

Standards and Incentives in Safety Regulation

Inaugural-Dissertation
zur Erlangung des Grades
Doctor oeconomiae publicae (Dr. oec. publ.)
an der Ludwig-Maximilians-Universität München

2010

vorgelegt von
Felix Reinshagen

Referent: Prof. Dr. Klaus M. Schmidt
Korreferent: Prof. Ray Rees
Promotionsabschlussberatung: 1. Juni 2011

Acknowledgements

First and foremost, I would like to thank my supervisor, Klaus M. Schmidt, for the support and counsel that he has given me while writing this dissertation. I am also very thankful to Ray Rees and Sven Rady, who encouraged me – especially during the difficult start – and agreed to act as second and third examiner.

I have received valuable feedback and direction from a number of people. Among others, I would like to mention Jennifer Arlen, Benno Bühler, Matthias Dischinger, Florian Englmaier, Georg Gebhardt, Daniel Göller, Sandra Ludwig, Martin Schneider, Caspar Siegert, Anno Stolper, Christina Strassmair and Robert Ulbricht.

Thanks go to my numerous officemates, who – at some time or another – happily shared one of the large number of offices that I occupied during my doctoral career, namely Silvia Appelt, Joachim Klein, Nicolas Klein, Johannes Maier, Jan Schikora, Caspar Siegert and Anno Stolper.

I am also grateful to the Department of Economics for admitting me to the doctoral program, despite my unusual background as a lawyer. I received financial support from the Deutsche Forschungsgemeinschaft (DFG) through GRK 801, which I gratefully acknowledge.

Many thanks go to my family, who supported me during my long education, but especially to my late grandmother, Irmgard Gauss, who generously “topped up” my scholarship grant. Finally, I wish to thank Christina Strassmair – her love, support and advice was invaluable during the last years.

Contents

Preface	6
1 Standards and Incentives under Moral Hazard with Limited Liability	11
1.1 Introduction	11
1.2 Setup of the Model	15
1.3 Benchmark Case: Incentives only	16
1.4 Joint Use of Incentives and Standards	18
1.5 Conclusion	22
1.6 Appendix	22
2 Regulation and Liability as Complements and Substitutes	28
2.1 Introduction	28
2.2 Setup of the Model	32
2.3 The First-Best Care Levels	34
2.4 Care Levels with Liability only	35
2.5 The Optimal Combination of Liability and Regulation	37
2.6 Conclusion	42
2.7 Appendix	43
3 Minimum Education Requirements for Professions	47
3.1 Introduction	47
3.2 Objective Function of the Agent	51
3.3 Market Equilibrium	53

<i>Acknowledgements</i>	4
3.4 Regulatory Interventions	56
3.5 Welfare Analysis	59
3.5.1 Achieving the First-Best	60
3.5.2 Comparisons in the Second Best	62
3.6 Conclusion	64
3.7 Appendix	67
3.7.1 Proof of Proposition 3.6	67
3.7.2 Derivation of equation 3.5.2	69
3.7.3 Example with cost function $c(k) = \frac{a}{k}$	69
Bibliography	71

List of Figures

2.1	Observable Effort with Substitutes	40
2.2	Observable Effort with Complements	41
3.1	Timing of the model	52
3.2	Payoffs without regulation	57
3.3	Payoffs with minimum human capital	59
3.4	Payoffs with minimum human capital (not costless)	62

Preface

The optimal regulation of safety is a popular topic, both for academics and for the wider public. On the one hand, technological progress has multiplied the destructive potential of human action. The debates about nuclear technology or genetically modified organisms give daily proof to this statement. On the other hand, the same progress has increased human capability to guard against dangers while rising incomes have increased the demand for safety. What former generations might have mourned as a “stroke of fate” will today be denounced as a “preventable tragedy”.

The three chapters in this dissertation all deal with the optimal regulation of safety. By safety we mean the prevention of accidents, i.e. events which cause significant harm but occur with a low probability. The last criterion distinguishes accidents from bad outcomes that are considered “normal”. For example, the failure of a company to produce an important innovation will not be considered as an accident. Still, this definition of accidents might be broader than the ordinary use of the word. For example, consider the case of a physician who fails to perform an important test. This failure might cause a bad outcome for the patient, albeit with a low a probability. If the bad outcome is realised, our definition of an accident is fulfilled.

Consider an agent whose activity might lead to an accident. There are usually many actions that the agent can take to lower the probability of an accident. According to the usual criterion of cost-benefit analysis, an action should be taken if the benefits of this action – the decrease in expected damage from accidents – are greater than or at least equal to the costs of the action.

If the performance of these actions can be observed by a regulator or a court, the optimal regulation of safety is straightforward. All actions that pass the cost-benefit test should be prescribed by law or regulation. We will call this approach “standard setting”. But in many cases, the observation of these actions will be impossible or at least very expensive. This means that a strategy of “standard setting” will be infeasible or at least very expensive. In such a situation the regulator can still implement all preventive actions which pass the cost-benefit test if he is able to impose on the agent all costs that arise from a potential accident. This strategy, the use of “incentives”, makes the agent the “residual claimant” on the risk. It will induce the agent to voluntarily implement all actions that pass the cost-benefit test. The legal institution of “strict liability”, either on a statutory or a contractual basis, can be seen as an attempt to make the agent the residual claimant.

In practice it is often very difficult to shift all potential costs of an accident to the agent. The financial assets of the agent might be far below the amount that is necessary to pay for all accidental damages in the worst case. It might also be very difficult to find all victims of an accident and to accurately estimate the damages they have suffered. Finally, the agent might be able to escape the reach of the law. If the liability of the agent is limited in such a way, he is no longer induced to take all preventive activities that pass the cost-benefit test.

So in practice neither standard setting nor incentives will work perfectly. Standard setting can only affect actions which are easily observable, while incentives can shift only parts of the costs of accidents to the agent. We will therefore investigate the optimal combination of these two approaches. We will assume that preventive actions can be separated into two distinct categories, observable actions and unobservable actions. Examples of actions that we think of as observable are conformance with technical norms in the construction of a plant, or whether a physician has access to adequate equipment. Another example of an observable action is the acquisition of human capital by the persons who will control the risk, e.g. whether he has a relevant academic degree or has passed a prescribed examination. Examples of actions that we think of

as unobservable are the mental alertness of the person who controls the risk or whether safety procedures are actually followed in day-to-day work.

Given this distinction, it will be a natural outcome of our models that most of the time it is optimal to use standards and incentives at the same time. Our main interest will be the optimal level of standard setting in such a situation. In the following chapters we show that inadequate incentives for unobservable actions will also influence the optimal level of standards for observable actions. Often this optimal level will differ from the optimal use of observable actions in a first-best world. This has the consequence that a cost-benefit analysis of standards that ignores the interaction between observable and unobservable preventive actions might be flawed.

The first chapter is an application of this idea to contract theory. We consider a model of moral hazard with limited liability of the agent and effort that is two-dimensional. One dimension of the agent's effort is observable and the other is not. The principal can thus make the contract conditional not only on outcome but also on observable effort. In this chapter we make the assumption that there is no interaction between the costs or returns of the two kinds of effort; nevertheless, the limited liability of the agent will influence the levels of both kinds of effort. The principal's optimal contract gives the agent no rent and – in contrast to the first-best solution – uses too much observable effort and too little unobservable effort. This distortion in the relative use of the two kinds of effort increases if the agent's liability becomes more limited.

The second chapter integrates two-dimensional care into the model of accident prevention first formulated in Shavell (1980). In analogy to the first chapter, one dimension of care is observable while the other is not.¹ We consider a situation where the firm's liability is limited and analyze the use of strict liability combined with ex-ante regulation of observable care. In comparison to other approaches to model the combination of liability and regulation, joint use of the two instruments is a natural outcome in our model. In contrast to the first chapter, we allow for technical interactions between the two dimen-

¹In accordance with much of the literature in Law & Economics, we use the term "care" instead of effort.

sions of care.² If the two kinds of care are independent (the level of one kind of care does not influence the marginal cost or the marginal return of the other kind of care), it is optimal to regulate observable care at the (socially optimal) first-best level. If the two dimensions of care are complements or substitutes, the optimal level of regulation is influenced by two considerations. An increase in observable care has direct benefits because it decreases the probability of an accident but has also indirect effects because it influences the firm's incentive to take unobservable care. We show that if the moral hazard problem is serious enough, regulated observable care will be below its first-best level in the case of complements and above its first-best level in the case of substitutes.

The third chapter deals with a controversial instrument to prevent accidents, namely minimum education requirements for professional services. Many critics allege that these requirements are unnecessary or at least excessively strict. We provide a partial justification for these requirements, by demonstrating how they can induce high quality of work, even though they only regulate education, which is just one input for this work. In our model, minimum education requirements serve as an hostage to ensure high quality. The threat of losing their occupation-specific human capital makes professionals more sensitive to the punishment of being excluded from the profession, but the same human capital makes it also easier for them to do high quality work. We compare education requirements with an alternative method that works in a similar way, namely quantitative entry restrictions. We show that under certain parameters minimum education requirements achieve the first-best solution. Furthermore, when social welfare is given by consumer surplus, education requirements are always preferable to quantitative entry restrictions.

The reader will find that limited liability is a crucial assumption in all three chapters. This poses the question whether there is an easy way to lift such a limit by another mechanism. For example, a regulator could demand that an agent who wants to undertake the dangerous activity "posts a bond", i.e. deposits an amount of money equal the potential damage with the regulator.

²In addition, the social welfare function is different: the principal does also care about the payoff of the agent.

For such “bonding” to work perfectly, the deposited money has to come out of the private wealth of the person who controls the risk. For example, if the agent is a corporation, the deposit should be paid by the executive who controls the corporation’s safety policy. If someone else, like a bank, an insurance company or the corporation’s shareholders provide the money, a new principal-agent problem arises between these financiers and the person who controls the risk. The first-best will only be implemented if a financier can perfectly observe this person’s safety actions. If he cannot do so, the financier will probably impose his own safety standards on the agent – the financier will become the regulator. For example, an insurance company might make coverage conditional on the insuree’s compliance with certain technical regulations or his proof of an adequate education level.

If the bond has to come out of the private wealth of the person who controls the risk, many occupations will be reserved for the rich, which raises serious issues of social justice and the allocation of talent in society. In the extreme, a position like CEO of a firm operating nuclear power stations might not be filled at all. So while bonding can play an important role in alleviating problems of moral hazard, it does not render the mechanisms discussed here irrelevant.

Chapter 1

Standards and Incentives under Moral Hazard with Limited Liability

1.1 Introduction

Consider a principal-agent relationship with moral hazard. There will probably be many actions that the agent can take to further the principal's project. Some of these actions will be observable, some not. In the following, we will subsume all actions that are observable under the term observable effort, and all actions that are not observable under the term unobservable effort. In the first-best, without moral hazard, the optimal mix of efforts will in general include a mix of both kinds of effort. The contract that is usually assumed in situations of moral hazard is conditional on the observed *outcome* only. In this chapter we will look at a contract that is also conditional on the level of observable *effort*. This means that the contract will stipulate a specific level of observable effort and the principal will only pay if he observes at least this level of observable effort.

Our main interest in this chapter is the level of the contractually specified observable effort and its relation to the induced level of unobservable effort. We assume that there is no direct interaction between the costs or returns of the

two kinds of effort; nevertheless, the limited liability of the agent will influence the levels of both kinds of effort. Moral hazard problems with limited liability of the agent usually have the following outcome: if the principal cannot extract the whole surplus at the first-best level of effort, he will lower the implemented effort below the first-best level.¹ In contrast, in our model the specified level of observable effort will be above the first-best level, while unobservable effort will still be below the first-best level. This will also mean that the combination of observable and unobservable effort will not be cost-minimizing, i.e. the given amount of total effort is produced with too much observable effort and too little unobservable effort. In other words, the agent would be able to produce the same level of total effort with lower costs.

For an application, think about a situation where the principal wants the agent to undertake a project that can fail with catastrophic consequences. Consider a government that licenses a firm to operate an hazardous technology, like a chemical factory or a nuclear reactor. The government wants the firm to undertake effort that increases the probability that the firm operates safely. Some of this effort, like the compliance with technical regulations for the construction of the plant, or the education level of the operating personnel can be controlled rather easily. But other elements essential to safe operation will be very hard to observe, like the workload and alertness of the personnel or whether the firm's management exerts pressure on them to "bend the rules". The "regulatory contract" in such situations usually includes both standards for observable effort ("regulation") and monetary payments that depend on the outcome of the project ("fines" and "liability"). The compliance with the standards can and will be enforced ex-ante, while the ex-post payments give the firm incentives to undertake unobservable effort. A similar problem exists if a big firm subcontracts part of a project to a small firm. If the small firm produces bad quality, the damage for the big firm might be immense. Contractual arrangements in such situations will usually not only include payments that are conditional on final outcomes but will also authorise the big firm to monitor whether the work of the small firm is in compliance with contractual standards. In addition, the

¹This may or may not imply a rent for the agent.

big firm might demand that the small firm will have its operations “certified” by a third party.

Our result suggest that in such situations the principal will set standards that demand observable effort which is above the first-best level, but the level of total effort will be below the first-best. For example, the work of a small sub-contractor will be more oriented toward observable effort compared to the case where the big firm would do the work itself. To generalize, we suggest a possible inefficiency existing under moral hazard with limited liability, which does not lie in the amount of total effort but in the way this effort is produced. This inefficiency has seen scant attention in theory but is often complained about in practice.

Many employees of big organizations complain about “bureaucracy”. They feel that their work is inefficiently organized – it would be more productive if there were fewer regulations to observe and more time could be spend on doing “real work”. Regulatory regimes for hazardous activities are criticized for putting too much emphasis on compliance with technical standards rather than on soft factors like “safety culture”. And many observers question whether a firm’s decision to seek certification for use of a “quality management systems” is mainly motivated by customer pressure, while the real effect on quality is questionable.²

This work is related to a number of papers which all exploit a similar effect: if the solution to the moral hazard problem calls for granting the agent a rent, the principal will try to expropriate this rent by forcing the agent to undertake some other activity that benefits the principal. This activity might be socially inefficient, but because its costs come out of the agent’s rent, it is still advantageous for the principal to implement it. The activity in question might be another principal-agent project (Laux, 2001), reporting activities like “paperwork”

²The question whether firms introducing ISO 9000 quality management systems are mainly motivated by external reasons (customer pressure etc.) or by internal reasons (concern for quality and cost improvements) has been the subject of numerous studies, which have come to conflicting results. A overview of previous studies can be found in Heras Saizarbitoria et al. (2006); the Delphi study described in their paper finds that external reasons are dominating. In a similar vein, Buttle (1997) describes a survey of ISO 9000 certified firms; the highest scoring motivation for certification is “anticipated demand from future customers for ISO 9000”.

(Strausz, 2006) or the effort in a preceding period of the principal-agent relationship (Kräkel and Schöttner, 2010). Our model is the first that applies this effect to the choice between observable effort and unobservable effort. This setting is not only of great practical importance, it does also allow for a sharp characterization of the trade-off that is responsible for the implementation of a socially inefficient activity.

In the “Law & Economics” literature, Bhole and Wagner (2008) analyze a setting where a firm can take observable effort as well as unobservable effort to prevent an accident.³ They find that in many situations only the combined use of both liability and regulation will lead to optimal levels of effort in both dimensions. There are two important differences to our approach. First, in a tort law setting the principal has a different objective function (total welfare) and usually a restricted choice of policy measures. Second, Bhole and Wagner only consider a binary choice of observable effort; because in their model a high level of observable effort is first-best, the question of excessive regulation of observable effort is ruled out by assumption.

Multi-dimensional effort has been studied in number of other settings in the literature. In the most prominent treatment by Holmstrom and Milgrom (1991), different dimensions of effort interact through the agent’s cost function. In our model, there is no such interaction; observable effort and unobservable effort influence each other only because of the shared limited liability constraint.

The rest of the chapter is structured as follows: Section 1.2 sets up the model. In section 1.3, we discuss a benchmark case, namely a contract that is conditional on outcome only. The main part of the chapter is section 1.4, which analyzes a contract that does also regulate the agents effort, while section 1.5 concludes. Proofs can be found in the Appendix.

