

Minimum Payments, Incentives, and Markets*

Ohad Kadan[†] and Jeroen Swinkels[‡]
First Draft, October 2009

September 7, 2010

Abstract

We consider the effect of an increase in the minimum feasible payment (driven by a minimum wage, limited liability, social norms or a legal restriction that pay cannot be negative) in markets with performance pay. We build a model in which a principal can adjust both the incentives of agents *and* the number of agents it employs. While a principal with a fixed number of agents will typically reduce the amount of effort it chooses to implement when the minimum payment increases, a principal that can also adjust its employment level will often instead *increase* the effort it induces. When embedded in a simple market setting, a typical outcome is that a moderate minimum payment makes all participants in the economy, *including agents who keep their job*, worse off. On the other hand, a considerable minimum payment, one that makes the issue of employee retention irrelevant, can help agents who keep their jobs.

1 Introduction

Many employees are subject to a minimum wage, but also receive incentive pay. Waiters fall in this category, as do many retail employees and sales people.¹ The design of incentive contracts for much more highly compensated workers often also have relevant minimum feasible payments (caused by limited liability or by legal constraint). The smallest payment (largest penalty) in a contract between a firm and a supplier is similarly limited by what courts will enforce, and by the net worth of the supplier.

*We thank Phil Reny, Charlie Brown, Sandra Black, Toomas Hinnosaar, Justin Ng, and Jonathan Lhost for helpful comments. We also thank seminar audiences at UT Austin, Northwestern, Princeton, and Michigan.

[†]Olin Business School, Washington University in St. Louis

[‡]Kellogg School of Management, Northwestern University

¹The current minimum wage in the U.S. is \$7.25/hour. This minimum is different for some sectors than others. However, even in sectors such as the restaurant industry, in which the mandatory cash minimum wage is, as of 2009, \$2.13, the employer must assure that the combination of bases plus tip income exceeds \$7.25/hour. Under U.S. law, a salesperson can be paid “straight commission” only if they are employed primarily offsite, and satisfy a number of other criteria. And, even here, a legal minimum of zero applies: a salesperson cannot be charged for the opportunity to approach a given set of targets, although that can easily be the contract that would prevail absent legal restrictions.

What happens in such settings when the minimum feasible payment changes? How do firms choose to change incentives in response to a change in the minimum wage?² Do the resultant incentives lead workers to work harder or less hard than before? What happens to employment? Are agents, even if they keep their job, better off? And, along the same line, does one ask more, or less, of a supplier with shallower pockets from which to pay penalties?

To explore this, we extend the standard moral hazard problem in two important ways. First, in the real world, Nordstroms can choose both how motivated each one of their sales people is *and* how many of them walk the floor. A firm can divide work across multiple suppliers. So, we build a model in which the principal can adjust both the number of agents it has, and how hard they work.

Second, in the standard model, the outside option of the agent is exogenous, there is no question of whether or not the principal chooses to operate, and the benefit to the principal of effort is exogenously determined. In reality, all of these things are determined as part of a market outcome. So, we embed our principal-agent model in the simplest possible market setting that allows us to endogenize these key features.

We next discuss each of these in a little more depth.

1.1 A Moral Hazard Model with Multiple Agents

Consider a principal that can employ n identical agents, each of whom in equilibrium exerts some effort level e , and that the principal cares about the total effort exerted, $E = ne$. The principal sees a signal about the effort of any given agent independently across agents, and rewards each agent on the basis of this signal. For the moment, take the outside option of the agent as exogenous.

We distinguish between the case in which n is fixed and given exogenously, and the case in which the principal can adjust both n and e . The first case may correspond to a market in which adjustment costs of labor are very high, resulting in a rigid labor market (as is the case in some European countries). Alternatively, the first case with $n = 1$ may correspond to the CEO job, where employing more than one is not an option. The second case corresponds to markets with relatively low adjustment costs for labor.

If n is exogenous, then the heart of the analysis from the single-agent case carries through. In particular, Kadan and Swinkels (2009a,2009b,2010) show for the one agent per principal case that as the minimum allowable payment goes up, the effort the principal wishes to implement generally goes down. Intuitively, increasing the minimum payment reduces the ability of the principal to punish low outcomes. But, eliciting effort by rewarding high outcomes is particularly costly when the effort to be induced is high, and so a higher minimum payment m leads to a higher marginal cost of inducing effort.³

²Even a restaurant can change the share of tips that gets shared among the kitchen staff, with the chef, or with the bartender, and can change who pays for uniforms.

³Formally, if C is the minimized cost of inducing effort, the set of papers referenced explore when $C_{me} > 0$. Conditions for a theoretical result are non-trivial. The referenced papers present a variety of sets of such conditions, and supplement them by numerical exploration.

In the standard principal agent problem (with one or a fixed number of workers), one can decouple the problem into first solving for the least cost way of inducing any given effort level e per agent, and only then looking for the optimal e . When both n and e can be optimally chosen by the principal, a very useful further decoupling of the problem occurs: at an optimum, the principal is indifferent between eliciting one more unit of effort from existing agents and hiring one more agent at the existing effort level. Thus, having found the least cost way of inducing any given e , the principal can then find the e^* that minimizes the average cost a per unit of effort induced. Effectively, e^* is the “minimum efficient scale” for any given agent, and, because the principal can work on both margins, the marginal cost of inducing total effort E is nothing but the average cost per unit of effort a . Having solved for a , the principal chooses E to maximize profits given an effective marginal cost of effort a . Thus, it does not matter whether the word processing pool of a firm is typing up the notes of autobody claims adjusters or medical researchers, the right thing to do is to begin by minimizing the cost per page accurately typed.

Because of this decoupling, it is immediate that when m rises, the optimal E falls. This is so simply because when m rises a constraint is being tightened, and so the minimized average cost a rises as well. Hence, while in the single agent model, signing the change in effort is hard, here it is easy.

Total effort can fall by more than one path. One way is to have as many (or more) agents, but have them move a little slower. Certainly this is what our earlier results suggest: since it is more expensive to motivate the agent, it makes sense to ask a little less of them. But, because the total cost of employing the agent goes up in m as well, there is also a force in the direction of inducing higher effort: if both the average total cost and marginal cost of a production technology are raised, the effect on minimum efficient scale depends on the relative magnitudes of the effects. So, another possibility is that agents are given incentives to work harder than ever, but there are sufficiently fewer of them that the overall effort level still falls. One way or the other, service will be worse at your favorite restaurant.

We next turn to exploring the relative magnitudes of the average and marginal cost effects. We have the surprising result that for a substantial set of cases, the average total cost effect dominates the marginal cost effect, and the firm chooses to induce *higher* effort from its employees when m rises: E falls, but n falls even faster. The result is clearest when the combination of the minimum payment and performance incentives is high enough that worker retention ceases to be relevant (i.e., when the individual rationality constraint, IR , has ceased to bind), but holds beyond this case as well. Thus, even though agents are homogeneous, and so minimum payments have no selection effects, an increase in the minimum payment, by causing the principal to optimally face agents with more intense incentives, leads to a measured productivity increase.⁴ We also identify cases where the marginal cost effect dominates the average total cost effect, so that minimum efficient scale falls. In such a setting, whether n

⁴This is not the substitution effect. The principal is not shifting away from paying agents to some factor the price of which has not changed, but rather, re-optimizing given that both the cost of having an agent around and the cost of inducing effort at the margin has risen.

falls depends on how quickly E falls compared to e , and n could even rise if the firm has very inflexible demand for E .

These results are phrased not in terms of the primitives of the problem, but rather in terms of how the induced distribution over payments changes as the firm implements, and the agent accepts, different effort levels. We discuss the (obvious) disadvantages and (considerable) advantages of this approach in Section 7.2. At the heart is a requirement that as a worker works harder, he is having more of a positive effect on his pay at higher levels than at low. We show that numerical simulation is quite supportive of the sensibility of this condition. When the condition is satisfied, the question of what happens to minimum efficient scale depends on a curvature condition on the utility function.

The question of what happens to minimum efficient scale when the minimum payment changes is very tightly connected to the question of what happens to minimum efficient scale as the outside option changes (whether or not there is a relevant minimum payment). We show for a large set of cases that as the outside option rises, so does the optimal induced effort, and of course, the utility of the agent within the relationship. We thus have the phenomenon of highly compensated, high outside option, but very hard-working employees, as is typical of many successful professionals (see Aguiar and Hurst (2007) and Gicheva (2009) for supporting evidence).

1.2 The Moral Hazard Model in a Market Context

Thus far, we have taken the benefit to the principal of effort, the choice of the principal to operate, and the outside option of the agents as exogenous. But, for many important questions, this is clearly inadequate. To address this, we next embed the moral hazard problem in the simplest possible market framework. We consider this step of our analysis as being most sensible when thinking about principals as firms and agents as workers, and so adopt this lexicon. Firms convert E into homogeneous output through a concave production function.

Workers are identical with respect to their productivity and signal process, but have heterogeneous outside options. Labor supply thus is given by an exogenous upward sloping supply function relating the number of workers willing to work to the utility being offered. The price of output is determined competitively, and depends on an exogenously modeled demand curve, and on the hiring and production decisions of firms. Firms have heterogeneous fixed costs.

Recall that in the case of an exogenous outside option and market price, effort per employee typically falls when the number of agents n per firm is fixed, but often rises when n is variable. Interestingly, once equilibrium effects are accounted for, much of this distinction disappears: if the optimal effort level per worker e^* rises with the minimum wage m for firms that can adjust their output, then the market outcome, even for firms that cannot adjust n , will typically be one with lower output, but higher output per employee. Employment thus takes a double hit in either case.

The mechanism, however, is different. When firms can adjust the workforce, an increase in the minimum wage leads to lower output per firm produced by a smaller number of workers each of whom works harder than before. Thus, lower employment

is a result of firms actively reducing their workforce, and of decreased overall output at the new higher market price. When the number of workers per firm is rigid, workers again work harder, but the increase in the cost of inducing effort increases the average costs of production more than the associated increase in equilibrium price. Consequently, the number of active firms in the sector declines and employment falls. Thus, while firms cannot adjust their workforce, price adjustments in the product market lead to a similar outcome in terms of employment.

The fact that labor demand falls with m creates a strong suggestion in our model that moderate minimum wages in incentive-pay driven sectors hurt everyone, *including those who keep their jobs*, and that successive increases in m cause further harm. This is so because at moderate minimum wages, IR will bind – the utility of workers in the performance pay sector will be just high enough to attract enough workers into the sector. This required utility falls when the sector shrinks, and firms optimally re-adjust the entire incentive scheme to take advantage of this.⁵

If m is high enough, then it may be that the combination of m and performance incentives is more than enough to attract employees into the sector (the firm may not wish to lower compensation, as the only way to do that also lowers effort). Employment takes the same double hit as before. But, because IR no longer governs, under reasonable conditions, workers who keep their jobs have higher utility even though (in fact, partly because) they are working harder.⁶

Worker utility is thus typically U shaped in m , with the trough at the point where IR just ceases to bind. Only if m is substantially higher than this point is there any hope that it leaves even those who keep their jobs better off. So, a low minimum wage hurts everyone, while a sufficiently high minimum wage may help those within the performance pay sector, but only at the cost of workers in the other sectors, and of consumers as well. In contrast, note that helping workers to lift themselves into better careers aids them directly, and, by removing workers from the sector, improves the outside option of the marginal worker remaining in the sector, and hence the utility of all workers in the sector.

This discussion is, at its heart, demand and supply. A minimum wage is designed to move the market away from the intersection of demand and supply curves. In the standard analysis, this helps those who keep their jobs, but hurts other market participants. Here, the employer has an easy tool to grab any utility above that required to bring enough workers into the sector, and, when IR is binding, will optimally do so.

Our results have implications far beyond, say, the service sector. In many markets,

⁵It may at first seem counter-intuitive that the firm can move to paying less bonus compensation, but elicit higher effort. One way to achieve this in many settings is to reward “decent” performance less, but “excellent” performance more. Under, for example, *MLRP*, this can simultaneously lower the agent’s expected utility from bonus payments, but increase the marginal incentive to exert effort.

