

# Minimum Payments that Reward Honest Reputation Feedback

Radu Jurca  
Ecole Polytechnique Fédérale de Lausanne  
(EPFL)  
Artificial Intelligence Laboratory  
CH-1015 Lausanne, Switzerland  
radu.jurca@epfl.ch

Boi Faltings  
Ecole Polytechnique Fédérale de Lausanne  
(EPFL)  
Artificial Intelligence Laboratory  
CH-1015 Lausanne, Switzerland  
boi.faltings@epfl.ch

## ABSTRACT

Online reputation mechanisms need honest feedback to function effectively. Self interested agents report the truth only when explicit rewards offset the cost of reporting and the potential gains that can be obtained from lying. Side-payment schemes (monetary rewards for submitted feedback) can make truth-telling rational based on the correlation between the reports of different buyers.

In this paper we use the idea of automated mechanism design to construct the payments that minimize the budget required by an incentive-compatible reputation mechanism. Such payment schemes are defined by a linear optimization problem that can be solved efficiently in realistic settings. Furthermore, we investigate two directions for further lowering the cost of incentive-compatibility: using several reference reports to construct the side-payments, and filtering out reports that are probably false.

## Categories and Subject Descriptors

I.2.11 [Artificial Intelligence]: Distributed Artificial Intelligence

## General Terms

Algorithms, Design, Economics

## Keywords

honest feedback, reputation mechanisms, mechanism design

## 1. INTRODUCTION

Online buyers increasingly resort to reputation forums for obtaining information about the products or services they intend to purchase. The testimonies of previous buyers disclose hidden, *experience-related* [13], product attributes (e.g., quality, reliability, ease of use, etc.) that can only be

observed after the purchase. This previously unavailable information allows the buyers to take better, more efficient decisions.

Quality-based differentiation of products is also beneficial for the sellers. High quality, when recognizable by the buyers, brings higher revenues. Manufacturers can therefore optimally plan the investment in their products, such that the difference between the higher revenues of a better product, and the higher cost demanded by the improved quality, is maximized. Honest reputation feedback is thus essential for establishing an efficient market.

Human users exhibit high levels of honest behavior (and truthful sharing of information) without explicit incentives. However, in a future e-commerce environment dominated by rational agents, reputation mechanism designers need to make sure that sharing truthful information is in the best interest of the reporter.

Two factors make this task difficult. First, feedback reporting is usually costly. Most forums still require a conscious effort to formulate and submit feedback: buyers need to understand the rating scale (e.g., five star ratings – where one star is the lowest score, five star is the highest score – or “top five” preferences where one is the best score and five is the lowest), they need to manually fill in forms, and supervise the submission of the report. As feedback reporting does not bring direct benefits, many agents only report when they have ulterior motives, thus leading to a biased sample of reputation information.

Second, truth-telling is not always in the best interest of the reporter. In some settings, for instance, false denigration decreases the reputation of a product and allows the reporter to make a future purchase for a lower price. In other contexts, providers can offer monetary compensations in exchange for favorable feedback: e.g., doctors get gifts for recommending new drugs, authors ask their friends to write positive reviews about their latest book [6, 19]. One way or another, external benefits can be obtained from lying and selfish agents will exploit them.

Both problems can be addressed by a payment scheme that explicitly rewards honest feedback by a sufficient amount  $\Delta$  to offset both the cost of reporting and the gains that could be obtained through lying. Seminal work in the mechanism design literature [5, 4] shows that side payments can be designed to create the incentive for agents to report their private opinions truthfully, a property called *incentive compatibility*. The best such payment schemes have been constructed based on “proper scoring rules” [11, 8, 2], and ex-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

EC'06, June 11–15, 2006, Ann Arbor, Michigan, USA.  
Copyright 2006 ACM 1-59593-236-4/06/0006 ...\$5.00.

plot the correlation between the observations of different buyers about the same good.

Miller, Resnick and Zeckhauser (henceforth referred to as MRZ) [12] present a payment mechanism based on proper scoring rules that is particularly well suited for online feedback forums. In their mechanism, a central processing facility “scores” every submitted feedback by comparing it with another report (called the *reference* report) about the same good. The score does not reflect the agreement with the reference report; instead it *measures* the quality of the probability distribution for the reference report, induced by the submitted feedback. Payments directly proportional to these scores make honest reporting a Nash equilibrium. The payments can then be scaled so that in equilibrium, the return when reporting honestly is better by at least a margin  $\Delta$ . However, this scaling can lead to arbitrarily high feedback payments. This can be a problem because the payments cause a loss to the reputation mechanism that must be made up in some way, either by sponsorship or by charges levied on the users of the reputation information.

In this paper, we use the idea of *automated mechanism design* [3, 14] and compute optimal payments that minimize the budget required to achieve a certain margin  $\Delta$ . We thus lose the simplicity of a closed-form scoring rule, but gain in efficiency of the mechanism. Specifically, we derive the optimal payment scheme such that:

- given a required margin  $\Delta$  to offset reporting and honesty costs, the expected budget required for feedback payments is minimized; or, conversely,
- given certain budget constraints, the margin  $\Delta$  is maximized.

Using the framework for computing optimal feedback payment schemes, we then investigate two complementary methods that can be used to further decrease the cost of incentive-compatibility. The first requires the use of several reference reports to score feedback. We formally prove that the expected budget required by the reputation mechanism decreases with the number of employed reference reports. The second method adapts probabilistic filtering techniques to eliminate the reports that are probably false. We experimentally show that such filters are successful in decreasing the lying incentives, without greatly distorting the information provided by the reputation mechanism.

Section 2 formally introduces our setting. Section 3 describes the algorithm for computing the optimal payments, followed by an analysis of the computational complexity. In Section 4 we extend the mechanism by considering a) several reference reports, and b) a probabilistic filter to eliminate reports that are probably false. We conclude with a discussion and future work.

## 2. THE SETTING

Similar to [12], we consider an online market where a number of rational buyers (or “agents”) experience the same product or service. The quality of the product remains fixed, and defines the product’s (unknown) *type*.  $\Theta$  is the finite set of possible types, and  $\theta$  denotes a member of this set. We assume that all buyers share a common belief regarding the prior probability  $Pr[\theta]$ , that the product is of type  $\theta$ .  $\sum_{\theta \in \Theta} Pr[\theta] = 1$ .

After purchasing the product, the buyer perceives a noisy signal about the quality (i.e., true type) of the product.  $O^i$  denotes the random signal observed by the buyer  $i$ , and  $\mathcal{S} = \{s_1, s_2, \dots, s_M\}$  is the set of possible values for  $O^i$ . The observations of different buyers are conditionally independent, given the type of the product. Let  $f(s_j|\theta) = Pr[O^i = s_j|\theta]$  be the probability that a buyer observes the signal  $s_j$  when the true type of the product is  $\theta$ .  $f(\cdot|\cdot)$  is assumed common knowledge, and  $\sum_{j=1}^M f(s_j|\theta) = 1$  for all  $\theta \in \Theta$ .

A central reputation mechanism asks every buyer to submit feedback. Let  $a^i = (a_1^i, \dots, a_M^i)$  be the reporting strategy of buyer  $i$ , such that the buyer announces  $a_j^i \in \mathcal{S}$  whenever she observes the signal  $s_j$ . The honest reporting strategy is  $\bar{a} = (s_1, \dots, s_M)$ , when the buyer always declares the truth.

