

# Lectures on the Theory of Contracts

Lars A. Stole\*

First version: September 1993  
Current version: September 1999

## 1 Preface

The contours of contract theory as a field are difficult to define. Many would argue that contract theory is a subset of Game Theory which is defined by the notion that one party to the game (typically called the principal) is given all of the bargaining power and so can make a take-it-or-leave-it offer to the other party or parties (i.e., the agent(s)). In fact, the techniques for screening contracts were largely developed by pure game theorists to study allocation mechanisms and game design. But then again, carefully defined, everything is a subset of game theory.

Others would argue that contract theory is an extension of price theory in the following sense. Price theory studies how actors interact where the actors are allowed to choose prices, wages, quantities, etc. and studies partial or general equilibrium outcomes. Contract theory extends the choice spaces of the actors to include richer strategies (i.e. contracts) rather than simple one-dimensional choice variables. Hence, a firm can offer a nonlinear price menu to its customers (i.e., a screening contract) rather than a simple uniform price and an employer can offer its employee a wage schedule for differing levels of stochastic performance (i.e., an incentives contract) rather than a simple wage.

Finally, one could group contract theory together by the substantive questions it asks. How should contracts be developed between principals and their agents to provide correct incentives for communication of information and actions. Thus, contract theory seeks to understand organizations, institutions, and relationships between productive individuals when there are differences in personal objectives (e.g., effort, information revelation, etc.). It is this later classification that probably best defines contract theory as a field, although many interesting questions such as the optimal design of auctions and resource allocation do not fit this description very well but comprise an important part of contract theory nonetheless.

---

\*©1993, Lars A. Stole.

The notes provided are meant to cover the rough contours of contract theory. Much of their substance is borrowed heavily from the lectures and notes of Mathias Dewatripont, Bob Gibbons, Oliver Hart, Serge Moresi, Klaus Schmidt, Jean Tirole, and Jeff Zwiebel. In addition, helpful, detailed comments and suggestions were provided by Rohan Ptichford, Adriano Rampini, David Roth, Jennifer Wu and especially by Daivd Martimort. I have relied on many outside published sources for guidance and have tried to indicate the relevant contributions in the notes where they occur. Financial support for compiling these notes into their present form was provided by a National Science Foundation Presidential Faculty Fellowship and a Sloan Foundation Fellowship. I see the purpose of these notes as (i) to standardize the notation and approaches across the many papers in the field, (ii), to present the results of later papers building upon the theorems of earlier papers, and (iii) in a few cases present my own intuition and alternative approaches when I think it adds something to the presentation of the original author(s) and differs from the standard paradigm. Please feel free to distribute these notes in their entirety if you wish to do so.

## 2 Moral Hazard and Incentives Contracts

### 2.1 Static Principal-Agent Moral Hazard Models

#### 2.1.1 The Basic Theory

**The Model** We now turn to the consideration of moral hazard. The workhorse of this literature is a simple model with one principal who makes a take-it-or-leave-it offer to a single agent with outside reservation utility of  $\underline{U}$  under conditions of symmetric information. If the contract is accepted, the agent then chooses an action,  $a \in \mathcal{A}$ , which will have an effect (usual stochastic) on an outcome,  $x \in \mathcal{X}$ , of which the principal cares about and is typically “informative” about the agent’s action. The principal may observe some additional signal,  $s \in \mathcal{S}$ , which may also be informative about the agent’s action. The simplest version of this model casts  $x$  as monetary profits and  $s = \emptyset$ ; we will focus on this simple model for now ignoring information besides  $x$ .

We will assume that  $x$  is observable and *verifiable*. This latter term is used to indicate that enforceable contracts can be written on the variable,  $x$ . The nature of the principal’s contract offer will be a wage schedule,  $w(x)$ , according to which the agent is rewarded. We will also assume for now that the principal has full commitment and will not alter the contract  $w(x)$  later – even if it is Pareto improving.

The agent takes a hidden action  $a \in \mathcal{A}$  which yields a random monetary return  $\tilde{x} = x(\tilde{\theta}, a)$ ; e.g.,  $\tilde{x} = \tilde{\theta} + a$ . This action has the effect of stochastically improving  $x$  (e.g.,  $E_{\theta}[x_a(\tilde{\theta}, a)] > 0$ ) but at a monetary disutility of  $\psi(a)$ , which is continuously differentiable, increasing and strictly convex. The monetary utility of the principal is  $V(x - w(x))$ , where  $V' > 0 \geq V''$ . The agent’s net utility is separable in cost of effort and money:

$$U(w(x), a) \equiv u(w(x)) - \psi(a),$$

where  $u' > 0 \geq u''$ .

Rather than deal with the stochastic function  $\tilde{x} = x(\tilde{\theta}, a)$ , which is referred to as the *state-space representation* and was used by earlier incentive papers in the literature, we will find it useful to consider instead the density and distribution induced over  $x$  for a given

action; this is referred to as the *parameterized distribution characterization*. Let  $\mathcal{X} \equiv [\underline{x}, \bar{x}]$  be the support of output; let  $f(x, a) > 0$  for all  $x \in \mathcal{X}$  be the density; and let  $F(x, a)$  be the cumulative distribution function. We will assume that  $f_a$  and  $f_{aa}$  exist and are continuous. Furthermore, our assumption that  $E_\theta[x_a(\hat{\theta}, a)] > 0$  is equivalent to  $\int_{\underline{x}}^{\bar{x}} F_a(x, a) dx < 0$ ; we will assume a stronger (but easier to use) assumption that  $F_a(x, a) < 0 \quad \forall x \in (\underline{x}, \bar{x})$ ; i.e., effort produces a first-order stochastic dominant shift on  $\mathcal{X}$ . Finally, note that since our support is fixed,  $F_a(\underline{x}, a) = F_a(\bar{x}, a) = 0$  for any action,  $a$ . The assumption that the support of  $x$  is fixed is restrictive as we will observe below in remark 2.

**The Full Information Benchmark** Let's begin with the full information outcome where effort is observable and verifiable. The principal chooses  $a$  and  $w(x)$  to satisfy

$$\max_{w(\cdot), a} \int_{\underline{x}}^{\bar{x}} V(x - w(x)) f(x, a) dx,$$

subject to

$$\int_{\underline{x}}^{\bar{x}} u(w(x)) f(x, a) dx - \psi(a) \geq \underline{U},$$

where the constraint is the agent's participation or individual rationality (IR) constraint. The Lagrangian is

$$\mathcal{L} = \int_{\underline{x}}^{\bar{x}} [V(x - w(x)) + \lambda u(w(x))] f(x, a) dx - \lambda \psi(a) - \lambda \underline{U},$$

where  $\lambda$  is the Lagrange multiplier associated with the IR constraint; it represents the shadow price of income to the agent in each state. Assuming an interior solution, we have as first-order conditions,

$$\frac{V'(x - w(x))}{u'(w(x))} = \lambda, \quad \forall x \in \mathcal{X},$$

$$\int_{\underline{x}}^{\bar{x}} [V(x - w(x)) + \lambda u(w(x))] f_a(x, a) dx = \lambda \psi'(a),$$

and the IR constraint is binding.

The first condition is known as the Borch rule: the ratios of marginal utilities of income are equated across states under an optimal insurance contract. Note that it holds for every  $x$  and not just in expectation. The second condition is the choice of effort condition.

**Remarks:**

1. Note that if  $V' = 1$  and  $u'' < 0$ , we have a risk neutral principal and a risk averse agent. In this case, the Borch rule requires  $w(x)$  be constant so as to provide perfect insurance to the agent. If the reverse were true, the agent would perfectly insure the principal, and  $w(x) = x + k$ , where  $k$  is a constant.

2. Also note that the first-order condition for the choice of effort can be re-written as follows:

$$\begin{aligned}
\psi'(a) &= \int_{\underline{x}}^{\bar{x}} [V(x-w(x))/\lambda + u(w(x))] f_a(x, a) dx, \\
&= [V(x-w(x))/\lambda + u(w(x))] F_a(x, a) \Big|_{\underline{x}}^{\bar{x}} \\
&\quad - \int_{\underline{x}}^{\bar{x}} [V'(x-w(x))(1-w'(x))/\lambda + u'(w(x))w'(x)] F_a(x, a) dx, \\
&= - \int_{\underline{x}}^{\bar{x}} u'(w(x)) F_a(x, a) dx.
\end{aligned}$$

Thus, if the agent were risk neutral (i.e.,  $u' = 1$ ), integrating by parts one obtains

$$\int_{\underline{x}}^{\bar{x}} x f_a(x, a) dx = \psi'(a).$$

I.e.,  $a$  maximizes  $E[x|a] - \psi(a)$ . So even if effort cannot be contracted on as in the full-information case, if the agent is risk neutral then the principal can “sell the enterprise to the agent” with  $w(x) = x - k$ , and the agent will choose the first-best level of effort.

**The Hidden Action Case** We now suppose that the level of effort cannot be contracted upon and the agent is risk averse:  $u'' < 0$ . The principal solves the following program

$$\max_{w(\cdot), a} \int_{\underline{x}}^{\bar{x}} V(x-w(x)) f(x, a) dx,$$

subject to

$$\begin{aligned}
&\int_{\underline{x}}^{\bar{x}} u(w(x)) f(x, a) dx - \psi(a) \geq \underline{U}, \\
a &\in \arg \max_{a' \in \mathcal{A}} \int_{\underline{x}}^{\bar{x}} u(w(x)) f(x, a') dx - \psi(a'),
\end{aligned}$$

where the additional constraint is the incentive compatibility (IC) constraint for effort. The IC constraint implies (assuming an interior optimum) that

$$\begin{aligned}
&\int_{\underline{x}}^{\bar{x}} u(w(x)) f_a(x, a) dx - \psi'(a) = 0, \\
&\int_{\underline{x}}^{\bar{x}} u(w(x)) f_{aa}(x, a) dx - \psi''(a) \leq 0,
\end{aligned}$$

which are the local first- and second-order conditions for a maximum.

The *first-order approach* (FOA) to incentives contracts is to maximize subject to the first-order condition rather than IC, and then check to see if the solution indeed satisfies IC ex post. Let’s ignore questions of the validity of this procedure for now; we’ll return to the problems associated with its use later.

Using  $\mu$  as the multiplier on the effort first-order condition, the Lagrangian of the FOA program is

$$\mathcal{L} = \int_{\underline{x}}^{\bar{x}} V(x - w(x))f(x, a)dx + \lambda \int_{\underline{x}}^{\bar{x}} u(w(x))f(x, a)dx - \psi(a) - \underline{U} \\ + \mu \int_{\underline{x}}^{\bar{x}} u(w(x))f_a(x, a)dx - \psi'(a) .$$

Maximizing  $w(x)$  pointwise in  $x$  and simplifying, the first-order condition is

$$\frac{V'(x - w(x))}{u'(w(x))} = \lambda + \mu \frac{f_a(x, a)}{f(x, a)}, \quad \forall x \in \mathcal{X}. \quad (1)$$

We now have a modified Borch rule: the marginal rates of substitution may vary if  $\mu > 0$  to take into account the incentives effect of  $w(x)$ . Thus, risk-sharing will generally be inefficient.

Consider a simple two-action case in which the principal wishes to induce the high action:  $\mathcal{A} \equiv \{a_L, a_H\}$ . Then the IC constraint implies the inequality

$$\int_{\underline{x}}^{\bar{x}} u(w(x))[f(x, a_H) - f(x, a_L)]dx \geq \psi(a_H) - \psi(a_L).$$

The first-order condition for the associated Lagrangian is

$$\frac{V'(x - w(x))}{u'(w(x))} = \lambda + \mu \frac{f(x, a_H) - f(x, a_L)}{f(x, a_H)}, \quad \forall x \in \mathcal{X}.$$

In both cases, providing  $\mu > 0$ , the agent is rewarded for outcomes which have higher relative frequency under high effort. We now prove that this is indeed the case.

**Theorem 1** (Holmst m, [1979]) Assume that the FOA program is valid. Then at the optimum of the FOA program,  $\mu > 0$ .

**Proof:** The proof of the theorem relies upon first-order stochastic dominance  $F_a(x, a) < 0$  and risk aversion  $u'' < 0$ . Consider  $\frac{\partial \mathcal{L}}{\partial a} = 0$ . Using the agent's first-order condition for effort choice, it simplifies to

$$\int_{\underline{x}}^{\bar{x}} V(x - w(x))f_a(x, a)dx + \mu \int_{\underline{x}}^{\bar{x}} u(w(x))f_{aa}(x, a)dx - \psi''(a) = 0.$$

Suppose that  $\mu \leq 0$ . By the agent's second-order condition for choosing effort we have

$$\int_{\underline{x}}^{\bar{x}} V(x - w(x))f_a(x, a)dx \leq 0.$$

Now, define  $w_\lambda(x)$  as the solution to refmh-star when  $\mu = 0$ ; i.e.,

$$\frac{V'(x - w_\lambda(x))}{u'(w_\lambda(x))} = \lambda, \quad \forall x \in \mathcal{X}.$$

Note that because  $u'' < 0$ ,  $w_\lambda$  is differentiable and  $w'_\lambda(x) \in [0, 1)$ . Compare this to the solution,  $w(x)$ , which satisfies

$$\frac{V'(x - w(x))}{u'(w(x))} = \lambda + \mu \frac{f_a(x, a)}{f(x, a)}, \quad \forall x \in \mathcal{X}.$$

When  $\mu \leq 0$ , it follows that  $w(x) \leq w_\lambda(x)$  if and only if  $f_a(x, a) \geq A_s$ . Thus,

$$V(x - w(x))f_a(x, a) \geq V(x - w_\lambda(x))f_a(x, a), \quad \forall x \in \mathcal{X},$$

and as a consequence,

$$\int_{\underline{x}}^{\bar{x}} V(x - w(x))f_a(x, a)dx \geq \int_{\underline{x}}^{\bar{x}} V(x - w_\lambda(x))f_a(x, a)dx.$$

The RHS is necessarily positive because integrating by parts yields

$$V(x - w_\lambda(x))F_a(x, a)|_{\underline{x}}^{\bar{x}} - \int_{\underline{x}}^{\bar{x}} V'(x - w_\lambda(x))(1 - w'_\lambda(x))F_a(x, a)dx > 0.$$

But then this implies a contradiction.  $\square$

**Remarks:**

1. Note that we have assumed that the agent's chosen action is in the interior of  $a \in \mathcal{A}$ . If the principal wishes to implement the least costly action in  $\mathcal{A}$ , perhaps because the agent's actions have little effect on output, then the agent's first-order condition does not necessarily hold. In fact, the problem is trivial and the principal will supply full insurance if  $V'' = 0$ . In this case an inequality for the agent's corner solution must be included in the maximization program rather than a first-order condition; the associated multiplier will be zero:  $\mu' = 0$ .
2. The assumption that the support of  $x$  does not depend upon effort is crucial. If the support "shifts", it may be possible to obtain the first best, as some outcomes may be perfectly informative about effort.
3. Commitment is important in the implementation of the optimal contract. In equilibrium the principal knows the agent took the required action. Because  $\mu > 0$ , the principal is imposing unnecessary risk upon the agent, *ex post*. Between the time the action is taken and the uncertainty resolved, there generally exist Pareto-improving contracts to which the parties could mutually renegotiate. The principal must commit not to do this to implement the optimal contract above. We will return to this issue when we consider the case in which the principal cannot commit not to renegotiate.
4. Note that  $f_a/f$  is the derivative of the log-likelihood function,  $\log f(x, a)$ , and hence is the gradient for a MLE estimate of  $a$  given observed  $x$ . In equilibrium, the principal "knows"  $a = a^*$  even though he commits to a mechanism which rewards based upon the informativeness of  $x$  for  $a = a^*$ .

5. The above FOA program can be extended to models of adverse selection and moral hazard. That is, the agent observes some private information,  $\theta$ , before contracting (and choosing action) and the principal offers a wage schedule,  $w(x, \hat{\theta})$ . The difference now is that  $\mu(\hat{\theta})$  will generally depend upon the agent's announced information and may become nonpositive for some values. See Holmstrom[1979], section 6.
6. Note that although  $w'(x) \in [0, 1)$ , we haven't demonstrated that the optimal  $w(x)$  is  $\lambda$  monotonic. Generally, first-order stochastic dominance is not enough. We will need a stronger property.