³In an article on liability for nuclear accidents, Trebilcock and Winter (1997) sketch a tort-law model with observable and unobservable effort but do not fully solve it.

1.2 Setup of the Model

There are two kinds of effort, observable effort $o \in [0, o_{max}]$ and unobservable effort $u \in [0, u_{max}]$ with $o_{max}, u_{max} > 0$ and $o_{max} + u_{max} \leq 1$. The agent's project has two outcomes, it can either succeed or fail, $s \in \{0, 1\}$. The probability of success ($s = 1$) depends on the agents effort and is given by $p(o, u) = o + u$. At times we will denote this probability as *total effort*. If the agent exerts effort, he suffers costs of $c_o(o) + c_u(u)$. Note that under this setup there is no direct interaction between the two kinds of effort: the level of one kind of effort does not influence the marginal cost or the marginal return of the other kind of effort.⁴ We further need the following technical assumptions for the cost functions:

Assumption 1.1. $c_o(o)$ and $c_u(u)$ are continuous, three times differentiable, strictly increasing and strictly convex.

Assumption 1.2. $c_o(o_{max}) = c_u(u_{max}) = \infty$.

Assumption 1.3. $c'_o(0) = c'_u(0) = 0$.

Assumption 1.4. $c'''_o(o), c'''_u(u) > 0$.

Assumption 1.5. $c_o(0) = c_u(0) = 0$.

Assumptions 1.2 and 1.3 ensure that the agent's problem has an interior solution, while Assumption 1.4 makes the principal's problem concave (the condition on $c'''_o(o)$ is only needed for the benchmark case).

The benefit for the principal if the project succeeds is set to $B > 0$. Both parties are risk neutral. To induce effort, the principal will write a contract that specifies a transfer scheme $t(s, o)$ that can depend on the outcome of the project and the observed effort. The agent faces a liability limit $L \geq 0$, which can either be interpreted as the maximum fine that can be imposed on the agent ex-post, or the maximum bond that can be posted by the agent ex-ante.⁵ This liability limit is expressed by:

⁴In reality those direct interaction will often exist, making the two kinds of efforts either complements or substitutes. In this chapter, we assume no direct interaction to isolate those effects that are due to limited liability.

⁵We assume that the liability limit does not depend on the level of efforts.

Assumption 1.6. $t(s, o) \geq -L \forall s \in \{0, 1\}, o \in [0, o_{max}]$.

We have to distinguish two concepts. On the one hand, we have the socially optimal *first-best* effort levels o^* and u^* , which are given by $c'_o(o^*) = B$ and $c'_u(u^*) = B$. On the other hand, for a given level of total effort \bar{p} , we can find the least expensive combination of observable and unobservable effort that produces \bar{p} . Such a *cost-minimizing* combination of efforts will be characterized by $c'_o(o) = c'_u(u)$.⁶ It is easy to see that first-best effort levels are also a cost-minimizing combination of efforts, but that there are also many other cost-minimizing combinations of efforts that are not first-best.

1.3 Benchmark Case: Incentives only

To establish a benchmark case, we will first consider a contract that conditions only on *outcome*. This contract can be described by the transfer scheme:

$$t(s, o) = \begin{cases} b + w & \text{if } s = 1 \\ w & \text{if } s = 0 \end{cases}$$

It has the usual property that the principal sets a base wage w and a bonus b . It follows that the profit function of the principal is given by $\Pi(o, u, b, w) = (B - b) \cdot p(o, u) - w$, while the payoff function of the agent is $V(o, u, b, w) = bp(o, u) + w - c_o(o) - c_u(u)$. The principal has to solve the problem:

$$\begin{aligned} & \max_{o, u, b, w} \Pi(o, u, b, w) \\ & \text{subject to:} \\ & V(o, u, b, w) \geq 0 \quad \text{PC} \\ & w \geq -L, w + b \geq -L \quad \text{LLCs} \\ & (o, u) \in \underset{(o, u)}{\operatorname{argmax}} V(o, u, b, w) \quad \text{IC} \end{aligned} \tag{1.1}$$

⁶This condition results from $\min_{o, u} c_o(o) + c_u(u)$, subject to $p(o, u) = \bar{p}$. Formally, the marginal rate of technical substitution between these two kinds of effort must be equal to the ratio of respective marginal costs.

The fact that u is unobservable does not necessarily mean that the first-best will not be implemented. In fact, if the principal sets $b = B$, the agent will deliver effort levels o^* and u^* . The wage w^* that extracts all the agent's surplus is then given by $V(o^*, u^*, B, w^*) = 0$, which can be written as $w^* = c_o(o^*) - c_u(u^*) - Bp(o^*, u^*)$.

But this extraction of surplus is feasible only if $w^* \geq -L$; in this case, the principal can “sell the project” to the agent. If $w^* < -L$, the principal faces a tradeoff between incentivizing effort and extracting rent. In the following, we will always assume that the first-best will not be implemented, namely

Assumption 1.7. $w^* < -L$.

We will find the optimal effort levels o_{bm} and u_{bm} by using the so-called first-order approach. The following proposition shows that this approach is valid in our setting because the agent's optimal choice of effort levels is at a stationary point.

Proposition 1.1. *The optimal solution to (1.1) has $b > 0$ and o_{bm}, u_{bm} will be given by the agent's first-order order conditions $b - c'_o(o) = 0$ and $b - c'_u(u) = 0$, with $o_{bm} \in (0, o_{max})$, $u_{bm} \in (0, u_{max})$ and total effort $p(o, u) > 0$.*

We can therefore replace the incentive constraint with the agent's first-order conditions. Additionally, because $b > 0$, one of the limited liability constraints, $w + b \geq -L$, is superfluous. The Lagrangian for the principal's problem can now be written as:

$$\begin{aligned} \mathcal{L}(o, u, b, w, \lambda, \eta, \mu_o, \mu_u) = \\ (B - b) \cdot p(o, u) - w + \lambda (b \cdot p(o, u) + w - c_o(o) - c_u(u)) \\ + \eta (w + L) + \mu_o (b - c'_o(o)) + \mu_u (b - c'_u(u)) \end{aligned} \quad (1.2)$$

In the optimal solution, the limited liability constraint $w \geq -L$ will always be binding, while the participation constraint may be binding or not.

Proposition 1.2. *The optimal solution to (1.2) has $w = -L$ and*

$$b = B - \frac{(1 - \lambda)p(o, u)}{\frac{1}{c''_o(o)} + \frac{1}{c''_u(u)}} \quad (1.3)$$

with $0 \leq \lambda < 1$. If the agent will get a rent, we have $\lambda = 0$.

The optimal contract can be found by trying out two cases. In the first case with $\lambda = 0$, the optimal effort levels are given by the trade-off between the costs of incentives and the principal's benefit from having more effort, ignoring the PC (this will usually mean a rent for the agent). But if those effort levels and $w = -L$ do not satisfy the PC, we have the case $\lambda > 0$. The principal sets $w = -L$ and chooses the unique level of b that makes the PC binding. This will mean higher effort levels than in the first case and no rent for the agent.⁷

In both cases we will have $c'_o(o_{bm}) = c'_u(u_{bm}) = b < B$. This implies that both kinds of effort are below the first-best level ($o_{bm} < o^*$ and $u_{bm} < u^*$), but because $c'_o(o_{bm}) = c'_u(u_{bm})$, they form a cost-minimizing combination.

1.4 Joint Use of Incentives and Standards

We now look at a contract that makes the principal's payments conditional not only on outcome, but also on observable effort. At first glance the problem of finding the optimal contract looks quite simple: set the observable effort to o^* and optimize over u (because we assume $p(o, u) = o + u$, there is no interaction between the two kinds of effort). But it will turn out that the optimal contract will have a level of observable effort that is above o^* .

We consider contracts of the following form:⁸

$$t(s, o) = \begin{cases} b + w & \text{if } s = 1 \text{ and } o \geq \underline{o} \\ w & \text{if } s = 0 \text{ and } o \geq \underline{o} \\ -L & \text{if } o < \underline{o} \end{cases}$$

⁷Which case obtains depends on the severity of the liability limit. Define L^* by $V(o^*, u^*, B, -L^*) = 0$ and \tilde{L} by $V(o_{bm}, u_{bm}, B, -\tilde{L}) = 0$ (where o_{bm} and u_{bm} are given by (1.3) with $\lambda = 0$). If $0 \leq L < \tilde{L}$ the agent gets a rent, if $\tilde{L} \leq L < L^*$ there will be no rent.

⁸The principal cannot improve his profit by using a more general contract that distinguishes between more levels of o , because, besides his effort level, the agent has no other private information.

where \underline{o} is contractually specified level of observable care. The principal's expected profit is given by

$$\Pi_{\underline{o}}(o, u, b, w) = \begin{cases} (B - b) \cdot p(o, u) - w & \text{if } o \geq \underline{o} \\ B \cdot p(o, u) + L & \text{if } o < \underline{o} \end{cases}$$

while the agent's payoff has the form:

$$V_{\underline{o}}(o, u, b, w) = \begin{cases} bp(o, u) + w - c_o(o) - c_u(u) & \text{if } o \geq \underline{o} \\ -L - c_o(o) - c_u(u) & \text{if } o < \underline{o} \end{cases}$$

The principal's problem is given by:

$$\begin{aligned} & \mathbf{max}_{o, u, b, w, \underline{o}} \Pi_{\underline{o}}(o, u, b, w) \\ & \mathbf{subject\ to:} \\ & V_{\underline{o}}(o, u, b, w) \geq 0 \quad \text{PC} \\ & w \geq -L, \quad w + b \geq -L \quad \text{LLCs} \\ & (o, u) \in \underset{(o, u)}{\operatorname{argmax}} V_{\underline{o}}(o, u, b, w) \quad \text{IC} \end{aligned} \tag{1.4}$$

Denote by \hat{o} and \hat{u} the effort levels that are implemented in the optimum. The first problem is again to show that the first-order approach is valid here.

Proposition 1.3. *The optimal solution to (1.4) has $\hat{o} = \underline{o}$ and $b > 0$. Effort level \hat{u} will be given by the agent's first-order order condition $b - c'_u(u) = 0$, with $\hat{o} \in (0, o_{max})$, $\hat{u} \in (0, u_{max})$ and total effort $p(\hat{o}, \hat{u}) > 0$.*

We can again use the agent's first order condition for u and ignore the constraint $w + b \geq 0$. The Lagrangian for the principal's problem can be written as:

$$\begin{aligned} \mathcal{L}(o, u, b, w, \lambda, \eta, \mu) = & \\ & (B - b) \cdot p(o, u) - w + \lambda (b \cdot p(o, u) + w - c_o(o) - c_u(u)) \\ & + \eta (w + L) + \mu (b - c'_u(u)) \end{aligned} \tag{1.5}$$

Proposition 1.4. *In the optimal solution to (1.5), both the participation constraint $V_{\underline{o}}(o, u, b, w) \geq 0$ and the limited liability constraint $w \geq -L$ are binding. The optimal effort levels \hat{o} and \hat{u} are given by*

$$c'_o(o) = B + \frac{1-\lambda}{\lambda}(B - c'_u(u)) \quad (1.6)$$

$$\text{and } c'_u(u) = B - (1-\lambda) \cdot p(o, u) \cdot c''_u(u) \quad (1.7)$$

with $0 < \lambda < 1$.

It is quite intuitive that the principal will not give the agent a rent. Suppose the principal would choose some o and some $b < B$ so that the agent gets a rent. The principal could then increase observable effort and get a marginal benefit of $B - b$ while letting the agent take the additional costs out of his rent. So the principal will transform the agent's rent into his own benefit.

From (1.6) and (1.7) and $0 < \lambda < 1$ we can conclude that $c'_u(u) < B$ and $c'_o(o) > B$. This implies that $\hat{o} > o^*$ and $\hat{u} < u^*$, so observable effort is above and unobservable effort is below the first-best level. We also note that \hat{o} and \hat{u} are not a cost-minimizing combination of efforts (because $c'_o(\hat{o}) \neq c'_u(\hat{u})$), meaning that $p(\hat{o}, \hat{u})$ could be produced more cheaply by a different combination of efforts. It is also clear that $\hat{o} > o_{bm}$, but we cannot tell whether \hat{u} is greater or smaller than u_{bm} . In fact, numerical simulations show that both cases can occur.

The principal is willing to set observable effort above the first-best level because stipulating more observable effort has the additional benefit of inducing more unobservable effort. When the principal demands additional observable effort, he must compensate the agent for the additional cost (because the PC is binding), but does so by increasing b , thereby increasing the agent's incentive for providing unobservable effort. This can be seen if we combine the two implicit equations (1.6) and (1.7) by eliminating λ :

$$c'_o(o) - B = (B - b) \cdot \frac{1}{c''_u(u)} \cdot \frac{1}{p(o, u)} (c'_o(o) - b) \quad (1.8)$$

Equation (1.8) can be interpreted as the trade-off facing the principal at the margin when he increases \hat{o} beyond o^* . The term on the left-hand-side is the princi-

pal's cost of increasing observable effort further above the first-best level. Because the PC is binding, he has to compensate the agent for the marginal cost of additional effort but receives additional expected benefit of only B (which is smaller than $c'_o(o)$ because $\hat{o} > o^*$). The right hand side is his marginal benefit and can be interpreted as follows (read from right to left): if o is increased, the agent has marginal costs of $c'_o(o)$ but receives a marginal increase in expected payoff of only b . To compensate the agent for a small loss in payoff, the principal has to marginally increase b by $\frac{1}{p(o,u)}$. An marginal increase in b will increase unobservable effort by $\frac{1}{c''_u(u)}$, while a marginal increase in u will give the principal an marginal benefit of $B - b$. These effects can be labeled as follows:

$$c'_o(o) - B = (B - b) \cdot \underbrace{\frac{1}{c''_u(u)}}_{\frac{d\Pi}{du}} \cdot \underbrace{\frac{1}{p(o,u)}}_{\frac{db}{dV}} \underbrace{(c'_o(o) - b)}_{-\frac{dV}{do}}$$

In our model, the agent's limited liability causes a combination of the two kinds of effort that is not cost-minimizing, namely too much observable and too little unobservable effort. This suggests that a decrease in L – the problem of limited liability becomes worse – will increase this distortion. The next proposition shows that this is indeed the case.

Proposition 1.5. *If L decreases (the agent's liability becomes more limited), \hat{o} increases and \hat{u} decreases.*

This result looks more obvious than it is. Because if L decreases, it changes not only the optimal combination of o and u that implements a given level of $p(o, u)$ (substitution effect), but it may also change the level of $p(o, u)$ that is optimal for the principal to implement (scale effect).⁹ Proposition 1.5 shows that the first effect dominates. This result also suggests a possible way to test our theory: for agents with a stricter liability limit we should observe standards that prescribe a higher level of observable effort.

⁹The terminology is taken from Nagatani (1978). It can be shown that if L decreases, the substitution effect is positive for o and negative for u . But if the optimal $p(o, u)$ decreases, the scale effect will be negative for both kinds of effort.

1.5 Conclusion

The chapter analyzes a model of moral hazard with limited liability of the agent where the agent's effort has one observable and one unobservable dimension. For simplicity, we only consider the case where the two kinds of efforts do not interact with each other. We consider different contracts with regard to two questions: whether each of the two kinds of effort is above or below its first-best level and whether the two levels form a cost-minimizing combination. With a contract that is conditional on outcome only, both kinds of effort are below their first-best levels but they form a cost-minimizing combination. With a contract that is conditional on both outcome and observable effort, unobservable effort will still be below its first best level while observable effort will be above the first-best level. This combination of efforts will not be cost-minimizing. The distortion between the two kinds of efforts increases if the agent's liability becomes more limited.

