.1037/pspp0000099

Masters, G. N. (1982). A rasch model for partial credit scoring. *Psychometrika*, *47*(2), 149–174. https://doi.org/10.1007/bf02296272

Mateo, M. J. C., Diaz-Morales, J. F., Barreno, C. E., Prieto, P. D., & Randler, C. (2012). Morningness–eveningness and sleep habits among adolescents: Age and gender differences. *Psicothema*, *24*(3), 410–415.

Meehl, P. E. (2004). What's in a taxon? *Journal of Abnormal Psychology*, *113*(1), 39. https://doi.org/10.1037/0021-843X.113.1.39

Min, J.-K., Doryab, A., Wiese, J., Amini, S., Zimmerman, J., & Hong, J. I. (2014). Toss'n'turn: Smartphone as sleep and sleep quality detector. *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, 477–486. https://doi.org/10.1145/2556288.2557220

Natale, V., & Cicogna, P. (2002). Morningness–eveningness dimension: Is it really a continuum? *Personality and Individual Differences*, *32*(5), 809–816. https://doi.org/10.1016/S0191-8869(01)00085-X

Ohayon, M., Wickwire, E. M., Hirshkowitz, M., Albert, S. M., Avidan, A., Daly, F. J., ..., & Hazen, N. (2017). National sleep foundation's sleep quality recommendations: First report. *Sleep Health*, *3*(1), 6–19. https://doi.org/10.1016/j.sleh.2016.11.006

Preckel, F., Fischbach, A., Scherrer, V., Brunner, M., Ugen, S., Lipnevich, A. A., & Roberts, R. D. (2019). Circadian preference as a typology: Latent-class analysis of adolescents' morningness/eveningness, relation with sleep behavior, and with academic outcomes. *Learning and Individual Differences*, *78*, 101725. https://doi.org/10.1016/j.lindif.2019.03.007

Putilov, A. A. (2017). Owls, larks, swifts, woodcocks and they are not alone: A historical review of methodology for multidimensional self-assessment of individual differences in sleep-wake pattern. *Chronobiology international*, *34*(3), 426–437. https://doi.org/10.1080/07420528.2017.1278704

Putilov, A. A., Donskaya, O. G., & Verevkin, E. G. (2015). How many diurnal types are there? a search for two further "bird species". *Personality and Individual Differences*, *72*, 12–17. https://doi.org/10.1016/j.paid.2014.08.003

R Core Team (2018). *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing. https://www.R-project.org/

Harari, G. M., Lane, N. D., Wang, R., Crosier, B. S., Campbell, A. T., & Gosling, S. D. (2016). Using smartphones to collect behavioral data in psychological science. *Perspectives on Psychological Science*, *11*(6), 838–854. https://doi.org/10.1177/1745691616650285

Kafadar, K. (2003). John tukey and robustness. *Statistical Science*, *18*(3), 319–331. https://doi.org/10.1214/ss/1076102419

Lauderdale, D. S., Knutson, K. L., Yan, L. L., Liu, K., & Rathouz, P. J. (2008). Self-reported and measured sleep duration. *Epidemiology*, *19*(6), 838–845. https://doi.org/10.1097/ede.0b013e318187a7b0

Lin, Y.-H., Wong, B.-Y., Lin, S.-H., Chiu, Y.-C., Pan, Y.-C., & Lee, Y.-H. (2019). Development of a mobile application (app) to delineate "digital chronotype" and the effects of delayed chronotype by bedtime smartphone use. *Journal of Psychiatric Research*, *110*, 9–15. https://doi.org/10.1016/j.jpsychires.2018.12.012

Newman, D. A. (2014). Missing data. *Organizational Research Methods*, *17*(4), 372–411. https://doi.org/10.1177/1094428114548590

Randler, C. (2007). Gender differences in morningness–eveningness assessed by self-report questionnaires: A meta-analysis. *Personality and Individual Differences*, *43*(7), 1667–1675. https://doi.org/10.1016/j.paid.2007.05.004

Randler, C. (2008). Morningness–eveningness, sleep–wake variables and big five personality factors. *Personality and Individual Differences*, *45*(2), 191–196. https://doi.org/10.1016/j.paid.2008.03.007

Randler, C., Díaz-Morales, J. F., Rahafar, A., & Vollmer, C. (2016). Morningness–eveningness and amplitude–development and validation of an improved composite scale to measure circadian preference and stability (messi). *Chronobiology international*, *33*(7), 832–848. https://doi.org/10.3109/07420528.2016.1171233

Randler, C., & Engelke, J. (2019). Gender differences in chronotype diminish with age: A meta-analysis based on morningness/chronotype questionnaires. *Chronobiology International*, *36*(7), 888–905. https://doi.org/10.1080/07420528.2019.1585867

Randler, C., Schredl, M., & Göritz, A. S. (2017). Chronotype, sleep behavior, and the big five personality factors. *Sage Open*, *7*(3), 1–9. https://doi.org/10.1177/2158244017728321

Revelle, W. (2018). *psych: Procedures for psychological, psychometric, and personality research*. R package version 1.8.12, https://CRAN.R-project.org/package=psych

Roenneberg, T. (2015). Having trouble typing? What on earth is chronotype? *Journal of Biological Rhythms*, *30*(6), 487–491. https://doi.org/10.1177/0748730415603835

Roenneberg, T., Allebrandt, K. V., Merrow, M., & Vetter, C. (2012). Social jetlag and obesity. *Current Biology*, *22*(10), 939–943. https://doi.org/10.1016/j.cub.2012.03.038

Roenneberg, T., Keller, L. K., Fischer, D., Matera, J. L., Vetter, C., & Winnebeck, E. C. (2015). Human activity and rest in situ. *Methods in Enzymology*, *552*, 257–283. https://doi.org/10.1016/bs.mie.2014.11.028

Roenneberg, T., Kuehnle, T., Juda, M., Kantermann, T., Allebrandt, K., Gordijn, M., & Merrow, M. (2007). Epidemiology of the human circadian clock. *Sleep Medicine Reviews*, *11*(6), 429–438. https://doi.org/10.1016/j.smrv.2007.07.005

Roenneberg, T., Wirz-Justice, A., & Merrow, M. (2003). Life between clocks: Daily temporal patterns of human chronotypes. *Journal of Biological Rhythms*, *18*(1), 80–90. https://doi.org/10.1177/0748730402239679

Roepke, S. E., & Duffy, J. F. (2010). Differential impact of chronotype on weekday and weekend sleep timing and duration. *Nature and Science of Sleep*, *2*, 213–220. https://doi.org/10.2147/NSS.S12572

Rousseeuw, P. J. (1987). Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, *20*, 53–65. https://doi.org/10.1016/0377-0427(87)90125-7

Rousseeuw, P. J., & Croux, C. (1993). Alternatives to the median absolute deviation. *Journal of the American Statistical Association*, *88*(424), 1273–1283. https://doi.org/10.2307/2291267

Schoedel, R., Au, Q., Völkel, S. T., Lehmann, F., Becker, D., Bühner, M., Bischl, B., Hussmann, H., & Stachl, C. (2018). Digital footprints of sensation seeking. *Zeitschrift für Psychologie*, *226*(4), 232–245. https://doi.org/10.1027/2151-2604/a000342

Schoedel, R., Pargent, F., Au, Q., Völkel, S., Schuwerk, T., Bühner, M., & Stachl, C. (2020). To challenge the morning lark and the night owl: Using smartphone sensing data to investigate day-night behavior patterns. *European Journal of Personality*, *34*, 733–752. https://doi.org/10.1002/per.2258

Schuwerk, T., Kaltefleiter, L. J., Au, J.-Q., Hoesl, A., & Stachl, C. (2019). Enter the wild: Autistic traits and their relationship to mentalizing and social interaction in everyday life. *Journal of Autism and Developmental Disorders*, *49*(10), 4193–4208. https://doi.org/10.1007/s10803-019-04134-6

Sella, E., Carbone, E., Toffalini, E., & Borella, E. (2020). Personality traits and sleep quality: The role of sleep-related beliefs. *Personality and Individual Differences*, *156*, 109770. https://doi.org/10.1016/j.paid.2019.109770

Shah, A. D., Bartlett, J. W., Carpenter, J., Nicholas, O., & Hemingway, H. (2014). Comparison of random forest and parametric imputation models for imputing missing data using mice: A caliber study. *American Journal of Epidemiology*, *179*(6), 764–774. https://doi.org/10.1093/aje/kwt312

Shambroom, J. R., Fábregas, S. E., & Johnstone, J. (2012). Validation of an automated wireless system to monitor sleep in healthy adults. *Journal of Sleep Research*, *21*(2), 221–230. https://doi.org/10.1111/j.1365-2869.2011.00944.x

Simonsohn, U., Simmons, J. P., & Nelson, L. D. (2015). *Specification curve: Descriptive and inferential statistics o all reasonable specifications*. SSRN, https://doi.org/10.2139/ssrn.2694998

Spears, S. K., Montgomery-Downs, H. E., Steinman, S. A., Duggan, K. A., & Turiano, N. A.(2019). Sleep: A pathway linking personality to mortality risk. *Journal of Research in Personality*, *81*, 11–24. https://doi.org/10.1016/j.jrp.2019.04.007

Stachl, C., Au, Q., Schoedel, R., Buschek, D., Völkel, S., Schuwerk, T., …, & Bühner, M. (2019). *Behavioral patterns in smartphone usage predict big five personality traits*. https://doi.org/10.31234/osf.io/ks4vd

Stachl, C., Hilbert, S., Au, J.-Q., Buschek, D., De Luca, A., Bischl, B., Hussmann, H., & Bühner, M. (2017). Personality traits predict smartphone usage. *European Journal of Personality*, *31*(6), 701–722. https://doi.org/10.1002/per.2113

Stachl, C., Schoedel, R., Au, Q., Völkel, S., Buschek, D., Hussmann, H., …, & Bühner, M. (2018). *The phonestudy project*. https://doi.org/10.17605/osf.io/ut42y

Steegen, S., Tuerlinckx, F., Gelman, A., & Vanpaemel, W. (2016). Increasing transparency through a multiverse analysis. *Perspectives on Psychological Science*, *11*(5), 702–712. https://doi.org/10.1177/1745691616658637

Sutin, A. R., Gamaldo, A. A., Stephan, Y., Strickhouser, J. E., & Terracciano, A. (2019). Personality traits and the subjective and objective experience of sleep. *International Journal of Behavioral Medicine*. https://doi.org/10.1007/s12529-019-09828-w

Takano, K., Sakamoto, S., & Tanno, Y. (2014). Repetitive thought impairs sleep quality: An experience sampling study. *Behavior Therapy*, *45*(1), 67–82. https://doi.org/10.1016/j.beth.2013.09.004

Tan, P.-N., Steinbach, M., & Kumar, V. (2006). *Cluster analysis: Basic concepts and algorithms*.

Terman, J. S., Terman, M., Lo, E.-S., & Cooper, T. B. (2001). Circadian time of morning light administration and therapeutic response in winter depression. *Archives of General Psychiatry*, *58*(1), 69–75. https://doi.org/10.1001/archpsyc.58.1.69

Tibshirani, R., & Walther, G. (2005). Cluster validation by prediction strength. *Journal of Computational and Graphical Statistics*, *14*(3), 511–528. https://doi.org/10.1198/106186005X59-243

Tonetti, L., Pascalis, V. D., Fabbri, M., Martoni, M., Russo, P. M., & Natale, V. (2016). Circadian typology and the alternative five-factor model of personality. *International Journal of Psychology*, *51*(5), 332–339. https://doi.org/10.1002/ijop.12170

Tsaousis, I. (2010). Circadian preferences and personality traits: A meta-analysis. *European Journal of Personality*, *24*(4), 356–373. https://doi.org/10.1002/per.754

Ushey, K., McPherson, J., Cheng, J., Atkins, A., & Allaire, J. (2018). *packrat: A dependency management system for projects and their r package dependencies*. R package version 0.4.9-3, https://CRAN.R-project.org/package=packrat

van Buuren, S., & Groothuis-Oudshoorn, K. (2011). mice: Multivariate Imputation by Chained Equations in R. *Journal of Statistical Software*, *45*(3), 1–67. https://www.jstatsoft.org/v45/i03/

Vela-Bueno, A., Fernandez-Mendoza, J., Olavarrieta-Bernardino, S., Vgontzas, A. N., Bixler, E. O., de la Cruz-Troca, J. J., Rodriguez-Muñoz, A., & Oliván-Palacios, J. (2008). Sleep and behavioral correlates of napping among young adults: A survey of first-year university students in madrid, spain. *Journal of*

American College Health*, *57*(2), 150–158. https://doi.org/10.3200/jach.57.2.150-158

Vitale, J. A., Roveda, E., Montaruli, A., Galasso, L., Weydahl, A., Caumo, A., & Carandente, F.(2015). Chronotype influences activity circadian rhythm and sleep: Differences in sleep quality between weekdays and weekend. *Chronobiology International*, *32*(3), 405–415. https://doi.org/10.3109/07420528.2014.986273

Wei, T., & Simko, V. (2017). *R package corrplot: Visualization of a correlation matrix*. R package version 0.84, https://github.com/taiyun/corrplot

Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. New York: Springer-Verlag. http://ggplot2.org

Wickham, H., François, R., Henry, L., & Müller, K. (2019). *Dplyr: A grammar of data manipulation*. R package version 0.8.0.1, https://CRAN.R-project.org/package=dplyr

Wittmann, M., Dinich, J., Merrow, M., & Roenneberg, T. (2006). Social jetlag: misalignment of biological and social time. *Chronobiology International*, *23*(1-2), 497–509. https://doi.org/10.1080/07420520500545979

Zavada, A., Gordijn, M. C., Beersma, D. G., Daan, S., & Roenneberg, T. (2005). Comparison of the Munich Chronotype Questionnaire with the Horne--Östberg's Morningness–Eveningness Score. *Chronobiology International*, *22*(2), 267–278. https://doi.org/10.1081/CBI-200053536

Zimmermann, L. K. (2011). Chronotype and the transition to college life. *Chronobiology International*, *28*(10), 904–910. https://doi.org/10.3109/07420528.2011.618959

## APPENDIX A: SUPPLEMENTAL METHOD

### Measures

Table A1. Description of the algorithm for detecting nightly inactivity

| Step | Description |
| --- | --- |
| 1 | Exclude passive smartphone events (GPS logs, notifications, and related screen events) |
| 2 | Exclude active usage events lasting shorter than two minutes and label them as checking behavior |
| 3 | Search for the maximum distance between consecutive events |
| 4 | Label the starting point of the maximum distance as last event of the day and the end point as first event of the next day |

To avoid longer periods of inactivity being detected during the day, the time frame for maximum distance detection was limited to 6.00 PM to 2.00 PM of the following day. We defined and filtered checking behaviour, because we wanted to exclude less significant actions like checking the clock or notification texts.