We now turn to the question of monotonicity.

**Definition 1** The **monotone likelihood ratio property (MLRP)** is satisfied for a distribution  $F$  and its density  $f$  iff

$$\frac{d}{dx} \frac{f_a(x, a)}{f(x, a)} \geq 0.$$

Note that when effort is restricted to only two types so  $f$  is non-differentiable, the analogous MLRP condition is that

$$\frac{d}{dx} \frac{f(x, a_H) - f(x, a_L)}{f(x, a_H)} \geq 0.$$

Additionally it is worth noting that MLRP implies that  $F_a(x, a) < 0$  for  $x \in (\underline{x}, \bar{x})$  (i.e., first-order stochastic dominance). Specifically, for  $x \in (\underline{x}, \bar{x})$

$$F_a(x, a) = \int_{ux}^x \frac{f_a(s, a)}{f(s, a)} f(s, a) ds < 0,$$

where the latter inequality follows from MLRP (when  $x = \bar{x}$ , the integral is 0; when  $x < \bar{x}$ , the fact that the likelihood ratio is increasing in  $s$  implies that the integral must be strictly negative).

We have the following result.

**Theorem 2** (Holmstrom 1979], Shavell [1979]). Under the first-order approach, if  $F$  satisfies the monotone likelihood ratio property, then the wage contract is increasing in output.

The proof is immediate from the definition of MLRP and our first-order conditions above.

**Remarks:**

1. Sometimes you do not want monotonicity because the likelihood ratio is u-shaped. For example, suppose that only two effort levels are possible,  $\{a_L, a_H\}$ , and only three output levels can occur,  $\{x_1, x_2, x_3\}$ . Let the probabilities be given by the following  $f(x, a)$ :

$f(x, a)$	$x_1$	$x_2$	$x_3$
$a_H$	0.4	0.1	0.5
$a_L$	0.5	0.4	0.1
L-Ratio	-0.25	-3	0.8

If the principal wishes to induce high effort, the idea is to use a non-monotone wage schedule to punish moderate outputs which are most indicative of low effort, reward high outputs which are quite informative about high effort, and provide moderate income for low outputs which are not very informative.

2. If agents can freely dispose of the output, monotonicity may be a constraint which the solution must satisfy. In addition, if agent's can trade output amongst themselves (say the agents are sharecroppers), then only a linear constraint is feasible with lots of agents; any nonlinearities will be arbitrated away.
3. We still haven't demonstrated that the first-order approach is correct. We will turn to this shortly.

**The Value of Information** Let's return to our more general setting for a moment and assume that the principal and agent can enlarge their contract to include other information, such as an observable and verifiable signal,  $s$ . When should  $w$  depend upon  $s$ ?

**Definition 2**  $x$  is **sufficient** for  $\{x, s\}$  with respect to  $a \in \mathcal{A}$  iff  $f$  is multiplicatively separable in  $s$  and  $a$ ; i.e.

$$f(x, s, a) \equiv y(x, a)z(x, s).$$

We say that  $s$  is **informative** about  $a \in \mathcal{A}$  whenever  $x$  is not sufficient for  $\{x, s\}$  with respect to  $a \in \mathcal{A}$ .

**Theorem 3** (Holmst in, [1979], Shavell [1979]). Assume that the FOA program is valid and yields  $w(x)$  as a solution. Then there exists a new contract,  $w(x, s)$ , that strictly Pareto dominates  $w(x)$  iff  $s$  is informative about  $a \in \mathcal{A}$ .

**Proof:** Using the FOA program, but allowing  $w(\cdot)$  to depend upon  $s$  as well as  $x$ , the first-order condition determining  $w$  is given by

$$\frac{V'(x - w(x, s))}{u'(w(x, s))} = \lambda + \mu \frac{f_a(x, s, a)}{f(x, s, a)},$$

which is independent of  $s$  iff  $s$  is not informative about  $a \in \mathcal{A}$ .  $\square$

The result implies that without loss of generality, the principal can restrict attention to wage contracts that depend only upon a set of sufficient statistics for the agent's action. Any other dependence cannot improve the contract; it can only increase the risk the agent faces without improving incentives. Additionally, the result says that any informative signal about the agent's action should be included in the optimal contract!



**Application: Insurance Deductibles.** We want to show that under reasonable assumptions it is optimal to offer insurance policies which provide full insurance less a fixed deductible. The idea is that if conditional on an accident occurring, the value of the loss is uninformative about  $a$ , then the coverage of optimal insurance contract should also not depend upon actual losses – only whether or not there was an accident. Hence, deductibles are optimal. To this end, let  $x$  be the size of the loss and assume  $f_a(0, a) = 1 - p(a)$  and  $f(x, a) = p(a)g(x)$  for  $x < 0$ . Here, the probability of an accident depends upon effort in the obvious manner:  $p'(a) < 0$ . The amount of loss  $x$  is independent of  $a$  (i.e.,  $g(x)$  is independent of  $a$ ). Thus, the optimal contract is characterized by a likelihood ratio of  $\frac{f_a(x, a)}{f(x, a)} = \frac{p'(a)}{p(a)} < 0$  for  $x < 0$  (which is independent of  $x$ ) and  $\frac{f_a(x, a)}{f(x, a)} = \frac{-p'(a)}{1-p(a)} > 0$  for  $x = 0$ . This implies that the final income allocation to the agent is fixed at one level for all  $x < 0$  and at another for  $x = 0$ , which can be implemented by the insurance company by offering full coverage less a deductible.

**Asymptotic First-best** It may be that by making very harsh punishments with very low probability that the full-information outcome can be approximated arbitrarily closely. This insight is due to Mirrlees [1974].

**Theorem 4** (Mirrlees, [1974].) Suppose  $f(x, a)$  is the normal distribution with mean  $a$  and variance  $\sigma^2$ . Then if unlimited punishments are possible, the first-best can be approximated arbitrarily closely.

**Sketch of Proof:** We prove the theorem for the case of a risk-neutral principal. We have

$$f(x, a) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-a)^2}{2\sigma^2}},$$

so that

$$\frac{f_a(x, a)}{f(x, a)} = \frac{d}{da} \log f(x, a) = \frac{(x-a)}{\sigma^2}.$$

That is, detection is quite efficient for  $x$  very small.

The first-best contract has a constant  $w^*$  such that  $u(w^*) = \underline{U} + \psi(a^*)$ , where  $a^*$  is the first-best action. The approximate first-best contract offers  $w^*$  for all  $x \geq x_o$  ( $x_o$  very small), and  $w = k$  ( $k$  very small) for  $x < x_o$ . Choose  $k$  low enough such that  $a^*$  is optimal for the agent (IC) at a given  $x_o$ :

$$\int_{-\infty}^{x_o} u(k) f_a(x, a^*) dx + \int_{x_o}^{\infty} u(w^*) f_a(x, a^*) dx = \psi'(a^*).$$

We want to show that with this contract, the agent's IR constraint can be satisfied arbitrarily closely as we lower the punishment region. Note that the loss with respect to the first-best is

$$\Delta \equiv \int_{-\infty}^{x_o} (u(w^*) - u(k)) f(x, a^*) dx,$$

for the agent. Define  $M(x_o) \equiv \frac{f_a(x_o, a^*)}{f(x_o, a^*)}$ . Because the normal distribution satisfies MLRP, for all  $x < x_o$ ,  $f_a/f < M$  or  $f > f_a/M$ . This implies that the difference between the agent's

utility and  $\underline{U}$  is bounded above by

$$\frac{1}{M} \int_{-\infty}^{x_o} (u(w^*) - u(k)) f_a(x, a^*) dx \geq \Delta.$$

But by the agent's IC condition, this bound is given by

$$\frac{1}{M} \int_{-\infty}^{\infty} u(w^*) f_a(x, a^*) dx - \psi'(a^*) \quad ,$$

which is a constant divided by  $M$ . Thus, as punishments are increased, i.e,  $M$  decreased, we approach the first best.  $\square$

The intuition for the result is that in the tails of a normal distribution, the outcome is very informative about the agent's action. Thus, even though the agent is risk averse and harsh random punishments are costly to the principal, the gain from informativeness dominates at any punishment level.

**The Validity of the First-order Approach** We now turn to the question of the validity of the first-order approach.

The approach of first finding  $w(x)$  using the relaxed FOA program and then checking that the principal's selection of  $a$  maximizes the agent's objective function is logically invalid without additional assumptions. Generally, the problem is that when the second-order condition of the agent is not be globally satisfied, it is possible that the solution to the unrelaxed program satisfies the agent's first-order condition (which is necessary) but not the principal's first-order condition. That is, the principal's optimum may involve a corner solution and so solutions to the unrelaxed direct program may not satisfy the necessary Kuhn-Tucker conditions of the relaxed FOA program. This point was first made by Mirrlees [1974].

There are a few important papers on this concern. Mirrlees [1976] has shown that MLRP with an additional convexity in the distribution function condition (CDFC) is sufficient for the validity of the first-order approach.

**Definition 3** A distribution satisfies the **Convexity of Distribution Function Condition (CDFC)** iff

$$F(x, \gamma a + (1 - \gamma)a') \leq \gamma F(x, a) + (1 - \gamma)F(x, a'),$$

for all  $\gamma \in [0, 1]$ . (I.e.,  $F_{aa}(x, a) \geq 0$ .)

A useful special case of this CDF condition is the *linear distribution function condition*:

$$f(x, a) \equiv a \bar{f}(x) + (1 - a) \underline{f}(x),$$

where  $\bar{f}(x)$  first-order stochastically dominates  $\underline{f}(x)$ .

Mirrlees' theorem is correct but the proof contains a subtle omission. Independently, Rogerson [1985] determines the same sufficient conditions for the first-order approach using a correct proof. Essentially, he derives conditions which guarantee that the agent's objective function will be globally concave in action for any selected contract.

In an earlier paper, Grossman and Hart [1983] study the general unrelaxed program directly rather than a relaxed program. They find, among other things, that MLRP and CDFC are sufficient for the monotonicity of optimal wage schedules. Additionally, they show that the program can be reduced to an elegant linear programming problem. Their methodology is quite nice and a significant contribution to contract theory independent of their results.

**Rogerson [1985]:**

**Theorem 5** (Rogerson, [1985].) The first-order approach is valid if  $F(x, a)$  satisfies the MLRP and CDF conditions.

We first begin with a simple commonly used, but incorrect, “proof” to illustrate the subtle circularity of proving the validity of the FOA program.

**“Proof”:** One can rewrite the agent’s payoff as

$$\begin{aligned} \int_{\underline{x}}^{\bar{x}} u(w(x)) f(x, a) dx - \psi(a) &= u(w(x)) F(x, a) \Big|_{\underline{x}}^{\bar{x}} \\ &- \int_{\underline{x}}^{\bar{x}} u'(w(x)) \frac{dw(x)}{dx} F(x, a) dx - \psi(a) \\ &= u(w(\bar{x})) - \int_{\underline{x}}^{\bar{x}} u'(w(x)) \frac{dw(x)}{dx} F(x, a) dx - \psi(a), \end{aligned}$$

where we have assumed for now that  $w(x)$  is differentiable. Differentiating this with respect to  $a$  twice yields

$$- \int_{\underline{x}}^{\bar{x}} u'(w(x)) \frac{dw(x)}{dx} F_{aa}(x, a) dx - \psi''(a) < 0,$$

for every  $a \in \mathcal{A}$ . Thus, the agent’s second-order condition is globally satisfied in the FOA program if  $w(x)$  is differentiable and nondecreasing. Under MLRP,  $\mu > 0$ , and so the first-order approach yields a monotonically increasing, differentiable  $w(x)$ ; we are done.  $\square$

**Note:** The mistake is in the last line of the proof which is circular. You cannot use the FOA  $\mu > 0$  result, without first proving that the first-order approach is valid. (In the proof of  $\mu > 0$ , we implicitly assumed that the agent’s second-order condition was satisfied).

Rogerson avoids this problem by focusing on a doubly-relaxed program where the first-order condition is replaced by  $\frac{d}{da} E[U(w(x), a)|a] \geq 0$ . Because the constraint is an inequality, we are assured that the multiplier is nonnegative:  $\delta \geq 0$ . Thus, the solution to the doubly-relaxed program implies a nondecreasing, differentiable wage schedule under MLRP. The second step is to show that the solution of the doubly-relaxed program satisfies the constraints of the relaxed program (i.e., the optimal contract satisfies the agent’s first-order condition with equality). This result, combined with the above “Proof”, provides a complete proof of the theorem by demonstrating that the double-relaxed solution satisfies the unrelaxed constraint set. This second step is provided in the following lemma.

**Lemma 1** (Rogerson, [1985].) At the doubly-relaxed program solution,

$$\frac{d}{da}E[U(w, a)|a] = 0.$$

**Proof:** To see that  $\frac{d}{da}E[U(w, a)|a] = 0$  at the solution doubly-relaxed program, consider the inequality constraint multiplier,  $\delta$ . If  $\delta > 0$ , the first-order condition is clearly satisfied. Suppose  $\delta = 0$  and necessarily  $\lambda > 0$ . This implies the optimal risk-sharing choice of  $w(x) = w_\lambda(x)$ , where  $w'_\lambda(x) \in [0, 1)$ . Integrating the expected utility of the principal by parts yields

$$E[V(x - w(x))|a] = V(\bar{x} - w_\lambda(\bar{x})) - \int_{\underline{x}}^{\bar{x}} V'(x - w_\lambda(x))[1 - w'_\lambda(x)]F(x, a)dx.$$

Differentiating with respect to action yields

$$\frac{\partial E[V|a]}{\partial a} = - \int_{\underline{x}}^{\bar{x}} V'(x - w_\lambda(x))[1 - w'_\lambda(x)]F_a(x, a)dx \geq 0,$$

where the inequality follows from  $F_a \leq 0$ . Given that  $\lambda > 0$ , the first-order condition of the doubly-relaxed program for  $a$  requires that

$$\frac{d}{da}E[U(w(x), a)|a] \leq 0.$$

This is only consistent with the doubly-relaxed constraint set if

$$\frac{d}{da}E[U(w(x), a)|a] = 0,$$

and so the first-order condition must be satisfied.  $\square$

**Remarks:**

1. Due to Mirrlees' [1976] initial insight and Rogerson's [1985] correction, the MLRP-CDFC sufficiency conditions are usually called the Mirrlees-Rogerson sufficient conditions.
2. The conditions of MLRP and CDFC are very strong. It is difficult to think of many distributions which satisfy them. One possibility is a generalization of the uniform distribution (which is a type of  $\beta$ -distribution):

$$F(x, a) \equiv \frac{x - \underline{x}}{\bar{x} - \underline{x}}^{\frac{1}{1-a}},$$

where  $\mathcal{A} = [0, 1)$ .

3. The CDF condition is particularly strong. Suppose for example that  $\tilde{x} \equiv a + \tilde{\varepsilon}$ , where  $\tilde{\varepsilon}$  is distributed according to some cumulative distribution function. Then the CDF condition requires that the density of the cumulative distribution is increasing in  $\varepsilon$ ! Jewitt [1988] provides a collection of alternative sufficient conditions on  $F(x, a)$  and  $u(w)$

which avoid the assumption of CDFC. Examples include CARA utility with either (i) a Gamma distribution with mean  $\alpha a$  (i.e.,  $f(x, a) = a^{-\alpha} x^{\alpha-1} e^{-x/a} \Gamma(\alpha)^{-1}$ ), (ii) a Poisson distribution with mean  $a$  (i.e.,  $f(x, a) = a^x e^{-a} / \Gamma(1+x)$ ), or (iii) a Chi-squared distribution with  $a$  degrees of freedom (i.e.,  $f(x, a) = \Gamma(2a)^{-1} 2^{-2a} x^{2a-1} e^{-x/2}$ ). Jewitt also extends the sufficiency theorem to situations in which there are multiple signals as in Holmst in [1979]. With CARA utility for example, MLRP and CDFC are sufficient for the validity of the FOA program with multiple signals. See also Sinclair-Desgagne [1994] for further generalizations for use of the first-order approach and multiple signals.

