

Sender-Receiver Games with Cheap Talk

Let us consider games where player 1 (the sender) has private information about his type in some set T_1 , and player 2 (the receiver) chooses an action in some set C_2 . Each player i has a given utility functions $u_i: C_2 \times T_1 \rightarrow \mathbb{R}$, and player 2's prior belief about 1's type is described by a given probability distribution $p \in \Delta(T_1)$. We assume that player 1 can send messages (with a large potential vocabulary) to player 2 before she chooses her action, but player 1's messages are just "cheap talk" which do not affect anybody's payoff except to the extent that player 2 responds to the messages sent by player 1.

With cheap talk, there is always a babbling equilibrium where 1's message is independent of his type and 2's action is her ex-ante optimal action independently of 1's message. But sometimes better equilibria can be found where substantive communication occurs.

For example, consider a game (from section 6.7 in Game Theory by Myerson) where $T_1 = \{1a, 1b\}$, $C_2 = \{x, y, z\}$, the prior probability distribution has $p(1a) = p(1b) = 0.5$, and the utility payoffs (u_1, u_2) depend on 2's action and 1's type as follows:

	$c_2 = x$	$c_2 = y$	$c_2 = z$	
$t_1 = 1a$	2, 3	0, 2	-1, 0	$(p(1a) = 0.5)$
$t_1 = 1b$	1, 0	2, 2	0, 3	$(p(1b) = 0.5)$

In the babbling equilibrium, player 2's optimal choice with prior beliefs is y , which would be the best outcome for player 1 when his type is $1b$. When 1's type is $1a$, however, he would like to tell player 2 that x is actually better for both of them. But in direct face-to-face communication without noise, there cannot be any equilibrium in which player 1 encourages player 2 to choose x with positive probability. Notice first that no belief about t_1 could ever make 2 willing to randomize between x and z . (Action x is optimal for 2 when the probability of $1a$ is $2/3$ or more, y is optimal for 2 when the probability of $1a$ is between $1/3$ and $2/3$, and z is optimal for 2 when the probability of $1a$ is $1/3$ or less.) If there were any message that player 1 could announce that would make player 2 willing to choose x (either for sure or in a randomization between x and y), then type $1a$ would always want to announce such a message (to maximize the probability of x), but then the absence of such a message would convince player 2 that 1's type is $1b$ and so would cause her to choose z ; but this in turn would imply that even type $1b$ should want to send the same message that type $1a$ would send. Thus, with direct communication, this game has no equilibrium other than the trivial babbling equilibrium.

But player 1 can send a credible message with noise. Imagine that player 1 has a carrier pigeon which, if sent, would reach player 2 with probability 0.4. There is an equilibrium in which 1 sends the carrier pigeon (with a note saying "I am type $1a$, please do x ") if $t_1 = 1a$ but not if $t_1 = 1b$. If the pigeon does not arrive, player 2's posterior belief about the probability of type $1b$ is $0.5 \times 1 / (0.5 \times 1 + 0.5 \times (1 - 0.4)) = 5/8$, and so player 2 is still prefer to choose y , not z . So noise can help player 1 here to send messages that credibly guide 2's action.

The discovery that noise can actually help to support credible communication that might not be possible without noise raises the question: What else is possible? We now develop a general framework for answering this question, for any given sender-receiver game.

Given a sender-receiver game (T_1, C_2, u_1, u_2, p) , a general coordination mechanism or mediation plan can be described by a function $\mu: T_1 \rightarrow \Delta(C_2)$. We can think of this plan as being implemented by a trustworthy mediator who first asks player 1 to confidentially report his type, and then, depending on this report, the mediator will recommend an action to player 2. For any $t_1 \in T_1$ and $c_2 \in C_2$, and The number $\mu(c_2 | t_1)$ denotes the probability that the mediator will recommend action c_2 if player 1 reports type t_1 . So

$$[1] \quad \sum_{d_2 \in C_2} \mu(d_2|t_1) = 1 \quad \text{and} \quad \mu(c_2|t_1) \geq 0, \quad \forall c_2 \in C_2, \quad \forall t_1 \in T_1.$$

The plan of sending a carrier pigeon corresponds to the mediation plan with $\mu(c_2|t_1)$ as follows:

	$c_2=x$	$c_2=y$	$c_2=z$
$t_1=1a$	0.4	0.6	0
$t_1=1b$	0	1	0

Consider the expected payoffs that will result from a plan μ if player 1 reports his type honestly and player 2 chooses her action obediently under this plan μ . Player 2's expected payoff is

$$U_2(\mu) = \sum_{t_1 \in T_1} p(t_1) \sum_{c_2 \in C_2} \mu(c_2|t_1) u_2(c_2, t_1).$$

If player 1's type is t_1 then his expected payoff is

$$U_1(\mu|t_1) = \sum_{d_2 \in C_2} \mu(d_2|t_1) u_1(c_2, t_1).$$

But if player 1 with type t_1 dishonestly reported type s_1 to the mediator, then he could get expected payoff

$$\hat{U}_1(\mu, s_1|t_1) = \sum_{d_2 \in C_2} \mu(d_2|s_1) u_1(c_2, t_1).$$

On the other hand, if player 2 planned to choose action d_2 when a particular action c_2 is recommended, then the net change in player 2's expected payoff would be

$$\sum_{t_1 \in T_1} p(t_1) \mu(c_2|t_1) (u_2(d_2, t_1) - u_2(c_2, t_1)).$$

Then, honest reporting by player 1 and obedient action by player 2 can be equilibrium behavior with mediation plan μ if and only if

$$[2] \quad U_1(\mu|t_1) \geq \hat{U}_1(\mu, s_1|t_1), \quad \forall t_1 \in T_1, \quad \forall s_1 \in T_1; \quad \text{and}$$

$$[3] \quad \sum_{t_1 \in T_1} p(t_1) \mu(c_2|t_1) (u_2(d_2, t_1) - u_2(c_2, t_1)) \leq 0, \quad \forall c_2 \in C_2, \quad \forall d_2 \in C_2.$$

We say that μ is incentive compatible iff μ satisfies these constraints [2] and [3]. Here [2] are the informational incentive constraints saying that player 1 should not want to lie about his type, and [3] are the strategic incentive constraints saying that player 2 should not want to disobey her recommendations.

A mechanism μ is (strictly) interim dominated by another mechanism ν iff every type of player 1 would expect to better under ν than μ , and player 2 would also expect to do better under ν than μ ; that is,

$$U_1(\nu|t_1) > U_1(\mu|t_1) \quad \forall t_1 \in T_1, \quad \text{and} \quad U_2(\nu) > U_2(\mu).$$

A mechanism μ is (weakly) incentive-efficient iff it is incentive compatible and it is not interim dominated by any other incentive-compatible mechanism. That is, μ is not incentive-efficient if a social planner can find some other incentive-compatible mechanism ν such that, at the point in time when player 1 knows his type but player 2 does not, we could be sure that both players 1 and 2 would definitely prefer to implement ν rather than μ . A mechanism μ is incentive-efficient if and only if there exist some nonnegative utility weights $\lambda_1(t_1) \geq 0$ for each t_1 in T_1 and $\lambda_2 \geq 0$ such that these weights are not all zero and μ is an optimal solution to the problem of maximizing

$$[4] \quad \lambda_2 U_2(\mu) + \sum_{t_1 \in T_1} \lambda_1(t_1) U_1(\mu|t_1)$$

over μ subject to the probability constraints [1] and the incentive constraints [2] and [3].

To characterize optimal solutions to this maximization problem, we consider the Lagrangean:

$$L(\mu, \lambda, \alpha) = \lambda_2 U_2(\mu) + \sum_{t_1 \in T_1} \lambda_1(t_1) U_1(\mu|t_1) + \\ + \sum_{t_1 \in T_1} \sum_{s_1 \in T_1} \alpha_1(s_1|t_1) (U_1(\mu|t_1) - \hat{U}_1(\mu, s_1|t_1)) + \\ + \sum_{c_2 \in C_2} \sum_{d_2 \in C_2} \alpha_2(d_2|c_2) \sum_{t_1 \in T_1} p(t_1) \mu(c_2|t_1) (u_2(c_2, t_1) - u_2(d_2, t_1)).$$

Now let us define (λ, α) -virtual utility functions v_1 and v_2 as follows:

$$v_1(c_2, t_1, \lambda, \alpha) = [(\lambda_1(t_1) + \sum_{s_1 \in T_1} \alpha_1(s_1|t_1)) u_1(c_2, t_1) - \sum_{s_1 \in T_1} \alpha_1(t_1|s_1) u_1(c_2, s_1)] / p(t_1),$$

$$v_2(c_2, t_1, \lambda, \alpha) = \lambda_2(t_1) u_2(c_2, t_1) + \sum_{d_2 \in C_2} \alpha_2(d_2|c_2) (u_2(c_2, t_1) - u_2(d_2, t_1)).$$

These definitions are constructed to give us the following equation

$$L(\mu, \lambda, \alpha) = \sum_{t_1 \in T_1} p(t_1) \mu(c_2|t_1) (v_1(c_2, t_1, \lambda, \alpha) + v_2(c_2, t_1, \lambda, \alpha)).$$

That is, the Lagrangean $L(\mu, \lambda, \alpha)$ is just the expected sum of the players' (λ, α) -virtual utilities.