1.6 Appendix

Proof of Proposition 1.1

We first show that all $b \leq 0$ give the principal the same profit. If the principal sets $b \leq 0$, the agent will always choose $o = 0$ and $u = 0$, and $w = 0$ will make the PC binding. This will give the principal a profit $\Pi = 0$, for all $b \leq 0$. Thus to show that $b \leq 0$ is not optimal it is sufficient to show that $b = 0$ is not optimal

We now show that for a given $b \geq 0$ and w , the maximum of $V(o, u, b, w) = bp(o, u) + w - c_o(o) - c_u(u)$ will be characterized by the first-order conditions $b - c'_o(o) = 0$ and $b - c'_u(u) = 0$. Because $V(o, u, b, w)$ is strictly concave in o and u , an interior maximum will be characterized by the first-order conditions. As regards to corner solutions, $o = o_{max}$ or $u = u_{max}$ cannot be a maximum because the costs would be infinite, so zero effort would be better. A possible corner solution with $o = 0$ and $u = 0$ would have $b - c'_o(0) \leq 0$. Because of $b \geq 0$ and $c'_o(0) = 0$ this implies $b = 0$ and this maximum would also fulfill the first-order condition with equality.

Now we show that $b = 0$ cannot be an optimum. Suppose otherwise: then the agent would choose $o = 0$ and $u = 0$, and $w = 0$ would make the PC binding. If the principal would marginal increase b he would get:

$$\frac{d\Pi}{db} = -p(o, u) + (B - b) \left(\frac{do}{db} + \frac{du}{db} \right) \quad (1.9)$$

$$= B \left(\frac{1}{c''_o(o)} + \frac{1}{c''_u(u)} \right) > 0 \quad (1.10)$$

where the values of $\frac{do}{db}$ and $\frac{du}{db}$ come from implicitly differentiating the agent's first-order conditions; at the same time, at this point, $\frac{dV}{db} = p(0, 0) - c'_o(0) \frac{do}{db} - c'_u(0) \frac{du}{db} = 0$ so the PC will still be satisfied. Because this implies $o_{bm}, u_{bm} > 0$, we must have $p(o, u) > 0$. \square

Proof of Proposition 1.2

The first order conditions for a maximum are:

$$\frac{\partial \mathcal{L}}{\partial o} = (B - b) + \lambda(b - c'_o(o)) - \mu_o c''_o(o) = 0 \quad (1.11)$$

$$\frac{\partial \mathcal{L}}{\partial u} = (B - b) + \lambda(b - c'_u(u)) - \mu_u c''_u(u) = 0 \quad (1.12)$$

$$\frac{\partial \mathcal{L}}{\partial b} = -p(o, u) + \lambda p(o, u) + \mu_o + \mu_u = 0 \quad (1.13)$$

$$\frac{\partial \mathcal{L}}{\partial w} = -1 + \lambda + \eta = 0 \quad (1.14)$$

$$\lambda, \eta, \mu_o, \mu_u \geq 0 \text{ (with complementary slackness)}$$

It cannot be the case that both the PC and the LLC are slack. In this case, the principal could always increase his profit by decreasing w (formally, equation (1.14) can never be fulfilled). Furthermore, it cannot be the case that the PC is binding and the LLC is not binding: Because this would imply $\eta = 0$, and from (1.14) we would get $\lambda = 1$. Then (1.13) and complementary slackness give us $\mu_o = 0, \mu_u = 0$ and from this we get $o = o^*$ and $u = u^*$ (using equations (1.11) and (1.12)). But this contradicts Assumption 1.7.

By using the agent's first order conditions we can simplify the equations to

$$b + \mu_o c''_o(o) = B$$

$$b + \mu_u c''_u(u) = B$$

$$\mu_o + \mu_u = (1 - \lambda)p(o, u)$$

with $0 \leq \lambda < 1$. Solving this system of equations for b yields:

$$b = B - \frac{(1 - \lambda)p(o, u)}{\frac{1}{c_o''(o)} + \frac{1}{c_u''(u)}}.$$

This implies $b < B$ and $\mu_o, \mu_u > 0$.

For this solution to be a maximum, the Lagrange function evaluated with the Lagrange-multipliers found above must be concave. This function is given by:

$$\mathcal{L}^*(o, u, b, w) = Bp(o, u) + (1 - \lambda)L - \lambda[c_o(o) + c_u(u)] - \mu_o c_o'(o) - \mu_u c_u'(u)$$

which is concave in o, u, b and w (because $c_o'''(o), c_u'''(u) > 0$). \square

Proof of Proposition 1.3

We first show that $o = \underline{o}$ by contradiction. Consider the case $o < \underline{o}$. If the agent would like to disobey the contract, his optimal choice of efforts is $o = 0$ and $u = 0$, which would give him a payoff of $-L < 0$. But this cannot be optimal for the agent because obeying and delivering $o = \underline{o}$ would give him a non-negative payoff (because the principal has to fulfill the PC).

Now consider the case $o > \underline{o}$. This would be optimal for the agent if the effort level given by $c_o'(o) = b$ is higher than \underline{o} , or $c_o'(\underline{o}) < b$. To show the opposite first note that in the principal's optimum it must be the case that $c_o'(\underline{o}) \geq B$. If not, the principal could marginally increase o while holding the agent's payoff constant by increasing w . This would increase the principal's profit marginally by $B - c_o'(o) > 0$. Second, it cannot be optimal for the principal to set $b > B$. Consider

$$\frac{d\Pi}{db} = -p(o, u) + (B - b)\frac{1}{c_u''(u)}$$

which is negative for $b > B$. If the LLC is binding, this shows that decreasing b will increase Π . If the LLC is not binding, the principal could extract the increase in the agent's surplus

$$\frac{dV}{db} = p(o, u) + b\frac{1}{c_u''(u)} - c_u'(u)\frac{1}{c_u''(u)} = p(o, u)$$

where the last equality results from using the agent's first-order condition. So the principal's profit would increase by $\frac{d\Pi}{db} + \frac{dV}{db} = (B - b)\frac{1}{c_u''(u)}$ which is still negative for $b > B$. So we have $c_o'(\underline{o}) \geq B$ and $B \geq b$ which implies $c_o'(\underline{o}) \geq b$ which means that in the optimum \underline{o} will be greater than the effort level implied by $c_o'(o) = b$.

We next show that all $b \leq 0$ give the principal the same profit. Suppose the principal chooses some $o \in [0, o_{max}]$ and sets $b \leq 0$. Then the agent will choose $u = 0$. For the PC to hold the principal has to set $w = -bo + c_o(o) > 0 \geq -L$. This will give him the same profit $\Pi = Bo - c_o(o)$ for all $b \leq 0$.

The proof that for all $b \geq 0$ the optimal u will be given by the agent's first-order condition $b - c'_u(u) = 0$ is analogous to the argument in the proof of Proposition 1.1. Because $\hat{o} > o^* > 0$ we will also have $p(o, u) > 0$. \square

Proof of Proposition 1.4

The first order conditions for a maximum are:

$$\frac{\partial \mathcal{L}}{\partial o} = (B - b) + \lambda (b - c'_o(o)) = 0 \quad (1.15)$$

$$\frac{\partial \mathcal{L}}{\partial u} = (B - b) + \lambda (b - c'_u(u)) - \mu c''_u(u) = 0 \quad (1.16)$$

$$\frac{\partial \mathcal{L}}{\partial b} = -p(o, u) + \lambda p(o, u) + \mu = 0 \quad (1.17)$$

$$\frac{\partial \mathcal{L}}{\partial w} = -1 + \lambda + \eta = 0 \quad (1.18)$$

$$\lambda, \eta, \mu \geq 0 \text{ (with complementary slackness)}$$

We show that a solution to these conditions must have both the PC and the LLC binding. If the PC is not binding, we will have $\lambda = 0$. Then (1.15) gives us $b = B$. But from (1.17) we get $\mu = p(o, u)$ and plugging into (1.16) gives us $b = B - p(o, u)c''_u(u) < B$, a contradiction. If only the PC is binding but the LLC is not, we will have $\eta = 0$. From (1.18) we get $\lambda = 1$ and from (1.17) we get $\mu = 0$. Plugging into (1.15) and (1.16) gives us $c'_o(o) = B$ and $c'_u(u) = B$ respectively. This implies, that the first best can be achieved with a bonus contract without violating the LLC. But this contradicts Assumption 1.7.

From (1.18) we get $\lambda = 1 - \eta$ and with $\lambda, \eta > 0$, we must have $0 < \lambda < 1$. From (1.17) we get $\mu = (1 - \lambda)p(o, u) > 0$. Substituting for μ into (1.16) and rearranging gives us

$$c'_u(u) = b = B - (1 - \lambda) \cdot p(o, u) \cdot c''_u(u) < B$$

and substituting $b = c'_u(u)$ into (1.15) gives us:

$$c'_o(o) = B + \frac{1 - \lambda}{\lambda} (B - c'_u(u)) > B.$$

For this solution to be a maximum, the Lagrange function evaluated with the Lagrange-multipliers found above must be concave. This function is given by:

$$\mathcal{L}^*(o, u, b, w) = Bp(o, u) + (1 - \lambda)L - \lambda[c_o(o) + c_u(u)] - \mu c'_u(u)$$

which is concave in o, u, b and w (because $c''_u(u) > 0$). \square

Proof of Proposition 1.5

We have to show that $\frac{do}{dL} < 0$ and $\frac{du}{dL} > 0$. The optimal values for o, u, b and w are given by the solution to the four equations:

$$b \cdot p(o, u) + w - c_o(o) - c_u(u) = 0 \quad (1.19)$$

$$L + w = 0 \quad (1.20)$$

$$b - c'_u = 0 \quad (1.21)$$

$$(c'_o - B)c''_u \cdot p(o, u) + (b - B)(c'_o - b) = 0 \quad (1.22)$$

where (1.22) is a rewritten form of (1.8). If we differentiate these four equations with respect to L , we get:

$$p(o, u) \frac{db}{dL} + 1 \cdot \frac{dw}{dL} + (b - c'_o) \frac{do}{dL} + (b - c'_u) \frac{du}{dL} = 0 \quad (1.23)$$

$$1 + \frac{dw}{dL} = 0 \quad (1.24)$$

$$\frac{db}{dL} - c''_u \frac{du}{dL} = 0 \quad (1.25)$$

$$\begin{aligned} & ((c'_o - b) + (B - b)) \frac{db}{dL} + (c''_o c''_u p(o, u) + (c'_o - B)c''_u + (b - B)c''_o) \frac{do}{dL} \\ & + ((c'_o - B)c'''_u p(o, u) + (c'_o - B)c''_u) \frac{du}{dL} = 0 \end{aligned} \quad (1.26)$$

We can now solve (1.24) for $\frac{dw}{dL} = -1$ and (1.25) for $\frac{db}{dL} = c''_u \frac{du}{dL}$. Plugging these results into (1.23) and using (1.21) gives us

$$p(o, u) c''_u \frac{du}{dL} + (b - c'_o) \frac{do}{dL} = 1 \quad (1.27)$$

while plugging the results into (1.26) gives us:

$$\begin{aligned} & (c''_o c''_u p(o, u) + (c'_o - B)c''_u + (b - B)c''_o) \frac{do}{dL} \\ & + ((c'_o - B)c'''_u p(o, u) + 2(c'_o - b)c''_u) \frac{du}{dL} = 0 \end{aligned} \quad (1.28)$$

To simplify calculations, we make the following substitutions:

$$e = b - c'_o$$

$$f = p(o, u)c''_u$$

$$g = [c''_o c''_u p(o, u) + (c'_o - B)c''_u + (b - B)c''_o]$$

$$h = [(c'_o - B)c'''_u p(o, u) + 2(c'_o - b)c''_u]$$

The two equations can then be written as

$$\begin{bmatrix} e & f \\ g & h \end{bmatrix} \cdot \begin{bmatrix} \frac{do}{dL} \\ \frac{du}{dL} \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

Using Cramer's Rule we can solve for

$$\frac{do}{dL} = \frac{h}{eh - fg} \quad \text{and} \quad \frac{du}{dL} = \frac{-g}{eh - fg}.$$

Because at the optimum $c'_o > B > b$, we will have $e < 0$, $f > 0$ and $h > 0$. To sign g , we rewrite (1.22) and get

$$\frac{B - b}{c''_u p(o, u)} = \frac{c'_o - B}{c'_o - b} < 1$$

where the inequality follows again from $c'_o > B > b$. Because $c''_u p(o, u) > 0$, this implies $c''_u p(o, u) > B - b$. Now we can easily show $g > 0$. These results imply $eh - fg < 0$ and finally $\frac{do}{dL} < 0$, $\frac{du}{dL} > 0$. \square

Chapter 2

Regulation and Liability as Complements and Substitutes

2.1 Introduction

The relationship between liability and regulation as instruments for controlling hazardous technologies has been a long-standing question for Law & Economics scholars. The literature has mainly concentrated on two questions: first, whether – if one uses only one instrument – it is better to use regulation or liability, second, whether the joint use of regulation and liability is better than using only one of these instruments.

In this chapter, it will always be optimal to use both regulation and liability.¹ Our main interest is the stringency of regulation, namely whether the care demanded by regulation differs from the one that is implemented in the first-best. The main feature of our model – in comparison to most of the literature – is that the care that can be exercised to lower the probability of an accident is two-dimensional, where only one dimension of care can be subject to regulation.

For some action to be subject to regulation, it must fulfill the following requirements. First, the action must be describable in a regulation. Second, the regulator must be able to measure whether the firm has complied with the regulation

¹In the case where liability alone can implement the first-best, using regulation is superfluous but will also do no harm.

and third, a possible breach must be verifiable in a court of law.² Even if an action could be subject to regulation in principle, it might often not be practical to do so, because the regulator's cost of enforcing the regulation are prohibitively high. For example, a certain safety procedure might be easily describable, but might only be enforceable if every worker has a safety inspector at his side.³

In this chapter we consider a firm that operates a hazardous technology. To decrease the probability of an accident, the firm can take care in two dimensions. One dimension of care can be regulated and will be called "observable care". The other dimension of care cannot be regulated and will be called "unobservable care". An example of observable care is the compliance with a technical regulation that can easily be controlled by inspections. Examples of unobservable care are the workload and alertness of the operators, or whether management exerts pressure on staff to "bend the rules". We also assume that in case of an accident the legal system is not able to impose the full costs of the accident on the firm. So the use of liability alone will not induce the firm to take the optimal levels of observable and unobservable care.

In such a situation it will always be optimal to use both liability and regulation. Our main focus will be the optimal level of regulation. Existing regulations are often criticized for being excessively strict, meaning that their costs outweigh their benefits, and also for encouraging carelessness in matters not subject to regulation.

If both dimensions of care are independent, meaning that the level of care in one dimension does not influence the marginal cost or the marginal return of care in the other dimension, regulation should simply impose the first-best

²For example, experts analyzing disasters like the Chernobyl accident or the destruction of the Space Shuttle Columbia have stressed deficits in the "safety culture" of the responsible organizations (INSAG, 1992; CAIB, 2003). According to one definition, the safety culture of an organization is "the product of individual and group values, attitudes, perceptions, competencies and patterns of behavior that determine the commitment to, and the style and proficiency of, an organization's health and safety management" (ACSNI, 1993). Because "values", "attitudes" and "perceptions" are hard to measure, it seems difficult to enforce a regulation that simply prescribes a good safety culture.

³A solution to this problem are "stochastic inspections". But even an unannounced inspection might not uncover all day-to-day behavior.

level of observable care. Things become more complicated if we allow for interactions between the marginal returns of the two dimensions of care.⁴

For example, consider the case where the two dimensions of care are complements, meaning that the marginal return of one dimension of care increases in the level of the other dimension of care. Because the level of unobservable care will be too low compared to the first-best, a natural intuition in this case suggests that the regulated level of observable care should be higher than in the first best, to induce the firm to take more unobservable care. But there is a countervailing effect. Because the level of unobservable care is too low and both dimensions are complements, the marginal return of observable care will be lower than in the first-best, making a high level of observable care less beneficial. Which of the two effects dominates depends on the curvature of the cost curve for unobservable care. If the moral hazard problem is serious enough, we can show that regulated observable care will be below its first-best level.

The model in this chapter is most applicable to the question of avoiding environmental accidents, like major oil spills or the uncontrolled emission of radioactive or toxic material. The victims of those accidents will be numerous and widespread, making Coasian bargaining very costly. Furthermore, there are not many cost-effective avoidance activities which can be undertaken by victims, so we can ignore victims' incentives. But it should be noted that the idea of two-dimensional care should also be relevant to bilateral accidents, where victims' care decisions are important. Environmental accidents are also distinct from other environmental problems like the regulation of certain, observable effluents like carbon dioxide, which is at the center of current debates or problems where the emission is observable but uncertainty exists whether a certain damage has been caused by this emission.

