# See No Evil: How a company can benefit from *not* observing behavior

By: Daniël Springer, Erasmus University Rotterdam

**Summary:**

This thesis proposes a model that shows how companies may be better off when not observing certain behavior from certain employees, even if that behavior contains relevant information. The reason for this is that employees may start distorting their behavior so actively to attain future rewards, that the added value of the information is exceeded by the costs of these distortions.

# Chapter 1: Introduction

When Jeremy Bentham conceived his "ideal prison", the Panopticon, he stated it would provide "a new mode of obtaining power of mind over mind". The idea was that by making sure all the behavior of prisoners could be observed, prisoners would stop misbehaving. A similar reaction can be seen by anyone who has found himself in the presence of a police car. At once, even the greatest road hog will start driving excessively well-mannered. In both cases, the *potential* for the guards/police to use what they see against the perpetrator is enough to (positively) influence behavior.

Obviously, the same logic holds in companies. In environments where output can be (almost) perfectly observed (e.g. tomato picking), companies often manage to significantly increase production by offering some sort of contingent wage. Many authors (e.g. Mills, 1985) have found, however, that companies don't always use all available information when deciding on whom to promote or whom to give a bonus. More often than would be expected, seniority, not performance, decides who gets a certain promotion or bonus.

In my thesis, I offer a potential explanation for this behavior. In all the examples I gave above, the "distortion" of the behavior of the person who is being observed is positive. However, this is not necessarily always true. In an environment with imperfect information, it could be the case that the company can induce information from certain behavior. A famous example of this would be "completing university" (behavior) which can signal "intelligence" (the information) (Spence, 1973).

## 1.1: Signaling games and companies

Obviously, these kinds of "signaling games" are constantly played in companies. For instance, employees may work longer to signal commitment to the company. Potentially, however, when certain signals are rewarded, employees could start producing a suboptimally large amount of them. In that case, it would indeed be optimal for the company not to observe the signal, or at least not use it. Promoting based on seniority could be way to do this.

At the same time, my thesis tries to answer a question that has personally bugged me for a long time. While I was working at an accounting department, my senior would routinely refuse to delegate certain tasks to me that I felt were quite simple and a waste of his time. I never understood this behavior. However, in the logic of the model I propose, my senior might just have been signaling that he was a hard worker to *his* boss.

In spirit, my thesis is close to literature on the distorting effects reputational concerns can have. Specifically, the effect I try to show is quite similar to that described by Morris (2001). In his paper, a policy advisor with the same preferences as the policy maker sometimes distorts his signal to "distance" himself from a "bad" type of policy advisor, the two of which the policy maker can't distinguish. On the other hand, the "bad" type of advisor is induced to truthfully reveal his signal more often.

The same dynamic is shown in Ottaviani and Sorensen (2006), who show that when experts receive pay-offs according to their reputation, truthfully revealing their signals is not an equilibrium. Also here, the process of optimizing utility over more than one period when some outside observer uses the information provided can lead to a suboptimal outcome.

## 1.2: Career concerns

In the specific model I present, the reason why the agent cares about his reputation is concern for his career. The information provided by the agent's actions ultimately serves to

determine whether or not a promotion is given. In this way, my thesis is related to the broad range of articles on career concerns in environments of imperfect information. Holmström (1979) initially posited that increased amounts of information could ameliorate the effects of the principal agent problem. However, in Holmström (1999), the same author shows that as the quality of the signal a principal receives about the quality of an agent increases, the agent may work less. The reasoning behind this result is that lower uncertainty decreases the agent's ability to influence the principal's perception of his quality by working harder.

Crémer (1994) shows that high-powered incentives, such as firing unproductive employees, may become unfeasible as information about the agent becomes better. The rationale behind this result is that if there is no full correlation between output and effort, it might be ex-ante optimal for the principal to judge agents only on output. Ex-post, however, sufficiently "cheap" information about effort could induce the principal to renegotiate on the "fire when output is low" policy. Knowing this, the agent will presume such a move and not (fully) act on the intended incentive. In such cases, keeping monitoring costs high can be optimal for the firm.

These studies share the result that with extra information, sorting (discriminating between types of agent) is improved, but discipline (effort) is reduced. My thesis focuses not so much on reduced effort, but more on a distortion in decision making by the agent. Along these lines, Prat (2005) proposes a model that shows a similar effect. In his model, an agent with a private signal may disregard that signal in order to increase his reputation with his principal. Prat, like my thesis, suggests that these issues could be solved by not observing the action of the agent, but merely its consequences.

My thesis differs from Prat's paper in numerous ways. Firstly, my model shows dynamics that occur in *more than two* hierarchical levels in the organization and proposes a link to delegation behavior of mid-level managers. Secondly, agents (seniors) in my model differ in the private costs they associate with certain behavior, not their ability. Thirdly, in my model, an agent has to balance between the effect his decision has on his wage through company profits and the benefits of a potential promotion, whereas in Prat, pay-offs for the agent only depend on his reputation. Lastly, my model introduces possible promotions as a means of rewarding the "good" type of employees.

In general, my model adds several distinct characteristics that I feel replicate quite a few real life situations. It is more detailed than previous models and therefore reaches some more nuanced conclusions about whether or not observing behavior is optimal. On the other hand, it also replicates effects found in previous model studies in a substantially different model, thereby lending more credibility to those models.

My model can be characterized as a two-period signaling game with costly signals.


# Chapter 2: model and results

## 2.1: Model

### *2.1.1: Relevant actors*

In this model, three actors are relevant. These are: the manager/owner of the company, denoted as M, a senior employee, denoted as S, and a junior employee, denoted as J. All actors are rational and optimize their own (expected) utility. Since M is the owner of the company, the terms "M" and "the company" are used interchangeably. There are two types of S, $S_G$ and $S_B$, which differ in their effort averseness, which is expressed through their private costs of effort, $c_P$. I will discuss this difference when I explain the decision variables. The ex-ante probability of a random S being $S_W$ ($P(S=S_W)$) equals $\mu$, so that $P(S=S_L) = 1-\mu$.

J's also differ from each other in one aspect: their ability to assess the profitability of a certain project. This ability is denoted as $\lambda$, with $\lambda=\frac{1}{2}$ meaning "J has a 50% chance of assessing the pay-off from a certain project correctly". $\lambda$ is uniformly distributed on the interval $[\frac{1}{2},1]$.

## *2.1.2.: Available information*

This model is characterized by information asymmetry between M and S. There are in essence two important pieces of information in this model, which are: S's type and J's ability ($\lambda$). Prior to all decisions made, M has no information about either of these variables, other than the ex-ante probabilities and distributions.

