

# Putting Peer Prediction Under the Micro(economic)scope and Making Truth-Telling Focal

Yuqing Kong<sup>1</sup>(✉), Katrina Ligett<sup>2,3</sup>, and Grant Schoenebeck<sup>1</sup>

<sup>1</sup> University of Michigan, Ann Arbor, USA  
{yuqkong, schoeneb}@umich.edu

<sup>2</sup> California Institute of Technology, Pasadena, USA  
katrina@caltech.edu

<sup>3</sup> Hebrew University, Jerusalem, Israel

**Abstract.** Peer-prediction [19] is a (meta-)mechanism which, given any proper scoring rule, produces a mechanism to elicit private information from self-interested agents. Formally, truth-telling is a strict Nash equilibrium of the mechanism. Unfortunately, there may be other equilibria as well (including uninformative equilibria where all players simply report the same fixed signal, regardless of their true signal) and, typically, the truth-telling equilibrium does not have the highest expected payoff. The main result of this paper is to show that, in the symmetric binary setting, by tweaking peer-prediction, in part by carefully selecting the proper scoring rule it is based on, we can make the truth-telling equilibrium focal—that is, truth-telling has higher expected payoff than any other equilibrium.

Along the way, we prove the following: in the setting where agents receive binary signals we (1) classify all equilibria of the peer-prediction mechanism; (2) introduce a new technical tool for understanding scoring rules, which allows us to make truth-telling pay better than any other informative equilibrium; (3) leverage this tool to provide an optimal version of the previous result; that is, we optimize the gap between the expected payoff of truth-telling and other informative equilibria; and (4) show that with a slight modification to the peer-prediction framework, we can, in general, make the truth-telling equilibrium focal—that is, truth-telling pays more than any other equilibrium (including the uninformative equilibria).

**Keywords:** Information elicitation · Peer prediction · Crowdsourcing

---

Y. Kong—Supported by National Science Foundation Career Award 1452915 and CCF Award 1618187.

K. Ligett—Supported in part NSF grants 1254169 and 1518941, US-Israel Binational Science Foundation Grant 2012348, the Charles Lee Powell Foundation, a subcontract through the Brandeis project, a grant from the HUJI Cyber Security Research Center, and a startup grant from Hebrew University's School of Computer Science.

G. Schoenebeck—Supported by National Science Foundation Career Award 1452915 and Algorithms in the Field Award 1535912.

# 1 Introduction

From Facebook.com’s “What’s on your mind?” to Netflix’s 5-point ratings, from innumerable survey requests in one’s email inbox to Ebay’s reputation system, user feedback plays an increasingly central role in our online lives. This feedback can serve a variety of important purposes, including supporting product recommendations, scholarly research, product development, pricing, and purchasing decisions. With increasing requests for information, agents must decide where to turn their attention. When privately held information is elicited, sometimes agents may be intrinsically motivated to both participate and report the truth. Other times, self-interested agents may need incentives to compensate for costs associated with truth-telling and reporting: the effort required to complete the rating (which could lead to a lack of reviews), the effort required to produce an accurate rating (which might lead to inaccurate reviews), foregoing the opportunity to submit an inaccurate review that could benefit the agent in future interactions [11] (which could, e.g., encourage negative reviews), or a potential loss of privacy [8] (which could encourage either non-participation or incorrect reviews).

To overcome a lack of (representative) reviews, a system could reward users for reviews. However, this can create perverse incentives that lead to inaccurate reviews. If agents are merely rewarded for participation, they may not take time to answer the questions carefully, or even meaningfully.

To this end, explicit reward systems for honest ratings have been developed. If the ratings correspond to objective information that will be revealed at a future date, this information can be leveraged (e.g., via prediction markets) to incentive honesty. In this paper, we study situations where this is not the case: the ratings cannot be independently verified either because no objective truth exists (the ratings are inherently subjective) or an objective truth exists, but is not observable.

In such cases, it is known that correlation between user types can be leveraged to elicit truthful reports by using side payments [1, 2, 4, 5]. Miller et al. [19] propose a particular such (meta-)mechanism for truthful feedback elicitation, known as *peer prediction*. Given any proper scoring rule (a simple class of payment functions we describe further below), and a prior where each agent’s signal is “stochastically relevant” (informative about other agents’ signals), the corresponding peer prediction mechanism has truth-telling as a strict Bayesian-Nash equilibrium.

There is a major problem, however: alternative, non-truthful equilibria may have higher payoff for the agents than truth-telling. This is the challenge that our work addresses.

*Our Results.* The main result of this paper is to show that by tweaking peer prediction, in part by specially selecting the proper scoring rule it is based on, we can make the truth-telling equilibrium focal—that is, truth-telling has higher expected payoff than any other equilibrium.

Along the way we prove the following: in the setting where agents receive binary signals we (1) classify all equilibria of the peer prediction mechanism; (2) introduce a new technical tool, the best response plot, and use it to show that we can find proper scoring rules so the truth-telling pays more, in expectation, than any other informative equilibrium; (3) we provide an optimal version of the previous result, that is we optimize the gap between the expected payoff of truth-telling and other informative equilibrium; and (4) we show that with a slight modification to the peer prediction framework, we can, in general, make the truth-telling equilibrium focal—that is, truth-telling pays more than any other equilibrium (including the uninformative equilibria).