A number of papers have modeled the problem of joint use of regulation and liability in a framework that was first developed by Shavell (1984). In this framework, care is one-dimensional and observable by the regulator, but regulation does not implement the social optimum because the firm has some private information. In Shavell's setup, this private information is about the potential

⁴To simplify, we do not consider interactions in the marginal costs.

harm that can be caused by an accident. But because liability is uncertain (underenforcement) and tortfeasor's assets are limited, liability alone does not implement the social optimum either. Shavell shows that it can be optimal to use both regulation and liability. In this case the optimal level of regulation is below the level that would be optimal if regulation is used alone. Schmitz (2000) shows that the optimality of joint use disappears if there is no underenforcement; then, either regulation or liability should be used. But joint use can become optimal again if there is heterogeneity in asset limitations. Rouillon (2008) extends these results to the case where agents differ in their probability of getting sued, while Hiriart et al. (2004) embed the problem in a contract-theoretic setting where the regulator gives the firm incentives to reveal its private information about potential harm. A different framework is used in Kolstad et al. (1990). In contrast to Shavell, they assume that the firm is only liable if it is found to have acted with negligence; regulation is welfare increasing because it reduces the firm's uncertainty about the legal standard that determines negligence.

There are some papers that model care as two-dimensional. In an article on liability for nuclear accidents, Trebilcock and Winter (1997) sketch a model with observable and unobservable care but do not fully solve it. In the setting analyzed by Bhole and Wagner (2008), a firm can take observable care as well as unobservable care to prevent an accident. They find that in many situations only the combined use of both liability and regulation will lead to optimal levels of effort in both dimensions. In contrast to us, Bhole and Wagner do only consider a binary choice of observable care; because in their model a high level of observable care is socially optimal, the question of excessive regulation is ruled out by assumption. In Hutchinson and van 't Veld (2005), unobservable care will lower the probability of an accident, while observable care will only lower the damage if an accident occurs. Finally, Bartsch (1997) considers two-dimensional care where one dimension is perfectly observable and the other only imperfectly and analyzes the firm's choice of care under a negligence regime.

The rest of the chapter is structured as follows. In section 2.2 we describe the setup of the model and distinguish the possible interactions between the two kinds of care. Section 2.3 characterizes the socially optimal care levels, while section 2.4 shows that liability alone will induce care levels that are too low, if the liability limit is below potential harm. The main part is section 2.5, where we describe the regulator's optimization problem under joint use of liability and regulation. In that section we also analyze whether observable effort will be above or below the first-best level, while section 2.6 concludes. All proofs can be found in the appendix.

2.2 Setup of the Model

There is a firm that uses a production process that could cause an accident. The firm can invest in two dimensions of care, observable care $o \in [0, o_{max})$ and unobservable care $u \in [0, u_{max})$ with $o_{max}, u_{max} < 1$. Define $X := [0, o_{max}) \times [0, u_{max})$.

The firm's costs for those investments in care are given by $c_o(o) + c_u(u)$. The probability that the accident is *avoided* depends on both levels of care and is given by $p(o, u)$. In the following we will sometimes call this probability the "level of safety". If an accident happens, the damage to society is D .

We make the following assumptions for $p(o, u)$ (subscripts denote partial derivatives):

Assumption 2.1. $0 \leq p(o, u) \leq 1$ for all $(o, u) \in X$, $p(0, 0) = 0$.

Assumption 2.2. $p_o, p_u > 0$.

Assumption 2.3. $p_{oo}, p_{uu} = 0$.

Assumption 2.4. $p_{ou} \in [-1, 1]$.

If $p_{ou} = 0$, we say that the two dimensions of care are independent, if $p_{ou} > 0$ we say that they are complements, and if $p_{ou} < 0$ we say that they are substitutes.

A functional form that fulfills the assumptions above is given by:

$$p(o, u) = \alpha_o o + \alpha_u u + \alpha_{ou} ou \quad (2.1)$$

with $0 \leq \alpha_o \leq 1$, $0 \leq \alpha_u \leq 1$, $\alpha_o + \alpha_u \leq 1$ and $|\alpha_{ou}| < \min\{\alpha_o, \alpha_u\}$.

We can interpret this functional form as follows: Because both o and u are between zero and unity, they can be interpreted as probabilities. We can associate each with some distinct piece of equipment or process involved in avoiding an accident. Then o and u are the probabilities that the respective process work properly, while $1 - o$ and $1 - u$ are the respective probabilities of failure. It is now easy to interpret three extreme cases:

independence contingent on the state of the world, one or the other process has to work to avoid the accident; the accident is avoided with probability $\alpha_o o + \alpha_u u$ with $\alpha_o, \alpha_u \geq 0$ and $\alpha_o + \alpha_u \leq 1$;

pure complements both processes are needed to avoid an accident; the accident is avoided with probability $o \cdot u$;

pure substitutes at least one process has to work to avoid the accident; the accident is avoided with probability $o + u - o \cdot u$.

The functional form given in (2.1) can be understood as a convex combination of the three functions, where the case of pure complements ($p(o, u) = o \cdot u$) is not allowed to happen.

We make the following technical assumptions for the cost functions:

Assumption 2.5. $c_o(o)$ and $c_u(u)$ are continuous, three times differentiable, strictly increasing and strictly convex.

Assumption 2.6. $c_o(o_{max}) = c_u(u_{max}) = \infty$.

Assumption 2.7. $c'_o(0) = c'_u(0) = 0$.

Assumption 2.8. $c_o(0) = c_u(0) = 0$.

These conditions are rather standard, e.g. Assumptions 2.6 and 2.7 ensure that the firm's problem has an interior solution. The next assumption is a little bit more intricate:

Assumption 2.9. For all $u \in [0, u_{max})$ we have $p_o c''_u(u) > |p_{ou}| \cdot c'_u(u)$. For all $o \in [0, o_{max})$ we have $p_u c''_o(o) > |p_{ou}| \cdot c'_o(o)$.

This assumption rules out cases where the complementary or substitutive nature of o and u is “too strong”. Note that this assumption implies that the function $\tilde{p}(k_o, k_u) := p(c_o^{-1}(k_o), c_u^{-1}(k_u))$ is strictly concave. Furthermore, it rules out the possibility that one dimension of care is an “inferior factor”, meaning that the socially optimal use of this factor decreases when D is increased.

2.3 The First-Best Care Levels

The socially optimal levels of o and u minimize the expected social costs of the hazardous technology, formally

$$\min_{(o,u) \in X} C_D(o, u) = (1 - p(o, u))D + c_o(o) + c_u(u). \quad (2.2)$$

Proposition 2.1. The problem given by (2.2) has an unique, interior solution, where the socially optimal care levels o^* and u^* are given by $p_o D = c'_o(o^*)$ and $p_u D = c'_u(u^*)$.

If we compare other outcomes with this first-best, we have to distinguish two concepts. On the one hand, we have the socially optimal *first-best* care levels o^* and u^* , which are given by $c'_o(o^*) = p_o D$ and $c'_u(u^*) = p_u D$. On the other hand, for a given level of safety \bar{p} , we can find the least expensive combination of observable and unobservable care that produces \bar{p} . Such a *cost-minimizing* combination of care will have $\frac{p_o}{p_u} = \frac{c'_o(o)}{c'_u(u)}$.⁵ It is easy to see that first-best care levels are also a cost-minimizing combination of care, but that there are also

⁵This condition results from $\min_{o,u} c_o(o) + c_u(u)$, subject to $p(o, u) = \bar{p}$. Formally, the marginal rate of technical substitution between these two kinds of care must be equal to the ratio of respective marginal costs.

many other cost-minimizing combinations of care that are not first-best. How can we interpret the cost-minimizing combination? Consider a safety specialist who is given the task of implementing a specific level of safety and is not aware of any incentive problem. He will recommend this combination of the two levels of care. If he finds – for example – that $\frac{c'_o(o)}{p_o} > \frac{c'_u(u)}{p_u}$ he will diagnose an over-reliance on the observable dimension of care.

2.4 Care Levels with Liability only

This section considers the outcome if the firm is only subject to strict liability. Strict liability means that in case of an accident the firm has to pay out compensation to the victim's of the accident, regardless of the level of care it has taken. The amount of compensation is given by L , with $0 \leq L \leq D$.

There exist a number of institutional reasons why L might be below damages D . First, there exists the possibility that if the accident occurs, the firm has insufficient funds to cover all damages.⁶ Second, often the legal system itself limits the amount of compensation. This might happen either explicitly, with numerical liability caps,⁷ or implicitly, if certain categories of damages, like “pain and suffering” are excluded or underestimated when computing damage amounts. Third, there might be some uncertainty whether in case of an accident the firm can be sued successfully (by all victims). Let q be the probability of a successful law-suit, then $L = qD$ can be interpreted as the expected damage payment.⁸ These considerations suggest another interpretation, namely that L is the amount of compensation that the firm has to pay in expectation,

⁶A regulator might take measures to prevent this from happening, like requiring insurance coverage, which we will not consider. But note that this leads to the question of moral hazard in the relation between firm and insurer. The insurer's problem might be very similar to regulator's problem considered here.

⁷For example in the US, the Oil Pollution Act of 1990 limits the liability for natural resource and economic damages to \$75 million per offshore oil spill; only direct cleanup cost are expected. However, there may be ways for plaintiffs to suspend those limits (see Richardson, 2010).

⁸Theoretically, this problem could be remedied by the use of “damage multipliers”.

conditional that an accident has occurred. Our model requires that this expectation does not depend on $p(o, u)$.⁹

The firm which is faced with strict liability and a liability limit L minimizes its expected cost of accidents:

$$\min_{(o,u) \in X} C_L(o, u) = (1 - p(o, u))L + c_o(o) + c_u(u). \quad (2.3)$$

Proposition 2.2. *For all $0 < L \leq D$, the problem given by (2.3) has an unique, interior solution, where the firm's privately optimal care levels \hat{o} and \hat{u} are given by $p_o L = c'_o(\hat{o})$ and $p_u L = c'_u(\hat{u})$.*

If $L = D$, strict liability will implement the first-best care levels, because then the firm's problem is identical to the government's. What if $L < D$? Unfortunately, this question cannot be answered by simply inspecting the first-order conditions. Consider $p_o L = c'_o(\hat{o})$; if L decreases, it will lower $p_o L$ directly but a change in u might also increase p_o if $p_{ou} \neq 0$. We have to do comparative statics with respect to L . Define $\hat{p}(L) := p(\hat{o}(L), \hat{u}(L))$, the privately optimal safety level as a function of the liability limit L .

Proposition 2.3. *The privately optimal safety level $\hat{p}(L)$ is strictly increasing in L .*

So a higher liability level causes a higher level of safety implemented by the firm.

Proposition 2.4. *For all $L < D$, the social costs of accidents $(1 - p(o, u))D + c_o(o) + c_u(u)$ are decreasing in L .*

So the higher liability level, the lower will be the social costs of the hazardous technology. The higher liability level induces more care by the firm, but because care levels are still lower than socially optimal, the savings in accident cost will outweigh the increased cost of care.

⁹In contrast, in Hutchinson and van 't Veld (2005) observable care reduces the amount of damage and therefore influences the amount of liability.

2.5 The Optimal Combination of Liability and Regulation

We now consider the problem of a regulator who has to choose the optimal level of regulation, given the firm's liability limit L . Regulation takes the form of postulating a minimum amount of observable care \underline{o} that the firm has to exert. Given this amount of regulated observable care and its liability L , the firm will choose the amount of unobservable care that minimizes its private costs. This restricted problem seems to be quite realistic for many regulatory agencies. Those agencies can do little about the general legal framework that governs liability but often have a great amount of discretion regarding the technical details of regulation.

Formally the regulator has to solve:

$$\begin{aligned} \min_{\underline{o} \in [0, o_{max}], (o, u) \in X} & (1 - p(o, u))D + c_o(o) + c_u(u) \\ \text{subject to: } & (o, u) \in \underset{(o, u) \in X}{\operatorname{argmin}} [1 - p(o, u)]L + c_o(o) + c_u(u) \\ & \text{subject to: } o \geq \underline{o}. \end{aligned} \quad (2.4)$$

As we show below, the firm's minimization problem has a unique solution that is characterized by the first-order condition $p_u L = c'_u(u)$. So we can replace the constraint above with this first order condition and use the so-called first-order approach.

Proposition 2.5. *For all $0 < L \leq D$, a solution to the problem given by (2.4) will have $o = \underline{o} \geq \hat{o}$ and $u > 0$. Such a solution is characterized by the first-order condition:*

$$c'_o(o) = p_o D + p_u (D - L) \frac{p_{ou} L}{c''_u(u)}.$$

The intuition for this result is as follows. If the government chooses a level of regulation below \hat{o} , the regulation is not a binding constraint on the firm and it will choose its privately optimal levels \hat{o} and \hat{u} . So we can assume that the

government sets regulatory level of at least \hat{o} . We can now show that the *firm's* maximization problem has a unique solution given by $p_u L - c'_u(u) = 0$. So we can define $u(o)$ as the firm's best response in unobservable care for every level of regulation $\underline{o} \geq \hat{o}$. If we differentiate the firm's first order condition with respect to o we get:

$$\frac{du}{do} = \frac{p_{ou}L}{c''_u(u)}.$$

Then a marginal increase in observable care changes the social costs of accidents as follows

$$\frac{dC_H}{do} = c'_o(o) - p_o D + (c'_u(u) - D p_u) \frac{du}{do} \quad (2.5)$$

$$= c'_o(o) - p_o D - p_u(D - L) \frac{p_{ou}L}{c''_u(u)} \quad (2.6)$$

This is negative at \hat{o} so it will always be worthwhile to increase observable effort above the level which is privately optimal for the firm. At the optimal point, the cost of additional effort will be just equal to its direct and indirect benefits.

If $p_{ou} = 0$, the regulator demands the first-best level of observable care while unobservable care is at the same level as without regulation. If $p_{ou} \neq 0$, inspecting the first order-conditions will not give us definitive information whether o and u are above or below the first-best values, because p_o and p_u will be different from their values at the first-best. We can only make the following local statements:

- because $p_u D > c'_u(u)$, for the given amount of o , the level of u is too low, so if the firm would voluntarily increase u marginally, social welfare would increase;
- if o and u are complements, we have $p_o D < c'_o(o)$, if they are substitutes, we have $p_o D > c'_o(o)$; so a safety expert who evaluates the regulation only with regard to its direct costs and benefits – excluding its effect on unobservable care – would conclude that the regulation is too strict in the case of complements and too lenient in the case of substitutes.

To get information on the absolute level of o and u we need to do a comparative statics analysis with regard to L . There are three countervailing effects that happen if L decreases. Let's consider the case of complements. The decrease in liability means that the firm has less powerful incentives to exert unobservable care, so u will fall. How should the regulator react to this development? Because the two dimension of care are complements, the regulator might consider increasing the level of o by regulation, thereby increasing the marginal return to u and the firm's incentives to take unobservable care. But the decrease in u will also lower the marginal return to o , which suggests to lower the level of regulation. Finally, the lower level of L will decrease the effectiveness of the channel whereby more o gives more incentives for u , which also suggests lowering regulation of o . The change in o if L decreases depends on the sum of these three effects, which depends on the curvature of the cost curve for unobservable care. Define

$$M(u) = 2 \frac{c_u''(u)}{c_u'(u)} - \frac{c_u'''(u)}{c_u''(u)}.$$

Proposition 2.6. *Assume that if problem (2.4) has more than one solution, the regulator chooses the one which has the lowest o . If L decreases, the change of o depends on the signs of p_{ou} and $M(u)$.*

1. *If o and u are complements ($p_{ou} > 0$), a decrease in L will cause o to decrease if $M(u) > 0$ and to increase if $M(u) < 0$.*
2. *If o and u are substitutes ($p_{ou} < 0$), a decrease in L will cause o to increase if $M(u) > 0$ and to decrease if $M(u) < 0$.*

To clarify this proposition, assume $M(u) > 0$ on $[0, u^*]$. Then if $L < D$ the level of regulation will be as follows:

- if o and u are **complements** ($p_{ou} > 0$), the level of regulated observable care will be **below** the first-best level;
- if o and u are **substitutes** ($p_{ou} < 0$), the level of regulated observable care will be **above** the first-best level.