### Clustering

#### K-Means Algorithm

For clustering, we used the *k-means* algorithm, which is one of the most frequently used algorithms for clustering (Tan et al., 2006). In the following section, we only describe the basic principles behind *k-means* clustering and refer the interested reader to Tan et al. (2006) for a detailed explanation. After the user has defined the expected number of clusters $k$, $k$ points in the sample data are randomly determined and represent initial centroids. In the second step, all remaining data

points are assigned to the centroid for which the Euclidean distance is lowest. Afterward, the centroids in each of the *k* clusters are updated by calculating the arithmetic mean of all points in the respective clusters. Step by step, the procedures are repeated as long as the centroids do not change anymore, which indicates that the grouping structure in the data has been identified. As the centroid represents the data points within the clusters, *k-means* clustering is also often referred to as prototype-based or partitional clustering (Tan et al., 2006).

*Evaluation Metrics*
To ensure cluster validity, we took several steps to find nonrandom structures in our data. The first step is to determine the appropriate number of clusters. Tibshirani and Walther (2005) proposed to reframe clustering as a supervised prediction problem by splitting the data into a training and a test set and estimating the number of pairwise cases that are assigned to the same cluster in the test set based on centroids of the training set. The associated prediction strength measure defined by Tibshirani and Walther (2005) can be used to determine an optimal number of clusters. Another important aspect is cluster stability (Hennig, 2007). If clusters disappear when data are slightly modified, they are not regarded as stable and consequently might reflect only random structure. Hennig (2007) therefore suggests bootstrapping the data and considering the Jaccard coefficient (*JC*) for each cluster separately. The *JC* gives the proportion of data points (participants) that are assigned to the same cluster across the bootstrapped iterations, thus expressing the similarity of cluster solutions across bootstrapped data sets on a cluster-wise basis (Hennig, 2007). Further descriptive measures of cluster stability are the criteria of *recovery* and *dissolution,* which count how often each cluster has been successfully recovered and dissolved across all bootstrap iterations (Hennig, 2007; 2008). As recommended by Hennig (2018), we used 100 bootstrap replications and interpreted clusters as stable if the *JC* exceeded values above 0.85.

*Imputation of Missing Values*
Based on a variable-by-variable procedure, missings are replaced by values of a conditional distribution, which results from estimating imputation models using the remaining variables of the data set (van Buuren & Groothuis-Oudshoorn, 2011). We chose the random forest as an imputation algorithm as it has been proven useful for complex, incomplete data problems (Shah et al., 2014). To reduce the imputation bias caused by stochastic variation, we specified 50 imputation models. For each of the resulting 50 data sets, we performed a separate cluster analysis and report the mean/modus of the performance coefficients and cluster membership across data sets (Basagaña et al., 2013).

**Multilevel Modelling**

*Decisions in the Multiverse*
For constructing the data multiverse (Steegen et al., 2016), we considered the following decisions concerning preprocessing steps:

*Descision 1: Coding of Weekend.* In an earlier draft of the manuscript, we defined the weekend not as a period from *Friday to Sunday*, but from *Friday to Monday*. We found it challenging to decide whether Sunday evening and the following night still belong to the weekend or whether it is more of a weekday in terms of sleep–wake behaviour. In sleep research, the nights from Friday to Saturday and from Saturday to Sunday are considered as weekends traditionally. Because on Monday, one usually has to attend to social obligations again, sleep behaviour during the night from Sunday to Monday is assumed not to be chosen as freely and used to balance the weekly sleep deficit as the other two weekend nights (Roenneberg et al., 2007). Despite the standard in sleep research, we want to include both variants in our multilevel modelling and thus make our research process transparent.

*Decision 2: Number of Weeks.* We considered the number of repeated measurements to be plausible as both *3* and *4 weeks* because we noticed during the aggregation of the raw timestamped event data that some participants had only partially participated in the last weekend (e.g. only on Saturday, no longer on Sunday).

*Decision 3: Outliers.* For the handling of outliers, we found two points of view plausible. First, smartphone sensing-derived variables are usually susceptible to distortion due to data errors, which do not matter if enough data are aggregated using robust measures over a longer period. However, as for week-based variables, only a few single data points can be summarized; outliers due to data errors are more problematic. Therefore, we identified outliers as cases *deviating more than three times the mean absolute deviation from the median* and replaced them by the *person-specific median* of the corresponding variable. Second, the identification of outliers arising from the underlying smartphone usage behaviour can be emphasized. In this case, it would be plausible to use a method for outlier handling that limits the variability of the smartphone indicators less than using the median. To cover this aspect, we used *winsorization* as the second alternative.

*Decision 4: Missing Values.* Dealing with missing values in multilevel models is a challenging task. Traditionally, *listwise deletion* has been used, which uses only complete observations for estimating the model (e.g. Newman, 2014). Besides the disadvantage of the reduced sample and power, results are likely to be biased if the incomplete observations differ systematically from complete observations (Newman, 2014; Grund et al., 2018). An alternative approach to deal with missing data is to apply *multiple imputation*. However, in the context of multilevel models, this is not a trivial task as the imputation model itself should consider the multilevel structure. Current methods and software implementations are reaching their limits if more complicated use cases like random slopes or cross-level interactions are included in the model (Grund et al., 2018). For our analyses, we used the multivariate imputation by chained equations technique and implemented a random slope imputation model with group-level variables as proposed by Grund et al. (2018). Please note the imputation bias because we were unable to

integrate cross-level interactions with existing software implementations. In addition, Grund et al. (2018) point out that this area of research is still ongoing and that there are no clear recommendations for dealing with missing data in use cases such as ours.

### Model Description

To comprehensibly illustrate the multilevel model used for the multiverse analysis, we present the pseudo-model equation using the lmer syntax of the *lme4 package* in R (Bates et al., 2015). We specified a *random-intercept-random-slope model* predicting the mean duration of nightly inactivity on weekends based on the mean nightly inactivity duration during the previous week (level 1). Chronotype, the big five traits, age, and gender were included as level 2 predictors. The level 1 predictor duration of nightly inactivity during the week was person centred and the individual mean was entered as level 2 predictor (Curran & Bauer, 2011). Besides, the cross-level interaction of the mean nightly inactivity duration during the previous week and chronotype was added:

$$
\begin{aligned}
\text{Nightly Inactivity}_{\text{weekend}} \sim {} & 1 \\
& + \text{Nightly Inactivity}_{\text{week}}(\text{L1, z, pc}) \\
& + \text{Chronotype (L2, z, gc)} \\
& + \text{Nightly Inactivity}_{\text{week}}(\text{L2, z, gc}) \\
& + \text{Openness (L2, z, gc)} \\
& + \text{Conscientiousness (L2, z, gc)} \\
& + \text{Extraversion (L2, z, gc)} + \text{Agreeableness (L2, z, gc)} \\
& + \text{Emotional Stability (L2, z, gc)} + \text{Age (L2, z, gc)} \\
& + \text{Gender (L2, dc)} + \text{Nightly Inactivity}_{\text{week}}(\text{L1, z, pc}) \\
& * \text{Chronotype (L2, z, gc)} + (1 + \text{Nightly Inactivity}_{\text{week}} \\
& (\text{L1, z, pc}) | \text{ userid}),
\end{aligned}
\tag{A1}
$$

where L1 denotes predictors on level 1, L2 denotes predictors on level 2, z denotes that predictors were *z*-standardized, pc denotes that predictors were person-mean-centred, gc denotes that predictors were grand-mean-centred, and dc denotes that gender was dummy coded (0 = male, 1 = female).

## APPENDIX B: SUPPLEMENTAL RESULTS

### Big Five Personality Traits

Table B1. Descriptive statistics of personality factors and facets

| Variable | M | SD | alpha CI95% |
|---|---|---|---|
| Openness | -0.05 | 0.71 | [0.93, 0.94] |
| O1: Openness to imagination | 1.28 | 1.41 | [0.84, 0.87] |
| O2: Openness to aesthetics | 0.37 | 1.29 | [0.85, 0.88] |
| O3: Openness to feelings | 2.05 | 2.09 | [0.91, 0.93] |
| O4: Openness to actions | 1.35 | 1.4 | [0.84, 0.87] |
| O5: Openness to ideas | 1.65 | 1.42 | [0.82, 0.86] |
| O6: Openness to value/norm system | 0.9 | 1.02 | [0.73, 0.79] |
| Conscientiousness | -0.09 | 0.74 | [0.95, 0.96] |
| C1: Competence | 0.84 | 1.22 | [0.76, 0.82] |
| C2: Love of order | 1.1 | 1.58 | [0.87, 0.90] |
| C3: Sense of duty | 1.93 | 1.41 | [0.80, 0.85] |
| C4: Ambition | 1.83 | 1.68 | [0.86, 0.89] |
| C5: Discipline | 1.45 | 1.46 | [0.81, 0.86] |
| C6: Caution | 1.51 | 1.34 | [0.80, 0.84] |
| Extraversion | -0.01 | 0.74 | [0.95, 0.96] |
| E1: Friendliness | 1.45 | 1.29 | [0.80, 0.84] |
| E2: Sociableness | 1.3 | 1.74 | [0.89, 0.92] |
| E3: Assertiveness | 0.45 | 1.38 | [0.84, 0.87] |
| E4: Dynamism | 1.2 | 1.59 | [0.85, 0.88] |
| E5: Adventurousness | 0.45 | 1.49 | [0.88, 0.91] |
| E6: Cheerfulness | 1.97 | 1.64 | [0.86, 0.89] |
| Ageeableness | -0.06 | 0.75 | [0.92, 0.94] |
| A1: Willingness to trust | 0.4 | 1.43 | [0.86, 0.89] |
| A2: Genuineness | 1.01 | 0.94 | [0.61, 0.70] |
| A3: Helpfulness | 1.65 | 1.38 | [0.77, 0.82] |
| A4: Obligingness | 1.17 | 1.31 | [0.81, 0.85] |
| A5: Modesty | 0.77 | 1.13 | [0.79, 0.84] |
| A6: Good naturedness | 2.1 | 1.77 | [0.84, 0.88] |
| Emotional Stability | -0.03 | 0.71 | [0.93, 0.94] |
| ES1: Carefreeness | 0.12 | 1.3 | [0.82, 0.86] |
| ES2: Equanimity | 0.57 | 1.07 | [0.78, 0.83] |
| ES3: Positive mood | 0.95 | 1.43 | [0.84, 0.88] |
| ES4: Self consciousness | 0.66 | 1.18 | [0.83, 0.86] |
| ES5: Self control | 0.64 | 1 | [0.74, 0.81] |
| ES6: Emotional robustness | 0.65 | 1.19 | [0.80, 0.85] |

$N$ = 597; Alpha CI95% = 95% bootstrapped confidence intervals for Cronbach alpha coefficients.

# Enter the Wild: Autistic Traits and Their Relationship to Mentalizing and Social Interaction in Everyday Life

This article investigates the relationship between mentalizing, autistic traits, and social behavior in everyday life.

### *Contributing article*

Schuwerk, T., Kaltefleiter, L. J., Au, J.-Q., Hoesl, A., and Stachl, C. (2019). Enter the wild: Autistic traits and their relationship to mentalizing and social interaction in everyday life. *Journal of Autism and Developmental Disorders*, 49(10):4193–4208 This publication was part of the PhoneStudy project (see Chapter 7).

### *Copyright information*

See included licence agreement.

### *Declaration of contributions*

The preliminary work carried out by Jiew-Quay Au on other projects within the PhoneStudy project played a crucial role in this article. The PhD candidate also took charge of data preprocessing and extracting relevant variables for the data analysis. Furthermore, the doctoral candidate made valuable contributions to the review and editing process of the manuscript.

#### *Contribution of the coauthors*

Tobias Schuwerk developed the study concept and wrote most of the manuscript as the main author. Larissa Kaltefleiter played a key role by preparing the app usage data and also contributed by writing some parts of the methodology section. The other coauthors assisted in setting up the data collection process and helped in revising the manuscript.

| | |
|---|---|
| License Number | 5606570242398 |
| License date | Aug 12, 2023 |
| Licensed Content Publisher | Springer Nature |
| Licensed Content Publication | Journal of Autism and Developmental Disorders |
| Licensed Content Title | Enter the Wild: Autistic Traits and Their Relationship to Mentalizing and Social Interaction in Everyday Life |
| Licensed Content Author | Tobias Schuwerk et al |
| Licensed Content Date | Jul 4, 2019 |
| Type of Use | Thesis/Dissertation |
| Requestor type | academic/university or research institute |
| Format | print and electronic |
| Portion | full article/chapter |
| Will you be translating? | no |
| Circulation/distribution | 1 - 29 |
| Author of this Springer Nature content | yes |
| Title | PhD Thesis: Challenges in Machine Learning for Predicting Psychological Attributes from Smartphone Data |
| Institution name | Ludwig Maximilians Universität München |
| Expected presentation date | Oct 2023 |
| Requestor Location | Jiew-Quay Au<br>Am Prangebach 17<br><br>Gelsenkirchen, 45896<br>Germany<br>Attn: Jiew-Quay Au |

Terms and Conditions

**Springer Nature Customer Service Centre GmbH Terms and Conditions**

The following terms and conditions ("Terms and Conditions") together with the terms specified in your [RightsLink] constitute the License ("License") between you as Licensee and Springer Nature Customer Service Centre GmbH as Licensor. By clicking 'accept' and completing the transaction for your use of the material ("Licensed Material"), you confirm your acceptance of and obligation to be bound by these Terms and Conditions.

**1. Grant and Scope of License**

1. 1. The Licensor grants you a personal, non-exclusive, non-transferable, non-sublicensable, revocable, world-wide License to reproduce, distribute, communicate to the public, make available, broadcast, electronically transmit or create derivative works using the Licensed Material for the purpose(s) specified in your RightsLink Licence Details only. Licenses are granted for the specific use requested in the order and for no other use, subject to these Terms and Conditions. You acknowledge and agree that the rights granted to you under this License do not include the right to modify, edit, translate, include in collective works, or create derivative works of the Licensed Material in whole or in part unless expressly stated in your RightsLink Licence Details. You may use the Licensed Material only as permitted under this Agreement and will not reproduce, distribute, display, perform, or otherwise use or exploit any Licensed Material in any way, in whole or in part, except as expressly permitted by this License.

1. 2. You may only use the Licensed Content in the manner and to the extent permitted by these Terms and Conditions, by your RightsLink Licence Details and by any applicable laws.

1. 3. A separate license may be required for any additional use of the Licensed Material, e.g. where a license has been purchased for print use only, separate permission must be obtained for electronic re-use. Similarly, a License is only valid in the language selected and does not apply for editions in other languages unless additional translation rights have been granted separately in the License.

1. 4. Any content within the Licensed Material that is owned by third parties is expressly excluded from the License.

1. 5. Rights for additional reuses such as custom editions, computer/mobile applications, film or TV reuses and/or any other derivative rights requests require additional permission and may be subject to an additional fee. Please apply to journalpermissions@springernature.com or bookpermissions@springernature.com for these rights.