4. Jewitt also has a nice elegant proof that the solution to the relaxed program (valid or invalid) necessarily has  $\mu > 0$  when  $v'' = 0$ . The idea is to show that  $u(w(x))$  and  $1/u'(w(x))$  generally have a positive covariance and note that at the solution to the FOA program, the covariance is equal to  $\mu\psi'(a)$ . Specifically, note that 1 is equivalent to

$$f_a(x, a) = \frac{1}{u'(w(x))} - \lambda \frac{f(x, a)}{\mu},$$

which allows us to rewrite the agent's first-order condition for action as

$$\int_x^{\bar{x}} u(w(x)) \frac{1}{u'(w(x))} - \lambda \int_x^{\bar{x}} f(x, a) dx = \mu\psi'(a).$$

Since the expected value of each side of 1 is zero, the mean of  $\frac{1}{u'(w(x))}$  is  $\lambda$  and so the righthand side of the above equation is the covariance of  $u(w(x))$  and  $\frac{1}{u'(w(x))}$ ; and since both functions are increasing in  $w(x)$ , the covariance is nonnegative implying that  $\mu \geq 0$ . This proof could have been used above instead of either Holmst in's or Rogerson's results to prove a weaker theorem applicable only to risk-neutral principals.

### Grossman-Hart [1983]:

We now explicitly consider the unrelaxed program following the approach of Grossman and Hart.

Following the framework of G-H, we assume that there are only a finite number of possible outputs,  $x_1 < x_2 < \dots, x_N$ , which occur with probabilities:  $f(x_i, a) \equiv \text{Prob}[\tilde{x} = x_i | a] > 0$ . We assume that the principal is risk neutral:  $V'' = 0$ . The agent's utility function is slightly more general:

$$U(w, a) \equiv K(a)u(w) - \psi(a),$$

where every function is sufficiently continuous and differentiable. This is the most general utility function which preserves the requirement that the agent's preference ordering over income lotteries is independent of action. Additionally, if either  $K'(a) = 0$  or  $\psi'(a) = 0$ , the agent's preferences over action lotteries is independent of income. We additionally assume that  $\mathcal{A}$  is compact,  $u' > 0 > u''$  over the interval  $(\underline{I}, \infty)$ , and  $\lim_{w \rightarrow \underline{I}} K(a)u(w) = -\infty$ . [The latter bit excludes corner solutions in the optimal contract. We were implicitly assuming

this above when we focused on the FOA program.] Finally, for existence of a solution we need that for any action there exists a  $w \in (\underline{I}, \infty)$  such that  $K(a)u(w) - \psi(a) \geq \underline{U}$ .

When the principal cannot observe  $a$ , the second-best contract solves

$$\max_{w,a} \quad f(x_i, a)(x_i - w_i),$$

subject to

$$a \in \arg \max_{a'} \quad f(x_i, a')[K(a')u(w_i) - \psi(a')],$$

$$f(x_i, a)[K(a)u(w_i) - \psi(a)] \geq \underline{U}.$$

We proceed in two steps: first, we solve for the least costly way to implement a given action. Then we determine the optimal action to implement.

What's the least costly way to implement  $a^* \in \mathcal{A}$ ? The principal solves

$$\min_w \quad f(x_i, a^*)w_i,$$

subject to

$$f(x_i, a^*)[K(a^*)u(w_i) - \psi(a^*)] \geq \quad f(x_i, a)[K(a)u(w_i) - \psi(a)], \quad \forall a \in \mathcal{A},$$

$$f(x_i, a^*)[K(a^*)u(w_i) - \psi(a^*)] \geq \underline{U}.$$

This is not a convex programming problem amenable to the Kuhn-Tucker theorem. Following Grossman and Hart, we convert it using the following transformation: Let  $h(u) \equiv u^{-1}(u)$ , so  $h(u(w)) = w$ . Define  $u_i \equiv u(w_i)$ , and use this as the control variable. Substituting yields

$$\min_u \quad f(x_i, a^*)h(u_i),$$

subject to

$$f(x_i, a^*)[K(a^*)u_i - \psi(a^*)] \geq \quad f(x_i, a)[K(a)u_i - \psi(a)], \quad \forall a \in \mathcal{A},$$

$$f(x_i, a^*)[K(a^*)u_i - \psi(a^*)] \geq \underline{U}.$$

Because  $h$  is convex, this is a convex programming problem with a linear constraint set. Grossman and Hart further show that if either  $K'(a) = 0$  or  $\psi'(a) = 0$ , (i.e., preferences over actions are independent of income), the solution to this program will have the IR constraint binding. In general, the IR constraint may not bind when wealth effects (i.e., leaving money to the agent) induce cheaper incentive compatibility. From now on, we assume that  $K(a) = 1$  so that the IR constraint binds.

Further note that when  $\mathcal{A}$  is finite, we will have a finite number of constraints, and so we can appeal to the Kuhn-Tucker theorem for necessary and sufficient conditions. We will do this shortly. For now, define

$$C(a^*) \equiv \begin{cases} \inf_i \{ f(x_i, a^*)h(u_i) \} & \text{if } w \text{ implements } a^*, \\ \infty & \text{otherwise.} \end{cases}$$

Note that some actions cannot be feasibly implemented with any incentive scheme. For example, the principal cannot induce the agent to take a costly action that is dominated:  $f(x_i, a) = f(x_i, a') \forall i$ , but  $\psi(a) > \psi(a')$ .

Given our construction of  $C(a)$ , the principal's program amounts to choosing  $a$  to maximize  $B(a) - C(a)$ , where  $B(a) \equiv \sum_i f(x_i, a)x_i$ . Grossman and Hart demonstrate that a (second-best) optimum exists and that the inf in the definition of  $C(a)$  can be replaced with a min.

### Characteristics of the Optimal Contract

1. Suppose that  $\psi(a_{FB}) > \min_{a'} \psi(a')$ ; (i.e., the first-best action is not the least cost action). Then the second-best contract produces less profit than the first-best. The proof is trivial: the first-best requires full insurance, but then the least cost action will be chosen.
2. Assume that  $\mathcal{A}$  is finite so that we can use the Kuhn-Tucker theorem. Then we have the following program:

$$\max_u \quad - \sum_i f(x_i, a^*)h(u_i),$$

subject to

$$f(x_i, a^*)[u_i - \psi(a^*)] \geq \sum_i f(x_i, a_j)[u_i - \psi(a_j)], \quad \forall a_j \neq a^*,$$

$$\sum_i f(x_i, a^*)[u_i - \psi(a^*)] \geq \underline{U}.$$

Let  $\mu_j \geq 0$  be the multiplier on the  $j$ th IC constraint;  $\lambda \geq 0$  the multiplier on the IR constraint. The first-order condition for  $u_i$  is

$$h'(u_i) = \lambda + \sum_{a_j \in A, a_j \neq a^*} \mu_j \frac{f(x_i, a^*) - f(x_i, a_j)}{f(x_i, a^*)}.$$

We know from Grossman-Hart that the IR constraint is binding:  $\lambda > 0$ . Additionally, from the Kuhn-Tucker theorem, providing  $a^*$  is not the minimum cost action,  $\mu_j > 0$  for some  $j$  where  $\psi(a_j) < \psi(a^*)$ .

3. Suppose that  $\mathcal{A} \equiv \{a_L, a_H\}$  (i.e., there are only two actions), and the principal wishes to implement the more costly action,  $a_H$ . Then the first-order condition becomes:

$$h'(u_i) = \lambda + \mu_L \frac{f(x_i, a_H) - f(x_i, a_L)}{f(x_i, a_H)}.$$

Because  $\mu_L > 0$ ,  $w_i$  increases with  $\frac{f(x_i, a_L)}{f(x_i, a_H)}$ . The condition that this likelihood ratio increase in  $i$  is the MLRP condition for discrete distributions and actions. Thus, MLRP is sufficient for monotonic wage schedules when there are only two actions.

4. Still assuming that  $\mathcal{A}$  is finite, with more than two actions all we can say about the wage schedule is that it cannot be decreasing *everywhere*. This is a very weak result.

The reason is clear from the first-order condition above. MLRP is not sufficient to prove that  $\prod_{a_j \in A} \mu_j \frac{f(x_i, a_j)}{f(x_i, a^*)}$  is nonincreasing in  $i$ . Combining MLRP with a variant of CDFC (or alternatively imposing a *spanning condition*), however, Grossman and Hart show that monotonicity emerges.

Thus, while before we demonstrated that MLRP and CDFC guarantees the FOA program is valid and yields a monotonic wage schedule, Grossman and Hart's direct approach also demonstrates that MLRP and CDFC guarantee monotonicity directly. Moreover, as Grossman and Hart discuss in section 6 of their paper, many of their results including monotonicity generalize to the case of a risk averse principal.

### 2.1.2 Extensions: Moral Hazard in Teams

We now turn to an analysis of multi-agent moral hazard problems, frequently referred to as moral hazard in teams (or partnerships). The classic reference is Holmstrom[1982].<sup>1</sup> Holmstrom makes two contributions in this paper. First, he demonstrates the importance of a budget-breaker. Second, he generalizes the notion of sufficient statistics and informativeness to the case of multi-agent situations and examines relative performance evaluation. We consider each contribution in turn.

#### The Importance of a Budget Breaker.

The canonical multi-agent model, has one risk neutral principal and  $N$  possibly risk-averse agents, each who privately choose  $a_i \in \mathcal{A}_i$  at a cost given by a strictly increasing and convex cost function,  $\psi_i(a_i)$ . We assume as before that  $a_i$  cannot be contracted upon. The output of the team of agents,  $x$ , depends upon the agents' efforts  $a \equiv (a_1, \dots, a_N) \in \mathcal{A} \equiv \mathcal{A}_1 \times \dots \times \mathcal{A}_N$ , which we will assume is deterministic for now. Later, we will allow for a stochastic function as before in the single-agent model.

A contract is a collection of wage schedules  $w = (w_1, \dots, w_N)$ , where each agent's wage schedule indicates the transfer the agent receives as a function of verifiable output; i.e.,  $w_i(x) : \mathcal{X} \rightarrow \mathbb{R}^N$ .

The timing of the game is as follows. At stage 1, the principal offers a wage schedule for each agent which is observable by everyone. At stage two, the agents reject the contract or accept and simultaneously and noncooperatively select their effort levels,  $a_i$ . At stage 3, the output level  $x$  is realized and participating agents are paid appropriately.

We say we have a **partnership** when there is effectively no principal and so the agents

---

<sup>1</sup>Mookherjee [1984] also examines many of the same issues in the context of a Grossman-hart [1983] style moral hazard model, but with many agents. His results on sufficient statistics mirrors those of Holmström's, but in a discrete environment.



split the output amongst themselves; i.e., the wage schedule satisfies budget balance:

$$\sum_{i=1}^N w_i(x) = x, \quad \forall x \in \mathcal{X}.$$

An important aspect of a partnership is that the budget (i.e., transfers) are always exactly balanced on and off the equilibrium path.

Holmst in [1982] points out that one frequently overlooked benefit of a principal is that she can break the budget while a partnership cannot. To illustrate this principle, suppose that the agents are risk neutral for simplicity.

**Theorem 6** Assume each agent is risk neutral. Suppose that output is deterministic and given by the function  $x(a)$ , strictly increasing and differentiable. If aggregate wage schedules are allowed to be less than total output (i.e.,  $\sum_i w_i < x$ ), then the first-best allocation  $(x^*, a^*)$  can be implemented as a Nash equilibrium with all output going to the agents (i.e.,  $\sum_i w_i(x^*) = x^*$ ), where

$$a^* = \arg \max_a x(a) - \sum_{i=1}^N \psi_i(a_i)$$

and  $x^* \equiv x(a^*)$ .

**Proof:** The proof is by construction. The principal accomplishes the first best by paying each agent a fixed wage  $w_i(x^*)$  when  $x = x^*$  and zero otherwise. By carefully choosing the first-best wage profile, agent's will find it optimal to produce the first best. Choose  $w_i(x^*)$  such that

$$w_i(x^*) - \psi_i(a_i^*) \geq 0 \quad \text{and} \quad \sum_{i=1}^N w_i(x^*) = x^*.$$

Such a wage profile can be found because  $x^*$  is optimal. Such a wage profile is also a Nash equilibrium. If all agents other than  $i$  choose their respective  $a_j^*$ , then agent  $i$  faces the following tradeoff: expend effort  $a_i^*$  so as to obtain  $x^*$  exactly and receive  $w_i(x^*) - \psi_i(a_i^*)$ , or shirk and receive a wage of zero. By construction of  $w_i(x^*)$ , agent  $i$  will choose  $a_i^*$  and so we have a Nash equilibrium.  $\square$

Note that although we have budget balance on the equilibrium path, we do not have budget balance off the equilibrium path. Holmst in further demonstrates that with budget balance one and off the equilibrium path (e.g., a partnership), the first-best cannot be obtained. Thus, in theory the principal can play an important role as a budget-breaker.

**Theorem 7** Assume each agent is risk neutral. Suppose that  $a^* \in \text{int } \mathcal{A}$  (i.e., all agent's provide some effort in the first-best allocation) and each  $\mathcal{A}_i$  is a closed interval,  $[\underline{a}_i, \bar{a}_i]$ . Then there do not exist wage schedules which are balanced and yield  $a^*$  as a Nash equilibrium in the noncooperative game.

Holmst in [1982] provides an intuitive in-text ‘‘proof’’ which he correctly notes relies upon a smooth wage schedule (we saw above that the optimal wage schedule may be discontinuous). The proof given in the appendix is more general but also complicated so I’ve

provided an alternative proof below which I believe is simpler.

**Proof:** Define  $a_j(a_i)$  by the relation  $x(a_{-j}^*, a_j) \equiv x(a_{-i}^*, a_i)$ . Since  $x$  is continuous and increasing and  $a^* \in \text{int } \mathcal{A}$ , a unique value of  $a_j(a_i)$  exists for  $a_i$  sufficiently close to  $a^*$ . The existence of a Nash equilibrium requires that for such an  $a_i$ ,

$$w_j(x(a^*)) - w_j(x(a_{-j}^*, a_j(a_i))) \equiv w_j(x(a^*)) - w_j(x(a_{-i}^*, a_i)) \geq \psi_j(a_j^*) - \psi_j(a_j(a_i)).$$

Summing up these relationships yields

$$\sum_{j=1}^N (w_j(x(a^*)) - w_j(x(a_{-i}^*, a_i))) \geq \sum_{j=1}^N (\psi_j(a_j^*) - \psi_j(a_j(a_i))).$$

Budget balance implies that the LHS of this equation is  $x(a^*) - x(a_{-i}^*, a_i)$ , and so

$$x(a^*) - x(a_{-i}^*, a_i) \geq \sum_{j=1}^N (\psi_j(a_j^*) - \psi_j(a_j(a_i))).$$

Because this must hold for all  $a_i$  close to  $a_i^*$ , we can divide by  $(a_i^* - a_i)$  and take the limit as  $a_i \rightarrow a_i^*$  to obtain

$$x_{a_i}(a^*) \geq \sum_{j=1}^N \psi'_j(a_j^*) \frac{x_{a_i}(a^*)}{x_{a_j}(a^*)}.$$

But the assumption that  $a^*$  is a first-best optimum implies that  $\psi'_j(a_j^*) = x_{a_j}(a^*)$ , which simplifies the previous inequality to  $x_{a_i}(a^*) \geq N x_{a_i}(a^*)$  – a contradiction because  $x$  is strictly increasing.  $\square$

### Remarks:

1. Risk neutrality is not required to obtain the result in Theorem 6. But with sufficient risk aversion, we can obtain the first best even in a partnership if we consider random contracts. For example, if agents are infinitely risk averse, by randomly distributing the output to a single agent whenever  $x = x(a^*)$ , agents can be given incentives to choose the correct action. Here, randomization allows you to break the “utility” budget, even though the wage budget is satisfied. Such an idea appears in Rasmusen [1987] and Legros and Matthews [1993]
2. As indicated in the proof, it is important that  $x$  is continuous over an interval  $\mathcal{A} \subset \mathbb{R}$ . Relaxing this assumption may allow us to get the first best. For example, if  $x$  were informative as to who cheated, the first best could be implemented by imposing a large fine on the shirker and distributing it to the other agents.
3. It is also important for the proof that  $a^* \in \text{int } \mathcal{A}$ . If, for example,  $a_i^* = \min_{a' \in \mathcal{A}_i} \psi_i(a')$  for some  $i$  (i.e.,  $i$ 's efficient contribution is to shirk), then  $i$  can be made the “principal” and the first best can be implemented.