Clearly the sign of $M(u)$ depends on the curvature of $c_u(u)$. But we can make more general statements than that. Because $c'_u(0) = 0$, we have $M(u) \rightarrow \infty$ if $u \rightarrow 0$. So we will have $M(u) > 0$ if u is low enough. We can also consider the situation at $L = 0$: in case of substitutes the optimal level of observable care will be higher than in the first-best and in case of complements it will be lower than in the first-best.¹⁰ This means that – in the case of substitutes – even if there exists a neighbourhood where a lower L causes a decrease in o , this decrease will be more than reversed as L goes to zero (an analogous argument applies to complements). This suggests that $M(u) > 0$ can be treated as the “normal” case.¹¹

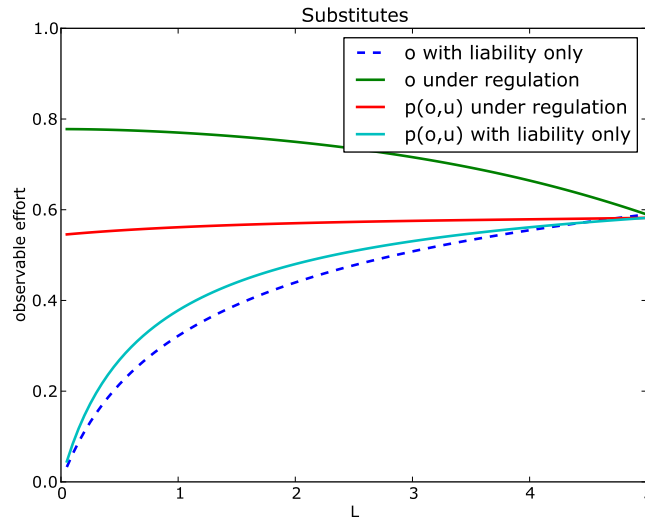


Figure 2.1: Observable Effort with Substitutes

To exemplify the interaction between the two dimensions of care, we show the result of numerical simulations.¹² Figures 2.1 and 2.2 both refer to “normal” case as defined above. They show the level of o and $p(o, u)$ under liability only

¹⁰In this situation, the regulator simply sets the optimal o given that u is zero.

¹¹A cost function that fulfills our assumption and has $M(u) > 0$ is given by $c_u(u) = \ln\left(\frac{1}{1-u}\right) - u$.

¹²We use the cost function $c(x) = \ln\left(\frac{1}{1-x}\right) - x$ for both $c_o(o)$ and $c_u(u)$. For the case of complements we use $p(o, u) = 0.1 \cdot o + 0.3 \cdot u + 0.3 \cdot o \cdot u$ and for the case of substitutes $p(o, u) = 0.7 \cdot (o + u - o \cdot u)$.

and under regulation and liability, for the case of substitutes and complements respectively.

In the case of substitutes (Figure 2.1), a lower liability limit L will not affect $p(o, u)$ too much because the decrease in u will be largely compensated by an increase in o (this increase will lead to a “second-round effect” where u will be decreased even more, and so on). So with low levels of L we will see a high level of regulation which will have crowded out much of the unobservable care. Critics may complain about this bureaucratic regime but given the low liability it is the best the regulator can do.

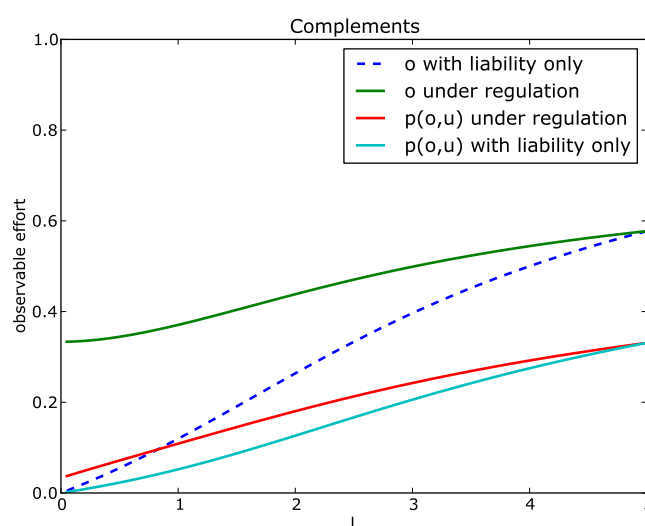


Figure 2.2: Observable Effort with Complements

In the case of complements (Figure 2.2), regulation will counteract a decrease in liability only to a small extent. If L is very low, unobservable care will be low but regulation will not be of much use (because observable care does not work well without unobservable care). It follows that the level of safety may become very low. In such a situation, it might make sense to explore other policy instruments, even if they are very costly, like trying to increase the liability limit or restricting the hazardous technology to wealthy individuals. In addition, given that the private cost of the activity will be much below its social cost, there will

be a tendency to use this activity above the socially optimal level. A regulator might therefore consider restricting the use of the activities to cases where it is of high social value;¹³ as an alternative, mandatory insurance might not raise the level of care but insurance premia will cause firms to internalize the expected damage costs and curtail excessive use of the activity.

2.6 Conclusion

In this chapter we have analyzed the optimal use of regulation and liability in a situation where only one dimension of care can be made subject to regulation. The optimal level of regulation in this dimension will only be equal to the first-best level if there is no interaction between the two dimensions of care. If the two kinds of care are substitutes or complements, the optimal level of regulation depends on the curvature of the cost curve for unobservable care. If the liability limit is low enough, it will always be the case that the optimal level of regulation will be above the first-best level for substitutes and below the first-best level for complements.

This result may appear somewhat counterintuitive if one considers only the incentive effect. But considered from the perspective of a regulator who “consumes” observable and unobservable care, it actually seems quite natural: if the regulator gets less unobservable care (because of lower liability) he wants to consume more observable care if this is a substitute and less observable care if this is a complement.

¹³An example for such a “needs test” can be found in the German law on gun ownership.

2.7 Appendix

Proof of Proposition 2.1

We first show that we can restrict the search for a minimum to

$$\tilde{X} = \{(o, u) \mid o, u \geq 0, c_o(o) + c_u(u) \leq D\},$$

which is a compact subset of X . For a choice of care levels (o, u) that is in X but not in \tilde{X} we will have

$$\begin{aligned} C_D(o, u) &= (1 - p(o, u))D + c_o(o) + c_u(u) \\ &> (1 - p(o, u))D + D \\ &> D. \end{aligned}$$

On the other hand, choosing $(o, u) = (0, 0)$ which is in \tilde{X} would give us social costs of D , so a minimum of social costs cannot be outside of \tilde{X} .

We note that $C_D(o, u)$ is continuous on \tilde{X} so the Weierstraß Theorem ensures that social costs will have a minimum on \tilde{X} . This shows the existence of a solution to problem (2.2).

We next show that we can rule out a solution on the border of \tilde{X} . A point on the border $c_o(o) + c_u(u) = D$ cannot be a minimum, because choosing $(o, u) = (0, 0)$ would be better (see above). We now consider a point on the border $o = 0$; at such a point, increasing o would change social costs by

$$\frac{dC_D}{do} = -p_o D + c'_o(0) < 0$$

(because $p_o > 0$ and $c'_o(0) = 0$). So increasing o would decrease social costs which rules out a minimum on this border. By an analogous argument we can rule out a minimum on the border $u = 0$. Therefore, the minimum must be in the interior of \tilde{X} and therefore in the interior of X .

As an interior minimum, the solution to problem (2.2) will be given by the first-order conditions

$$\frac{\partial C_D}{\partial o} = -p_o D + c'_o(o) = 0 \text{ and } \frac{\partial C_D}{\partial u} = -p_u D + c'_u(u) = 0$$

and because of because $D = \frac{c'_o(o)}{p_o} = \frac{c'_u(u)}{p_u}$ and Assumption 2.9 this minimum is unique. \square

Proof of Proposition 2.2

By analogy to the proof of Proposition 2.1. \square

Proof of Proposition 2.3

The variables \hat{o} , \hat{u} and \hat{p} are defined implicitly by the equations

$$\begin{aligned} p(o, u) - \hat{p} &= 0 \\ c'_o(o) - p_o L &= 0 \\ c'_u(u) - p_u L &= 0 \end{aligned}$$

Differentiating with respect to L gives us

$$\begin{aligned} p_o \frac{do}{dL} + p_u \frac{du}{dL} - \frac{d\hat{p}}{dL} &= 0 \\ c''_o(o) \frac{do}{dL} - L p_{ou} \frac{du}{dL} - p_o &= 0 \\ c''_u(u) \frac{du}{dL} - L p_{ou} \frac{do}{dL} - p_u &= 0 \end{aligned}$$

In matrix form:
$$\begin{pmatrix} -1 & p_o & p_u \\ 0 & c''_o(o) & -p_{ou}L \\ 0 & -p_{ou}L & c''_u(u) \end{pmatrix} \begin{pmatrix} \frac{d\hat{p}}{dL} \\ \frac{do}{dL} \\ \frac{du}{dL} \end{pmatrix} = \begin{pmatrix} 0 \\ p_o \\ p_u \end{pmatrix}$$

Using Cramer's rule:
$$\frac{d\hat{p}}{dL} = \frac{\begin{vmatrix} 0 & p_o & p_u \\ p_o & c''_o(o) & -p_{ou}L \\ p_u & -p_{ou}L & c''_u(u) \end{vmatrix}}{\begin{vmatrix} -1 & p_o & p_u \\ 0 & c''_o(o) & -p_{ou}L \\ 0 & -p_{ou}L & c''_u(u) \end{vmatrix}}.$$

The denominator has the opposite sign as $\begin{vmatrix} c''_o(o) & -p_{ou}L \\ -p_{ou}L & c''_u(u) \end{vmatrix}$ which is positive because $L = \frac{c'_o(o)}{p_o} = \frac{c'_u(u)}{p_u}$ and Assumption 2.9, so the denominator is negative. We can also conclude that $\begin{pmatrix} c''_o(o) & -p_{ou}L \\ -p_{ou}L & c''_u(u) \end{pmatrix}$ is positive definite, i.e. the quadratic form $x' \begin{pmatrix} c''_o(o) & -p_{ou}L \\ -p_{ou}L & c''_u(u) \end{pmatrix} x$ is positive for all $x \in \mathbb{R}^2$. So it must be positive for all x with $\begin{pmatrix} p_o \\ p_u \end{pmatrix}' x = 0$. But this implies that the bordered matrix given in the numerator is negative. Because both denominator and numerator are negative, $\frac{d\hat{p}}{dL}$ is positive. \square

Proof of Proposition 2.4

$$\begin{aligned}
\frac{dC_D}{dL} &= \frac{d}{dL} [D(1 - p(o, u)) + c_o(o) + c_u(u)] \\
&= -D \left(p_o \frac{do}{dL} + p_u \frac{du}{dL} \right) + c'_o(o) \frac{do}{dL} + c'_u(u) \frac{du}{dL} \\
&= -((D - L) + L) \left(p_o \frac{do}{dL} + p_u \frac{du}{dL} \right) + c'_o(o) \frac{do}{dL} + c'_u(u) \frac{du}{dL} \\
&= -(D - L) \left(p_o \frac{do}{dL} + p_u \frac{du}{dL} \right) + (c'_o(o) - p_o L) \frac{do}{dL} + (c'_u(u) - p_u L) \frac{du}{dL} \\
&= -(D - L) \frac{d\hat{p}}{dL} \\
&< 0
\end{aligned}$$

□

Proof of Proposition 2.5

If $D = L$, the regulator can achieve the first-best by setting $\underline{o} = \hat{o}$; this solution fulfills the first-order condition. In the following we assume $L < D$. Let again \hat{o} be the privately optimal observable care level at L . If the regulator sets $\underline{o} < \hat{o}$, the regulation is not binding and the firm will choose its privately optimal care levels (\hat{o}, \hat{u}) . So the social cost of accidents are equal for all $\underline{o} \leq \hat{o}$, so we can restrict the search for a maximum to values $o \geq \hat{o}$.

Given \underline{o} , the firm minimizes its private costs. The result is given by the first-order condition $-p_u L + c'_u(u) = 0$. Applying the implicit function theorem gives us:

$$\frac{du}{do} = -\frac{-p_{ou}L}{c''_u(u)} = \frac{p_{ou}L}{c''_u(u)}$$

So an additional marginal amount of o gives us :

$$\begin{aligned}
\frac{dC_D}{do} &= (-p_o D + c'_o(o)) + (-p_u D + c'_u(u)) \frac{du}{do} \\
&= (-p_o D + c'_o(o)) + (-p_u D + c'_u(u)) \frac{p_{ou}L}{c''_u(u)} \\
&= c'_o(o) - p_o D - p_u(D - L) \frac{p_{ou}L}{c''_u(u)}
\end{aligned}$$

If we evaluate $\frac{dC_H}{do}$ at the lower border \hat{o} , where $c'_o(o) = p_o L$, we get

$$\begin{aligned}\frac{dC_D(\hat{o})}{do} &= -p_o(D-L) - p_u(D-L) \frac{p_{ou}L}{c''_u(u)} \\ &= -(D-L) \left(p_o + p_u \frac{p_{ou}L}{c''_u(u)} \right) \\ &= -(D-L) \left(p_o + p_{ou} \frac{c'_u(u)}{c''_u(u)} \right)\end{aligned}$$

From Assumption 2.9 we know that $c''_u(u) + p_{ou} \frac{c'_u(u)}{p_o} > 0$, so $\frac{dC_D(\hat{o})}{do} < 0$. On the other hand, the optimal o must be smaller than \hat{o} , with $c_o(\hat{o}) = D$. This means that the government's problem has an interior solution in (\hat{o}, \bar{o}) , which will fulfill the first order condition:

$$c'_o(o) - p_o D = p_u(D-L) \frac{p_{ou}L}{c''_u(u)} \quad (2.7)$$

□

Proof of Proposition 2.6

Define $\left. \frac{du}{dL} \right|_{o \text{ fixed}}$ as the change in u induced by a change in L while holding o fixed. Using the implicit function theorem on $p_u L - c'_u(u) = 0$ we get:

$$\left. \frac{du}{dL} \right|_{o \text{ fixed}} = \frac{p_u}{c''_u(u)}$$

Now we consider how $\frac{dC_D}{do}$ (marginal social costs) change, if L increases.¹⁴

$$\begin{aligned}\frac{\partial}{\partial L} \frac{dC_D}{do} &= p_u \frac{p_{ou}L}{c''_u(u)} - p_u(D-L) \frac{p_{ou}}{c''_u(u)} + \left(-p_{ou}D - p_u(D-L) \frac{p_{ou}L}{(c''_u(u))^2} (-1) c'''_u(u) \right) \left. \frac{du}{dL} \right|_{o \text{ fixed}} \\ &= p_u \frac{p_{ou}L}{c''_u(u)} - p_u(D-L) \frac{p_{ou}}{c''_u(u)} + \left(-p_{ou}D + p_u(D-L) \frac{p_{ou}L}{(c''_u(u))^2} c'''_u(u) \right) \frac{p_u}{c''_u(u)} \\ &= \frac{p_u p_{ou}}{c''_u(u)} \left((L-D) - (D-L) + (D-L) \frac{c'_u(u) c'''_u(u)}{(c''_u(u))^2} \right) \\ &= (-p_{ou}) \frac{p_u}{c''_u(u)} (D-L) \left(2 - \frac{c'_u(u) c'''_u(u)}{(c''_u(u))^2} \right)\end{aligned}$$

Consider the case with $M(u) > 0$; this implies $2 - \frac{c'_u(u) c'''_u(u)}{(c''_u(u))^2} > 0$. Now the sign of the change in marginal social costs depends on the sign of p_{ou} . If $p_{ou} > 0$, marginal social costs decrease with L so the optimal o will be higher. If $p_{ou} < 0$, marginal social cost increase with L so the optimal o will be lower. □

¹⁴We use the method of increasing marginal returns, see Edlin and Shannon (1998).