**2. Reservation of Rights**

Licensor reserves all rights not expressly granted to you under this License. You acknowledge and agree that nothing in this License limits or restricts Licensor's rights in or use of the Licensed Material in any way. Neither this License, nor any act, omission, or statement by Licensor or you, conveys any ownership right to you in any Licensed Material, or to any element or portion thereof. As between Licensor and you, Licensor owns and retains all right, title, and interest in and to the Licensed Material subject to the license granted in Section 1.1. Your permission to use the Licensed Material is expressly conditioned on you not impairing Licensor's or the applicable copyright owner's rights in the Licensed Material in any way.

**3. Restrictions on use**

3. 1. Minor editing privileges are allowed for adaptations for stylistic purposes or formatting purposes provided such alterations do not alter the original meaning or intention of the Licensed Material and the new figure(s) are still accurate and representative of the Licensed Material. Any other changes including but not limited to, cropping, adapting, and/or omitting material that affect the meaning, intention or moral rights of the author(s) are strictly prohibited.

3. 2. You must not use any Licensed Material as part of any design or trademark.

3. 3. Licensed Material may be used in Open Access Publications (OAP), but any such reuse must include a clear acknowledgment of this permission visible at the same time as the figures/tables/illustration or abstract and which must indicate that the Licensed Material is not part of the governing OA license but has been reproduced with permission. This may be indicated according to any standard referencing system but

must include at a minimum 'Book/Journal title, Author, Journal Name (if applicable), Volume (if applicable), Publisher, Year, reproduced with permission from SNCSC'.

**4. STM Permission Guidelines**

4. 1. An alternative scope of license may apply to signatories of the STM Permissions Guidelines ("STM PG") as amended from time to time and made available at https://www.stm-assoc.org/intellectual-property/permissions/permissions-guidelines/.

4. 2. For content reuse requests that qualify for permission under the STM PG, and which may be updated from time to time, the STM PG supersede the terms and conditions contained in this License.

4. 3. If a License has been granted under the STM PG, but the STM PG no longer apply at the time of publication, further permission must be sought from the Rightsholder. Contact journalpermissions@springernature.com or bookpermissions@springernature.com for these rights.

**5. Duration of License**

5. 1. Unless otherwise indicated on your License, a License is valid from the date of purchase ("License Date") until the end of the relevant period in the below table:

| | |
|---|---|
| Reuse in a medical communications project | Reuse up to distribution or time period indicated in License |
| Reuse in a dissertation/thesis | Lifetime of thesis |
| Reuse in a journal/magazine | Lifetime of journal/magazine |
| Reuse in a book/textbook | Lifetime of edition |
| Reuse on a website | 1 year unless otherwise specified in the License |
| Reuse in a presentation/slide kit/poster | Lifetime of presentation/slide kit/poster. Note: publication whether electronic or in print of presentation/slide kit/poster may require further permission. |
| Reuse in conference proceedings | Lifetime of conference proceedings |
| Reuse in an annual report | Lifetime of annual report |

| | |
|---|---|
| Reuse in training/CME materials | Reuse up to distribution or time period indicated in License |
| Reuse in newsmedia | Lifetime of newsmedia |
| Reuse in coursepack/classroom materials | Reuse up to distribution and/or time period indicated in license |

**6. Acknowledgement**

6. 1. The Licensor's permission must be acknowledged next to the Licensed Material in print. In electronic form, this acknowledgement must be visible at the same time as the figures/tables/illustrations or abstract and must be hyperlinked to the journal/book's homepage.

6. 2. Acknowledgement may be provided according to any standard referencing system and at a minimum should include "Author, Article/Book Title, Journal name/Book imprint, volume, page number, year, Springer Nature".

**7. Reuse in a dissertation or thesis**

7. 1. Where 'reuse in a dissertation/thesis' has been selected, the following terms apply: Print rights of the Version of Record are provided for; electronic rights for use only on institutional repository as defined by the Sherpa guideline (www.sherpa.ac.uk/romeo/) and only up to what is required by the awarding institution.

7. 2. For theses published under an ISBN or ISSN, separate permission is required. Please contact journalpermissions@springernature.com or bookpermissions@springernature.com for these rights.

7. 3. Authors must properly cite the published manuscript in their thesis according to current citation standards and include the following acknowledgement: '*Reproduced with permission from Springer Nature'.*

**8. License Fee**

You must pay the fee set forth in the License Agreement (the "License Fees"). All amounts payable by you under this License are exclusive of any sales, use, withholding, value added or similar taxes, government fees or levies or other assessments. Collection

and/or remittance of such taxes to the relevant tax authority shall be the responsibility of the party who has the legal obligation to do so.

**9. Warranty**

9. 1. The Licensor warrants that it has, to the best of its knowledge, the rights to license reuse of the Licensed Material. **You are solely responsible for ensuring that the material you wish to license is original to the Licensor and does not carry the copyright of another entity or third party (as credited in the published version).** If the credit line on any part of the Licensed Material indicates that it was reprinted or adapted with permission from another source, then you should seek additional permission from that source to reuse the material.

9. 2. EXCEPT FOR THE EXPRESS WARRANTY STATED HEREIN AND TO THE EXTENT PERMITTED BY APPLICABLE LAW, LICENSOR PROVIDES THE LICENSED MATERIAL "AS IS" AND MAKES NO OTHER REPRESENTATION OR WARRANTY. LICENSOR EXPRESSLY DISCLAIMS ANY LIABILITY FOR ANY CLAIM ARISING FROM OR OUT OF THE CONTENT, INCLUDING BUT NOT LIMITED TO ANY ERRORS, INACCURACIES, OMISSIONS, OR DEFECTS CONTAINED THEREIN, AND ANY IMPLIED OR EXPRESS WARRANTY AS TO MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. IN NO EVENT SHALL LICENSOR BE LIABLE TO YOU OR ANY OTHER PARTY OR ANY OTHER PERSON OR FOR ANY SPECIAL, CONSEQUENTIAL, INCIDENTAL, INDIRECT, PUNITIVE, OR EXEMPLARY DAMAGES, HOWEVER CAUSED, ARISING OUT OF OR IN CONNECTION WITH THE DOWNLOADING, VIEWING OR USE OF THE LICENSED MATERIAL REGARDLESS OF THE FORM OF ACTION, WHETHER FOR BREACH OF CONTRACT, BREACH OF WARRANTY, TORT, NEGLIGENCE, INFRINGEMENT OR OTHERWISE (INCLUDING, WITHOUT LIMITATION, DAMAGES BASED ON LOSS OF PROFITS, DATA, FILES, USE, BUSINESS OPPORTUNITY OR CLAIMS OF THIRD PARTIES), AND WHETHER OR NOT THE PARTY HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES. THIS LIMITATION APPLIES NOTWITHSTANDING ANY FAILURE OF ESSENTIAL PURPOSE OF ANY LIMITED REMEDY PROVIDED HEREIN.

**10. Termination and Cancellation**

10. 1. The License and all rights granted hereunder will continue until the end of the applicable period shown in Clause 5.1 above. Thereafter, this license will be

terminated and all rights granted hereunder will cease.

10. 2. Licensor reserves the right to terminate the License in the event that payment is not received in full or if you breach the terms of this License.

**11. General**

11. 1. The License and the rights and obligations of the parties hereto shall be construed, interpreted and determined in accordance with the laws of the Federal Republic of Germany without reference to the stipulations of the CISG (United Nations Convention on Contracts for the International Sale of Goods) or to Germany′s choice-of-law principle.

11. 2. The parties acknowledge and agree that any controversies and disputes arising out of this License shall be decided exclusively by the courts of or having jurisdiction for Heidelberg, Germany, as far as legally permissible.

11. 3. This License is solely for Licensor's and Licensee's benefit. It is not for the benefit of any other person or entity.

**ORIGINAL PAPER**

# Enter the Wild: Autistic Traits and Their Relationship to Mentalizing and Social Interaction in Everyday Life

**Tobias Schuwerk**[1] · **Larissa J. Kaltefleiter**[1] · **Jiew-Quay Au**[2] · **Axel Hoesl**[3] · **Clemens Stachl**[4]

## Abstract

Theories derived from lab-based research emphasize the importance of mentalizing for social interaction and propose a link between mentalizing, autistic traits, and social behavior. We tested these assumptions in everyday life. Via smartphone-based experience sampling and logging of smartphone usage behavior we quantified mentalizing and social interaction in our participants' natural environment. Mentalizing occurred less frequently than reasoning about actions and participants preferred to mentalize when alone. Autistic traits were negatively correlated with communication via smartphone. Yet, they were not associated with social media usage, a more indirect way of getting in touch with others. Our findings critically inform recent theories on social cognition, social behavior, and the role of autistic traits in these phenomena.

**Keywords** Autism · Experience sampling method · Mentalizing · Mobile sensing · Theory of mind

"Why is she not texting me back?" A large part of everyday social life consists of trying to answer questions like this to make sense of others' behavior. Mentalizing is a powerful cognitive tool to explain and predict behavior. It is the ability to impute mental states such as beliefs, desires or intentions to others and ourselves. Mentalizing is considered essential for social interaction.

Theories on the cognitive basis of autism spectrum conditions (hereafter 'autism') are in line with this view by suggesting a causal link between altered social cognition and

reciprocal social interaction and communication in autism (Frith 2012; Tager-Flusberg 1999). The autism spectrum is characterized by a set of autistic traits, such as problems with balanced and reciprocal social interaction, rigid behavior patterns, difficulties in adapting to change, strong attention to details, or a strong focus of attention. These autistic traits reach a clinical level in autism but are also prevalent in non-clinical samples (Baron-Cohen et al. 2001). In line with the idea that autistic traits manifest as a continuum, relatives of autistic people also show an increased—yet subclinical—level of autistic traits (Sasson et al. 2013) and a heritability of autistic traits has been documented in the general population (Hoekstra et al. 2007). Further, based on the examination of autistic traits profiles of over 6900 individuals without autism, Ruzich et al. (2015) concluded that autistic traits are continuously distributed in the general population.

To date, central pillars of theories suggesting the importance of mentalizing for everyday social interaction and a link between mentalizing, autistic traits, and actual social behavior remain under-researched. On the one hand, knowledge about mentalizing in people with and without autism stems almost exclusively from lab-based research (c.f., Atherton et al. 2018). On the other hand, social interaction outside the lab is usually assessed indirectly via interviews or questionnaires (e.g., Kreider et al. 2016). Only a handful of studies addressed the impact of social cognitive deficits of individuals with autism on their everyday social

✉ Tobias Schuwerk
  Tobias.Schuwerk@psy.lmu.de

1 Department of Psychology, Developmental Psychology, Ludwig-Maximilians-Universität München, Leopoldstr. 13, 80802 Munich, Germany

2 Department of Statistics, Computational Statistics, Ludwig-Maximilians-Universität München, Munich, Germany

3 Media Informatics Group, Ludwig-Maximilians-Universität München, Munich, Germany

4 Department of Psychology, Methods and Assessment, Ludwig-Maximilians-Universität München, Munich, Germany

🖄 Springer

life (e.g., Begeer et al. 2010; Chen et al. 2016; Frith et al. 1994). Consequently, there is a large gap between the solid empirical basis of mentalizing characteristics in the lab and knowledge about actual social interaction in everyday lives of people with and without autism. Central questions that remain unanswered are: When and how do we mentalize? Is there a relationship between autistic traits and the amount and quality of mentalizing, the amount of social interaction, and more generally the extent of exposure to the social world and social network size in everyday life?

In this study, we assessed autistic traits, social cognition in everyday life, and actual social behavior. The conceptualization of autism as a dimensional condition and the prevalence of autistic traits in the general population made it possible to address the questions above in a non-autistic sample (Landry and Chouinard 2016). Our strategy was two-fold: First, we employed the experience sampling method (ESM), a way to capture moment-to-moment cognitive processing in an everyday context (Hektner et al. 2007), to measure the amount and quality of mentalizing outside the lab. Second, we measured the amount of communication and exposure to the social world via logging of smartphone usage behavior. Both measures were then compared with the participants' level of autistic traits, controlled for Big Five personality dimensions, social anxiety, and verbal intelligence.

One other study previously used ESM to quantify the extent to which we mentalize. Bryant et al. (2013) sampled thoughts of 30 participants during a period of 10 h. The main finding was that overall, adults think more about actions (e.g. "I will send her another text tonight") than about mental states (e.g. "She probably thinks I am busy at work"). However, this pattern was context-sensitive: they thought more about actions than mental states when they were interacting, but more about mental states than actions when they were alone.

In the present study, participants answered ESM surveys over a period of 30 days via their smartphones. First, we aimed to replicate Bryant et al. (2013) findings in a larger sample over a longer sampling period. Second, we added new thought categories that are crucial to understand what mentalizing is used for in everyday life. Specifically, we were able to investigate whether the participants' mental state thoughts referred to the past, present, or future, and whether it referred more to themselves, someone else, or both. Third, derived from the notion that autism is associated with a reduced use of mentalizing (cf., Frith et al. 1994), we hypothesized a negative relationship between autistic traits and the overall amount of mental state thoughts (but, see Begeer et al. 2010).

It has been argued that people with autism find social interactions little rewarding and that they have a diminished motivation to engage with others (Chevallier et al. 2012; Kohls et al. 2012; Shultz et al. 2015). A large body of

lab-based findings supports the social motivation theory of autism (e.g., Chawarska et al. 2013; Clements et al. 2018). However, observations from outside the lab are often not consistent with this theory. For example, previous research showed that people with autism are interested in social interactions and do experience loneliness when this desire is not sufficiently satisfied (Howard et al. 2006; Locke et al. 2010). Moreover, by presenting alternative explanations of experimental findings and by drawing on autistic testimony, Jaswal and Akhtar (2018) recently challenged the key notion that autistic individuals lack social interest. Here, we tested the hypothesis that autistic traits are related to the emotional valence experienced while cognitively engaging with the social world. We expected to find a negative relationship between the level of autistic traits and the emotional valence while thinking about others. This hypothesis resulted from a conclusion based on two premises: First, according to the social motivation theory, autistic symptomatology should be linked to a reduced reward value of thinking about others. Second, observable autistic symptoms are tied to autistic traits, which are also prevalent in the general population.

In the second part of our study, we investigated links between autistic traits and actual social interaction in everyday life. To this end, we tracked our participants smartphone usage behavior. Smartphones are ubiquitous and record an abundance of our everyday life behavior. Crucially, a main purpose of smartphones is to communicate. Therefore, smartphone usage profiles can be used to study links between psychological phenomena and behavior in an ecologically valid and non-disruptive way (Miller 2012; Stachl et al. 2017).

Results from initial studies remain ambiguous about the extent to which people with autism use electronic devices to get in touch with the social world (Mazurek 2013; Mazurek et al. 2012; van Schalkwyk et al. 2017). Here, we were able to distinguish between communication (e.g., using a messaging app) and social media usage, a way to connect to the social world without the need to directly communicate. Considering that maintaining reciprocity in interaction is challenging in autism, the former might be particularly difficult for people with high autistic traits, whereas the latter might provide a low-threshold opportunity to participate in social life.