An incentive-compatible mechanism μ solves the maximization problem [4] if and only if there exists a vector α of nonnegative Lagrange multipliers that satisfy the complementary slackness conditions

$$\alpha_1(s_1|t_1) \geq 0 \text{ and } \alpha_1(s_1|t_1) (U_1(\mu|t_1) - \hat{U}_1(\mu, s_1|t_1)) = 0, \quad \forall t_1 \in T_1, \forall s_1 \in T_1;$$

$$\alpha_2(d_2|c_2) \geq 0 \text{ and } \alpha_2(d_2|c_2) \sum_{t_1 \in T_1} p(t_1) \mu(c_2|t_1) (u_2(c_2, t_1) - u_2(d_2, t_1)) = 0, \quad \forall c_2 \in C_2, \forall d_2 \in C_2;$$

and μ maximizes the the expected sum of the players (λ, α) -virtual utilities over all mechanisms $\mu: T_1 \rightarrow \Delta(C_2)$ without regard to the incentive constraints. This virtual maximization holds when μ satisfies the optimal-support condition:

$$\{c_2 | \mu(c_2|t_1) > 0\} \subseteq \operatorname{argmax}_{c_2 \in C_2} (v_1(c_2, t_1, \lambda, \alpha) + v_2(c_2, t_1, \lambda, \alpha)), \quad \forall t_1 \in T_1.$$

The complementary slackness conditions say that a Lagrange multiplier can be strictly positive only if its corresponding incentive constraint is satisfied by μ as a binding equality. When $\alpha_1(t_1|s_1) > 0$ (which indicates some difficulty in deterring 1 from reporting t_1 when s_1 is true), we may say that type s_1 jeopardizes type t_1 for player 1. With this terminology, player 1's (λ, α) -virtual utility is a positive multiple of his true utility minus positive multiples of the utility that he would get with other possible types that jeopardize his true type. When $\alpha_2(d_2|c_2) > 0$ (which indicates some difficulty in deterring 2 from choosing d_2 when c_2 is recommended), we may say that action d_2 jeopardizes action c_2 for player 2. Player 2's (λ, α) -virtual utility is a positive multiple of her actual utility minus positive multiples of the utility that she would get with other actions that jeopardize her recommended action.

When people enter into a relationship or transaction, the problem of getting them to act appropriately is called moral hazard, and the problem of getting them to share information appropriately is called adverse selection. In our sender-receiver games, the sender (player 1) is subject to an adverse selection problem, represented by the informational incentive constraints [2], and the receiver (player 2) is subject to a moral hazard problem, represented by the strategic incentive constraints [3].

Let us apply these conditions to the question of what incentive-compatible mediation plan would be best for the receiver (player 2) in our example with the payoffs (u_1, u_2) :

	$c_2=x$	$c_2=y$	$c_2=z$	
$t_1=1a$	2, 3	0, 2	-1, 0	$(p(1a)=0.5)$
$t_1=1b$	1, 0	2, 2	0, 3	$(p(1b)=0.5)$

In this case we have $\lambda_2=1$ and $\lambda_1(1a)=\lambda_1(1b)=0$, because we are just trying to maximize $U_2(\mu)$ subject to the probability constraints and incentive constraints. This is a linear programming problem which can be efficiently solved by many computer programs, including the Solver program in MS Excel.

But we can solve the problem by hand if we can guess which probabilities in μ will be positive and which Lagrange multipliers in α will be positive.

In this case the correct guesses are that

$$\mu(x|1a)>0, \mu(y|1a)>0, \mu(y|1b)>0, \mu(z|1b)>0, \alpha_1(1a|1b)>0, \text{ and } \alpha_2(z|y)>0.$$

That is, the mediator randomizes between recommending x or y if player 1 reports 1a, and the mediator randomizes between recommending y or z if player 1 reports 1b; the binding incentive constraints are that player 1 should not report 1a when 1b is true, and that player 2 should not do z when y is recommended. To simplify our notation, let $\beta = \alpha_1(1a|1b)$, $\gamma = \alpha_2(z|y)$, $p = \mu(y|1a)$, and $q = \mu(y|1b)$, so that $\mu(x|1a) = 1-p$ and $\mu(z|1b) = 1-q$. With these Lagrange multipliers, the virtual utilities (v_1, v_2) become

	$c_2=x$	$c_2=y$	$c_2=z$
$t_1=1a$	$-(1)\beta/0.5, 3$	$-(2)\beta/0.5, 2+(2-0)\gamma$	$-(0)\beta/0.5, 0$
$t_1=1b$	$(1)\beta/0.5, 0$	$(2)\beta/0.5, 2+(2-3)\gamma$	$(0)\beta/0.5, 3$

The optimal-support conditions require: $-(1)\beta/0.5 + 3 = -(2)\beta/0.5 + 2+(2-0)\gamma \geq (0)\beta/0.5 + 0$, and $(1)\beta/0.5 + 0 \leq (2)\beta/0.5 + 2+(2-3)\gamma = (0)\beta/0.5 + 3$.

The two equations are satisfied when $\beta=0.5$ and $\gamma=1$, and then the two inequalities are also satisfied.

The binding incentive constraints ($\alpha_1(1a|1b)>0$ and $\alpha_2(z|y)>0$) require

$$2q+0(1-q) = (1)(1-p)+2p \text{ and } (0.5)p(2)+(0.5)q(2) = (0.5)p(0)+(0.5)q(3),$$

and these equations are satisfied by $p=1/3$, $q=2/3$. So we get a mechanism with $\mu(c_2|t_1)$ as follows:

	$c_2=x$	$c_2=y$	$c_2=z$
$t_1=1a$	$2/3$	$1/3$	0
$t_1=1b$	0	$2/3$	$1/3$

It is straightforward to verify that this mechanism μ satisfies all other incentive constraints (with slack), and so it satisfies all the Lagrangean conditions for maximizing player 2's expected payoff among all incentive-compatible mechanisms.

Motivating an agent with a linear type drawn from a continuous distribution on an interval.

Suppose that the agent's type \tilde{t} is a random variable drawn from an interval $[A,B]$.

The agent's type t is his cost of effort, and his utility for income w and effort q is $w - tq$.

Consider any contract $(w(\bullet),q(\bullet))$ where the terms of trade for each type θ would be $(w(\theta),q(\theta))$.

Let $U(w,q|t) = w(t) - tq(t)$ denote the expected utility of type t under this contract.

For any pair of possible types t and s in $[A,B]$, the $(s|t)$ -informational incentive constraint says

$$U(w,q|t) = w(t) - tq(t) \geq w(s) - tq(s) = U(w,q|s) + (s-t)q(s).$$

Similarly, the $(t|s)$ -incentive constraint implies $U(w,q|s) \geq U(w,q|t) + (t-s)q(t)$

So the $(t|s)$ and $(s|t)$ constraints together imply $(s-t)q(t) \geq U(w,q|t) - U(w,q|s) \geq (s-t)q(s)$.

So when $s > t$ we must have $q(t) \geq q(s)$, and so $q(t)$ is a decreasing function of the cost-type t .

These inequalities over many small steps from t up to B yield the information-rent equation:

$$U(w,q|t) = U(w,q|B) + \int_t^B q(r) dr.$$

With all $q(r) \geq 0$, the high cost-type B has the least gain from trade: $U(w,q|B) = \min_{t \in [A,B]} U(w,q|t)$.