The main technical tool we use is a best response plot, which allows us to easily reason about the payoffs of different equilibria. We first prove that no asymmetric equilibria exist. The naive approach then would be to simply plot the payoffs of different symmetric strategies. However, for even the simplest proper-scoring rules, these payoff curves are paraboloid, and hence difficult to analyze directly. The best response plot differs from this naive approach in two ways: first, instead of plotting the strategies of agents explicitly, the best response plot aggregates the results of these actions; second, instead of plotting the payoffs of all agents, the best response plot analyzes the payoff of one distinguished agent which, given the strategies of the remaining agents, plays her best response. This makes the plot piece-wise linear for all proper scoring rules, which makes analysis tractable. We hope that the best response plot will be useful in future work using proper scoring rules.

## 1.1 Related Work

Since the seminal work of Miller et al. introducing peer prediction [19], a host of results in closely related models have followed (see, e.g., [9, 11, 12, 14]), primarily motivated by opinion elicitation in online settings where there is no objective ground truth.

Recent research [7] indicates that individuals in lab experiments do not always truth-tell when faced with peer prediction mechanisms; this may in part be related to the issue of equilibrium multiplicity. Gao et al. [7] ran studies over Mechanical Turk using two treatments: in the first they compensated the participants according to peer prediction payments, and in the second they gave them a flat reward for participation. In their work, the mechanism had complete knowledge of the prior. The participants responded truthfully more often when the payoffs were fixed than in response to the peer prediction payments. However, it should be noted that the task the agents were asked to perform took little effort (report the received signal), and the participants were not primed with any information about the truthful equilibrium of the peer prediction mechanism (they were only told the payoffs)—an actual surveyor would have incentive to prime the participants to report truthfully.

The most closely related work is a series of papers by Jurca and Faltings [12, 14], which studies collusion between the reporting agents. In a weak model

of collusion, the agents may be able to coordinate ahead of time (before receiving their signals) to select the equilibrium with the highest payoff. Jurca and Faltings use techniques from algorithmic mechanism design to design a mechanism where, in most situations, the only symmetric *pure* Nash equilibria are truth-telling. They explicitly state the challenge of analysing mixed-Nash equilibrium as an open question, and show challenges to doing this in their algorithmic mechanism design framework [12, 14]. Our techniques, in contrast, allow us to analyse all Nash equilibria of the peer prediction mechanism including both mixed-strategy and asymmetric equilibria. Instead of eliminating equilibria, we enforce that they have a lower expected payoff than truth-telling. Additionally, the algorithmic mechanism design framework used by Jurca and Faltings sacrifices “the simplicity of specifying the payments through closed-form scoring rules” [12] that was present in the peer prediction paper. Our work recovers a good deal of that simplicity.

Jurca and Faltings further analyze other settings where colluding agents can make transfer payments, or may collude after receiving their signals. In particular, they again use automated mechanism design to show that in the case where agents coordinate after receiving their signals that even without transfer payments, there will always be multiple equilibria; in this setting, they pose the question of whether the truth-telling equilibrium can be endowed with the highest expected payoff. We do not deal with this setting explicitly, but in the settings we consider, we show that even in the face of multiple equilibria, we can ensure that the truth-telling equilibrium has the highest expected payoff and no other equilibrium is paid the same with truth-telling.

In a different paper [11], Jurca and Faltings show how to minimize payments in the peer prediction framework. Their goal is to discover how much “cost” is associated with a certain marginal improvement of truth-telling over lying. In this paper, they also consider generalizations of peer prediction, where more than one other agent’s report is used as a reference. Our work takes this to the extreme (as did [8] before us) using *all* of the *other* agents’ reports as references.

A key motivation of one branch of the related work is removing the assumption that the mechanism knows the common prior [3, 6, 10, 13, 15, 18, 20–22, 24, 25]. Dasgupta and Ghosh [3], Kamble et al. [15], Kong and Schoenebeck [16], Shnayder et al. [23] study a different setting where agents are asked to answer several a priori similar questions. Our results can be applied even when there is just a single questions (thus we do not need to assume any relation between questions). Kamble et al. [15]’s mechanism applies to both homogeneous and heterogeneous population but requires a large number of a priori similar tasks. However, Kamble et al. [15]’s mechanism contains non-truthful equilibria that are paid higher than truth-telling. Dasgupta and Ghosh [3]’s mechanism has truth-telling as the equilibrium with the highest payoff, but contains a non-truthful equilibrium that is paid as much as truth-telling. Prelec [20] shows that in his Bayesian Truth Serum (BTS), truth-telling maximizes each individual’s expected “Information-score” across all equilibria. However, this guarantee is not strict, and requires the number of agents to be infinite, even to just have truth-telling be an equilibrium.

Moreover, it is hard to classify the equilibria or optimize mechanism in Prelec’s setting. Another drawback of BTS is that it requires agents to report “prediction” while our mechanism only requires agents to report a single signal. Radanovic and Faltings [21]’s mechanism solves this drawback but that mechanism is in a sensing scenario and needs to compare the information of an sensor’s local neighbours with the information of global sensors while our mechanism does not require this local/global structure. Moreover, like BTS, Radanovic and Faltings [21]’s mechanism does not have the strictness guarantee and requires the number of agents to be infinite even to have truth-telling as an equilibrium. In addition, Lambert and Shoham [18] provide a mechanism such that no equilibrium pays more than truth-telling, but here all equilibria pay the same amount; and while truth-telling is a Bayesian Nash equilibrium, unlike in peer prediction it generally is not a strict Bayesian Nash equilibrium. *Minimal Truth Serum (MTS)* [22] is a mechanism where agents have the option to report or not report their predictions, and also lacks analysis of non-truthful equilibria. MTS uses a typical zero-sum technique such that all equilibria are paid equally.