⁶Note that while it can be that workers receive more than their outside option, this paper is not part of the efficiency-wage literature begun by Shapiro and Stiglitz (1984). In particular, we do not consider firing as part of the incentive scheme (doing so would clearly be an interesting extension). Rather, firms provide total compensation beyond that needed to attract workers because of the combination of a desire to provide incentives and the requirement to pay the minimum wage.

with or without explicit bonuses, changes in job title are part of the motivational scheme, and so a change in the minimum wage can require adjustment of pay for a considerable range of higher paid positions. If starting pay on a given career track is the minimum wage, then arguably every job on the normal progression of that career track is affected by the minimum wage, and our results, while for a static model, are suggestive of how wages along that progression should change. In particular, our results suggest that a firm faced with a higher minimum wage will find it more expensive to induce effort, and so will use less effort in total across its workforce, but will, in many circumstances, choose to induce more effort for each worker it keeps. If the *IR* constraint binds, the effect of this on total labor demand can result in a situation where, as viewed from the beginning of their career path, worker utility is lower.

Second, a very real “minimum wage” is extant in most managerial situations, and, indeed for many professional workers with strong incentive pay: it is essentially impossible to pay less than zero, no matter how disastrous the observable outcome associated with that employee.⁷ In fact, the base salary of CEOs sets a lower bound on their pay. While potentially, such base salary can be adjusted by the firm, in practice base salaries are quite rigid, generating a de-facto minimum payment for high level management. In many cases this minimum payment is set close to \$1 million, following the tax law that does not permit US firms to expense non-incentive pay totaling more than \$1 million.

Through changes in tax law, accounting practice, and via academic research, one could perhaps facilitate incentive schemes where the bonus for today’s outcome is revokable if tomorrow’s outcome is bad, again making negative payments more feasible. Our results suggest that such changes can be desirable, even from the point of view of the utility of the agent being required to put more skin in the game. Similarly, in partnerships, buy-in is common, and, to some extent relaxes this constraint.⁸

1.3 Outline of the Paper

Section 2 reviews related literature. Section 3 gives the model. Section 4 gives some key results when the number of agents is fixed. Section 5 discusses the model when the principal can adjust the number of agents. We begin by discussing the effect of a change in the minimum payment on the total effort the principal will choose to induce, and then, in Section 6 turn to the more difficult issue of effort per agent, with the technical analysis concentrated in Section 7.⁹ We supplement our theoretical results with numerical analysis. In Section 8 we turn to a model with firms as principals and workers as agents in a simple market context, and examine the effect of a change in

⁷See Kadan and Swinkels (2008) for an extensive exploration of the role of minimum payment constraints in understanding the composition of executive pay.

⁸Buy-in is common in architectural firms. While it may seem remarkable in light of subsequent events, Goldman routinely asked managers to co-invest in the funds they managed.

⁹While we hope the reader will give Section 7 a chance, we provide guidance there for readers who wish to skip the technical details.

the minimum wage under some natural-seeming assumptions on how a change in the minimum wage changes the demand and supply for labor. In Section 9, we use our results to see just when these natural-seeming assumptions in fact hold. Section 10 concludes. Appendix I contains various proofs, while Appendix II contains details of our numerical explorations of the problem.

2 Related Literature

Lazear (2000) conducts an empirical study of Safelite Glass Corporation, and their introduction of a performance based scheme. His theoretical structure concentrates on a world with observable effort. He emphasizes the role of heterogeneity of worker ability in explaining the use of performance pay in such a setting, and notes that in evaluating the effects of performance pay on productivity, one must think about both incentive effects and selection effects.

Lemieux *et al.* study a sorting model along the lines of Lazear (1986), arguing that most of the observable increase in wage inequality between the late 1970's and early 1990's was driven by an increase in the use of performance pay, and attributing this to a change in the underlying return to skill, which in turn increased the returns to sorting. They find performance pay to be more prevalent in higher paid professions, with, for example, 78% of sales workers receiving performance pay, and 65% of those in finance, insurance, and real estate.¹⁰

The selection effects highlighted by these papers are important. We view examining the implications of a model of homogeneous workers with effort incentives as complementary to these efforts. A full understanding presumably incorporates both sorts of effects. It would be interesting to know whether unobserved issues like the difficulty of a given installation, the time spent rounding up parts for a given vehicle, and the amount of hand-holding any specific customer needs make even the Safelite example one where a moral hazard model gives extra useful insight.

We are not aware of empirical work that directly addresses the implications of minimum wages for productivity. Part of this may affect data limitations (see Griliches (1994)). As Bartel, Ichniowski, and Shaw (2004) argue “inside econometrics” – econometric study that takes place at the level of a single firm, and that involves a much more in depth look at that particular firm, is likely to be an illuminating route here. Two studies by Baker, Gibbs, and Holmström (1994a,b) follow this route (but do not consider minimum wages). An advantage of our theoretical technique is that the antecedents to our theorems seem very amenable to observation.

De Fraja (1999) highlights the importance of paying attention to the change in working conditions that a change in the minimum wage will bring about, and he points to the fact that a change in the minimum wage can have consequences for workers paid more than the minimum wage as well. In his model, effort is observable,

¹⁰The average total hourly pay of those receiving performance pay (including performance pay) in their sample is just over \$10 an hour in adjusted 1979 dollars. By comparison, the minimum wage for most of the sample period was \$3.35, and so represented about 1/3 of pay for those receiving performance pay at the beginning of the sample period (but fell in real terms during that time).

workers differ in their willingness to exert effort, and an increase in the minimum wage changes the menu of (effort, wage) pairs that firms will offer.

An active debate on the effect of minimum wage increases on employment levels begins with Card and Krueger (1994), who look at employment at fast food restaurants and argue against a negative impact on employment. Neumark and Wascher (2000) disagree with this conclusion, arguing that the data collection technique employed by Card and Kreuger skewed their results, and directly collect payroll data instead.¹¹ A good review of the burgeoning literature in the field, with a strong focus on international data, is provided by Neumark and Wascher (2007).

Several other papers embed an agency model in a market setting. Prescott and Townsend (1984) were the first to study the extent to which competitive equilibria can be applied to environments with moral hazard and adverse selection problems. They establish a general structure that embeds both moral hazard and adverse selection problems as special cases, allowing for randomized contracts. Rustichini and Siconolfi (2008) explore efficiency properties of equilibria in this setting, and establish cases of non-existence. Zame (2007) develops a general equilibrium model in which firms, incentive contracts, prices for inputs and outputs, and consumption are all determined endogenously. Baranchuk, MacDonald, and Yang (2009) study an agency model with agents of varying levels of ability. Firm size, entry and exit of firms, and incentives are all determined in equilibrium. Firms that hire more able managers provide stronger incentives, and generate larger output. None of the above papers studies the effects of minimum wages on incentives and employment, which is the focus of this paper.

3 Model and Preliminaries

A principal employs n agents. This may be fixed,¹² or n may be a choice variable. Each agent i chooses an effort level $e \in [0, 1]$, following which the principal observes a signal x_i according to distribution $F(\cdot|e)$. It is not important to us that x_i be one dimensional or satisfy, for example, *MLRP*. Components of the signal might be gross sales, an output measure, a stock price, customer satisfaction surveys, or an evaluation by a supervisor. Conditional on effort levels, the signals x_i are assumed to be independent across agents. Each agent has differentiable disutility of effort $c(e)$, and increasing, strictly concave, twice continuously differentiable utility function $u(w)$.¹³

Each agent is paid according to a contract $\pi(x_i)$, where we assume that i 's compensation is based only on x_i , and not, for example, on aggregate performance or on the signals of other agents. This is one feature that distinguishes our model from the

¹¹ Payroll data before and after a change in the desired effort level is precisely the sort of information we need to sign our theoretical predictions.

¹² At, for example, $n = 1$, yielding the standard moral hazard model (e.g. Mirrlees (1976, 1999) and Holmström (1979)).

¹³ Along the lines of Kadan, Reny, and Swinkels (2010), one could also allow for other components of rewards, such as promotions, or non-cash compensation. Our results can be generalized to that setting as long as there is always a cash component to rewards that enters utility in an additively separable way. We avoid the additional layer of complexity involved with this generalization.

literature on moral hazard in teams originated by Holmström (1982).¹⁴ A minimum payment constraint (M) is modeled by assuming that $\pi(x_i)$ is constrained to be at least $m \geq 0$ for all x_i and for all i . The agent chooses his effort level to maximize his utility. That is,

$$e \in \arg \max_{e'} \int_0^{\bar{x}} u(\pi(x_i)) dF(x_i|e') - c(e'). \quad (IC)$$

Finally, there is an individual rationality constraint (IR): the agent must receive at least utility u_0 .¹⁵ We assume symmetry so that all agents are motivated by the same incentive contract $\pi(\cdot)$ and therefore, in equilibrium, each agent exerts the same effort level e .

The expected payoff to the principal before labor costs is represented by a function $B(E)$ that depends only on the total effort E exerted by the agents. Given the symmetry assumed, $E = ne$. We assume that $B(\cdot)$ is twice differentiable, increasing and strictly concave. In Section 8, when we embed the model in a market setting, $B(\cdot)$ will get a particular functional form which has as a parameter the price in the product market.

The principal's problem is to choose the number of agents to employ, n , the effort to induce from each agent, e , and the incentive contract, $\pi(\cdot)$, to maximize

$$B(ne) - n \int_0^{\bar{x}} \pi(x) dF(x|e),$$

subject to the individual rationality, (IR), incentive compatibility (IC), and minimum payment (M) constraints.

For any given (e, m, u_0) , let

$$C(e, m, u_0) = \min_{\pi(\cdot)} \int_0^{\bar{x}} \pi(x) dF(x|e)$$

subject to IR , IC , and M . That is, $C(e, m, u_0)$ is the cost of inducing effort level e from any given agent subject to minimum payment m and outside option u_0 . To focus the analysis, we assume that $C(e, m, u_0)$ is twice continuously differentiable, with $C_e > 0$ and $C_{ee} > 0$ so that costs are increasing and convex in effort. Primitives for these conditions are of interest, but not the subject of this paper.¹⁶ Note that the

¹⁴Our main justification for this is simplicity and realism, but in many contexts, it is reasonable that, for example the final output of a firm is not observable at the time compensation is made, or that if there are many workers, then firm-level performance has essentially no information about e_i beyond that already provided by x_i .

¹⁵When we embed the moral hazard model into a market setting, we will allow u_0 to vary with i , but keep $c(\cdot)$ and $u(\cdot)$ independent of i .

¹⁶For existence, monotonicity in e , and uniqueness of $C(e, m, u_0)$ in a one dimensional setting when the first order approach holds, see Jewitt, Kadan, and Swinkels (2008). That paper also presents sufficient conditions on, *inter alia*, $c(\cdot)$, to assure that $C_{ee} > 0$. The optimal incentive contract in these cases will typically be "option like," with m being paid with a strictly positive probability. For existence in a quite general setting, see Kadan, Reny, and Swinkels (2010), although note that in a general setting, conditions for $C_e > 0$ and $C_{ee} > 0$ are not well understood, and so these have to be viewed as simplifying assumptions.

convexity of C implies that it is optimal to induce the same effort level from each agent. Thus, the symmetry assumption is without loss of generality.

The principal's problem can now be rephrased as

$$\max_{n,e} B(ne) - nC(e, m, u_0). \quad (1)$$

Given our assumptions, (1) has a unique solution (n^*, e^*) . We will further assume that this solution is interior, and hence given by the first order conditions of the problem.

4 A Fixed Number of Agents

When n is exogenous, we have essentially the standard moral hazard setup. Having first solved the cost minimization problem for each e , yielding $C(\cdot, m, u_0)$, the principal can choose e to solve (1). This is the standard “two-step” for such problems (Grossman and Hart (1983)), where first one finds the cost minimizing contract for each effort level, and then one solves for the optimal effort level.