We hypothesized that an association between autistic traits and the amount of everyday life communication via smartphone should become evident in a negative relationship between autistic traits and the amount our participants used their smartphones to communicate. Further, if autistic traits are related to a reduced participation in the social world (Mazurek et al. 2012), we should find a negative relationship between autistic traits and social media app usage. Previous research reported a smaller social network size in people with autism (Kreider et al. 2016), and an association

of autistic traits in the general population with loneliness, as well as with the duration and quality of friendships (Jamil et al. 2017; Wainer et al. 2013). Thus, we hypothesized that the level of autistic traits should be negatively correlated with the number of contacts saved on the participants smartphones.

## Method

The pre-registration, material, and data of this study can be found at OSF (https://osf.io/39tvf/). We report how we determined the sample size, all data exclusions, all manipulations, and all measures in this study. For the sake of brevity, deviations from the pre-registration protocol are described in the Supplemental Material. The demographic information is not shared as it cannot be guaranteed that it is impossible to identify individual data sets.

## Participants

In total, 234 adults (51% female) between the age of 18 and 50 years of age ($M = 22.70$, $SD = 3.85$) took part in this study. They were mainly recruited via university mailing lists and campus bulletins. The participants gave informed written consent and received €25 for their participation. If they managed to complete 50 out of 60 ESM surveys, they received an extra €1 for each additionally filled survey (max. €35). They further took part in a lottery to win a smartphone or tablet worth €400. On average, the participants answered 41 surveys ($SD = 9$). The ethics committee of the Department of Psychology and Education of LMU Munich approved this study. Participants were included if they used an Android smartphone and reported no history of psychiatric or neurological condition. In the debriefing questionnaire, $n = 0$ participants reported that they were aware of a family member with autism. Notably, judging from the prevalence of autism, this is statistically unlikely. An explanation could be that our participants who had a family member with autism did not know about the diagnosis. Further, there are still comparably few clinical institutions that offer evidence-based autism diagnostics in Germany. Many people receive their diagnosis late in life or remain undiagnosed. German native speakers or people with equivalent language skills were able to participate in the study. Forty-three additional adults signed up for the study but had to be excluded because they did not show up for the post-sampling lab appointment ($n = 14$), they had technical problems with the application on their smartphone ($n = 18$), the data was lost irrecoverably (e.g., the smartphone broke, $n = 5$), they neither filled enough ESM surveys nor enough smartphone

usage data was sampled ($n = 6$, criteria below). Data collection started in August 2016 and ended in August 2017.

The participants (74% were currently enrolled students) stem from various fields of studies or occupation (40% social/medical, 25% mathematics/physics/engineering, 7% humanities, 3% law, 12% business/economics, 0.43% arts, 1% multiple subjects/occupations, 12% other). A total of 64% held a secondary degree, 34% held a postsecondary degree, and 2% had other degrees. A list of the participants' smartphone types and Android versions can be found at the OSF. Most of our participants were in their 20s. Only ten participants were 30 years or older. The pattern of results remained stable when running the analyses without participants from the upper age tail (participants older than 3 SD's from the mean age). Thus, we kept these older participants included in the final sample.

The sample size was determined based on an a priori power analysis. For a weak correlation ($r = 0.2$), with $\alpha$ (two-tailed) set to 0.05 and $(1 - \beta)$ set to 0.8, a minimum of 193 participants was required. For the analysis of the ESM surveys, and of the smartphone usage behavior analysis, we ended up with two different—yet largely overlapping—subsamples ($n = 220$ for the ESM analysis and $n = 223$ for the smartphone usage data analysis, 209 participants were in both subsamples). In some cases, we received data for one, but not the other analysis (e.g., if a participant did not fill enough ESM surveys, but sufficient smartphone usage data was collected). The analyses of the ESM data and the smartphone usage data were run with the respective subsample.

## Measures and Analysis

### Autistic Traits Questionnaires

We assessed the level of autistic traits via the three most commonly used and validated self-report questionnaires for adults. These questionnaires sensitively assess the prevalence of autistic traits in the general population, each one tapping into slightly different aspects of autistic personality traits. For the analyses in this study, individual scores in these three questionnaires were combined in a single compound score of autistic traits (mean of z-transformed scores of each questionnaire). All questionnaires (including the control questionnaires) were filled via PCs in the lab.

#### Autism-Spectrum Quotient

The *Autism-Spectrum Quotient* (AQ; Baron-Cohen et al. 2001) is a 50-item self-report questionnaire that measures the level of autism-associated traits in the five subscales *social skills*, *attention switching*, *attention to detail*, *communication*, *imagination*. The sum score ranges between 0

and 50 (the higher the score, the more autistic traits were reported). In a meta-analysis, Ruzich et al. (2015) showed that AQ scores are continuously distributed in the general population. In a typical nonclinical sample, the mean score is approximately 17 (*SD* range 0.8–9.7). For this study, we used a German adaption (Freitag et al. 2007).

### Empathy Quotient

The *Empathy Quotient* (EQ; Baron-Cohen and Wheelwright 2004) assesses cognitive and affective aspects of empathic traits with 40 items. A high EQ score (range 0–80) indicates a high level of empathy. Previous research showed that individuals with autism score significantly lower in the EQ than individuals without autism (Baron-Cohen and Wheelwright 2004). Baron-Cohen and Wheelright reported a mean EQ score of 42.1 (*SD* = 10.6) in a general population sample. On average, women score higher than men. We employed the German translation retrieved from http://www.autismresearchcentre.com/arc_tests. For the calculation of the compound score we used reverse scoring of the z-transformed EQ scores.

### Broader Autism Phenotype

The broader autism phenotype questionnaire (BAP; Hurley et al. 2007) measures a set of personality traits and language characteristics that are qualitatively similar to core symptoms of autism. It was initially developed to assess the prevalence of these characteristics in families of people with autism. The BAP consists of 36 items and the three subscales *aloof* (lack of interest/joy in social interactions), *rigid* (change aversion) and *pragmatic* (communication difficulties due to deviations in social aspects of language use). A mean score is calculated for each subscale and over all items (the higher the score, the more autistic traits were reported). In the study by Hurley and colleagues, the general population sample had a mean total score of 2.74 (*SD* = 0.55). The German version created for this study can be found at the OSF.

### Control Questionnaires

To ensure that possible effects can be attributed to the variation in autistic traits, and not to other potentially confounding factors, we assessed several control measures.

### Social Interaction Anxiety and Social Phobia

Social anxiety and social phobia (SPS) are highly prevalent comorbidities of autism (MacNeil et al. 2009). Further, these are also strongly related phenomena in the general population (Liew et al. 2015). Yet, a recent study also reported differential effects of social anxiety and autistic traits on social

attention, suggesting that these phenomena might be—at least partly—distinct (Kleberg et al. 2017). In this study, we included the *Social Interaction Anxiety Scale* and the SPS *Scale* (SIAS and SPS respectively; Mattick and Clarke 1998; German version by Stangier et al. 1999) to identify the variance that is attributable to social interaction, anxiety, and SPS. Each scale consists of 20 items. Mattick et al. reported a mean SPS score of 14.1 (*SD* = 10.2), and a mean SIAS score of 19.0 (*SD* = 10.1) in an undergraduate sample (N = 482).

### Verbal Intelligence

We employed a German multiple-choice vocabulary test as a rough estimate of verbal intelligence [*Mehrfachwahl-Wortschatz-Intelligenztest*, MWT-B; Lehrl (2005)]. The aim was to control for a potential influence of verbal intelligence on performance in our measures of interest (ESM and smartphone usage data), which are both inherently language-dependent.

### Big Five Personality

The German version of the Big Five Structure Inventory was employed to obtain Big Five personality scores (BFSI; Arendasy 2009). We used the person parameter of the partial credit model (PCM; see Masters 1982). The self-report questionnaire consists of 300 items. The participants are asked to evaluate how typically/untypically an adjective or a short phrase describes how they are. The response is provided using a four-point Likert scale ranging from *untypical for me* to *typical for me*. The Big Five personality dimensions (*Openness to Experience, Conscientiousness, Extraversion, Emotional Stability*/Absence of *Neuroticism*, and *Agreeableness*) are measured on the factor- and the facet-level.

### Debriefing Questionnaire

A short debriefing questionnaire, completed by the participants at the end of the study, assessed (1) the pleasantness of study participation, (2) how difficult it was to identify the respective thoughts for the ESM surveys, (3) whether the participant's daily life during the study was typical or not, (4) if, and if so how, the study had an influence on the way they used their smartphone, and (5) how many hours a day they usually interact with others (face-to-face and via technical devices).

### Experience Sampling Method

We integrated an ESM extension into an already existing version of the *PhoneStudy* Android logging application (made available for Android 4.0 or higher; see also Stachl

et al. 2017). The participants completed 60 surveys in 30 days. The timing of the surveys was pseudo-randomized and unpredictable for the participants. The participants were instructed that, on average, they will receive 2 (0–4) surveys per day, and that the surveys will only be scheduled between 10 am and 8 pm. A status screen, accessible via the navigation drawer, informed the participants how many surveys they already completed, and how many surveys they will receive this day. Participants who completed less than 33% of the ESM surveys (20 out of 60), were excluded from the analysis. The participants who were excluded due to this criterion ($n = 14$) did not differ in their level of autistic traits (compound score) from the final sample, $t(14.35) = -0.85$, $p = .411$, $g = -0.26$, $CI_{95\%} = [-0.8, 0.28]$.

The current ESM measure was closely adapted from a study by Bryant et al. (2013). All 60 surveys were identical and consisted of five multiple-choice questions in a fixed sequence. The first question referred to the type of thought: "What were you thinking of just before the beep?" (response options: mental state/action/miscellaneous/I cannot tell exactly right now). The second question asked about the direction of the thought: "Who was involved in this thought?" (response options: I/someone else/I and someone else/miscellaneous/I cannot tell exactly right now). The third question addressed the time reference of the thought: "What was the timeline of the thought?" (response options: past/present/future/none of these options). The fourth question referred to the participant's mood while thinking this thought: "How did you feel while having this thought?" (response options: pleasant/neutral/unpleasant/I cannot tell exactly right now). The fifth question asked whether participants were interacting while having the thought: "Were you engaged with others while having this thought?" (response options: yes/no).

The ESM surveys popped up as visual notifications on the lock screen, accompanied by a beep and a haptic feedback (vibration). To answer the survey, participants had to touch the notification. Once opened, they had 10 min to fill the survey, after that the notification disappeared and the survey was counted as missed. Participants were instructed to answer as many surveys as promptly as possible, without putting themselves in danger by doing so (e.g., if they were currently driving). At the beginning of the study, the participants completed a standardized instruction and training, implemented in the *PhoneStudy* app (for details see material at OSF). In a standardized step-by-step procedure, the application instructed the participants on how to adequately respond to the ESM prompts. For example, for the first question on the type of thought, it was crucial to explain the meaning of the terms *mental state* and *action*. The participants were instructed that mental states only exist in their own or another person's head. Examples for mental states are opinions, beliefs, desires, or feelings. An *action* was defined

as something that they or others are doing. All definitions were accompanied by examples (e.g., I think Sarah is still at work, I will brush my teeth before I go to bed). The other questions were explained accordingly (see OSF for details). A potential disadvantage of fixed response categories as compared to free text responses could be a wrong or imprecise categorization of the thought of interest. Yet, comparing both response formats, Bryant et al. (2013) found the same pattern of results. Based on cost-effectiveness considerations and the difficulty to unambiguously categorize free text, we decided to use multiple-choice responses.

Following the instruction, the participants completed a training session (referred to as "quiz" in the app). It consisted of 36 example thoughts that had to be categorized correctly (4 question types × 9 example thoughts). For example, the thought "I want to eat chocolate although I shouldn't" had to be categorized correctly as mental state that refers to the participant him- or herself and to the present. For the question addressing the participant's mood, any option was counted as correct. The training session was only passed if all questions were answered correctly. Incorrectly answered questions were repeated until the correct response was provided. Throughout the whole test period, the instruction and the training were available via the navigation drawer.

At the end of the study, participants provided feedback about the ESM methodology in a short debriefing questionnaire. In the current sample, 17% rated the ESM procedure as pleasant, 73% as neither pleasant nor unpleasant, and 10% as unpleasant. The debriefing questionnaire showed that participants were sufficiently able to identify a respective thought (7% always, 73% most of time, and 20% half of the time). Note that the participants were instructed to select the option "I cannot tell exactly right now" in situations in which they were not able to unambiguously identify a respective thought.

## Social Interaction Via Smartphone

Smartphone usage behavior was automatically recorded via the *PhoneStudy* Android mobile sensing application (Stachl et al. 2017). The app uses background services to monitor a wide range of smartphone usage behavior, such as app usage, communication (calls, SMSs), mobility assessed via geolocation, listened music tracks, Bluetooth/Wifi connections, battery-charging events, and boot events. For the planned analyses of the current study, we focussed on the following variables as indicators of social interaction via smartphone: number and duration of incoming and outgoing calls, number and total length of received and sent SMSs, and number and duration of events in which participants used apps for social interaction (e.g., WhatsApp, Facebook, Twitter, etc.). Further, the number of contacts at the end of the logging period was recorded as an indicator of social

network size. The *PhoneStudy* app neither tracks the content of written text nor does it record spoken words. Contacts are hashed. In a first anonymization step, we assured that personal information and logged data were never jointly stored. After the second anonymization step, neither the experimenters nor the participants were able to link personal information to a data set. Because the collected raw data is still sensitive (e.g., via geolocation in combination with the usage of certain apps), the possibility that a person could be identified cannot be excluded. Therefore, we saved this data inaccessible to the public, adhering to data storage guidelines of the local university.

The smartphone usage events were logged as a list of timestamp-sorted actions. Each event was a row that contained information about the time of the event (e.g., "1488966198449"), geolocation (e.g., "48.156024, 11.582928″), application name (e.g., "WhatsApp"), and package name (e.g., "com.whatsapp"). The service assessed the currently running app every 2 s, creating a log entry if it had changed. Devices operating on newer versions of the Android operating system supported reading the app usage history directly. On capable devices, our app thus automatically switched to this method, retrieving the latest history every 15 min. The participants were instructed to regularly transfer the collected data to our server, using SSL encryption. Additionally, the final database was automatically transferred to the server once the logging period ended.

In a first processing step, we filtered out events that did not reflect usage behavior. These events were produced by apps that run in the background and are not voluntarily controlled by the participant (e.g., the launch and functioning of a manufacturer-specific keyboard). Those background apps vary between manufacturer types and Android versions. A list of all filtered background apps that were at work in the current sample can be found at the OSF. Subsequently, we identified and categorized usage events of apps for social interaction. Due to the multitude of relevant apps and because some apps could not be unambiguously categorized whether they are used for social interaction or not, we had to individually decide in which category an app fitted best. A source for these decisions were descriptions of the applications' purpose that are available at the Google Play Store.