The w function can be determined from the q function and the value $U(w,q|B)$ by:

$$w(t) = U(w,q|t) + tq(t) = U(w,q|B) + tq(t) + \int_t^B q(r) dr$$

With $q(\bullet)$ weakly decreasing, these functions (w,q) will satisfy incentive compatibility because

$$U(w,q|t) - [w(s) - tq(s)] = U(w,q|t) - U(w,q|s) - (s-t)q(s) = \int_t^s q(r) dr - (s-t)q(s) = \int_t^s [q(r) - q(s)] dr \geq 0.$$

Suppose the principal's beliefs about the agent's type are described by the cumulative distribution

$F(t) = P(\tilde{t} \leq t)$, and $f(t) = F'(t)$ is the continuous probability density of this distribution,

with $f(t) > 0$ for all t in $[A,B]$. Here $F(B)=1$, $F(A)=0$, and $P(a \leq \tilde{t} \leq b) = F(b) - F(a) = \int_a^b f(t) dt$

whenever $a \leq b$. Then the expected wage bill is

$$\int_A^B w(t) f(t) dt = \int_A^B [U(w,q|t) + tq(t)] f(t) dt = \int_A^B [U(w,q|B) + tq(t) + \int_t^B q(r) dr] f(t) dt$$

$$= U(w,q|B) + \int_A^B tq(t) f(t) dt + \int_A^B \int_A^r f(t) dt q(r) dr$$

$$= U(w,q|B) + \int_A^B tq(t) f(t) dt + \int_A^B F(r) q(r) dr = U(w,q|B) + \int_A^B q(t) [t + F(t)/f(t)] f(t) dt.$$

So the incentive-compatible expected wage $E(w(\tilde{t}))$ looks like what the principal would have to pay without incentive constraints if the cost of each type t were increased to a virtual cost $t + F(t)/f(t)$.

This virtual-cost formula expresses the fact that, when we ask more effort from any type t , we increase the amount that we must pay all types below t , because of incentive constraints.

Example: Akerlof's Lemons.

The "agent" is the seller of a unique object, of which the "principal" is the only potential buyer. The seller's type is the value of the object to him, which depends on his unverifiable private information about its quality. Then $q(t)$ can be reinterpreted as the probability of his selling the good if he acts like type t , which must satisfy $0 \leq q(t) \leq 1$, and $w(t)$ is his expected revenue from selling if he acts like type t .

Suppose \tilde{t} is drawn from a Uniform distribution on the interval from 0 to 100, but the value of the object to the buyer also depends on the quality (which the buyer would learn only after the transaction) and would be $1.5\tilde{t}$. So the object would always be worth 50% more to the buyer.

If (w,q) satisfies the incentive constraints and $U(w,q|t) \geq 0$, the buyer's expected gain from trade is

$$\int_0^{100} [1.5tq(t) - w(t)] f(t) dt = \int_0^{100} [1.5t - t - F(t)/f(t)] q(t) f(t) dt - U(w,q|100)$$

$$= \int_0^{100} [1.5t - 2t] q(t) dt / 100 - U(w,q|100) \leq 0. \text{ The buyer can only expect to lose if any } q(t) > 0.$$

Equilibrium in markets with adverse selection

Let Y denote the set of possible contracts, and let T denote the set of possible types of consumers. For simplicity, we may sometimes assume that Y and T are nonempty finite sets. (When they are infinite sets, some sums below may need to be rewritten as integrals). A given utility function $U: \mathbb{R} \times Y \times T \rightarrow \mathbb{R}$ specifies the utility $U(p, x, t)$ that any consumer of type t would get from buying contract x at price p . A given cost function $C: Y \times T \rightarrow \mathbb{R}$ specifies the expected cost $C(x, t)$ for a firm to fulfill a contract x for a consumer of type t . A given probability distribution μ in $\Delta(T)$ specifies the fraction $\mu(t) > 0$ of consumers who have each type t in the general population. Each consumer must buy exactly one contract. (We could include the no-trade option as an $x=0$ contract with cost 0 and utility 0 for all types.)

In a Rothschild-Stiglitz model of insurance markets, a consumer's type $t \in [0, 1]$ would denote his probability of suffering some loss $\ell > 0$ from some given initial wealth W , and the contract parameter $x \in [0, 1]$ would denote the fraction of this loss to be covered by an insurance policy. Then, given some concave increasing utility function $u(\bullet)$ for monetary wealth, we would get

$$U(p, x, t) = t u(W - p - (1-x)\ell) + (1-t) u(W - p).$$

In general, we assume here that, for each $x \in Y$ and each $t \in T$, $U(p, x, t)$ is strictly decreasing and continuous in the price p . Also, for each $x \in Y$, $t \in T$, $p \in \mathbb{R}$, and $y \in Y$, we assume that there exists some price $\varphi(p, x, t, y)$ such that $U(p, x, t) = U(\varphi(p, x, t, y), y, t)$. (This says that money is important enough for a price adjustment to change any consumer's preference over any pair of contracts.) So $\varphi(p, x, t, y)$ denotes the price of y which would make a type- t consumer indifferent between buying y and buying x at price p .

In this market, each consumer will buy exactly one contract. We also assume that there are multiple firms which can sell any of these contracts, and any one of these firms could serve the entire population of consumers. Thus, in a competitive equilibrium, prices should be such that firms expect zero profits for every contract.

We consider here a simplified version of Azevedo and Gottlieb's (2017) definition of competitive equilibrium for markets with adverse selection. We define a competitive equilibrium to be a pair (q, γ) such that $q = (q(x))_{x \in Y}$ is a price vector in \mathbb{R}^Y , $\gamma = (\gamma(x|t))_{x \in Y, t \in T}$ is an allocation vector in $\Delta(Y)^T$, and the following conditions are satisfied:

- [0] $\sum_{y \in Y} \gamma(y|t) = 1$ and $\gamma(x|t) \geq 0$, $\forall x \in Y, \forall t \in T$;
- [1] $\sum_{x \in Y} \gamma(x|t) U(q(x), x, t) = \max_{x \in Y} U(q(x), x, t)$, $\forall t \in T$;
- [2] $\sum_{t \in T} \mu(t) \gamma(x|t) (C(x, t) - q(x)) = 0$, $\forall x \in Y$; and
- [3] $\forall y \in Y, \exists t \in T$ such that $U(q(y), y, t) = \max_{x \in Y} U(q(x), x, t)$ and $q(y) \leq C(y, t)$.

Condition [0] must be satisfied because each number $\gamma(x|t)$ denotes the fraction of type- t consumers who choose contract x . Condition [1] is an optimality condition for consumers, saying that consumers of each type only choose contracts that maximize their utility at the given q prices. Condition [2] is a zero-profit condition for firms, saying that firms expect to break even on any contract that attracts a positive fraction of the consumers. In this case, when $\sum_{s \in T} \mu(s) \gamma(x|s) > 0$, condition [2] implies that $q(x)$ equals the average cost of all consumers who choose contract x

$$q(x) = \sum_{t \in T} C(x, t) \mu(t) \gamma(x|t) / \sum_{s \in T} \mu(s) \gamma(x|s).$$

Condition [3] says that, for each contract, there is at least one type which is willing to choose this contract but would not be profitable for the competitive firms at the given price. Condition [3] is actually implied by condition [2] for any contract that has positive demand in the equilibrium, but requiring condition [3] for all contracts adds a further restriction on the pricing of contracts that have zero demand under γ . Without this restriction, we could get equilibria that satisfy [2] trivially for any one contract x simply by setting its price $q(x)$ so high that all $\gamma(x|t) = 0$. With condition [3], a real possibility of unprofitable sales can explain what deters firms from trying to increase demand for a contract by shading its price slightly.

When all costs $C(x,t)$ are nonnegative, any positively traded contract must have a nonnegative price in equilibrium; but the definition of competitive equilibrium here allows the possibility that some untraded contracts might have negative prices, as if there were some small fund for subsidizing sales of the contract. Of course, a contract that has zero demand at a negative price would have the same zero demand at the higher price of zero. The point of imputing a negative price here is only to identify which types would be willing to buy the contract with the least subsidy, so that we can verify that the imputed price is not greater than the costs of serving these types.

To simplify some equations below, let us introduce the notation: $\bar{U}(q,t) = \max_{x \in Y} U(q(x),x,t)$.

Fact. When Y and T are finite sets and U satisfies the assumptions that are listed above, a competitive equilibrium must exist.

Proof. For each x in Y , let $\bar{c}(x) = \max_{t \in T} C(x,t)$, and $\underline{c}(x) = \min_{t \in T} C(x,t)$.

Then let $c_0(x) = \min_{t \in T} \min_{y \in Y} \varphi(\underline{c}(y),y,t,x) - 1$.

Here we have $c_0(x) < \underline{c}(x)$, because $\varphi(\underline{c}(x),x,t,x) = \underline{c}(x)$, and we also have $\underline{c}(x) \leq \bar{c}(x)$. Now, for any ε such that $0 < \varepsilon < 1$, consider a modified market in which the fraction of each type t in T is $(1-\varepsilon)\mu(t)$ and, on for each contract x , a fraction $\varepsilon/\#Y$ of the consumers are a new artificial type that only buys contract x and has cost $c_0(x)$. (Here $\#Y$ is the number of contracts in the finite set Y .)

Now for any price vector q in $\times_{x \in Y} [c_0(x),\bar{c}(x)]$ and any allocation vector γ in $\Delta(Y)^T$ such that, consider a mapping that selects a new price vector q' and a new allocation vector γ' as follows.