Since $E = ne$, we have from (1) that the optimal effort level $E^* = ne^*$ is determined by

$$B'(E^*) = C_e\left(\frac{E^*}{n}, m, u_0\right).$$

So, since $B'' < 0$, the effect of a change in the minimum wage on E^* satisfies

$$E_m^* \underset{s}{=} -C_{me}, \quad (2)$$

where $\underset{s}{=}$ means “has the same sign as.” That is, if effort is harder to implement when m is higher, optimally induced effort will decline with m . Since n is exogenous, this trivially implies that the optimal effort per agent, e^* , will also decline in m .

The sign of C_{me} is not obvious. Kadan and Swinkels (2009a) show a variety of conditions under which in the one dimensional setting, $C_{me} > 0$. They also numerically explore (2009b) a large set of examples, and have no counter-example to this. Finally, they explore conditions (2010) under which $C_{me} > 0$ which are of the same flavor as those we use here, in that the conditions depend not on the precise structure of the primitives, but rather on how changes in the desired level of induced effort affect the realized distribution over pay. If, for example, inducing higher effort involves a first order stochastic shift in realized pay, then $C_{me} > 0$ will always hold.

Similarly,

$$E_{u_0}^* \underset{s}{=} -C_{u_0e}.$$

Kadan and Swinkels (2010) show that C_{u_0e} , the effect of the outside option on the marginal cost of inducing effort can go either way, but especially when pay exceeds the minimum required only when performance is fairly good, reasonable conditions on primitives yield $C_{u_0e} < 0$, so that an agent who has a higher outside option has a lower cost of being motivated. The intuition here is that one is making payments above the minimum for the dual purposes of yielding the agent enough utility, and providing

incentives. As the implicit value of providing pay for the purpose of providing utility goes up, the implicit cost of doing so for incentive reasons goes down. On the other hand, because an employee with a higher outside option is being paid more on average, providing incentives can often involve extra pay at places where pay is already high, and so the marginal utility of income is low, thus requiring a larger increase in pay for any given increase in incentives. This provides a force in the direction of $C_{u_0e} > 0$.

5 A Varying Number of Agents

5.1 A “Three-Step”

Now, we turn to the model where n is endogenous. We will for simplicity treat n as a continuous variable from this point on.¹⁷ Note that using $E = ne$ we can rewrite (1) as

$$\max_{E,e} B(E) - E \frac{C(e, m, u_0)}{e}. \quad (3)$$

An immediate implication of (3) is that an optimizing principal will minimize $\frac{C(e, m, u_0)}{e}$.

Remark 1 *If e^* solves (1) then e^* solves*

$$\min_e \frac{C(e, m, u_0)}{e}.$$

This is intuitive, as the benefit to the principal $B(\cdot)$ relies on total effort only, and so maximizing (3) for any given E is equivalent to minimizing $\frac{C(e, m, u_0)}{e}$.

Let $e^*(m, u_0)$ be the optimal individual per-agent effort level given m and u_0 .¹⁸ Let

$$a(m, u_0) = \frac{C(e^*(m, u_0), m, u_0)}{e^*(m, u_0)}$$

be the minimized average cost of effort. We can then further rewrite (3) as

$$\max_E B(E) - a(m, u_0) E,$$

for which the first order condition is

$$B'(E^*) = a(m, u_0). \quad (4)$$

Intuitively, at the optimum the principal is indifferent between eliciting one more unit of effort from each existing agent, and hiring an additional agent at the current effort level. Hence, the marginal cost of total effort is just the average cost of inducing effort

¹⁷This allows us to replace a pair of difference equations (saying that the principal neither wishes to add or to eliminate an employee) by a single derivative. The approximation is small when n is large, and the results can be adapted if one wishes at the cost of some transparency.

¹⁸This is unique given that $C_{ee} > 0$.

from each existing agent, $a(m, u_0)$. At the optimum total effort E^* , this marginal cost is equal to the marginal benefit $B'(E^*)$.

Thus, while in the case of exogenous n , one can use the standard two-step approach, here we can use a three-step approach: First find the optimal incentive contract for each e and the cost function $C(\cdot)$. Second, find the individual effort level e^* that minimizes $\frac{C(e, m, u_0)}{e}$. Finally, introduce B and find E^* , the total effort that maximizes total profit.

5.2 Comparative Statics in Total Effort

By the envelope theorem,

$$a_m(m, u_0) = \frac{C_m(e^*(m, u_0), m, u_0)}{e^*(m, u_0)},$$

and so at any (m, u_0) where the minimum wage constraint M binds,¹⁹

$$a_m(m, u_0) > 0. \tag{5}$$

Intuitively, since $C(e, m, u_0)$ goes up for every e , its minimized average value $a(m, u_0)$ goes up as well. Similarly,

$$a_{u_0}(m, u_0) > 0 \tag{6}$$

when IR is binding at e^* .²⁰ From (4), (5), and (6), we have the following simple conclusion.

Theorem 1 *If M binds at $(e^*(m, u_0), m, u_0)$ then $E_m^*(m, u_0) < 0$. If IR binds at $(e^*(m, u_0), m, u_0)$ then, $E_{u_0}^*(m, u_0) < 0$.*

That is, the principal chooses to induce less total effort when a binding minimum payment or outside option rises. Note the contrast here to the case with a fixed number of agents. There, while we certainly believe that C_{me} is generally positive, the analysis that supports that conclusion is difficult, and the conditions involved are non-trivial. Here, the result stems simply from the fact that tightening a binding constraint increases costs, and thereby also increases the minimized average cost. For C_{u_0e} , where there are reasonable examples in the fixed number case where a higher u_0 generates a higher e , the contrast is even starker.

The key reason why the result here is so much more robust is that because the principal can adjust both the effort level per agent *and* the number of agents, the

¹⁹For the purposes of this paper, we shall say that M binds at (e, m, u_0) if m is paid with positive probability at the (unique) optimal contract implementing e given m and u_0 . As Jewitt, Kadan and Swinkels discuss, it can be that in the primal problem, the principal chooses an e for which the least cost implementation never pays M , but would, were it not for M , implement a different effort level and earn higher profits. They argue that the two senses of the term agree when $C_{ee} > 0$ for the problem without the M constraint.

²⁰If $C_{ee} > 0$ did not hold, then, $e(m, u_0)$ would not necessarily be unique, but $a_m(m, u_0) > 0$ would continue to hold if M was binding at all average cost minimizing e , and would hold as a weak inequality otherwise, and similarly for $a_{u_0}(m, u_0)$.

question of whether the principal will want to adjust E^* by changing the resources expended on providing motivation at the margin is replaced by the much simpler question of whether the principal wants to adjust E^* by adjusting the number of agents.

6 “Minimum Efficient Scale” and Effort Per Agent

For many questions, knowing what happens to total effort E suffices. But, for others, as for example what happens to total employment, we also need to know what happens to e and n . So, we now turn to the question of what happens to the effort level of the individual agent as m changes. The key observation is that because e^* solves the problem of minimizing $\frac{C}{e}$, one can think of e^* as minimum efficient scale for the activity of deriving effort from a given agent. Understanding e^* will be useful both here, and later when we turn to understanding the implications of embedding the model with *either* a fixed or variable number of agents, into a market context.

We begin with a lemma.

Lemma 1 *At any m, u_0 and at $e^*(u_0, m)$,*

$$e_m^* = \frac{C_m}{s} \frac{C}{C} - \frac{C_{me}}{C_e},$$

and

$$e_{u_0}^* = \frac{C_{u_0}}{s} \frac{C}{C} - \frac{C_{u_0e}}{C_e}.$$

This is in fact very intuitive. The minimum efficient scale e^* is defined by the intersection point of the average cost $\frac{C}{e}$ and the marginal cost C_e curves. Thus, whether e^* rises or falls when the cost structure changes depends on whether the percentage by which the marginal cost curve goes up, here $\frac{C_{me}}{C_e}$ (or $\frac{C_{u_0e}}{C_e}$) is smaller or larger than the percentage by which the average cost curve moves up, here $\frac{C_m}{C}$ (or $\frac{C_{u_0}}{C}$).

We begin with a simple observation.

Remark 2 *If at m, u_0 and $e^*(m, u_0)$, $C_{me} \leq 0$ and M binds, then since $C_m > 0$, $e_m^* > 0$. Similarly, if $C_{u_0e} \leq 0$ and IR binds then $e_{u_0}^* > 0$*

Thus, if raising m lowers the marginal cost of inducing effort, the principal will certainly induce more of it. Hence, any time that the exogenous n model would have predicted an increase in effort, the endogenous n model does as well.

As discussed, we suspect that the typical case has $C_{me}(m, u_0) > 0$. But, by Lemma 1, even when $C_{me} > 0$ it can still be that $e_m^* > 0$. It turns out (see Appendix II) that it is fairly easy to generate examples where $C_{u_0e} < 0$. But, again, even when $C_{u_0e} > 0$, we can still have $e_{u_0}^* > 0$. In each case, the question is whether the increase in marginal costs is matched by at least as large a change in average costs. Section 7 looks at when the relative size of these effects can be compared, shows that e_m^* and $e_{u_0}^*$ are tightly linked, and derives plausible (but not vacuous) conditions under which the net effect is indeed in the direction of larger effort per agent.

6.1 Effect on Employment

When $e_m^* > 0$ or $e_{u_0}^* > 0$, the principal responds to an increase in the minimum payment or outside option by *increasing* the effort it asks from an individual agent. Since (by Theorem 1) $E_m^* \leq 0$ and $E_{u_0}^* \leq 0$ (and strictly so when the respective constraints bind) it follows that

$$n^*(m, u_0) = \frac{E^*(m, u_0)}{e^*(m, u_0)}$$

takes a double hit – less overall effort is demanded, but each remaining agent supplies more of it.

On the other hand, when e_m^* or $e_{u_0}^*$ is negative, the question comes down to the relative magnitude of the two terms. It can easily be shown that for any given e_m^* or $e_{u_0}^*$, when B'' is sufficiently small (in absolute value), the drop in E^* is larger than that in e^* . That is, the demand by the principal for E is sufficiently elastic as to overwhelm the decrease in minimum efficient scale, and total employment falls. On the other hand, if B'' is large (which means that the principal has an inflexible demand for E), and if e_m^* or $e_{u_0}^*$ is negative, then we can have the effect that an increase in m or u_0 raises the total number of agents employed. Note that what is happening in this case is that each agent exerts less effort, while the principal demands roughly the same total effort.

6.2 An Extension and a Caveat

The following simple extension will turn out to be very useful when we embed the model in a market setting. Consider the case where there are significant fixed costs associated with the employment of each additional agent. Then the result in Lemma 1 can be overturned. That is, assume that the principal maximizes

$$B(ne) - n(T + C(e, m, u_0)),$$

where $T > 0$ is some fixed expense associated with an extra agent.²¹ Then, the principal sets e to solve

$$\min_e \frac{T + C(e, m, u_0)}{e}.$$

The analysis showing that E falls with m or u_0 is as before. But, the conclusion of Lemma 1 will fail for T large enough. In particular, if one replicates the analysis from above, one arrives at

$$e_m^* = \frac{C_e}{T + C} - \frac{C_{me}}{C_m}, \tag{7}$$

and so for large enough T , $e_m^* \geq 0$ need no longer hold. Thus, when the hiring of an extra agent imposes large new costs independent of the compensation of the agent, it becomes more likely that $e_m^* < 0$.

²¹Pace Econ 101, there are distinctions to be made here about how T is properly thought of in shorter and longer run contexts.

7 Signing the Change in Effort Level

This section dives into the mathematics of the moral hazard problem to see what the forces are that drive e_m^* and $e_{u_0}^*$. The reader who is willing to trust that e_m^* and $e_{u_0}^*$ can reasonably be positive even when $C_{me} > 0$ or $C_{u_0e} > 0$, can move directly to the market setting of Section 8 without loss of continuity. Subsection 7.1 states the results and 7.2 interprets the conditions. Subsection 7.3 presents a more general result from which these results flow, and shows the proof for the simpler of our two cases. A discussion of our numerical exploration of the problem is contained in Section 7.4 and in Appendix II.

