Towards Social Norm Design for Crowdsourcing Markets

Chien-Ju Ho

University of California, Los Angeles cjho@cs.ucla.edu

Jennifer Wortman Vaughan

University of California, Los Angeles jenn@cs.ucla.edu

Abstract

Crowdsourcing markets, such as Amazon Mechanical Turk, provide a platform for matching prospective workers around the world with tasks. However, they are often plagued by workers who attempt to exert as little effort as possible, and requesters who deny workers payment for their labor. For crowdsourcing markets to succeed, it is essential to discourage such behavior. With this in mind, we propose a framework for the design and analysis of incentive mechanisms based on social norms, which consist of a set of rules that participants are expected to follow, and a mechanism for updating participants' public reputations based on whether or not they do. We start by considering the most basic version of our model, which contains only homogeneous participants and randomly matches workers with tasks. The optimal social norm in this setting turns out to be a simple, easily comprehensible incentive mechanism in which market participants are encouraged to play a tit-for-tat-like strategy. This simple mechanism is optimal even when the set of market participants changes dynamically over time, or when some fraction of the participants may be irrational. In addition to the basic model, we demonstrate how this framework can be applied to situations in which there are heterogeneous users by giving several illustrating examples. This work is a first step towards a complete theory of incentive design for crowdsourcing systems. We hope to build upon this framework and explore more interesting and practical aspects of real online labor markets in our future work.

Introduction

Online labor markets have emerged as a popular platform for matching prospective workers around the world with paying work. While some online labor markets, like oDesk and Elance, focus on matching skilled laborers with relatively long-term projects, others, like Clickworker and Amazon Mechanical Turk, are designed to match workers with short, simple *micro-tasks*. A typical micro-task might involve captioning a picture or transcribing an audio message.

In principle, these online labor markets for micro-tasks, or *crowdsourcing markets*, could revolutionize the way in which projects are completed by giving individuals immediate access to large, diverse, and flexible pools of workers any Yu Zhang

University of California, Los Angeles yuzhang@ucla.edu

Mihaela van der Schaar University of California, Los Angeles mihaela@ee.ucla.edu

time of day or night. In practice, they are plagued by workers who attempt to exert as little effort as possible (Ipeirotis, Provost, and Wang 2010), and requesters who advertise spammy tasks or deny workers payment for tasks completed.

To address this problem, we advocate incorporating an explicit incentive mechanism into the design of crowdsourcing markets. In particular, we propose a class of incentive mechanisms based on *social norms* (Kandori 1992). A social norm consists of a set of rules that participants are expected to follow, and a mechanism for updating participants' public reputations based on whether or not they do. For example, in a crowdsourcing market, workers may be encouraged to exert high effort only for tasks posted by requesters who have made payments on time in the past. Our goal is to develop a formal framework to help platform designers identify optimal social norms for their applications, taking into account parameters about the environment and the participants.

In this work, we take into account several innate features of crowdsourcing markets:

- *The two-sided nature*. Unlike P2P systems in which all participants play similar roles, in crowdsourcing markets, the set of workers accepting tasks is typically mostly disjoint from the set of requesters posting them.
- *The difficulty of quality assessments.* For many types of micro-tasks (for example, translation from an obscure language), the quality of submitted work is difficult to determine. In some cases, requesters may need to decide whether or not to pay a worker before they are able to accurately assess the quality of his work. Even with additional time, it may not be possible to assess the quality of submitted work with certainty.
- *Anonymity.* Market participants could potentially create new identities to erase their history and start fresh. For example, in Mechanical Turk, workers are able to create new identities with only an email address.
- *Dynamic changes in the population*. In any online community, the set of participants may change over time as new individuals discover the community and current participants lose interest.
- *The existence of irrational participants.* While it is useful analytically to make the natural assumption that users are self-interested and rational, any practical incen-

Copyright © 2012, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

tive mechanism must be robust to the existence of some nonstrategic participants.

• *Platform fees.* To make a profit, the operator of a crowd-sourcing market generally charges a fee for every payment made. This fee impacts the incentives of both market participants and the market designer.

As a demonstration of our framework, we start by considering the large and general class of social norms with threshold-based social strategies. In a threshold-based social strategy, workers are expected to exert high effort on tasks posted by any requester whose reputation exceeds a specified value. Similarly, requesters are expected to pay those workers whose reputation exceeds a specified value. We show that under the class of maximum punishment reputation update rules, in which users who violate the social strategy have their reputation reset to the minimum possible value, the optimal threshold-based social strategy from the platform designer's perspective is a tit-for-tat-like strategy. In particular, it says that participants should punish other participants who did not adhere to the social strategy in the previous round. In addition to the basic model, we also demonstrate how to apply our framework when there are heterogeneous users by giving illustrative examples.