For each contract x , $q'(x)$ is the average cost

$$q'(x) = (c_0(x)\varepsilon/\#Y + \sum_{t \in T} C(x,t)\mu(t)\gamma(x|t)) / (\varepsilon/\#Y + \sum_{t \in T} \mu(t)\gamma(x|t)).$$

For each type t , $\gamma'(\bullet|t)$ can be any probability distribution over Y such that

$$\{x \in Y | \gamma'(x|t) > 0\} \subseteq \operatorname{argmax}_{x \in Y} U(q(x),x,t), \quad \forall t \in T.$$

By the Kakutani fixed-point theorem, we can find a pair $(q^\varepsilon, \gamma^\varepsilon)$ in $(\times_{x \in Y} [c_0(x),\bar{c}(x)]) \times \Delta(Y)^T$ that maps to itself under this mapping. This $(q^\varepsilon, \gamma^\varepsilon)$ will satisfy conditions [0] and [1] and also a perturbed version of the zero-profit condition [2] with the new artificial types included as an ε fraction of the population (which ensure that every contract gets strictly positive demand).

Now consider a sequence of numbers ε that converge to 0, so that the artificial types become an infinitesimal fraction of the population as we go to the limit. By compactness of the domain $(\times_{x \in Y} [c_0(x),\bar{c}(x)]) \times \Delta(Y)^T$, there exists a subsequence such that all the $q^\varepsilon(x)$ and $\gamma^\varepsilon(x|t)$ converge to some limits $q^*(x)$ and $\gamma^*(x|t)$ that are also in this domain. We can now show that this (q^*, γ^*) is a competitive equilibrium.

The fact that conditions [0] and [1] and an approximate version of [2] are satisfied for each $(q^\varepsilon, \gamma^\varepsilon)$ implies, by continuity, that these conditions are satisfied for the limit (q^*, γ^*) . Furthermore, the limit of $\varepsilon \rightarrow 0$ implies that every contract x that has positive γ^* -demand must have $q^*(x)$ that is a weighted average of $C(x,t)$ costs and so satisfies $\underline{c}(x) \leq q^*(x) \leq \bar{c}(x)$. Each type t must be generating positive demand $\gamma^*(x|t) > 0$ for at least one of its optimal contracts x , for which we then get that t 's optimal utility $\bar{U}(q^*,t)$ must satisfy $\bar{U}(q^*,t) = U(q^*(x),x,t) \leq U(\underline{c}(x),x,t)$. (Utility is decreasing in the price.)

Now let us show that condition [3] is satisfied for any contract y . Since $c_0(y) < \underline{c}(y)$, the price $q^*(y)$ could not equal $c_0(y)$ unless demand for y was zero. But $c_0(y)$ was defined so that any type would strictly prefer y at price $c_0(y)$ over any other contract x priced at $\underline{c}(x)$. Thus, for every contract y , we must have $q^*(y) > c_0(y)$, which implies $q^\varepsilon(y) > c_0(y)$ for all sufficiently small ε , which in turn implies that $\gamma^\varepsilon(y|t) > 0$ for at least one real type t in T . (The demand $\gamma^\varepsilon(y|t)$ would be small, of order ε , but it must be strictly positive.) As there are only finitely many types, we can choose a subsequence (if necessary) so that there is some type t in T such that $\gamma^\varepsilon(y|t) > 0$ for all ε , and so we get $U(q^\varepsilon(y),y,t) = \bar{U}(q^\varepsilon,t)$ for all ε . Among such types t , we can pick the one with the highest cost in contract y , and then we must also get $q^\varepsilon(y) \leq C(y,t)$ for all ε , because the price $q^\varepsilon(y)$ is a

weighted average of the real types that choose y in γ^ε and the lower-cost artificial type for y . So in the limit, the price $q^*(y)$ satisfies condition [3], $U(q^*(y), y, t) = \bar{U}(q^*, t)$ and $q^*(y) \leq C(y, t)$, with this type t .

Thus (q^*, γ^*) is a competitive equilibrium. *QED*

For markets where the set of possible contracts Y is infinite, one can still prove an existence of competitive equilibria with some additional continuity assumptions which are developed by Azevedo and Gottlieb (2017). (The basic idea of the proof is to consider the limits of competitive equilibria that we would get for an increasing sequence of finite subsets of Y which in the limit become dense in the whole set.)

Now let us consider a class of models that include Rothschild-Stiglitz insurance markets with two type. Let $T = \{L, H\}$ where L is the low type and H is the high type. The type may be interpreted as the probability of suffering a loss, with $0 < L < H < 1$. The set of possible contracts Y will be the interval $[0, 1]$ (or some subset of this interval), where any x in $[0, 1]$ denotes the fraction of the insured loss that will be covered by the insurance company. In addition to the previous assumptions about U , we now add several further assumptions about the U and C functions in our model, all of which will be satisfied by a Rothschild-Stiglitz insurance model:

$C(x, t)$ is continuous in x and $U(p, x, t)$ is continuous in x and p , $\forall t \in T$;

$C(0, L) = C(0, H) = 0$, but $0 < C(x, L) < C(x, H) \forall x > 0$;

$U(C(x, t), x, t)$ is strictly increasing in x , $\forall t \in T$;

if $y > x$ then $\varphi(p, x, L, y) < \varphi(p, x, H, y)$, but if $y < x$ then $\varphi(p, x, L, y) > \varphi(p, x, H, y)$.

The second assumption says that the high types are more expensive to serve. The third assumption says that any type of consumer would prefer full insurance at a price that is actuarially fair for this type. The fourth assumption is a "single-crossing" condition. It says that, when consumers are asked about their willingness to switch from a contract x at a price p to some other contract y , if y is greater than x then type H would be willing to switch than type L , in the sense that type H would accept the increased coverage y for a higher price than type L would accept; but if $y < x$ then the type L would be more willing to switch than type H . If U is continuously differentiable in both p and x , then the single-crossing property is implied by a condition that, when we normalize the marginal utility of money ($-p$) across types, type H has greater marginal utility for increasing coverage x .

$$(\partial U(p, x, L) / \partial x) / (-\partial U(p, x, L) / \partial p) < (\partial U(p, x, H) / \partial x) / (-\partial U(p, x, H) / \partial p).$$

Fact. When $Y = [0, 1]$, $T = \{L, H\}$, and U and C satisfy the above assumptions, the unique competitive equilibrium is the best separating plan, which is defined as follows. Type H chooses the contract $x_H = 1$ which has price $q(1) = C(1, H)$. Type L chooses a contract x_L at the price $q(x_L) = C(x_L, L)$, where x_L is the solution to the equation $U(C(x_L, L), x_L, H) = U(C(1, H), 1, H)$. Any other contract y has the price $q(y) = \varphi(C(x_L, L), x_L, L, y)$ if $y < x_L$, and $q(y) = \varphi(C(1, H), 1, H, y)$ if $y > x_L$.

Proof. Every contract x in the interval $[0, 1]$ must be priced so that at least one type is willing to buy it, so that either $U(q(x), x, L) = \bar{U}(q, L)$ or $U(q(x), x, L) = \bar{U}(q, H)$. The pricing function $q(x)$ cannot be discontinuous because, at any point of discontinuity, the type that is (supposedly) willing to choose contracts that converge to the highest price limit would actually strictly prefer to switch to a nearby contract that has a price close to the lowest price limit. Thus, the set of contracts that any type is willing to choose at the q prices is a closed set, and so the connected interval $Y = [0, 1]$ must include some contract y^* that both types are willing to choose. That is, we have some y^* such that $U(q(y^*), y^*, L) = \bar{U}(q, L)$ and $U(q(y^*), y^*, H) = \bar{U}(q, H)$.

By the single-crossing property, this y^* must be unique, and only type L is willing to choose any contract $x < y^*$. (If $x < y^*$ satisfied $U(q(x), x, H) = \bar{U}(q, H) = U(q(y^*), y^*, H)$, then type L would strictly prefer x at price $q(x)$ over y^* at price $q(y^*)$, contradicting the fact that y^* is optimal for both types.) Similarly, if $x > y^*$, then only type H is willing to choose the contract x at price $q(x)$.