4. Legros and Matthews [1993] extend Holmst in’s budget-breaking argument and establish necessary and sufficient conditions for a partnership to implement the efficient set of actions.

In particular, their conditions imply that partnerships with finite action spaces and a generic output function,  $x(a)$ , can implement the first best. For example, let  $N = 3$ ,  $\mathcal{A}_i = \{0, 1\}$ ,  $\psi_i = \psi$ , and  $a^* = (1, 1, 1)$ . Genericity of  $x(a)$  implies that  $x(a) = x(a')$  if  $a = a'$ . Genericity therefore implies that for any  $x = x^*$ , the identity of the shirker can be determined from the level of output. So letting

$$w_i(x) = \begin{cases} \frac{1}{3}x & \text{if } x = x^*, \\ \frac{1}{2}(F + x) & \text{if } j = i \text{ deviated,} \\ -F & \text{if } i \text{ deviated.} \end{cases}$$

For  $F$  sufficiently large, choosing  $a_i^*$  is a Nash equilibrium for agent  $i$ .

Note that determining the identity of the shirker is not necessary; only the identity of a non-shirker can be determined for any  $x = x^*$ . In such a case, the non-shirker can collect sufficiently large fines from the other players so as to make  $a^*$  a Nash equilibrium. (Here, the non-shirker acts as a budget-breaking principal.)

5. Legros and Matthews [1993] also show that asymptotic efficiency can be obtained if (i)  $\mathcal{A}_i \subset \mathbb{R}$ , (ii)  $\underline{a}_i \equiv \min_a \mathcal{A}_i$  and  $\bar{a}_i \equiv \max_a \mathcal{A}_i$  exist and are finite, and (iii),  $a_i^* \in (\underline{a}_i, \bar{a}_i)$ .

This result is best illustrated with the following example. Suppose that  $N = 2$ ,  $\mathcal{A}_i \equiv [0, 2]$ ,  $x(a) \equiv a_1 + a_2$ , and  $\psi_i(a_i) \equiv \frac{1}{2}a_i^2$ . Here,  $a^* = (1, 1)$ . Consider the following strategies. Agent 2 always chooses  $a_2 = a_2^* = 1$ . Agent 1 randomizes over the set  $\{\underline{a}_1, a_1^*, \bar{a}_1\} = \{0, 1, 2\}$  with probabilities  $\{\delta, 1 - 2\delta, \delta\}$ , respectively. We will construct wage schedules such that this is an equilibrium and show that  $\delta$  may be made arbitrarily small.

Note that on the equilibrium path,  $x \in [1, 3]$ . Use the following wage schedules for  $x \in [1, 3]$ :  $w_1(x) = \frac{1}{2}(x-1)^2$  and  $w_2(x) = x - w_1(x)$ . When  $x \notin [1, 3]$ , set  $w_1(x) = x + F$  and  $w_2(x) = -F$ . Clearly agent 2 will always choose  $a_2 = a_2^* = 1$  providing agent 1 plays his equilibrium strategy and  $F$  is sufficiently large. But if agent 2 plays  $a_2 = 1$ , agent 1 obtains  $U_1 = 0$  for any  $a_1 \in [0, 2]$ , and so the prescribed randomization strategy is optimal. Thus, we have a Nash equilibrium.

Finally, it is easy to verify that as  $\delta$  goes to zero, the first best allocation is obtained in the limit. One difficulty with this asymptotic mechanism is that the required size of the fine is  $F \geq \frac{1-2\delta+3\delta^2}{2\delta(2-\delta)}$ , so as  $\delta \rightarrow 0$ , the magnitude of the required fine explodes,  $F \rightarrow \infty$ . Another difficulty is that the strategies require very “unnatural” behavior by at least one of the agents.

6. When output is a stochastic function of actions, the first-best action profile may be sustained if the actions of the agents can be differentiated sufficiently and monetary transfers can be imposed as a function of output. Williams and Radner [1988] and Legros and Matsushima [1991] consider these issues.

## Sufficient Statistics and Relative Performance Evaluation.

We now suppose more realistically that  $x$  is stochastic. We also assume that actions affect a vector of contractible variables,  $y \in Y$ , via a distribution parameterization,  $F(y, a)$ . Of course, we allow the possibility that  $x$  is a component of  $y$ .

Assuming that the principal gets to choose the Nash equilibrium which the agents play, the principal's problem is to

$$\max_{a,w} \int_Y \left( E[x|a, y] - \sum_{i=1}^N w_i(y) \right) f(y, a) dy,$$

subject to  $(\forall i)$

$$\int_Y u_i(w_i(y)) f(y, a) - \psi_i(a_i) \geq \underline{U}_i,$$

$$a_i \in \arg \max_{a'_i \in \mathcal{A}_i} \int_Y u_i(w_i(y)) f(y, (a'_i, a_{-i})) - \psi_i(a'_i).$$

The first set of constraints are the IR constraints; the second set are the IC constraints. Note that the IC constraints imply that the agents are playing a Nash equilibrium amongst themselves. [Note: This is our first example of the principal designing a game for the agent's to play! We will see much more of this when we explore mechanism design.] Also note that when  $x$  is a component of  $y$  (e.g.,  $y = (x, s)$ ), we have  $E[x|a, y] = x$ .

Note that the actions of one agent may induce a distribution on  $y$  which is informative for the principal for the actions of a different agent. Thus, we need a new definition of statistical sufficiency to take account of this endogeneity.

**Definition 4**  $T_i(y)$  is **sufficient for  $y$  with respect to  $a_i$**  if there exists an  $g_i \geq 0$  and  $h_i \geq 0$  such that

$$f(y, a) \equiv h_i(y, a_{-i}) g_i(T_i(y), a), \quad \forall (y, a) \in Y \times \mathcal{A}.$$

$T(y) \equiv \{T_i(y)\}_{i=1}^N$  is **sufficient for  $y$  with respect to  $a$**  if  $T_i(y)$  is sufficient for  $y$  with respect to  $a_i$  for every agent  $i$ .

The following theorem is immediate.

**Theorem 8** If  $T(y)$  is sufficient for  $y$  with respect to  $a$ , then given any wage schedule,  $w(y) \equiv \{w_i(y)\}_i$ , there exists another wage schedule  $\tilde{w}(T(y)) \equiv \{w_i(T_i(y))\}_i$  that weakly Pareto dominates  $w(y)$ .

**Proof:** Consider agent  $i$  and take the actions of all other agents as given (since we begin with a Nash equilibrium). Define  $\tilde{w}_i(T_i)$  by

$$u_i(\tilde{w}_i(T_i)) \equiv \int_{\{y|T_i(y)=T_i\}} u_i(w_i(y)) \frac{f(y, a)}{g_i(T_i, a)} dy = \int_{\{y|T_i(y)=T_i\}} u_i(w_i(y)) h_i(y, a_{-i}) dy.$$

The agent's expected utility is unchanged under the new wage schedule,  $\tilde{w}_i(T_i)$ , and so IC and IR are unaffected. Additionally, the principal is weakly better off as  $u'' < 0$ , so by Jensen's inequality,

$$\tilde{w}_i(T_i) \leq \int_{\{y|T_i(y)=T_i\}} w_i(y) h_i(y, a_{-i}) dy.$$

Integrating over the set of  $T_i$ ,

$$\int_Y \tilde{w}_i(T_i(y)) f(y, a) dy \leq \int_Y w_i(y) f(y, a) dy.$$

This argument can be repeated for  $N$  agents, because in equilibrium the actions of  $N - 1$  agents can be taken as fixed parameters when examining the  $i$  agent.  $\square$

The intuition is straightforward: we constructed  $\tilde{w}$  so as not to affect incentives relative to  $w$ , but with improved risk-sharing, hence Pareto dominating  $w$ .

We would like to have the converse of this theorem as well. That is, if  $T(y)$  is not sufficient, we can strictly improve welfare by using additional information in  $y$ . We need to be careful here about our statements. We want to define a notion of insufficiency that pertains for all  $a$ . Along these lines,

**Definition 5**  $T(y)$  is **globally sufficient** iff for all  $a$ ,  $i$ , and  $T_i$

$$\frac{f_{a_i}(y, a)}{f(y, a)} = \frac{f_{a_i}(y', a)}{f(y', a)}, \text{ for almost all } y, y' \in \{y | T_i(y) = T_i\}.$$

Additionally,  $T(y)$  is **globally insufficient** iff for some  $i$  the above statement is false for all  $a$ .

**Theorem 9** Assume  $T(y)$  is globally insufficient for  $y$ . Let  $\{w_i(y) \equiv \tilde{w}_I(T_i(y))\}_i$  be a collection of non-constant wage schedules such that the agent's choices are unique in equilibrium. Then there exist wage schedules  $\hat{w}(y) = \{\hat{w}_i(y)\}_i$  that yield a strict Pareto improvement and induces the same equilibrium actions as the original  $w(y)$ .

The proof involves showing that otherwise the principal could do better by altering the optimal contract over the positive-measure subset of outcomes in which the condition for global sufficiency fails. See Holmst oin, [1982] page 332.

The above two theorems are useful for applications in agency theory. Theorem 8 says that randomization does not pay if the agent's utility function is separable; any uninformative noise should be integrated out. Conversely, Theorem 9 states that if  $T(y)$  is not sufficient for  $y$  at the optimal  $a$ , we can do strictly better using information contained in  $y$  that is not in  $T(y)$ .

**Application: Relative Performance Evaluation.** An important application is *relative performance evaluation*. Let's switch back to the state-space parameterization where  $\tilde{x} = x(a, \tilde{\theta})$  and  $\tilde{\theta}$  is a random variable. In particular, let's suppose that the information system of the principal is rich enough so that

$$\tilde{x}(a, \theta) \equiv \int_i x_i(a_i, \tilde{\theta}_i)$$

and *each*  $x_i$  is contractible. Ostensibly, we would think that each agent's wage should depend only upon its  $x_i$ . But the above two theorems suggest that this is not the case when the  $\theta_i$  are not independently distributed. In such a case, the output of one agent may be informative about the effort of another. We have the following theorem along these lines.

**Theorem 10** Assume the  $x_i$ 's are monotone in  $\theta_i$ . Then the optimal sharing rule of agent  $i$  depends on the individual  $i$ 's output alone if and only if the  $\theta_i$ 's are independently distributed.

**Proof:** If the  $\theta_i$ 's are independent, then the parameterized distribution satisfies

$$f(x, a) = \prod_{i=1}^N f_i(x_i, a_i).$$

This implies that  $T_i(x) = x_i$  is sufficient for  $x$  with respect to  $a_i$ . By theorem 8, it will be optimal to let  $w_i$  depend upon  $x_i$  alone.

Suppose instead that  $\theta_1$  and  $\theta_2$  are dependent but that  $w_1$  does not depend upon  $x_2$ . Since in equilibrium  $a_2$  can be inferred, assume that  $x_2 = \theta_2$  without loss of generality and subsume  $a_2$  in the distribution. The joint distribution of  $x_2 = \theta_2$  conditional on  $a_1$  is given by

$$f(x_1, \theta_2, a_1) = \tilde{f}(x_1^{-1}(a_1, x_1), \theta_2),$$

where  $\tilde{f}(\theta_1, \theta_2)$  is the joint distribution of  $\theta_1$  and  $\theta_2$ . It follows that

$$\frac{f_{a_1}(x_1, \theta_2, a_1)}{f(x_1, \theta_2, a_1)} = \frac{\tilde{f}_{\theta_1}(x_1^{-1}(a_1, x_1), \theta_2)}{\tilde{f}(x_1^{-1}(a_1, x_1), \theta_2)} \frac{\partial x_1^{-1}(a_1, x_1)}{\partial a_1}.$$

Since  $\theta_1$  and  $\theta_2$  are dependent,  $\frac{\tilde{f}_{\theta_1}}{\tilde{f}}$  depends upon  $\theta_2$ . Thus  $T$  is globally insufficient and theorem 9 applies, indicating that  $w_1$  should depend upon information in  $x_2$ .  $\square$

### Remarks on Sufficient Statistics and Relative Performance:

1. The idea is that competition is not useful per se, but only as a way to get a more precise signal of an agent's action. With independent shocks, relative performance evaluation only adds noise and reduces welfare.
2. There is a literature on tournaments by Lazear and Rosen [1981] which indicates how basing wages on a tournament among workers can increase effort. Nalebuff and Stiglitz [1983] and Green and Stokey [1983] have made related points. But such tournaments are generally suboptimal as only with very restrictive technology will ordinal rankings be a sufficient statistic. The benefit of tournaments is that it is less susceptible to output tampering by the principal, since in any circumstance the wage bill is invariant.
3. All of this literature on teams presupposes that the principal gets to pick the equilibrium the agents will play. This may be unsatisfactory as better equilibria (from the agents' viewpoints) may exist. Mookherjee [1984] considers the multiple-equilibria problem in his examination of the multi-agent moral hazard problem and provides an illuminating example of its manifestation.

## General Remarks on Teams and Partnerships:

1. Itoh [1991] has noted that sometimes a principal may want to reward agent  $i$  as a function of agent  $j$ 's output, even if the outputs are independently distributed, when "teamwork" is desirable. The reason such dependence may be desirable is that agent  $i$ 's effort may be a vector which contains a component that improves agent  $j$ 's output. Itoh's result is not in any way at odds with our sufficient statistic results above; the only change is that efforts are multidimensional and so the principal's program is more complicated. Itoh characterizes the optimal contract in this setting.
2. It is possible that agents may get together and collude over their production and write (perhaps implicit) contracts amongst themselves. For example, some of the risk which the principal imposes for incentive reasons may be reduced by the agents via risk pooling. This obviously hurts the principal as it places an additional constraint on the optimal contract: the marginal utilities of agents will be equated across states. This cost of collusion has been noted by Itoh [1993]. Itoh also considers a *benefit* of collusion: when efforts are mutually observable by agents, the principal may be better off. The idea is that through the principal's choice of wage schedules, the agents can be made to police one another and increase effort through their induced side contracts. Or more precisely, the set of implementable contracts increases when agents can contract on each other's effort. Thus, collusion may (on net) be beneficial. The result that "collusion is beneficial to the principal" must be taken carefully, however. We know that when efforts are mutually observable by the agents there exist revelation mechanisms which allow the principal to obtain the first best as a Nash equilibrium (where agents are instructed to report shirking to the principal). We normally think that such first-best mechanisms are problematic *because of* collusion or coordination on other detrimental equilibria. It is not the collusion of Itoh's model that is beneficial – it is the mutual observation of efforts by the agents. Itoh's result may be more appropriately stated as "mutual effort observation may increase the principal's profits *even if* agents collude."

### 2.1.3 Extensions: A Rationale for Linear Contracts

We now briefly turn to an important paper by Holmstøin and Milgrom [1987] which provides an economic setting and conditions under which contracts will be linear in aggregates. This paper is fundamental both in explaining the simplicity of real world contracts and providing contract theorists with a rationalization for focusing on linear contracts. Although the model of the paper is dynamic in structure, its underlying stationarity (i.e., CARA utility and repeated environment) generates a static form:

The optimal dynamic incentive scheme can be computed as if the agent were choosing the mean of a normal distribution only once and the principal were restricted to offering a linear contract.

We thus consider Holmstøin and Milgrom's [1987] contribution here as an examination of static contracts rather than dynamic contracts.