While our current model does not include all the important features in real crowdsourcing systems, we hope to provide a formal framework and gain insights on how to design incentive schemes in crowdsourcing markets. Our next steps would be to generalize our framework and include more interesting and practical aspects of social norm design in crowdsourcing markets. Some examples include the matching mechanism between workers and requesters, the full heterogeneity of the users, and the convergence of user distributions when users are learning.

Related Work

Crowdsourcing markets have recently attracted attention as a powerful tool for harnessing human labor to complete tasks that are notoriously difficult for computers. However, it is well-known that Turk users are not always trustworthy, especially when they can gain by being lazy or even cheating (Ipeirotis, Provost, and Wang 2010). Different approaches have been proposed to deal with the poor quality of work, such as redundantly assigning the same tasks to several workers and cleverly averaging their responses (Ipeirotis, Provost, and Wang 2010; Karger, Oh, and Shah 2011). However, redundancy leads to wasted effort and cannot be applied to tasks for which averaging isn't meaningful, such as creating a paragraph of original text. Although Mechanical Turk currently tracks the success rate of each worker (i.e., the fraction of tasks on which a worker has been paid), the idea of embedding a more sophisticated incentive mechanism directly in the market is not often discussed.

There are a variety of options one might consider when designing an incentive mechanism, such as using virtual currency (Kash, Friedman, and Halpern 2009) or virtual points (Jain, Chen, and Parkes 2009). For crowdsourcing systems that already involve real cash payments, we believe it is most natural to consider a mechanism based on *repu*-

tation (Resnick et al. 2000; Friedman, Resnick, and Sami 2007; Dellarocas 2005). In a reputation-based mechanism, rewards and punishments are typically determined based on a differential service scheme, which might require that a user who behaved well in the past should receive more resources or better service than a user who did not. This preferential treatment provides an incentive for users to behave well.

The problem formulation we propose is an instance of the well-studied repeated Prisoner's Dilemma (Ellison 1994). However, instead of finding achievable payoffs and analyzing equilibria, in this work, we focus on how to optimally design a reputation mechanism within a design space with given parameters. The particular reputation mechanisms we propose are based on social norms (Kandori 1992), which consist of a set of prescribed rules that market participants are asked to follow, and a mechanism for updating reputations based on whether or not they do. One advantage of social norms over traditional reputation systems is that the equilibrium selection problem is easy by design. The market operator announces the strategy that everyone should follow, and participants only need to verify that following is in their own best interest. This differs from typical reputation systems which may have a large number of equilibria, some of which lead to low payoffs for all.

To the best of our knowledge, we are the first to consider incorporating social norms into crowdsourcing markets. The most relevant related work is the social norm design for P2P systems (Zhang, Park, and van der Schaar 2011; Zhang and van der Schaar 2012). In P2P systems, all users play similar roles, and only the transmitters take actions. This differs from crowdsourcing markets. The set of workers is mostly disjoint from the set of requesters, and both sets can take actions. Interactions in P2P systems are modeled as gift-giving games, while interactions in crowdsourcing markets are naturally modeled as instances of the Prisoner's Dilemma. These apparently small differences lead to dramatically different results; to achieve optimality in P2P systems, it is necessary to carefully tune parameters based on the environment, while in our setting, a simple and intuitive mechanism is optimal for a wide range of environments.

Problem Formulation

We study the problem of designing social norms for crowdsourcing markets. In this paper, we take the point of view of the *platform designer* (e.g., Amazon in the case of Mechanical Turk) who runs the market and typically receives a fixed percentage of all payments that are made. We therefore consider the goal of maximizing the total amount of money exchanged in the system, which can be viewed as a proxy for the platform designer's profit.

For clarity of presentation, we initially make several simplifying assumptions in our model. (Some of these assumptions are relaxed in the later discussion.) First, we assume that the user population is large and static, with the number of workers in the population equal to the number of requesters. Second, we assume that at each point in time, workers and requesters are randomly matched. The problem of assigning tasks to workers is of great importance, and is discussed in our section on next steps. Third, we assume that workers and requesters are homogeneous in the sense that all workers pay the same cost (in terms of time and effort) to complete any task, and all requesters receive the same utility for tasks completed. Under these assumptions, it is natural to set one fixed price for all tasks as well.