On the other hand, S knows his own type and knows with certainty the ability of his junior. For the purpose of this thesis, it is irrelevant whether or not J knows his own type, since J acts according to a fixed set of rules (see 2.1.4). Apart from these pieces of information (S's type and J's $\lambda$), all information is common knowledge.

## *2.1.3: Decisions and sequence*

The sequence of decisions and events in this model is as follows. Nature first assigns a type to S and an ability level $\lambda$ to J. S observes both his own type and J's $\lambda$. After learning these values, S has to decide on delegating the authority to make a certain decision. The decision variable is $I \in 0,1$. The pay-off in period 1 is denoted as $\pi \in \{ \pi_L, \pi_H \}$ with $\pi_L<0<\pi_H$. It depends on the state of the world $SoW \in \{ 0, 1 \}$ and the choice of I. When $\Sigma =$I, $\pi = \pi_H$ whereas when $\Sigma \neq$ I, $\pi = \pi_L$.

Denote the delegation decision as $D \in \{ 0, 1 \}$, with D=0 meaning "no delegation" (S decides) and D=1 meaning "delegation" (J decides). When S chooses D=1, J investigates the decision at hand and receives a (noisy) signal about the state of the world, $\sigma \in \{ 0, 1 \}$ which is correct with probability $\lambda$. Based on his signal, J then decides which value for I to choose.

When, however, S chooses D=0, he investigates the decision himself and learns the state of the world, SoW, *with certainty*. He can then decide on I. Having to investigate the decision himself does, however, entail certain costs. In all scenarios, there are costs to the company when D=0[1]. Denote these costs as $c_C(D)$, with $c_C(D=0)>0$ and $c_C(D=1)=0$. S also incurs private costs (of effort) when D=0, denoted as $c_P(S)$. $c_P$ is a function of S's type, with $c_P(S_G)=0$ and $c_P(S_B)>0$.

After D and I have been chosen, M observes these actions and updates his beliefs about S's type. He then decides whether or not to promote S, $P \in \{ 0, 1 \}$. M wants to promote S if, and only if, $S=S_G$ and is therefore better off if he can distinguish between the two types as much as possible. When promoted, S receives a higher wage, the expected value of which is denoted as w. I assume w to be beyond M's control, which corresponds to the case in which there is sufficient competition in the labor market to force M to offer a certain wage to retain employees. The reason for this assumption is that allowing M to choose w would mean an entirely new strategic dimension would be added, which is beyond the scope of my thesis.

---

[1] Think of these as the opportunity costs to the company of assigning the senior employee to project A instead of using his expertise on project B or the costs of paying overtime to S.

Sequence of events



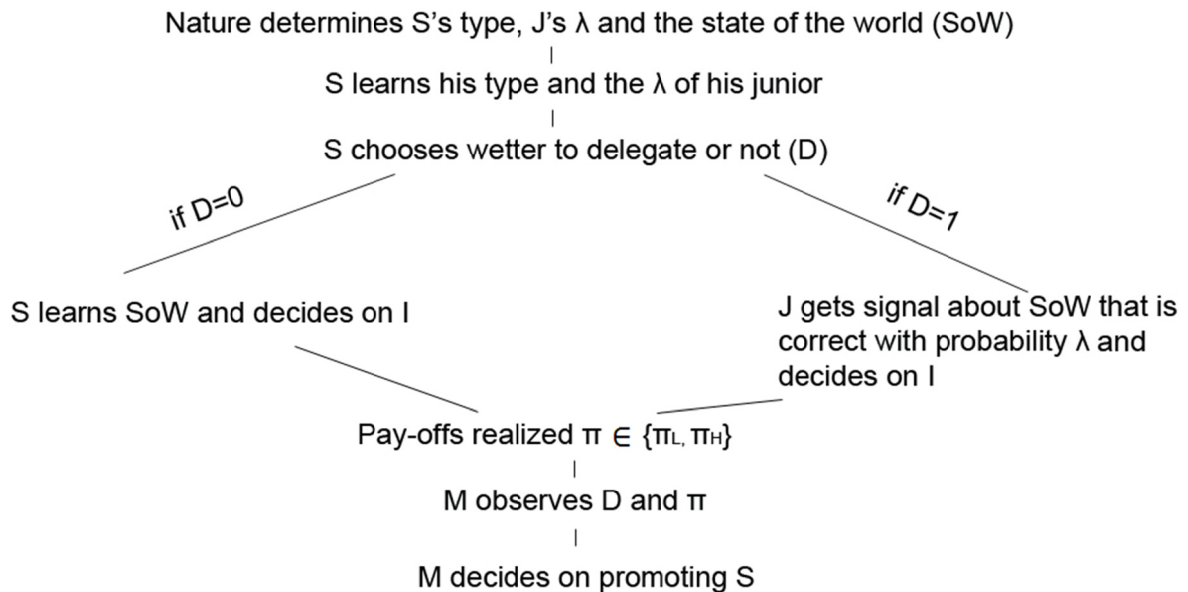Nature determines S's type, J's λ and the state of the world (SoW)
|
S learns his type and the λ of his junior
|
S chooses wetter to delegate or not (D)

*if D=0*

*if D=1*

S learns SoW and decides on I

J gets signal about SoW that is correct with probability λ and decides on I

Pay-offs realized π ∈ {πL, πH}
|
M observes D and π
|
M decides on promoting S

Fig 1: an overview of the sequence of events described above

### *2.1.4.: Utility functions and strategies*

I will now go through all actors in this model and assess their utility functions (where relevant) and strategies:

**J** has only one decision he might be asked to make, which is I if D=1. J makes his decision based on his signal about the state of the world (SoW), σ. His behavior can therefore be characterized by the following strategy mapping:

I: {0,1} → {0,1}

With I(σ) being the decision given the signal received. To focus on the relevant mechanism, assume that J's optimal strategy is:

$$I(\sigma) = \begin{cases} 0 \text{ if } \sigma = 0 \\ 1 \text{ if } \sigma = 1 \end{cases}$$

That is: J is not a strategic actor, but merely a "machine" that always acts according to his own signal. Whether or not J knows his own λ is therefore irrelevant.

**S** must always decide on D, and if D=0, also on I. This last decision is, obviously, not very relevant, since S knows the value of Σ with certainty and will simply choose I=1 if Σ=1 and I=0 if Σ=0. More importantly, S has to decide on D, based on his knowledge of λ. S wants to maximize his utility, which equals:

$$U(P,S,D) = g(\pi - c_C(D)) + w(P) - c_P(S)$$

With g ∈ [0,1] [2], $c_C$=0 if D=1 and w=0 if P=0. S's behavior can then be described by the following strategy mapping:

$$D: [0,1] \rightarrow \{0,1\}$$

With D(λ) being the delegation decision given the ability of the junior λ.