Equilibrium multiplicity is clearly a pervasive problem in this literature. While our present work only applies to the classical peer prediction mechanism, it provides an important step in addressing equilibrium multiplicity, and a new toolkit for reasoning about proper scoring rules.

*Subsequent Work.* Kong and Schoenebeck [17] show analogous results in the setting where mechanism does not know the prior; however, they also prove that results as strong as those in this paper are impossible in that setting.

## 2 Preliminaries, Background, and Notation

### 2.1 Game Setting

Consider a setting with  $n$  agents  $A$ . If  $A' \subseteq A$ , we let  $-A'$  denote  $A \setminus A'$ . Each agent  $i$  has a private signal  $\sigma_i \in \Sigma$ . We consider a game in which each agent  $i$  reports some signal  $\hat{\sigma}_i \in \Sigma$ . Let  $\boldsymbol{\sigma}$  denote the vector of signals and  $\hat{\boldsymbol{\sigma}}$  denote the vector of reports. Let  $\boldsymbol{\sigma}_{-i}$  and  $\hat{\boldsymbol{\sigma}}_{-i}$  denote the signals and reports excluding that of agent  $i$ ; we regularly use the  $-i$  notation to exclude an agent  $i$ .

We would like to encourage truth-telling, namely that agent  $i$  reports  $\hat{\sigma}_i = \sigma_i$ . To this end, agent  $i$  will receive some payment  $\nu_i(\hat{\sigma}_i, \hat{\boldsymbol{\sigma}}_{-i})$  from our mechanism. In this paper, the game will be *anonymous*, in that each player’s payoffs will depend only on the player’s own report and the *fraction* of other players giving each possible report  $\in \Sigma$ , and not on the identities of those players.

**Assumption 1 (Binary Signals).** *We will refer to the case when  $\Sigma = \{0, 1\}$  as the **binary signal** setting, and we focus on this setting in this paper.*

**Assumption 2 (Symmetric Prior).** *We assume throughout that the agents’ signals  $\boldsymbol{\sigma}$  are drawn from some joint **symmetric prior**  $Q$ : a priori, each agent’s signal is drawn from the same distribution. We in fact only leverage a weaker assumption, that  $\forall \sigma, \sigma'$ , and  $\forall i \neq j$  and  $k \neq l$ , we have  $\Pr[\sigma_j = \sigma' | \sigma_i = \sigma] = \Pr[\sigma_l = \sigma' | \sigma_k = \sigma]$ .*

That is, the inference your signal lets you draw about others' signals does not depend on your identity or on the identity of the other agent.

Given the prior  $Q$ , for  $\sigma \in \Sigma$ , let  $q(\sigma)$  be the fraction of agents that an agent expects will have  $\sigma_j = \sigma$  *a priori*. Let

$$q(\sigma'|\sigma) := \Pr[\sigma_j = \sigma' | \sigma_i = \sigma]$$

(where  $j \neq i$ ) be the fraction of other agents that a user  $i$  expects have received signal  $\sigma'$  given that he has signal  $\sigma$ .

**Assumption 3 (Signals Positively Correlated).** *We assume throughout that the prior  $Q$  is **positively correlated**, namely that  $q(\sigma|\sigma) > q(\sigma)$ , for all  $\sigma \in \Sigma$ .*

That is, once a player sees that his signal is  $\sigma$ , this strictly increases his belief that others will have signal  $\sigma$ , when compared with his prior. Notice that even after an agent receives his signal, he may still believe that he is in the minority. Thus, simply encouraging agent agreement is not sufficient to incentivize truthful reporting.

**Assumption 4 (Signal Asymmetric Prior).** *An additional assumption we will often use is that the prior is **signal asymmetric**. For binary signals, as we consider in this paper, this simply means that  $q(0) \neq q(1)$ .*

For a richer signal space, intuitively, a signal asymmetric prior is one that changes under a relabeling of the signals, so that lying can potentially be distinguishable from truth-telling.

We say that an agent plays *response*  $\sigma \rightarrow \hat{\sigma}$ , if the agent reports signal  $\hat{\sigma}$  when he receives signal  $\sigma$ . Let  $X$  be the set of all responses (e.g.  $X = \{0 \rightarrow 0, 0 \rightarrow 1, 1 \rightarrow 0, 1 \rightarrow 1\}$  when  $\Sigma = \{0, 1\}$ ). In a *pure-strategy* an agent chooses a response for each  $\sigma \in \Sigma$ , and thus there are  $|\Sigma|^{|\Sigma|}$  possible pure strategies. Let  $S$  be the set of pure strategies and let  $s_i \in S$  denote a pure-strategy for agent  $i$ . We will also consider mixed strategies  $\theta_i$ , where agent  $i$  randomizes over pure strategies. Here we write

$$\theta_i(\sigma'|\sigma) := \Pr[\hat{\sigma}_i = \sigma' | \sigma_i = \sigma].$$

A strategy profile  $\theta = (\theta_1, \dots, \theta_n)$  consists of a strategy for each agent.