When a worker and a requester are matched, the worker chooses to exert either high or low effort for the requester's task, and the requester chooses whether or not to pay the worker. Ideally, a requester would like to be able to carefully evaluate the worker's work before choosing whether or not to pay. However, in reality, the requester does not always have the ability to do this. Payment decisions are generally made relatively quickly, while it may take significant time for the requester to evaluate the work. For this reason, we assume that a requester must choose whether or not to pay a worker *before* the action of the worker is revealed. However, afterwords, when the requester has had a chance to evaluate the work, the worker and requester can report each other's actions to the platform designer. The platform designer can then use these reports to update reputations.

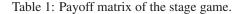
Even after time has passed, it may be difficult for a requester to accurately evaluate the quality of a worker's work, especially for advanced tasks. Therefore, we allow some noise in the reports. In particular, we assume that with probability ϵ_w (ϵ_r), the action of a worker (requester) is misreported. Both ϵ_w and ϵ_r are common knowledge. Note that in our setting, there is no incentive for workers or requesters to report dishonestly. These reports have no direct impact on the reputation of the reporter. Additionally, because we are considering a very large population, any indirect effects of a dishonest report would be negligible.

For the majority of this paper, we assume that all market participants are rational and would like to maximize their long-term discounted sum of utilities. When dealing with rational participants, we assume that all workers and requesters have the same discount factor, δ .

Single Stage Game

Since workers and requesters are unable to observe each other's action before choosing their own actions, it is natural to model their interaction as a normal-form game as in Table 1. In this game, C denotes a worker's cost for exerting high effort on the requester's task, V denotes the benefit a requester receives when a worker exerts high effort, and P denotes the monetary payment associated with the task. We implicitly assume the worker's cost for exerting low effort is 0, but this assumption is without loss of generality since our analysis depends only on the difference between the cost of exerting high effort and the cost of exerting low effort. Similarly, it is without loss of generality to assume the requester's value for receiving low effort work is 0.

		Requester	
		Pay	Don't Pay
Worker	High Effort	P-C, V-P	-C, V
	Low Effort	P, -P	0, 0



We are interested only in the natural setting in which the action pair (High Effort, Pay) leads to higher payoffs than the action pair (Low Effort, Don't Pay) for both the worker and the requester. If this were not the case, then social welfare would be maximized by everyone choosing Low Effort or Don't Pay, and there would be no need to design a social norm. Therefore, we assume that 0 < C < P < V.

Repeated Game with Social Norm

In the single-stage game described above, the only Nash equilibrium (Low Effort, Don't Pay) results in a low payoff for both the requester and the worker. Fortunately, in crowd-sourcing markets, workers and requesters typically participate in the market many times. Therefore, it is natural to model a crowdsourcing market as a repeated game, in which workers and requesters play the game above with a randomly chosen opponent at each point in time, as is common in the literature (Ellison 1994). The social norms that we propose will be for this repeated game setting. Note that in this work, we focus on how to *design* optimal social norms in the repeated game setting instead of analyzing equilibrium and finding achievable payoffs.

In the social norms that we design, both requesters and workers are assigned reputation values that are updated over time. The platform designer announces a prescribed social strategy (e.g., "don't exert effort for/pay participants with zero reputation") and updates participants' reputations based on how well they follow this strategy. Formally speaking, a social norm consists of two parts:

- The *social strategies* σ_w and σ_r are functions that define the prescribed actions for the worker and requester respectively when a worker with reputation θ_w is matched with a requester with reputation θ_r .
- The reputation update rules τ_w and τ_r define how to update the reputation of the worker and requester respectively after each transaction. We assume the reputation values are bounded between 0 and an upper bound L.¹

Designing a social norm involves specifying both the social strategies and the reputation update rules.

Optimal Social Norm Design

As a warm-up, in this section, we illustrate how to design optimal social norms under the most basic version of our model. We limit our attention to the most natural class of social strategies, threshold-based social strategies, paired with maximum punishment reputation update rules, both defined below. We first derive a set of conditions that can be used to determine whether or not a particular social norm is sustainable, i.e., whether or not each user has incentive to follow the social norm if she believes that everyone else is following. We then formalize the objective function of the platform designer and show how to find the optimal sustainable social norm with respect to this objective function.

A threshold-based social strategy is based on a pair of threshold values (k_r, k_w) . Workers are expected to exert

¹Setting different limits for workers and requesters would not affect the analysis. For simplicity, we use one parameter L.

high effort when they face requesters with reputation at least k_r , and exert low effort otherwise. Similarly, requesters are expected to pay the payment only when they are matched with workers with reputation at least k_w .