**M** has to decide on P: whether or not to promote S given the information he received about D and, if D=1, about π. Formally:

$$P: \{0,1\} \times \{\pi_L, \pi_H\} \rightarrow \{0,1\}$$

Where P(D,π) is the promotion decision given S's choice of D and first-period profits π. Assume that M's (simplified) utility function for the two periods is of the form:

$$U_M(P, S, D) = (1 - g)(\pi - c_C(D)) + P(h(S) - w)$$

Where h is some constant > 0 that depends on the type of S, with $h(S_G) > h(S_B)$. To make the promotion decision important, assume that $U_M(1, S_B) < 0$ and $U_M(1, S_G) > 0$. The interpretation of this is that the benefits M receives from promoting S are greater than the wage increase when S=$S_G$, but smaller when S=$S_B$.


## 2.2: Results without reputational concerns

First of all, I will consider a model in which there is no second period and therefore, no reputational concerns. These results are used as benchmarks against which to compare the results with a second period.

Proposition 1:
When S can't get a promotion, there are values of λ for which $S_B$ and $S_G$ are indifferent between choosing D=0 and D=1. This cut-off point is suboptimal from the company's perspective for S=$S_B$, but optimal when S=$S_G$.

Proof:
The fact that $S_B$ incurs costs when choosing D=0 that are irrelevant for the company ensures that he chooses D=1 more often than $S_G$. Since $S_G$'s strategy was optimal from the company's perspective, therefore, $S_B$'s strategy is necessarily not.

See appendix for algebraic proof.

This is a very intuitive result. $S_G$ has no private costs associated with D, but fully depends on the company outcome. With interest aligned in such a manner, it is hardly surprising for $S_G$ to act in the company's best interest. $S_B$, on the other hand, does have private costs associated with a certain choice of D, and will therefore be less inclined to choose D=0 than the company would like him too.

As a benchmark, it would be interesting to see what the first-best outcome of this game would be from a social welfare perspective (that is: taking the utility of both M and S into account).

---

[2] Ergo, S receives a certain percentage of the profit made under his supervision.

Proposition 2:
From a social welfare perspective, it would be optimal for $S_B$ to increase the threshold level above which he chooses D=1. $S_G$'s optimal strategy from a social welfare perspective is the same as under proposition 1.

Proof:
See appendix

## 2.3: Results with reputational concerns

Having looked at the rather straight-forward game without the influence of a second period, I will now look at whether or not these results also hold when M can indeed promote (certain) S's. First, I will look at the impact S's decision has on his reputation.

Proposition 3:
Choosing D=0 increases S's reputation, which is M's belief about the probability that $S=S_G$, while choosing D=1 decreases it as long as for at least one type of S, the optimal choice of D depends on the value of $\lambda$. When that condition hold and D=1, the outcome $\pi = \pi_H$ gives a higher reputation for S than the outcome $\pi = \pi_L$.

Proof:
See appendix

### 2.3.1: Model with commitment device

Before looking at behavior in equilibrium, I will now make a small side-step in the form of propositions 4 and 5. Proposition 4 will analyze the case in which M can ex-ante commit himself to a certain promotion policy. This would in effect reverse the decision-making process and allow M to determine the outcome of the game. This model is relevant for two reasons. First of all, it shows what would happen in case M can commit himself, which can be a reasonable assumption when contracts are sufficiently good. Secondly, the model shows the trade-off between making a more informed decision and distorting behavior in a potentially harmful way.

Proposition 4:
When M can credibly commit himself to a certain promotion policy, he faces a trade-off between not distorting first period behavior and being able to use the information conveyed by S's choice of D.

Proof:
See appendix

The distortion caused by using the information conveyed in S's choice of D is not necessarily negative. To see why, consider that both types of S will choose D=0 more often, which is (from M's perspective) bad if $S=S_G$ but could be either good or bad if $S=S_B$ (see proposition 1). Proposition 5 algebraically proves this intuition for one outcome of the model discussed in proposition 4. Note that the result also holds for proposition 7, since the distortion is the same in a non-commitment model with semi-separating equilibria[3].

Proposition 5:
The distortion caused by using the information conveyed by S's choice of D can either be positive or negative for M.

---

[3] However, in that case, M doesn't face a trade-off, since he can't credibly commit to a certain policy and therefore gets "forced" into a certain outcome

<u>Proof:</u>
See appendix

<u>*2.3.2: Model without commitment device*</u>

With propositions 4 and 5 in mind, I will now establish equilibrium behavior in the "original" model, in which M chooses his promotion policy after S has chosen D. I will limit myself to pure strategy equilibria, since they suffice for the effect I'm trying to demonstrate.

<u>Proposition 6:</u>
There are conditions under which there is a pooling equilibrium, in which S always chooses the same D and M always chooses the same P.

<u>Proof:</u>
See appendix

In the pooling equilibria, obviously, there is no information conveyed by S's choice of D. Therefore, the first-period result is the same as without a second period and M has no information to act on besides the ex-ante probability that $S=S_G$. I will now look at the case where for at least one type of S, the optimal choice of D depends on $\lambda$.

<u>Proposition 7</u>
There exist multiple semi-separating equilibria

<u>Proof</u>
See appendix

# Chapter 3: discussion of results

The results that are derived from the model have some interesting features. First of all, this model is obviously limited in its direct applicability, since it describes a very specific situation within a company. That doesn't mean, however, that some of the results don't indicate more generalizable tendencies.

The model shows that under imperfect information, when (mid-level) managers differ in some dimension, certain decisions they make may reveal information about them. In this specific case, the dimension is effort-averseness, and the decision relates to delegating to a subordinate. However, similar models could be constructed with, for instance, differences in ability.

When the behavior is not observed, or is observed but not acted on, the fact that the behavior contains information is not important. However, once a higher level manager / the company starts using the information provided by the behavior, this will in turn affect the behavior. The key point that this thesis shows, is that depending on the exact situation, this "distortion" in behavior could be either positive or negative. As an educated guess, I would say that in this particular model, it is negative for most reasonable sets of values.

That, however, is still not the end of the story. First of all, when the distortion is negative, that doesn't necessarily mean that the company would be better off by not observing the behavior. The reason for that is, as explained in proposition 3, that the information provided by the behavior could increase the chance of making the correct decision. In these cases, the company clearly faces a trade-off between acting on superior information and ensuring

that mid-level managers don't distort their behavior. The result of this trade-off depends on the values of all variables.