The reputation mechanism pays buyers for submitting feedback. The amount received by buyer  $i$  is computed by taking into account the signal announced by  $i$ , and the signal announced by another buyer,  $r(i)$ , called the *reference reporter* of  $i$ . Let  $\tau(a_j^i, a_k^{r(i)})$  be the payment received by  $i$  when she announces the signal  $a_j^i$  and the reference reporter announces the signal  $a_k^{r(i)}$ . The expected payment of the buyer  $i$  depends on the prior belief, on her observation  $O^i = s_j$ , and on the reporting strategies  $a^i$  and  $a^{r(i)}$ :

$$\begin{aligned} V(a^i, a^{r(i)}|s_j) &= E_{s_k \in \mathcal{S}} [\tau(a_j^i, a_k^{r(i)})] \\ &= \sum_{k=1}^M Pr[O^{r(i)} = s_k | O^i = s_j] \tau(a_j^i, a_k^{r(i)}); \end{aligned} \quad (1)$$

The conditional probability distribution for the signal observed by the buyer  $r(i)$  can be computed as:

$$Pr[s_k|s_j] = \sum_{\theta \in \Theta} f(s_k|\theta) Pr[\theta|s_j];$$

where  $Pr[\theta|s_j]$  is the posterior probability of the type  $\theta$  given the observation  $s_j$ , computed from Bayes’ Law:

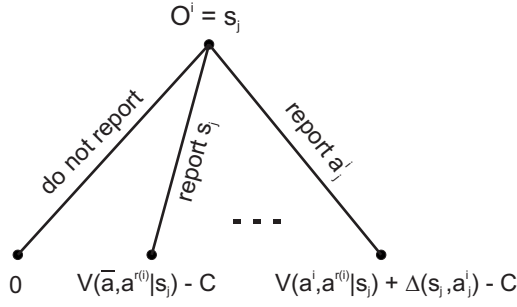
$$Pr[\theta|s_j] = \frac{f(s_j|\theta) Pr[\theta]}{Pr[s_j]}; \quad Pr[s_j] = \sum_{\theta \in \Theta} f(s_j|\theta) Pr[\theta];$$

Reporting is costly, and buyers can obtain external benefits from lying. Let  $C \geq 0$  be an upper bound for the feedback reporting cost of one buyer, and let  $\Delta(s_j, a_j^i)$  be an upper bound on the external benefit a buyer can obtain from falsely reporting the signal  $a_j^i$  instead of  $s_j$ . The cost of reporting  $C$  is assumed independent of the beliefs and observations of the buyer; moreover, for all signals  $s_j \neq s_k \in \mathcal{S}$ ,  $\Delta(s_j, s_j) = 0$  and  $\Delta(s_j, s_k) \geq 0$ .

## 3. COMPUTING THE OPTIMAL PAYMENT SCHEME

Let us consider the buyer  $i$  who purchases the product and observes the quality signal  $O^i = s_j$ . When asked by the reputation mechanism to submit feedback, the buyer can choose: (a) to honestly report  $s_j$ , (b) to report another signal  $a_j^i \neq s_j \in \mathcal{S}$  or (c) not to report at all. Figure 1 presents the buyer’s expected payoff for each of these cases, given the payment scheme  $\tau(\cdot, \cdot)$  and the reporting strategy  $a^{r(i)}$  of the reference reporter.

Truthful reporting is a Nash equilibrium (NEQ) if the buyer finds it optimal to announce the true signal, whenever



**Figure 1: Reporting feedback. Choices and Payoffs.**

the reference reporter also reports the truth. Formally, the honest reporting strategy  $\bar{a}$  is a NEQ if and only if for all signals  $s_j \in \mathcal{S}$ , and all reporting strategies  $a^* \neq \bar{a}$ :

$$\begin{aligned} V(\bar{a}, \bar{a}|s_j) &\geq V(a^*, \bar{a}|s_j) + \Delta(s_j, a_j^*); \\ V(\bar{a}, \bar{a}|s_j) &\geq C; \end{aligned}$$

When the inequalities are strict, honest reporting is a strict NEQ.

For any observed signal  $O^i = s_j \in \mathcal{S}$ , there are  $M - 1$  different dishonest reporting strategies  $a^* \neq \bar{a}$  the buyer can use: i.e., report  $a_j^* = s_h \in \mathcal{S} \setminus \{s_j\}$  instead of  $s_j$ . Using (1) to expand the expected payment of a buyer, the NEQ conditions become:

$$\begin{aligned} \sum_{k=1}^M Pr[s_k|s_j] (\tau(s_j, s_k) - \tau(s_h, s_k)) &> \Delta(s_j, s_h); \\ \sum_{k=1}^M Pr[s_k|s_j] \tau(s_j, s_k) &> C; \end{aligned} \quad (2)$$

for all  $s_j, s_h \in \mathcal{S}$ ,  $s_j \neq s_h$ .

Any payment scheme  $\tau(\cdot, \cdot)$  satisfying the conditions in (2) is incentive-compatible. MRZ [12] prove that such schemes exist.

Given the incentive-compatible payment scheme  $\tau(\cdot, \cdot)$ , the expected amount paid by the reputation mechanism to one buyer is:

$$W = E_{s_j \in \mathcal{S}} [V(\bar{a}, \bar{a}|s_j)] = \sum_{j=1}^M Pr[s_j] \left( \sum_{k=1}^M Pr[s_k|s_j] \tau(s_j, s_k) \right);$$

The optimal payment scheme minimizes the budget required by the reputation mechanism, and therefore solves the following linear program (i.e., linear optimization problem):

LP 1.

$$\begin{aligned} \min \quad & W = \sum_{j=1}^M Pr[s_j] \left( \sum_{k=1}^M Pr[s_k|s_j] \tau(s_j, s_k) \right) \\ \text{s.t.} \quad & \sum_{k=1}^M Pr[s_k|s_j] (\tau(s_j, s_k) - \tau(s_h, s_k)) > \Delta(s_j, s_h); \\ & \forall s_j, s_h \in \mathcal{S}, s_j \neq s_h; \\ & \sum_{k=1}^M Pr[s_k|s_j] \tau(s_j, s_k) > C; \quad \forall s_j \in \mathcal{S} \\ & \tau(s_j, s_k) \geq 0; \forall s_j, s_k \in \mathcal{S} \end{aligned}$$

The payment scheme  $\tau(\cdot, \cdot)$  solving LP 1 depends on the cost of reporting, on the external benefits from lying, and on the prior belief about the type of the product. To illustrate what these payments look like, the next subsection introduce a very simple example.

### 3.1 Example

Let us consider one buyer that needs the services of a plumber. The plumber can be either *Good* ( $G$ ) or *Bad* ( $B$ ): i.e.,  $\Theta = \{G, B\}$ . Since the plumber is listed on the Yellow Pages, the buyer believes that the plumber is probably good:  $Pr[G] = 0.8, Pr[B] = 0.2$ . However, even a good plumber can sometimes make mistakes and provide low quality service. Similarly, a bad plumber gets lucky from time to time and provides satisfactory service. Our buyer does not have the necessary expertise to judge the particular problem she is facing; she therefore perceives the result of the plumber's work as a random signal conditioned to the plumber's true type. We assume that the probability of a successful service (i.e., *high* quality) is 0.9 if the plumber is good, and 0.2 if the plumber is bad (the probabilities of a *low* quality service are 0.1 and 0.8 respectively). Following the notation in Section 2, we have:  $f(h|G) = 1 - f(l|G) = 0.9$  and  $f(h|B) = 1 - f(l|B) = 0.2$ .

Considering the prior belief, and the conditional distribution of quality signals, the buyer expects to receive high quality with probability:  $Pr[h] = 1 - Pr[l] = f(h|G)Pr[G] + f(h|B)Pr[B] = 0.76$ . After observing the plumber's work, the buyer updates her prior beliefs regarding the type of the plumber and can estimate the probability that the next buyer (i.e., the reference reporter) will get satisfactory service:  $Pr[h|h] = 1 - Pr[l|h] = 0.86$  and  $Pr[h|l] = 1 - Pr[l|l] = 0.43$ .