Definition 1 The threshold-based social strategies σ_w and σ_r with parameters k_r and k_w are defined as:

$$\sigma_{w}(\theta_{w}, \theta_{r}) = \begin{cases} \text{"High Effort"} & \text{if } \theta_{r} \ge k_{r}, \\ \text{"Low Effort"} & \text{otherwise.} \end{cases}$$
$$\sigma_{r}(\theta_{w}, \theta_{r}) = \begin{cases} \text{"Pay"} & \text{if } \theta_{w} \ge k_{w}, \\ \text{"Don't Pay"} & \text{otherwise.} \end{cases}$$

We initially consider a simple class of update rules. Each time the user follows the social strategy, her reputation is increased by 1 until it hits the limit L. If the user deviates from the social strategy, her reputation is reset to 0. We call this the *maximum punishment reputation update rule*.

Definition 2 The maximum punishment reputation update rules τ_w and τ_r with parameter L are defined as:

$$\begin{aligned} \tau_w(\theta_w, \theta_r, a) &= \begin{cases} \min\{\theta_w + 1, L\} & \text{if } a = \sigma_w(\theta_w, \theta_r), \\ 0 & \text{otherwise.} \end{cases} \\ \tau_r(\theta_w, \theta_r, a) &= \begin{cases} \min\{\theta_r + 1, L\} & \text{if } a = \sigma_r(\theta_w, \theta_r), \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

We first show that there exists a unique stationary reputation distribution for the maximum punishment reputation update rules. We then derive the sustainable conditions and find the optimal values of the parameters k_r , k_w , L, and P. We show that the resulting optimal social norm is both simple and intuitive for market participants.

The Stationary Reputation Distribution

Our first lemma states that there exists a stationary distribution of worker and requester reputations under the assumption that all users follow the prescribed social strategy. Recall that every worker's (requester's) action at a given time step is misreported with probability ϵ_w (ϵ_r).

Lemma 1 In the basic model, under the maximum punishment reputation update rule with any L > 0 paired with any social strategy, if all users follow the social strategy, then there exists a unique stationary distribution $\{\eta_w(\theta)\}$ over worker reputations and a unique stationary distribution $\{\eta_r(\theta)\}$ over requester reputations given by

$$\begin{split} \eta_w(\theta) &= (1 - \epsilon_w)^{\theta} \epsilon_w, \, \eta_r(\theta) = (1 - \epsilon_r)^{\theta} \epsilon_r \,, \, \text{for } \theta \in [0, L-1] \\ \eta_w(L) &= (1 - \epsilon_w)^L, \, \eta_r(L) = (1 - \epsilon_r)^L \end{split}$$

The proof of this lemma is simple. Assuming all workers follow the social strategy, we can write the transition of worker reputations over time as a set of linear equations. By treating the $\eta_w(\theta)$ and $\eta_r(\theta)$ as variables, we can get the above results by solving the linear equations.²

The stationary distributions do not depend on the social strategy used. Note that if the reporting errors are 0, all workers and requesters have the highest reputation value.

Sustainability Conditions

We next investigate the conditions under which users have incentive to follow the threshold-based social norm. If the social strategy states that requesters never pay and workers never exert high effort, it would be sustainable trivially, but this bad equilibrium is undesirable. We therefore focus on the non-trivial case where $k_w \leq L$ and $k_r \leq L$.

To check if a social norm is sustainable, we check whether a user has incentive to deviate from the social norm given all other users are following. Since we assume the population of users is large, the deviation of single agent will not influence the overall distribution. Therefore, we can calculate the long-term expected payoff using the stationary reputation distribution. Below we show the sustainable conditions for the threshold-based social norms.

Lemma 2 In the basic model, the threshold-based social strategy with parameters k_w and k_r paired with the maximum punishment reputation update rule with L > 0 is sustainable if and only if $k_w > 0$, $k_r > 0$, and

$$\frac{1}{\delta^{k_w}(1-\epsilon_w)^{k_w-1}(1-2\epsilon_w)}C \le P \le \delta^{k_r}(1-\epsilon_r)^{k_r-1}(1-2\epsilon_r)V.$$

The proof of this lemma uses the one-shot deviation principle of game theory and examines the conditions under which no user has incentive to deviate at any time step. To gain intuition about Lemma 2, consider the case in which the reports of behavior are noise-free, that is, $\epsilon_r = \epsilon_w = 0$. The lemma then implies that the social norm is sustainable if $C \leq \delta^{k_w} P$ and $P \leq \delta^{k_r} V$. These are precisely the conditions under which a worker or requester with reputation 0 has incentive to follow the social norm. The left-hand side of each equation is the current maximum additional payoff a worker or a requester can get by deviating from the social norm, and the right-hand side is the difference of expected future payoff she can get if she does not deviate from the social norm in the current stage. Since users with reputation 0 have the least incentive to follow (they have no reputation to lose), these conditions suffice to ensure that all users have incentive to follow.