However, my results also show that without any commitment devices, companies don't actually have this choice. The equilibria in proposition 7 show situations in which M's ex-post optimal reaction is to use the signal provided by S's choice of D, which doesn't mean this would also be his optimal outcome. Similarly, proposition 6 shows equilibria in which there is no information conveyed by S's choice, even though an ex-ante commitment to a certain promotion policy may induce S to act otherwise.

Since a company often chooses ex-post on whom to promote, the resulting equilibrium may be one in which the signal provided is used, even if this is sub-optimal from the company's point of view, or vice-versa. This could be due to imperfect contracts or external factors (law, unions) that prohibit contracts that commit the company to a specific promotion policy. When the rationality assumption is relaxed, it may also be due to a failure on the part of the company.to recognize the distorting effects a certain promotion policy has.

### 3.1: 'Political correctness'

As mentioned, the behavior that I show in a company context resembles that shown by Morris (2001) in his article on political correctness. When a certain behavior is more likely to be associated with a negative trait, an advisor/manager will refrain from this behavior more often. As in Morris' article, my results show that this situation can be harmful to the company (or in Morris: the policy maker). Note that, given proposition 2, the distortion is even more likely to negatively influence social welfare.

In addition to this logic, my paper illustrates that companies who expect high costs from this sort of distorting behavior may look for policies to reduce it. This would be the case in, for instance, industries where there exist large uncertainties about some quality of employees. In these cases, contractually agreeing to promote someone after a set amount of years could indeed be a solution. Another option would be to simply limit the amount of possible behaviors. In my model, for instance, M could force S to always delegate the decision to J.

All in all, I feel my thesis proposes a valid reason for companies to sometimes not use all the information at their disposal when assessing their employees. It also provides interesting insights into the effect of signaling on the work floor and serves as a warning for companies not to underestimate the effect that observation can have on behavior. Lastly, my thesis hints at an interesting potential role of delegation decisions, as a signaling device for employees.

### 3.2: Further research

To my knowledge, the effect that I describe and that has popped up in the literature in several forms has never been empirically tested. A first step to test my model would be to look at whether or not companies that face high uncertainty about their employee's abilities are more likely to promote based on something else than behavior or use means to not collect data on behavior in the first place.

In terms of theory, it would be interesting to try and fit the model even closer with reality. That would mean, for instance, that a promotion decision is probably not made based on one incident. When agents are judged not based on one action, but on a continuous stream of decisions they have to make, a reasonable guess would be that some of the problems indicated in the literature could evaporate, since the incremental effect of each decision on reputation is lowered.

Lastly, it would be interesting to see what happens when the Junior employee in my model becomes more than a machine-like transmitter of his signal, but also has strategic concerns about appearing competent. This would obviously require uncertainty on the Senior's behalf about J's type. The complexity of the model would greatly increase, but might be very apt at replicating situations like those found in hospitals, where all levels of employees face some uncertainty about the abilities of their subordinates and constantly have to make delegation (and promotion) decisions.

# Appendix

Proof for proposition 1

S's optimal strategy is found by finding the conditions under which his incentive constraint is met. More precisely, I look for the value of $\lambda$ for which S is indifferent between choosing D=0 and D=1.

Expected pay-off without delegation:

$$g(\pi_H - c_C(D = 0)) - c_P(S)$$

Expected pay-off with delegation

$$g(\lambda \pi_H + (1 - \lambda)\pi_L)$$

So S delegates when

$$g(\lambda \pi_H + (1 - \lambda)\pi_L) \geq g(\pi_H - c_C(D = 0)) - c_P(S)$$

Rearranging the incentive constraints gives the optimal strategy for S as a function of $\lambda$:

$$D(\lambda) = \begin{cases} 0 \text{ if } \lambda < \lambda^- \\ 1 \text{ if } \lambda \geq \lambda^- \end{cases}$$

where

$$\lambda^- = \frac{\pi_H - \pi_L - c_C(D = 0)}{\pi_H - \pi_L}$$

Using the company's incentive constraint, it can be shown that this is also the optimal strategy for S from the company's perspective. $S_B$'s (individually) optimal strategy, however, is

$$D(\lambda) = \begin{cases} 0 \text{ if } \lambda < \lambda^- \\ 1 \text{ if } \lambda \geq \lambda^- \end{cases}$$

where

$$\lambda^- = \frac{g(\pi_H - \pi_L - c_C(D = 0)) - c_P(S)}{g(\pi_H - \pi_L)}$$

Proof for proposition 2

Optimizing from a social welfare perspective basically requires adding up the utilities of both S and the company and finding the value of $\lambda$ for which D=0 and D=1 give the same amount of total utility. Since I've shown in the proof of proposition 1 that $S_G$ has the exact same

optimal strategy as the company, sticking to this strategy is obviously also optimal from a social welfare perspective.

For $S_B$, find the socially optimal strategy requires adding up the incentive constraints for S and the company. These are:

$$g(\lambda\pi_H + (1-\lambda)\pi_L) \geq g(\pi_H - c_C(D=0)) - c_P(S)$$
for $S_B$ and

$$(1-g)(\lambda\pi_H + (1-\lambda)\pi_L) \geq (1-g)(\pi_H - c_C(D=0))$$
for the company

In total:

$$\lambda\pi_H + (1-\lambda)\pi_L \geq g(\pi_H - c_C(D=0)) - c_P(S) + (1-g)(\pi_H - c_C(D=0))$$

Which means that from a social welfare perspective, $S_B$ should choose D=1 if:

$$\lambda^- \geq \frac{\pi_H - \pi_L - c_C(D=0) - c_P(S)}{\pi_H - \pi_L}$$

and D=0 otherwise. As can be expected, therefore, $S_B$ should choose D=0 more often if the social optimum is to be attained.

<u>Proof for proposition 3</u>

First of all, it should be noted that there are essentially five possible combinations of strategies of the two types of S. In two of these sets, neither type varies his choice of D according to $\lambda$, meaning they both either always delegate, or never delegate. Obviously, both these options ensure that S's decision reveals no information about his type, and is therefore irrelevant to his reputation.

Two interesting cases are those in which either both types choose D=0 for at least some values of $\lambda$ and D=1 for at least some values of $\lambda$, or only one type does. In both these cases, however, the cut-off point for $\lambda$, above which S wishes to delegate, will be higher for $S_G$ than for $S_B$. Formally: $\lambda_G^- \vee \lambda_B^- \in \left(\frac{1}{2}, 1\right) : 1 > \lambda_G^- > \lambda_B^- > 0$. To see why, just consider that this cut-off point reflects indifference between delegating and not delegating, and that $S_B$ has extra costs involved when choosing D=0 which $S_G$ doesn't.

