

CS269I: Incentives in Computer Science

Lecture #17: Scoring Rules and Peer Prediction (Incentivizing Honest Forecasts and Feedback)*

Tim Roughgarden[†]

November 28, 2016

1 Scoring Rules

1.1 Motivation

We next consider the goal of eliciting a good prediction of an uncertain event. For example:

1. You might ask a weather forecaster: will it rain tomorrow?
2. You might ask a political pundit: will a Democrat or a Republican win the next election?
3. You might ask a Microsoft employee: will the next version of MS Office get shipped on time, or will it be delayed?

It's deceptively tricky to evaluate the quality of a prediction. For example, suppose a weather forecaster declares “30% chance of rain tomorrow,” and then it rains. Was the forecast bad, or did the forecaster just get unlucky?

For example, one simplistic approach would be to call a prediction “good” if it puts more than 50% probability on the outcome that actually occurred, and “bad” otherwise. But this is a weird rule—a forecaster then has no reason to say anything other than “100%” or “0%.” (Since “51%” is treated the same as “100%,” anyways.) So why should a forecaster bother formulating a more nuanced prediction? Evidently, we want a less binary notion of a “good” prediction, with the quality of the prediction increasing with the amount of probability assigned to the outcome that occurred.

To formalize this idea, let X denote a finite set of all possible outcomes. All of our examples so far involve a binary event (“rain” or “shine,” “Democrat” or “Republican,” “on

*©2016, Tim Roughgarden.

[†]Department of Computer Science, Stanford University, 474 Gates Building, 353 Serra Mall, Stanford, CA 94305. Email: tim@cs.stanford.edu.

time” or “delayed”), meaning that X has two elements. It’s fine to think only about binary outcomes for the rest of this and next lecture, even though the results that we’ll cover apply more generally.

Definition 1.1 (Scoring Rule) A *scoring rule* is a real-valued function of the form $S(\mathbf{q}, i)$, where \mathbf{q} is a probability distribution over X (a prediction) and i is an outcome of X (the realized outcome).

For example, \mathbf{q} could be “30% chance of rain” and i could be “rain.” Our binary scoring rule above corresponds to setting $S(\mathbf{q}, i) = 1$ if $q_i > \frac{1}{2}$ and 0 otherwise (where q_i denotes the probability ascribed to outcome i in the distribution \mathbf{q}).

1.2 Strictly Proper Scoring Rules

Our next goal is a notion of a “truthful” scoring rule. To formalize this, we need to specify what motivates forecasters—i.e., their preferences. Here’s the model:

1. A forecaster has a “belief” \mathbf{p} , which is a probability distribution over X .¹
2. A forecaster wants to choose her prediction \mathbf{q} in order to maximize her score. (Maybe \mathbf{q} equals the forecaster’s true belief \mathbf{p} , or maybe not.) Actually this doesn’t quite make sense, since the forecaster’s score will depend on the realized outcome, which is not under the forecaster’s control. So we assume that the forecaster wants to maximize her *expected* score:

$$\max_{\mathbf{q}} \mathbf{E}_{i \sim \mathbf{p}}[S(\mathbf{q}, i)], \quad (1)$$

where the expectation is with respect to the distribution \mathbf{p} over outcomes that the forecaster believes is the true one.

Note we are assuming that a forecaster cares about her (expected) score. This could be because the score represents a dollar reward, or some other motivating currency like reputation points.

Given this model of what a forecaster wants, we are in a position to define “truthful” scoring rules (where setting $\mathbf{q} = \mathbf{p}$ is the best course of action). For historical reasons, such rules are called *strictly proper* scoring rules.

Here is the key definition to know about scoring rules:

Definition 1.2 (Strictly Proper Scoring Rule) A scoring rule S is *strictly proper* if, no matter what the true belief \mathbf{p} of the forecaster is, her unique best response is to report truthfully (i.e., to set $\mathbf{q} = \mathbf{p}$).