The buyer can submit one binary feedback (i.e.,  $l$  or  $h$ ) to an online reputation mechanism. Let the price of the plumber's work be fixed and normalized to 1, and the cost of formatting and submitting feedback be  $C = 0.01$ . The buyer has clear incentives to misreport:

- by reporting low quality when she actually received high quality, the buyer can hope to both decrease the price and increase the future availability of this (good) plumber. Assume that the external benefits of lying can be approximated as  $\Delta(h, l) = 0.06$
- by reporting high quality when she actually received low quality, the buyer can hope to decrease the relative reputation of other plumbers and thus obtain a faster (or cheaper) service from a better plumber in the future. Assume the lying incentive can be approximated as  $\Delta(l, h) = 0.02$

The optimal feedback payments solves the following problem:

	$\tau(h, h)$	$\tau(h, l)$	$\tau(l, h)$	$\tau(l, l)$	
min	0.65	0.11	0.10	0.14	
s.t.	0.86	0.14	-0.86	-0.14	> 0.06
	0.86	0.14			> 0.01
	-0.43	-0.57	0.43	0.57	> 0.02
			0.43	0.57	> 0.01
	$\geq 0$	$\geq 0$	$\geq 0$	$\geq 0$	

and are equal to:  $\tau(h, h) = 0.086$ ,  $\tau(l, l) = 0.1$ ,  $\tau(h, l) = \tau(l, h) = 0$ . The expected payment to a truth-telling buyer is 0.07 (i.e., 7% of the price of the service) for the reputation mechanism.

### 3.2 Unknown lying incentives

LP 1 reveals a strong correlation between the minimum expected cost and the external benefits obtained from lying: low lying incentives generate lower expected payments. When finding accurate approximations for the lying incentives is difficult, the mechanism designer might want to compute the payment scheme that satisfies certain budget constraints, and maximizes the tolerated misreporting incentives. The algorithm for computing these payments follows directly from LP 1: the objective function becomes a constraint (e.g., expected budget is bounded by some amount,  $\Gamma$ ) and the new objective is to maximize the worst case (i.e., minimum) expected payment loss caused by misreporting:

LP 2.

$$\begin{aligned}
 & \max \quad \Delta \\
 & \text{s.t.} \quad \sum_{j=1}^M Pr[s_j] \left( \sum_{k=1}^M Pr[s_k|s_j] \tau(s_j, s_k) \right) \leq \Gamma; \\
 & \quad \sum_{k=1}^M Pr[s_k|s_j] \left( \tau(s_j, s_k) - \tau(s_h, s_k) \right) > \Delta; \\
 & \quad \forall s_j, s_h \in \mathcal{S}, s_j \neq s_h; \\
 & \quad \sum_{k=1}^M Pr[s_k|s_j] \tau(s_j, s_k) > \Delta; \quad \forall s_j \in \mathcal{S} \\
 & \quad \tau(s_j, s_k) \geq 0; \forall s_j, s_k \in \mathcal{S}
 \end{aligned}$$

The resulting scheme guarantees that any buyer will report honestly when the reporting costs and external lying benefits are smaller than  $\Delta$ .

Coming back to the example in Section 3.1, let us assume that we cannot accurately approximate the benefits obtained from lying. Given the same limit on the expected budget (i.e.,  $\Gamma = 0.07$ ), we want to compute the payment scheme that maximizes the tolerance to lying. Solving LP 2 gives:  $\tau(h, h) = 0.077$ ,  $\tau(l, l) = 0.14$ ,  $\tau(h, l) = \tau(l, h) = 0$  and  $\Delta = 0.047$ .

### 3.3 Computational Complexity and Possible Approximations

The linear optimization problems LP 1 and LP 2 are similar in terms of size and complexity: LP 1 has  $M^2$  variables and  $M^2$  inequality constraints, LP 2 has  $M^2 + 1$  variables and  $M^2 + 1$  inequality constraints. We will therefore analyze the complexity (and runtime) of LP 1, and extend the conclusions to LP 2 as well.

The worst case complexity of linear optimization problems is  $O(n^4L)$ , where  $n = M^2$  is the number of variables, and  $L$  is the size of the problem (approximately equal to the total number of bits required to represent the problem). We experimentally evaluated the average time required to solve LP 1 by using the standard linear solver in the Optimization Toolbox of Matlab 7.0.4. For different sizes of the feedback set (i.e., different values of  $M$ ) we randomly generated 2000 settings, as described in Appendix A. Table 1 presents the average CPU time required to find the optimal payment scheme on an average laptop: e.g., 1.6 GHz Centrino processor, 1Gb RAM, WinXP operating system. Up to  $M = 16$  possible quality signals, general purpose hardware and software can find the optimal payment scheme in less than half a second.

M	CPU time [ms]	M	CPU time [ms]
2	11.16 ( $\sigma = 3.5$ )	10	92.79 ( $\sigma = 7.5$ )
4	19.24 ( $\sigma = 3.7$ )	12	174.81 ( $\sigma = 11.1$ )
6	29.22 ( $\sigma = 4.4$ )	14	316.63 ( $\sigma = 18.4$ )
8	55.62 ( $\sigma = 6.7$ )	16	521.47 ( $\sigma = 25.4$ )

**Table 1: Average CPU time (and standard deviation) for computing the optimal payment scheme.**

The optimal payment scheme depends on the prior belief regarding the type of the product, and therefore, must be recomputed after every submitted feedback. Although linear optimization algorithms are generally fast, frequent feedback reports could place unacceptable workloads on the reputation mechanism. Two solutions can be envisaged to ease the computational burden:

- publish batches of reports instead of individual ones. The beliefs of the buyers thus change only once for every batch, and new payments must be computed less frequently. The right size for the batch should be determined by considering the frequency of submitted reports and the tradeoff between computational cost, and the efficiency losses due to delayed information.
- approximate the optimal payments, either by closed form functions (e.g., scoring rules) or by partial solutions of the optimization problem. The rest of this section develops on these latter techniques.

The first approximation for the optimal incentive compatible payment scheme is provided by the MRZ mechanism [12]. They suggest the payment scheme:  $\tau(s_j, s_k) = R(s_k|s_j)$ , where  $R(\cdot)$  is a proper scoring rule. The three best known proper scoring rules are:

- the logarithmic scoring rule:

$$R(s_k|s_j) = \ln(Pr[s_k|s_j]);$$

- spherical scoring rule:

$$R(s_k|s_j) = \frac{Pr[s_k|s_j]}{\sqrt{\sum_{s_h \in \mathcal{S}} Pr[s_h|s_j]^2}};$$