We can think of each  $\theta$  as a linear transformation from a distribution over received signals to a distribution of reported signals. Given a set of agents  $A' \subset A$ , we define

$$\theta'_{A'}(\sigma'|\sigma) := E_{i \leftarrow A'}[\theta_i(\sigma'|\sigma)]$$

where  $i \leftarrow A'$  means  $i$  is chosen uniformly at random from  $A'$ . When discussing *symmetric strategy profiles* where all players use the same strategy, we will often abuse notation and use notation for one agent's strategy to denote the entire strategy profile.

A *Bayesian Nash equilibrium* consists of a strategy profile  $\theta = (\theta_1, \dots, \theta_n)$  such that no player wishes to change his strategy, given the strategies of the other

players and the information contained the prior and his signal: for all  $i$  and for all alternative strategies  $\theta'_i$  for  $i$ ,  $\mathbb{E}[\nu_i(\boldsymbol{\theta})] \geq \mathbb{E}[\nu_i(\theta'_i, \boldsymbol{\theta}_{-i})]$ , where the expectations are over the realizations of the randomized strategies and the prior  $Q$ . We call such an equilibrium **focal** if it provides a strictly larger payoff, in expectation, to each agent, than any other equilibrium and **weakly focal** if it provides a larger payoff (maybe not strictly).

Given a symmetric prior  $Q$  and strategy profile  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_n)$ , we define

$$\hat{q}_j(\sigma'|\sigma) := \Pr[\hat{\sigma}_j = \sigma' | \sigma_i = \sigma] = \sum_{\sigma'' \in \Sigma} q(\sigma''|\sigma)\theta_j(\sigma'|\sigma'')$$

for  $i \neq j$ . Intuitively,  $\hat{q}_j(\sigma'|\sigma)$  is the probability of player  $j$  reporting  $\sigma'$ , given that another player  $i$  sees signal  $\sigma$ ; note that this probability does not depend on the identity of  $i$ , by symmetry of the prior. Given a set of agents  $A' \subset A$ , we define

$$\hat{q}'_{A'}(\sigma'|\sigma) := E_{j \leftarrow A'} \hat{q}_j(\sigma'|\sigma)$$

where  $j \leftarrow A'$  means  $j$  is chosen uniformly at random from  $A'$  (again assuming that the implicit reference agent  $i \notin A'$ ). If  $\boldsymbol{\theta} = (\theta, \dots, \theta)$  is symmetric, we simplify our notation to  $\hat{q}(\sigma'|\sigma)$  because the referenced set of agents does not matter.

In the binary signal setting when  $\boldsymbol{\theta}$  is symmetric, we have:

$$\hat{q}(1|0) = \theta(1|0)q(0|0) + \theta(1|1)q(1|0) \tag{1}$$

$$\hat{q}(1|1) = \theta(1|0)q(0|1) + \theta(1|1)q(1|1) \tag{2}$$

Additionally, we observe that  $q(1|b) = 1 - q(0|b)$ ,  $\theta_i(1|b) = 1 - \theta_i(0|b) \forall i$ , and  $\hat{q}(1|b) = 1 - \hat{q}(0|b)$ . Note that we will typically use  $b$  instead of  $\sigma$  to refer to binary signals (bits).

There are four pure strategies for playing the game in the binary signal setting: always 0, always 1, truth-telling, lying:

$$S = \left\{ \begin{pmatrix} 0 \rightarrow 0 \\ 1 \rightarrow 0 \end{pmatrix}, \begin{pmatrix} 0 \rightarrow 1 \\ 1 \rightarrow 1 \end{pmatrix}, \begin{pmatrix} 0 \rightarrow 0 \\ 1 \rightarrow 1 \end{pmatrix}, \begin{pmatrix} 0 \rightarrow 1 \\ 1 \rightarrow 0 \end{pmatrix} \right\} = \{\mathbf{0}, \mathbf{1}, \mathbf{T}, \mathbf{F}\}.$$

We will denote mixed strategies as  $\begin{pmatrix} 0 \rightarrow \theta(1|0) \\ 1 \rightarrow \theta(1|1) \end{pmatrix}$ .

### 2.2 Proper Scoring Rules

A scoring rule  $PS : \Sigma \times \Delta_\Sigma \rightarrow \mathbb{R}$  takes in signal  $\sigma \in \Sigma$  and a distribution over signals  $\delta_\Sigma \in \Delta_\Sigma$  and outputs a real number. A scoring rule is *proper* if, whenever the first input is drawn from a distribution  $\delta_\Sigma$ , then the expectation of  $PS$  is maximized by  $\delta_\Sigma$ . A scoring rule is called *strictly proper* if this maximum is unique. We will assume throughout that the scoring rules we use are strictly proper. By slightly abusing notation, we can extend a scoring rule to be  $PS : \Delta_\Sigma \times \Delta_\Sigma \rightarrow \mathbb{R}$  by simply taking  $PS(\delta_\Sigma, \delta'_\Sigma) = \mathbb{E}_{\sigma \leftarrow \delta_\Sigma} (PS(\sigma, \delta'_\Sigma))$ .

In the case of scoring rules over binary signals, a distribution can be represented by a number in the unit interval, denoting the probability placed on the signal 1. In the binary signal setting, then, we extend proper scoring rules to be defined on  $[0, 1] \times [0, 1]$ .