Chapter 3

Minimum Education Requirements for Professions

3.1 Introduction

Many occupations cannot be undertaken by everybody; the classical “learned professions” like law or medicine have been restricted to persons with a prescribed education for a long time. Today, this policy has been extended to many other occupations. Recently, Kleiner and Krueger (2009) found that 29 % of the US workforce is required to have a government issued license to do their job; for 85 % of those jobs, a specific exam was necessary to get this licence. The defenders of such “minimum education requirements” usually argue that without such barriers to entry, the quality of professional services will deteriorate because consumers are not able to judge this quality at the time of consumption.

Minimum education requirements have been an object of economists’ critiques for a long time. The main objection seems to be that they restrict consumers’ choices, forcing them to buy high-quality services even if they would have preferred cheaper low-quality services.¹ As an alternative, some critics suggest, the government should restrict itself to certification of education levels,

¹The classical critique of minimum education requirements can be found in Friedman (1962), ch. 9. Recently, this question has come up in European competition policy, comp. Com-

allowing consumers to buy the quality level of their own choice. Other critics do not advocate abolishing minimum education requirements, but think that existing requirements are often excessively high and should be lowered.²

But there exists a problem with these requirements that is even more fundamental. They regulate only an input for the production of professional services, while a regulator will care about the quality of the output. So if unregulated professionals will produce low quality output, it needs to be explained why an input regulation will improve the situation. A well-educated professional might still produce low quality if he does not pay attention to his work or accepts too many clients, thereby “spreading himself too thin”. So it seems that the proper solution to the problem of professional quality does not lie in minimum education requirements, but in the direct regulation of output quality. If such a regulation is enforcing high quality, professionals have an incentive to voluntarily acquire the optimal amount of education, so minimum education requirements seem to be superfluous.

This chapter suggests a reason why minimum education requirements might still be necessary in addition to a direct regulation of output quality. In our model, professional education serves as an “hostage” that makes direct quality regulation enforceable by granting professionals a quasi-rent. We can show that minimum education requirements will sometimes implement the first-best. In other cases, they might still be preferable to the alternate policy of granting the professionals a pure rent by restricting entry to the profession by numerical limits.

We assume that quality is not observable during purchase and delivery of the service, but that low quality will – with some probability – lead to a bad outcome at a later point of time. If such an outcome is observed, the professional will be excluded from the profession. Minimum education requirements make professionals more sensitive to this punishment, because exclusion devalues their occupation-specific human capital by taking away the associated quasi-

mission of the European Communities (2004b) and Commission of the European Communities (2005a).

²Leland (1979) deals with the optimal level of regulation in a model where quality can be directly regulated.

rent. In addition, high education makes it easier to produce high quality. In other words, education that is specific to the profession is both an “hostage” to ensure high quality and makes such quality easier to accomplish.

We make the assumption that exclusion is the only possible punishment. While this is clearly unrealistic, in practice punishments will be restricted by wealth constraints and the need for marginal deterrence.³ So one might interpret our assumption in the sense that other punishments are not sufficient to ensure high quality.

The idea of using education as an hostage to ensure good behavior has been mentioned in the literature, for example by Svorny (1987) and Shapiro (1986), but – to our knowledge – has never been formally modeled. On the other hand, it is a well known property of moral hazard that “limited liability” of the agent often makes it optimal to grant him a rent. A famous example is found in Shapiro and Stiglitz (1984) with regard to the relation between workers and employers; indeed, their “efficiency wages” are in many ways comparable to the quantitative entry restrictions that we will consider as an alternative to education requirements. But in our case, the welfare loss is not due to unemployment but to increases in the price of the professional service. Scoppa (1997) considers the use of firm specific human capital as an “hostage” in a moral hazard model between workers and employers and compares it with the use of efficiency wages. But, in contrast to us, he concentrates on the problem of self-enforcing contracts, and does not deal with the question of the efficient amount of investment in human capital,⁴ which is central for us.

There exists a parallel literature that deals with the problem of ensuring quality in competitive markets, where this quality cannot be contracted upon, but customers learn about the quality provided by individual firms in the past. Klein and Leffler (1981) have pointed out the paradox that a firm will only provide high quality if it is earning a rent which it wants to protect, but that the existence of such a rent seems to be incompatible with free entry. They suggest a

³Marginal deterrence refers to the problem that in order to efficiently deter very harmful acts there might be upper bounds for the punishment of less harmful acts.

⁴In his model, human capital investment is assumed to be either efficient or inefficient on the relevant range.

number of possible solutions but do not model them formally.⁵ Shapiro (1983) addresses this paradox with a model in which firms invest in the asset of their reputation. The subsequent work of Shapiro (1986) is, in some ways, most closely related to ours. Like this study, he models professional quality provision as a moral hazard problem and interprets education requirements as an input regulation. But in contrast to us, customers *can* observe quality, but only after purchasing the good, so producers can acquire a reputation for good quality. Education requirements are imposed to make it more attractive for producers to earn a reputation for high-quality work. Shapiro shows that under certain parameters education requirements can be welfare improving compared to a policy of *laissez-faire*.

In contrast to Shapiro, Svorny (1987) stresses that the threat of losing the future returns to their professional education gives professionals an incentive for good behavior, which is a central idea of this chapter.⁶ We extend her work by formally modeling this idea; furthermore, we explicitly differentiate between “pure” entry restrictions, which only provide a stream of rents to practitioners, and education requirements, which do provide social benefits by decreasing the cost of providing good quality.

Another related paper is Donabedian (1995), which studies professional self-regulation. He stresses that the effectiveness of such regulation depends on the “exit costs” of leaving the profession; these costs are given by the rents and quasi-rents earned by practitioners. But in comparison to us, the human capital investment of the professional is not made explicit, while it is central to our approach. On the other hand, we abstract from choice between professional self-regulation and state regulation, which is central to Donabedian.

The rest of the chapter proceeds as follows: While section 3.2 describes the decision problem of the individual professional, section 3.3 derives the market equilibrium without regulation, which results in low quality work. Possible

⁵The solution nearest to ours are “nonsalvagable productive assets”.

⁶But she cannot find empirical evidence that stricter licensing requirements for physicians (in 1965) were quality enhancing. In a later paper (Svorny, 1992) she argues that subsequent recent changes in the US health care market have reduced the need for licensing and its incentive effects.

regulatory inventions, namely minimum education requirements and quantitative entry restrictions are discussed in section 3.4. In the following welfare analysis of section 3.5, we demonstrate the condition for minimum education requirements to achieve the first-best solution and compare minimum education requirements and quantitative entry restrictions in cases where the first best is not achievable. Section 3.6 concludes with a discussion of the policy relevance of our results.

3.2 Objective Function of the Agent

In the following we will call work producing high quality “good work” and work producing low quality “bad work”. Each agent faces two decisions: the kind of work he produces and level of investment in human capital. An agent produces one unit of a service by either doing good work or bad work, $w \in \{0, 1\}$, where $w = 1$ means good work. The agent receives price p when he sells the service on the market.

While the cost of bad work is zero, the cost of good work is given by $c(k) > 0$, where k is the investment in human capital. We assume that $c(k)$ is strictly decreasing in k , so more human capital makes good work easier. We further assume that $c(k)$ is continuously differentiable, has a second derivative with $c''(k) > 0$ for all $k \geq 0$. To ensure an interior solution we assume $\lim_{k \rightarrow 0} c'(k) < -r$ and $\lim_{k \rightarrow \infty} c'(k) = 0$.

There are two complementary interpretations for the assumption that the cost of good work decreases with k . First, greater knowledge about their field makes it easier for professionals to avoid errors which might cause a bad outcome. Second, professional education does also include indoctrination into the “professional” way of doing work, including professional ethics.⁷ Professionals who have gone through this kind of education may do good work unthinkingly or

⁷This indoctrination may happen through formal education but also – in practical training – through students observing and imitating practitioners.

may even experience psychological costs when doing bad work that counter-vail the greater exertion needed for good work.⁸

The agent is risk-neutral and maximizes expected present value of utility with discount rate $r \in (0, 1)$. The timing is as follows (see Figure 3.1): Time is discrete and goes from zero to infinity. In period $t = 0$ an agent decides how much to invest in human capital. In periods $t \in \{1, 2, 3, \dots\}$ the agent does either good or bad work. If an agent does bad work in some period t , a bad outcome occurs with probability ϕ after the period. If a bad outcome occurs, the bad work of the agent will be revealed, and the agent is prohibited from offering the service in the future and earns zero outside the profession. (Bad work in one period cannot lead to a bad outcome in a later period.)

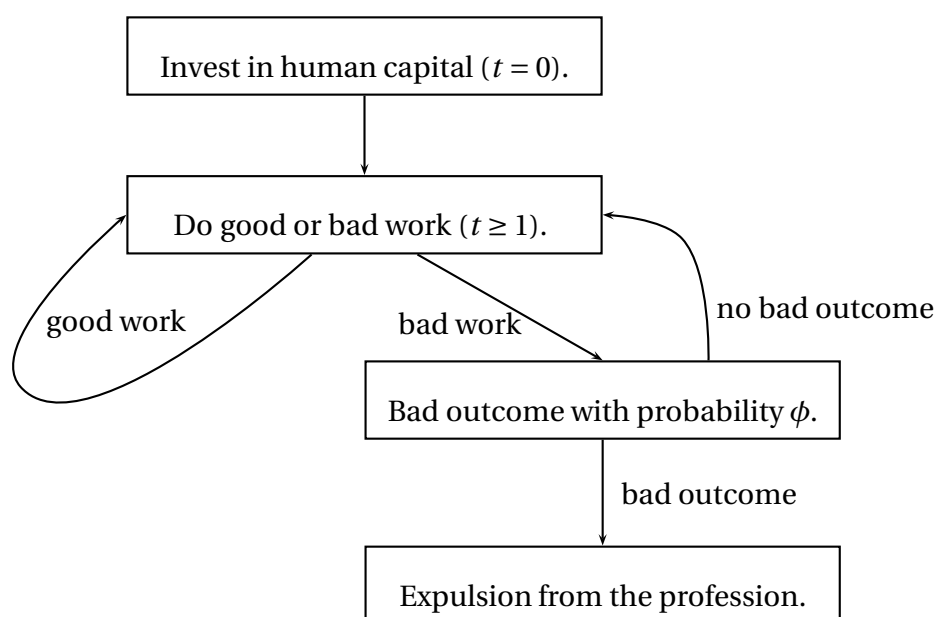


Figure 3.1: Timing of the model

⁸An explicit model with such psychological costs can be found in Akerlof and Kranton (2005).

The agent's payoff depends on his kind of work w , his human capital k , and the market price p . For an agent, in every period t greater than zero, the present value of being in the profession $V(p, k)$ is governed by the recursive equation:

$$V(p, k) = \max_w p - c(k) + \left(\frac{1}{1+r} \right) (wV(p, k) + (1-w)(1-\phi)V(p, k))$$

The value of being in the profession is this period's earnings plus the discounted *expected* value of being in the profession in the next period. Because this value does not depend on the number of the period, the optimal decision on $w \in \{0, 1\}$ must be the same for every period.

So the discounted sum of future earnings for good work is given by:

$$V_1(p, k) = p - c(k) + \left(\frac{1}{1+r} \right) V_1(p, k).$$

And the discounted sum of future earnings for bad work is given by:

$$V_0(p, k) = p - 0 + \left(\frac{1}{1+r} \right) (1-\phi)V_0(p, k).$$

If we solve these equations we get:

$$V_1(p, k) = \left(\frac{1+r}{r} \right) (p - c(k)) \text{ and } V_0(p, k) = \left(\frac{1+r}{r+\phi} \right) \cdot p.$$

The total payoff of an agent – taking account of the initial investment in human capital – is given by:

$$\Pi(w, p, k) = \left(\frac{1}{1+r} \right) V_w(p, k) - k.$$

3.3 Market Equilibrium

There exists an unlimited mass of identical agents willing to enter the profession. The market demand for the service is given by $D_e(p)$, where e is the quality of work expected by costumers. We make the usual assumptions $D_e(p) \geq 0$ and

$D'_e(p) \leq 0$; furthermore, $e' > e$ implies $D_{e'}(p) \geq D_e(p)$. We assume that agents have rational expectations. They only enter the profession if they expect to earn non-negative profits at the equilibrium price. Furthermore, we assume them to be price-takers.⁹ If some agents are forced to leave the profession, they are replaced with new agents. Note that those agents face the same decision problem as agents that have been there from the beginning.

A market equilibrium is characterized by values \hat{k}_0 , \hat{k}_1 , \hat{p} , \hat{w} , that satisfy the following conditions (all agents choose the same kind of work and the same level of investment):

1. Profit maximization, i.e. agents choose the optimal capital level (\hat{k}_0 or \hat{k}_1) for the kind of work (good/bad) they do:

$$\text{for both } w \in \{0, 1\}, \hat{k}_w = \underset{k}{\operatorname{argmax}} \Pi(w, \hat{p}, k).$$

2. Incentive compatibility, i.e. agents choose the kind of work that gives them a higher payoff:

$$\Pi(\hat{w}, \hat{p}, \hat{k}_{\hat{w}}) \geq \Pi(\hat{w} - 1, \hat{p}, \hat{k}_{\hat{w}-1}).$$

3. Free entry, i.e. agents earn expected profits of zero (given their investment in education):

$$\Pi(\hat{w}, \hat{p}, \hat{k}_{\hat{w}}) = 0.$$

Condition 1 ensures that all agents choose the optimal amount of human capital, given the expected market price and the kind of work they plan to do. Condition 2 makes sure that agents prefer to do work of quality \hat{w} , given the market price and the amount of human capital that is optimal for each kind of work.¹⁰ Condition 3 represents free entry – as long as positive profits can be earned, further agents will enter the market.

⁹They also take as given the expectation of customers about quality of work.

¹⁰Additional constraints that ensure that the chosen kind of work is optimal under a sub-optimal choice of capital are superfluous.

We begin by determining the optimal choice of human capital. This depends on the intended kind of work. If the agent plans to do bad work, his objective function is

$$\Pi(0, p, k) = \left(\frac{1}{1+r} \right) V_0(p, k) - k = \frac{p}{r+\phi} - k$$

and his optimal human capital is obviously zero, so $\hat{k}_0 = 0$.

The optimal \hat{k}_1 maximizes

$$\Pi(1, p, k) = \left(\frac{1}{1+r} \right) V_1(\hat{p}, k) - k = \frac{p-c(k)}{r} - k.$$

It is characterized by: $-c'(\hat{k}_1) = r$ (marginal savings from human capital investment equal the demanded return to capital). Our assumptions on $c(k)$ make sure that there is an interior solution.

Suppose, agents prefer to do good work. This means, the incentive constraint must hold, i.e.

$$\Pi(1, \hat{p}, \hat{k}_1) \geq \Pi(0, \hat{p}, \hat{k}_0)$$

that is

$$\left(\frac{1}{1+r} \right) V_1(\hat{p}, \hat{k}_1) - \hat{k}_1 \geq \left(\frac{1}{1+r} \right) V_0(\hat{p}, \hat{k}_0) - 0$$

Which is equivalent to:

$$\hat{p} \geq \left(1 + \frac{r}{\phi} \right) (c(\hat{k}_1) + r\hat{k}_1). \quad (3.1)$$

On the other hand, we have the zero profit condition:

$$\Pi(1, \hat{p}, \hat{k}_1) = \frac{1}{1+r} V(1, p) - \hat{k}_1 = 0$$

which implies

$$\hat{p} = c(\hat{k}_1) + r\hat{k}_1. \quad (3.2)$$

But because $\left(1 + \frac{r}{\phi}\right) > 1$, the incentive constraint and zero profit condition cannot be fulfilled at the same time. On the other hand, for $\hat{w} = 0$ and $\hat{p} = 0$ all three equilibrium conditions are fulfilled. So we get the following result:

Proposition 3.1. *There is only one market equilibrium, where $\hat{w} = 0$, $\hat{p} = 0$, $\hat{k}_0 = 0$ and $\hat{k}_1 > 0$.*

So without regulation, agents will acquire zero human capital and produce bad work. The intuition for this is easy. Define $C = c(k_1) + rk_1$ as the long-run costs of doing good work (i.e. the price that ensures zero profits). As long as the price is above $\left(1 + \frac{r}{\phi}\right)C$, it is optimal for agents to acquire human capital before entering the market and as long as p does not sink below C , it is profitable to do so. But entering *without* human capital is profitable for any positive price, and with $p \leq \left(1 + \frac{r}{\phi}\right)C$ it is optimal to do so. Therefore, agents expect that entry will drive down prices to zero. At this price, entry with human capital is no longer profitable, so only agents without human capital will enter in the first place.