Denoting the cut-off points for $\lambda$ as $\lambda_G^-$ for $S_G$ and $\lambda_B^-$ for $S_B$[4], I use Bayes' rule to determine M's updated beliefs about S's type when at least one type differentiates his choice of D according to the value of $\lambda$ (note that the value of $\pi$ conveys no information when D=0, so that I omit it in that case):

$$\Pr(S = S_G \,|\, D = 0) = \frac{\mu(2\lambda_G^- - 1)}{\mu(2\lambda_G^- - 1) + (1-\mu)(2\lambda_B^- - 1)} \qquad (1)$$

$$\Pr(S = S_G \,|\, D = 1 \text{ and } \pi = \pi_L) = \frac{\mu(\lambda_G^- - 1)^2}{\mu(\lambda_G^- - 1)^2 + (1-\mu)(\lambda_B^- - 1)^2} \qquad (2)$$

---

[4] Note that $\lambda_G^-, \lambda_B^- \in [\frac{1}{2}, 1]$

$$\Pr(S = S_G \mid D = 1 \text{ and } \pi = \pi_H) = \frac{\mu(1 - (\lambda_G^-)^2)}{\mu(1 - (\lambda_G^-)^2) + (1 - \mu)(1 - (\lambda_B^-)^2)} \qquad (3)$$

Since $\lambda_G^- > \lambda_B^-$ and $0 < \mu < 1$, it can easily be shown that (1)>µ, (2)<µ and (3)<µ. Therefore, choosing D=0 increases the perceived likelihood that S=$S_G$, whereas choosing D=1 decreases it, regardless of the actual profits made during the first period. Another very intuitive result, which is slightly less straightforward to prove, is that (3) > (2). To see why, consider that:

(3) − (2) = $\mu(1 - \mu)(1 - \lambda_G^-)(1 - \lambda_B^-)(2(\lambda_G^- - \lambda_B^-))$

Which is more than 0 since $\lambda_G^- > \lambda_B^-$. Thus, the fact that J made the correct decision after D=1 increases the perceived likelihood that S=$S_G$

The last case is that in which $S_G$ always chooses D=0 and $S_B$ always chooses D=1. In this case, there is perfect separation, so that

$$\Pr(S = S_G \mid D = 0) = \Pr(S = S_B \mid D = 1) = 1$$

Proof for proposition 4

First, consider the possible strategies for M. M has to decide on:

$P(D, \pi): \{0,1\} \times \{\pi_L, \pi_H\} \to \{0,1\}$

Which essentially gives him four options: always choose P=0 (strategy 1), always choose P=1 (strategy 2), choose

$P(D, \pi) = \begin{cases} 0 \text{ if } D = 1 \\ 1 \text{ if } D = 0 \end{cases}$
(strategy 3)

or choose

$P(D, \pi) = \begin{cases} 0 \text{ if } D = 1 \text{ and } \pi = \pi_L \\ 1 \text{ if } D = 1 \text{ and } \pi = \pi_H \\ \quad 1 \text{ if } D = 0 \end{cases}$
(strategy 4)

The first two options give the same result. When M essentially doesn't use the information conveyed in S's choice of D, S no longer has to take second period effects into account. This, obviously, means that S will simply adopt his one-period optimal strategy, as derived under proposition 1. Therefore, not using the information S's choice conveys ensures that S doesn't strategically distort his behavior when choosing D.

Secondly, I'll consider strategy 3. If this is M's strategy, then it's clear that he will promote if D=0 and not promote if D=1. To see why, consider that from M's utility function, it follows that he only wants to promote S if S=$S_G$. Since, as shown under proposition 3, $S_G$ is more likely to choose D=0 than $S_B$, the proposed strategy is the only rational way of using the information conveyed by S's decision.

When M chooses this strategy, S knows that by choosing D=0 he will get the benefits associated with promotion for sure, whereas he forfeits those benefits by choosing D=1. Therefore, his incentive constraint (for choosing D=1) becomes:

$$g(\lambda\pi_H + (1-\lambda)\pi_L) \geq g(\pi_H - c_C(D=0)) - c_P(S) + w$$

Which solves for

$$\lambda^- \geq \frac{g(\pi_H - \pi_L - c_C(D=0)) - c_P(S) + w}{g(\pi_H - \pi_L)}$$

Since by assumption, w>0, this means that S will choose D=0 more often on average because of the effect his decision has on his future wage. Therefore, choosing this strategy distorts S's first-period behavior.

The last possible strategy for M is strategy 4. Again, the results from proposition 3 make it clear that this is the only possible rational way of differentiating according to $\pi$, since the chance that $S=S_G$ when D=1 is larger when $\pi = \pi_H$ than when $\pi = \pi_L$. In this case, S's incentive constraint becomes:

$$g(\lambda\pi_H + (1-\lambda)\pi_L) + \lambda w \geq g(\pi_H - c_C(D=0)) - c_P(S) + w$$

Which solves for:

$$\lambda^- \geq \frac{g(\pi_H - \pi_L - c_C(D=0)) - c_P(S) + w}{g(\pi_H - \pi_L) + w}$$

Again, this means that S distorts his first-period strategy as compared to the scenario where the signal provided by his choice of D is not used by M.

Proof for proposition 5

Without loss of generality, I will focus on M's strategy 3 as described in proposition 4. I will first calculate the expected loss/gain from the changed behavior by both types of S:

$S_G$'s cut-off point for choosing D=1 under strategy 3 is:

$$\lambda^- > \frac{g(\pi_H - \pi_L - c_C(D=0)) + w}{g(\pi_H - \pi_L)}$$

Under strategy 1 or 2, this is:

$$\lambda^- > \frac{g(\pi_H - \pi_L - c_C(D=0))}{g(\pi_H - \pi_L)}$$

The expected damage to the company is the probability that $\lambda$ falls between the two values described above, multiplied by the average damage in that area.

The probability that $\lambda$ falls in this area equals:

$$2\left(\frac{g(\pi_H - \pi_L - c_C(D=0)) + w}{g(\pi_H - \pi_L)} - \frac{g(\pi_H - \pi_L - c_C(D=0))}{g(\pi_H - \pi_L)}\right) = \frac{2w}{g(\pi_H - \pi_L)}$$

With the remark than when one or more of the values are < 0,5 or >1, they replaced by 0,5 or 1 respectively.