*Example 1 (Example of Proper Scoring Rule).* The Brier Scoring Rule for predicting a binary event is defined as follows. Let  $I$  be the indicator random variable for the binary event to be predicted. Let  $q$  be the predicted probability of the event occurring. Then:

$$B(I, q) = 2I \cdot q + 2(1 - I) \cdot (1 - q) - q^2 - (1 - q)^2.$$

Note that if the event occurs with probability  $p$ , then the expected payoff of reporting a guess  $q$  is (abusing notation slightly):

$$B(p, q) = 2p \cdot q + 2(1 - p) \cdot (1 - q) - q^2 - (1 - q)^2 = 1 - 2(p - 2p \cdot q + q^2)$$

This is (uniquely) maximized when  $p = q$ , and so the Brier scoring rule is a strictly proper scoring rule. Note also that  $B(p, q)$  is a linear function in  $p$ . Hence, if  $p$  is drawn from a distribution, we have:  $\mathbb{E}_p[B(p, q)] = B(\mathbb{E}[p], q)$ , and so this is also maximized by reporting  $q = \mathbb{E}[p]$ .

### 2.3 Peer Prediction

Peer Prediction [19] with  $n$  agents receiving positively correlated binary signals  $\mathbf{b}$ , with symmetric prior  $Q$ , consists of the following mechanism  $\mathcal{M}(\hat{\mathbf{b}})$ :

1. Each agent  $i$  reports a signal  $\hat{b}_i$ .
2. Each agent  $i$  is uniformly randomly matched with an individual  $j \neq i$ , and is then paid  $PS(\hat{b}_j, q(1|\hat{b}_i))$ , where  $PS$  is a proper scoring rule.

That is, agent  $i$  is paid according to a proper scoring rule, based on  $i$ 's prediction that  $\hat{b}_j = 1$ , where  $i$ 's prediction is computed as either  $q(1|0)$  or  $q(1|1)$ , depending on  $i$ 's report to the mechanism. This can be thought of as having agent  $i$  bet on what agent  $j$ 's reported signal will be.

Notice that if agent  $j$  is truth-telling, then the Bayesian agent  $i$  would also be incentivized to truth-tell (strictly incentivized, if the proper scoring rule is strict). Agent  $i$ 's expected payoff (according to his own posterior distribution) for reporting his true type  $b_i$  has a premium compared to reporting  $-b_i$  of:

$$PS(\hat{b}_j, q(1|b_i)) - PS(\hat{b}_j, q(1|-b_i)) \geq 0$$

(strictly, for strict proper scoring rules) because we know that the expectation of  $PS(\hat{b}_j, \cdot)$  is (uniquely) maximized at  $q(1|b_i)$ . Now we introduce a convenient way to represent peer prediction mechanism.

**Definition 1 (Payoff Function Matrix).** *Each agent  $i$  who reports  $\hat{b}_i$  and is paired with agent  $j$  who reports  $\hat{b}_j$ , will be paid  $h_{\hat{b}_j, \hat{b}_i}$ . Then the peer prediction mechanism can be naturally represented as a  $2 \times 2$  matrix:*

$$\begin{pmatrix} h_{1,1} & h_{1,0} \\ h_{0,1} & h_{0,0} \end{pmatrix} = \begin{pmatrix} PS(1, q(1|1)) & PS(1, q(1|0)) \\ PS(0, q(1|1)) & PS(0, q(1|0)) \end{pmatrix}$$

which we call the payoff function matrix.



An example of a peer-prediction setting is included in the full version.

While truth-telling is always an equilibrium of the peer prediction mechanism, as we will see, it is not the only equilibrium. Two more equilibria are to always play 0 or always play 1. In Sect. 3.1, we further investigate equilibria of the peer prediction game. Based on the analysis of these multiple equilibria, we will develop a **modified peer prediction mechanism**, wherein players are paid according to the peer prediction based on a carefully-designed proper scoring rule, modulo some punishment imposed on the all playing 0 or all playing 1 strategy profiles. This modified mechanism will make the truth-telling equilibrium focal.

### 3 Summary of Main Results

In this section, we introduce our modified peer prediction mechanism and sketch the main theorem of this paper, that for almost any symmetric prior, there exists a modified peer prediction mechanism such that truth-telling is the focal equilibrium. Recall, we use the term *focal* to refer to an equilibrium with expected payoff strictly higher than that of any other Bayesian Nash equilibrium.

#### 3.1 Our Modified Peer Prediction Mechanism MPPM

Recall that modified peer prediction mechanism is the mechanism wherein players are paid according to peer prediction based on a carefully-designed proper scoring rule, modulo some punishment imposed on the all playing 0 or all playing 1 strategy profiles. So our approach differentiates between two types of equilibria:

**Definition 2 (Informative Strategy).** *We call always reporting 1 and always reporting 0 uninformative strategies; we call all other strategies (equilibria) informative.*

*Designing the Optimal Peer Prediction Mechanism.* We start to describe our modified peer prediction mechanism MPPM. We use two steps to design our MPPM. First we define the PPM:

**Definition 3.** *Given any binary, symmetric, positively correlated, and signal asymmetric prior  $Q$ , with  $q(1|1) > q(0|0)$  (the  $q(0|0) < q(1|1)$  case is analogous), we first design our peer prediction mechanism  $PPM(Q)$  and represent it as a payoff function matrix (See Definition 1).  $PPM(Q)$  depends on the region that  $Q$  belongs to, we defer the definitions of regions  $R_1, R_2, R_3$  to full version but provide Fig. 1 here to illustrate them.*