For our analyses, we formed two categories which served as dependent variables (a list of apps per category can be found at the OSF). The first category, termed *communication*, subsumed events of apps with the main purpose to communicate with others verbally or via text messages. These events were generated by pre-installed apps for phone calls and messaging, as well as by apps from other providers (e.g., WhatsApp, Signal, or Skype). For this analysis, we made no distinction between verbal communication and text messaging, because many of these apps offer both communication forms and this could not be differentiated in the logged

event. We did not consider e-mail apps for this category. First, because a substantial amount of e-mail traffic is related to contacting companies or agencies (e.g., for online shopping). Second, because the amount of work-related e-mails, a rather involuntary form of communication, could not be identified for filtering them out.

The second category, termed *social media usage*, grouped events of apps that connected the participants to the social world without the need to directly communicate. Although messaging can be a feature of these apps, the main reason to use these apps is not communication. Apps for classical social networks such as Facebook or Instagram are in this category. An important reason to use such an app is to address one's need to belong and/or one's need to self-represent (Nadkarni and Hofmann 2012). Further, browsing one's timeline can merely be used to gather news on individually-relevant topics. Another type of apps in this category is used to coordinate group tasks (e.g., shared calendars, apps that help to share costs between several people, or apps that can be used to manage a sports team). Dating apps were also included in this category. Although communication takes place in dating apps, their main purpose is to look at other people's profiles in order to find a matching person.

In the next processing step, the total number of events per app and category was calculated. Further, the total usage duration of apps of the two categories was calculated. This was done by computing the difference between the timestamp of an event of interest (e.g., the first occurence of a "WhatsApp" usage event) and the timestamp of the next event generated by the usage of a different app or operation (e.g., turning the screen off). Ten participants had to be excluded because usage data was missing for more than 3 days of their logging period. For nine participants, logging data was missing for less than 3 days. For these participants, we interpolated the number and duration of usage per app (via the rule of three, in total 0.17% of the data) to match the logging period of exactly 30 days. This criterion was set during data preprocessing, prior to data analysis.

Due to a logging issue, a systematic error was introduced to the number and duration of app usage events. In some situations, it was not logged when a participant turned off her screen, which led to implausibly long app usage events. For example, if a participant used Whatsapp before she went to bed and the event of turning off the screen was not logged, the whole time until the next event in the morning (e.g., alarm clock) was incorrectly counted as duration of WhatsApp usage. As the occurence of this logging error was related to the amount our participants used their smartphone, a simple exclusion of these events would have biased our data set. To solve this issue, we identified these events in the raw data and replaced them with the participant's mean usage duration of this app. The number of logging error events was added to the recorded total number of usage

events per app. Thus, the total number of app usage events could be accurately reconstructed. For the total duration of communication events, 9.07% of the data was interpolated. For the variable total duration of social media usage, 2.33% of the data was interpolated. Aggregated data before and after this correction is available at the OSF.

All data processing and analyses were performed with statistical software R 3.5.0 (R Core Team 2018). A full list of employed packages can be found at OSF.

## Procedure

The study was comprised of three parts. First, participants were invited to a pre-sampling lab appointment (based on the participant's schedule, those were individual or group sessions). In the morning of the same day, they received instructions via mail on how to install the app. At the beginning of the lab appointment, the experimenter made sure that everyone had successfully installed the app and provided help if necessary. Subsequently, participants completed the standardized ESM instruction and training. The experimenters answered any upcoming questions. After that, the participants completed the verbal intelligence questionnaire and the BFSI on a PC. Note that half of the participants filled the BFSI at the post-sampling lab appointment. Further, all participants additionally completed the BFSI on their smartphone either at the beginning or at the end of the 30 days. This data was used for an independent study: https://osf.io/h9pdb. The ESM period started one day after the first lab appointment. During the following 30 days, which constituted the second part of the study, the participants received

the 60 ESM surveys. During the same time, their smartphone usage behavior was recorded. For the third part, the participants were invited to a post-sampling lab appointment, in which they filled the autistic trait questionnaires, the social interaction anxiety and SPS questionnaires on a PC. Additionally, they completed the debriefing protocol (a paper-and-pencil version). Finally, they received their reimbursement, based on the amount of filled ESM surveys.

## Results

All confirmatory partial correlations on the relationship between the level of autistic traits and the other measures of interest (ESM surveys and smartphone usage behavior) were corrected for multiple comparisons using the Holm-Bonferroni adjustment. For all computed $t$ tests, Hedges $g$ was used as a measure of effect size.

### Autistic Traits, Control Measures and Debriefing

Table 1 provides descriptive statistics and reliability analysis parameters for the questionnaire results. In short, the means and standard deviations of the current sample are highly comparable to those reported for the general population in previous literature. A more detailed description of the distribution of the autistic traits questionnaire results is provided in the Supplemental Material (Fig. S1). Further, the Supplemental Material provides a correlation matrix of the questionnaire results (Fig. S2). The pattern of correlations speaks for the construct validity of the employed measures. In the debriefing questionnaire, the participants indicated that they

**Table 1** Descriptive statistics of questionnaire results

|  | M | SD | Range | Skew | Kurtosis | Cronbach's alpha | CI Cronbach's alpha |
|---|---|---|---|---|---|---|---|
| AQ | 16.27 | 5.93 | 28.00 | 0.52 | − 0.44 | 0.80 | [0.77, 0.84] |
| EQ | 40.68 | 10.99 | 60.00 | − 0.04 | − 0.32 | 0.79 | [0.76, 0.83] |
| BAP | 2.74 | 0.57 | 3.08 | 0.35 | − 0.35 | 0.90 | [0.88, 0.92] |
| SPS | 15.27 | 12.63 | 67.00 | 1.27 | 1.39 | 0.92 | [0.91, 0.94] |
| SIAS | 23.18 | 14.01 | 72.00 | 0.70 | − 0.07 | 0.92 | [0.91, 0.94] |
| Verbal IQ | 106.97 | 10.26 | 54.00 | 0.81 | 0.02 | | |
| BFSI: O | − 0.07 | 0.71 | 4.20 | 0.13 | 0.14 | 0.94 | [0.93, 0.95] |
| BFSI: C | − 0.13 | 0.69 | 4.07 | − 0.07 | 0.18 | 0.95 | [0.95, 0.96] |
| BFSI: E | − 0.15 | 0.68 | 3.89 | 0.10 | − 0.09 | 0.96 | [0.96, 0.97] |
| BFSI: A | − 0.04 | 0.72 | 3.97 | 0.41 | 0.54 | 0.94 | [0.92, 0.95] |
| BFSI: N | − 0.07 | 0.79 | 4.61 | 0.35 | 0.37 | 0.94 | [0.93, 0.95] |

The last two columns provide results of the reliability analyses of the personality trait questionnaires

*BFSI* values reflect person parameters of the PCM (Masters 1982)

*AQ* Autism-Spectrum Quotient, *EQ* Empathy Quotient, *BAP* Broader Autism Phenotype, *SPS* Social Phobia Scale, *SIAS* Social Interaction Anxiety Scale, Verbal IQ refers to the MWT-B, a German multiple choice vocabulary test, *BFSI* Big Five Structure Inventory, *O* openness to experience, *C* conscientiousness, *E* extraversion, *A* agreeableness, *N* emotional stability/absence of neuroticism

usually interact with others for about 7.05 h per day (face-to-face and via technical devices; $SD = 3.41$ h, range 1–16 h). Further, 68% of the participants indicated that their daily routine during the sampling period was typical ("as usual"), 18% stated their daily routine was untypical ("I did things I usually don't do"), and 14% could not decide whether their daily routine was typical or untypical. In total, 60% of the participants reported that the study had no influence on their smartphone usage behavior. Of the 40% who indicated an influence, 7% stated that they used their smartphone more often, 2% said they were more aware of their usage behavior. 16% looked more often on the phone, 7% took the phone more often with them, and only 1% stated that the study had some influence on their actual smartphone usage behavior (the remaining 7% provided no information on the nature of the specific influence).

## Experience Sampling

The ESM survey analysis is based on a sample of 220 participants. Descriptive statistics of the questionnaire results of this subsample can be found in the Supplemental Material.

### Confirmatory Analyses

We replicated the finding by Bryant et al. (2013) that participants think more about actions ($M_{action} = 0.56$, $SD_{action} = 0.18$) than about mental states ($M_{mental} = 0.28$, $SD_{mental} = 0.18$) in their everyday life, $t(219) = -12.92$, $p < .001$, $g = -0.87$, $CI_{95\%} = [-1.07, -0.67]$.

Further, we investigated whether the frequency of thoughts about mental states and actions was context-dependent. To this end, we calculated thought types (mental state, action, miscallaneous) relative to the context in which they occurred (interaction and alone) and performed a $2 \times 2$ repeated measures analysis of variance (ANOVA) with the within-participants factors thought type (mental state vs. action) and context (interaction vs. alone). See Fig. 1 for boxplots. Mirroring the finding of the $t$ test reported above, we found a significant main effect of thought type, $F(1,219) = 171.57$, $MSE = 0.11$, $p < .001$, $\hat{\eta}_G^2 = .349$. Due to the way frequency scores were calculated for this analysis (thought type relative to context), no main effect of context was observed, $F(1,219) = 1.41$, $MSE = 0.00$, $p = .236$, $\hat{\eta}_G^2 = .000$. Crucially, we found a significant interaction between thought type and context, $F(1,219) = 14.90$, $MSE = 0.03$, $p < .001$, $\hat{\eta}_G^2 = .011$. Bonferroni-corrected post hoc $t$ tests showed significant differences between all conditions. Action thoughts occured more frequently when the participants were interacting ($M = 0.59$, $SD = 0.22$) than when they were alone ($M = 0.54$, $SD = 0.19$), $t(219) = 3.75$, $p = .001$, $g = 0.25$, $CI_{95\%} = [0.06, 0.44]$. Conversely, mental state thoughts occured more often when the participants were alone ($M = 0.29$, $SD = 0.19$) than when they were interacting ($M = 0.26$, $SD = 0.20$), $t(219) = -3.39$, $p = .005$, $g = -0.23$, $CI_{95\%} = [-0.42, -0.04]$. Further, the post hoc $t$ tests showed that participants more frequently thought about actions than mental states when they were interacting $t(219) = -12.87$, $p < .001$, $g = -0.87$, $CI_{95\%} = [-1.06, -0.67]$. In parallel, when alone, participants also thought more frequently about



**Fig. 1** Mean frequency of thought type. This figure illustrates the mean frequency of thoughts about actions and mental states in percent

actions than about mental states $t(219) = -10.58$, $p < .001$, $g = -0.71$, $CI_{95\%} = [-0.91, -0.52]$.

Additionally, we addressed the question whether the participants' mental state thoughts referred more frequently to the past, present, or future in a one-way repeated measures ANOVA with the within-factor timeline (past, present, future). The respective boxplots are shown in Fig. 2a. This analysis revealed a significant difference between the times to which the participants' thoughts referred, $F(1.65, 360.54) = 241.10$, $MSE = 0.07$, $p < .001$, $\hat{\eta}_G^2 = .507$. Bonferroni-corrected post hoc $t$ tests showed that the participants' mental state thoughts referred more frequently to the present ($M = 0.59$, $SD = 0.25$) than to the past ($M = 0.12$, $SD = 0.15$), $t(219) = -20.32$, $p = .030$, $g = -1.37$, $CI_{95\%} = [-1.58, -1.16]$ and the future ($M = 0.20$, $SD = 0.20$), $t(219) = 14.45$, $p < .001$, $g = 0.97$, $CI_{95\%} = [0.77, 1.17]$. Further, their mental state thoughts more often referred to the future than to the past $t(219) = 14.45$, $p < .001$, $g = 0.97$, $CI_{95\%} = [0.77, 1.17]$.

We also analyzed whether the participants' mental state thoughts more frequently referred to themselves, others, or themselves and others. Boxplots can be found in Fig. 2b. A one-way repeated measures ANOVA with the factor direction (self, other, self and other) yielded a significant difference between the directions of mental state thoughts, $F(1.9, 416.76) = 22.46$, $MSE = 0.08$, $p < .001$, $\hat{\eta}_G^2 = .088$. Bonferroni corrected post hoc $t$ tests indicated that mental state thoughts referred more frequently to oneself ($M = 0.40$, $SD = 0.27$) than to others ($M = 0.26$, $SD = 0.22$), $t(219) = 4.76$, $p < .001$, $g = 0.32$, $CI_{95\%} = [0.13, 0.51]$, and to oneself and others ($M = 0.23$, $SD = 0.21$), $t(219) = 6.14$, $p < .001$, $g = 0.41$, $CI_{95\%} = [0.22, 0.6]$. There was no difference in the frequency of mental state thoughts referring to others versus oneself and others, $t(219) = 1.32$, $p = .566$, $g = 0.09$, $CI_{95\%} = [-0.1, 0.28]$.

Finally, we investigated whether our data would indicate an association of the level of autistic traits with the reported amount of mental state thoughts. The corresponding partial correlation was controlled for verbal IQ (MWT-B), SPS, social anxiety (SIAS), and Big Five personality dimensions (BFSI). We found no significant relation between the level of autistic traits and the amount of mental state thoughts in this analysis, $r = 0.02$, $p > .999$, $CI_{95\%} = [-0.12, 0.15]$, ($p_{uncorrected} = .786$).

To analyze the relationship between autistic traits and the emotional valence while cognitively engaging with the social world, we computed the mean valence of all thoughts that were (1) categorized as mental state or action and (2) that were directed to others (i.e. the categories "other" and "self and other"). The logged valence was coded as $-1$ (negative), 0 (neutral), or 1 (positive). The partial correlation between the level of autistic traits and the valence of thoughts that addressed the social world ($M = 0.21$, $SD = 0.28$; controlling for the same variables as above) revealed no significant relationship between these two variables, $r = 0.01$, $p = > .999$, $CI_{95\%} = [-0.12, 0.14]$, ($p_{uncorrected} = .884$).



**Fig. 2** **a** Mean percentage of the timeline of the thought. **b** Mean percentage of the direction of mental state thoughts

### Exploratory Analyses

It was previously described that people with autism use more conscious and explicit routes to reason about others' mental states in contrast to the comparably effortless mentalizing of people without autism (Hill and Frith 2003). This leads to the assumption that especially during social interaction, a situation which is challenging for many people with autism, they should be explicitly reasoning about mental states (cf., Begeer et al. 2010). Thus, people with autism might be more aware of their mental state reasoning and might use such an explicit form of mentalizing more frequently than people without autism. With our data, we can indirectly test this assumption by investigating whether higher autistic traits are associated with an increased frequency of mental state thoughts when our participants were interacting with others. However, we found no evidence for such a relationship in a partial correlation between the level of autistic traits and the amount of mental state thoughts during social interaction, while controlling for the influence of verbal IQ (MWT-B), SPS, social anxiety (SIAS), and Big Five personality dimensions (BFSI), $r = -0.01$, $p = .937$, $CI_{95\%}$ [−0.14, 0.13]. Note that the $p$ value of this exploratory analysis is uncorrected and should not be interpreted.