What is the average damage? In this area, S now chooses D=0 instead of D=1. The damage is therefore the difference in utility for M between S's decision under strategy 1 or 2 and his decision under strategy 3:

$$g(\lambda\pi_H + (1-\lambda)\pi_L) - g(\pi_H - c_C(D=0))$$

Since the influence of $\lambda$ is linear, I can just take the average value of $\lambda$ in this domain to get to the average damage.

$$\lambda = \frac{g(\pi_H - \pi_L - c_C(D=0)) + \frac{1}{2}w}{g(\pi_H - \pi_L)}$$

So the average damage is:

$$g\left(\frac{g(\pi_H - \pi_L - c_C(D=0)) + \frac{1}{2}w}{g(\pi_H - \pi_L)}\pi_H + \left(1 - \frac{g(\pi_H - \pi_L - c_C(D=0)) + \frac{1}{2}w}{g(\pi_H - \pi_L)}\right)\pi_L\right)$$
$$- g(\pi_H - c_C(D=0))$$

$$= \frac{w}{2}$$

So in total, the expected damage when S=S$_G$ is:

$$\frac{2w}{g(\pi_H - \pi_L)} * \frac{w}{2} = \frac{w^2}{g(\pi_H - \pi_L)}$$

Secondly, I will look at what happens when **S=S$_B$**

The probability remains the same, since again the domain is stretched by

$$\frac{w}{g(\pi_H - \pi_L)}$$

In the same manner as before, I calculate the average damage in this area:

$$\lambda = \frac{g(\pi_H - \pi_L - c_C(D=0)) - c_P(S) + \frac{1}{2}w}{g(\pi_H - \pi_L)}$$

$$\left(\frac{g(\pi_H - \pi_L - c_C(D=0)) - c_P(S) + \frac{1}{2}w}{g(\pi_H - \pi_L)}\pi_H\right.$$
$$\left. + \left(1 - \frac{g(\pi_H - \pi_L - c_C(D=0)) - c_P(S) + \frac{1}{2}w}{g(\pi_H - \pi_L)}\right)\pi_L\right) - g(\pi_H - c_C(D=0))$$
$$= \frac{w}{2} - c_P(S_B)$$

Taking into account the ex-ante probability that S=S$_G$, I arrive at the expected damage/gain for M when choosing strategy 3 instead of strategy 1 or 2:

$$\frac{2w}{g(\pi_H - \pi_L)} * (\frac{w}{2} - (1 - \mu)c_P(S_B))$$

Which means that if

$$w < 2 * (1 - \mu)c_P(S_B)$$

The distortion is positive for the company, and otherwise it is negative.

Proof for proposition 6

For a pooling equilibrium to exist, both types of S must be induced to always choose the same D. Let's first assume that M's optimal strategy is to either always choose P=0 or to always choose P=1. Denote these strategies as, respectively, strategies 1 and 2. In both these cases, S can't affect P by his choice of D, so he plays his first-period optimum as derived under proposition 1.

There are obviously two compatible strategies for S (that is: strategies that lead to a pooling equilibrium): always choose D=0 or always choose D=1. First, consider the case of D=0. Since $S_G$ is on average more likely to choose D=0 (i.e.: has a higher cut-off point for $\lambda$ above which he chooses D=1), the relevant constraint is that of $S_B$. $S_B$ always chooses D=0 if:

$$\frac{g(\pi_H - \pi_L - c_C(D = 0)) - c_P(S_B)}{g(\pi_H - \pi_L)} > 1$$

Secondly, consider the case of "always choose D=1". The exact same logic applies, in reverse, so that the relevant incentive constraint is that of $S_G$. $S_G$ always chooses D=1 if:

$$\frac{\pi_H - \pi_L - c_C(D = 0)}{\pi_H - \pi_L} < \frac{1}{2}$$

Next, I will check whether or not given these strategies, it is indeed optimal for M to always choose either strategy 1 and 2. If these are the only options for M, given that S is pooling, when will M choose strategy 1? He will if

$$\mu < \frac{w - h(S_B)}{h(S_G) - h(S_B)}$$

and choose strategy 2 otherwise. There are two alternative strategies for M, which are:

$$P(D, \pi) = \begin{cases} 0 \text{ if } D = 1 \\ 1 \text{ if } D = 0 \end{cases} \quad \text{(strategy 3)}$$

or

$$P(D, \pi) = \begin{cases} 0 \text{ if } D = 1 \text{ and } \pi = \pi_L \\ 1 \text{ if } D = 1 \text{ and } \pi = \pi_H \\ \quad 1 \text{ if } D = 0 \end{cases} \quad \text{(strategy 4)}$$

Since every type of S always chooses either D=0 or D=1, both strategy 3 and 4 are not acting on extra relevant information. Therefore, both these strategies are weakly dominated by the optimal strategy 1 or 2. To see why this is true, consider the case in which S always chooses D=1. Strategy 3 is then equivalent to strategy 1, which may or may not dominate strategy 2. Strategy 4 if D=1 means that the expected pay-off is:

$$\theta(\mu h(S_G) + (1 - \mu)h(S_B) - w)$$

Which is also weakly dominated by either strategy 1 or 2. So, for all information sets that are reached, strategy 1 or 2 is optimal. Now, the only other thing that I need to compute is appropriate out-of-equilibrium beliefs for both M and S. These must be such that both M and S are not induced to deviate from their strategies. The following figure shows the possible situations and the required, arbitrary out-of-equilibrium beliefs:

**Possible situations and requirements**

| ↓ M's strategy ➜ S's strategy | Always choose D=0 | Always choose D=1 |
|---|---|---|
| **Always choose P=0** | If S chooses D=1, M must still choose P=0. If M chooses P=1, S must still choose D=0. | If S chooses D=0, M must still choose P=0 If M chooses P=1, S must still choose D=1. |
| **Always choose P=1** | If S chooses D=1, M must still choose P=1. If M chooses P=0, S must still choose D=0. | If S chooses D=0, M must still choose P=1. If M chooses P=0, S must still choose D=1. |

**Required out-of-equilibrium beliefs**

| ↓ M's strategy ➜ S's strategy | Always choose D=0 | Always choose D=1 |
|---|---|---|
| **Always choose P=0** | $x < \dfrac{w - h(S_B)}{h(S_G) - h(S_B)}$ | $y < \dfrac{w - h(S_B)}{h(S_G) - h(S_B)}$ |
| **Always choose P=1** | $x > \dfrac{w - h(S_B)}{h(S_G) - h(S_B)}$ | $y > \dfrac{w - h(S_B)}{h(S_G) - h(S_B)}$ |

$x = Pr(S = S_G \mid D = 1) \mid y = Pr(S = S_G \mid D = 0)$

Fig 2: an overview of the required out-of equilibrium beliefs, both formally and informally

Note that the out-of-equilibrium beliefs only put restrictions on M's beliefs, not S's. The reason for that is that given the conditions for a pooling equilibrium exist, S's choice of D is completely independent of M's promotion policy. Therefore, the condition that S must not be induced to change his action when M does is by definition satisfied.