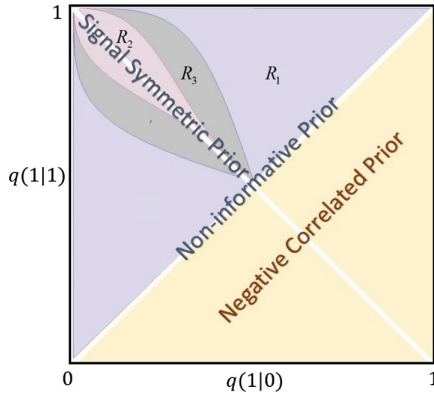
1. If  $Q \in R_1$ , then  $PPM(Q) = \mathcal{M}_1(Q)$
2. If  $Q \in R_2$ , then  $PPM(Q) = \mathcal{M}_2(Q)$
3. If  $Q \in R_3$ , then we pick a small number  $\epsilon > 0$  and  $PPM(Q, \epsilon) = \mathcal{M}_3(Q, \epsilon)$

where

$$\mathcal{M}_1(Q) = \begin{pmatrix} \zeta(Q) & 0 \\ 0 & 1 \end{pmatrix}, \mathcal{M}_2(Q) = \begin{pmatrix} 1 & 0 \\ 0 & \eta(Q) \end{pmatrix}, \mathcal{M}_3(Q, \epsilon) = \begin{pmatrix} \zeta(Q, \epsilon) & \delta(Q, \epsilon) \\ 0 & 1 \end{pmatrix}$$

and

$0 \leq \zeta(Q), \eta(Q) \leq 1$  are constants that only depend on common prior  $Q$ .  $0 \leq \zeta(Q, \epsilon), \delta(Q, \epsilon) \leq 1$  are constants that only depend on common prior  $Q$  and  $\epsilon > 0$ .<sup>1</sup>



**Fig. 1.** The regions  $R_1, R_2, R_3$  are good “priors” where we can make truth-telling focal when the number of agents is sufficient large. The white diagonals are “bad” priors we cannot make truth-telling focal. In the top-right to bottom-left diagonal,  $q(1|0) = q(1|1)$ , so the private signal does not have any information. We call this diagonal the set of non-informative priors. In the top-left to bottom-right diagonal,  $q(0|0) = q(1|1)$  (actually we can see  $q(0|0) = q(1|1)$  iff  $q(0) = q(1)$  via some calculations). This diagonal is the set of signal symmetric priors. The yellow region is the set of the negative correlated priors. (Color figure online)

Note that actually  $PPM(Q)$  is a quite simple mechanism. We use region  $R_1$  as example: if the prior belongs to region  $R_1$ , for every  $i$ , agent  $i$  will receive 0 payment if the agent paired with agent  $i$ , call him agent  $j$ , reports a different signal than him. If both agent  $i$  and agent  $j$  report 1, agent  $i$  will receive a payment of  $0 \leq \zeta(Q) \leq 1$ , if both agent  $i$  and agent  $j$  report 0, agent  $i$  will receive payment of 1.

Actually for regions  $R_1, R_2$ , the  $PPM(Q)$  we define here is the optimal peer prediction mechanism in that it maximizes the advantage of truth-telling over the informative equilibria which have the second largest expected payoff over all Peer-prediction mechanisms with payoffs in  $[0, 1]$ . For region  $R_3$ , the optimal peer prediction mechanism does not exist, but the advantage of the  $PPM(Q, \epsilon)$  we define approaches the optimal advantage as  $\epsilon$  goes to 0.

<sup>1</sup> Explicit statement in full version.

**Definition 4.** We define  $\Delta^*(Q)$  to be the supremum of the advantage of truth-telling over the informative equilibria which have the second largest expected payoff over all Peer-prediction mechanisms with payoffs in  $[0, 1]$ .

*Add Punishment.* In our  $PPM(Q)$ , an uninformative strategy can still obtain the highest payoff. For example, in mechanism  $\mathcal{M}_1$ , agents will receive maximal payment 1 by simply always reporting 0.

Our final  $MPPM(Q)$  Mechanism is the same as the  $PPM(Q)$  except that we add a punishment designed to hurt the all 0 or all 1 equilibria.

**Definition 5.** Our Modified Peer-Prediction Mechanism  $MPPM(Q)$  (or  $MPPM(Q, \epsilon)$  has payoffs identical to  $PPM(Q)$  (or  $MPPM(Q, \epsilon)$ ) except that, in the event all the other agents play all  $\mathbf{0}$  or all  $\mathbf{1}$ , it will issue an agent a punishment of  $p = \frac{1-t}{2(1-\epsilon_Q)} + \frac{\Delta^*(Q)}{2\epsilon_Q}$  where  $\epsilon_Q$  is the maximum probability that a fixed set of  $n - 1$  agents receive the same signal (either all  $\mathbf{0}$  or all  $\mathbf{1}$ );  $t$  is the expected of payoff of truth-telling  $\mathbf{T}$  in the  $PPM(Q)$ , and  $\Delta^*(Q)$  is as defined in Definition 4.

To make truth-telling focal, we would like to impose a punishment to the agents if everyone reports the same signal. However, such a punishment may distort the equilibria of the mechanism. To avoid this, we punish an agent by  $p$  when all the other agents report the same signal. Because an agent’s strategy does not influence his punishment, his marginal benefit for deviation remains the same and so the equilibrium remain the same. However, while all  $\mathbf{0}$  and all  $\mathbf{1}$  remain equilibrium, in them,  $MPPM(Q)$  will punish each agent by  $p$ .