This result gives us a feasible interval in which we can set the price P. If the ratio of value to cost or the discount factor is too small, (i.e., if $V/C \leq 1/\delta^{k_r+k_w}$), then there is no way to set the price to obtain a sustainable social norm.

The Optimal Social Norm

With the sustainability conditions in place, we are now ready to design the optimal threshold-based social norm. We define optimality in terms of the total of all payments, a proxy for the revenue of the platform designer. The designer has control of four parameters: the payment P, the threshold values k_r and k_w , and the maximum reputation L.

We first formally define the objective function of the design problem. Let $P(\theta_w, \theta_r) = P$ if the requester of reputation θ_r is asked to pay in a transaction with worker of reputation θ_w , i.e., if $\sigma_r(\theta_w, \theta_r) =$ "Pay", and $P(\theta_w, \theta_r) = 0$ otherwise. Assuming that the population has already converged to the unique stationary distri-

²Formal proofs of this and all results are omitted due to space constraints, but will be posted in a longer version of this paper.

bution, we can express the objective function U as $U = \sum_{\theta_r=0}^{L} \sum_{\theta_w=0}^{L} \eta_w(\theta_w) \eta_r(\theta_r) P(\theta_w, \theta_r).$

Note that $P(\theta_w, \theta_r)$ depends on the choice of k_w , but not k_r . If everyone follows the social strategy, the stationary distributions do not depend on k_w or k_r . We can therefore conclude that under the assumption that all users follow the social strategies:

- For fixed values P, k_r , and L, U is non-increasing in k_w .
- For fixed values P, k_w , and L, U is constant in k_r .

This allows us to derive the optimal parameter settings.

Theorem 1 In the basic model, restricting attention to threshold-based social strategies paired with maximum punishment reputation update rules, if

$$\frac{V}{C} \ge \frac{1}{\delta^2 (1 - 2\epsilon_w)(1 - 2\epsilon_r)}$$

then total payments are maximized by setting $k_r = k_w = L = 1$ and $P = \delta(1 - 2\epsilon_r)V$. The optimal value of the objective U is then $\delta(1 - \epsilon_w)(1 - 2\epsilon_r)V$. Otherwise, there is no sustainable social norm.

The optimal payment P is chosen to be the maximum value satisfying the condition in Lemma 2. In addition, if the social norm is not sustainable with $k_r = k_w = 1$, then it will not be sustainable with any setting of k_r and k_w . Since U is non-increasing both in k_r and k_w , setting $k_r = k_w = 1$ is optimal, and setting L = 1 does not affect the optimality.

With both of the thresholds and the maximum reputation L equal to 1, the optimal social norm from Theorem 1 is surprisingly simple and intuitive. In this social norm, there are only two different reputation values. Modulo the effects of noise, users who followed the social norm on the previous time step have reputation 1, while users who did not follow have reputation 0. The social strategy then says that users should play a tit-for-tat-like strategy, exerting high effort for or paying those users who obeyed the social norm on the previous time step, and punishing those who did not.

This result is both reassuring and somewhat surprising. It tells us that there is no need to construct complicated incentive mechanisms that are difficult for users to understand, or to heavily optimize parameters of the mechanism based on properties of the market participants. More interestingly, in the next section, we show that this intuitive social norm is still optimal even if we relax some of the assumptions of the basic model. This is in stark contrast to the P2P setting (Zhang, Park, and van der Schaar 2011; Zhang and van der Schaar 2012) in which in order to derive optimal social norms it is necessary to tune complex parameters based on properties of the environment.

Beyond the Basic Model

Dynamic Population and Whitewashing

Until now, we have assumed that the user population is static. However, in real crowdsourcing markets, participants enter and exit the market over time. This affects user incentives in several ways. First, if users have the ability to leave and rejoin the market, they may be tempted to engage in *whitewashing* (Friedman and Resnick 2001), creating a new identity to escape a bad reputation. Second, if a user knows he will not stay in the market forever, he may be less willing to work hard to earn a high reputation.

We first consider the problem of whitewashing. Luckily, it is easy to prevent whitewashing using social norms by setting the reputation of any new user to 0. It is clear that doing this removes any incentive that a user might have to erase his history. With this problem under control, we need only to consider how the turnover rate affects the stationary reputation distributions and users' expectations about their future payoffs. To simplify analysis, we assume the size of worker and requester populations stay the same at all times, but at every time step, some fraction of users leave the market and are replaced by new users (with reputation 0). Let $\alpha_w (\alpha_r)$ be the fraction of workers who leave and are replaced by new workers (requesters) at each time step.