So for any $\varepsilon > 0$, condition [3] requires $U(q(y^* - \varepsilon), y^* - \varepsilon, L) = \bar{U}(q, L)$ and $q(y^* - \varepsilon) \leq C(y^* - \varepsilon, L)$. Taking the limit as $\varepsilon \rightarrow 0$, we must have $q(y^*) \leq C(y^*, L)$. The contracts that L actually chooses must be in the interval $[0, y^*]$. If L actually had positive demand for some contract $x < y^*$, then the zero-profit condition for firms would require

that contract to be priced at $q(x)=C(x,L)$, but $U(C(x,L),x,L)$ is strictly increasing in x , and so we would get the contradiction

$$\bar{U}(q,L) = U(q(x),x,L) = U(C(x,L),x,L) < U(C(y^*,L),y^*,L) \leq U(q(y^*),y^*,L) = \bar{U}(q,L).$$

Thus type L cannot be actually buying any contract less than y^* , which is the highest contract that L is willing to choose. We found that $q(y^*) \leq C(y^*,L)$, but for firms to break even in selling y^* , its price cannot be strictly less than the low cost $C(y^*,L)$, and so we must have $q(y^*) = C(y^*,L)$, and only the type-L consumers are actually buying this contract. That is, y^* is the contract that was called x_L in the statement of the Fact above.

Then type H consumers must be actually buying some strictly higher contract $z > y^*$ which must satisfy the break-even price of $q(z)=C(z,H)$. But $U(C(z,H),z,H)$ is strictly increasing in z , and we know that the greatest contract 1, which only H is willing to choose, must in equilibrium satisfy $q(1) \leq C(1,H)$. Thus, if type H cannot be buying any contract $z < H$, because it would yield the contradiction $\bar{U}(q,H) = U(C(z,H),z,H) < U(C(1,H),1,H) \leq U(q(1),z,H) = \bar{U}(q,H)$. Thus, type H is buying the maximal contract 1 at the price $C(1,H)$.

QED

The definition of competitive equilibrium here is a simplified version of the equilibrium concept developed by Azevedo and Gottlieb [2017]. Their concept differs from the simple competitive equilibrium concept here in that they require that the prices of untraded contracts should be evaluated as limits of equilibrium prices from a sequence of perturbed models in which vanishingly small populations of artificial consumers are introduced to guarantee some positive low-cost demand for every contract. (Such a construction is used in the proof of equilibrium existence above.) This perturbational condition is analogous to conditions that are applied in perfect and proper refinements of Nash equilibrium. Examples can be constructed (with infinitely many types, including pairs of types that have identical utility functions but different costs) such that the perfect competitive equilibrium concept of Azevedo and Gottlieb helps to exclude some counter-intuitive equilibria that would be admitted by the simple competitive equilibrium concept developed here.

Markets with adverse selection on the sellers' side (e.g. Spence signaling in labor markets)

We can also model markets where competitive firms buy labor or other inputs from individuals with private information about themselves that might be relevant for their potential productivity in any labor-supply contract. Again in this case we can let Y denote the set of possible labor contracts, and let T denote the set of possible types of individual workers. A given utility function $U: \mathbb{R} \times Y \times T \rightarrow \mathbb{R}$ specifies the utility $U(p,x,t)$ that any consumer of type t would get from selling his labor under the terms of contract x at wage p . In this case, of course, the worker's utility $U(p,x,t)$ would be a strictly increasing function of the wage-price p . A given productivity function $V: Y \times T \rightarrow \mathbb{R}$ specifies the expected output value $V(x,t)$ for a firm from a labor contract x with a worker of type t . A given probability distribution μ in $\Delta(T)$ specifies the fraction $\mu(t) > 0$ of workers who have each type t in the general population. Each worker must enter into exactly one labor contract. (We could include the no-trade option as an $x=0$ contract with productivity 0 and utility 0 for all types.)

Our simple concept of competitive equilibrium can be directly extended to such models. Then a competitive equilibrium would be any pair (q,γ) such that $q=(q(x))_{x \in Y}$ is a price vector in \mathbb{R}^Y , $\gamma=(\gamma(x|t))_{x \in Y, t \in T}$ is an allocation vector in $\Delta(Y)^T$, and the following conditions are satisfied:

- [0'] $\sum_{y \in Y} \gamma(y|t) = 1$ and $\gamma(x|t) \geq 0, \forall x \in Y, \forall t \in T$;
- [1'] $\sum_{x \in Y} \gamma(x|t) U(q(x),x,t) = \max_{x \in Y} U(q(x),x,t), \forall t \in T$;
- [2'] $\sum_{t \in T} \mu(t) \gamma(x|t) (V(x,t) - q(x)) = 0, \forall x \in Y$; and
- [3'] $\forall y \in Y, \exists t \in T$ such that $U(q(y),y,t) = \max_{x \in Y} U(q(x),x,t)$ and $q(y) \geq V(y,t)$.

Thus, for markets with adverse selection on the sellers' side, we can similarly require that [2'] competitive firms should expect to break even on any contract that attracts a positive fraction of the informed sellers, and that, [3'] even for an untraded contract, there should be at least one seller-type which would be willing to choose this contract but would not be profitable for the competitive firms at the given price.

Moral hazard with constant risk tolerance, monetary effort cost, 2 actions.

A risk-neutral principal is designing an incentive plan for a risk-averse agent.

The agent must choose among two unobservable actions: a_H and a_L .

Each action has a cost to the agent, $c(a_L) = c_L < c(a_H) = c_H$. The agent could earn w_0 elsewhere.

Suppose the agent's utility from choosing action a and getting wage w would be

$u(w-c(a)) = -\text{EXP}(-(w-c(a))/T)$, where $T>0$ is the agent's constant risk tolerance.

[This model differs from standard textbook models in that effort cost here is monetary, subtracted from income before applying the utility function: $u(w-c(a))$ instead of $u(w)-c(a)$.]

The principal can only observe an outcome \tilde{y} that is a random variable which depends on the agent's action according to the conditional probability distribution $p(y|a_H)$ or $p(y|a_L)$.

Let Y denote the set of possible values of \tilde{y} , which we assume here to be a finite set.

The principal can promise the agent a wage $w(y)$ that depends on the observable outcome.

Suppose that the principal's expected payoff is much higher when the agent chooses a_H .

So the principal's problem is to design the wage-function $w(\bullet)$ to minimize the expected wage expense, subject to the constraints that the agent should not prefer the outside option w_0 or the lower action a_L :

$$\begin{aligned} \text{Choose } (w(y))_{y \in Y} \text{ to minimize } & \sum_{y \in Y} p(y|a_H) w(y) \text{ subject to} \\ & \sum_{y \in Y} p(y|a_H) u(w(y)-c(a_H)) \geq u(w_0) && \text{(participation constraint: } \lambda), \\ & \sum_{y \in Y} p(y|a_H) u(w(y)-c(a_H)) \geq \sum_{y \in Y} p(y|a_L) u(w(y)-c(a_L)) && \text{(moral-hazard constraint: } \mu). \end{aligned}$$

The participation constraint must be binding, or else the principal could reduce all $w(y)$.

If the moral-hazard constraint were not binding, then the optimal solution would be $w(y)=w_0+c_H$ for all outcomes y , but then the agent would prefer the lower-cost action a_L (with $c(a_L)<c(a_H)$).

So both constraints must be binding at the optimal solution. Let λ and μ denote the Lagrange multipliers of the participation and moral-hazard constraints respectively. The Lagrangean is

$$\begin{aligned} \mathcal{L}(w;\lambda,\mu) = & \sum_y p(y|a_H) w(y) - \lambda[\sum_y p(y|a_H) u(w(y)-c_H) - u(w_0)] \\ & - \mu[\sum_y p(y|a_H) u(w(y)-c_H) - \sum_y p(y|a_L) u(w(y)-c_L)]. \end{aligned}$$

The optimality conditions $0 = \partial \mathcal{L} / \partial w(y)$ yield

$$\forall y \in Y: \quad 0 = p(y|a_H) - (\lambda + \mu) p(y|a_H) u'(w(y)-c_H) + \mu p(y|a_L) u'(w(y)-c_L).$$

With constant risk tolerance, $u'(w-c) = \text{EXP}(-(w-c)/T)/T = -u(w-c)/T$.

Thus, summing the optimality conditions over all y in Y , we get

$$0 = 1 + (\lambda + \mu) \sum_y p(y|a_H) u(w(y)-c_H)/T - \mu \sum_y p(y|a_L) u(w(y)-c_L)/T,$$

which with the binding constraints yields $0 = 1 + (\lambda + \mu) u(w_0)/T - \mu u(w_0)/T = 1 + \lambda u(w_0)/T$.

So the participation constraint's Lagrange multiplier is $\lambda = 1/u'(w_0) = -T/u(w_0) = T \times \text{EXP}(w_0/T)$.

With constant risk tolerance T , $u'(w(y)-c_L) = \eta u'(w(y)-c_H)$, where $\eta = \text{EXP}(-(c_H-c_L)/T)$.

So the optimality conditions become $0 = 1 - u'(w(y)-c_H)[\lambda + \mu - \mu \eta p(y|a_L)/p(y|a_H)]$, $\forall y \in Y$.