A difficulty arises: if the number of agents is too small like 2 or 3, it is possible (and even probable) that all agents report their true signals, yet are still punished by the  $MPPM(Q)$  mechanism. Punishments like this might distort the payoffs among the informative equilibrium. However, if  $\epsilon_Q$  (the probability that  $n - 1$  agents receives the same signal) is sufficient small, this is no longer a problem. For most reasonable priors, as the number of agents increases,  $\epsilon_Q$  will go to zero. Formally we will need that the number of agents is large enough such that  $\epsilon_Q < \frac{\Delta^*(Q)}{1-t+\Delta^*(Q)}$ .

If the number of agents is too small such that  $\epsilon_Q \leq \frac{\Delta^*(Q)}{1-t+\Delta^*(Q)}$ , we cannot show that  $MPPM(Q)$  has truth-telling as a focal equilibrium.

In particular, we can see if  $\epsilon_Q \rightarrow 0$  (say as the number of agents increases), then at some point, truth-telling will be focal. We know that such a limit is necessary because, for example, with two agents making truth-telling focal is impossible.

Note that if the prior tells us the probability of a 1 event is concentrated far away from 0 and 1, the number of agents we need to make truth-telling focal will be very small since uninformative equilibria (all 1 and all 0) are far away from truth-telling.

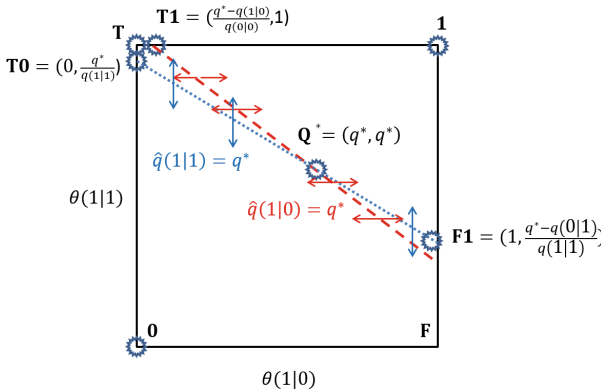
**Theorem 5.** (Main Theorem (Informal)) Let  $Q$  be a binary, symmetric, positively correlated and signal asymmetric prior, and let  $\epsilon_Q$  be the maximum probability that a fixed set of  $n - 1$  agents receive the same signal (either all  $\mathbf{0}$  or all  $\mathbf{1}$ ). Then

1. In our PPM, truth-telling has the largest expected payoff among all informative equilibria. Moreover, over the space of Peer-Prediction mechanisms, our PPM(Q) maximizes the advantage truth-telling has over the informative equilibrium which have the second largest expected payoff, over all Peer-prediction mechanisms with payoffs in [0, 1] for regions R<sub>1</sub>, R<sub>2</sub> and PPM(Q, ε) approaches the maximal advantage for region R<sub>3</sub> as ε goes to 0.
2. There exists a constant  $\xi_{q(1|1),q(1|0)}$  which only depends on  $q(1|1)$  and  $q(1|0)$  such that, if  $\epsilon_Q < \xi$ , our MPPM(Q) makes truth-telling focal.

Now we list all equilibria of the peer prediction mechanism in the below theorem (Fig. 2).

**Definition 6.** For a prior Q, proper scoring rule PS, and a binary signal space, we define  $q^*$  to be the fraction of other agents reporting 1 that would make an agent indifferent between reporting 0 or 1, i.e.,

$$q^* := \{p \mid PS(p, q(b|1)) = PS(p, q(b|0)), 0 \leq p \leq 1\}.$$



**Fig. 2.** Illustration of the 7 equilibria of a peer prediction mechanism under a specific scoring rule (see the full version). Note that to the right of the dashed red line where  $\hat{q}(1|0) = q^*$ , the best response is to increase  $\theta(1|0)$ ; to the left of the dashed red line, the best response is to decrease  $\theta(1|0)$ ; and on the line an agent is indifferent. Similarly, above the dotted blue line where  $\hat{q}(1|1) = q^*$ , the best response is to increase  $\theta(1|1)$ ; below the dotted blue line, the best response is to decrease  $\theta(1|1)$ ; and on the line an agent is indifferent. (Color figure online)

**Theorem 6.** Let Q be a symmetric and positively correlated prior on  $\{0, 1\}^n$ , and let M be a peer-prediction mechanism run with a strictly proper scoring rule with break-even  $q^*$  (Definition 6). Then there are no asymmetric equilibria. All equilibria are symmetric and depend only on  $q^*$ ; they are

$$0, 1, T, Q^* \triangleq \begin{pmatrix} 0 \rightarrow q^* \\ 1 \rightarrow q^* \end{pmatrix},$$

$$\mathbf{T0} \triangleq \begin{pmatrix} 0 \rightarrow 0 \\ 1 \rightarrow \frac{q^*}{q(1|1)} \end{pmatrix}, \mathbf{T1} \triangleq \begin{pmatrix} 0 \rightarrow \frac{q^* - q(1|0)}{q(0|0)} \\ 1 \rightarrow 1 \end{pmatrix}$$

and also conditionally include

$$\mathbf{F} \text{ if } q(0|1) \leq q^* \leq q(0|0) \tag{3}$$

$$\mathbf{F1} \triangleq \begin{pmatrix} 0 \rightarrow 1 \\ 1 \rightarrow \frac{q^* - q(0|1)}{q(1|1)} \end{pmatrix} \text{ if } q(0|1) \leq q^* \tag{4}$$

$$\mathbf{F0} \triangleq \begin{pmatrix} 0 \rightarrow \frac{q^*}{q(0|0)} \\ 1 \rightarrow 0 \end{pmatrix} \text{ if } q^* \leq q(0|0) \tag{5}$$

Due to the space limitation, we defer all proofs to our full version (see <https://arxiv.org/abs/1603.07319>).