7.1 The Conditions and Results

The central conditions behind our results depend on how the realized pay distribution of the agent changes when the principal changes the effort level it wishes to induce. In particular, fixing u_0 and m , let $G(p|e)$ be the probability that the agent receives p or below when the principal induces effort level e . That is, G is the measure of the set of x such that $\pi(x, e, m, u_0) \leq p$.²² Thus, G depends on both the contract and on the distribution over outcomes. As such, it is not a primitive. But, as discussed at more length in the next subsection, G is inherently both interpretable and observable in a way that assumptions on primitives are unlikely to be.

Let $\bar{p}(e)$ be the upper limit of the support of $G(p|e)$. We will require throughout that G has some basic continuity and differentiability properties.

Assumption 1 $G(p|e)$ is continuous in (p, e) on $[m, \infty) \times [0, 1]$. G is differentiable with respect to e with $G_e(\cdot|e)$ continuous on $[m, \infty) \times [0, 1]$.

This is satisfied, for example, when the First Order Approach (FOA) holds under the conditions studied in Jewitt, Kadan, and Swinkels (2008).

Our second condition is that $\frac{-G_e(p|e)}{1-G(p|e)}$ is increasing in p on the domain of p . This says that as the agent both faces and responds to a contract that makes her work harder, the probability of being above any given pay level (which is $1 - G(p|e)$) changes by more, in percentage terms, at higher pay levels. Further interpretation is provided in the next subsection.

Finally, our results will be delineated according to whether $\frac{1}{u'(p)}$ is concave or convex in p . For the CRRA utility function, $u(p) = \frac{p^{1-\gamma}}{1-\gamma}$, and so $\frac{1}{u'}$ is concave iff $\gamma \leq 1$, i.e., with log utility as the borderline case. For the negative exponential, $u(p) = -e^{-\gamma p}$ with $\gamma > 0$, we have $\frac{1}{u'(p)} = e^{\gamma p}$ and so $\frac{1}{u'}$ is convex.

Theorem 2 Assume that $\frac{-G_e(p|e)}{1-G(p|e)}$ is increasing in p on $[m, \bar{p}(e)]$, and that $\frac{1}{u'(p)}$ is concave. Let $e^*(m, u_0) \in (0, 1)$. Then,

²²If signals are one dimensional, and if π is increasing in x (as it will if the FOA and MLRP hold), then

$$G(p|e) = F(\pi^{-1}(p; e, m, u_0) | e).$$

- (i) if IR is non-binding at $(e^*(m, u_0), m, u_0)$, then $e_m^*(m, u_0) \geq 0$,
and
(ii) if M is not binding at $(e^*(m, u_0), m, u_0)$ (that is, if $G(m|e) = 0$), then $e_{u_0}^*(m, u_0) \geq 0$.

Each result works by signing the appropriate expression from Lemma 1. So, for at least two important sets of cases, we can unambiguously conclude that the marginal cost effect is dominated by the average cost effect: despite the fact that effort is more costly to induce on the margin, more effort is in fact induced. One interpretation of (ii) is that if better opportunities in an alternate sector give agents in a particular sector a better outside option, then principals in that sector should respond by facing agents with a contract that is more attractive, but also involves working harder.

A result delineating a set of cases where the marginal cost effect dominates the average cost effect is given by the following.

Theorem 3 Assume that $\frac{-G_e(p|e)}{1-G(p|e)}$ is increasing in p on $[m, \bar{p}(e)]$, and that $\frac{1}{w'(p)}$ is convex, with $\frac{1}{w'(0)} = 0$. Let $e^*(m, u_0) \in (0, 1)$. Then,

- (i) if IR is non-binding at $(e^*(m, u_0), m, u_0)$, then $e_m^*(m, u_0) \leq 0$,
and
(ii) if M is not binding at $(e^*(m, u_0), m, u_0)$, then $e_{u_0}^*(m, u_0) \leq 0$.

For cases where IR binds but m is still sometimes paid, the forces involved in sorting out e_m^* and $e_{u_0}^*$ become more complex. For a discussion of what we do and do not know at a theoretical level when IR binds, see Section 7.3, which presents a more general result of which Theorems 2 and 3 are an implication.

Note that this result does *not* depend on the first order approach (FOA). However, the ease with which one can check the key condition on G in examples certainly does, as when the FOA holds there is a quick numerical method for finding $C(\cdot)$. See Subsection 7.4.

7.2 Interpreting the Conditions

As noted, G is a function of both the primitives and the contract. As discussed at more length in Kadan and Swinkels (2010), we think there is significant value in assumptions stated in terms of G .

First, properties of G are often easy to interpret in a way that conditions on primitives may not be. The conditions behind Theorem 2 have economic content that conditions on for example F_{eee} do not. Second, even where results fully in terms of primitives are possible, they will tend to be very restrictive, because a large number of effects all need to be signed in the same direction. Results via G can help illuminate the balance of forces necessary. Note especially that by dint of a condition on G we dispense with the first order approach and with any need for one dimensional and ordered signals. Third, results via G can be a useful stepping stone towards results in terms of primitives.

Finally, it is not clear that the standard notion of “primitive” is exactly on target. It seems hard to imagine that one can meaningfully address, from an econometric point of view, whether, for example, the distribution over signals moves in any particular way with effort, especially when looking at more subtle properties of the distribution. It may be very difficult to even sort out what measures individual principals are using, and to get any sort of systematic access to how those measures move with effort.

But, G is inherently observable, and the requirements on the convexity of $\frac{1}{u'}$ are rather intuitive and interpretable. So, while G is partly endogenously determined, it may be that we can arrive at a sensible understanding of whether, when the principal asks more effort of its agents (and they respond), the distribution of pay moves in the appropriate sense. For example, one could imagine a combination of the right natural experiment and access to aggregate tax records as saying a lot about this issue. But then, our theorems allow one to make predictions about harder to observe things like the effect of a prospective policy change on employment and employee well-being.

To see an intuition for the condition that $\frac{-G_e(p|e)}{1-G(p|e)}$ be increasing in p , note that when signals are one-dimensional, and $MLRP$ holds,

$$\frac{-F_e(x|e)}{1-F(x|e)}$$

is increasing in x . So, the issue is that the changes in contract are not so severe as to confound this effect.

Remark 3 *Assume that signals are one dimensional, and that $\pi(\cdot, e, m, u_0)$ is increasing. Then, we can differentiate*

$$G(p|e) = F(\pi^{-1}|e)$$

by e to yield $G_e = F_e + \frac{d}{de}[\pi^{-1}] f$, and so

$$\frac{-G_e}{1-G} = \frac{-F_e(\pi^{-1}|e)}{1-F(\pi^{-1}|e)} + \left[-\frac{d}{de}[\pi^{-1}] \right] \frac{f(\pi^{-1}|e)}{1-F(\pi^{-1}|e)}. \quad (8)$$

For the case where the IR does not bind, a fairly typical outcome is that π_e is positive at all outputs at or beyond $\hat{x}(e, m)$, and so $\frac{d}{de}[\pi^{-1}]$ is negative (see Kadan and Swinkels (2009a) for further discussion). Under $MLRP$, $\frac{-F_e}{1-F}$ is increasing, as is $\frac{f}{1-F}$, and so it would thus be enough that $-\frac{d}{de}[\pi^{-1}]$ is increasing in p , or, equivalently, that $\pi_{ex} \geq 0$. So, if the contract increases and gets steeper in x as e increases, then our condition is satisfied.

In general, one has the intuition that if an increase in e causes the contract to get steeper, and single cross the old contract, then the condition should hold. A general proof is complicated by the fact at p below the crossing, one is dealing with the offsetting effects from the first and second terms on the RHS of (8).

7.3 A More General Result

Theorems 2 and 3 are corollaries of the following result.

Theorem 4 *Assume that $\frac{-G_e(p|e)}{1-G(p|e)}$ is increasing in p . If $\frac{1}{u'(p)}$ is concave, then for all $e \in (0, 1)$*

$$\frac{\frac{C_{me}}{u'(m)} + C_{u_0e}}{\frac{C_m}{u'(m)} + C_{u_0}} \leq \frac{C_e}{C}.$$

If $\frac{1}{u'(p)}$ is convex, with $\frac{1}{u'(0)} = 0$, then the inequality is reversed.

The first step to proving Theorem 4 is the following lemma.

Lemma 2 *For any e, m, u_0*

$$\frac{C_m}{u'(m)} + C_{u_0} = \int_0^{\bar{p}(e)} \frac{1}{u'(p)} dG(p|e). \quad (9)$$

See Kadan and Swinkels (2010) for a proof. The thought experiment is that, in utility terms, both the minimum payment and the outside option are raised by the same amount. Simply adding the same utility at all points in the contract restores *IC*, noting that adding a constant to utility at each signal does nothing to affect the relative attractiveness of any two effort levels to the agent, and, of course, restores *IR* as well. Raising utility at any given p costs $\frac{1}{u'(p)}$. The integral term is the expectation of this cost.

The proof of Theorem 4 begins from the fact that after integrating by parts,

$$C(e, m, u_0) = \int_0^{\bar{p}(e)} [1 - G(p|e)] dp$$

and so

$$C_e(e, m, u_0) = \int_m^{\bar{p}(e)} [-G_e(p|e)] dp,$$

Using this, we show how to relate $\frac{C_e}{C}$ to a conditional expectation of $\frac{-G_e(p|e)}{1-G(p|e)}$ with respect to the density given by

$$h_1(p) = \frac{[1 - G(p|e)]}{\int_0^{\bar{p}(e)} [1 - G(t|e)] dt} dp.$$

Similarly, using Lemma 2, we relate

$$\frac{\frac{C_{me}}{u'(m)} + C_{u_0e}}{\frac{C_m}{u'(m)} + C_{u_0}}$$

to a conductorial expectation of $\frac{-G_e(p|e)}{1-G(p|e)}$ with respect to the density

$$h_2(p) = \frac{\left[\frac{1}{u'(p)}\right]' [1 - G(p|e)] dp}{\int_0^{\bar{p}(e)} \left[\frac{1}{u'(t)}\right]' [1 - G(t|e)] dt}.$$

We show that when $\frac{1}{u'(p)}$ is concave, h_2 is stochastically dominated by h_1 , while when it is convex, h_2 stochastically dominates h_1 . Each result follows since $\frac{-G_e(p|e)}{1-G(p|e)}$ is increasing by assumption.

To see that Theorem 4 implies Theorem 2, note that when IR is non-binding, C_{u_0} and C_{u_0e} are both zero, while when m is paid with zero probability, C_{me} and C_m are zero. Theorem 4 similarly implies Theorem 3.²³

When both constraints are active, Theorem 2 tells us that a weighted average of $\frac{C_{me}}{C_m}$ and $\frac{C_{u_0e}}{C_{u_0}}$ is less than $\frac{C_e}{C}$ when $\frac{1}{u'}$ is concave, and conversely when $\frac{1}{u'}$ is convex with $\frac{1}{u'(0)} = 0$. But, it does not give us much guidance as to the individual terms. An approach of any form (using the G structure or otherwise) that gave theoretical results for the terms individually when both constraints are binding would be highly desirable.

7.4 Numerical Examples

As one path to a fuller understanding of the forces involved in e_m^* and $e_{u_0}^*$, in Appendix II we provide two detailed examples to illustrate the effect of changes in m and u_0 on e^* . To construct the examples we use the numerical algorithm introduced in Kadan and Swinkels (2009b). This algorithm is based on the constructive existence proof from Jewitt, Kadan, and Swinkels (2008). This allows us to easily check whether $\frac{-G_e(p|e)}{1-G(p|e)}$ is increasing. Further description of the numerical algorithm as well as additional numerical explorations can be found in Kadan and Swinkels (2009b).

The examples in Appendix II suggest that the requirement that $\frac{-G_e(p|e)}{1-G(p|e)}$ be increasing in p is satisfied in many cases of interest. In fact, in extensive exploration, we have yet to find an example where this condition fails. The examples show cases where M is binding but IR is not and the conditions of Theorem 2 are satisfied. So, while $C_{em} > 0$, we still have that $e_m^* > 0$. We also show an example in which both IR and M are binding, and where Theorem 4 applies to tell us that at least one of e_m^* and $e_{u_0}^*$ is positive. But, in that example, we in fact have $e_m^* < 0$, showing that e_m^* need not always be positive when both constraints are active, even if the other conditions for the result hold.