Then the optimal wage $w(y)$ can be determined from the equations:

$$T \times \text{EXP}((w(y)-c_H)/T) = 1/u'(w(y)-c_H) = \lambda + \mu - \mu \eta p(y|a_L)/p(y|a_H), \quad \forall y \in Y.$$

Thus, the optimal wage $w(y)$ is monotone decreasing in the likelihood ratio $p(y|a_L)/p(y|a_H)$.

With $\mu \leq \lambda/(\eta \max_{y \in Y} p(y|a_L)/p(y|a_H) - 1)$, μ is determined by the requirement that the moral-hazard constraint must be satisfied as a binding equality.

Moral hazard with a linear wage rule, Normal risks, and constant risk tolerance.

Suppose that a firm's net profits will be a random variable \tilde{y} drawn from a Normal distribution with with a mean that depends on an agent's effort, but with a variance σ^2 that is independent of effort.

If the agent's effort is high (a_H), then the agent's effort cost is c_H and profit has mean $E(\tilde{y}|a_H) = m_H$.

If the agent's effort is low (a_L), then the agent's effort cost is c_L and profit has mean $E(\tilde{y}|a_L) = m_L$.

Suppose $m_H > m_L$ and $c_H > c_L$.

Nobody else can observe the agent's effort, but his wage can depend on the observable profit \tilde{y} .

For now, let us assume that firm must specify the agent's wage as a linear function of observed profits, according to any linear formula $w(\tilde{y}) = \alpha + \beta\tilde{y}$.

(The likelihood-ratio result of the previous model suggests that such linear rules might be inferior to nonlinear rules. But Holmstrom and Milgrom [1987] provided a fundamental reason why such linear rules would actually be optimal solutions to the moral hazard problem: Suppose that Normally distributed returns accrue over some interval of time in a dynamic Brownian-motion process that has effort-dependent drift, and suppose that the agent can change his hidden effort at any point in time depending on the current state and past history of the observable Brownian process.)

Suppose that the agent has a constant risk tolerance T , but the owners of the firm are risk neutral, which means that they have linear utility for money or infinite risk tolerance.

The agent's best outside option is to earn w_0 elsewhere with zero effort costs.

What linear wage rule (α, β) is best for the firm if the firm wants the agent to choose high effort?

Under an (α, β) linear wage rule, the agent's net certainty equivalent with high effort is

$$E(\alpha + \beta\tilde{y}|a_H) - (0.5/T)\text{Var}(\alpha + \beta\tilde{y}|a_H) - c_H = \alpha + \beta m_H - (0.5/T)(\beta\sigma)^2 - c_H.$$

By choosing his outside option, the agent could instead get certainty equivalent w_0 .

With low effort, the agent could get the certainty equivalent $\alpha + \beta m_L - (0.5/T)(\beta\sigma)^2 - c_L$.

So for the agent to want to work here and exert high effort, (α, β) must satisfy the

participation constraint $\alpha + \beta m_H - (0.5/T)(\beta\sigma)^2 - c_H \geq w_0$, and the

moral-hazard incentive constraint $\alpha + \beta m_H - (0.5/T)(\beta\sigma)^2 - c_H \geq \alpha + \beta m_L - (0.5/T)(\beta\sigma)^2 - c_L$.

These constraints imply $\alpha + \beta m_H \geq w_0 + c_H + (0.5/T)(\beta\sigma)^2$ and $\beta \geq (c_H - c_L)/(m_H - m_L)$.

The risk-neutral owners of the firm want to minimize their expected wage bill $\alpha + \beta m_L$.

So the optimal linear rule should choose α (given any β) to make the participation constraint bind, and should choose β as small as possible to make the moral-hazard constraint bind.

That is, the optimal rule is $\beta = (c_H - c_L)/(m_H - m_L)$ and $\alpha = w_0 + c_H + (0.5/T)(\beta\sigma)^2 - \beta m_H$.

The firm's minimal expected wage bill is $w_0 + c_H + (0.5/T)(\beta\sigma)^2$, and the owner's expected net profit is thus $m_H - w_0 - c_H - (0.5/T)(\beta\sigma)^2$.

Given that the agent is risk averse and the owners are risk neutral, the optimal sharing rule without moral hazard would have $\beta=0$ (and so $\alpha=w_0+c_H$), if the agent could be forced to choose high effort.

So moral hazard provides a fundamental reason why risk-averse agents cannot insure away their professional risks: because their share of such risks is essential to motivate their efforts.

Credit rationing in a binary moral-hazard model with limited liability (Stiglitz-Weiss).

An agent can undertake a project which requires an amount I to be invested this period.

Next period, the project will return the amount RI if it succeeds, or 0 if it fails.

The agent must borrow the amount I but can offer collateral worth C for the loan.

If the agent manages the investment well, then its probability of success will be p_H .

But the agent could act badly, divert an amount γI , and reduce the probability of success to p_L .

This reduction in the probability of success is the only observable consequence of such malfeasance.

The agent is risk neutral but has limited liability in the sense that he cannot lose more than C .

Let w_0 denote the agent's alternative wage next period if he does not manage this project,

Let ρ denote the market rate of interest per period for risk free investments elsewhere.

Suppose that $p_H R > 1 + \rho > p_L R + \gamma$, so that the project could be worthwhile for risk-neutral investors if the agent manages it well, but not if he acts badly.

If the agent is charged an interest rate r to borrow I , then the agent will get the surplus $(R - (1+r))I$ if the project is a success, but he will default and lose the collateral C if the project fails.

Including the value of defaulted collateral, the investors' expected return is $p_H(1+r)I + (1-p_H)C$.

Per unit invested, the investors' expected rate of return is $p_H(1+r) + (1-p_H)C/I$.

The agent's participation constraint is $p_H(R - (1+r))I - (1-p_H)C \geq w_0$.

The agent's participation constraint implies an upper bound on the investors' rate of return:

$$p_H R - w_0/I \geq p_H(1+r) + (1-p_H)C/I.$$

The moral-hazard constraint to deter agential malfeasance is

$$p_H(R - (1+r))I - (1-p_H)C \geq \gamma I + p_L(R - (1+r))I - (1-p_L)C.$$

The moral-hazard constraint yields a lower bound on an agent's stake $C + (R - (1+r))I \geq \gamma I / (p_H - p_L)$,

and also an upper bound on investors' returns: $C/I + p_H R - p_H \gamma / (p_H - p_L) \geq p_H(1+r) + (1-p_H)C/I$.

Suppose that C and w_0 are small so that $C + w_0 < p_H \gamma I / (p_H - p_L)$.

[poor agent case]

This poor-agent condition implies that $C/I + p_H R - p_H \gamma / (p_H - p_L) < p_H R - w_0/I$.

Then the upper bound on investors' rate of return is determined by the moral-hazard constraint:

$$p_H(1+r) + (1-p_H)C/I \leq C/I + p_H R - p_H \gamma / (p_H - p_L).$$

The maximal interest rate r^* that can avoid moral hazard is $1 + r^* = R + C/I - \gamma / (p_H - p_L)$.

Suppose also that R satisfies $R \geq (1+\rho)/p_H + \gamma / (p_H - p_L) - C / (p_H I)$.

Then the investors are willing to lend at this rate r^* , because

$$p_H(1+r^*)I + (1-p_H)C = p_H[R + C/I - \gamma / (p_H - p_L)]I + (1-p_H)C \geq (1+\rho)I.$$

Even at this maximal interest rate, however, the agent's expected gain is strictly positive:

$$p_H(R - (1+r^*))I - (1-p_H)C - w_0 = p_H \gamma / (p_H - p_L) - C - w_0 > 0.$$

The quantity $p_H \gamma / (p_H - p_L)$ here may be called the expected moral-hazard rent that the agent must be allowed in this financial transaction, to deter his malfeasance.

Now suppose that there are many such agents, but only a limited supply of funds from investors who understand these projects well enough to enter into such financial deals.

Then the interest rate for such loans can rise only to the maximal rate r^* .

If the demand from agents at this rate exceeds the supply of funds, there must be credit rationing.

If the agent also discounts future payoffs by the discount factor $\delta = 1/(1+\rho)$ per period, then the agent's participation and moral hazard constraints become (in second-period values):

$$p_H(R - (1+r))I - (1-p_H)C \geq w_0/\delta,$$

$$p_H(R - (1+r))I - (1-p_H)C \geq \gamma I/\delta + p_L(R - (1+r))I - (1-p_L)C.$$

Then the resulting bounds on the investors' expected rate of return become:

$$\min\{p_H R - w_0/(\delta I), p_H R + C/I - p_H \gamma / (\delta(p_H - p_L))\} \geq p_H(1+r) + (1-p_H)C/I.$$

Long-term agency relationships and back-loaded moral hazard rents

Let us consider a model similar to the credit-rationing model, but now the projects are producing some special output. Many competitive firms are able to produce this output, but they must hire agents to manage the production projects. The size of a production project can be measured in terms of the value of the inputs I that are given to the agent for the project. When an agent is given inputs worth I (say, measured in dollars) to manage in one period, then next period the agent will produce either AI units of output (say, measured in lbs) if the project succeeds, or else output will be if the project fails. (Here A is a constant measured in lbs of output per dollar of input.) Let p_H denote the probability of success if the manager acts appropriately, but the manager could instead divert a fraction γ of the invest funds and reduce the probability of success to p_L . Here $p_L < p_H$.