### Social Interaction Via Smartphone

The analysis of social interaction via smartphone is based on a subsample of 223 participants. Descriptive statistics of the questionnaire results of this subsample are provided in the Supplemental Material. Table 2 gives an overview of the descriptive statistics of the logged smartphone usage behavior that served as measures of interest. On average, the participants used communication apps for 24 h in the 30-day-long logging period ($SD = 17$ h). This corresponds to a mean of 48 min per day ($SD = 34$ min). Social media apps were used, on average, for 15 h in the sampling period ($SD = 14$ h). This equals a mean social media duration of 29 min per day ($SD = 28$ min). On average, our participants had 189 contacts saved on their smartphone ($SD = 138$

**Table 2** Descriptive statistics of smartphone usage behavior

|  | M | SD | Range |
| --- | --- | --- | --- |
| Total number of communication events | 1981 | 1470 | 8609 |
| Total duration of communication events (in h) | 24.08 | 17.16 | 94.72 |
| Total number of social media events | 641 | 770 | 5276 |
| Total duration of social media events (in h) | 14.7 | 14.18 | 71.23 |
| Number of contacts | 189 | 138 | 1039 |

The total number of events and the total sum of event durations in the 30-day-long logging period is shown. The number of contacts was recorded at the end of the sampling period

contacts). These app usage rates reflect the previously reported so-called *application micro-usage* behavior (Ferreira et al. 2014). Our participants spent, on average, 48 s at a time using an app from the *communication* category ($SD = 26$ s). The average usage duration of apps from the *social media* category was 91 s at a time ($SD = 74$ s).

### Confirmatory Analyses

All partial correlations were again controlled for verbal IQ (MWT-B), SPS, social anxiety (SIAS), and Big Five personality dimensions (BFSI). Scatterplots displaying the relationship between the level of autistic traits and the the amount of communication via smartphone can be found in Fig. 3. A main aim of our study was to test whether the participants' level of autistic traits was associated with their amount of communication via smartphone. After correcting for multiple comparisons, we found a significant negative correlation between the level of autistic traits and the total number of communication events, $r = -0.18$, $p = .048$, $CI_{95\%} = [-0.31, -0.05]$, ($p_{uncorrected} = .007$). The negative correlation between the level of autistic traits and the total duration of communication events was not significant after the Holm-Bonferroni adjustment, $r = -0.16$, $p = .111$, $CI_{95\%} = [-0.29, -0.03]$, ($p_{uncorrected} = .019$).

We found no significant correlation between the level of autistic traits and exposure to the social world, operationalized via the total number of social media events, $r = -0.04$, $p > .999$, $CI_{95\%} = [-0.18, 0.09]$, ($p_{uncorrected} = .516$). Also the correlation between the level of autistic traits and the total duration of social media events was not significant, $r = -0.05$, $p > .999$, $CI_{95\%} = [-0.18, 0.09]$, ($p_{uncorrected} = .483$).

There was also no significant correlation between the level of autistic traits and the number of contacts saved on the participants' smartphones, $r = -0.04$, $p > .999$, $CI_{95\%} = [-0.17, 0.10]$, ($p_{uncorrected} = .583$).

### Exploratory Analyses

We ran a regression analysis to further explore the significant correlation between the level of autistic traits and the number of communication events. We were interested in the specific influence of the level of autistic traits on communication via smartphone. Previous literature suggested that social anxiety, SPS, and autistic traits are strongly related, but constitute still distinct phenomena (Kleberg et al. 2017; Liew et al. 2015), To better assess the differential contributions of each domain, we introduced social anxiety as well as the interaction between social anxiety and autistic traits as additional predictors into the model. The dimension extraversion from the Big Five personality inventory was added as a control variable.

**Fig. 3** Scatterplots showing the relationship between level of autistic traits and communication via smartphone. Correlation coefficients are from the partial correlation of the measures of interest, controlled for verbal IQ (MWT-B), social phobia (SPS), social anxiety (SIAS), and Big Five personality dimensions (BFSI). $p < .05*$ after correcting for multiple comparisons



For the confirmatory analyses, we used the level of autistic traits, a compound score of the participants' AQ, EQ, and BAP scores. However, a reliability analysis of these three z-transformed scores revealed that EQ scores were not a good predictor of AQ and BAP scores, implying that the EQ measured a different construct than the AQ and the BAP. With the EQ included, Cronbach's $\alpha$ was 0.78. When the EQ was left out, Cronbach's $\alpha$ increased to 0.86 (whe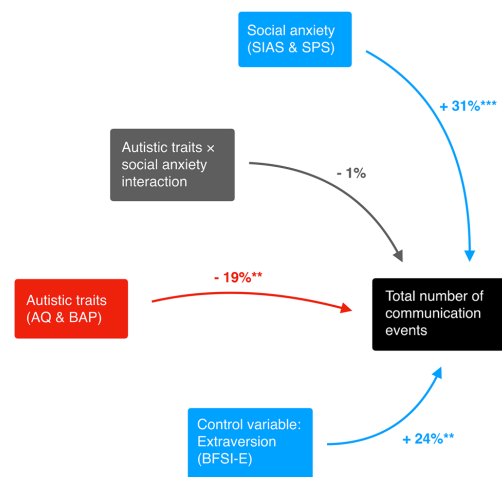n the AQ was dropped, Cronbach's $\alpha$ was 0.61, when the BAP was dropped, Cronbach's $\alpha$ was 0.59). Further, also the EQ's discriminatory power was the lowest of the three measures ($r_{EQ} = 0.46$, $r_{AQ} = 0.69$, $r_{BAP} = 0.71$). Based on these results, we excluded the EQ from this exploratory analysis and built a compound score only from z-transformed AQ and BAP scores to get a better estimate of the level of autistic traits.

A reliability analysis of the employed measures for social anxiety and social phobia (SPS and SIAS) revealed a Cronbach's $\alpha$ of 0.87 and a sufficient discriminatory power, $r = 0.77$. This fits well with the conceptualization of the SIAS and SPS as complementary measures of the same underlying construct (Mattick and Clarke 1998). Thus, for this exploratory analysis, both measures were combined into one score of social anxiety.

The distributions of the independent variables indicated that a negative binomial regression model is appropriate. Figure 4 illustrates the model and provides the percent ratio of the Incident Rate Ratio $[-100 \times (1 - Exp(b))]$. The level of autistic traits significantly predicted the total number of communication events, $b = -0.21$, $SE = 0.07$, $Z = -2.88$, $p = .004$. Holding the other predictors constant, an increase



**Fig. 4** Schematic illustration of the exploratory negative binomial regression of smartphone usage data. The values show the percent ratio of the Incident Rate Ratio $[-100 \times (1 - Exp(b))]$. Positive values indicate a positive, negative values a negative predictive relationship between the independent variables and the total number of communication events ($p < .001***$, $p < .01**$). Note that the $p$ values of this exploratory analysis are not corrected for multiple comparisons and the predictive relations should not be generalized without further cross-validation

of the level of autistic traits by one standard deviation was associated with a decrease by 19% of communication via smartphone, operationalized via the total number of communication events. In contrast, social anxiety showed a significant positive relation to the total number of communication events, $b = 0.27$, $SE = 0.07$, $Z = 3.71$, $p = < .001$. Keeping all other predictors constant, an increase of social anxiety in our aggregated score by one standard deviation was associated with a 31% increase of communication via smartphone. The interaction between the level of autistic traits and social anxiety did not significantly predict the communication via smartphone, $b = -0.01$, $SE = 0.05$, $Z = -0.24$, $p = .807$. Analoguous to social anxiety, the control variable extraversion was significantly positively related to the communication via smartphone, $b = 0.22$, $SE = 0.08$, $Z = 2.64$, $p = .008$. An increase of extraversion by one standard deviation led to a 24% increase of the communication via smartphone, when keeping the other predictors constant. It is important to note that due to the exploratory nature of this analysis, the found associations should not be readily generalized without further cross-validation in a new sample.

## Discussion

We investigated the nature of mentalizing and the links between autistic traits, mentalizing, and social interaction in everyday life. Corresponding to Bryant et al. (2013) findings, adults thought twice as much about actions than about mental states. Further, we found a similar context-specific variation. Our participants reported more thoughts about actions when they were interacting with others as compared to when they were alone and vice versa. Based on the idea that this form of mentalizing is effortful and resource-consuming and that our (neuro-)cognitive system works cost-efficiently (Bullmore and Sporns 2012; Fiebich and Coltheart 2015), we argue that overall, mental state thoughts occur less frequently than action thoughts because processing of mental states is cognitively costly. Rather, they occur preferably when we are alone, a situation in which cognitive resources are not occupied by the multitude of social information that has to be processed during interaction.

In our sample, mentalizing in everyday life was mainly used to process current mental states and only to a minor fraction dealt with past and future mental states. Further, paralleling Bryant et al. (2013), we found that most mental state thoughts revolved around one's own mental state. Yet, next to self- and other-directed thoughts, we introduced a third category to classify thoughts that referred to oneself and others because sometimes this cannot be disentangled. Our findings suggest that Bryant et al. underestimated the amount of thoughts that—at least partially—refer to others.

Our results show that about half of the mental state thoughts in our sample were directed to others or others and oneself.

In contrast to what can be postulated based on previous literature (cf., Frith et al. 1994), autistic traits were not related to a reduced use of mentalizing. Moreover, we found no relationship between autistic traits and the valence of thoughts that addressed the social world.

As hypothesized, autistic traits were negatively correlated with communication via phone calls or text messages. The exploratory regression analysis points to additional details on the nature of this relationship. An increase of autistic traits was associated with a decrease in communication via smartphone. Interestingly, there was no interaction between autistic traits and social anxiety, and social anxiety had a reverse effect on the amount of communication. This adds to evidence that both phenomena are overlapping but, yet, distinct (Kleberg et al. 2017; Liew et al. 2015).

Further, it allows for speculating about a potential compensatory use of computer-mediated communication for people with increased levels of autistic traits and social anxiety. A recent meta-analysis showed that social anxiety is positively correlated with comfort during computer-mediated interaction (Prizant-Passal et al. 2016). Additionally, research with adolescents suggests that communicating via technical devices may help to compensate for weak social skills when making contact with new people (Bonetti et al. 2010) and may support interaction with peers (Desjarlais and Willoughby 2010). For autism, empirical evidence on the use of computer-mediated communication is relatively sparse. Interviewing adults with autism, Burke et al. (2010) found that computer-mediated communication could help to reduce stress associated with interactions because, in contrast to real-life interactions, there is, e.g., no need to decode nonverbal signals and conversations are more pre-defined. However, the autistic adults also reported that they find it challenging to maintain relationships online due to trust issues or insecurity in the usage of social rules.

In sum, it seems that while for people with increased social anxiety communication via smartphone could serve a compensatory purpose, this may be different in the case of people with elevated autistic traits. Further research is necessary to follow up on this result. Van Schalkwyk et al. (2017) recently presented a short questionnaire that could help in the endeavor to systematically assess the role computer-mediated communication plays in the lives of autistic people.

Autistic traits were not associated with the amount of social media usage, a more indirect way of getting in touch with the social world. We also found no relation between autistic traits and social network size (Kreider et al. 2016). This suggests an interesting dissociation between different ways of engaging with the social world. The reduced communication could be related to difficulties with fast and

flexible social information processing, required for reciprocal social interactions. Unlike communication via smartphone, social media usage can be entirely passive and follows clear rules (e.g., liking, retweeting, …). Thus, it may be less challenging for people with difficulties in reciprocal interaction (cf., van Schalkwyk et al. 2017).

We argue that the lacking relationship between autistic traits and (1) the valence of thoughts that addressed the social world and (2) the amount of social media usage speaks against claims that can be derived from the social motivation theory. This theory holds that reduced joy in social situations, fewer friendships, or a reduced preference for collaborative activities results from a lack of interest in the social world because autistic individuals find social interactions little rewarding (Chevallier et al. 2012; Kohls et al. 2012; Shultz et al. 2015).

The empirical basis of this theory comes largely from lab-based experimental research (e.g., Chawarska et al. 2013; Clements et al. 2018). Notably, a rare example from outside the lab is also presented by Chevallier and colleagues (Chevallier et al. 2012). In their study, employing a questionnaire, adolescents with autism reported a reduced enjoyment of social situations. Further, the amount of reported pleasure in social interactions correlated with autism symptom severity. Yet, the authors did not ask the participants *why* they experienced little enjoyment during interaction.

One answer and a critical alternative explanation of such findings is that social interaction is not less rewarding or joyful in autism, but that it is more difficult and stressful for many autistic people, which often results in social withdrawal, fewer friendships, and reduced relationship quality. However, the source is not a lacking motivation of the autistic individual. A fast-growing body of evidence highlights the role non-autistic people play in successful social interaction (Morrison et al. 2019; Sasson et al. 2017). This suggests that, rather than attributing social interaction problems only to deficits within the autistic individual, we should focus on the mismatch between interaction partners to explain social impairments (Bolis et al. 2017).

Our finding that direct communication, but neither the self-reported emotional valence when cognitively engaging with others, nor the general tendency to get in touch with the social world, were associated with the level of autistic traits is in line with this point of view. Yet, it is important to stress the caveat that this interpretation rests on the assumption that due to their qualitative similarity, autistic traits in the general population can be employed to study clinically relevant autism symptomatology (Landry and Chouinard 2016). However, the conclusions based on the examination of autistic traits in the general population cannot be readily generalized to autism. For example, previous work suggests that the AQ taps the same latent traits in people with and without autism, but that the same test scores do not necessarily reflect the same level of autistic traits (Murray et al. 2014). A next step would be to run the current study in a sample of people with an autism diagnosis.

## Limitations

Some methodological factors should be considered in the evaluation of our findings. Compared to experimentally testing cognition in the lab, experience sampling introduces a considerable measurement error. It cannot be definitely determined to what extend the participants accurately classified their thoughts. In particular, the type of thought categorization inevitably left room for ambiguity. For example, it could have been difficult for the participants to distinguish between thoughts like "I will eat a cookie" (which should be categorized as an action) and "I want to eat a cookie" (which should be categorized as a mental state). To address this issue, we explicitly used thoughts like these to train participants to differentiate between mental state and action thoughts in the instruction and the subsequent quiz. Further, participants always had the possibility to use the response option "I cannot tell exactly right now" if they were not able to clearly categorize their thought. This option helped us to reduce potential noise in the data due to unclearly or ambiguously categorized thoughts in our categories of interest.

Further, interaction via smartphone constitutes only a part of our social life. Our conclusions cannot be directly expanded to other forms of interaction. However, from an experimental psychologist's point of view, given the difficulty to study cognition and behavior outside the lab, even with these limitations both measures can be considered being relatively valid means to capture these phenomena.

Due to a logging issue, we lost data for two of our four measures of social interaction via smartphone. The number of communication and social media events could be entirely reconstructed. Yet, we had to impute 9.07% of the communication event durations and 2.33% of the social media event durations, which adds a level of uncertainty to these two measures. However, it is important to note that the results of these two measures mirror the results of the respective usage event count data. Thus, we are confident to conclude that, if we introduced a bias to the data set with this interpolation, its effect is neglectable.