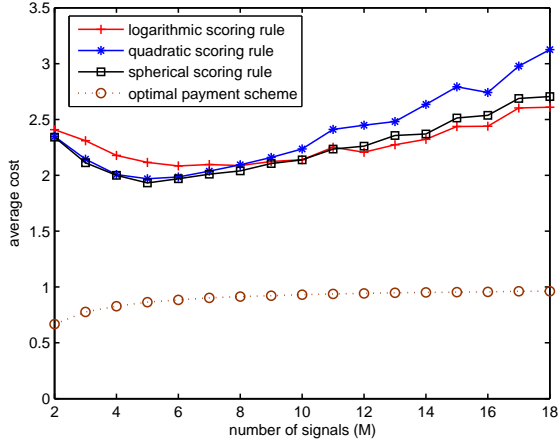
- quadratic scoring rule:

$$R(s_k|s_j) = 2Pr[s_k|s_j] - \sum_{s_h \in \mathcal{S}} Pr[s_h|s_j]^2;$$

The constraints from LP 1 can be satisfied by: (a) adding a constant to all payments such that they become positive: i.e.,  $\tau(s_j, s_k) = \tau(s_j, s_k) - \min_{s_h, s_l \in \mathcal{S}} \tau(s_h, s_l)$ , and (b) multiplying all payments with a constant such that the expected payment loss when lying outweighs external benefits: i.e.,  $\tau(s_j, s_k) = \alpha \cdot \tau(s_j, s_k)$  where:

$$\alpha = \max_{\substack{a_j^*, s_j \in \mathcal{S} \\ a_j^* \neq s_j}} \frac{\Delta(s_j, a_j^*)}{V(\bar{a}, \bar{a}|s_j) - V(a^*, \bar{a}|s_j)}; \quad (3)$$

For the example in Section 3.1, the payments computed based on scoring rules (properly scaled according to (3)) are the following:



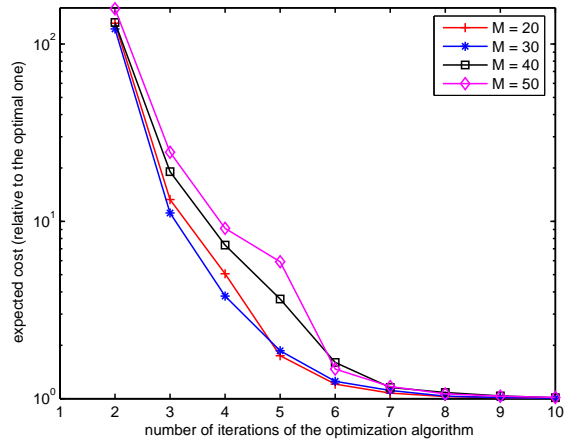
**Figure 2: Incentive-compatible payments based on proper scoring rules.**

- $\tau_l(h, h) = 0.27$ ,  $\tau_l(h, l) = 0$ ,  $\tau_l(l, h) = 0.17$ ,  $\tau_l(l, l) = 0.21$  and an expected cost of 0.22 for the logarithmic scoring rule;
- $\tau_s(h, h) = 0.2$ ,  $\tau_s(h, l) = 0$ ,  $\tau_s(l, h) = 0.11$ ,  $\tau_s(l, l) = 0.15$  and an expected cost of 0.17 for the spherical scoring rule;
- $\tau_q(h, h) = 0.23$ ,  $\tau_q(h, l) = 0$ ,  $\tau_q(l, h) = 0.13$ ,  $\tau_q(l, l) = 0.18$  and an expected cost of 0.19 for the quadratical scoring rule;

The payments based on scoring rules are two to three times more expensive than the optimal ones. The same ratio remains valid for more general settings. We investigated 2000 randomly generated settings (see Appendix A) for different number of quality signals. Figure 2 plots the average expected payment to one buyer when payments are computed using scoring rules.

Computational methods can also be used to obtain faster approximations of the optimal payment scheme. Most linear programming algorithms find the optimal solution by iterating through a set of feasible points that monotonically converge to the optimal one. Such algorithms are *just in time* algorithms as they can be stopped at any time, and provide a feasible solution (i.e., a payment scheme that is incentive-compatible, but maybe not optimal). The more time available, the better the feasible solution. The reputation mechanism can thus set a deadline for the optimization algorithm, and the resulting payment scheme makes it optimal for the buyers to report the truth.

Figure 3 plots the convergence of the Matlab linear programming algorithm for large problems (i.e., large number of signals) where approximations are likely to be needed. For 500 randomly generated settings, we plot (on a logarithmic scale) the average relative cost (relative to the optimal one) of the partial solution available after  $t$  iteration steps of the algorithm. As it can be seen, most of the computation time is spent making marginal improvements to the partial solution. For  $M = 50$  quality signals, the full optimization takes 20 steps on the average. However, the partial solution after



**Figure 3: Incentive-compatible payments based on partial solutions.**

6 steps generates expected costs that are only 40% higher on the average than the optimal ones.

Finally, the two techniques can be combined to obtain fast accurate approximations. As many linear programming algorithms accept initial solutions, one could use the scoring rules approximations to specify a starting point for an iterative optimization algorithm.

## 4. FURTHER LOWERING THE FEEDBACK PAYMENTS

The payment scheme computed in Section 3 generates, by definition, the lowest expected budget required by an incentive-compatible reputation mechanism, for a given setting. In this section we investigate two modifications of the mechanism itself, in order to lower the cost of reputation management even further. The first, proposes the use of several reference reporters for scoring one feedback. We formally prove that the higher the number of reference raters, the lower becomes the expected cost of reputation management.

The second idea is to reduce the potential lying incentives by filtering out false reports. Intuitively, the false feedback reports that bring important external benefits must significantly differ from the average reports submitted by honest reporters. Using a probabilistic filter that detects and ignores “abnormal” reports, lying benefits can be substantially reduced. The constraints on the optimal payments thus become more relaxed, and the optimal expected cost decreases.

### 4.1 Using several reference raters

The setting described in Section 2 is modified in the following way: we consider  $N$  (instead of only one) reference reports when computing the feedback payment due to an agent. By an abuse of notation we use the same  $r(i)$  to denote the set of  $N$  reference reporters of agent  $i$ . Let  $a^{r(i)} = (a^j)_{j \in r(i)}$  denote the vector of reporting strategies of the agents in  $r(i)$ , and let  $a_k^{r(i)}$  be a set of submitted reports. The set of possible values of  $a_k^{r(i)}$  is  $\mathcal{S}(N)$ .

As the signals observed by the agents are independent, the order in which the reference reports were submitted is not relevant. We take  $\mathcal{S}(N)$  to be the set of all unordered sequences of reports of length  $N$ .  $a_k^{r(i)}$  can be represented by a vector  $(n_1, \dots, n_M)$ , where  $n_j$  is the number of reference reporters announcing the signal  $s_j$ .  $\mathcal{S}(N)$  thus becomes:

$$\mathcal{S}(N) = \left\{ (n_1, \dots, n_M) \in \mathbb{N}^M \mid \sum_{j=1}^M n_j = N \right\};$$

The expected payment of agent  $i$  is:

$$V(a^i, a^{r(i)}|s_j) = \sum_{a_k^{r(i)} \in \mathcal{S}(N)} Pr[a_k^{r(i)}|s_j] \tau(a_j^i, a_k^{r(i)});$$

and the optimal payment scheme  $\tau(\cdot, \cdot)$  solves:

LP 3.

$$\begin{aligned} \min \quad & \sum_{j=1}^M Pr[s_j] \left( \sum_{a_k \in \mathcal{S}(N)} Pr[a_k|s_j] \tau(s_j, a_k) \right); \\ \text{s.t.} \quad & \sum_{a_k \in \mathcal{S}(N)} Pr[a_k|s_j] \left( \tau(s_j, a_k) - \tau(s_h, a_k) \right) > \Delta(s_j, s_h); \\ & \forall s_j, s_h \in \mathcal{S}, s_j \neq s_h, \\ & \sum_{a_k \in \mathcal{S}(N)} Pr[a_k|s_j] \tau(s_j, a_k) > C; \quad \forall s_j \in \mathcal{S} \\ & \tau(s_j, a_k) \geq 0; \forall s_j \in \mathcal{S}, a_k \in \mathcal{S}(N) \end{aligned}$$

The optimization problem LP 3 has  $M^2$  constraints and  $M \cdot |\mathcal{S}(N)|$  variables, where  $|\mathcal{S}(N)| = \binom{M-1}{N+M-1}$  is the cardinality of  $\mathcal{S}(N)$

**PROPOSITION 1.** *The minimum budget required by an incentive compatible reputation mechanism decreases as the number of reference reporters increases.*

**PROOF.** The proof is based on the observation that the number of constraints in the optimization problem LP 3 does not depend on the number  $N$  of reference reporters. Therefore, the number of variables of the dual of LP 3 does not depend on  $N$ . We define a sequence of primal and dual optimization problems,  $LP(N)$  and  $DP(N)$  respectively, characterizing the setting where  $N$  reference reports are considered. We show that any feasible solution of  $DP(N+1)$ , is also feasible in  $DP(N)$ .  $DP(N)$  is therefore “less constrained” than  $DP(N+1)$ , and consequently will have a higher maximal cost. From the Duality Theorem, it follows that the expected cost of the payment scheme defined by  $LP(N)$  is higher than the expected cost of the payments defined by  $LP(N+1)$ . Thus, the budget required by an incentive compatible reputation mechanism decreases as the number of reference reporters increases.