We assume here that $\gamma < (p_H/p_L + p_L/p_H - 2)$, a condition which will guarantee that competitive firms cannot expect to profit from employing agents who divert the γ fraction of inputs.

The inputs I that one individual agent can manage in any one period can be any amount in some wide interval, from some lower bound \underline{I} to some upper bound \bar{I} . We assume that this range of feasible project sizes is quite wide (so that the ratio \bar{I}/\underline{I} is large), but agents' alternative wages and collateral are negligible relative to even the minimal size \underline{I} of these projects (so that, in the notation of the previous credit-rationing model, we can let $w_0 = 0$ and $C = 0$).

An agent's career can span T periods, after which the agent will retire. We assume that agents and investors are risk-neutral and discount future income at some per-period discount factor δ .

Let π denote the price per unit of the outputs produced by these projects. When it is profitable for firms to hire agents to manage such production projects, then many competitive firms should be expected to do so, driving down the price of the good which is the output of these projects. So let us ask, what is the lowest output price π at which hiring agents under an optimal incentive contract can be just barely profitable for competitive firms.

Assuming limited liability for the agents, the terms of an agent's contract in any period must specify (I, u_S, u_F) , where $I \geq 0$ is the value of inputs to be managed by the agent this period, and next period the expected value of the agent's future compensation will be u_S if this project succeeds, or u_F if this project fails. For the agent to manage the plan appropriately, the plan must satisfy the moral-hazard constraint is $p_H \delta u_S + (1 - p_H) \delta u_F \geq \gamma I + p_L \delta u_S + (1 - p_L) \delta u_F$, and the limited-liability and investment-bound constraints are $u_S \geq 0$, $u_F \geq 0$, $\underline{I} \leq I \leq \bar{I}$. (Here the agent's participation constraint $p_H \delta u_S + (1 - p_H) \delta u_F \geq 0$ is trivially satisfied.)

The moral-hazard constraint here implies $u_S \geq u_F + \gamma I / (\delta(p_H - p_L))$.

In one period, the firm can minimize the cost per unit of output by choosing $u_F = 0$ and $u_S = \gamma I / (\delta(p_H - p_L))$. Under this one-period production plan, the expected present-discounted value of revenue from outputs next period would be $\delta p_H \pi AI$, and the expected present-discounted value of costs would be $I + \delta p_H \gamma I / (\delta(p_H - p_L)) = I + p_H \gamma I / (p_H - p_L)$.

So this one-period production plan would be profitable at an output price π above

$$\pi_1 = (I + p_H \gamma I / (p_H - p_L)) / (\delta p_H AI) = (1 / p_H + \gamma / (p_H - p_L)) / (\delta A).$$

If the agent were not promised large rewards worth $\gamma I / (\delta(p_H - p_L))$ for success, he would divert a γ fraction of the inputs to his own current consumption have have only p_L probability of success; but such an alternative would not be worthwhile at any price less than $1 / (\delta A p_L)$. So to ensure that competitive firms will pay to satisfy the moral-hazard constraint, we need the inequality $1 / (\delta A p_L) > (1 / p_H + \gamma / (p_H - p_L)) / (\delta A)$, which is equivalent to our parametric assumption $\gamma < (p_H / p_L + p_L / p_H - 2)$.

The rewards for the agent in this one-period have an expected present-discounted value of $\delta p_H u_S = p_H \gamma I / (p_H - p_L)$, which is the expected moral-hazard rent that the agent gets from the power that his job gives him over the input resources I . This moral-hazard rent is not needed to make the agent want to take this job (as we assumed that any positive compensation would satisfy the participation constraint); it is only needed to deter the agent from abusing the power of this job. So the opportunity to have such a job with its moral-hazard rent can itself be a reward that an agent should value. This fact implies that a firm should be able to reduce its expected costs per unit of revenue further by employing agents with a multi-period incentive contract where an agent in one period can be motivated by the prospect of managing larger projects in later periods if he can succeed in the current project.

For example, consider extending the one-period contract above to a two-period contract as follows. Let I_1 denote the value of inputs that the agent manages in his first period. If the project succeeds, instead of simply paying the agent $u_S = \gamma I_1 / (p_H - p_L)$ in cash, the firm could instead give the agent an opportunity to manage a larger project with input I_2 such that the expected moral-hazard rent $p_H \gamma I_2 / (p_H - p_L)$ is just equal to the promised value u_S . That is, I_2 should satisfy $p_H \gamma I_2 / (p_H - p_L) = \gamma I_1 / (\delta(p_H - p_L))$, which holds when

$$I_2 = I_1 / (\delta p_H).$$

Under this two-period incentive plan, the agent will be paid $\gamma I_2 / (\delta(p_H - p_L))$ in period 3 only if the agent's first and second projects both succeed; otherwise he gets no compensation; and the second project with inputs I_2 is implemented only if the first project with inputs I_1 is a success. So from the perspective of period 1, the expected present-discounted value of revenues is

$$\pi \delta p_H A I_1 + \pi (\delta p_H)^2 A I_2 = 2\pi \delta p_H A I_1$$

and the expected discounted value of costs is

$$I_1 + \delta p_H I_2 + (\delta p_H)^2 \gamma I_2 / (\delta(p_H - p_L)) = 2I_1 + p_H \gamma I_1 / (p_H - p_L).$$

Thus the two-period incentive contract would yield positive expected profits for the firm at any output-price greater than $\pi_2 = [1/p_H + (1/2) \gamma / (p_H - p_L)] / (\delta A)$.

More generally, a firm can operate profitably at even lower output prices with a T -period incentive contract as follows. In the first period of an agent's career, the agent is entrusted with the minimum feasible amount of inputs $I_1 = \underline{I}$. At any later period t in $\{2, \dots, T\}$, if the agent's past projects were all successful then the agent will be entrusted with inputs worth

$$I_t = \underline{I} / (\delta p_H)^{t-1}.$$

If the agent ever has an unsuccessful project then he will be dismissed without any further compensation, but if all his projects succeed then, in period $T+1$ he will retire with payment

$$w_{T+1} = \gamma I_T / (\delta(p_H - p_L)) = p_H \gamma \underline{I} / ((\delta p_H)^T (p_H - p_L)).$$

At each period t , these equations yield

$$p_H \gamma I_t / (p_H - p_L) = w_{T+1} (\delta p_H)^{T+1-t}.$$

So under this plan, the moral hazard rents for all responsibilities over the agent's entire career are back-loaded to one big payment at the end of his career, conditional on good performance throughout the agent's career. For this plan to be feasible, we must assume that the range of feasible project sizes is wide enough that

$$\bar{I} \geq \underline{I} / (\delta p_H)^{T-1}.$$

With this T -period plan, the expected period-1 discounted value of the firm's revenues is

$$\sum_{t \in \{1, \dots, T\}} \pi (\delta p_H)^t A I_t = T \pi \delta p_H A \underline{I}$$

and the expected discounted value of the firm's costs is

$$\sum_{t \in \{1, \dots, T\}} (\delta p_H)^{t-1} \underline{I} + (\delta p_H)^T w_{T+1} = T \underline{I} + p_H \gamma \underline{I} / (p_H - p_L).$$

So this plan (with moral-hazard rents paid only at period T+1 and then only if the agent has had a consistently good record over T periods) can be profitable at any output price greater than

$$\pi^* = [1/p_H + (1/T) \gamma / (p_H - p_L)] / (\delta A).$$

Optimality of this production plan can be verified by Lagrangean analysis. The simplest way to do this is to decompose the T-period contract-design problem into T one-period problems as follows. At any period t in $\{1, \dots, T\}$, given the expected (long-run) reward v_t that has been promised to the agent under the contract in the previous period, the firm wants to choose this period's terms of employment $(I_t, u_{S,t}, u_{F,t})$ so as to

$$\begin{aligned} & \text{maximize } (\delta p_H \pi_{t+1} A - 1) I_t - \delta \lambda_{t+1} [p_H u_{S,t} + (1 - p_H) u_{F,t}] && \text{(net profit)} \\ & \text{subject to } \delta p_H u_{S,t} + \delta (1 - p_H) u_{F,t} \geq v_t && \text{(promise-keeping, } \lambda_t) \\ & \delta p_H u_{S,t} + \delta (1 - p_H) u_{F,t} - [\gamma I_t + \delta p_L u_{S,t} + \delta (1 - p_L) u_{F,t}] \geq 0 && \text{(moral hazard, } \mu_t) \\ & u_{S,t} \geq 0, u_{F,t} \geq 0. && \text{(limited liability)} \end{aligned}$$