²³A simple possible extension is to consider the effect of a tax schedule on income, and let $t(p)$ be the agent's take-home pay as a function of gross pay p , where $t(0) = 0$, and t is concave. Then, we can effectively replace $u(\cdot)$ in all calculations by $u(t(\cdot))$, which is *a fortiori* concave. So, Theorem 4 continues to hold, where what now counts is the concavity or convexity of $\frac{1}{u'(t(p))t'(p)}$.

8 A Market Setting

Thus far, we have taken u_0 , the outside option of the agent, as given, as well as whether the principal is operating or not, and the function B giving the benefit to the principal of any given effort level. In this section, we embed the principal-agent model in a simple market setting, one in which the value to principals of extra effort, and the number of principals and agents, is endogenous, and consider some of the implications of doing so. The most natural application of our model is to a context where principals are firms, and agents are workers, and we hence adopt that nomenclature.

8.1 A Model

The Product Market We examine a single sector consisting of a set of firms who employ workers to produce a homogeneous output which is then sold into a competitive final goods market with inverse demand curve $P^D(q)$.

Workers All workers employed in the sector are as described in Section 3, and, for simplicity, are identical in terms of the sector (i.e., have the same $c(\cdot)$ and F). But, in contrast to the previous section, we assume that workers have heterogeneous outside options, so that there is an upward sloping supply function $N^S(u)$ relating the number of workers willing to work in the sector to the utility being offered in the sector.²⁴ All workers in the sector are indistinguishable from the point of view of firms, and so they all receive the same contract, with the IR constraint being determined by the utility of the marginal agent in the sector.

Firms For simplicity, we take the number of firms as continuous, and assume that when j firms are operating, the marginal firm has fixed cost FC_j of operation, where FC_j is strictly increasing and continuously differentiable in j . Conditional on operation, each firm faces the same strictly concave, increasing, and twice differentiable production function $y(E)$ relating total effort expended by its workforce to output, where $y(0) = 0$.²⁵ So, if the market price is p , then for firm j , the benefit function takes the form

$$B(E) = py(E) - FC_j.$$

Note that conditional on operating, all firms face effectively the same maximization problem, as their profits differ only by a constant, and so they all choose the same E and e .

²⁴As a primitive for this, assume workers have different outside options in another sector not subject to the minimum wage. The other sector could be unemployment, self-employment, the grey-market economy, or devoting more time to school. It could also be a sector with a different production and information technology where a legal minimum payment exists, but pay is flatter in performance (as, for example, constant), and no payment near m is in fact ever made. This last example is especially relevant when thinking about highly compensated individuals, where the choice may be between a safe job with somewhat variable compensation well above the minimum payment, or a risky job with minimum compensation running into a lower bound that exists for any of the reasons we have discussed.

²⁵One could also consider, for example a model where E translates into quality.

8.2 Market Clearing and Its Implications

We begin by noting that because agents receive a compensation package consisting of more than their minimum payment, “market clearing” in this model does not take place in wage space. Rather, it takes place in terms of the utility offered by the contract. That is, for any given minimum payment, we look at a demand curve for labor which gives, for any given utility, the number of workers firms wish to hire if they must offer at least that utility, and the number of workers who are willing to work given that the contract offers the given utility.

The demand curve for labor is in turn derived from the workings of the product market: we ask, for given utility level u and minimum payment m , what the induced costs of firms in the product market are, given that principals solve a moral hazard problem for each agent as described. This determines the supply curve in the product market. The supply curve interacts with product market demand to determine prices and quantities in the product market, and how much labor is demanded.

Call this derived labor demand $N^D(u, m)$, and for now, let us posit that $N^D(u, m)$ is weakly decreasing in u and in m . That is, assume that as either the utility to be provided or the minimum payment required increase, the demand for labor weakly decreases. Shortly, we turn to establishing when this is in fact reasonable, but for now, we see how market clearing should be defined given this assumption, and trace its implications.

Recall that, as we have seen, for some specifications of u and m , the principal, in choosing an optimal contract, may choose to provide utility greater than u . This will be reflected by a vertical portion of the demand curve for labor. For example, in Figure 1, for the black demand curve (corresponding to a minimum payment m_L), over the range of utilities below \hat{u}_L the principal optimally chooses to provide utility \hat{u}_L instead, and IR is non-binding. The green demand curve illustrates a demand curve for labor for a higher minimum payment, m_H . It lies to the left of the black demand curve, reflecting our assumption that $N^D(\cdot; m)$ decreases in m . The form in which it is drawn, in which $\hat{u}_H > \hat{u}_L$, is at a natural case, and we conjecture that it is *the* natural case. The intuition is that when the minimum wage is increased, the firm chooses to induce higher effort. To do this, payments will be increased at high performance levels. But, since IR is not binding, there are no payments at particularly low performance levels that can be correspondingly reduced.

Figure 1 also illustrates two possible supply curves for labor. The upper (red) supply curve intersects demand where IR binds. Hence, equilibrium quantities and utilities are determined by the crossings. When m_L is replaced by m_H , the equilibrium quantity of labor, and utility of those who keep their job, are both reduced. For the lower (blue) demand curve, the situation is different. For each of m_L and m_H , we are on the vertical portion of the demand curve, and IR does not bind in equilibrium. So, while the equilibrium quantity of labor falls when m_L is replaced by m_H , the equilibrium utility is the one that principals choose to provide subject to this minimum, and is thus \hat{u}_L for m_L , and \hat{u}_H for m_H . Thus, at the market clearing N and u , there is in fact an excess supply of labor.

Thus, when IR is binding, as for the red demand curve, an increase in the mini-

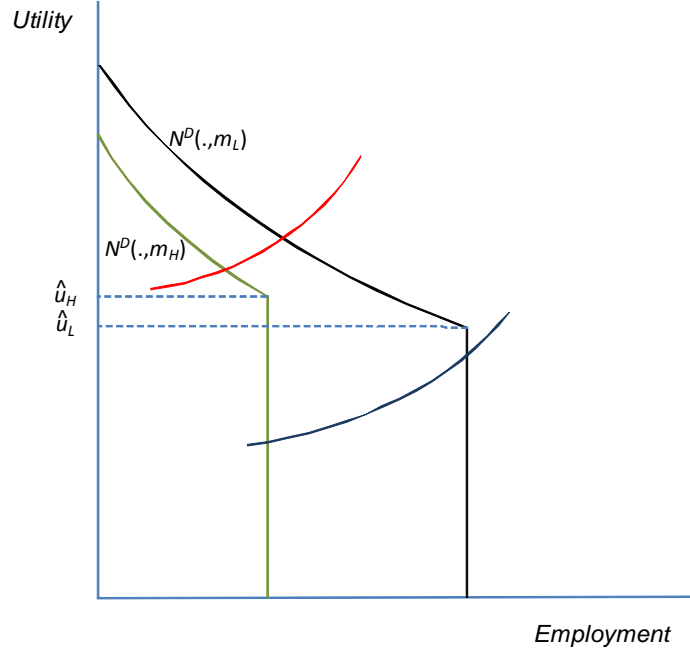


Figure 1: Demand and Supply for Labor under Different Levels of Minimum Pay

minimum payment *lowers* the utility *even of agents who keep their jobs*. This conclusion is trivial under the assumptions given. But, it is important! It says that minimum wages are likely to have very perverse effects in performance pay sectors where the *IR* constraint binds.

What is happening here is simple: when the *IR* constraint binds, the firm, when faced with a higher minimum wage, finds it optimal to re-shuffle incentives, employment, and induced effort in such a way as to put the marginal worker back on the *IR* constraint. Given that there are now fewer workers, this involves a lower utility than before, simply because the supply curve is upward sloping. This is very different from a standard model where the wage is the only free parameter. There workers who keep their jobs and earn m_H instead of m_L are trivially better off, as they are paid more and nothing else about their terms of employment changes.

Note also that this very simple conclusion holds for many settings other than the moral hazard problem at issue here. As long as the firm has multiple dimensions upon which it can adjust the contract of the worker when the minimum wage changes, it is less than clear that an increase in the minimum wage helps even workers who keep their jobs, as the firm will, in many settings have the incentive to adjust the other terms of employment in such a way as to return the agent to the *IR* constraint.²⁶

²⁶It is not for nothing that the actual law implementing the minimum wage is quite intricate, with many clauses directed precisely at attempting to prevent firms from “clawing back” a minimum wage increase by, *inter alia* requiring workers to buy their own uniform, changing the number of breaks they have, and requiring them to be on the factory floor by starting time rather than in the change

As long as the firm has the incentive to do so, as long as the effect of the increase in the minimum wage is to lower the number of workers demanded, and as long as the supply of workers is larger when the utility they are offered is higher, the result holds. The content of the analysis that follows is to verify when each of these is true in the specific setting of a moral hazard problem interacting with a minimum wage.²⁷

When the *IR* constraint is not binding, as for the blue demand curve, it remains the case that an increase in the minimum wage reduces equilibrium employment. But, when $\hat{u}_H > \hat{u}_L$, at least it is the case that workers who keep their positions are better off than before.

9 Unpacking the Model

Figure 1 embeds a number of assumptions. First, the demand curve for labor is drawn as sloping downward. Second, demand is assumed to move in when the minimum wage increases. Finally, as drawn, the minimum utility offered increases when *IR* is non-binding and the minimum wage increases. None of these is in fact obvious.

In this section, we analyze the product market more carefully to see how labor demand is derived. We shall see that indeed, it is natural for the signs to be as assumed. But, we will also see the forces involved that can lead to a reversal.

9.1 Marginal and Average Costs

We will need to understand the relationship between what happens to the marginal and average costs of output for the marginal firm. For the variable number of workers case, the average total cost (per unit of output, not per unit of effort!) of any given firm j using E units of effort is

$$ATC_j(y(E), m, u_0) = \frac{FC_j + a(m, u_0)E}{y(E)}, \quad (10)$$

and the marginal cost is

$$MC(y(E), m, u_0) = \frac{a(m, u_0)}{y'(E)}. \quad (11)$$

For the fixed number of workers per firm case, we have

$$ATC_j(y(E), m, u_0) = \frac{FC_j + nC(E/n, m, u_0)}{y(E)}, \quad (12)$$

and

$$MC(y(E), m, u_0) = \frac{C_e(E/n, m, u_0)}{y'(E)}. \quad (13)$$

room.

²⁷A survey of employers (Converse et al, 1981, Table 17, page 279) found that 20.8% of all firms with minimum wage employees self-reported that they increased the responsibilities of workers in response to the 1980 change in the U.S. federal minimum wage. Tantalizingly, for our purposes, the same figure when restricted to employees with some tip income was 38.3%.

9.2 Output and Labor Demand

Firms that operate choose E so that $p = MC(y(E), m, u_0)$, and choose to operate if $p \geq ATC_j(y(E), m, u_0)$ for the optimal E . Let $j(p, m, u_0)$ be the number of active firms, and $E(p, m, u_0)$ be the total effort employed by each active firm (recall that the optimal output and employment will be the same for all active firms, since they differ only in their fixed costs). In both the variable and fixed n cases, it follows that output per firm and the number of operating firms is continuous and increasing in p . Hence, for any given m and u_0 , the product market has a well defined upward sloping supply curve, and so there is a unique market clearing price, $p(m, u_0)$ and quantity $q(m, u_0)$. Let $P^S(q, m, u_0)$ be the inverse supply curve. Labor demand with a variable number of workers is

$$N^D(m, u_0) = j(p(m, u_0), m, u_0) \frac{E(p(m, u_0), m, u_0)}{e^*(m, u_0)},$$

while for the fixed number of workers case, it is simply

$$N^D(m, u_0) = nj(p(m, u_0), m, u_0).$$

9.3 Effects of an Increase in m or u_0 .

Let us consider the effect of an increase in m on E and N^D . The analysis for an increase in u_0 is essentially identical.