### One-period Model:

There are  $N + 1$  possible outcomes:  $x_i \in \{x_0, \dots, X_N\}$ , with probability of occurrence given by the vector  $p = (p_0, \dots, p_N)$ . We assume that the agent directly chooses  $p \in \Delta(N)$  at a cost of  $c(p)$ . The principal offers a contract  $w(x_i) = \{w_0, \dots, w_N\}$  as a function of outcomes. Both principal and agent have exponential utility functions (to avoid problems of wealth effects).

$$U(w - c(p)) \equiv -e^{-r(w-c(p))},$$

$$V(x - w) \equiv \begin{cases} -e^{-R(x-w)} & \text{if } R > 0, \\ x - w & \text{if } R = 0. \end{cases}$$

Assume that  $R = 0$  for now. The principal solves

$$\max_{w, p} \sum_{i=0}^N p_i(x_i - w_i) \quad \text{subject to}$$

$$p \in \arg \max_p \sum_{i=0}^N p_i U(w - c(p)),$$

$$\sum_{i=0}^N p_i U(w - c(p)) \geq U(\underline{w}),$$

where  $\underline{w}$  is the certainty equivalent of the agent's outside opportunity. (I will use generally underlined variables to represent certainty equivalents.)

Given our assumption of exponential utility, we have the following result immediately.

**Theorem 11** Suppose that  $(w^*, p^*)$  solves the principal's one-period program for some  $\underline{w}$ . Then  $(w^* + \underline{w}' - \underline{w}, p^*)$  solves the program for an outside certainty equivalent of  $\underline{w}'$ .

**Proof:** Because utility is exponential,

$$\sum_{i=0}^N p_i U(w^*(x_i) - c(p^*)) = -U(-\underline{w} + \underline{w}') \quad \sum_{i=0}^N p_i U(w^*(x_i) + \underline{w}' - \underline{w}).$$

Thus,  $p^*$  is still incentive compatible and the IR constraint is satisfied for  $U(\underline{w}')$ . Similarly, given the principal's utility is exponential, the optimal choice of  $p^*$  is unchanged.  $\square$

The key here is that there are absolutely no wealth effects in this model. This will be an important ingredient in our proofs below.

### T-period Model:

Now consider the multi-period problem where the agent chooses a probability each period after having observed the history of outputs up until that time. Let superscripts



denote histories of variables; i.e.,  $X^t = \{x_0, \dots, x_t\}$ . The agent gets paid at the end of the period,  $w(X^t)$  and has a combined cost of effort equal to  $\sum_{\tau=0}^t c(p_\tau)$ . Thus,

$$U(w, \{p_t\}_t) = -e^{-r(w - \sum_{\tau=0}^t c(p_\tau))}.$$

Because the agent observes  $X^{t-1}$  before deciding upon  $p_t$ , for a given wage schedule we can write  $p_t(X^{t-1})$ . We want to first characterize the wage schedule which implements an arbitrary  $\{p_t(X^{t-1})\}_t$  effort function. We use dynamic programming to this end.

Let  $\mathcal{U}_t$  be the agent's expected utility (ignoring past effort costs) from date  $t$  forward. Thus,

$$\mathcal{U}_t(X^t) \equiv E U \left( w(X^T) - \sum_{\tau=t+1}^T c(p_\tau) \mid X^t \right).$$

Note here that  $\mathcal{U}_t$  differs from a standard value function by the constant  $U(-\sum_{\tau=t}^T c(p_\tau))$ . Let  $\underline{w}_t(X^t)$  be the certain equivalent of income of  $\mathcal{U}_t$ . That is,  $U(\underline{w}_t(X^t)) \equiv \mathcal{U}_t(X^t)$ . Note that  $\underline{w}_t(X^{t-1}, x_{it})$  is the certain equivalent for obtaining output  $x_i$  in period  $t$  following a history of  $X^{t-1}$ .

To implement  $p_t(X^{t-1})$ , it must be the case that

$$p_t(X^{t-1}) \in \arg \max_{p_t} \prod_{i=0}^N p_{it} U(\underline{w}_t(X^{t-1}, x_{it}) - c(p_t)),$$

where we have dropped the irrelevant multiplicative constant.

Our previous theorem 11 applies:  $p_t(X^{t-1})$  is implementable and yields certainty equivalent  $\underline{w}_{t-1}(X^{t-1})$  iff  $p_t(X^{t-1})$  is also implemented by

$$\tilde{w}_t(x_{it} | p_t(X^{t-1})) \equiv \underline{w}_t(X^{t-1}, x_{it}) - \underline{w}_{t-1}(X^{t-1})$$

with a certainty equivalent of  $\underline{w} = 0$ .

Rearranging the above relationship,

$$\underline{w}_t(X^{t-1}, x_{it}) = \tilde{w}_t(x_{it} | p_t(X^{t-1})) + \underline{w}_{t-1}(X^{t-1}).$$

Integrating this difference equation from  $t = 0$  to  $T$  yields

$$w(X^T) \equiv \underline{w}_T(X^T) = \sum_{t=0}^T \tilde{w}_t(x_{it} | p_t(X^{t-1})) + \underline{w}_0,$$

or in other words, the end of contract wage is the sum of the individual single-period wage schedules for implementing  $p_t(X^{t-1})$ .

Let  $\tilde{w}_t(p_t(X^{t-1}))$  be an  $N + 1$  vector over  $i$ . Then rewriting,

$$w(X^T) = \sum_{t=1}^T \tilde{w}_t(p_t(X^{t-1})) \cdot (A^t - A^{t-1}) + \underline{w}_0,$$

where  $A^t = (A_0^t, \dots, A_N^t)$  and  $A_i^t$  is an account that gives the number of times outcome  $i$  has occurred up to date  $t$ .

We thus have characterized a wage schedule,  $w(X^T)$ , for implementing  $p_t(X^{t-1})$ . Moreover, Holmst om and Milgrom show that if  $c$  is differentiable and  $p_t \in \text{int}\Delta(N)$ , such a wage schedule is uniquely defined. We now wish to find the optimal contract.

**Theorem 12** The optimal contract is to implement  $p_t(X^{t-1}) = p^* \forall t$  and offer the wage schedule

$$w(X^T) = \sum_{t=1}^T w(x_t, p^*) = w(p^*) \cdot A^T.$$

**Proof:** By induction. The theorem is true by definition for  $T = 1$ . Suppose that it holds for  $T = \tau$  and consider  $T = \tau + 1$ . Let  $\mathcal{V}_T^*$  be the principal's value function for the  $T$ -period problem. The value of the contract to the principal is

$$-e^{Rw_0} E[V(x_{t=1} - w_{t=1})E[V(\sum_{t=2}^{\tau+1} (x_t - w_t) | X^1)]] \leq -e^{Rw_0} E[V(x_{t=1} - w_{t=1})\mathcal{V}_\tau^*] \leq -e^{Rw_0} \mathcal{V}_1^* \mathcal{V}_\tau^*.$$

At  $p_t = p^*$ ,  $w_t = w(x_t, p^*)$ , this upper bound is met.  $\square$

**Remarks:**

1. Note very importantly that the optimal contract is *linear in accounts*. Specifically,

$$w(X^T) = \sum_{t=1}^T w(x_t, p^*) = \sum_{i=0}^N w(x_i, p^*) \cdot A_i^T,$$

or letting  $\alpha_i \equiv w(x_i, p^*) - w(x_0, p^*)$  and  $\beta \equiv T \cdot w(x_0, p^*)$ ,

$$w(X^T) = \sum_{i=1}^N \alpha_i A_i^T + \beta.$$

This is not generally linear in profits. Nonetheless, many applied economists typically take Holmst in and Milgrom's result to mean linearity in profits for the purposes of their applications.

2. If there are only two accounts, such as success or failure, then wages are linear in "profits" (i.e., successes). From above we have

$$w(X^T) = \alpha A_1^T + \beta.$$

Not surprisingly, when we take the limit as this binomial process converges to unidimensional Brownian motion, we preserve our linearity in profits result. With more than two accounts, this is not so. Getting an output of 50 three times is not the same as getting the output of 150 once and 0 twice.

3. Note that the history of accounts is irrelevant. Only total instances of outputs are important. This is also true in the continuous case. Thus,  $A^T$  is "sufficient" with respect to  $X^T$ . This is not inconsistent with Holmst in [1979] and Shavell [1979]. Sufficiency notions should be thought of as sufficient information regarding the binding constraints. Here, the binding constraint is shifting to another constant action, for which  $A^T$  is sufficient.

4. The key to our results are stationarity which in turn is due exclusively to time-separable CARA utility and an i.i.d. stochastic process.

### Continuous Model:

We now consider the limit as the time periods become infinitesimal. We now want to ask what happens if the agent can continuously vary his effort level and observe the realizations of output in real time.

### Results:

1. In the limit, we obtain a linearity in accounts result, where the accounts are movements in the stochastic process. With unidimensional Brownian motion, (i.e., the agent controls the drift rate on a one-dimensional Brownian motion process), we obtain linearity in profits.
2. Additionally, in the limit, if only a subset of accounts can be contracted upon (specifically, a linear aggregate), then the optimal contract will be linear in those accounts. Thus, if only profits are contractible, we will obtain the linearity in profits result in the limit – even when the underlying process is multinomial Brownian motion. This does not happen in the discrete case. The intuition roughly is that in the limit, information is lost in the aggregation process, while in the discrete case, this is not the case.
3. If the agent must take all of his actions simultaneously at  $t = 0$ , then our results do not hold. Instead, we are in the world of static nonlinear contracts. In a continuum, Mirrlees's example would apply, and we could obtain the first best.

### The Simple Analytics of Linear Contracts:

To see the usefulness of Holmström and Milgrom's [1987] setting for simple comparative statics, consider the following model. The agent has exponential utility with a CARA parameter of  $r$ ; the principal is risk neutral. Profits (excluding wages) are  $x = \mu + \varepsilon$ , where  $\mu$  is the agent's action choice (the drift rate of a unidimensional Brownian process) and  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ . The cost of effort is  $c(\mu) = \frac{k}{2}\mu^2$ .

Under the full information first-best contract,  $\mu^{FB} = \frac{1}{k}$ , the agent is paid a constant wage to cover the cost of effort,  $w^{FB} = \frac{1}{2k}$ , and the principal receives net profits of  $\pi = \frac{1}{2k}$ .

When effort is not contractible, Holmström and Milgrom's linearity result tells us that we can restrict attention to wage schedules of the form  $w(x) = \alpha x + \beta$ . With this contract,

the agent's certainty equivalent<sup>2</sup> upon choosing an action  $\mu$  is

$$\alpha\mu + \beta - \frac{k}{2}\mu^2 - \frac{r}{2}\alpha^2\sigma^2.$$

The first-order condition is  $\alpha = \mu k$  which is necessary and sufficient because the agent's utility function is globally concave in  $\mu$ .

It is very important to note that the utilities possibility frontier for the principal and agent is *linear* for a given  $(\alpha, \mu)$  and independent of  $\beta$ . The independence of  $\beta$  is an artifact of CARA utility (that's our result from Theorem refce above), and the linearity is due to the combination of CARA utility and normally distributed errors (the latter of which is due to the central limit theorem).

As a consequence, the principal's optimal choice of  $(\alpha, \mu)$  is independent of  $\beta$ ;  $\beta$  is chosen solely to satisfy the agent's IR constraint. Thus, the principal solves

$$\max_{\alpha, \mu} \mu - \frac{k}{2}\mu^2 - \frac{r}{2}\alpha^2\sigma^2,$$

subject to  $\alpha = \mu k$ . The solution gives us  $(\alpha^*, \mu^*, \pi^*)$ :

$$\alpha^* = (1 + rk\sigma^2)^{-1},$$

$$\mu^* = (1 + rk\sigma^2)^{-1}k^{-1} = \alpha^* \mu^{FB} < \mu^{FB},$$

$$\pi^* = (1 + rk\sigma^2)^{-1}(2k)^{-1} = \alpha^* \pi^{FB} < \pi^{FB}.$$

The simple comparative statics are immediate. As either  $r$ ,  $k$ , or  $\sigma^2$  decrease, the power of the optimal incentive scheme increases (i.e.,  $\alpha^*$  increases). Because  $\alpha^*$  increases, effort and profits also increase closer toward the first best. Thus when risk aversion, the uncertainty in measuring effort, or the curvature of the agent's effort function decrease, we move toward the first best. The intuition for why the curvature of the agent's cost function matters can be seen by totally differentiating the agent's first-order condition for effort. Doing so, we find that  $\frac{d\mu}{d\alpha} = \frac{1}{C''(\mu)} = \frac{1}{k}$ . Thus, lowering  $k$  makes the agent's effort choice more responsive to a change in  $\alpha$ .

### Remarks:

1. Consider the case of additional information. The principal observes an additional signal,  $y$ , which is correlated with  $\varepsilon$ . Specifically,  $E[y] = 0$ ,  $V[y] = \sigma_y^2$ , and  $Cov[\varepsilon, y] = \rho\sigma_y\sigma_e$ . The optimal wage contract is linear in both aggregates:  $w(x, y) = \alpha_1 x + \alpha_2 y + \beta$ . Solving for the optimal schemes, we have

$$\alpha_1^* = (1 + rk\sigma_\varepsilon^2(1 - \rho^2))^{-1},$$

---

<sup>2</sup>Note that the moment generating function for a normal distribution is  $M_x(t) = e^{\mu t + \frac{1}{2}\sigma^2 t^2}$  and the defining property of the m.g.f. is that  $E_x[e^{tx}] = M_x(t)$ . Thus,

$$E_\varepsilon[e^{-r(\alpha\mu + \alpha\varepsilon + \beta - C(\mu))}] = e^{-r(\alpha\mu + \beta - C(\mu)) + \frac{1}{2}\alpha^2 r^2 \sigma^2}.$$

Thus, the agent's certainty equivalent is  $\alpha\mu + \beta - C(\mu) - \frac{r}{2}\alpha^2\sigma^2$ .

$$\alpha_2^* = -\alpha_1 \frac{\sigma_y}{\sigma_\varepsilon} \rho.$$

As before,  $\mu^* = \alpha_1^* \mu^{FB}$  and  $\pi^* = \alpha_1^* \pi^{FB}$ . It is as if the outside signal reduces the variance on  $\varepsilon$  from  $\sigma_\varepsilon^2$  to  $\sigma_\varepsilon^2(1 - \rho^2)$ . When either  $\rho = 1$  or  $\rho = -1$ , the first-best is obtainable.

2. The allocation of effort across tasks may be greatly influenced by the nature of information. To see this, consider a symmetric formulation with two tasks:  $x_1 = \mu_1 + \varepsilon_1$  and  $x_2 = \mu_2 + \varepsilon_2$ , where  $\varepsilon_i \sim \mathcal{N}(0, \sigma_i^2)$  and are independently distributed across  $i$ . Suppose also that  $C(\mu) = \frac{1}{2}\mu_1^2 + \frac{1}{2}\mu_2^2$  and the principal's net profits are  $\pi = x_1 + x_2 - w$ . If only  $x = x_1 + x_2$  were observed, then the optimal contract has  $w(x) = \alpha x + \beta$ , and the agent would equally devote his attention across tasks. Additionally, if  $\sigma_1 = \sigma_2$  and the principal can contract on both  $x_1$  and  $x_2$ , the optimal contract has  $\alpha_1 = \alpha_2$  and so again the agent equally allocates effort across tasks.

Now suppose that  $\sigma_1 < \sigma_2$ . The resulting first-order conditions imply that  $\alpha_1^* > \alpha_2^*$ . Thus, optimal effort allocation may be entirely determined by the information structure of the contracting environment. The intuition here is that the “price” of inducing effort on task 1 is lower for the principal because information is more informative. Thus, the principal will “buy” more effort from the agent on task 1 than task 2.

#### 2.1.4 Extensions: Multi-task Incentive Contracts

We now consider more explicitly the implications of multiple tasks within a firm using the linear contracting model of Holmström and Milgrom [1987]. This analysis closely follows Holmström and Milgrom [1991].