The categorization of smartphone apps was a source of ambiguity in this study. Deciding whether an app belonged to the *communication* or the *social media* category was difficult in some cases. We therefore specified clear and justified criteria (detailed in the Methods section; a complete list of all apps per category can be found at OSF). We always judged from the app's main usage purpose. For example, we categorized dating apps as *social media*, although the inbuilt messaging function is an important part of these

apps. Nonetheless, these apps are mainly used to find dates. A consequence of the application of this coding scheme is that some of our *social media* events include direct communication. Yet, if this would have substantially affected our data, we would have observed also a significant negative correlation between social media usage and autistic traits, just as it was the case with communication events. In turn, this strengthens the finding of the significant negative relationship between direct communication and autistic traits, a key finding of our study.

## Conclusion

Our data provide evidence that thinking about others and our own actions and mental states makes up most of our conscious cognitive processing. We were able to show that elevated autistic traits are associated with reduced computer-mediated communication, potentially because reciprocal direct interaction is difficult for people with high autistic traits. Yet, autistic traits were unrelated to the general tendency to get in touch with the social world and with the social network size, indirectly supporting findings that people with autism seek social participation via technology (Mazurek 2013).

### Compliance with Ethical Standards

**Conflict of interest** The authors declare that they have no conflict of interest.

**Ethical Approval** The ethics committee of the Department of Psychology and Education of LMU Munich approved the study.

**Informed Consent** The participants gave informed written consent to participate in the study.

## References

Arendasy, M. (2009). *BFSI: Big-Five Struktur-Inventar (Test & Manual)*. Mödling: SCHUHFRIED GmbH.

Atherton, G., Lummis, B., Day, S. X., & Cross, L. (2018). What am i thinking? Perspective-taking from the perspective of adolescents with autism. *Autism, 23*(5), 1186–1200. https://doi.org/10.1177/1362361318793409.

Baron-Cohen, S., & Wheelwright, S. (2004). The broad autism phenotype questionnaire. *Journal of Autism and Developmental Disorders, 34*(2), 163–175. https://doi.org/10.1023/B:JADD.0000022607.19833.00.

Baron-Cohen, S., Wheelwright, S., Skinner, R., Martin, J., & Clubley, E. (2001). The autism-spectrum quotient (AQ): Evidence from asperger syndrome/high-functioning autism, malesand females, scientists and mathematicians. *Journal of Autism and Developmental Disorders, 31*(1), 5–17. https://doi.org/10.1023/A:1005653411471.

Begeer, S., Malle, B. F., Nieuwland, M. S., & Keysar, B. (2010). Using theory of mind to represent and take part in social interactions: Comparing individuals with high-functioning autism and typically developing controls. *European Journal of Developmental Psychology, 7*(1), 104–122. https://doi.org/10.1080/17405620903024263.

Bolis, D., Balsters, J., Wenderoth, N., Becchio, C., & Schilbach, L. (2017). Beyond autism: Introducing the dialectical misattunement hypothesis and a bayesian account of intersubjectivity. *Psychopathology, 50,* 355–372. https://doi.org/10.1159/000484353.

Bonetti, L., Campbell, M. A., & Gilmore, L. (2010). The relationship of loneliness and social anxiety with children's and adolescents' online communication. *Cyberpsychology, Behavior, and Social Networking, 13*(3), 279–285. https://doi.org/10.1089/cyber.2009.0215.

Bryant, L., Coffey, A., Povinelli, D. J., & Pruett, J. R. (2013). Theory of Mind experience sampling in typical adults. *Consciousness and Cognition, 22*(3), 697–707. https://doi.org/10.1016/j.concog.2013.04.005.

Bullmore, E., & Sporns, O. (2012). The economy of brain network organization. *Nature Reviews Neuroscience, 13,* 336–349. https://doi.org/10.1038/nrn3214.

Burke, M., Kraut, R., & Williams, D. (2010). Social use of computer-mediated communication by adults on the autism spectrum. In Proceedings of the 2010 ACM *conference on Computer supported cooperative work*, ACM, New York, USA (pp. 425–434). https://doi.org/10.1145/1718918.1718991

Chawarska, K., Macari, S., & Shic, F. (2013). Decreased spontaneous attention to social scenes in 6-month-old infants later diagnosed with autism spectrum disorders. *Biological Psychiatry, 74*(3), 195–203. https://doi.org/10.1016/j.biopsych.2012.11.022.

Chen, Y.-W., Bundy, A., Cordier, R., Chien, Y.-L., & Einfeld, S. (2016). The experience of social participation in everyday contexts among individuals with autism spectrum disorders: An experience sampling study. *Journal of Autism and Developmental Disorders, 46*(4), 1403–1414. https://doi.org/10.1007/s10803-015-2682-4.

Chevallier, C., Grèzes, J., Molesworth, C., Berthoz, S., & Happé, F. (2012a). Brief report: Selective social anhedonia in high functioning autism. *Journal of Autism and Developmental Disorders, 42*(7), 1504–1509. https://doi.org/10.1007/s10803-011-1364-0.

Chevallier, C., Kohls, G., Troiani, V., Brodkin, E. S., & Schultz, R. T. (2012b). The social motivation theory of autism. *Trends in Cognitive Sciences, 16*(4), 231–239. https://doi.org/10.1016/j.tics.2012.02.007.

Clements, C. C., Zoltowski, A. R., Yankowitz, L. D., Yerys, B. E., Schultz, R. T., & Herrington, J. D. (2018). Evaluation of the social motivation hypothesis of autism: A systematic review and meta-analysis. *JAMA Psychiatry, 75*(8), 797–808. https://doi.org/10.1001/jamapsychiatry.2018.1100.

Desjarlais, M., & Willoughby, T. (2010). A longitudinal study of the relation between adolescent boys and girls' computer use with friends and friendship quality: Support for the social compensation or the rich-get-richer hypothesis? *Computers in Human Behavior, 26*(5), 896–905. https://doi.org/10.1016/j.chb.2010.02.004.

Ferreira, D., Goncalves, J., Kostakos, V., Barkhuus, L., & Dey, A. K. (2014). Contextual experience sampling of mobile application micro-usage. *Proceedings of the 16th International Conference on Human-Computer Interaction with Mobile Devices; Services*, (pp. 91–100). https://doi.org/10.1145/2628363.2628367

Fiebich, A., & Coltheart, M. (2015). Various ways to understand other minds: Towards a pluralistic approach to the explanation of social understanding. *Mind & Language, 30,* 235–258. https://doi.org/10.1111/mila.12079.

Freitag, C. M., Retz-Junginger, P., Retz, W., Seitz, C., Palmason, H., Meyer, J. et al. (2007). Evaluation der deutschen Version des Autismus-Spektrum-Quotienten (AQ)—die Kurzversion. *Zeitschrift Für Klinische Psychologie Und Psychotherapie,* ACM, (pp. 280–289). https://doi.org/10.1026/1616-3443.36.4.280

Frith, U. (2012). Why we need cognitive explanations of autism. *The Quarterly Journal of Experimental Psychology, 65*(11), 2073–2092. https://doi.org/10.1080/17470218.2012.697178.

Frith, U., Happé, F., & Siddons, F. (1994). Autism and theory of mind in everyday life. *Social Development, 3*(2), 108–124. https://doi.org/10.1111/j.1467-9507.1994.tb00031.x.

Hektner, J. M., Schmidt, J. A., & Csikszentmihalyi, M. (2007). *Experience sampling method: Measuring the quality of everyday life.* Thousand Oaks, CA: Sage.

Hill, E. L., & Frith, U. (2003). Understanding autism: Insights from mind and brain. *Philosophical Transactions of the Royal Society of London B: Biological Sciences, 358*(1430), 281–289. https://doi.org/10.1098/rstb.2002.1209.

Hoekstra, R. A., Bartels, M., Verweij, C. J. H., & Boomsma, D. I. (2007). Heritability of autistic traits in the general population. *Archives of Pediatrics and Adolescent Medicine, 161*(4), 372–377. https://doi.org/10.1001/archpedi.161.4.372.

Howard, B., Cohn, E., & Orsmond, G. I. (2006). Understanding and negotiating friendships: Perspectives from an adolescent with asperger syndrome. *Autism, 10*(6), 619–627. https://doi.org/10.1177/1362361306068508.

Hurley, R. S. E., Losh, M., Parlier, M., Reznick, J. S., & Piven, J. (2007). The broad autism phenotype questionnaire. *Journal of Autism and Developmental Disorders, 37*(9), 1679–1690. https://doi.org/10.1007/s10803-006-0299-3.

Jamil, R., Gragg, M. N., & DePape, A.-M. (2017). The broad autism phenotype: Implications for empathy and friendships in emerging adults. *Personality and Individual Differences, 111,* 199–204. https://doi.org/10.1016/j.paid.2017.02.020.

Jaswal, V. K., & Akhtar, N. (2018). Being vs. appearing socially uninterested: Challenging assumptions about social motivation in autism. *Behavioral and Brain Sciences.* https://doi.org/10.1017/s0140525x18001826.

Kleberg, J. L., Högström, J., Nord, M., Bölte, S., Serlachius, E., & Falck-Ytter, T. (2017). Anxiety in children and adolescents with autism spectrum disorders. *Journal of Autism and Developmental Disorders, 47*(12), 3814–3821. https://doi.org/10.1007/s10803-016-2978-z.

Kohls, G., Chevallier, C., Troiani, V., & Schultz, R. T. (2012). Social "wanting" dysfunction in autism: Neurobiological underpinnings and treatment implications. *Journal of Neurodevelopmental Disorders, 4*(1), 10. https://doi.org/10.1186/1866-1955-4-10.

Kreider, C. M., Bendixen, R. M., Young, M. E., Prudencio, S. M., McCarty, C., & Mann, W. C. (2016). Social networks and participation with others for youth with learning, attention, and autism spectrum disorders. *Canadian Journal of Occupational Therapy, 83*(1), 14–26. https://doi.org/10.1177/0008417415583107.

Landry, O., & Chouinard, P. A. (2016). Why we should study the broader autism phenotype in typically developing populations. *Journal of Cognition and Development, 17*(4), 584–595. https://doi.org/10.1080/15248372.2016.1200046.

Lehrl, S. (2005). *Mehrfachwahl-Wortschatz-Intelligenztest MWT-B.* Balingen: Spitta-Verlag.

Liew, S. M., Thevaraja, N., Hong, R. Y., & Magiati, I. (2015). Development and validation of measures of social phobia scrutiny fear and social interaction anxiety. *Journal of Autism and Developmental Disorders, 45*(3), 858–872. https://doi.org/10.1007/s10803-014-2238-z.

Locke, J., Ishijima, E. H., Kasari, C., & London, N. (2010). Loneliness, friendship quality and the social networks of adolescents with high-functioning autism in an inclusive school setting. *Journal of Research in Special Educational Needs, 10*(2), 74–81. https://doi.org/10.1111/j.1471-3802.2010.01148.x.

MacNeil, B. M., Lopes, V. A., & Minnes, P. M. (2009). Anxiety in children and adolescents with autism spectrum disorders. *Research in Autism Spectrum Disorders, 3*(1), 1–21. https://doi.org/10.1016/j.rasd.2008.06.001.

Masters, G. N. (1982). A rasch model for partial credit scoring. *Psychometrika, 47*(2), 149–174. https://doi.org/10.1007/BF02296272.

Mattick, R. P., & Clarke, J. (1998). Development and validation of measures of social phobia scrutiny fear and social interaction anxiety. *Behaviour Research and Therapy, 36*(4), 455–470. https://doi.org/10.1016/S0005-7967(97)10031-6.

Mazurek, M. O. (2013). Social media use among adults with autism spectrum disorders. *Computers in Human Behavior, 29*(4), 1709–1714. https://doi.org/10.1016/j.chb.2013.02.004.

Mazurek, M. O., Shattuck, P. T., Wagner, M., & Cooper, B. P. (2012). Prevalence and correlates of screen-based media use among youths with autism spectrum disorders. *Journal of Autism and Developmental Disorders, 42*(8), 1757–1767. https://doi.org/10.1007/s10803-011-1413-8.

Miller, G. (2012). The smartphone psychology manifesto. *Perspectives on Psychological Science, 7*(3), 221–237. https://doi.org/10.1177/1745691612441215.

Morrison, K. E., DeBrabander, K. M., Faso, D. J., & Sasson, N. J. (2019). Variability in first impressions of autistic adults made by neurotypical raters is driven more by characteristics of the rater than by characteristics of autistic adults. *Autism.* https://doi.org/10.1177/1362361318824104.

Murray, A. L., Booth, T., McKenzie, K., Kuensberg, R., & O'Donnell, M. (2014). Are autistic traits measured equivalently in individuals with and without an autism spectrum disorder? An invariance analysis of the autism spectrum quotient short form. *Journal of Autism and Developmental Disorders, 44,* 55–64. https://doi.org/10.1007/s10803-013-1851-6.

Nadkarni, A., & Hofmann, S. (2012). Why do people use facebook? *Personality and Individual Differences, 52,* 243–249.

Prizant-Passal, S., Shechner, T., & Aderka, I. M. (2016). Social anxiety and internet use—a meta-analysis: What do we know? What are we missing? *Computers in Human Behavior, 62,* 221–229. https://doi.org/10.1016/j.chb.2016.04.003.

R Core Team. (2018). *R: A language and environment for statistical computing.* Vienna, Austria: R Foundation for Statistical Computing. Retrieved September 1, 2018, from https://www.R-project.org/.

Ruzich, E., Allison, C., Smith, P., Watson, P., Auyeung, B., Ring, H., et al. (2015). Measuring autistic traits in the general population: A systematic review of the autism-spectrum quotient (AQ) in a nonclinical population sample of 6,900 typical adult males and females. *Molecular Autism, 6,* 2. https://doi.org/10.1186/2040-2392-6-2.

Sasson, N. J., Faso, D. J., Nugent, J., Lovell, S., Kennedy, D. P., & Grossman, R. B. (2017). Neurotypical peers are less willing to interact with those with autism based on thin slice judgments. *Scientific Reports, 7,* 40700. https://doi.org/10.1038/srep40700.

Sasson, N. J., Lam, K. S. L., Childress, D., Parlier, M., Daniels, J. L., & Piven, J. (2013). The broad autism phenotype questionnaire: Prevalence and diagnostic classification. *Autism Research, 6*(2), 134–143. https://doi.org/10.1002/aur.1272.

Shultz, S., Jones, W., & Klin, A. (2015). Early departures from normative processes of social engagement in infants with autism spectrum disorder. In A. Puce & B. I. Bertenthal (Eds.), *The many*

*faces of social attention: Behavioral and neural measures* (pp. 157–177). Cham: Springer International Publishing. https://doi.org/10.1007/978-3-319-21368-2_6.

Stachl, C., Hilbert, S., Au, J.-Q., Buschek, D., De Luca, A., Bischl, B., et al. (2017). Personality traits predict smartphone usage. *European Journal of Personality, 31*(6), 701–722. https://doi.org/10.1002/per.2113.