Proof for proposition 7

The following image illustrates the circle of events that has to hold for a semi-separating equilibrium to exist:
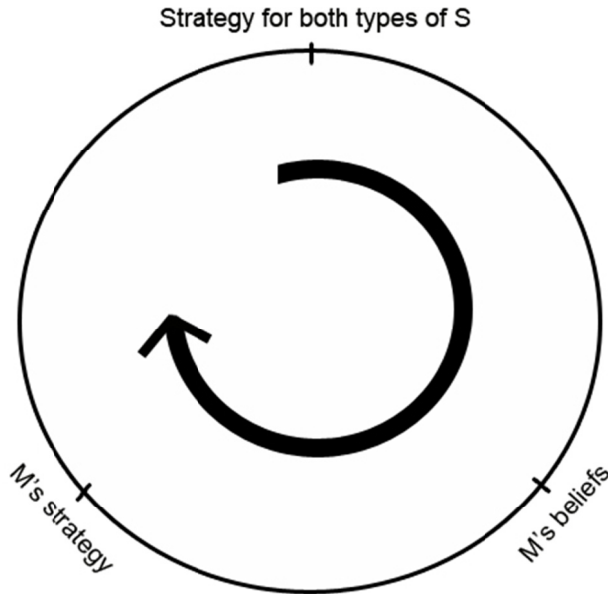
Fig. 3: the equilibrium circle

I will somewhat uncomfortable start from M's beliefs. Obviously, M can't have any beliefs without knowing S's strategy, but the alternative of starting with S's strategy poses the problem that there is an infinite amount of those. Therefore, I ask the reader to imagine that somehow, M has arrived at a set of beliefs.

M's optimal actions depend on his beliefs about S's type. There are three relevant beliefs, which are:
$\Pr(S = S_G \mid D = 0)$ (1)
$\Pr(S = S_G \mid D = 1 \text{ and } \pi = \pi_L)$ (2)
$\Pr(S = S_G \mid D = 1 \text{ and } \pi = \pi_H)$ (3)

It is optimal for M to choose P=0 if his belief about S's type exceeds a certain threshold, which is:

$$\frac{w - h(S_B)}{h(S_G) - h(S_B)} \ (4)$$

The following situations are possible (refer to proposition 3 to see why). For each situation, I have described M's optimal action:

(1) > (3) > (2) > (4) -> always promote (strategy 1)
(4) > (1) > (3) > (2) -> never promote (strategy 2)
(1) > (4) > (3) > (2) -> promote if D=0 (strategy 3)
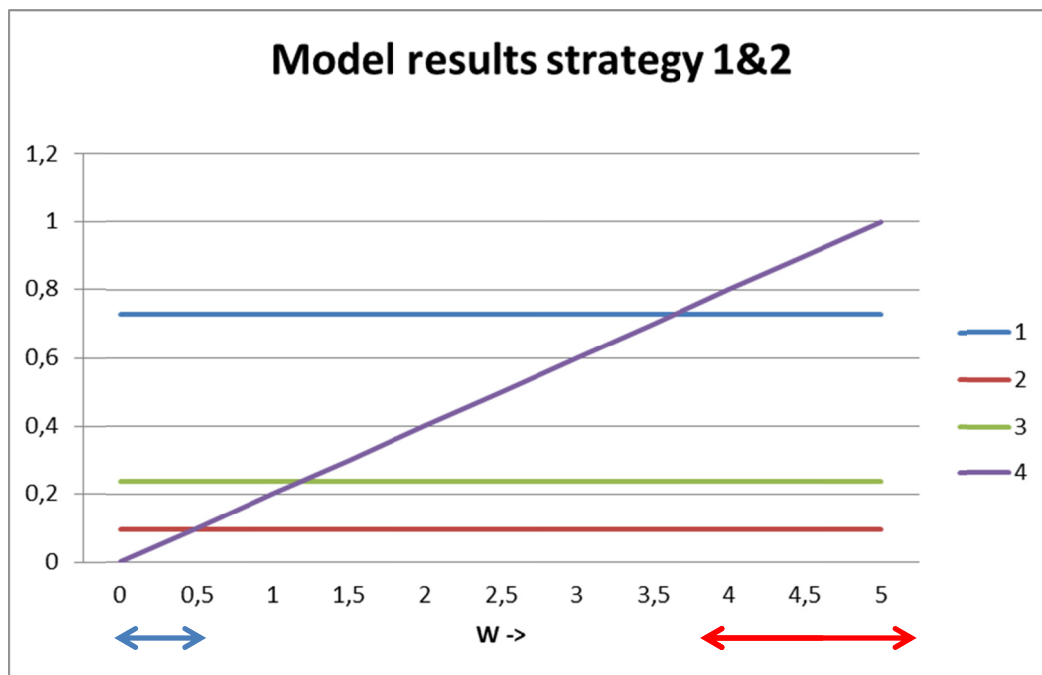(1) > (3) > (4) > (2) -> promote if D=0 or D=1 and $\pi = \pi_H$ (strategy 4)

Now that I've established M's strategy for each belief set, what remains is checking for each strategy what S's optimal strategy is and, in turn, if S's strategy could lead to the assumed beliefs for M. Performing this last step algebraically leads to expressions that are too large to interpret, due to the sheer amount of variables. I will therefore resort to giving numerical examples that prove the existence of a certain type of equilibrium.

*Strategy 1 and 2*

I'll first look at the case where M's beliefs are such that strategy 1 or 2 is optimal. In this case, S can't influence M's decision and will therefore revert to his first-period optimum. Therefore, he will choose D=0 if:

$$\lambda < \frac{g(\pi_H - \pi_L - c_C(D = 0)) - c_P(S)}{g(\pi_H - \pi_L)}$$

And D=1 otherwise. If, in turn, these two cut-off points (for each type of S) are such that the mentioned condition required for M to choose strategy 1 or 2 is met, there is an equilibrium.