##### The Basic Linear Model with Multiple Tasks:

The principal can contract on the following  $k$  vector of aggregates:

$$x = \mu + \varepsilon,$$

where  $\varepsilon \sim \mathcal{N}(0, \Sigma)$ . The agent chooses a vector of efforts,  $\mu$ , at a cost of  $C(\mu)$ .<sup>3</sup> The agent's utility is exponential with CARA parameter of  $r$ . The principal is risk neutral, offers wage schedule  $w(x) = \alpha'x + \beta$ , and obtains profits of  $B(\mu) - w$ . [Note,  $\alpha$  and  $\mu$  are vectors;  $B(\mu)$ ,  $\beta$  and  $w(x)$  are scalars.]

As before, CARA utility and normal errors implies that the optimal contract solves

$$\max_{\alpha, \mu} B(\mu) - C(\mu) - \frac{r}{2} \alpha' \Sigma \alpha,$$

such that

$$\mu \in \arg \max_{\tilde{\mu}} \alpha' \tilde{\mu} - C(\tilde{\mu}).$$

Given the optimal  $(\alpha, \mu)$ ,  $\beta$  is determined so as to meet the agent's IR constraint:

$$\beta = \underline{w} - \alpha' \mu + C(\mu) + \frac{r}{2} \alpha' \Sigma \alpha.$$

---

<sup>3</sup>Note that Holmström and Milgrom [1991] take the action vector to be  $t$  where  $\mu(t)$  is determined by the action. We'll concentrate on the choice of  $\mu$  as the primitive.

The agent's first-order condition (which is both necessary and sufficient) satisfies

$$\alpha_i = C_i(\mu), \quad \forall i,$$

where subscripts on  $C$  denote partial derivatives with respect to the indicated element of  $\mu$ . Comparative statics on this equation reveal that

$$\frac{\partial \mu}{\partial \alpha} = [C_{ij}(\mu)]^{-1}.$$

This implies that in the simple setting where  $C_{ij} = 0 \quad \forall i = j$ , that  $\frac{d\mu_i}{d\alpha_i} = \frac{1}{C_{ii}(\mu)}$ . Thus, the marginal affect of a change in  $\alpha$  on effort is inversely related to the curvature of the agent's cost of effort function.

We have the following theorem immediately.

**Theorem 13** The optimal contract satisfies

$$\alpha^* = (I + r[C_{ij}(\mu^*)]\Sigma)^{-1} B'(\mu^*).$$

**Proof:** Form the Lagrangian:

$$\mathcal{L} \equiv B(\mu) - C(\mu) - \frac{r}{2}\alpha'\Sigma\alpha + \lambda(\alpha - C'(\mu)),$$

where  $C'(\mu) = [C_i(\mu)]$ . The  $2k$  first-order conditions are

$$\begin{aligned} B'(\mu^*) - C'(\mu^*) - \lambda[C_{ij}(\mu^*)] &= 0, \\ -r\Sigma\alpha^* + \lambda &= 0. \end{aligned}$$

Substituting out  $\lambda$  and solving for  $\alpha^*$  produces the desired result.  $\square$

**Remarks:**

1. If  $\varepsilon_i$  are independent and  $C_{ij} = 0$  for  $i = j$ , then

$$\alpha_i^* = B_i(\mu^*)(1 + rC_{ii}(\mu^*)\sigma_i^2)^{-1}.$$

As  $r$ ,  $\sigma_i$ , or  $C_{ii}$  decrease,  $\alpha_i^*$  increases. This result was found above in our simple setting of one task.

2. Given  $\mu^*$ , the cross partial derivatives of  $B$  are unimportant for the determination of  $\alpha^*$ . Only cross partials in the agent's utility function are important (i.e.,  $C_{ij}$ ).

**Simple Interactions of Multiple Tasks:**

Consider the setting where there are two tasks, but where the effort of only the first task can be measured:  $\sigma_2 = \infty$  and  $\sigma_{12} = 0$ . A motivating example is a teacher who teaches basic skills (task 1) which is measurable via student testing and higher-order skills such as

creativity, etc. (task 2) which is inherently unmeasurable. The question is how we want to reward the teacher on the basis of basic skill test scores.

Suppose that under the optimal contract  $\mu^* > 0$ ; that is, both tasks will be provided at the optimum.<sup>4</sup> Then the optimal contract satisfies  $\alpha_2^* = 0$  and

$$\alpha_1^* = \frac{B_1(\mu^*) - B_2(\mu^*) \frac{C_{12}(\mu^*)}{C_{22}(\mu^*)}}{1 + r\sigma_1^2} \left[ C_{11}(\mu^*) - \frac{C_{12}(\mu^*)^2}{C_{22}(\mu^*)} \right]^{-1}.$$

Some interesting conclusions emerge.

1. If effort levels across tasks are substitutes (i.e.,  $C_{12} > 0$ ), the more positive the cross-effort effect (i.e., more substitutable the effort levels), the lower is  $\alpha_1^*$ . (In our example, if a teacher only has 8 hours a day to teach, the optimal scheme will put less emphasis on basic skills the more likely the teacher is to substitute away from higher-order teaching.) If effort levels are complements, the reverse is true, and  $\alpha_1^*$  is increased.
2. The above result has a flavor of the public finance results that when the government can only tax a subset of goods, it should tax them more or less depending upon whether the taxable goods are substitutes or complements with the un-taxable goods. See, for example, Atkinson and Stiglitz [1980 Ch. 12] for a discussion concerning taxation on consumption goods when leisure is not directly taxable.
3. There are several reasons why the optimal contract may have  $\alpha_1^* = 0$ .
  - Note that in our example,  $\alpha_1^* < 0$  if  $B_1 < B_2 C_{12}/C_{22}$ . Thus, if the agent can freely dispose of  $x_1$ , the optimal constrained contract has  $\alpha_1^* = 0$ . No incentives are provided.
  - Suppose that technologies are otherwise symmetric:  $C(\mu) = c(\mu_1 + \mu_2)$  and  $B(\mu_1, \mu_2) \equiv B(\mu_2, \mu_1)$ . Then  $\alpha_1^* = \alpha_2^* = 0$ . Again, no incentives are provided.
  - Note that if  $C_i(0) > 0$ , there is a fixed cost to effort. This implies that a corner solution may emerge where  $\alpha_i^* = 0$ . A final reason for no incentives.

### Application: Limits on Outside Activities.

Consider the principal's problem when an additional control variable is added: the set of allowable activities. Suppose that the principal cares only about effort devoted to task 0:  $\pi = \mu_0 - w$ . In addition, there are  $N$  potential tasks which the agent could spend effort on and which increase the agent's personal utility. We will denote the set of these tasks by  $K = \{1, \dots, N\}$ . The principal has the ability to exclude the agent from any subset of these activities, allowing only tasks or activities in the subset set  $A \subset K$ . Unfortunately, the principal can only contract over  $x_0$ , and so  $w(x) = \alpha x_0 + \beta$ .

It is not always profitable, even in the full-information setting, to exclude these tasks from the agent, because they may be efficient and therefore reduce the principal's wage

---

<sup>4</sup>Here we need to assume something like  $C_2(\mu_1, \mu_2) < 0$  for  $\mu_2 \leq 0$  so that without any incentives on task 2, the agent still allocates some effort on task 2. In the teaching example, absent incentives, a teacher will still teach some higher-order skills.

bill. As a motivating example, allowing an employee to make personal calls on the company WATTS line may be a cheap perk for the firm to provide and additionally lowers the necessary wage which the firm must pay. Unfortunately, the agent may then spend all of the day on the telephone rather than at work.

Suppose that the agent's cost of effort is

$$C(\mu) = c \mu_0 + \sum_{i=1}^N \mu_i - \sum_{i=1}^N u_i(\mu_i).$$

The  $u_i$  functions represent the agent's personal utility from allocating effort to task  $i$ ;  $u_i$  is assumed to be strictly concave and  $u_i(0) = 0$ . The principal's expected returns are simply  $B(\mu) = p\mu_0$ .

We first determine the principal's optimal choice of  $A^*(\alpha)$  for a given  $\alpha$ , and then we solve for the optimal  $\alpha^*$ . The first-order condition which characterizes the agent's optimal  $\mu_0$  is

$$\alpha = c' + \sum_{i=0}^N \mu_i,$$

and (substituting)

$$\alpha = v'_i(\mu_i), \quad \forall i.$$

Note that the choice of  $\mu_i$  depends only upon  $\alpha$ . Thus, if the agent is allowed an additional personal task,  $k$ , the agent will allocate time away from task 0 by an amount equal to  $v_k^{-1}(\alpha)$ . The benefit of allowing the agent to spend time on task  $k$  is  $v_k(\mu_k(\alpha))$  (via a reduced wage) and the (opportunity) cost is  $p\mu_k(\alpha)$ . Therefore, the optimal set of tasks for a given  $\alpha$  is

$$A^*(\alpha) = \{k \in K | v_k(\mu_k(\alpha)) > p\mu_k(\alpha)\}.$$

We have the following results for a given  $\alpha$ .

**Theorem 14** Assume that  $\alpha$  is such that  $\mu(\alpha) > 0$  and  $\alpha < p$ . Then the optimal set of allowed tasks is given by  $A^*(\alpha)$  which is monotonically expanding in  $\alpha$  (i.e.,  $\alpha \leq \alpha'$ , then  $A^*(\alpha) \subset A^*(\alpha')$ ).

**Proof:** That the optimal set of allowed tasks is given by  $A^*(\alpha)$  is true by construction. The set  $A^*(\alpha)$  is monotonically expanding in  $\alpha$  iff  $v_k(\mu_k(\alpha)) - p\mu_k(\alpha)$  is increasing in  $\alpha$ . I.e.,

$$[v'_k(\mu_k(\alpha)) - p] \frac{d\mu_k(\alpha)}{d\alpha} = [\alpha - p] \frac{1}{v''_k(\mu_k(\alpha))} > 0.$$

□

**Remarks:**

1. The fundamental premise of exclusion is that incentives can be given by either increasing  $\alpha$  on the relevant activity or decreasing the opportunity cost of effort (i.e., by reducing the benefits of substitutable activities).



2. The theorem indicates a basic proposition with direct empirical content: *responsibility (large  $\alpha$ ) and authority (large  $A^*(\alpha)$ ) should go hand in hand*. An agent with high-powered incentives should be allowed the opportunity to expend effort on more personal activities than someone with low-powered incentives. In the limit when  $\sigma \rightarrow 0$  or  $r \rightarrow 0$ , the agent is residual claimant  $\alpha^* = 1$ , and so  $A^*(1) = K$ . Exclusion will be more frequently used the more costly it is to supply incentives.
3. Note that for  $\alpha$  small enough,  $\mu(\alpha) = 0$ , and the agent is not hired.
4. The set  $A^*(\alpha)$  is independent of  $r$ ,  $\sigma$ ,  $C$ , etc. These variables only influence  $A^*(\alpha)$  via  $\alpha$ . Therefore, an econometrician can regress  $\|A^*(\alpha)\|$  on  $\alpha$ , and  $\alpha$  on  $(r, \sigma, \dots)$  to test the multi-task theory. See Holmstrom and Milgrom [1994].

Now, consider the choice of  $\alpha^*$  given the function  $A^*(\alpha)$ .

**Theorem 15** Providing that  $\mu(\alpha^*) > 0$  at the optimum,

$$\alpha^* = p \left( 1 + r\sigma^2 \left( \frac{1}{c''(\mu_i(\alpha^*))} + \sum_{k \in A^*(\alpha^*)} \frac{1}{v_k''(\mu_k(\alpha^*))} \right) \right)^{-1}.$$

The proof of this theorem is an immediate application of our first multi-task characterization theorem. Additionally, we have the following implications.

**Remarks:**

1. The theorem indicates that when either  $r$  or  $\sigma$  decreases,  $\alpha^*$  increases. [Note that this implication is not immediate because  $\sigma^*$  appears on both sides of the equation; some manipulation is required. With quadratic cost and benefit functions, this is trivial.] By our previous result on  $A^*(\alpha)$ , the set of allowable activities also increases as  $\alpha^*$  increases.
2. Any personal task excluded in the first-best arrangement (i.e.,  $v_k'(0) < p$ ) will be excluded in the second-best optimal contract given our construction of  $A^*(\alpha)$  and the fact that  $v_k$  is concave. This implies that there will be more constraints on agent's activities when performance rewards are weak due to a noisy environment.
3. Following the previous remark, one can motivate rigid rules which limit an agent's activities (seemingly inefficiently) as a way of dealing with substitution possibilities. Additionally, when the "personal" activity is something such as rent-seeking (e.g., inefficiently spending resources on your boss to increase your chance of promotion), a firm may wish to restrict an agent's access to such an activity or withdraw the bosses discretion to promote employees so as to reduce this inefficient activity. This idea was formalized by Milgrom [1988] and Milgrom and Roberts [1988].
4. This activity exclusion idea can also explain why firms may not want to allow their employees to "moonlight". Or more importantly, why a firm may wish to use an

internal sales force which is not allowed to sell other firms' products rather than an external sales force whose activities vis-a-vis other firms cannot be controlled.

**Application: Task Allocation Between Two Agents.**

Now consider two agents,  $i = 1, 2$ , who are needed to perform a continuum of tasks indexed by  $t \in [0, 1]$ . Each agent  $i$  expends effort  $\mu_i(t)$  on task  $t$ ; total cost of effort is  $C(\int \mu_i(t)dt)$ . The principal observes  $x(t) = \mu(t) + \varepsilon(t)$  for each task, where  $\sigma^2(t) > 0$  and  $\mu(t) \equiv \mu_1(t) + \mu_2(t)$ . The wages paid to the agents are given by:

$$w_i(x) = \int_0^1 \alpha_i(t)x(t)dt + \beta_i.$$

By choosing  $\alpha_i(t)$ , the principal allocates agents to the various tasks. For example, when  $\alpha_1(.4) > 0$  but  $\alpha_2(.4) = 0$ , only agent 1 will work on task .4.

Two results emerge.

1. For any required effort function  $\mu(t)$  defined on  $[0, 1]$ , it is never optimal to assign two agents to the same task:  $\alpha_1^*(t)\alpha_2^*(t) \equiv 0$ . This is quite natural given the teams problem which would otherwise emerge.
2. More surprisingly, suppose that the principal must obtain a uniform level of effort  $\mu(t) = 1$  across all tasks. At the optimum, if  $\int \mu_i(t)dt < \int \mu_j(t)dt$ , then the hardest to measure tasks go to agent  $i$  (i.e., all tasks  $t$  such that  $\sigma(t) \geq \bar{\sigma}$ .) This results because you want to avoid the multi-task problems which occur when the various tasks have vastly different measurement errors. Thus, the principal wants information homogeneity. Additionally, the agent with the hard to measure tasks exerts lower effort and receives a lower “normalized commission” because the information structure is so noisy.

**Application: Common Agency.** Bernheim and Whinston [1986] were the first to undertake a detailed study of the phenomena of common agency with moral hazard. “Common agency” refers to the situation in which several principals contract with the same agent in common. The interesting economics of this setting arise when one principal’s contract imposes an externality on the contracts of the others.

Here, we follow Dixit [1996] restricting attention to linear contracts in the simplest setting of  $n$  independent principals who simultaneously offer incentive contracts to a single agent who controls the  $m$ -dimensional vector  $t$  which in turn effects the output vector  $x \in \mathbb{R}^m$ . Let  $x = t + \varepsilon$ , where  $x, t \in \mathbb{R}^m$  and  $\varepsilon \in \mathbb{R}^m$  is distributed normally with mean vector 0 and covariance matrix  $\Sigma$ . Cost of effort is a quadratic form with a positive definite matrix,  $C$ .

1. The first-best contract (assuming  $t$  can be contracted upon) is simply  $t = C^{-1}b$ .

2. The second-best *cooperative* contract. The combined return to the principals from effort vector  $t$  is  $b't$ , and so the total expected surplus is

$$b't - \frac{1}{2}t' Ct - \frac{r}{2}\alpha' \Sigma \alpha.$$

This is maximized subject to  $t = C^{-1}\alpha$ . The first-order condition for the slope of the incentive contract,  $\alpha$ , is

$$C^{-1}b - [C^{-1} + r\Sigma]\alpha = 0,$$

or  $b = [I + rC\Sigma]\alpha$  or  $b - \alpha = rC\Sigma\alpha > 0$ .