Theorem 2 *Restricting attention to threshold-based social* strategies paired with maximum punishment reputation update rules, for any turnover rates α_w and $\alpha_r \in [0, 1)$, if

$$\frac{V}{C} \ge \frac{1}{\delta^2 (1 - \alpha_r)(1 - \alpha_w)(1 - 2\epsilon_w)(1 - 2\epsilon_r)}$$

then total payments are maximized by setting $k_r = k_w = L = 1$ and $P = \delta(1 - \alpha_r)(1 - 2\epsilon_r)V$. Otherwise, there is no sustainable social norm.

The analysis is similar to Theorem 1 with two modifications. First, the stationary distribution changes in this setting. Second, users will discount their expected payoff by the turnover rate, since they expect to leave the market with probability α_w or α_r in the next time step.

Nonstrategic Users

Until now, we have assumed that all users are rational and strategic. While this assumption is natural and useful analytically, in real markets there exist nonstrategic users. For example, a system may have *altruists*, who always perform "good" actions (e.g., making payments), *malicious* users, who always perform "bad" actions (e.g., not paying), or users who choose actions ignoring social strategies and their own utility. Any reputation mechanisms should be robust to the existence of nonstrategic users to be useful in practice.

The following theorem applies whenever some fraction of the population is nonstrategic, but choose actions independent of their opponents' reputation at each point in time. Call such a user *oblivious-nonstrategic*. We assume that the fraction of nonstrategic users is known, but the identities of the nonstrategic users are not.

Theorem 3 Restricting attention to threshold-based social strategies paired with maximum punishment reputation update rules, for any $f_r, f_w \in [0, 1)$, if a fraction f_r of requesters and a fraction f_w of workers are oblivious-nonstrategic, then if

$$\frac{V}{C} \ge \frac{1}{\delta^2 (1 - f_w)(1 - f_r)(1 - 2\epsilon_r)(1 - 2\epsilon_w)}$$

then total payments are maximized by setting $k_r = k_w = L = 1$ and $P = \delta(1 - f_w)(1 - 2\epsilon_r)V$. Otherwise, there is no sustainable social norm.

The existence of nonstrategic users shortens the interval of feasible prices, i.e., the gain of future payoff by following the social norm is discounted by the fraction of nonstrategic users. However, it does not affect the social norm design; the tit-for-tat-like social norm is still optimal.

Transaction Fees

We now discuss how to explicitly incorporate the platform designer's fee into the analysis. We assume the platform designer always takes a fixed portion m from the payment as the transaction fee. Since the objective function is 1/m of the total payment, introducing this fee does not change the design problem in terms of the objective. However, there is now a difference between the payment paid by requesters (P) and the payment received by workers (P(1 - m)). The optimal social norm is characterized as follows.

Theorem 4 Suppose that the platform designer takes a fixed portion $m \in [0,1)$ of all payments. Restricting attention to threshold-based social strategies paired with maximum punishment reputation update rules, if

$$\frac{V}{C} \geq \frac{1}{\delta^2(1-2\epsilon_w)(1-2\epsilon_r)(1-m)}$$

then total fees are maximized by setting $k_r = k_w = L = 1$, $P = \delta(1-2\epsilon_r)V$, and $m = 1 - \frac{1}{\delta^2(1-2\epsilon_w)(1-2\epsilon_r)}\frac{C}{V}$. Otherwise, there is no sustainable social norm.

Heterogeneous Users

In this section, we illustrate how our framework can be applied to situations in which there are heterogeneous users. These results are meant as illustrative examples only, and the important problem of designing optimal social norms for settings with heterogeneous users is left for future work.

We focus on the case in which requesters are still homogeneous with value V, but workers have different costs. In this case, information is *asymmetric* in the sense that workers know their own costs, but the requesters with whom they are matched know only the distribution over worker costs and not the cost of the particular worker. To simplify the analysis and exposition of results, we restrict our attention to the case of threshold-based social strategies with binary reputation values (i.e., L = 1). With L = 1, it is natural to set the thresholds $k_w = k_r = 1$. Deriving the optimal social norm within this simplified setting reduces to finding the payment P that maximizes total money exchanged.