Stangier, U., Heidenreich, T., Berardi, A., Ulrike, G., & Hoyer, J. (1999). Development and validation of measures of social phobia scrutiny fear and social interaction anxiety. *Zeitschrift Für Klinische Psychologie Und Psychotherapie, 28,* 28–36. https://doi.org/10.1026//0084-5345.28.1.28.

Tager-Flusberg, H. (1999). A psychological approach to understanding the social and language impairments in autism. *International Review of Psychiatry, 11*(4), 325–334. https://doi.org/10.1080/09540269974203.

van Schalkwyk, G. I., Marin, C. E., Ortiz, M., Rolison, M., Qayyum, Z., McPartland, J. C., et al. (2017). Social media use, friendship quality, and the moderating role of anxiety in adolescents with autism spectrum disorder. *Journal of Autism and Developmental Disorders, 47*(9), 2805–2813. https://doi.org/10.1007/s10803-017-3201-6.

Wainer, A. L., Block, N., Donnellan, M. B., & Ingersoll, B. (2013). The broader autism phenotype and friendships in non-clinical dyads. *Journal of Autism and Developmental Disorders, 43*(10), 2418–2425. https://doi.org/10.1007/s10803-013-1789-8.

# PART III - Conclusion

# Concluding Remarks and Outlook

In summary, the research conducted in this thesis has yielded noteworthy advancements in the field of machine learning, specifically in the context of predicting psychological attributes through smartphone data. Essential learning tasks have been successfully implemented in Chapter 5, primarily focusing on the development of multilabel classification algorithms designed for binary target variables. A natural next step is to explore multivariate regression and multi-output prediction by incorporating these algorithms into the newer R-package `mlr3`. Additionally, applying these algorithms to predict personality traits and other attributes such as gender, or age simultaneously could offer an exciting path for further investigation.

In Chapter 4 the methods required pre-defined groups of features while data driven techniques of finding groups were not explored deeper. Future research could focus on developing innovative methods that integrate domain expertise with data-driven insights, potentially leading to more effective feature grouping strategies Furthermore, the methods introduced here, currently lack a packaged format, such as an R-package, for easy accessibility. Consideration can be given to integrating these methods into established packages like `mlr3` or `iml`. Such an implementation would provide other researchers and practitioners with straightforward access to these techniques.

In Chapter 8 personality traits could reliably be predicted from smartphone data. The predictive performance, as measured by the Pearson correlation coefficient, is in line with previous work such as Kosinski et al. (2013) for many personality factors. However, it should be noted that the factor *Agreeableness* could not be reliably predicted. This limitation could potentially stem from the lack of sufficiently informative features. During the PhoneStudy, only the event of a message or app usage was logged, without capturing information of the content of the messages. One of the primary challenges lies in extracting meaningful features from the data that also adhere to stringent privacy standards (Dwork, 2011).

The results from Chapter 9 highlight the potential of combining smartphone data with traditional self-reporting to gain deeper insights into sensation seeking behaviors. While the accuracy of predicting self-reported sensation seeking scores from smartphone logs may be limited, it indicates that mobile devices can capture behavioral clues related to this psychological trait. Chapters 10, 11, and 12 collectively show the diverse research opportunities enabled by smartphone sensor data and emphasize the significant role played by the PhoneStudy. Each chapter offers a

unique perspective, highlighting different aspects of human behavior and psychology through smartphone-derived information. The integration of physiological data, such as heart rate variability and sleep patterns, could further enrich the predictive power of smartphone-based models (Lakens, 2013).

# References

Alexander, D. L. J., Tropsha, A., and Winkler, D. A. (2015). Beware of R$^2$: Simple, unambiguous assessment of the prediction accuracy of QSAR and QSPR models. *Journal of Chemical Information and Modeling*, 55(7):1316–1322.

Apley, D. W. and Zhu, J. (2020a). Visualizing the effects of predictor variables in black box supervised learning models. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 82(4):1059–1086.

Apley, D. W. and Zhu, J. (2020b). Visualizing the effects of predictor variables in black box supervised learning models. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 82(4):1059–1086.

Au, Q., Herbinger, J., Stachl, C., Bischl, B., and Casalicchio, G. (2022). Grouped feature importance and combined features effect plot. *Data Mining and Knowledge Discovery*, 36(4):1401–1450.

Bischl, B., Lang, M., Kotthoff, L., Schiffner, J., Richter, J., Studerus, E., Casalicchio, G., and Jones, Z. M. (2016). mlr: Machine learning in r. *The Journal of Machine Learning Research*, 17(1):5938–5942.

Bischl, B., Mersmann, O., Trautmann, H., and Weihs, C. (2012). Resampling methods for meta-model validation with recommendations for evolutionary computation. *Evolutionary computation*, 20(2):249–275.

Borchani, H., Varando, G., Bielza, C., and Larrañaga, P. (2015). A survey on multi-output regression. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 5(5):216–233.

Breiman, L. (2001). Random forests. *Machine learning*, 45:5–32.

Breiman, L., Friedman, J. H., Olshen, R. A., and Stone, C. J. (2017). *Classification And Regression Trees*. Routledge.

Caruana, R., Lou, Y., Gehrke, J., Koch, P., Sturm, M., and Elhadad, N. (2015). Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '15, page 1721–1730, New York, NY, USA. Association for Computing Machinery.

Chittaranjan, G., Blom, J., and Gatica-Perez, D. (2013). Mining large-scale smartphone data for personality studies. *Personal and Ubiquitous Computing*, 17:433–450.

de Montjoye, Y.-A., Quoidbach, J., Robic, F., and Pentland, A. (2013). Predicting personality using novel mobile phone-based metrics. In *Social Computing, Behavioral-Cultural Modeling and Prediction: 6th International Conference, SBP 2013, Washington, DC, USA, April 2-5, 2013. Proceedings 6*, pages 48–55. Springer.

Doshi-Velez, F. and Kim, B. (2017). Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*.

Dwork, C. (2011). A firm foundation for private data analysis. *Communications of the ACM*, 54(1):86–95.

Fisher, A., Rudin, C., and Dominici, F. (2019). All models are wrong, but many are useful: Learning a variable's importance by studying an entire class of prediction models simultaneously. *Journal of Machine Learning Research*, 20(177):1–81.

Friedman, J. H. (2001a). Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232.

Friedman, J. H. (2001b). Greedy function approximation: A gradient boosting machine. *The Annals of Statistics*, 29(5).

Fuller, D., Shareck, M., and Stanley, K. (2017). Ethical implications of location and accelerometer measurement in health research studies with mobile sensing devices. *Social Science & Medicine*, 191:84–88.

Goldberg, L. R., Sweeney, D., Merenda, P. F., and Hughes Jr, J. E. (1998). Demographic variables and personality: The effects of gender, age, education, and ethnic/racial status on self-descriptions of personality attributes. *Personality and Individual differences*, 24(3):393–403.

Goldstein, A., Kapelner, A., Bleich, J., and Pitkin, E. (2015). Peeking inside the black box: Visualizing statistical learning with plots of individual conditional expectation. *Journal of Computational and Graphical Statistics*, 24(1):44–65.

Harari, G. M., Müller, S. R., Aung, M. S., and Rentfrow, P. J. (2017). Smartphone sensing methods for studying behavior in everyday life. *Current opinion in behavioral sciences*, 18:83–90.

Harari, G. M., Müller, S. R., Stachl, C., Wang, R., Wang, W., Bühner, M., Rentfrow, P. J., Campbell, A. T., and Gosling, S. D. (2020). Sensing sociability: Individual differences in young adults' conversation, calling, texting, and app use behaviors in daily life. *Journal of Personality and Social Psychology*, 119(1):204–228.

James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013). *An Introduction to Statistical Learning: with Applications in R*. Springer.

John, O. P., Naumann, L. P., and Soto, C. J. (2008). Paradigm shift to the integrative big five trait taxonomy. *Handbook of personality: Theory and research*, 3(2):114–158.

Kosinski, M., Stillwell, D., and Graepel, T. (2013). Private traits and attributes are predictable from digital records of human behavior. *Proceedings of the National Academy of Sciences*, 110(15):5802–5805.

Lakens, D. (2013). Using a smartphone to measure heart rate changes during relived happiness and anger. *IEEE transactions on affective computing*, 4(2):238–241.

Lang, M., Binder, M., Richter, J., Schratz, P., Pfisterer, F., Coors, S., Au, Q., Casalicchio, G., Kotthoff, L., and Bischl, B. (2019). mlr3: A modern object-oriented machine learning framework in R. *Journal of Open Source Software*.

Lundberg, S. M. and Lee, S.-I. (2017a). A unified approach to interpreting model predictions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, page 4768–4777, Red Hook, NY, USA. Curran Associates Inc.

Lundberg, S. M. and Lee, S.-I. (2017b). A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30.

Molnar, C. (2022). *Interpretable Machine Learning*. 2 edition.

Molnar, C., Bischl, B., and Casalicchio, G. (2018). iml: An r package for interpretable machine learning. *JOSS*, 3(26):786.

Mønsted, B., Mollgaard, A., and Mathiesen, J. (2018). Phone-based metric as a predictor for basic personality traits. *Journal of Research in Personality*, 74:16–22.

Montag, C., Błaszkiewicz, K., Lachmann, B., Andone, I., Sariyska, R., Trendafilov, B., Reuter, M., and Markowetz, A. (2014). Correlating personality and actual phone usage. *Journal of Individual Differences*.

Montag, C. and Elhai, J. D. (2019). A new agenda for personality psychology in the digital age? *Personality and Individual Differences*, 147:128–134.

Montgomery, D. C., Peck, E. A., and Vining, G. G. (2021). *Introduction to linear regression analysis*. John Wiley & Sons.

Omheni, N., Mazhoud, O., Kalboussi, A., and HadjKacem, A. (2014). Prediction of human personality traits from annotation activities. In *Proceedings of the 10th International Conference on Web Information Systems and Technologies - Volume 2: WEBIST,*, pages 263–269. INSTICC, SciTePress.

Probst, P., Au, Q., Casalicchio, G., Stachl, C., and Bischl, B. (2017). Multilabel Classification with R Package mlr. *The R Journal*, 9(1):352–369.

Qureshi, M. N. I., Min, B., Jo, H. J., and Lee, B. (2016). Multiclass classification for the differential diagnosis on the ADHD subtypes using recursive feature elimination and hierarchical extreme learning machine: Structural MRI study. *PLOS ONE*, 11(8):e0160697.

Ribeiro, M. T., Singh, S., and Guestrin, C. (2016). "why should i trust you?": Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '16, page 1135–1144, New York, NY, USA. Association for Computing Machinery.

Rudin, C. (2019). Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206–215.

Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323(6088):533–536.

Saeb, S., Zhang, M., Karr, C. J., Schueller, S. M., Corden, M. E., Kording, K. P., and Mohr, D. C. (2015). Mobile phone sensor correlates of depressive symptom severity in daily-life behavior: an exploratory study. *Journal of medical Internet research*, 17(7):e175.

Schoedel, R., Au, Q., Völkel, S. T., Lehmann, F., Becker, D., Bühner, M., Bischl, B., Hussmann, H., and Stachl, C. (2018). Digital footprints of sensation seeking. *Zeitschrift für Psychologie*, 226(4):232–245.

Schoedel, R., Pargent, F., Au, Q., Völkel, S. T., Schuwerk, T., Bühner, M., and Stachl, C. (2020). To challenge the morning lark and the night owl: Using smartphone sensing data to investigate day–night behaviour patterns. *European Journal of Personality*, 34(5):733–752.

Schütze, M., de Souza Costa, D., de Paula, J. J., Malloy-Diniz, L. F., Malamut, C., Mamede, M., de Miranda, D. M., Brammer, M., and Romano-Silva, M. A. (2018). Use of machine learning to predict cognitive performance based on brain metabolism in neurofibromatosis type 1. *PLOS ONE*, 13(9):e0203520.

Schuwerk, T., Kaltefleiter, L. J., Au, J.-Q., Hoesl, A., and Stachl, C. (2019). Enter the wild: Autistic traits and their relationship to mentalizing and social interaction in everyday life. *Journal of Autism and Developmental Disorders*, 49(10):4193–4208.

Servia-Rodríguez, S., Rachuri, K. K., Mascolo, C., Rentfrow, P. J., Lathia, N., and Sandstrom, G. M. (2017). Mobile sensing at the service of mental well-being. In *Proceedings of the 26th International Conference on World Wide Web*. International World Wide Web Conferences Steering Committee.

Shapley, L. S. (1953). A value for n-person games. *Contributions to the theory of games*, 2(28):307–317.

Stachl, C., Au, Q., Schoedel, R., Gosling, S. D., Harari, G. M., Buschek, D., Völkel, S. T., Schuwerk, T., Oldemeier, M., Ullmann, T., Hussmann, H., Bischl, B., and Bühner, M. (2020). Predicting personality from patterns of behavior collected with smartphones. *Proceedings of the National Academy of Sciences*, 117(30):17680–17687.

Stachl, C., Hilbert, S., Au, J., Buschek, D., Luca, A. D., Bischl, B., Hussmann, H., and Bühner, M. (2017). Personality traits predict smartphone usage. *European Journal of Personality*, 31(6):701–722.

Stachl, C., Schoedel, R., Au, Q., Völkel, S., Buschek, D., Hussmann, H., Bischl, B., and Bühner, M. (2022). The phonestudy project.

Štrumbelj, E. and Kononenko, I. (2013). Explaining prediction models and individual predictions with feature contributions. *Knowledge and Information Systems*, 41(3):647–665.

Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):267–288.

Waegeman, W., Dembczyński, K., and Hüllermeier, E. (2018). Multi-target prediction: a unifying view on problems and methods. *Data Mining and Knowledge Discovery*, 33(2):293–324.

Wan, S., Qi, L., Xu, X., Tong, C., and Gu, Z. (2020). Deep learning models for real-time human activity recognition with smartphones. *Mobile Networks and Applications*, 25:743–755.

Wilmer, H. H., Sherman, L. E., and Chein, J. M. (2017). Smartphones and cognition: A review of research exploring the links between mobile technology habits and cognitive functioning. *Frontiers in Psychology*, 8.

Zbiciak, A. and Markiewicz, T. (2023). A new extraordinary means of appeal in the polish criminal procedure: the basic principles of a fair trial and a complaint against a cassatory judgment. *Access to Justice in Eastern Europe*, 6(2):1–18.

Zhang, M.-L. and Zhou, Z.-H. (2013). A review on multi-label learning algorithms. *IEEE transactions on knowledge and data engineering*, 26(8):1819–1837.

Zuckerman, M. (1994). *Behavioral expressions and biosocial bases of sensation seeking*. Cambridge university press.

Zuckerman, M. (2002). Zuckerman-kuhlman personality questionnaire (zkpq): an alternative five-factorial model.

# Eidesstattliche Versicherung

(Siehe Promotionsordnung vom 12. Juli 2011, §8 Abs. 2 Pkt. 5)

Hiermit erkläre ich an Eides statt, dass die Dissertation von mir selbstständig, ohne unerlaubte Beihilfe angefertigt ist.

München, den 01.02.2024                                   Jiew-Quay Au