3. The second-best *non-cooperative* un-restricted contract. Each principal's return is given by the vector  $b^j$ , where  $b = \sum_j b^j$ . Suppose each principal  $j$  is unrestricted in choosing its wage contract; i.e.,  $w^j = \alpha^{j'} + \beta^j$ , where  $\alpha^j$  is a full  $m$ -dimensional vector. Define  $A^{-j} \equiv \sum_{i=j} \alpha^i$  and  $B^{-j} \equiv \sum_{i=j} \beta^i$ . From principal  $j$ 's point of view, absent any contract from himself  $t = C^{-1}A^{-j}$  and the certainty equivalent is  $\frac{1}{2}A^{-j'}[C^{-1} - r\Sigma]A^{-j} + B^{-j}$ . The aggregate incentive scheme facing the agent is  $\alpha = A^{-j} + \alpha^j$  and  $\beta = B^{-j} + \beta^j$ . Thus the agent's certainty equivalent *with* principal  $j$ 's contract is

$$\frac{1}{2}(A^{-j} + \alpha^j)'[C^{-1} - r\Sigma](A^{-j} + \alpha^j) + B^{-j} + \beta^j.$$

The incremental surplus to the agent from the contract is therefore

$$A^{-j'}(C^{-1} - r\Sigma)\alpha^j + \frac{1}{2}\alpha^{j'}[C^{-1} - r\Sigma]\alpha^j + \beta^j.$$

As such, principal  $j$  maximizes

$$b^{j'} C^{-1} A^{-j} - r A^{-j'} \Sigma \alpha^j + b^{j'} C^{-1} \alpha^j - \frac{1}{2} \alpha^{j'} [C^{-1} + r \Sigma] \alpha^j.$$

The first-order condition is

$$C^{-1}b^j - [C^{-1} + r\Sigma]\alpha^j - r\Sigma A^{-j} = 0.$$

Simplifying,  $b^j = [I + rC\Sigma]\alpha^j + rC\Sigma A^{-j}$ . Summing across all principals,

$$b = [I + rC\Sigma]\alpha + rC\Sigma(n-1)\alpha = [I + nrC\Sigma]\alpha,$$

or  $b - \alpha = nrC\Sigma\alpha > 0$ . Thus, the distortion has increased by a factor of  $n$ . Intuitively, it is as if the agent's risk has increased by a factor of  $n$ , and so therefore incentives will be reduced on every margin. *Hence, unrestricted common agency leads to more effort distortions.*

Note that  $b^j = \alpha^j - rC\Sigma\alpha$ , so substitution provides

$$\alpha^j = b^j - rC\Sigma[I + nrC\Sigma]^{-1}b.$$

To get some intuition for the increased-distortion result, suppose that  $n = m$  and that each principal cares only about output  $j$ ; i.e.,  $b_i^j = 0$  for  $i \neq j$ , and  $b_j^j > 0$ . In such a case,

$$\alpha_i^j = -rC\Sigma[I + nrC\Sigma]^{-1}b < 0,$$

so each principal finds it optimal to pay the agent not to produce on the other dimensions!

4. The second-best *non-cooperative restricted* contract. We now consider the case in which each principal is restricted in its contract offerings so as to not pay the agent for output on the other principals' dimensions. Specifically, let's again assume that  $n = m$  and that each principal only cares about  $x_j$ :  $b_i^j = 0$  for  $i \neq j$ , and  $b_j^j > 0$ . The restriction is that  $\alpha_i^j = 0$  for  $j \neq i$ . In such a setting, Dixit demonstrates that the equilibrium incentives are higher than in the un-restricted case. Moreover, if efforts are perfect substitutes across agents,  $\alpha^j = b^j$  and first-best efforts are implemented.

## 2.2 Dynamic Principal-Agent Moral Hazard Models

There are at least three sets of interesting questions which emerge when one turns attention to dynamic settings.

1. Can efficiency be improved in a long-run relationship with full commitment long-term contract?
2. When can short-term contracts perform as well as long-term contracts?
3. What is the effect of renegotiation (i.e., lack of commitment) between the time when the agent takes an action and the uncertainty of nature is revealed?

We consider each issue in turn.

### 2.2.1 Efficiency and Long-Run Relationships

We first focus on the situation in which the principal can commit to a long-run contractual relationship. Consider a simple model of repeated moral hazard where the agent takes an action,  $a_t$ , in each period, and the principal pays the agent a wage,  $w_t(x^t)$ , based upon the history of outputs,  $x^t \equiv \{x_1, \dots, x_t\}$ .

There are two standard ways in which the agent's intertemporal preferences are modeled. First, there is time-averaging.

$$U = \frac{1}{T} \sum_{t=1}^T (u(w_t) - \psi(a_t)).$$

Alternatively, there is the discounted representation.

$$U = (1 - \delta) \sum_{t=1}^T \delta^{t-1} (u(w_t) - \psi(a_t)).$$

Under both models, it has been shown that repeated moral hazard relationships will achieve the first-best arbitrarily close as either  $T \rightarrow \infty$  (in the time-averaging case) or  $\delta \rightarrow 1$  (in the discounting case).

Radner [1985] shows the first-best can be approximated as  $T$  becomes large using the weak law of large numbers. Effectively, as the time horizon grows large, the principal observes the realized distribution and can punish the agent severally enough for discrepancies to prevent shirking.

Fudenberg, Holmst in, and Milgrom [1990] and Fudenberg, Levine, and Maskin [1994], and others have shown that in the discounted case, as  $\delta$  approaches 1, the first best is closely approximated. The intuition for their approach is that when the agent can save in the capital market, the agent is willing to become residual claimant for the firm and will smooth income across time. Although this result uses the agent’s ability to use the capital market in its proof, the first best can be approximately achieved with short-term contracts (i.e.,  $w_t(x_t)$  and not  $w_t(x^t)$ ). See remarks below.

Abreu, Milgrom and Pearce [1991] study the repeated partnership problem in which agents play trigger strategies as a function of past output. They produce some standard folk-theorem results with a few interesting findings. Among other things, they show that taking the limit as  $r$  goes to zero is not the same as taking the limit as the length of each period goes to zero, as the latter has a negative information effect. They also show that there is a fundamental difference between information that is good news (i.e., sales, etc.) and information that is bad news (i.e., accidents). Providing information is generated by a Poisson process, a series of unfavorable events is much more informative about shirking when news is “the number of bad outcomes” (e.g., a high number of failures) than when news is “the number of good outcomes” (e.g., a low number of successes). This latter result has to do with the likelihood function generated from a Poisson process. It is not clear that it generalizes.

### 2.2.2 Short-term versus Long-term Contracts

Although we have seen that when a relationship is repeated infinitely and the discount factor is close to 1 that we can achieve the first-best, what is the structure of long-term contracts when asymptotic efficiency is not attainable? Do short-run contracts perform worse than long-run contracts?

Lambert [1983] and Rogerson [1985] consider the setting in which a principal can freely utilize capital markets at the interest rate of  $r$ , *but agents have no access*. This is fundamental as we will see. In this situation, long-term contracts play a role. Just as cross-state insurance is sacrificed for incentives, intertemporal insurance is also less than first-best efficient. Wage contracts have memory (i.e., today’s wage schedule depends upon yesterday’s output) in order to reduce the incentive problem of the agent. The agent is generally left with a desire to use the credit markets so as to self-insure across time.

We follow Rogerson’s [1985] model here. Let there be two periods,  $T = 2$ , and a finite set of outcomes,  $\{x_1, \dots, x_N\} \equiv \mathcal{X}$ , each of which have a probability  $f(x_i, a) > 0$  of occurring during a period in which action  $a$  was taken. The principal offers a long-term contract of the form  $w \equiv \{w_1(x_i), w_2(x_i, x_j)\}$ , where the wage subscript denotes the period in which the wage schedule is in affect and the output subscripts denote the realized output. (The first argument of  $w_2$  is the first period output; the second argument is the second period output.) An agent either accepts or rejects this contract for the duration of the relationship. If accepted, the agent chooses an action in period 1,  $a_1$ . After observing the period 1 outcome, the agent takes the period 2 action,  $a_2(x_i)$ , that is optimal given the wage schedule in operation. Let  $a \equiv \{a_1, a_2(\cdot)\}$  denote a temporally optimal strategy by the agent for a given wage structure,  $w$ . The principal and the agent both discount the future at rate  $\delta = \frac{1}{1+r}$ . (Identical discount factors is not important for the memory result).

The agent's utility therefore is

$$U = u(w_1) - \psi(a_1) + \delta[u(w_2) - \psi(a_2)].$$

The principal maximizes

$$\sum_{i=1}^N f(x_i, a_1) \left( [x_i - w_1(x_i)] + \delta \sum_{j=1}^N f(x_j, a_2(x_i)) [x_j - w_2(x_i, x_j)] \right).$$

We have two theorems.

**Theorem 16** If  $(a, w)$  is optimal, then  $w$  must satisfy

$$\frac{1}{u'(w_1(x_j))} = \sum_{k=1}^N \frac{f(x_k, a_2(x_j))}{u'(w_2(x_j, x_k))},$$

for every  $j \in \{1, \dots, N\}$ .

**Proof:** We use a variation argument constructing a new contract  $w^*$ . Take the previous contract,  $w$  and change it along the  $x_j$  contingent as follows:

$$\begin{aligned} w_1^*(x_i) &= w_1(x_i) \text{ for } i = j, \\ w_2^*(x_i, x_k) &= w_2(x_i, x_k) \text{ for } i = j, k \in \{1, \dots, N\}, \end{aligned}$$

but

$$\begin{aligned} u(w_1^*(x_j)) &= u(w_1(x_j)) - \Delta, \\ u(w_2^*(x_j, x_k)) &= u(w_2(x_j, x_k)) + \frac{\Delta}{\delta}, \text{ for } k \in \{1, \dots, N\}. \end{aligned}$$

Notice that  $w^*$  only differs from  $w$  following the first-period  $x_j$  branch. By construction, the optimal strategy  $a$  is still optimal under  $w^*$ . To see this note that nothing changes following  $x_i$ ,  $i = j$  in the first period. When  $x_j$  occurs in the first period, the *relative* second-period wages are unchanged, so  $a_2(x_j)$  is still optimal. Finally, the expected present value of the  $x_j$  branch also remains unchanged, so  $a_1$  is still optimal.

Because the agent's expected present value is identical under both  $w$  and  $w^*$ , a necessary condition for the principal's optimal contract is that  $w$  minimizes the expected wage bill over the set of perturbed contracts. Thus,  $\Delta = 0$  must solve the variation program

$$\min_{\Delta} u^{-1}(u(w_1(x_j)) + \Delta) + \delta \sum_{k=1}^N f(x_k, a_2(x_j)) u^{-1} \left( u(w_2(x_j, x_k)) - \frac{\Delta}{\delta} \right).$$

The necessary condition for this provides the condition in the theorem.  $\square$

The above theorem provides a Borch-like condition. It says that the marginal rates of substitution between the principal and the agent should be equal across time *in expectation*. This is not the same as full insurance because of the expectation component. With the theorem above, we can easily prove that long-term contracts will optimally depend upon previous output levels. That is, contracts have memory.

**Theorem 17** If  $w_1(x_i) = w_1(x_j)$  and if the optimal second period effort conditional on period 1 output is unique, then there exists a  $k \in \{1, \dots, N\}$  such that  $w_2(x_i, x_k) = w_2(x_j, x_k)$ .

**Proof:** Suppose not. Let  $w$  have  $w_2(x_i, x_k) = w_2(x_j, x_k)$  for all  $k$ . Then the agent has as an optimal strategy  $a_2(x_i) = a_2(x_j)$ , which implies that  $f(x_k, a_2(x_i)) = f(x_k, a_2(x_j))$  for every  $k$ . But this violates the condition in theorem 16.  $\square$

**Remarks:**

1. Another way to understand Rogerson's result is to consider a rather loose approach using a Lagrangian and continuous outputs (this is loose because we will not concern ourselves with second-order conditions and the like):

$$\max_{w_1(x_1), w_2(x_1, x_2)} \int_{\underline{x}}^{\bar{x}} (x_1 - w_1(x_1))f(x_1, a_1)dx_1 + \int_{\underline{x}}^{\bar{x}} \int_{\underline{x}}^{\bar{x}} (x_2 - w_2(x_1, x_2))f(x_1, a_1)f(x_2, a_2)dx_1dx_2,$$

subject to

$$\int_{\underline{x}}^{\bar{x}} \int_{\underline{x}}^{\bar{x}} (u(w_1(x_1)) + \delta u(w_2(x_1, x_2)))f_a(x_1, a_1)f(x_2, a_2)dx_1dx_2 - \psi'(a_1) = 0,$$

$$\int_{\underline{x}}^{\bar{x}} u(w_2(x_1, x_2))f_a(x_1, a_2)dx_2 - \psi'(a_2) = 0,$$

$$\int_{\underline{x}}^{\bar{x}} u(w_1(x_1))f(x_1, a_1)dx_1 - \psi(a_1) + \int_{\underline{x}}^{\bar{x}} \int_{\underline{x}}^{\bar{x}} \delta u(w_2(x_1, x_2))f(x_1, a_1)f(x_2, a_2)dx_1dx_2 - \psi(a_2) \geq 0.$$

Let  $\mu_1$ ,  $\mu_2(x_1)$  and  $\lambda$  represent the multipliers associated with each constraint (note that the second constraint – the IC constraint for period 2 – depends upon  $x_1$ , and so there is a separate constraint and multiplier for each  $x_1$ ). Differentiating and simplifying one obtains

$$\frac{1}{u'(w_1(x_1))} = \lambda + \mu_1 \frac{f_a(x_1, a_1)}{f(x_1, a_1)}, \quad \forall x_1$$

$$\frac{1}{u'(w_2(x_1, x_2))} = \lambda + \mu_1 \frac{f_a(x_1, a_1)}{f(x_1, a_1)} + \mu_2(x_1) \frac{f_a(x_2, a_2)}{f(x_2, a_2)}, \quad \forall (x_1, x_2).$$

Combining these expressions, we have

$$\frac{1}{u'(w_2(x_1, x_2))} = \frac{1}{u'(w_1(x_1))} + \mu_2(x_1) \frac{f_a(x_2, a_2)}{f(x_2, a_2)}.$$

Because  $\int_{\underline{x}}^{\bar{x}} f_a(x_2, a_2)dx_2 = 0$ , the expectation of  $1/u'(w_2)$  is simply  $1/u'(w_1)$ . Hence,  $\mu_2(x)f_a(x_2, a_2)/f(x_2, a_2)$  represents the deviation of date 2 marginal utility from the first period as a function of both periods' output.

2. Fudenberg, Holmst in and Milgrom [1990] demonstrate the importance of the agent’s credit market restriction. They show that if all public information can be used in contracting and recontracting takes place with common knowledge about technology and preferences, then agent’s perfect access to credit markets results in short-term contracts being as equally effective as long-term contracts. (An additional technical assumption is needed regarding the slope of the expected utility frontier of IC contracts; this is true, for example, when preferences are additively separable over time and utility is unbounded below.)

Together with the folk-theorem results of Radner and others, this consequently implies that short term contracts in which agents have access to credit markets in a repeated setting can obtain the first best arbitrarily closely as  $\delta \rightarrow 1$ . F-H-M [1990] present a construction of such such-term contracts (their theorem 6). Essentially, the agent self-insures by saving in the initial periods and then smoothing income over time.

3. Malcomson and Spinnewyn [1988] show similar results to FHM [1990] where long-term contracts can be duplicated by a sequence of loan contracts.
4. Rey and Salanie [1990] consider three contracts of varying lengths: one period contracts (they call these spot contracts), two-period overlapping contracts (they call these short-term contracts), and multi-period contracts (i.e., long-term contracts). They show that for many contracting situations (including Rogerson’s [1985] model), a sequence of two-period contracts that are renegotiated each period can mimic a long-term contract. Thus, even without capital market access, long-term contracts are not necessary if two-period contracts can be written and renegotiated period by period. The intuition is that a two-period contract mimics a loan/savings contract with the principal.