Two Worker Types

To get our feet wet and gain some intuition about how we might deal with workers with different costs, we first consider the most simple case in which there are two types of workers with costs C_1 and C_2 . Without loss of generality, we assume $C_1 \ge C_2$. Let f_1 be the fraction of workers with cost C_1 (type 1 workers) and f_2 be the fraction with cost C_2 (type 2 workers). We begin by deriving the conditions under which workers of each type would follow the social norm *if they believed* that all requesters were following. **Lemma 3** Consider the two-worker-type setting. Suppose that all workers believe that all requesters will always follow the social norm. For $i \in \{1, 2\}$, if $P \ge C_i/(\delta(1-2\epsilon_w))$, then workers of type *i* maximize their expected utility by always following the social norm regardless of their own reputation. Otherwise, workers of type *i* maximize their expected utility by always exerting low effort.

Let $P_w(C_i) = C_i/(\delta(1 - 2\epsilon_w))$. From this lemma, we know that this is the minimum payment for type *i* workers to follow the social norm, given the belief that all requesters are following. If $P < P_w(C_i)$, type *i* workers will always choose to exert low effort. The key insight is that in this case, type *i* workers exhibit precisely the same behavior as the malicious oblivious-nonstrategic users.

Similarly, let $P_r(V) = \delta(1 - 2\epsilon_r)V$ be the maximum payment requesters are willing to pay if they believe that all workers always follow the social norm. If requesters instead believe that a fraction f of workers are oblivious nonstrategic and all others follow the social norm, then by an argument similar to the one used in Theorem 3, requesters would be willing to pay a maximum of $(1 - f)P_r(V)$.

Combining these ideas, we formally state how to set payment P to maximize the total payment.

Theorem 5 Consider the two-worker-type setting.

- *I.* If $V \geq \frac{C_1}{\delta^2(1-2\epsilon_w)(1-2\epsilon_r)}$, then total payments are maximized by setting $P = \delta(1-2\epsilon_r)V$. In this case, all workers and requesters follow the social norm.
- 2. If $\frac{C_2}{\delta^2(1-2\epsilon_w)(1-2\epsilon_r)(1-f_1)} \leq V < \frac{C_1}{\delta^2(1-2\epsilon_w)(1-2\epsilon_r)}$, then total payments are maximized by setting $P = \delta(1 2\epsilon_r)(1-f_1)V$. In this case, requesters and type 2 workers follow the social norm, while type 1 workers do not.
- 3. Otherwise, there is no sustainable social norm.

General Distributions Over Worker Types

Building on the intuition gained in the two-worker-type setting, we now consider heterogeneous workers with costs $C \ge 0$ drawn from an arbitrary distribution with probability density function f(C). By the same argument used in Lemma 3, if the payment is P, given the belief that all requesters always follow the social norm, workers with any cost $C \le \delta(1 - 2\epsilon_w)P$ would maximize their expected utility by following the social norm. Therefore, given any payment P, if workers believe that requesters follow the social norm, the fraction of workers who would maximize their expected utility by following the social norm is $F(P) = \int_0^{\delta(1-2\epsilon_w)P} f(C)dC$.

By treating users who don't follow the social norm as oblivious-nonstrategic users, we know requesters will maximize their expected utility by following the social norm only if the payment $P \le F(P)P_r(V)$.

Proposition 1 Given homogeneous requesters with value V and heterogeneous workers with costs $C \ge 0$ drawn from a distribution with probability density function f(C), total payments are maximized by setting P to the largest positive value satisfying $P \le P_r(V) \int_0^{\delta(1-2\epsilon_w)P} f(C) dC$. In this case, all requesters and all workers with $C \leq \delta(1-2\epsilon_w)P$ would follow the social norm, while all workers with $C > \delta(1-2\epsilon_w)P$ would not. If no such P exists, then no social norm is sustainable.

As an illustrative example of how this result can be applied, we examine the case in which costs are drawn uniformly from an interval $[C_{min}, C_{max}]$.

Theorem 6 Given homogeneous requesters with value V and heterogeneous workers with costs C > 0 drawn uniformly from the interval $[C_{min}, C_{max}]$, if

$$\frac{V}{C_{max}} \geq \frac{1}{\delta^2(1-2\epsilon_w)(1-2\epsilon_r)},$$

then the total payment is maximized by setting $P = \delta(1 - 2\epsilon_r)V$, in which case all users follow the social norm. Otherwise, there is no sustainable social norm.

This result might seem surprising at first glance. It shows if the social norm is not sustainable with homogeneous workers with cost C_{max} , the social norm is also not sustainable when costs are uniformly drawn from $[C_{min}, C_{max}]$, even though some workers have lower costs.