Formally, we associate the dual variables  $y_j^h$  and  $y_j^j$  respectively, to the constraints:

$$\begin{aligned} \sum_{a_k \in \mathcal{S}(N)} Pr[a_k|s_j] \left( \tau(s_j, a_k) - \tau(s_h, a_k) \right) > \Delta(s_j, s_h); \\ \sum_{a_k \in \mathcal{S}(N)} Pr[a_k|s_j] \tau(s_j, a_k) > C; \end{aligned}$$

The dual problem  $DP(N)$  thus becomes:

$$\begin{aligned} \max \quad & \sum_{j=1}^M \left( C \cdot y_j^j + \sum_{h=1}^M \Delta(s_j, s_h) \cdot y_j^h \right); \\ \text{s.t.} \quad & \sum_{h=1}^M y_h^m Pr[a_k|s_m] - \sum_{\substack{h=1 \\ h \neq m}}^M y_h^m Pr[a_k|s_h] < Pr[s_m] Pr[a_k|s_m]; \\ & \forall s_m \in \mathcal{S}, a_k \in \mathcal{S}(N) \\ & y_j^h \geq 0; \forall j, h \in \{1, \dots, M\} \end{aligned}$$

Let  $y$  be a feasible solution of  $DP(N+1)$ . For any  $s_m \in \mathcal{S}$ , let  $s_j = \arg \min_{s \in \mathcal{S}} Pr[s|s_m]$ . For any  $a_k = (n_1, \dots, n_M) \in \mathcal{S}(N)$ , it is possible to find  $a_k^* = (n_1, \dots, n_j + 1, \dots, n_M) \in \mathcal{S}(N+1)$  such that  $N$  reference reporters announce  $a_k$ , and the remaining one reports  $s_j$ . For all  $s_h \in \mathcal{S}$ , we have:

$$\begin{aligned} Pr[a_k|s_h] &= N! \prod_{k=1}^M \frac{Pr[s_k|s_h]^{n_k}}{n_k!}; \\ Pr[a_k^*|s_h] &= Pr[s_j|s_h] Pr[a_k|s_h] \frac{N+1}{n_j+1}; \end{aligned}$$

$y$  is a feasible solution of  $DP(N+1)$ , therefore:

$$Pr[s_m] Pr[a_k^*|s_m] > \sum_{h=1}^M y_h^m Pr[a_k^*|s_m] - \sum_{\substack{h=1 \\ h \neq m}}^M y_h^m Pr[a_k^*|s_h];$$

Because:

$$\begin{aligned} Pr[a_k^*|s_m] &= \frac{N+1}{n_j+1} Pr[s_j|s_m] Pr[a_k|s_m]; \\ Pr[a_k^*|s_h] &= \frac{N+1}{n_j+1} Pr[s_j|s_h] Pr[a_k|s_h] \\ &\leq \frac{N+1}{n_j+1} Pr[s_j|s_m] Pr[a_k|s_h]; \end{aligned}$$

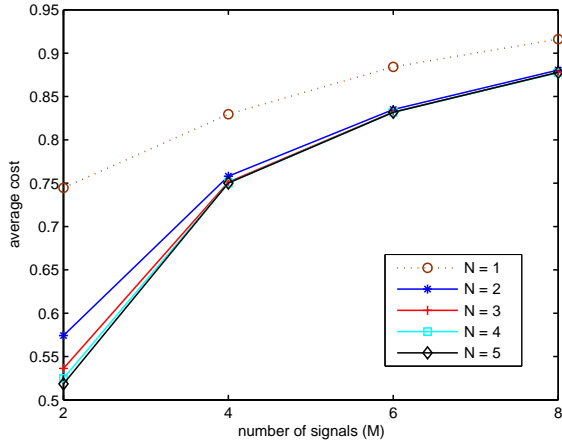
for all  $s_h \neq s_m$ ,  $y$  also satisfies:

$$Pr[s_m] Pr[a_k|s_m] > \sum_{h=1}^M y_h^m Pr[a_k|s_m] - \sum_{\substack{h=1 \\ h \neq m}}^M y_h^m Pr[a_k|s_h];$$

and is feasible in  $DP(N)$ . The “cost” of  $DP(N)$  is therefore greater or equal to the cost of  $DP(N+1)$ ; consequently, the budget required by a reputation mechanism using  $N$  reference reports is higher or equal to the budget required by a mechanism using  $N+1$  reference reports. ■

Using several reference reports decreases the cost of reputation management, but also increases the complexity of the algorithm defining the optimal payment scheme. We experimentally studied the quantitative effect of several reference reports on the budget of the reputation mechanism. For 2000 randomly generated settings (details available in Appendix A) Figure 4 plots the average cost as the number of reference reports is increased from 1 to 5. Significant savings (approx. 25% for a setting with  $M=2$  quality signals, and 4% for a setting with  $M=8$  quality signals) are mainly obtained from the second and third reference reports. As a good tradeoff between cost and computational complexity, practical systems can therefore use between 2 and 4 reference reports, depending on the number of quality signals in the set  $\mathcal{S}$ .

For the example in Section 3.1, using 2 reference reports gives the payment scheme:  $\tau_2(l, ll) = 0.11$ ,  $\tau_2(h, hh) =$



**Figure 4: Average expected payment to one agent when several reference reports are used.**

0.083, the rest of payments equal to 0, and an expected cost of 0.055. Using 3 reference reports gives the payment scheme:  $\tau_3(l, ll) = 0.15$ ,  $\tau_3(h, hhh) = 0.094$ , the rest of payments equal to 0, and an expected cost of 0.052.

## 4.2 Filtering out false reports

The feedback payments naturally decrease when the reporting and honesty costs become smaller. The cost of reporting can be decreased by software tools that help automate as much as possible the process of formatting and submitting feedback. On the other hand, the external incentives for lying can be reduced by filtering out the reports that are likely to be false.

“Truth filters” can be constructed based on statistical analysis. When all agents report truthfully, their reports follow a common distribution given by the product’s true type. Reports that stand out from the common distribution are either particularly unlikely, or dishonest. Either way, by filtering them out with high probability, the reputation information does not usually suffer significant degradation.