### 2.2.3 Renegotiation of Risk-sharing

#### Interim Renegotiation with Asymmetric Information:

Fudenberg and Tirole [1990] and Ma [1991] consider the case of moral hazard contracts where the principal has the opportunity to (i.e., the inability to commit not to) offer the agent a new Pareto-improving contract in the interim stage: *after the agent has supplied effort but before the outcome of the stochastic process is revealed*. Their first result is clear.

**Theorem 18** *Choosing any effort level other than the lowest cannot be a pure-strategy equilibrium for the agent in any PBE in which renegotiation is allowed by the principal.*

The proof is straightforward. If the theorem were not true, at the interim stage the principal and agent will be symmetrically informed (on the equilibrium path) and so the principal will offer full insurance. But the agent will intuit that the incentive contract will be renegotiated to a full-insurance contract, and so will supply only the lowest possible effort. As a consequence, high effort chosen with certainty cannot be implemented at any cost. Given the result, the authors naturally turn to mixed-strategy equilibria to study the optimal renegotiation-constrained contract.



We will sketch Fudenberg and Tirole's [1990] analysis of the optimal mixed-strategy renegotiation proof contract. Their insight was that in a mixed-strategy equilibrium (where the agent chooses a distribution over the possible effort choices) a moral hazard setting at the ex ante stage is converted to an adverse selection setting at the interim stage in which an agent's type is chosen action. They show that it is without loss of generality for the principal to offer the agent a renegotiation-proof contract which specifies an optimal mixed strategy for the agent to follow at the ex ante stage and an a menu of contracts for the agent to choose from at the interim stage such that the principal will not wish to renegotiate the contract. Because the analysis which follows necessarily relies upon some knowledge of screening contracts (which is covered in detail in Chapter 2), the reader unfamiliar with these techniques may wish to read the first few pages of chapter 2 (up through section 2.2.1).

Consider the following setting. There are two actions, high and low. To make things interesting, we suppose that the principal wishes to implement to high effort action. The relative cost of supplying high effort to low effort is  $\psi$  and the agent is risk averse in wages. That is,

$$U(w, H) = u(w) - \psi,$$

$$U(w, L) = u(w).$$

Following Grossman and Hart [1983], let  $h(\cdot)$  be the inverse of  $u$ . Thus,  $h(u(w)) = w$ , where  $h$  is strictly increasing and convex.

A high effort generates a distribution  $\{p, 1 - p\}$  on the profit levels  $\{\bar{x}, \underline{x}\}$ ; a low effort generates a distribution  $\{q, 1 - q\}$  on  $\{\bar{x}, \underline{x}\}$ , where  $\bar{x} > \underline{x}$  and  $p > q$ . The principal is risk neutral. Let  $\mu$  be the probability that the agent chooses the high action at the ex ante stage. An incentive contract provides wages as a function of the outcome and the agent's reported type at the interim stage:  $w \equiv \{(\bar{w}_H, \underline{w}_H), (\bar{w}_L, \underline{w}_L)\}$  Then,

$$V(w, \mu) = \mu[p\bar{w}_H + (1 - p)\underline{w}_H] + (1 - \mu)[q\bar{w}_L + (1 - q)\underline{w}_L].$$

The procedure we follow to solve for the optimal renegotiation-proof contract uses backward induction. Begin at the interim stage where the principal's beliefs are some arbitrary  $\mu$  and the agent's expected utility in absence of renegotiation is  $\{\bar{U}, \underline{U}\}$ . We solve for the optimal renegotiation contract  $w$  as a function of  $\mu$ . Then we consider the ex ante stage and maximize ex ante profits over the set of  $(\mu, w)$  pairs which are generate an interim optimal contract.

**Optimal Interim Contracts:** In the interim period, following standard revealed-preference tricks, one can show that the incentive compatibility constraint for the low-effort agent and the individual rationality constraint for the high-effort agent will be binding while the other constraints will slack. Let  $\bar{u}_H \equiv u(\bar{w}_H)$ ,  $\underline{u}_H \equiv u(\underline{w}_H)$ ,  $\bar{u}_L \equiv u(\bar{w}_L)$ , and  $\underline{u}_L \equiv u(\underline{w}_L)$ . Then the principal solves the following convex programming problem:

$$\max_w -\mu[p\bar{u}_H + (1 - p)\underline{u}_H] - (1 - \mu)[q\bar{u}_L + (1 - q)\underline{u}_L],$$

subject to

$$q\bar{u}_L + (1 - q)\underline{u}_L \geq q\bar{u}_H + (1 - q)\underline{u}_H,$$

$$p\bar{u}_H + (1-p)\underline{u}_H \geq \bar{U}.$$

Let  $\gamma$  be the IC multiplier and  $\lambda$  be the IR multiplier. Then by the Kuhn-Tucker theorem, we have a solution which satisfies the following four necessary first-order conditions.

$$\begin{aligned} \mu p h'(\bar{u}_H) &= p\lambda - q\gamma, \\ \mu(1-p)h'(\underline{u}_H) &= (1-p)\lambda - (1-q)\gamma, \\ (1-\mu)h'(\bar{u}_L) &= \gamma, \\ (1-\mu)h'(\underline{u}_L) &= \gamma. \end{aligned}$$

Combining the last two equations implies  $\bar{u}_L = \underline{u}_L = u_L$  (i.e., complete insurance for the low-effort agent), and therefore  $\gamma = (1-\mu)h'(u_L)$ . Using this result for  $\gamma$  in the first two equations, and substituting out  $\lambda$ , yields

$$\frac{\mu}{1-\mu} = \frac{h'(u_L)}{h'(\bar{u}_H) - h'(\underline{u}_H)} \frac{p-q}{p(1-p)}.$$

This equation is usually referred to as the renegotiation-proofness constraint. For any given wage schedule  $w \equiv \{(\bar{w}_H, \underline{w}_H), (\bar{w}_L, \underline{w}_L)\}$  (or alternatively, a utility schedule  $u \equiv \{(\bar{u}_H, \underline{u}_H), (\bar{u}_L, \underline{u}_L)\}$ ), there exists a unique  $\mu^*(u)$  which satisfies the above equation, and which provides the upper bound on feasible ex ante effort choices (i.e.,  $\forall \mu \leq \mu^*$ , the ex ante contract is renegotiation proof).

**Optimal Ex ante Contracts:** Now consider the optimal ex ante contract offer  $\{\mu, u\}$ . The principal solves

$$\begin{aligned} \max_{\mu, (\bar{u}_H, \underline{u}_H, u_L)} \quad & \mu[p(\bar{x} - h(\bar{u}_H)) + (1-p)(\underline{x} - h(\underline{u}_H))] \\ & + (1-\mu)[q(\bar{x} - h(u_L)) + (1-q)(\underline{x} - h(u_L))], \end{aligned}$$

subject to

$$\begin{aligned} u_L &= p\bar{u}_H + (1-p)\underline{u}_H - \psi, \\ u_L &\geq 0, \\ \frac{\mu}{1-\mu} &= \frac{h'(u_L)}{h'(\bar{u}_H) - h'(\underline{u}_H)} \frac{p-q}{p(1-p)}. \end{aligned}$$

The first constraint is the typical binding IC constraint that the high effort agent is just willing to supply high effort. More importantly, indifference also guarantees that the agent is willing to randomize according to  $\mu$ , and so it is a *mixing* constraint as well. The second constraint is the IR constraint for the low effort choice (and by the first equation, for the high-effort choice as well). The third equation is our renegotiation-proofness (RP) constraint. Note that if the principal attempts to choose  $\mu$  arbitrarily close to 1, then the RP constraint implies that  $\bar{u}_H \approx \underline{u}_H = u_H$ . Interim IC in turn implies that  $u_H \approx u_L$ . But this in turn will violate the ex ante IC (mixing) constraint for  $\psi > 0$ . Thus, it is not feasible to choose  $\mu$  arbitrarily close to 1.

**Remarks:**

1. In general, the RP constraint both lessens the principal's expected profit as well as reduces the principal's feasible set of implementable actions. Thus, non-commitment has two effects.
2. In general, it is not the case that the ex ante IR constraint for the agent will bind. Specifically, it is possible that by leaving the agents some rents that the RP constraint will be weakened. However, if utility displays non-increasing absolute risk aversion, the IR constraint will bind.
3. The above results are extended to a continuum of efforts and two outcomes by Fudenberg and Tirole [1990]. They also show that although it is without loss of generality to consider RP contracts, one can show that with equilibrium renegotiation, the principal can uniquely implement the optimal RP contract.
4. Our results about the cost of renegotiation were predicated upon the principal (who is uninformed at the interim stage) making a take-it-or-leave-it offer to the agent. According to Fudenberg and Tirole (who refer to Maskin and Tirole's [1992] paper on "The Principal-Agent Relationship with an Informed Principal, II: Common values"), not only is principal-led renegotiation simpler, but the same conclusions are obtained if the agent leads the renegotiation providing one requires that the original contract be "strongly renegotiation proof" (i.e., RP in *any* PBE at the interim stage). Nonetheless, in a paper by Ma [1994] it is shown that providing one is prepared to use a refinement on the principal's beliefs regarding off-the-equilibrium-path beliefs at the interim stage, agent-led renegotiation has no cost. That is, the second-best incentive contract remains renegotiation proof. Of course, according to Maskin and Tirole [1992], such a contract cannot be strongly renegotiation proof; i.e., there must be other equilibria where renegotiation does occur and is costly to the principal at the ex ante stage.
5. Matthews [1995].

### **Interim Renegotiation with Symmetric Information:**

The previous analysis assumed that at the interim (renegotiation) stage, one party was asymmetrically informed. Hermalin and Katz [1991] demonstrate that this is the source of *costly* renegotiation. With symmetrically informed parties, full insurance can be provided and first best effort can be implemented. Hence, it is possible that "renegotiation" can improve welfare, because it can change the terms of a contract to reflect new information about the agent's performance.

Consider the following variation on standard principal-agent timing. A contract is offered by the principal to the agent at the ex ante stage, and the agent immediately takes an action. Before the interim (renegotiation), however, both parties observe a signal  $s$  regarding the agent's chosen action. This signal is observable, but not verifiable (i.e., it cannot be part of any enforceable contract). The parties can now renegotiate. Following renegotiation, the verifiable output  $x$  is observed, and the agent is rewarded according to the contract in force and the realized output (which is contractible).

For simplicity, assume that both parties observe the chosen action,  $a$ , at the renegotiation stage. Then almost trivially it follows that the principal cannot be made worse off with renegotiation when the principal makes the interim offers. The set of implementable contracts remains unchanged. Generally, however, the principal can do better by using renegotiation to her benefit.

The the following proposition for agent-led renegotiation follows immediately.

**Theorem 19** Suppose that  $s$  perfectly reveals  $a$  and the agent makes a take-it-or-leave-it renegotiation offer at the interim stage. Then the first-best action is implementable at the full information cost.

**Proof:** By construction. The principal sells the firm to the agent at the ex ante stage for a price equal to the firm's first-best expected profit (net of effort costs). The agent is willing to purchase the firm, exert the first-best level of effort to maximize its resale value, and then sell the firm back to the principal at the interim stage, making a profit of zero. This satisfies IR and IC, and all rents go to the principal.  $\square$

A similar result is true if we give all of the bargaining power to the principal. In this case, the principal will offer the agent his certainty equivalent of the ex ante contract at the interim stage. Assume that  $a^{FB}$  is implementable with certainty equivalent equal to the agent's reservation utility. (This will be possible, for example, if the distribution vector  $p(a)$  is not an element of the convex hull of  $\{p(a')|a' = a\}$  for any  $a$ .) A principal would not normally want to do this because the risk premium the principal must give to the agent to satisfy IR is too great. But with the interim stage of renegotiation, the principal observes  $a$  and can renegotiate away all of the risk. Thus we have ...

**Theorem 20** Suppose that  $a^{FB}$  is implementable and the principal makes a take-it-or-leave-it offer at the renegotiation stage. Then the first-best is implementable at the full-information cost.

Hermalin and Katz go even further to show that under some mild technical assumptions that the first-best full information allocation is obtainable with any arbitrary bargaining game at the interim stage.

**Remarks:**

1. Hermalin and Katz extend their results to the case where  $a$  is only imperfectly observable and find that principal-led renegotiation is still beneficial providing that the commonly observable signal  $s$  is a sufficient statistic for  $x$  with respect to  $(s, a)$ . That is, having observed  $s$ , the principal can predict  $x$  as well as the agent can. Although the first-best cannot generally be achieved, renegotiation is beneficial.
2. Juxtaposing Hermalin and Katz's results with those of Fudenberg and Tirole's [1990], we find some interesting connections. In H&K, the terms of trade at the interim stage are a direct function of observed effort,  $a$ ; in F&T, the dependence is obtained only through costly interim incentive compatibility conditions. Renegotiation can be bad because it undermines commitment in absence of information on  $a$ ; it can be good if

renegotiation can be made conditional on  $a$ . Thus, the main differences are whether renegotiation takes place between asymmetrically or (sufficiently) symmetrically informed parties.

### 3 Mechanism Design and Self-selection Contracts

#### 3.1 Mechanism Design and the Revelation Principle

We consider a setting where the principal can offer a mechanism (e.g., contract, game, etc.) which her agents can play. The agent's are assumed to have private information about their preferences. Specifically, consider  $I$  agents indexed by  $i \in \{1, \dots, I\}$ .

- Each agent  $i$  observes only its own preference parameter,  $\theta_i \in \Theta_i$ . Let  $\theta \equiv (\theta_1, \dots, \theta_I) \in \Theta \equiv \prod_{i=1}^I \Theta_i$ .
- Let  $y \in Y$  be an allocation. For example, we might have  $y \equiv (x, t)$ , with  $x \equiv (x_1, \dots, x_I)$  and  $t \equiv (t_1, \dots, t_I)$ , and where  $x_i$  is agent  $i$ 's consumption choice and  $t_i$  is the agent's payment to the principal. The choice of  $y$  is generally controlled by the principal, although she may commit to a particular set of rules.
- Utility for  $i$  is given by  $U_i(y, \theta)$ ; note general interdependence of utilities on  $\theta_{-i}$  and  $y_{-i}$ . The principal's utility is given by the function  $V(y, \theta)$ . In a slight abuse of notation, if  $y$  is a distribution of outcomes, then we'll let  $U_i$  and  $V$  represent the value of expected utility after integrating with respect to the distribution.
- Let  $p(\theta_{-i}|\theta_i)$  be  $i$ 's probability assessment over the possible types of other agents given his type is  $\theta_i$  and let  $p(\theta)$  be the common prior on possible types.

Suppose that the principal has all of the bargaining power and can commit to playing a particular game or mechanism involving her agent(s). Posed as a mechanism design question, the principal will want to choose the game (from the set of all possible games) which has the best equilibrium (to be defined) for the principal. But this set of all possible games is enormous and complex. The revelation principle, due to Green and Laffont [1977], Myerson [1979], Harris and Townsend [1981], Dasgupta, Hammond, and Maskin [1979], et al., allows us to simplify the problem dramatically.

**Definition:** A communication *mechanism* or game,

$$\Gamma^c \equiv \{\mathcal{M}, \Theta, p, U_i(y(m), \theta)_{i=1, \dots, I}\},$$

is characterized by a message (i.e., strategy) space for each agent,  $\mathcal{M}_i$ , and an allocation  $y$  for each possible message profile,  $m \equiv (m_1, \dots, m_I) \in \mathcal{M} \equiv (\mathcal{M}_1, \dots, \mathcal{M}_I)$ ; i.e.,  $y : \mathcal{M} \mapsto Y$ . For generality, we will suppose that  $\mathcal{M}_i$  includes all possible mixtures over messages; thus,  $m_i$  may be a probability distribution. When no confusion would result, we sometimes indicate a mechanism  $\Gamma$  by the pair  $\{\mathcal{M}, y\}$ .

The timing of the communication mechanism game is as follows:

MIT OpenCourseWare  
<https://ocw.mit.edu>

14.124 Microeconomic Theory IV  
Spring 2017

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.