One can also define (weakly) proper scoring rules, where truthful reporting is one best response, perhaps among many. Do you see why this is not an interesting definition? Because even a constant function (like $S(\mathbf{q}, i) = 0$ for all \mathbf{q} and i) satisfies this definition. (Our first, binary, scoring rule is also weakly proper.)

¹This belief plays the same role as the “prior distribution” that we adopted last lecture to think about revenue-maximizing auction design.

1.3 A Non-Example

Let's try to find a strictly proper scoring rule. Recall that we want to reward a forecaster according to the strength of their prediction on the outcome that actually occurred. Maybe the first thing to try is a reward linear in the prediction probability:

$$S(\mathbf{q}, i) = q_i,$$

where q_i denotes the probability assigned to outcome $i \in X$ in the prediction \mathbf{q} . (In general, think of \mathbf{q} as a nonnegative vector, indexed by the outcomes of X , with $\sum_{i \in X} q_i = 1$.)

Is this linear scoring rule strictly proper? Actually, it's not even a weaker proper scoring rule! This rule incentivizes a forecaster to put 100% of her prediction on the outcome that she thinks is the most likely (see Exercise Set #9).

1.4 The Quadratic Scoring Rule

Uh oh — we've seen a number of impossibility results in this class, could there be another one lurking here? Fortunately not, as there are several natural strictly proper scoring rules. The first one is the *quadratic scoring rule*.² Its definition is:³

$$S(\mathbf{q}, i) = q_i - \frac{1}{2} \sum_{j \in X} q_j^2. \quad (2)$$

Thus the quadratic rule includes the same linear term as before, but also a quadratic penalty term designed to penalize extreme reports. (Extreme reports were the problem with the linear rule, remember?) The key difference is that while the score is still increasing with the probability assigned to the realized outcome, the rate of increase decreases as the report gets more extreme.

For example, if the forecaster assigns some event 100% probability and the predicted event occurs, her score is $\frac{1}{2}$. If the predicted event doesn't occur, then her score is $-\frac{1}{2}$. If the forecaster reports the uniform distribution, then no matter what outcome happens, her score is $\frac{1}{n} - \frac{1}{2} \cdot n \cdot \frac{1}{n^2} = \frac{1}{2n}$.

Proposition 1.3 *The quadratic scoring rule is strictly proper.*

Proof: The proof is just calculus. Fix a forecaster with a belief \mathbf{p} . For this fixed \mathbf{p} , the expected score (1) of the forecaster when reporting \mathbf{q} evaluates to

$$\sum_{i \in X} p_i q_i - \frac{1}{2} \sum_{i \in X} p_i \sum_{j \in X} q_j^2. \quad (3)$$

²Discovered by Brier in 1950 [2] who, believe it or not, really was interested in how to incentivize weather forecasters to produce the best-possible reports. (The article appeared in the *Monthly Weather Review*.)

³Often the quadratic rule is defined as twice the quantity in (2). This scaling factor has no effect on the rule's important properties; see also Proposition 1.4.

This function is strictly concave in \mathbf{q} (the first term is linear in \mathbf{q} , and the second term, as a negative quadratic, is strictly concave). This means that the function has a unique maximizer. So what is it? To identify it, let's examine the partial derivatives of the function (3). For $h \in X$, we have

$$\frac{d}{dq_h} = p_h - \sum_{i \in X} p_i q_h.$$

Since $\sum_{i \in X} p_i = 1$, the partial derivative is simply $p_h - q_h$. This means that setting $\mathbf{q} = \mathbf{p}$ zeroes out all of the partial derivatives (and is the only way to zero them all out), implying that \mathbf{q} is the unique maximum of (3). ■

Several comments on the proof. First, checking that all derivatives are zero (a “first-order condition”) is generally only *necessary* for optimality, and need not be sufficient. But for a strictly concave function like (3), the first-order conditions are also sufficient for optimality. Second, the proof shows something stronger than what is claimed in Proposition 1.3. Namely, setting $\mathbf{q} = \mathbf{p}$ is the unique maximizer of (3) over *all real-valued vectors* \mathbf{q} , not just over probability distributions \mathbf{q} . This is because we only used the first-order conditions of the unconstrained version of the problem of maximizing (3), rather than the constrained version which only considers probability distributions \mathbf{q} . Of course, if $\mathbf{q} = \mathbf{p}$ is optimal over all real-valued vectors, it is optimal in particular among all probability distributions.