This situation is demonstrated in Figure 3.2, where Π_0 is the expected payoff for $w = 0$ and $k = \hat{k}_0 = 0$ and Π_1 the expected payoff for $w = 1$ and $k = \hat{k}_1$. At the point where good work becomes unprofitable – the Π_1 line crosses the zero-line – bad work (the Π_0 line) will still give a positive payoff.

3.4 Regulatory Interventions

A regulator may intervene in this situation in several possible ways. He can restrict entry so that prices reach a level that makes it optimal to produce good work. He could also impose a level of minimum human capital to achieve the same aim.

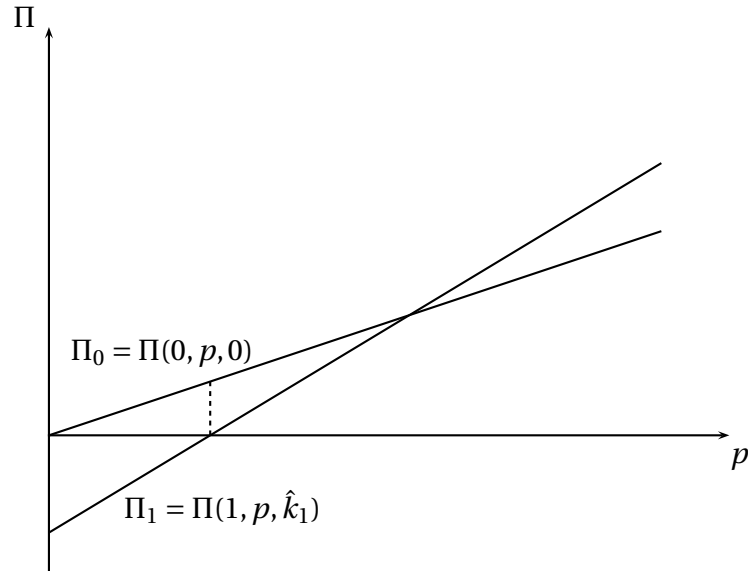


Figure 3.2: Payoffs without regulation

The regulator can simply restrict the entry to the profession so that prices remain high enough for agents to prefer good work over bad work, i.e. set entry numbers so that:

$$\bar{p} \geq \left(1 + \frac{r}{\phi}\right) (c(\hat{k}_1) + r\hat{k}_1).$$

Then it is optimal to acquire human capital and produce good work.

Proposition 3.2. *The regulator can ensure good work by restricting entry to the profession at a level so that*

$$\bar{p} = \left(1 + \frac{r}{\phi}\right) \cdot C.$$

The intuition behind this result is that if entry restrictions set prices to the high level of \bar{p} , it is very lucrative to be in the profession, so agents do not want to risk expulsion by doing bad work.

The other alternative is to set a minimal level of education k^* that agents must acquire before they enter the market; the level has to be set so that it is optimal to provide good work (the incentive constraint must hold) and must be consistent with free entry, which implies zero profits.

In this case the regulator solves the problem

$\min_{p, k^*} k^*$, subject to

$$\Pi(1, p, k^*) \geq \Pi(0, p, k^*)$$

and:

$$\Pi(1, p, k^*) = \frac{1}{1+r} V(p, k^*) - k^* = 0 \quad (\text{zero profits})$$

The two constraints imply:

$$\frac{p - c(k^*)}{r} \geq \frac{p}{r + \phi} \quad \text{and} \quad \frac{p - c(k^*)}{r} - k^* = 0$$

which boils down to:

$$\phi k^* \geq c(k^*).$$

The minimal k^* that fulfills this condition is given by: $\phi k^* = c(k^*)$. Because $c(k)$ is decreasing, such a k^* exists, which establishes the following result:

Proposition 3.3. *The regulator can ensure good work by setting a minimum human capital level of k^* , with $\phi k^* \geq c(k^*)$.*

The intuition behind this result is that agents do not want to risk leaving the profession, because they would lose their “quasi-rents”; doing bad work is unattractive because agents have to invest in “useless” education first. Note that under entry restrictions, agents earn economic rents. Under minimum education requirements, they earn only “quasi-rents” – the market price will be above short-run cost, but entry will drive it down to long-run cost.

Figure 3.3 illustrates the situation under minimum education requirements. The Π_0 -line has been shifted down because all agents are required to invest in human capital, even if it is useless for them because they want to provide bad work. (In this example Π_1 -line is unchanged but this need not be the case.)

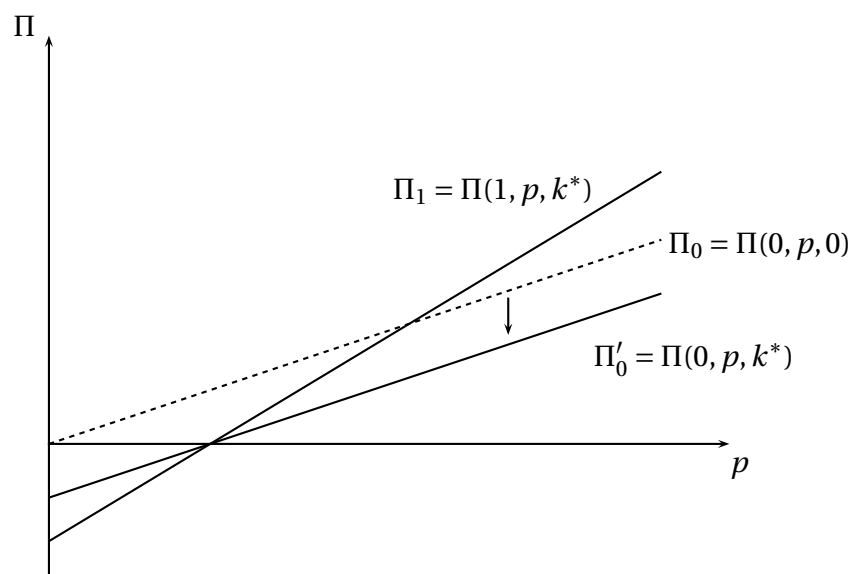


Figure 3.3: Payoffs with minimum human capital

3.5 Welfare Analysis

In the following we make a welfare comparison of entry restrictions and minimum human capital requirements. We assume that it is optimal to ensure high quality, and only investigate into the least costly way to do so.

It is *not* obvious which welfare standard is most appropriate to compare the two policies. The standard approach in industrial organization is to use total welfare (consumer surplus plus producer surplus). But one could argue that using only consumer surplus is more appropriate here. If the regulator imposes a numerical limit of places in the profession and this leads to the professionals

earning economic rents, applicants may spend resources to get one of those places. Because of this “rent dissipation” (Posner, 1975), a great part of the producer surplus might get lost. Nevertheless, the total welfare seems appropriate to define the first-best.

Proposition 3.4. *Under the total welfare standard the first-best-outcome is:*

$$p = C = c(\hat{k}_1) + r\hat{k}_1 \quad (\text{long-run costs of good work})$$

$$k = \hat{k}_1 \quad (\text{optimal education decision}).$$

Because it is optimal to ensure high quality (by assumption), the first-best w is good work while the optimal human capital level is given by cost minimization. The service should be provided to all consumer willing to pay at least the long-run cost of its provision. There are two relevant kinds of distortion from the first-best. The price of the service might be higher than the long run costs of good work and there might be excessive investment in education.

3.5.1 Achieving the First-Best

Under certain parameters, minimum education requirements can achieve good work costlessly. If the parameters are such that $\hat{k}_1 \geq k^*$, the socially optimal human capital is greater than the prescribed minimum human capital. Entrants to the profession, given that they are forced to acquire at least a capital level of k^* , will in their own interest choose the higher and socially optimal capital level of \hat{k}_1 . So the minimum education requirement does not distort the human capital decision and free entry will drive down prices to long-run costs. The following results shows when this is the case:

Proposition 3.5. *If $c(\hat{k}_1) \leq \phi\hat{k}_1$, the first-best can be implemented with education requirements.*

Proof: From $c(\hat{k}_1) \leq \phi\hat{k}_1$ follows $\phi\hat{k}_1 - c(\hat{k}_1) \geq 0$. Because $c(k)$ is decreasing, $\phi k - c(k)$ is strictly increasing in k . We know that at k^* we have $\phi k^* - c(k^*) = 0$. So we must have $\hat{k} \geq k^*$. This implies that the first-best is achievable. \square

In comparison, entry restrictions will never implement the first-best because they imply prices that are higher than long-run costs.

This can also be seen graphically. As long as the imposed minimum capital level is below the optimal capital level for good work, changes in minimum capital only affect the Π_0 -line (because for agents doing good work, the minimum education requirement is not binding). In these cases, if minimum human capital increases, only the Π_0 -line shifts down. Indeed, Figure 3.3 depicts the situation where education requirements implement the first-best, because only the Π_0 -line gets shifted down.

If $c(\hat{k}_1) > \phi \hat{k}_1$, it will not be sufficient to shift down only the Π_0 -line. Because in this case, when the imposed minimum human capital reaches the optimal human capital, it is still preferable to do bad work. And if the imposed minimum capital level exceeds the optimal capital level for good work, changes in minimum human capital affect both lines. Both the Π_0 -line and the Π_1 -line are shifted down, but the shift is greater for the Π_0 -line, because increases in human capital even above the optimal level do still decrease the costs of good work but are useless if the agent does bad work. This is the situation depicted in Figure 3.4.

To see how the condition $\hat{k}_1 \geq k^*$ depends on the parameters r and ϕ , we make a comparative statics analysis. Using the first order condition for profit maximizing human capital, $-c'(\hat{k}_1) = r$, we get:

$$\frac{d\hat{k}_1}{dr} = -\frac{1}{c''(\hat{k})} < 0$$

$$\frac{d\hat{k}_1}{d\phi} = 0.$$

And with the condition for the optimal minimum education requirement,

$k^* = \frac{c(k^*)}{\phi}$, we get

$$\frac{dk^*}{dr} = 0$$

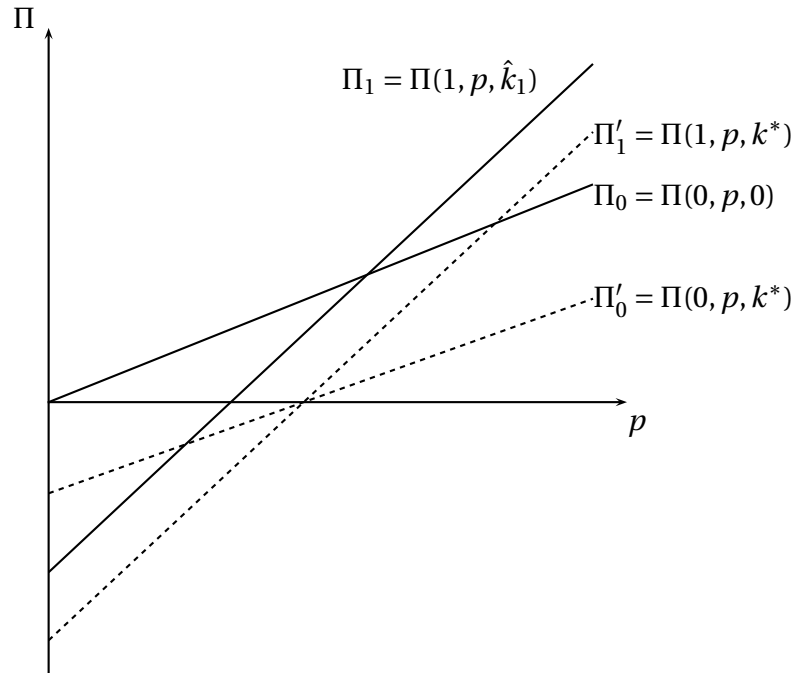


Figure 3.4: Payoffs with minimum human capital (not costless)

$$\frac{dk^*}{d\phi} = \frac{k^*}{c'(k^*) - \phi} < 0.$$

With these results we can see that the condition for the first-best being achievable, $\hat{k}_1 \geq k^*$, is more likely to be fulfilled the greater ϕ (higher probability of a bad outcome) and the lower r (greater patience). This can be clearly seen in the example where we have the cost function $c(k) = \frac{a}{k}$. In this case, the first-best can be implemented if $r \leq \phi$.

3.5.2 Comparisons in the Second Best

In cases where the first-best can not be achieved, i.e. $c(\hat{k}_1) > \phi \hat{k}_1$, the results depend on the welfare standard employed. If we use the consumer-surplus-standard, welfare is a decreasing function of price only. It can be shown that

education requirements imply a lower price than entry restrictions, so we can state:

Proposition 3.6. *If $c(\hat{k}_1) > \phi\hat{k}_1$, and we use the consumer-surplus-standard, minimum education requirements are at least weakly preferred to entry restrictions.*

(Proof in the Appendix.)

If we use the welfare standard of total surplus, things are no longer so clear. Because in this case the social welfare function does also include the profits earned by the professionals, so social welfare is no longer a function of price only.

In the first-best, the profits earned by professionals are zero, so social welfare is given by the area under the demand curve, i.e. $\int_{\hat{p}}^{\infty} D_1(\tilde{p})d\tilde{p}$. Under minimum education requirements, the profits are zero again, so social welfare is again given by consumer demand, albeit at the higher price of p^* , i.e. $\int_{p^*}^{\infty} D_1(\tilde{p})d\tilde{p}$. But under entry restrictions, social welfare consists of consumer surplus plus the profits earned by professionals. The profit per professional is \bar{p} minus the cost of good work under optimal capital, which is equal to \hat{p} . So social welfare is given by:

$$\int_{\bar{p}}^{\infty} D_1(\tilde{p})d\tilde{p} + (\bar{p} - \hat{p})D_1(\bar{p})$$

To compare the two policies, it is convenient to look at the respective welfare losses relative to the first-best case. Minimum education requirements are better, if their welfare loss is smaller than the welfare loss of entry restrictions, i.e.:

$$\int_{\hat{p}}^{p^*} D_1(\tilde{p})d\tilde{p} < \int_{\hat{p}}^{\bar{p}} D_1(\tilde{p})d\tilde{p} - (\bar{p} - \hat{p})D_1(\bar{p})$$

To make further comparison, we need to be more specific about the demand function. It is fairly intuitive that education requirements are better if demand is very elastic while entry restrictions are preferred if demand is very inelastic. This is best seen in extreme cases.

- Demand is totally inelastic between \hat{p} and \bar{p} . This means that $D_1(\tilde{p})$ is constant and equal to $D_1(\bar{p})$ and the welfare loss of entry restrictions is $(\bar{p} - \hat{p})D_1(\bar{p}) - (\bar{p} - \hat{p})D_1(\bar{p}) = 0$. On the other hand, the welfare loss for education requirements is $D_1(\hat{p})(p^* - \hat{p})$, which is strictly positive as long as $D_1(\hat{p}) > 0$. Thus entry restrictions are at least weakly preferred.
- Demand is very elastic between p^* and \bar{p} , i.e. $D_1(p^*) > 0$ and $D_1(\bar{p}) = 0$. In this case, the welfare loss of entry restrictions is $\int_{\hat{p}}^{\bar{p}} D_1(\tilde{p})d\tilde{p}$. The difference between this and the welfare loss of education requirements is $\int_{p^*}^{\bar{p}} D_1(\tilde{p})d\tilde{p}$ which is greater equal zero if $D_1(p^*) > 0$. Education requirements are strictly preferred as long as $D_1(\tilde{p}) > 0$ in $[p^*, p^* + \epsilon]$ for some $\epsilon > 0$.