9.3.1 The variable number of workers case

From (11), for any given E , m , and u_0 ,

$$MC_m(y(E), m, u_0) = \frac{a_m(m, u_0)}{y'(E)}, \quad (14)$$

while by (10),

$$ATC_m(y(E), m, u_0) = \frac{a_m(m, u_0) E}{y}. \quad (15)$$

But, by (5), $a_m(m, u_0) > 0$, and so, since y is increasing and concave, and $y(0) = 0$,

$$MC_m - ATC_m = \frac{1}{s} \frac{1}{y'(E)} - \frac{E}{y(E)} > 0.$$

Pick any given \hat{m} and \hat{u}_0 . Let the equilibrium market output level be \hat{q} and let \hat{E} be the associated effort level per firm. Then, it must be that

$$ATC_m(y(\hat{E}), \hat{m}, \hat{u}_0) < P_m^S(\hat{q}, \hat{m}, \hat{u}_0) < MC_m(y(\hat{E}), \hat{m}, \hat{u}_0).$$

That is, the supply curve moves up by more than the change in ATC of a representative firm, but by less than MC . This holds since if $P_m^S \leq ATC_m$, then an increase in m would result in a weak decrease in the number of active firms, and a strict decrease

in output per firm, a contradiction to \hat{q} being produced, and similarly, if $P_m^S \geq MC_m$, then we would have a weak increase in the number of firms, and a strict increase in output per firm, again a contradiction.

Since demand slopes down, it follows *a fortiori* that the equilibrium price $p(\hat{m}, \hat{u}_0)$ satisfies

$$0 < p_m(\hat{m}, \hat{u}_0) < MC_m(y(E(p(\hat{m}, \hat{u}_0), \hat{m}, \hat{u}_0))),$$

so that the equilibrium price goes up, but by less than the marginal cost of an active firm at the original equilibrium output. It follows that each active firm produces less than before. But then, since y is concave, labor has higher average productivity. Since output has fallen, it follows that total effort at the industry level has fallen. Thus, if effort per worker has gone up (that is if $e_m^* \geq 0$), then it must be that employment fell. More generally, employment can go up only if both total effort in the industry level did not fall too severely (inelastic demand for the final good, and inelastic supply of firms) and e fell considerably. We summarize our discussion.

Theorem 5 *In the variable n case, total effort demanded falls if either m or u_0 rises.*

Corollary 1 *In the variable n case, if $e_m^*(m, u_0) \geq 0$, then $N_m^D(m, u_0) \leq 0$ and if $e_{u_0}^*(m, u_0) \geq 0$, then $N_{u_0}^D(m, u_0) \leq 0$, with each inequality strict if the appropriate constraint is binding.*

9.3.2 The fixed number of workers case

Now, let us turn to the fixed n case. Here,

$$MC_m(y(E), m, u_0) = \frac{1}{y'(E)} C_{me}(E/n, m, u_0)$$

while

$$ATC_m(y(E), m, u_0) = \frac{1}{y(E)} n C_m(E/n, m, u_0).$$

So, as before

$$P_m^S \in (\min(MC_m, ATC_m), \max(MC_m, ATC_m)),$$

and

$$0 < p_m(m, u_0) < \max(MC_m(y(E(p(m, u_0), m, u_0))), ATC_m(y(E(p(m, u_0), m, u_0)))).$$

So, the question comes to the relative magnitudes of MC_m and ATC_m . Note that

$$ATC_m - MC_m = \frac{1}{y(E)} n C_m(E/n, m, u_0) - \frac{1}{y'(E)} C_{me}(E/n, m, u_0).$$

For the marginal firm \hat{j} , we have

$$\frac{C_e(E/n, m, u_0)}{y'(E)} = \frac{FC_{\hat{j}} + nC(E/n, m, u_0)}{y(E)}.$$

We can thus substitute and rearrange to arrive at

$$\begin{aligned} ATC_m - MC_m &= \frac{nC_m}{s} - \frac{C_{me}}{C_e} \\ &= \frac{C_e}{\frac{FC_j}{n} + C} - \frac{C_{me}}{C_m}. \end{aligned}$$

But, note that this is of exactly the same form as (7). Hence, if it turns out that the effort level which minimizes average cost per unit of effort rises with m (i.e. $e_m^* \geq 0$), then we have that $ATC_m > MC_m$, so that $p_m(m, u_0) < ATC_m$. Hence, the original marginal firm is no longer profitable, and we have that $j(p(m, u_0), m, u_0)$, and hence labor demand once again falls with m .

On the other hand, as fixed costs grow, or when we are in the cases where $e_m^* < 0$, it becomes more likely that we get

$$ATC_m < p_m(m, u_0) < MC_m.$$

Thus, in these cases, and with sufficiently inelastic demand for the final good, we could see labor demand actually rise with the minimum wage.

9.3.3 Comparing the Cases

Interestingly, in each case, the key condition to ensure that labor demand falls with m (or analogously with u_0) is that the effort level that minimizes costs per unit of effort expended rises with m (or with u_0). But, the mechanism is different.

Assume that $e_m^* > 0$, and consider first the variable number of workers per firm case. As one follows the supply curve up with m while holding fixed output, each firm gets smaller in terms of output (and hence total effort demanded). This is so as the marginal cost of the marginal firm rises by more than the average cost, and thus the supply curve moves up by more than the average cost, but less than the marginal cost. In particular, if the demand curve were vertical, output would remain the same, but would be produced by more firms, each of whom produces less. If the demand curve has slope, then *a fortiori* firms become *smaller* in terms of effort employed (and hence output). Since firms produce less, their average productivity (output per unit of effort) improves, since y is concave, and effort demanded falls. Since firms find it optimal to induce more effort from each worker, they thus reduce the workforce for both reasons. The number of firms in the sector can either rise or fall.

In the fixed number of workers per firm case, as one again follows the supply curve up for a given output, each firm gets bigger. This is so as now the marginal cost curve of the marginal firm moves up by less than the average cost curve, and so the supply curve moves up by more than the marginal cost, but less than average cost. Hence, in the event of a vertical demand curve, we would have more production concentrated on fewer firms, which again has the effect of asking for more from each worker. If market demand has slope, then *a fortiori* the number of firms falls, with output per firm now potentially going either way, depending on the relative magnitude of the forces.

The reason why the market reaction in terms of output per firm is opposite comes down to the difference between the “two-step” and the “three-step”. In the case where n can be varied, the firm solves the question of the optimal employee effort separately from the question of the scale at which the firm should be operating. The concavity of $y(\cdot)$ implies that as m goes up, the right way for the market to adjust is by re-optimizing in the direction of smaller firms thus economize on the amount of effort used. But, because of the separability, this is perfectly consistent with the amount of effort being asked from each remaining worker increasing. In the case where n is not variable, if the right way to rearrange production involves a higher output per worker, this can only be achieved by increasing output per firm, and, in the face of sufficiently inelastic market demand, that is in fact what happens.

9.4 Comparing \hat{u}_H and \hat{u}_L

Finally, consider the question of \hat{u}_H versus \hat{u}_L . Note that the general flow of our results is that firms will choose to implement a higher effort level at m_H versus m_L . In Kadan and Swinkels (2009a), we show that when m increases and the IR constraint is not binding, the agent, for any given e , is paid more at all outputs, and hence is better off. Further, every numerical example there and in (2009b) has the feature that once IR is non-binding, as the implemented e rises, so does the utility of the agent. A very strong intuition for this is that when IR is non-binding, the principal already always pays the minimum at poor outcomes. To induce higher effort, he thus has to add incentives at good outcomes, and cannot reduce them elsewhere. As usual, a proof of this is complicated by the fact that the information structure changes as e changes (see a discussion in Kadan and Swinkels (2010) on this issue). Our examples and intuition suggest, however, that the conclusion holds quite generally.²⁸ Thus, both forces, of the effect on the workers’ utility of the change in m holding fixed effort, and of the effect on the workers’ utility of the change in the effort level the principal wishes to implement, push in the direction of $\hat{u}_H > \hat{u}_L$.

9.5 Introducing a Minimum Wage Sector Without Performance Pay

Our model is implicitly one of two sectors, Sector 1 with a minimum wage and performance pay, and Sector 2 with no binding minimum wage. In reality, there are also many parts of the economy where there is a binding minimum payment and no performance pay. To capture this, assume a Sector 3 with a binding minimum wage \hat{m} but no performance pay. If one changes the minimum wage in only the performance pay based sector (i.e., if one changes m but not \hat{m}), our analysis is unaffected by this (and U.S. labor law does treat the minimum wage to, for example, tip-based worker separately from the minimum wage in other sectors), and workers for whom the IR constraint binds continue to be worse off.

But, one might hope that in the event of the minimum wage \hat{m} in Sector 3 being increased, workers in Sector 1 would be better off via the effect on their outside option. This is less than clear. In particular, if an increase in \hat{m} reduces demand for labor in

²⁸Here again, better theoretical results would be desirable.

Sector 3, then the pool of agents in Sector 2 is increased, holding fixed the number in Sector 1. If the relevant IR constraint in Sector 1 is modeled as keeping the marginal agent in Sector 1 from jumping to Sector 2, then agents in Sector 1 are again hurt, simply because the addition of workers to Sector 2 means that the outside option of the marginal worker has become worse. On the other hand, it may be that agents have some observable heterogeneity (something not allowed by our simple formal model), and that firms in Sector 1 would much rather take their marginal agent from Sector 3 than from Sector 2 (think of Sector 2 as the unemployed). Then, an increase in \hat{m} helps agents in Sector 1. Finally, one could imagine the IR constraint in Sector 1 as reflecting some lottery over utility in Sectors 2 and 3. We leave a fuller exploration of these considerations for future research.

10 Conclusion

We study the effect of minimum wages in sectors with incentive pay. To this end, we extend the standard moral hazard problem to allow the principal to employ multiple workers. We show that the solution to this problem consists of three steps: (i) solve for the optimal incentive contract and find the cost of implementing any given effort level; (ii) find the optimal effort level minimizing cost per unit of effort; and (iii) find the optimal number of employees. Since the optimal effort level minimizes the average cost of implementing effort, it can be viewed as the minimum efficient scale of the activity of eliciting effort from a given agent.

When employment is flexible, an increase in minimum wage results in less total effort being exerted. But, the effect on the effort exerted by each individual agent and on employment is not obvious. We show that under reasonable conditions, an increase in the minimum wage causes the average cost of inducing effort to go up more than the associated marginal cost. As a result, in the face of an increased minimum wage, the principal chooses to reduce employment while asking each agent to work *harder*. By contrast, when employment is rigid, a typical result is that the marginal cost of inducing effort goes up, and hence each agent exerts less effort.

We embed the model in a simple production economy where effort serves as a production factor. The demand for labor is derived from equilibrium in the product market. Interestingly, while in the setting without market equilibration we typically had opposite effects, here both in the case where firms can freely adjust their labor force and in the case where employment is rigid, employment typically falls when minimum wages rise. However, the mechanisms for achieving this outcome are different. When labor is adjustable, firms find it optimal when minimum wages rise, to ask more from each individual worker, while reducing the overall workforce. When the number of workers per firm is rigid, the increase in marginal cost of inducing effort tends to increase average costs of production more than the associated increase in equilibrium price. As a result, the number of firms that choose to operate declines, and employment falls.

The effect of an increase in the minimum wage on workers' welfare depends on whether the IR constraint is binding or not. When the IR constraint binds, an

increase in minimum wage typically results in lower employment and, given that the supply curve is upward sloping, the outside option of the marginal worker falls. But then, firms optimally adjust compensation in such a way that even workers who retain their job are rendered worse off despite the increase in their base pay. By contrast, if the IR constraint is not binding then minimum wages typically benefit workers who retain their jobs.

The paper has both policy and empirical implications. On the policy side the paper suggests that moderate levels of minimum wages, in which IR is likely to bind, are not beneficial. Not only does employment fall (as classical models would suggest) but also the workers who keep their job are worse off. There is a strong suggestion that attempts to aid the working poor ought to focus instead on their outside option, or on the relative magnitude of demand and supply. Note for example, that a jobs-training program that moves some people into more attractive jobs out of the sector increases the well-being not just of those who make the transfer, but, via the equilibrium outside option, helps those within the sector as well. On the empirical side, the paper provides guidance for the growing debate on the effect of minimum wages. It provides directions for empirical exploration of the effect of minimum wages on both the labor market and the associated product market in sectors where incentive pay is common.