The above graph shows the values of M's three beliefs (1, 2 and 3) and the cut-off point for M to choose P=1 (4) when M chooses strategy 1 or 2 (the resulting beliefs are the same) as a function of w. The other variables values are: $\pi_H = 10$; $\pi_L = -20$; $\mu = 0{,}4$; $g = 0{,}4$; $c_C(D = 0) = 5$, $c_P(S_G) = 0$; $c_P(S_B) = 3$; $h(S_B) = 0$; $h(S_G) = 5$
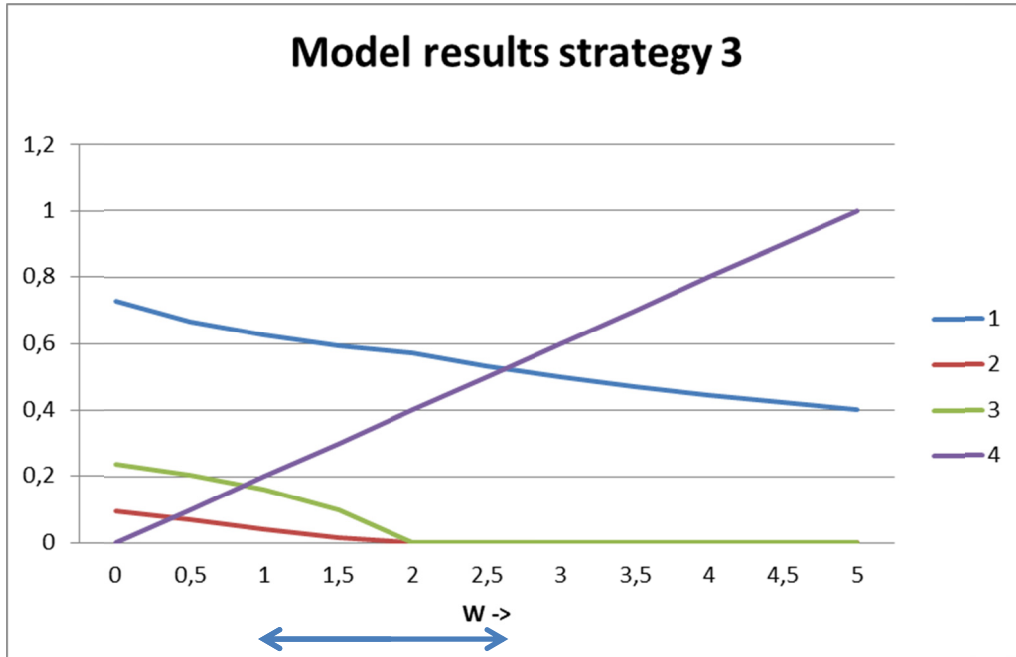
The blue arrow indicates the approximate value range of w for which M's beliefs support choosing strategy 1, the red arrow indicates the range which supports choosing strategy 2. In these cases, there is an equilibrium

*Strategy 3*

M's beliefs are such that strategy 3 is optimal. If M chooses strategy 3, S will choose D=0 if:

$$\lambda < \frac{g(\pi_H - \pi_L - c_C(D = 0)) - c_P(S) + w}{g(\pi_H - \pi_L)}$$

and D=1 otherwise. If, in turn, these two cut-off points (for each type of S) are such that the mentioned condition required for M to choose strategy 3 is met, there is an equilibrium.
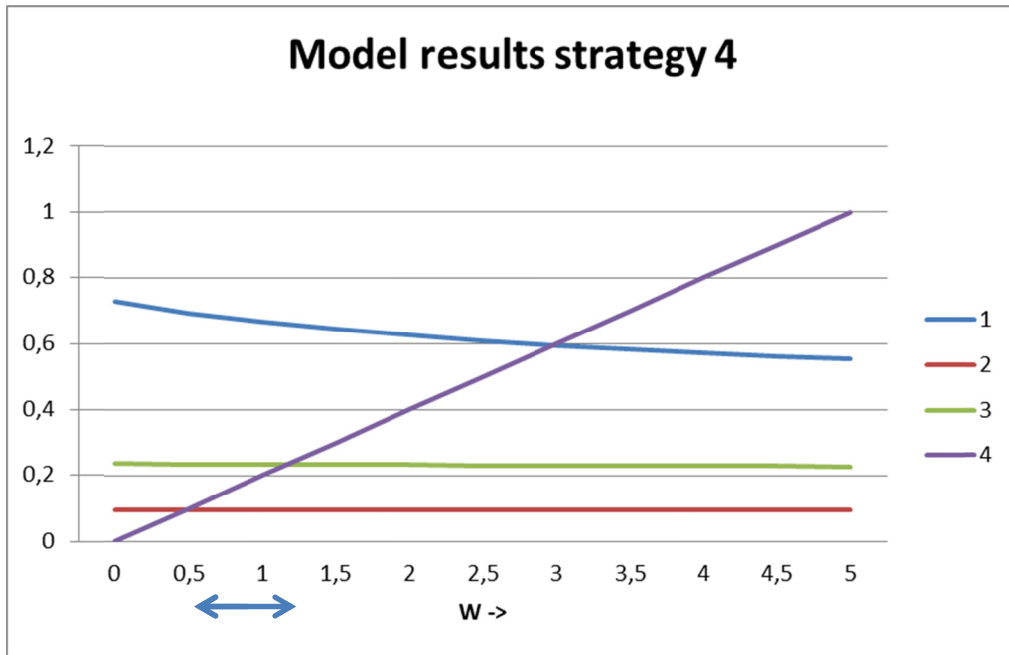
**Model results strategy 3**

The same parameters are used as under strategy 1 and 2. The blue arrow indicates the approximate range of values of w that support the equilibrium.

*Strategy 4*

M's beliefs are such that strategy 4 is optimal. If M chooses strategy 4, S will choose D=0 if:

$$\lambda < \frac{g(\pi_H - \pi_L - c_C(D = 0)) - c_P(S) + w}{g(\pi_H - \pi_L) + w}$$

and D=1 otherwise. Again, if the two cut-off points then meet the criteria for M to choose strategy 4, there is an equilibrium. I use the same parameters as for strategy 3, while substituting S's strategy.

Model results strategy 4

Again, the blue arrow shows the range of values of w for which S's strategy supports the beliefs that support the equilibrium. Note that because of the particular strategy for both types of S in this case, beliefs (2) and (3) are independent of w.

*Four types*

The above proves that with each of the four possible strategies for M, there exists a semi-separating weak sequential equilibrium.

# References

Crimer, J. (1995), "Arm's Length Relationships.", Quarterly Journal of Economics, 110(2), 275-95.

Holmström, B. (1979). "Moral Hazard and Observability.", Bell Journal of Economics, 10(1), 74-91.

Holmström, B. (1999). "Managerial Incentive Problems: A Dynamic Perspective." Review of Economic Studies, 66(1),169-82.

Morris, S. (2001), "Political Correctness", The Journal of Political Economy, 109(2):231-265

Prat, A. (2005), "The Wrong Kind of Transparency", The American Economic Review, 95(3), 862-877

Quin Mills, D. (1985), "Seniority Versus Ability in Promotion Decisions", Industrial and Labor Relations Review, 38(3):421-425

Spence, A.M. (1973), "Job Market Signaling", Quarterly Journal of Economics. 87:355-374

Ottaviani, M. and Sørensen, P. (2006), "Professional advice", Journal of Economic Theory. 126(1):120-14