Below we offer an intuition to explain why the results makes sense. Recall that $P_r(V)$ is the maximum amount that requesters are willing to pay if they believe that all workers always follow the social norm, and $P_w(C)$ is the minimum payment for which workers with cost C are willing to work if they believe all requesters always follow the social norm. Consider the nontrivial case in which $P_r(V) < P_w(C_{max})$. In this case, if we were to set $P = P_r(V)$, a portion of the workers would not be willing to work and would become uncooperative, behaving like malicious oblivious-nonstrategic users. But if some workers were uncooperative, requesters would not be willing to pay $P_r(V)$, and we would need to set the price P lower to incentivize requesters to follow the social norm. However, if we set the price lower, more workers would become uncooperative, and we would need to set the price even lower to satisfy requesters. When worker costs are drawn from a uniform distribution, this process can be repeated until all workers are uncooperative, and we can never find a value of P to make the social norm sustainable.

Fortunately, in most real-world situations, we would not expect worker costs to be uniformly distributed. Suppose instead that worker costs are drawn from a normal distribution with mean C_{avg} and standard deviation σ . Using Proposition 1, we can find the optimal payment using the properties of normal distributions. For example, if $0.841P_r(V) \ge (C_{avg} + \sigma)/(\delta(1 - 2\epsilon_w))$, we can find a payment which is at least $0.841P_r(V)$ and sustains the social norm.

Conclusion and Future Work

We introduced a framework for the design and analysis of incentive schemes for crowdsourcing markets based on social norms, and described a general technique that can be used to derive the optimal social norm from within a class of interest. We illustrated the use of this technique to derive the optimal social norm from within the natural class of threshold-based social strategies paired with maximum punishment reputation update rules, and showed that the optimal norm in this class is simple to implement and understand. Furthermore, the optimal social strategy does not depend on features of the environment such as the turnover rate of the population or the fraction of non-strategic users, making it applicable in a variety of settings. While the main results are proved under the assumption of homogeneous users, we also provided a selection of illustrative examples demonstrating how our framework can be used to analyze social norms for heterogeneous users as well.

This work is a first step towards a complete, robust theory of incentive design for crowdsourcing systems. An important next step in developing this theory is to build upon our illustrative examples of heterogeneity to derive techniques for obtaining optimal social norms for any distribution over worker costs and requester values. Allowing full heterogeneity would also introduce interesting questions related to the problem of optimally matching workers and requesters. Finally, developing a full theory would require digging into issues related to learning and convergence to the stationary reputation distribution. Our hope is that the framework and techniques in this paper will provide the necessary groundwork for future progress towards this goal.

Acknowledgement

This research was partially supported by the National Science Foundation under grant IIS-1054911.

References

Dellarocas, C. 2005. Reputation mechanism design in online trading environments with pure moral hazard. *Information Systems Research* 16(2):209–230.

Ellison, G. 1994. Cooperation in the prisoner's dilemma with anonymous random matching. *Review of Economic Studies* 61(3):567–588.

Friedman, E. J., and Resnick, P. 2001. The social cost of cheap pseudonyms. *Journal of Economics and Management Strategy* 10(2):173–199.

Friedman, E.; Resnick, P.; and Sami, R. 2007. Manipulationresistant reputation systems. In Nisan, N.; Roughgarden, T.; Tardos, E.; and Vazirani, V., eds., *Algorithmic Game Theory*. Cambridge University Press.

Ipeirotis, P. G.; Provost, F.; and Wang, J. 2010. Quality management on Amazon Mechanical Turk. In *HCOMP*.

Jain, S.; Chen, Y.; and Parkes, D. C. 2009. Designing incentives for online question and answer forums. In *ACM EC*.

Kandori, M. 1992. Social norms and community enforcement. *Review of Economic Studies* 59(1):63–80.

Karger, D. R.; Oh, S.; and Shah, D. 2011. Iterative learning for reliable crowdsourcing systems. In *NIPS*.

Kash, I. A.; Friedman, E. J.; and Halpern, J. Y. 2009. Manipulating scrip systems: Sybils and collusion. In *AMMA*.

Resnick, P.; Kuwabara, K.; Zeckhauser, R.; and Friedman, E. 2000. Reputation systems. *Communications of the ACM* 43:45–48.

Zhang, Y., and van der Schaar, M. 2012. Peer-to-peer multimedia sharing based on social norms. *Elsevier Journal Signal Processing: Image Communication*. To appear.

Zhang, Y.; Park, J.; and van der Schaar, M. 2011. Social norms for networked communities. Preprint available at http://arxiv.org/abs/1101.0272.