**Acknowledgments.** We thank David Parkes and Paul Resnick for helpful suggestions.

## References

1. Crémer, J., McLean, R.P.: Optimal selling strategies under uncertainty for a discriminating monopolist when demands are interdependent. *Econometrica* **53**(2), 345–361 (1985)
2. Crémer, J., McLean, R.P.: Full extraction of the surplus in bayesian and dominant strategy auctions. *Econometrica* **56**(6), 1247–1257 (1988)
3. Dasgupta, A., Ghosh, A.: Crowdsourced judgement elicitation with endogenous proficiency. In: Proceedings of the 22nd International Conference on World Wide Web, pp. 319–330. International World Wide Web Conferences Steering Committee (2013)
4. d’Aspremont, C., Gérard-Varet, L.A.: Bayesian incentive compatible beliefs. *J. Math. Econ.* **10**(1), 83–103 (1982)
5. d’Aspremont, C., Grard-Varet, L.A.: Incentives and incomplete information. *J. Publ. Econ.* **11**(1), 25–45 (1979)
6. Faltings, B., Jurca, R., Pu, P., Tran, B.D.: Incentives to counter bias in human computation. In: Second AAAI Conference on Human Computation and Crowdsourcing (2014)
7. Gao, X.A., Mao, A., Chen, Y., Adams, R.P.: Trick or treat: putting peer prediction to the test. In: Proceedings of the Fifteenth ACM Conference on Economics and Computation, pp. 507–524. ACM (2014)
8. Ghosh, A., Ligett, K., Roth, A., Schoenebeck, G.: Buying private data without verification. In: Proceedings of the Fifteenth ACM Conference on Economics and Computation, pp. 931–948. ACM (2014)
9. Goel, S., Reeves, D.M., Pennock, D.M.: Collective revelation: a mechanism for self-verified, weighted, and truthful predictions. In: Proceedings of the 10th ACM Conference on Electronic Commerce (EC 2009) (2009)
10. Jurca, R., Faltings, B.: Incentives for answering hypothetical questions. In: Proceedings of the 1st Workshop on Social Computing and User Generated Content (SC 2011). ACM (2011)

11. Jurca, R., Faltings, B.: Minimum payments that reward honest reputation feedback. In: Proceedings of the 7th ACM Conference on Electronic Commerce (EC 2006) (2006)
12. Jurca, R., Faltings, B.: Collusion-resistant, incentive-compatible feedback payments. In: Proceedings of the 8th ACM Conference on Electronic Commerce, pp. 200–209. ACM (2007)
13. Jurca, R., Faltings, B.: Incentives for expressing opinions in online polls. In: Proceedings of the 9th ACM Conference on Electronic Commerce (EC 2008) (2008)
14. Jurca, R., Faltings, B.: Mechanisms for making crowds truthful. *J. Artif. Int. Res.* **34**(1), 209 (2009)
15. Kamble, V., Shah, N., Marn, D., Parekh, A., Ramachandran, K.: Truth serums for massively crowdsourced evaluation tasks. arXiv preprint [arXiv:1507.07045](https://arxiv.org/abs/1507.07045)
16. Kong, Y., Schoenebeck, G.: A framework for designing information elicitation mechanisms that reward truth-telling. ArXiv e-prints, May 2016
17. Kong, Y., Schoenebeck, G.: Equilibrium selection in information elicitation without verification via information monotonicity. ArXiv e-prints, March 2016
18. Lambert, N., Shoham, Y.: Truthful surveys. In: Papadimitriou, C., Zhang, S. (eds.) WINE 2008. LNCS, vol. 5385, pp. 154–165. Springer, Heidelberg (2008). doi:[10.1007/978-3-540-92185-1\\_23](https://doi.org/10.1007/978-3-540-92185-1_23)
19. Miller, N., Resnick, P., Zeckhauser, R.: Eliciting informative feedback: the peer-prediction method. *Manag. Sci.* **51**(9), 1359–1373 (2005)
20. Prelec, D.: A Bayesian truth serum for subjective data. *Science* **306**(5695), 462–466 (2004)
21. Radanovic, G., Faltings, B.: Incentive schemes for participatory sensing. In: Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems, pp. 1081–1089. International Foundation for Autonomous Agents and Multiagent Systems (2015)
22. Riley, B.: Minimum truth serums with optional predictions. In: Proceedings of the 4th Workshop on Social Computing and User Generated Content (SC14) (2014)
23. Shnayder, V., Agarwal, A., Frongillo, R., Parkes, D.C.: Informed truthfulness in multi-task peer prediction. ArXiv e-prints, March 2016
24. Witkowski, J., Parkes, D.C.: Peer prediction without a common prior. In: Proceedings of the 13th ACM Conference on Electronic Commerce, pp. 964–981. ACM (2012)
25. Witkowski, J., Parkes, D.C.: Learning the prior in minimal peer prediction. In: Proceedings of the 3rd Workshop on Social Computing and User Generated Content at the ACM Conference on Electronic Commerce, p. 14. Citeseer (2013)