Appendix I: Proofs

Proof of Lemma 1 From Lemma 1, a necessary condition is that at $(e(u, m), m, u_0)$

$$\frac{d}{de} \left(\frac{C(e, m, u_0)}{e} \right) = 0, \quad (16)$$

and so

$$\frac{C}{e} = C_e. \quad (17)$$

Implicitly differentiating (17) by m and using the envelope theorem to eliminate the effect of e on the left hand side gives

$$\frac{C_m}{e} = C_{ee}e_m^* + C_{em}.$$

Rearranging, and using $C_{ee} > 0$, we have

$$e_m^* = \frac{C_m}{e} - C_{em}. \quad (18)$$

By (17), we can substitute for e at an optimum, and then rearrange to achieve

$$e_m^* = \frac{C_m}{C} - \frac{C_{me}}{C_e}.$$

One can similarly derive that

$$e_{u_0}^* = \frac{C_{u_0}}{C} - \frac{C_{u_0e}}{C_e}. \quad \blacksquare$$

Proof of Theorem 4 Note that

$$C(e, m, u_0) = \int_0^{\bar{p}(e)} p dG(p|e),$$

where recall that $\bar{p}(e)$ defines the top end of the support of $G(\cdot|e)$. Integrate this by parts (using $\frac{d}{dp}G(p|e) = \frac{d}{dp}(-[1 - G(p|e)])$) to arrive at

$$\begin{aligned} C(e, m, u_0) &= -p[1 - G(p|e)]\Big|_0^{\bar{p}(e)} + \int_0^{\bar{p}(e)} [1 - G(p|e)] dp \\ &= \int_0^{\bar{p}(e)} [1 - G(p|e)] dp. \end{aligned} \quad (19)$$

So

$$C_e(e, m, u_0) = \int_0^{\bar{p}(e)} [-G_e(p|e)] dp,$$

using Leibniz's rule (which is valid using Assumption 1), and noting that the term involving the derivative with respect to the upper limit of integration is zero since $1 - G(\bar{p}(e)|e) = 0$ by definition.

But then,

$$C_e(e, m, u_0) = \int_m^{\bar{p}(e)} \frac{-G_e(p|e)}{1 - G(p|e)} [1 - G(p|e)] dp \quad (20)$$

where we can restrict the range of integration (and avoid division by zero) since $G_e = 0$ for $p < m$. Let h_1 be the density given by

$$h_1(p) = \frac{1 - G(p|e)}{\int_0^{\bar{p}(e)} [1 - G(t|e)] dt},$$

and let H_1 be its cumulative. Then, from (19) and (20),

$$\begin{aligned} \frac{C_e(e, m, u_0)}{C(e, m, u_0)} &= \int_m^{\bar{p}(e)} \frac{-G_e(p|e)}{1 - G(p|e)} \frac{[1 - G(p|e)]}{\int_0^{\bar{p}(e)} [1 - G(t|e)] dt} dp \\ &= \int_m^{\bar{p}(e)} \frac{-G_e(p|e)}{1 - G(p|e)} h_1(p) dp, \end{aligned}$$

and so we have

$$\frac{C_e(e, m, u_0)}{C(e, m, u_0)} = [1 - H_1(m)] \int_m^{\bar{p}(e)} \frac{-G_e(p|e)}{1 - G(p|e)} \frac{h_1(p)}{1 - H_1(m)} dp. \quad (21)$$

Similarly, the right-hand side of (9) integrates by parts to give

$$\frac{C_m(e, m, u_0)}{u'(m)} + C_{u_0}(e, m, u_0) = \frac{1}{u'(0)} + \int_0^{\bar{p}(e)} \left[\frac{1}{u'(p)} \right]' [1 - G(p|e)] dp.$$

So,

$$\begin{aligned} \frac{C_{me}(e, m, u_0)}{u'(m)} + C_{u_0e}(e, m, u_0) &= \int_0^{\bar{p}(e)} \left[\frac{1}{u'(p)} \right]' [-G_e(p|e)] dp \\ &= \int_m^{\bar{p}(e)} \frac{-G_e(p|e)}{1-G(p|e)} \left[\frac{1}{u'(p)} \right]' [1-G(p|e)] dp. \end{aligned}$$

again appealing to Leibniz's rule, and to the fact that $G_e = 0$ below m . Define

$$h_2(p) = \frac{\left[\frac{1}{u'(p)} \right]' [1-G(p|e)]}{\int_0^{\bar{p}(e)} \left[\frac{1}{u'(t)} \right]' [1-G(t|e)] dt},$$

and let H_2 be the associated cumulative. Then,

$$\begin{aligned} \frac{\frac{C_{me}(e, m, u_0)}{u'(m)} + C_{u_0e}(e, m, u_0)}{\frac{C_m(e, m, u_0)}{u'(m)} + C_{u_0}(e, m, u_0)} &= \frac{\int_m^{\bar{p}(e)} \frac{-G_e(p|e)}{1-G(p|e)} \left[\frac{1}{u'(p)} \right]' [1-G(p|e)] dp}{\frac{1}{u'(0)} + \int_0^{\bar{p}(e)} \left[\frac{1}{u'(p)} \right]' [1-G(p|e)] dp} \\ &\leq \frac{\int_m^{\bar{p}(e)} \frac{-G_e(p|e)}{1-G(p|e)} \left[\frac{1}{u'(p)} \right]' [1-G(p|e)] dp}{\int_0^{\bar{p}(e)} \left[\frac{1}{u'(p)} \right]' [1-G(p|e)] dp} \\ &= \int_m^{\bar{p}(e)} \frac{-G_e(p|e)}{1-G(p|e)} h_2(p) dp, \end{aligned}$$

since $\frac{1}{u'(0)} \geq 0$, and hence we have

$$\frac{\frac{C_{me}(e, m, u_0)}{u'(m)} + C_{u_0e}(e, m, u_0)}{\frac{C_m(e, m, u_0)}{u'(m)} + C_{u_0}(e, m, u_0)} \leq [1 - H_2(m)] \int_m^{\bar{p}(e)} \frac{-G_e(p|e)}{1-G(p|e)} \frac{h_2(p)}{1-H_2(m)} dp. \quad (22)$$

Consider first the case that $\frac{1}{u'}$ is concave. Then, note that

$$\frac{\frac{h_2(p)}{1-H_2(m)}}{\frac{h_1(p)}{1-H_1(m)}}$$

is proportional to $\left[\frac{1}{u'(p)} \right]'$ which is decreasing by assumption. Hence $\frac{h_2}{1-H_2(m)}$ is stochastically dominated by $\frac{h_1}{1-H_1(m)}$,²⁹ and so, since $\frac{-G_e(p|e)}{1-G(p|e)}$ is increasing on $[m, \bar{p}(e)]$ by assumption,

$$\int_m^{\bar{p}(e)} \frac{-G_e(p|e)}{1-G(p|e)} \frac{h_1(p)}{1-H_1(m)} dp \geq \int_m^{\bar{p}(e)} \frac{-G_e(p|e)}{1-G(p|e)} \frac{h_2(p)}{1-H_2(m)} dp. \quad (23)$$

²⁹To see this, note that for any two densities z_1 and z_2 on $[0, \infty)$ with $\frac{z_2}{z_1}$ decreasing, we have that z_1 single crosses z_2 from below. So, $(1 - Z_1(t)) - (1 - Z_2(t))$ is single peaked. But, since the expression is 0 at $t = 0$ and $t = \infty$, it is weakly positive everywhere in between.

We claim that $1 - H_2(m) \leq 1 - H_1(m)$. To see this, note that

$$\begin{aligned} 1 - H_2(m) &= \frac{\int_m^{\bar{p}(e)} \left[\frac{1}{u'(p)} \right]' [1 - G(p|e)] dp}{\int_0^{\bar{p}(e)} \left[\frac{1}{u'(p)} \right]' [1 - G(p|e)] dp} \\ &= \frac{\int_m^{\bar{p}(e)} \left[\frac{1}{u'(p)} \right]' [1 - G(p|e)] dp}{\int_0^m \left[\frac{1}{u'(p)} \right]' [1 - G(p|e)] dp + \int_m^{\bar{p}(e)} \left[\frac{1}{u'(p)} \right]' [1 - G(p|e)] dp} \end{aligned}$$

which is of the form

$$\frac{X}{Y + X}.$$

Since both X and Y are positive, this expression is increasing in X and decreasing in Y . Since $\frac{1}{u'(p)}$ is concave, we thus have

$$X \leq \left[\frac{1}{u'(m)} \right]' \int_m^{\bar{p}(e)} [1 - G(p|e)] dp,$$

and

$$Y \geq \left[\frac{1}{u'(m)} \right]' \int_0^m [1 - G(p|e)] dp,$$

and hence

$$\begin{aligned} 1 - H_2(m) &= \frac{X}{Y + X} \\ &\leq \frac{\left[\frac{1}{u'(m)} \right]' \int_m^{\bar{p}(e)} [1 - G(p|e)] dp}{\left[\frac{1}{u'(m)} \right]' \int_0^m [1 - G(p|e)] dp + \left[\frac{1}{u'(m)} \right]' \int_m^{\bar{p}(e)} [1 - G(p|e)] dp} \\ &= 1 - H_1(m). \end{aligned}$$

So from (21), (22), and (23), it follows that

$$\frac{\frac{C_{me}(e, m, u_0)}{u'(m)} + C_{u_0e}(e, m, u_0)}{\frac{C_m(e, m, u_0)}{u'(m)} + C_{u_0}(e, m, u_0)} \leq \frac{C_e}{C}.$$

For the case where $\frac{1}{u'(0)} = 0$ note that (22) holds as an equality. But then, when $\left[\frac{1}{u'(p)} \right]'$ is increasing, we have that $\frac{h_2}{1 - H_2(m)}$ stochastically dominates $\frac{h_1}{1 - H_1(m)}$. And, it similarly follows that $1 - H_2(m) \geq 1 - H_1(m)$ and so we have

$$\frac{\frac{C_{me}(e, m, u_0)}{u'(m)} + C_{u_0e}(e, m, u_0)}{\frac{C_m(e, m, u_0)}{u'(m)} + C_{u_0}(e, m, u_0)} \geq \frac{C_e}{C}. \quad \blacksquare$$

Appendix II: Numerical Examples

In this appendix we provide two detailed examples to illustrate the effect of changes in m and u_0 on e^* . To construct the examples we use the numerical algorithm introduced in Kadan and Swinkels (2009a). This algorithm is based on the constructive existence proof from Jewitt, Kadan, and Swinkels (2008). It enables us to numerically find cost minimizing contracts for any given effort level, which, in turn, allows us to calculate $C(e, m, u_0)$, and solve for the optimal effort level in Lemma 1. Using the contract $\pi(x)$ and the exogenously given distribution $F(x|e)$, it is also easy to calculate the distribution of pay $G(p|e) = F(\pi^{-1}(p)|e)$. This allows us to check whether $\frac{-G_e(p|e)}{1-G(p|e)}$ is increasing in p .

The examples use the FGM copula distribution to express the relation between effort and signals: $f(x|e) = 1 + 0.5(1 - 2x)(1 - 2e)$. This distribution is convenient since it satisfies both *MLRP* and *CDFC*, implying that the FOA is valid. The validity of the FOA is useful since the algorithm depends on this feature. Under this specification, the cost minimizing contract for a given effort level is given by (see Jewitt, Kadan, and Swinkels (2008)):

$$\pi(x) = \max \left\{ m, u^{t-1} \left(\frac{1}{\lambda + \mu \frac{f_e(x|e)}{f(x|e)}} \right) \right\}, \quad (24)$$

where λ and μ are multipliers calculated using the algorithm mentioned above.