Probabilistic filters of false reports have been widely used in decentralized and multi-agent systems. Vu, Hauswirth and Aberer [16] use clustering techniques to isolate lying agents in a market of web-services. Their mechanism is based on a small number of “trusted” reports that provide the baseline for truthful information. The technique shows very good experimental results when lying agents use probabilistic strategies and submit several reports. Buchegger and Le Boudec [1] use Bayesian methods to detect free riders in a wireless ad-hoc network. Nodes consider both direct and second-hand information, however, second-hand information is taken into account only when it does not conflict with direct observations (i.e., second hand reports do not trigger significant deviations in the agent’s beliefs). In peer-to-peer reputation mechanisms (e.g., TRAVOS [15], CRE-DENCE [17] and [20, 7, 18]) agents weigh the evidence from peers by the distance from the agent’s direct experience.

However, all of the above cited results rely on two important assumptions: a) every agent submits several reports, b) according to some *probabilistic* lying strategy. Self-

interested agents can strategically manipulate their reports to circumvent the filtering mechanisms and take profit from dishonest reporting<sup>1</sup>. When all buyers are self-interested and submit only one feedback report, filtering methods based entirely on similarity metrics can never be accurate enough to filter out effectively *all* lying strategies without important losses of information.

In this section, we present an alternative filtering method that also exploits the information available to the agents. The intuition behind our method is simple: the probability of filtering out the report  $a_j^i$  submitted by agent  $i$  should not only depend on how well  $a_j^i$  fits the distribution of peer reports, but also on the benefits that  $a_j^i$  could bring to the reporter if it were false. When  $\Delta(s_j, a_j^i)$  is big (i.e. the agent has strong incentives to report  $a_j^i$  whenever her true observation was  $s_j$ ), the filtering mechanism should be more strict in accepting  $a_j^i$  given that peer reports make the observation of  $s_j$  probable. On the other hand, when  $\Delta(s_j, a_j^i)$  is small, filtering rules can be more relaxed, such that the mechanism does not lose too much information. In this way, the filter adapts to the particular context and allows an optimal tradeoff between diminished costs and loss of information.

Concretely, let  $Pr(\theta)$ ,  $\theta \in \Theta$  describe the current common belief regarding the true type of the product, let  $s_j, a_j^i \in \mathcal{S}$  be the signals observed, respectively announced by agent  $i$ , and let  $a_k \in \mathcal{S}(N)$  describe the set of  $N$  reference reports. The publishing of the report submitted by agent  $i$  is delayed until the next  $\hat{N}$  reports (i.e., the filtering reports) are also available. A filtering mechanism is formally defined by the table of probabilities  $\pi(a_j^i, \hat{a}_k)$  of accepting the report  $a_j^i \in \mathcal{S}$  when the filtering reports take the value  $\hat{a}_k \in \mathcal{S}(\hat{N})$ . With probability  $1 - \pi(a_j^i, \hat{a}_k)$  the report  $a_j^i$  will not be published by the reputation mechanism, and therefore not reflected in the reputation information. Note, however, that all reports (including dropped ones) are paid for as described in the previous sections.

The payment scheme  $\tau(\cdot, \cdot)$  and the filtering mechanism  $\pi(\cdot, \cdot)$  are incentive compatible if and only if for all signals  $s_j, s_h \in \mathcal{S}$ ,  $s_j \neq s_h$ , the expected payment loss offsets the expected gain obtained from lying:

$$\sum_{a_k \in \mathcal{S}(N)} Pr[a_k | s_j] (\tau(s_j, a_k) - \tau(s_h, a_k)) > \hat{\Delta}(s_j, s_h) \quad (4)$$

$$\hat{\Delta}(s_j, s_h) = \sum_{\hat{a}_k \in \mathcal{S}(\hat{N})} Pr[\hat{a}_k | s_j] \cdot \pi(s_h, \hat{a}_k) \cdot \Delta(s_j, s_h);$$

where  $\hat{\Delta}(\cdot, \cdot)$  is obtained by discounting  $\Delta(\cdot, \cdot)$  with the expected probability that a false report is recorded by the reputation mechanism.

Naturally, the feedback payments decrease with decreasing probabilities of accepting reports. However, a useful reputation mechanism must also limit the loss of information. As a metric for information loss we chose the number (or percentage) of *useful* reports that are dropped by the mechanism. A feedback report is *useful*, when given the true type of the product and a prior belief on the set of possible types, the posterior belief updated with the report is closer to the true type than the prior belief.

<sup>1</sup>it is true, however, that some mechanisms exhibit high degrees of robustness towards such lying strategies: individual agents can profit from lying, but as long as the big majority of agents reports honestly, the liars do not break the properties of the reputation mechanism



For the example in Section 3.1, when the plumber is actually good, recording a *high* quality report is *useful* (because the posterior belief is closer to reality than the prior belief), while recording a *low* quality report is not. Conversely, when the plumber is bad, recording a *low* quality report is useful, while recording a *high* quality report is not. The notion of usefulness captures the intuition that some reports can be filtered out in some contexts without any loss of information for the buyers (on the contrary, the community has more accurate information without the report).

Formally, information loss can be quantified in the following way. Given the true type  $\theta^* \in \Theta$  and the prior belief  $Pr(\cdot)$  on the set of possible types, the report  $s_j$  is useful if and only if  $Pr(\theta^* | s_j) < Pr(\theta^*)$ : i.e. the posterior belief updated with the signal  $s_j$  is closer to the true type than the prior belief. Given the filtering mechanism  $\pi(\cdot, \cdot)$ , and the true type  $\theta^*$ , the expected probability of dropping  $s_j$  is:

$$Pr[\text{drop } s_j | \theta^*] = 1 - \sum_{\hat{a}_k \in \mathcal{S}(\hat{N})} Pr[\hat{a}_k | \theta^*] \pi(s_j, \hat{a}_k); \quad (5)$$

where  $Pr[\hat{a}_k | \theta^*]$  is the probability that the filtering reports take the value  $\hat{a}_k$ , when the true type of the product is  $\theta^*$ . To limit the loss of information, the reputation mechanism must insure that given the current belief, whatever the true type of the product, no useful report is dropped with a probability greater than a given threshold,  $\gamma$ :

$$\forall s_j \in \mathcal{S}, \theta \in \Theta, \quad Pr[\theta] < Pr[\theta | s_j] \Rightarrow Pr[\text{drop } s_j | \theta] < \gamma; \quad (6)$$

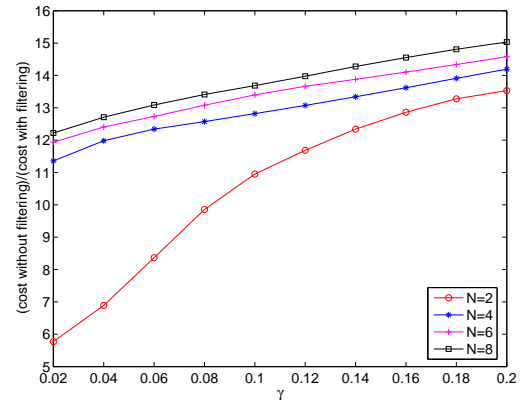
We can now define the incentive-compatible payment mechanism (using  $N$  reference reports) and filtering mechanism (using  $\hat{N}$  filtering reports) that minimize the expected cost:

LP 4.