Here λ_{t+1} is the expected discounted cost (in terms of dollars at period t+1) of promising the agent a reward worth u at period t+1. The price of output at period t+1 is denoted here by π_{t+1} . The promised reward v_t at any period t would be previous period's $u_{S,t-1}$ or $u_{F,t-1}$, depending on whether the previous project was a success or failure. At the final employment period $t=T$, we must have $\lambda_{T+1}=1$, because the reward will have to be paid in cash at period T+1. At any earlier period t, λ_{t+1} should be equal to the Lagrange multiplier of the promise-keeping constraint, which we call λ_t in our analysis of the period-t problem. With μ_t denoting the Lagrange multiplier of the moral-hazard constraint, Lagrangean analysis can verify the optimality of a solution with I_t strictly between its upper and lower bounds and with $u_{S,t} > 0$. The Lagrangean optimality conditions (below) imply

$$\begin{aligned} \mu_t &= (\delta p_H \pi_{t+1} A - 1) / \gamma \quad (\text{so that } \partial \mathcal{L} / \partial I_t = 0), \text{ and} \\ \lambda_t &= \lambda_{t+1} - \mu_t (p_H - p_L) / p_H = \lambda_{t+1} - (\delta p_H \pi_{t+1} A - 1) (p_H - p_L) / (\gamma p_H) \quad (\text{so that } \partial \mathcal{L} / \partial u_{S,t} = 0). \end{aligned}$$

With these multipliers, we get $\partial \mathcal{L} / \partial u_{F,t} = \delta (\lambda_t - \lambda_{t+1}) < 0$, so that the optimal solution has $u_{F,t} = 0$, and both constraints are satisfied as equalities with the solution $u_{S,t} = v_t / (\delta p_H)$ and $I_t = v_t (p_H - p_L) / (\gamma p_H)$.

But there is one difficulty with this solution, as the firm starts at period 1 with no contractual obligation to the agent, so that $v_1 = 0$. In any contract that involves positive production, the requirement of positive expected moral-hazard rents implies that the promise-keeping constraint at period 1 must be satisfied as a strict inequality with positive slack. For such a contract to be optimal, the Lagrange multiplier at period 1 must be zero. With the above recursive formula for λ_t , $\lambda_1 = 0$ and $\lambda_{T+1} = 1$ together imply

$$1 - \sum_{t \in \{1, \dots, T\}} (\delta p_H \pi_{t+1} A - 1) (p_H - p_L) / (\gamma p_H) = 0.$$

This equation is satisfied when all π_{t+1} are equal to the price π^* that we found above. In fact, it is satisfied whenever the average of the π_{t+1} prices is π^*

$$\sum_{t \in \{1, \dots, T\}} \pi_{t+1} / T = \pi^*$$

With this π^* -average price of output, our T-period plan satisfies all the Lagrangean conditions for optimality, provided also that each $\pi_{t+1} \geq 1 / (\delta p_H A)$, so that our solution has $\mu_t \geq 0$.

Under this T-period plan with all motivating wages back-loaded to retirement after a completely successful career, the responsibilities of a successful agent $I_t = \underline{I} / (\delta p_H)^{t-1}$ increase by a factor $1 / (\delta p_H)$ every period. That is $I_{t+1} / I_t = 1 / (\delta p_H)$. In a large economy, however, there may be many thousands of young agents in the new cohort that starts their career in any period, and we assume that each has a risk of failure that is independent of all the others. So among those agents in the original cohort who have not had any failure before period t in their careers, a fraction p_H should be expected to continue their successful careers through the next period t+1, while the others will fail and get no future responsibilities. So if we let J_t denote the aggregate investment managed by all agents in this cohort

in period t of their careers, then $J_{t+1}/J_t = p_H I_{t+1}/I_t = 1/\delta$. That is, if some very large number Ω of agents in the cohort begin managing I in the first period of their careers, then (by the law of large numbers) at any period t in their careers the cohort should be managing a total amount

$$J_t = (p_H)^{t-1} \Omega I_t = (p_H)^{t-1} \Omega \underline{I} / (\delta p_H)^{t-1} = \Omega \underline{I} / \delta^{t-1} = J_1 / \delta^{t-1}.$$

Let $W = \Omega w_{T+1} p_H^T$ denote the expected total wages that will be paid to successful members of this cohort when they retire. (There is no aggregate uncertainty about this amount, as the cohort consists of many agents with independent failure risks.) Then at each period t , we have

$$p_H \gamma J_t / (p_H - p_L) = p_H \gamma (p_H)^{t-1} \Omega \underline{I} / (p_H - p_L) = (p_H)^{t-1} \Omega w_{T+1} (\delta p_H)^{T+1-t} = W \delta^{T+1-t}.$$

That is, the expected moral-hazard rents that are associated with the investments managed by the cohort at the t 'th period of their careers are just equal to the current discounted value of the cohort's expected end-of-career rewards.

Now suppose that there are many competitive firms that are capable of hiring agents to produce this output under such T -period contracts. Then we may expect them to hire agents and increase production until the equilibrium price of output falls to π^* , or to some fluctuating sequence of prices that will average to π^* over any agent's career.

With the output price π^* , each firm just expects to break even on its T -period relationship with any agent. So at any point in time after the revenues from the first-period project have been realized, the firm's expected discounted value of net profits from future revenues and costs under the contract will be strictly negative. Thus, for any mid-career agent with a successful past record, there may be a temptation for the owners of the firm to falsely find fault with the agent's performance, so that he can be dismissed and replaced by a younger agent. So the π^* equilibrium here depends on the assumption that there exist two or more competitive firms which can make a credible commitment to pay back-loaded moral-hazard rents. In this economy, a firm's reputation for reliably judging and rewarding its responsible agents is an essential asset, without which the firm could not do business. More broadly, the productivity of this economy depends substantially on the social and legal structures that enable firms to commit to judging and rewarding their agents appropriately under long-term contracts.

Lagrangean analysis of the recursive problem at period t :

$$\begin{aligned} \mathcal{L}(I_t, u_{S,t}, u_{F,t}; \lambda_t, \mu_t) &= (\delta p_H \pi_{t+1} A - 1) I_t - \delta \lambda_{t+1} [p_H u_{S,t} + (1 - p_H) u_{F,t}] \\ &\quad + \lambda_t [\delta p_H u_{S,t} + \delta (1 - p_H) u_{F,t} - v_t] \\ &\quad + \mu_t [\delta p_H \delta u_{S,t} + \delta (1 - p_H) u_{F,t} - (\gamma I_t + \delta p_L u_{S,t} + \delta (1 - p_L) u_{F,t})] \end{aligned}$$

Optimality conditions for a solution with $\underline{I} < I_t < \bar{I}$, $u_{S,t} > 0$, $u_{F,t} = 0$:

$$\begin{aligned} \lambda_t &\geq 0 \quad \text{and} \quad \delta p_H u_{S,t} + \delta (1 - p_H) u_{F,t} - v_t \geq 0, \quad \text{with at least one equality;} \\ \mu_t &\geq 0 \quad \text{and} \quad p_H u_{S,t} + (1 - p_H) u_{F,t} - (\gamma I_t + p_L u_{S,t} + (1 - p_L) u_{F,t}) \geq 0, \quad \text{with at least one equality;} \\ 0 &= \partial \mathcal{L} / \partial I_t = (\delta p_H \pi_{t+1} A - 1) - \mu_t \gamma, \\ 0 &= \partial \mathcal{L} / \partial u_{S,t} = -\lambda_{t+1} \delta p_H + \lambda_t \delta p_H + \mu_t \delta (p_H - p_L), \\ 0 &\geq \partial \mathcal{L} / \partial u_{F,t} = -\lambda_{t+1} \delta (1 - p_H) + \lambda_t \delta (1 - p_H) - \mu_t \delta (p_H - p_L). \end{aligned}$$

The third condition implies $\mu_t = (\delta p_H \pi_{t+1} A - 1) / \gamma$.

The fourth condition implies $\lambda_t = \lambda_{t+1} - \mu_t (p_H - p_L) / p_H$ and

$$\partial \mathcal{L} / \partial u_{F,t} = -\lambda_{t+1} \delta (1 - p_H) + \lambda_t \delta (1 - p_H) - (\lambda_{t+1} - \lambda_t) p_H \delta = \delta (\lambda_t - \lambda_{t+1}).$$