The relation can also be easily be shown in the case of linear demand. If we assume $D_1(p) = 1 - bp$, for every p demand is more elastic the greater b . For simplicity, we assume that the parameters are such that $D_1(\hat{p}) \geq 0$. In this case, it can be shown that education requirements are preferred if

$$b > \frac{p^* - \hat{p}}{\frac{1}{2}\bar{p}^2 - \bar{p}\hat{p} + \frac{1}{2}p^{*2}}$$

(Derivation in the Appendix.)

3.6 Conclusion

In this chapter we have shown that minimum education requirements can serve as an instrument to make quality regulation enforceable. The profession-specific human capital that agents acquire to fulfill the requirement will provide them with a quasi-rent. The fear of losing this quasi-rent will make professionals reluctant to disobey the quality regulation. If the human capital level necessary to create this quasi-rent is lower than the level which is optimal to produce good work, minimum education requirements can implement the first-best. If the optimal level of human capital is too low to achieve the first-best, minimum education requirements can still be superior to numerical entry limits for the

profession, especially if the objective function of the regulator puts little weight on professionals' profits or the demand for the service is very elastic.

With regard to policy conclusions, our aim in this chapter is rather modest: to show that minimum education requirements can be a feasible instrument for ensuring high quality, but not whether or to what extent their use is optimal. Nevertheless, the level of real-world minimum education requirements should not only be evaluated with regard to the knowledge which is necessary to undertake a certain occupation but also with regard to the need to protect a profession from "fly by night operators" who have "nothing to lose". Occupations which require little physical but much human capital seem therefore most appropriate for minimum education requirements. They are very susceptible to opportunistic entry but minimum education requirements will not distort the optimal capital level too much.

Our models assumes that bad work is only detected at a very late stage, when a bad outcome has realised. Customers are assumed to have no other way of learning about the quality of work. So our model seems appropriate for situations where the bad effects of low quality services are realised with low probability or only after some time has passed. Otherwise, the professional's fear for his reputation might serve as an adequate instrument to ensure high quality. Our model seems also appropriate for situations where the professional's direct costumers do not fully internalize the quality of work.¹¹ This can be the case if the quality of the work creates an externality. For example the quality of a financial audit does not only affect the audited company – the direct customer of the accountant – but also outside financial investors. In other professions, it might be the professional's duty to protect customers from their own self-destructive impulses. For example, in many countries certain medications can only be bought if prescribed by a physician or a nurse. The aim of this rules is to guard against addiction to medical drugs. A pharmacist who refuses to sell a drug without the necessary prescription will produce "high quality" but might

¹¹In most cases those professionals are "gatekeepers" (Kraakman, 1986) in the sense that their cooperation is necessary to engage in some other regulated activity. For example, a firm who wants to raise capital on the stock market needs the cooperation of an auditor who certifies its accounts.

create a dissatisfied customer. In such a situation, reputation might not only be ineffective but could actually be counterproductive: a pharmacist who develops a reputation for bending the rules might be able to acquire additional customers.¹²

In cases where minimum education requirements do not achieve the first-best, one may ask whether there are other institutions which can enforce quality regulations at lower costs. A possible policy would be for the regulator to require the cost-minimizing level of human capital and demand an additional monetary deposit by the aspiring professional. If a bad outcome occurs, this deposit will be seized by the regulator, giving the professional additional incentives for good work. This policy seems very attractive, because “posting a bond” is a purely financial transaction that will not use up real resources. But a full evaluation should also include negative selection effects that occur if we restrict the profession to those who can finance this initial monetary investment. Of course the requirement of an additional investment in human capital will also exclude some applicants. But there may be significant differences between the class of applicants who can easily invest in additional human capital (because of high ability) and the class of people who can easily finance the deposit.

¹²The same concern applies to physician’s prescription behaviour. Svorny (1992) points out that licensing boards in the US focus their enforcement on physicians who prescribe narcotics inappropriately or who have a drug/alcohol problem of their own.

3.7 Appendix

3.7.1 Proof of Proposition 3.6

Under the consumer surplus standard the regulator's objective function is given by:

$$CS(p) = \int_p^{\infty} D_1(\tilde{p}) d\tilde{p}$$

The effect of price changes on consumer surplus is given by:

$$\frac{dCS}{dp} = -D_1(p)$$

which is negative if $D_1(p) > 0$. So consumer surplus is a decreasing function of price, except for prices so high that demand is zero.

That means that under the consumer surplus a policy is at least weakly preferable if it implies a lower price.

The price is greater under entry restrictions if:

$$\left(1 + \frac{r}{\phi}\right) (c(\hat{k}) + r\hat{k}) > c(k^*) + rk^*$$

Define

$$F(r) = \left(1 + \frac{r}{\phi}\right) [c(\hat{k}(r)) + r\hat{k}(r)] - [c(k^*) + rk^*]$$

with

$$\hat{k}(r) = \underset{k}{\operatorname{argmin}} c(k) + rk$$

The implicit function theorem implies that $\hat{k}(r)$ exists and is differentiable.

We have to show that $F(r) > 0$.

$$\begin{aligned}\frac{dF(r)}{dr} &= \frac{1}{\phi}[c(\hat{k}(r)) + r\hat{k}(r)] + \left(1 + \frac{r}{\phi}\right)\hat{k}(r) - k^* \\ &= \frac{1}{\phi}c(\hat{k}(\tilde{r})) - k^* + \left(1 + \frac{2r}{\phi}\right)\hat{k}(r)\end{aligned}$$

(Because of the envelope theorem we can ignore the effect of r on \hat{k}).

Now choose r^* so that $\hat{k}(r^*) = k^*$. Then

$$F(r) = F(r^*) + \int_{r^*}^r \frac{dF(\tilde{r})}{d\tilde{r}} d\tilde{r}$$

Note that r^* must be smaller than r , because we have $k^* > \hat{k}$ and $\hat{k}(r)$ is decreasing in r .

1. Step: We show $F(r^*) > 0$. Because $\hat{k}(r^*) = k^*$, we have

$$F(r^*) = \left(1 + \frac{r}{\phi}\right)[c(k^*) + rk^*] - [c(k^*) + rk^*] = \frac{r}{\phi}[c(k^*) + rk^*] > 0$$

2. Step: We show $\int_{r^*}^r \frac{dF(\tilde{r})}{d\tilde{r}} d\tilde{r} \geq 0$. This will be the case if $\frac{dF(\tilde{r})}{d\tilde{r}} d\tilde{r} > 0$ (monotonicity of the integral). Because the other terms in $\frac{dF(\tilde{r})}{d\tilde{r}}$ are all strictly positive, this will be the case if $\frac{1}{\phi}c(\hat{k}(\tilde{r})) - k^* \geq 0$. But as we can see:

$$\begin{aligned}\tilde{r} &\geq r^* \\ \Rightarrow \hat{k}(\tilde{r}) &\leq \hat{k}(r^*) = k^* \\ \Rightarrow c(\hat{k}(\tilde{r})) &\geq c(k^*) \\ \Rightarrow \frac{1}{\phi}c(\hat{k}(\tilde{r})) &\geq k^* \quad \text{because with } 0 < \phi < 1 \text{ we have } \frac{1}{\phi} > 1 \\ \Rightarrow \frac{1}{\phi}c(\hat{k}(\tilde{r})) - k^* &\geq 0\end{aligned}$$

□

3.7.2 Derivation of equation 3.5.2

$$[p - \frac{1}{2}bp^2]_{p^*}^{\bar{p}} - (\bar{p} - \hat{p})(1 - b\bar{p}) > 0$$

$$\bar{p} - \frac{1}{2}b\bar{p}^2 - p^* + \frac{1}{2}bp^{*2} - (\bar{p} - \hat{p}) + b\bar{p}(\bar{p} - \hat{p}) > 0$$

$$\frac{1}{2}bp^{*2} - \frac{1}{2}b\bar{p}^2 + b\bar{p}(\bar{p} - \hat{p}) > p^* - \bar{p} + (\bar{p} - \hat{p})$$

$$b > \frac{p^* - \hat{p}}{\bar{p}(\bar{p} - \hat{p}) - \frac{1}{2}\bar{p}^2 + \frac{1}{2}p^{*2}}$$

$$b > \frac{p^* - \hat{p}}{\frac{1}{2}\bar{p}^2 - \bar{p}\hat{p} + \frac{1}{2}p^{*2}}$$

3.7.3 Example with cost function $c(k) = \frac{a}{k}$

In this case, \hat{k}_1 is given by

$$-\left(-\frac{a}{\hat{k}_1^2}\right) = r$$

so

$$\hat{k}_1 = \sqrt{\frac{a}{r}}$$

and k^* is given by

$$\phi k^* = \frac{a}{k^*}$$

so

$$k^* = \sqrt{\frac{a}{\phi}}$$

for $k^* \leq \hat{k}_1$ to be true, we need

$$\sqrt{\frac{a}{\phi}} \leq \sqrt{\frac{a}{r}}$$

So this boils down to $r \leq \phi$.

Bibliography

- ACSNI (1993): *Third report: Organising for safety. Advisory Committee on the Safety of Nuclear Installations*, ACSNI Study Group on Human Factors.
- AKERLOF, G. A. AND R. E. KRANTON (2005): "Identity and the Economics of Organizations," *The Journal of Economic Perspectives*, 19, pp. 9–32.
- BARTSCH, E. (1997): "Environmental liability, imperfect information, and multidimensional pollution control," *International Review of Law and Economics*, 17, 139–146.
- BHOLE, B. AND J. WAGNER (2008): "The joint use of regulation and strict liability with multidimensional care and uncertain conviction," *International Review of Law and Economics*, 28, 123–132.
- BUTTLE, F. (1997): "ISO 9000: marketing motivations and benefits," *International Journal of Quality & Reliability Management*, 14, 936–947.
- CAIB (2003): *Report, Volume I*, Columbia Accident Investigation Board.
- Commission of the European Communities (2004b): *Report on Competition in Professional Services*, Commission of the European Communities.
- (2005a): *Professional Services – Scope for more reform*, Commission of the European Communities.
- DONABEDIAN, B. (1995): "Self-Regulation and the Enforcement of Professional Codes," *Public Choice*, 85, 107–18.

- EDLIN, A. S. AND C. SHANNON (1998): "Strict Monotonicity in Comparative Statics," *Journal of Economic Theory*, 81, 201 – 219.
- FRIEDMAN, M. (1962): *Capitalism and Freedom*, The University of Chicago Press.
- HERAS SAIZARBITORIA, I., G. ARANA LANDÍN, AND M. CASADESÚS FA (2006): "A Delphi study on motivation for ISO 9000 and EFQM," *International Journal of Quality & Reliability Management*, 23, 807–827.
- HIRIART, Y., D. MARTIMORT, AND J. POUYET (2004): "On the optimal use of ex ante regulation and ex post liability," *Economics Letters*, 84, 231–235.
- HOLMSTROM, B. AND P. MILGROM (1991): "Multitask Principal-Agent Analyses: Incentive Contracts, Asset Ownership, and Job Design," *Journal of Law, Economics and Organization*, 7, 24–52.
- HUTCHINSON, E. AND K. VAN 'T VELD (2005): "Extended liability for environmental accidents: what you see is what you get," *Journal of Environmental Economics and Management*, 49, 157–173.
- INSAG (1992): *The Chernobyl accident: updating of INSAG-1: report by the International Nuclear Safety Advisory Group*, International Atomic Energy Agency.
- KLEIN, B. AND K. B. LEFFLER (1981): "The Role of Market Forces in Assuring Contractual Performance," *Journal of Political Economy*, 89, 615–41.
- KLEINER, M. M. AND A. B. KRUEGER (2009): "Analyzing the Extent and Influence of Occupational Licensing on the Labor Market," Working Paper 14979, National Bureau of Economic Research.
- KOLSTAD, C. D., T. S. ULEN, AND G. V. JOHNSON (1990): "Ex Post Liability for Harm vs. Ex Ante Safety Regulation: Substitutes or Complements?" *American Economic Review*, 80, 888–901.
- KRAAKMAN, R. H. (1986): "Gatekeepers: The Anatomy of a Third-Party Enforcement Strategy," *Journal of Law, Economics and Organization*, 2, 53–104.

- KRÄKEL, M. AND A. SCHÖTTNER (2010): "Minimum wages and excessive effort supply," *Economics Letters*, 108, 341 – 344.
- LAUX, C. (2001): "Limited-Liability and Incentive Contracting with Multiple Projects," *RAND Journal of Economics*, 32, 514–26.
- LELAND, H. E. (1979): "Quacks, Lemons, and Licensing: A Theory of Minimum Quality Standards," *Journal of Political Economy*, 87, 1328–46.
- NAGATANI, K. (1978): "Substitution and Scale Effects in Factor Demands," *Canadian Journal of Economics*, 11, 521–27.
- POSNER, R. A. (1975): "The Social Costs of Monopoly and Regulation," *Journal of Political Economy*, 83, 807–27.
- RICHARDSON, N. (2010): "Deepwater Horizon and the Patchwork of Oil Spill Liability Law," Background, Resources for the Future.
- ROUILLON, S. (2008): "Safety regulation vs. liability with heterogeneous probabilities of suit," *International Review of Law and Economics*, 28, 133–139.
- SCHMITZ, P. W. (2000): "On the joint use of liability and safety regulation," *International Review of Law and Economics*, 20, 371–382.
- SCOPPA, V. (1997): "Firm Specific Investment in Human Capital as an Enforcement Mechanism Alternative to Efficiency Wages," Quaderni del Dipartimento di Economia Politica 221, Università degli Studi di Siena.
- SHAPIRO, C. (1983): "Premiums for High Quality Products as Returns to Reputations," *The Quarterly Journal of Economics*, 98, 659–79.
- (1986): "Investment, Moral Hazard, and Occupational Licensing," *Review of Economic Studies*, 53, 843–62.
- SHAPIRO, C. AND J. E. STIGLITZ (1984): "Equilibrium Unemployment as a Worker Discipline Device," *The American Economic Review*, 74, 433–444.
- SHAVELL, S. (1980): "Strict Liability versus Negligence," *The Journal of Legal Studies*, 9, pp. 1–25.

- (1984): “A Model of the Optimal Use of Liability and Safety Regulation,” *RAND Journal of Economics*, 15, 271–280.
- STRAUSZ, R. (2006): “Buried in paperwork: Excessive reporting in organizations,” *Journal of Economic Behavior & Organization*, 60, 460–470.
- SVORNY, S. V. (1987): “Physician Licensure: A New Approach to Examining the Role of Professional Interests,” *Economic Inquiry*, 25, 497–509.
- (1992): “Should We Reconsider Licensing Physicians?” *Contemporary Policy Issues*, 10, 31–38.
- TREBILCOCK, M. AND R. A. WINTER (1997): “The economics of nuclear accident law,” *International Review of Law and Economics*, 17, 215–243.

Eidesstattliche Versicherung

Ich versichere hiermit eidesstattlich, dass ich die vorliegende Arbeit selbständig und ohne fremde Hilfe verfasst habe. Die aus fremden Quellen direkt oder indirekt übernommenen Gedanken sowie mir gegebene Anregungen sind als solche kenntlich gemacht. Die Arbeit wurde bisher keiner anderen Prüfungsbehörde vorgelegt und auch noch nicht veröffentlicht.

Datum: 17.12.2010

Felix Reinshagen

Lebenslauf

- seit 2010 wissenschaftlicher Mitarbeiter
Seminar für Wirtschaftstheorie,
Prof. Klaus M. Schmidt, LMU München
- 2009 - 2010 Lehrassistent
Volkswirtschaftliche Fakultät der LMU München
- 2006 - 2010 Doktorandenstudium der Volkswirtschaftslehre
Munich Graduate School of Economics, LMU München
- 2003 - 2006 wissenschaftlicher Mitarbeiter
Geschäftsstelle der Monopolkommission, Bonn
- 2002 - 2003 Aufbaustudium der Volkswirtschaftslehre
Abschluss: Postgraduate Diploma in Economics and Econometrics
University of Essex, Colchester (UK)
- 2002 Zweites Juristisches Staatsexamen in Potsdam
- 2000 Erstes Juristisches Staatsexamen in München
- 1994 - 2000 Studium der Rechtswissenschaft
in Konstanz, Cardiff (UK) und München
- 1994 Abitur am Michaeli-Gymnasium, München
- 1. Mai 1975 Geboren in München