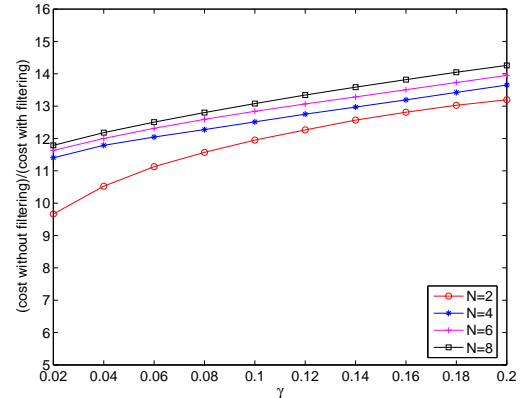
$$\begin{aligned} \min \quad & \sum_{j=1}^M Pr[s_j] \left( \sum_{a_k \in \mathcal{S}(N)} Pr[a_k | s_j] \tau(s_j, a_k) \right); \\ \text{s.t.} \quad & \sum_{a_k \in \mathcal{S}(N)} Pr[a_k | s_j] \left( \tau(s_j, a_k) - \tau(s_h, a_k) \right) > \hat{\Delta}(s_j, s_h); \\ & \forall s_j, s_h \in \mathcal{S}, s_j \neq s_h, \\ & \hat{\Delta}(s_j, s_h) \text{ is defined in (4)}; \\ & \sum_{a_k \in \mathcal{S}(N)} Pr[a_k | s_j] \tau(s_j, a_k) > C; \quad \forall s_j \in \mathcal{S} \\ & Pr[\theta] < Pr[\theta | s_j] \Rightarrow Pr[\text{drop } s_j | \theta] < \gamma; \quad \forall \theta, \forall s_j \\ & \tau(s_j, a_k) \geq 0, \pi(s_j, \hat{a}_k) \in [0, 1] \quad \forall s_j, \forall a_k, \forall \hat{a}_k; \end{aligned}$$

The effect of using probabilistic filtering of reports was experimentally studied on 500 randomly generated settings, for different number of filtering reports (i.e.,  $\hat{N}$ ), different number of quality signals (i.e.,  $M$ ) and different values of the parameter  $\gamma$ . Figure 5 plots the tradeoff between cost reduction (i.e. the ratio between the optimal cost without probabilistic filtering and the optimal cost with probabilistic filtering) and information loss for  $M = 3$  and  $M = 5$  quality signals. When  $M = 3$ , and we accept to lose 2% of the useful reports, the cost decreases 6 times by using  $\hat{N} = 2$  filtering reports, and 12 times by using  $\hat{N} = 8$  filtering reports. As intuitively expected, the cost decreases when we can use more filtering reports, and accept higher probabilities of losing useful feedback.

As a next experiment, we study the accuracy of the reputation information published by a mechanism that filters out



(a)  $M = 3$ , using  $N$  filtering reports;



(b)  $M = 5$ , using  $N$  filtering reports;

**Figure 5: Tradeoff between cost decrease and information loss.**

reports. For each of the random settings generated above, we also generate 200 random sequences of 20 feedback reports corresponding to a randomly chosen type. For different parameters (i.e., number of signals,  $M$ , number of filtering reports,  $\hat{N}$ , and threshold probability  $\gamma$ ), Figure 6 plots the mean square error of the reputation information<sup>2</sup> published by a mechanism that filters, respectively doesn't filter submitted reports. As expected, filtering out reports does not significantly alter the convergence of beliefs; on the contrary, filtering out reports may sometimes help to focus the beliefs on the true type of the product.

Finally, we illustrate the use of filtering mechanisms on the example from Section 3.1. Using one reference report, 3 filtering reports, and a threshold for dropping useful reports of 2%, gives the following mechanism: the payments:  $\tau(h, h) = 0.028$ ,  $\tau(h, l) = \tau(l, h) = 0$ ,  $\tau(l, l) = 0.04$  and the filtering probabilities:  $\pi(h, hhh) = 1 = \pi(h, hhl) = \pi(l, hll) = \pi(l, llh)$ ,  $\pi(l, llh) = \pi(l, llh) = 0.87$  and  $\pi(h, llh) = 0.3$ . The expected payment to one agent is 0.02

<sup>2</sup>The mean square error after  $i$  submitted reports is defined as:  $\epsilon_i = \sum_{\theta \in \Theta} (Pr[\theta | i] - I(\theta))^2$ , where  $Pr[\cdot | i]$  describes the belief of the agents regarding the type of the product after  $i$  submitted reports,  $I(\theta) = 1$  for  $\theta = \theta^*$  (the true type of the product), and  $I(\theta) = 0$  for  $\theta \neq \theta^*$ .



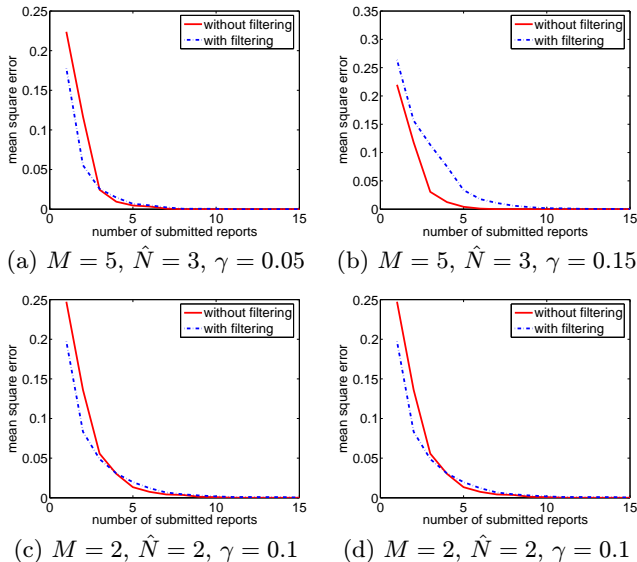


Figure 6: Convergence of reputation information.

## 5. DISCUSSION

The framework described in this paper can be used to address the collusion between clients and product manufacturers. The collusion can happen when the benefit obtained by providers from a false report offsets the payment returned to the client in exchange for lying. The feedback payments presented in this paper make sure that no provider can afford to buy false reports from rational clients.

Providers could, however, create fake buyer identities, (or “bribe” real buyers that never purchased the product) in order to bias reputation information. The problem can be addressed by security mechanisms that connect feedback reports to real transaction IDs. For example, a site like Expedia.com can make sure that no client leaves feedback about a hotel without actually paying for a room. Hotels could, of course, create fake bookings, but the repeated payment of Expedia commission fees makes the manipulation of information very expensive. On the other hand, social norms and legislation (i.e., providers that try to bribe clients risk being excluded from the market) could further avoid provider side manipulation.

One of the major assumptions behind the payment mechanism is that all buyers are risk-neutral. Honest reporting is rational because, in expectation, brings higher revenues. However, risk-averse buyers prefer the “sure” benefit from lying to the probabilistic feedback payment, and thus misreport. Fortunately, our mechanism can be adapted to any risk-model of the buyers, by modifying (1) to reflect the real expected payment. Of course, highly risk-averse buyers require significantly higher payments.

The motivation to minimize feedback payments might not be clear when the budget of the mechanism is covered by buyer subscription fees. Higher fees will be matched (in expectation) by higher feedback payments. However, this holds only for risk-neutral buyers. Real users (probably risk-averse) will regard the fixed fees as more expensive than the revenue expected from honest reporting; and higher subscription fees are increasingly more expensive. Moreover, no

real-world system that we know of, charges users for reputation information. Introducing reputation information fees could seriously deter participation.

An interesting question is what happens as more and more reports are recorded by the reputation mechanism. When the type of the product does not change, the beliefs of the buyers rapidly converge towards the true type (Figure 6). As more information becomes available to buyers, the private quality signal they observe triggers smaller and smaller changes of the prior belief. As a consequence, the probability distributions for the reference reports conditional on the private observation (i.e.,  $Pr[a_k|O^i]$ ) become closer, and the payments needed to guarantee a minimal expected loss from lying increase. Fortunately, external benefits from lying also decrease: the effect of a false report on reputation information tends to 0. Depending on the particular context, lying incentives decrease faster, respectively slower than the distance between the conditional probability distributions for the reference reports, and thus, truth-telling becomes easier respectively harder to guarantee. Whatever the case, it makes sense to stop collecting feedback when the beliefs of the buyers are sufficiently precise

In real settings, however, the true type of a product or service does actually change in time: e.g., initial bugs get eliminated, the technology improves, etc. Existing incentive compatible payment scheme rely on the fact that future buyers obtain exactly the same thing as the present ones. Depending on the time horizon, this assumption might not hold. Designing payment schemes that take into account the timely variations of the product’s type remains a challenge for future work.