Example 1 *A case where IR is not binding and M is binding.*

Assume $u(p) = \frac{1}{2}\sqrt{p}$, for which $\frac{1}{u}$ is concave. Assume also that $u_0 = 2$, $m = 2$, and $c(e) = e^2$. Figure 2 depicts several graphs describing the behavior of the cost minimizing contracts under this specification. First, the top left graph describes the cost minimizing contracts for $e = 0.2$ (blue), $e = 0.35$ (red), and $e = 0.5$ (green). It can be seen that the contract is “option like” in line with (24). In this example, *IR* is not binding at any of the effort levels ($\lambda = 0$). When effort goes up, pay goes up at all effort levels. Consequently, the distribution of pay for higher effort dominates that for a lower effort in the sense of *FOSD*. This can be easily observed from the top right graph, which depicts the distribution of pay $G(p|e)$ for the corresponding effort levels. The bottom left graph plots $\frac{-G_e(p|e)}{1-G(p|e)}$. This magnitude is monotone increasing in all three cases.

Since *IR* is not binding, we know from part (i) of Theorem 2 that $e_m^* > 0$. That is, the effort level that minimizes $\frac{C(e, m, u_0)}{e}$ goes up, when minimum wage goes up. To illustrate this point, the bottom right graph of Figure 2 shows the effect of changes in m on the minimum efficient scale effort level. Recall that this effort level can be plotted as the intersection point between the marginal cost $C_e(e, m, u_0)$ and the average cost $\frac{C(e, m, u_0)}{e}$ curves. We plot these two curves for two levels of minimum wage $m = 2$ (blue) and $m = 3$ (red). Note that marginal cost is higher for the higher minimum wage reflecting that $C_{me} > 0$. Despite this, as expected from Theorem 2, the crossing point moves to the right when m is increased.

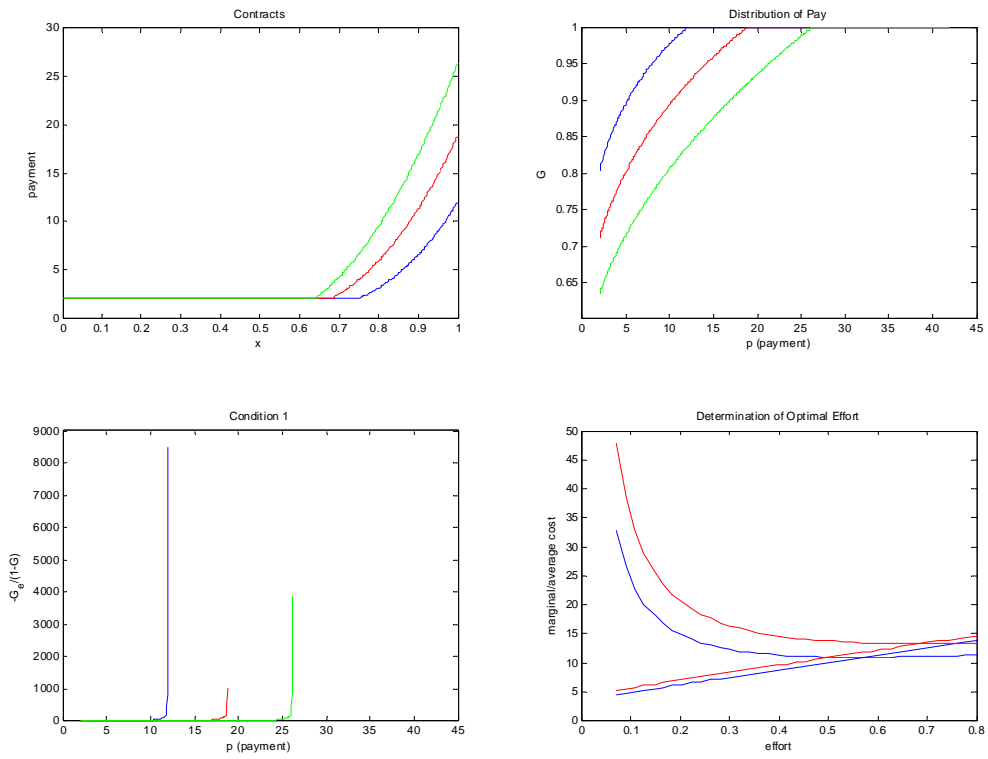


Figure 2: Graphs for Example 1

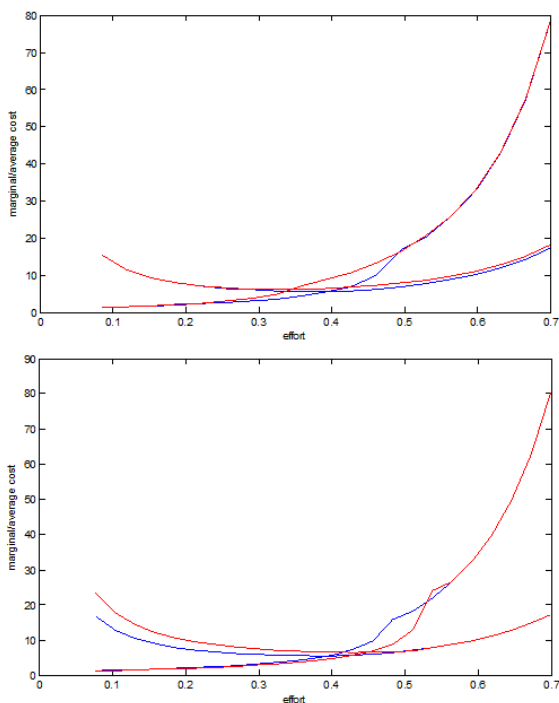


Figure 3: Graphs for Example 2

Example 2 *A case where both M and IR are binding.*

Assume $u(p) = p + \ln(p)$, for which $\frac{1}{u'}$ is concave. Assume also that $u_0 = 0.5$, $m = 0.05$ and $c(e) = \frac{e}{1-e}$. In this case, both the IR and the M constraints are binding at the relevant effort levels. Figure 3 illustrates what happens to the optimal effort level when either m or u_0 are increased. Each of the graphs depicts the marginal and average cost curves for given parameter values. The crossing point of this curves is the optimal effort level. In the top graph, the minimum wage is increased from 0.05 (blue) to 0.2 (red). It can be seen that the crossing point moves to the left, implying that optimal effort has declined. Note that this does not contradict Theorem 2 since both M and IR are binding. In the bottom graph, the outside option is increased from 0.5 (blue) to 1.5 (red), keeping m constant. The marginal cost of inducing effort then goes down at the relevant range, and optimal effort goes up. Note the kinks in the cost functions. These reflect effort levels where one of the constraints (M or IR) starts or ceases to bind.

References

- [1] Aguiar Mark and Erik Hurst (2007; *Measuring Trends in Leisure: The Allocation of Time over Five Decades*, **The Quarterly Journal of Economics**, 122, 969-

1006.

- [2] Baker, George, Michael Gibbs, and Bengt Holmström (1994a); *The Internal Economics of the Firm: Evidence from Personnel Data*, **The Quarterly Journal of Economics**, 109, 881-919.
- [3] Baker, George, Michael Gibbs, and Bengt Holmström (1994b); *The Wage Policy of a Firm*, **The Quarterly Journal of Economics**, 109, pp. 921-955.
- [4] Baranchuk, Nina, Glenn MacDonald, and Jun Yang (2009); *The Economics of Super Managers*, Working paper, Washington University in St. Louis.
- [5] Bartel, Ann, Casey Ichniowski and Kathryn Shaw (2004); *Using “Insider Econometrics” to Study Productivity* **American Economic Review**, 94, 217-223.
- [6] Card, David, and Alan B. Krueger (1994); *Minimum Wages and Employment: A Case Study of the Fast-Food Industry in New Jersey and Pennsylvania*, **American Economic Review** 84, 772-93.
- [7] Card, David, and Alan Krueger (1995); *Myth and Measurement; the New Economics of the Minimum Wage*, **Princeton University Press**.
- [8] Converse, Muriel, Richard Coe, Mary Corcoran, Maureen Kallick and James Morgan (1981); *The Minimum Wage: An Employer Survey*, **Report of the Minimum Wage Study Commission**, U.S. Government Printing Office, Volume VI.
- [9] De Fraja, Gianni (1999); *Minimum Wage Legislation, Productivity and Employment*, **Economica**, 66, 473-488.
- [10] Gicheva, Dora (2009); *Working Long Hours and Early Career Outcomes in the High-End Labor Market*, working paper, Yale University.
- [11] Griliches, Zvi (1994) *Productivity, R&D, and the Data Constraint*. **American Economic Review**, 84, 1-23.
- [12] Grossman, S., and O. Hart, (1983); *An Analysis of the Principal Agent Problem*, **Econometrica**, 51, 7-45
- [13] Holmström, Bengt, (1979); *Moral Hazard and Observability*, **Bell Journal of Economics**, 10, 74-91.
- [14] Holmström, Bengt, (1982); *Moral Hazard in Teams*, **Bell Journal of Economics**, 13, 324-340.
- [15] Holmström, Bengt, (1979); *Moral Hazard and Observability*, **Bell Journal of Economics**, 10, 74-91.
- [16] Jewitt, Ian (1988); *Justifying the First Order Approach to Principal-Agent Problems*, **Econometrica**, 56, 1177-1190.

- [17] Jewitt, Ian, Kadan, Ohad, and Jeroen M. Swinkels, 2008; *Moral Hazard with Bounded Payments*, **Journal of Economic Theory**, 143, 59–82.
- [18] Kadan, Ohad, and Jeroen M. Swinkels, (2008); *Stocks or Options? Moral Hazard, Firm Viability, and the Design of Compensation Contracts*, **Review of Financial Studies** 21, 451–482.
- [19] Kadan, Ohad, and Jeroen M. Swinkels, (2009a); *Minimum Payments and Effort Choices in Moral Hazard Models*. Available at http://www.kellogg.northwestern.edu/Faculty/Directory/Swinkels_Jeroen.aspx.
- [20] Kadan, Ohad, and Jeroen M. Swinkels, (2009b); *Numerical Exploration of the Moral Hazard Problem with Bounded Payments*. Available at http://www.kellogg.northwestern.edu/Faculty/Directory/Swinkels_Jeroen.aspx.
- [21] Kadan, Ohad, and Jeroen M. Swinkels, (2010); *On the Moral Hazard Problem without the First-Order Approach*. Available at http://www.kellogg.northwestern.edu/Faculty/Directory/Swinkels_Jeroen.aspx.
- [22] Lazear, Edward, (1986); *Salaries and Piece Rates*, **Journal of Business**, 59, 405–431.
- [23] Lazear, Edward (2000); *Performance Pay and Productivity*, **American Economic Review**, 90, 1346–1361.
- [24] Lemieux, Thomas, W. Bentley Macleod and Daniel Parent (2009); *Performance Pay and Wage Inequality*, **Quarterly Journal of Economics** 141, 1–49.
- [25] Mirrlees, J., (1976); *The Optimal Structure of Incentives and Authority within an Organization*, **Bell Journal of Economics**, 7, 105–131.
- [26] Mirrlees, J., (1999); *The Theory of Moral Hazard and Unobservable Behavior: Part I*, **Review of Economic Studies**, 66, 3–21.
- [27] Neumark, David and William Wascher (2000); *Minimum Wages and Employment: A Case Study of the Fast-Food Industry in New Jersey and Pennsylvania*, **American Economic Review**, 90, 1362–1396.
- [28] Neumark, David and William Wascher (2007); *Minimum Wages and Employment*, Discussion Paper No. 2570, IZA.
- [29] Prescott, Edward C., and Robert M. Townsend (1984), *Pareto Optima and Competitive Equilibria with Adverse Selection and Moral Hazard*, **Econometrica**, 52, 21-45.
- [30] Rogerson, William P., (1985); *The First-Order Approach to Principal-Agent Problems*, **Econometrica**, 53, 1357–1367.
- [31] Rustichini, Aldo, and Paolo Siconolfi (2008); *General Equilibrium in Economies with Adverse Selection*, **Economic Theory**, 37, 1-29.

- [32] Shapiro, Carl, and Joseph E. Stiglitz (1984); *Equilibrium Unemployment as a Worker Discipline Device*, **The American Economic Review**, 74, 433-444.
- [33] Zame, R. William (2007); *Incentives, Contracts, and Markets: A General Equilibrium Theory of Firms*, **Econometrica**, 75, 1453-1500.