The honest reporting Nash Equilibrium is unfortunately not unique. Other lying equilibria exist, and some of them generate higher expected payoffs for reporters than the truthful one. In a previous result, [9] we show that a small number of *trusted reports* (i.e., feedback reports that are true with high probability) can eliminate (or render unattractive) lying Nash equilibria. As future work, we intend to extend the presented framework to also account for multiple equilibria.

Collusion between buyers remains a problem for this class of incentive-compatible mechanisms. As explained by MRZ, buyers can synchronize their possibly false reports in order to increase their revenue. Choosing randomly the reference report for every submitted feedback can help eliminate small coalitions: only large coalitions are rational, when the probability of having a reference report from the same coalition is high enough. Another safeguard against reporting coalitions is to use trusted reports. In some settings [10], a small number of trusted reports can make collusion irrational.

Last, but not least, the mechanism we have presented might not be incentive compatible when buyers possess private information about the true type of the product. As future work we plan to relax the common knowledge requirement on the prior belief of the buyers, and extend the current framework to work for a range of acceptable prior beliefs.

## 6. CONCLUSION

Honest feedback is essential for the effectiveness of online reputation mechanisms. When feedback reporters are self-interested, explicit payments can make truthful reporting rational. Most of the existing incentive-compatible payment schemes are constructed based on proper scoring rules.

In this paper we use the idea of automated mechanism design to construct payment schemes that offset both the cost of reporting and the external gains an agent could obtain from lying. We show how a linear optimization problem can define the optimal payments that minimize the expected budget required by an incentive-compatible reputation mechanism. Experiments show that such payments can efficiently be computed by existing optimization algorithms.

We also investigate two methods that can further decrease incentive-compatible feedback payments. The first requires the use of several reference reports. We prove that higher numbers of reference reports lead to lower costs; however, experiments show that little benefit can be obtained by using more than 4 reference reports.

Finally, we show how probabilistic filtering mechanism can be used to filter out some of the reports that are probably false. By considering both the information available to the agents, and the similarity between peer reports, we were able to derive the filtering mechanism that significantly reduces the lying incentives while bounding the loss of information. The cost of incentive-compatible mechanisms that use such filters can thus be lowered by one order of magnitude.

## 7. ACKNOWLEDGMENTS

We thank Jean-Cédric Chappelier, Dan Jurca and the anonymous reviewers for helpful comments and suggestions.

## 8. REFERENCES

- [1] S. Buchegger and J.-Y. Le Boudec. Self-Policing Mobile Ad-Hoc Networks by Reputation. *IEEE Communication Magazine*.
- [2] R. T. Clemen. Incentive contracts and strictly proper scoring rules. *Test*, 11:167–189, 2002.
- [3] V. Conitzer and T. Sandholm. Complexity of mechanism design. In *Proceedings of the Uncertainty in Artificial Intelligence Conference (UAI)*, 2002.
- [4] J. Crémer and R. P. McLean. Optimal Selling Strategies under Uncertainty for a Discriminating Monopolist When Demands Are Interdependent. *Econometrica*, 53(2):345–61, 1985.
- [5] C. d’Aspremont and L.-A. Gauthier. Incentives and Incomplete Information. *Journal of Public Economics*, 11:25–45, 1979.
- [6] A. Harmon. Amazon Glitch Unmasks War of Reviewers. *The New York Times*, February 14, 2004.
- [7] R. Ismail and A. Jøsang. The Beta Reputation System. In *Proceedings of the 15th Bled Conf. on E-Commerce*, 2002.
- [8] S. Johnson, J. Pratt, and R. Zeckhauser. Efficiency Despite Mutually Payoff-Relevant Private Information: The Finite Case. *Econometrica*, 58:873–900, 1990.
- [9] R. Jurca and B. Faltings. Enforcing Truthful Strategies in Incentive Compatible Reputation Mechanisms. In *Internet and Network Economics*, volume 3828 of *LNCS*, pages 268 – 277. 2005.
- [10] R. Jurca and B. Faltings. Reputation-based Service Level Agreements for Web Services. In *Service Oriented Computing (ICSOC - 2005)*, volume 3826 of *LNCS*, pages 396 – 409. 2005.
- [11] M. Kandori and H. Matsushima. Private observation, communication and collusion. *Econometrica*, 66(3):627–652, 1998.
- [12] N. Miller, P. Resnick, and R. Zeckhauser. Eliciting Informative Feedback: The Peer-Prediction Method. Forthcoming in *Management Science*, 2005.
- [13] A. Parasuraman, V. Zeithaml, and L. Berry. A Conceptual Model of Service Quality and Its Implications for Future Research. *Journal of Marketing*, 49:41–50, 1985.
- [14] T. Sandholm. Automated mechanism design: A New Application Area for Search Algorithms. In *Proceedings of the International Conference on Principles and Practice of Constraint Programming*, 2003.
- [15] L. Teacy, J. Patel, N. Jennings, and M. Luck. Coping with Inaccurate Reputation Sources: Experimental Analysis of a Probabilistic Trust Model. In *Proceedings of AAMAS*, Utrecht, The Netherlands, 2005.
- [16] L.-H. Vu, M. Hauswirth, and K. Aberer. QoS-based Service Selection and Ranking with Trust and Reputation Management. In *Proceedings of the International Conference on Cooperative Information Systems (CoopIS 2005)*, 2005.
- [17] K. Walsh and E. Sirer. Fighting Peer-to-Peer SPAM and Decoys with Object Reputation. In *Proceedings of P2PECON*, Philadelphia, USA, 2005.
- [18] A. Whitby, A. Jøsang, and J. Indulska. Filtering out Unfair Ratings in Bayesian Reputation Systems. In *Proceedings of the 7th Intl. Workshop on Trust in Agent Societies*, 2004.
- [19] E. White. Chatting a Singer Up the Pop Charts. *The Wall Street Journal*, October 15, 1999.
- [20] B. Yu and M. Singh. Detecting Deception in Reputation Management. In *Proceedings of the AAMAS*, Melbourne, Australia, 2003.

## APPENDIX

### A. GENERATING RANDOM SETTINGS

We consider settings where  $M$  possible product types are each characterized by one quality signal: i.e., the sets  $\mathcal{S}$  and  $\Theta$  have the same number of elements, and every type  $\theta_j \in \Theta$  is characterized by one quality signal  $s_j \in \mathcal{S}$ . The conditional probability distribution for the signals observed by the buyers is computed as:

$$f(s_k|\theta_j) = \begin{cases} 1 - \epsilon & \text{if } k = j; \\ \epsilon/(M - 1) & \text{if } k \neq j; \end{cases}$$

where  $\epsilon$  is the probability that a buyer misinterprets the true quality of the product (all mistakes are equally likely). We take  $\epsilon = 10\%$ .

The prior belief is randomly generated in the following way: for every  $\theta_j \in \Theta$ ,  $p(\theta_j)$  is a random number, uniformly distributed between 0 and 1. The probability distribution over types is then computed by normalizing these random numbers:

$$Pr[\theta_j] = \frac{p(\theta_j)}{\sum_{\theta \in \Theta} p(\theta)};$$

The external benefits from lying are randomly uniformly distributed between 0 and 1.