Given one strictly proper scoring rule, one can construct others by shifting and scaling (i.e., by affine transformations).

Proposition 1.4 *If S is a strictly proper scoring rule, $a > 0$, and $b \in \mathbb{R}$, then $aS + b$ is also a strictly proper scoring rule.*

The reason is just that neither shifting everything by a constant b nor scaling everything by a positive constant $a > 0$ has any effect on a forecaster's best response. So if S is strictly proper (i.e., the unique best response is truthful), then so is $aS + b$. Of course, the score actually assigned to the forecaster varies with the choice of a and b .

One application of the above shifting trick is to shift a scoring rule so that forecasters are guaranteed nonnegative utility (maybe in expectation, or maybe always). For example, adding $\frac{1}{2}$ to the quadratic scoring rule ensures that it is always nonnegative.

1.5 The Logarithmic Scoring Rule

One curiosity about the quadratic scoring rule (2) is that the score assigned to a prediction depends not only on the report q_i on the outcome that actually occurred, but also on the distribution of \mathbf{q} on the outcomes that did not occur. It is arguably unnatural to do this—it's not clear why a scoring rule should have jurisdiction over the different predictions on the outcomes that did not occur. Is this oddity an inevitable consequence of the strongly proper condition?

Our second rule, discovered by Good in 1952 [6], is the *logarithmic scoring rule*. It is simply:

$$S(\mathbf{q}, i) = \ln q_i. \tag{4}$$

Several comments. First, note that the score assigned to a prediction depends *only* on the probability that the forecaster assigned to the outcome that occurred, and not on the probabilities assigned to the other outcomes. Second, note that the base of the logarithm doesn't really matter, since different logarithms differ by a constant factor (and strict properness is preserved by scaling). Third, the scoring rule as defined in (4) is never positive, so typically a shifted version of it is used. For example, adding the constant $\ln |X|$ to the rule in (4) ensures that a forecaster can guarantee herself nonnegative utility (by reporting the uniform distribution). Finally, note that the logarithmic scoring rule is not bounded below. If a forecaster assigns 0 to some outcome and that outcome actually transpires, then the forecaster's score is $-\infty$. (Of course, a forecaster who is 100% convinced that an outcome is impossible won't care what score she receives in that (impossible) case.) If this is undesirable, the rule can be modified by imposing a small but positive lower bound on all reported probabilities.

The logarithmic scoring rule is also a strictly proper rule.

Proposition 1.5 *The logarithmic scoring rule is strictly proper.*

Proof: Again, the proof is just calculus; we provide a sketch of the argument. Fix a forecaster with belief \mathbf{p} . (For clarity, assume that $p_i > 0$ for every $i \in X$, even though this is not necessary for the proof.) We again begin by writing down the expected score of a given report \mathbf{q} (cf., (3)):

$$\mathbf{E}_{i \sim \mathbf{p}}[S(\mathbf{q}, i)] = \sum_{i \in X} p_i \ln q_i. \quad (5)$$

Like (3), this expression is strictly concave in \mathbf{q} and hence has a unique maximizer. So what is it?

Again, we consider the partial derivatives, which are

$$\frac{d}{dq_i} = \frac{p_i}{q_i} \quad (6)$$

for all $i \in X$. Unlike in the proof of Proposition 1.3, there is no value of \mathbf{q} that zeroes out these derivatives. (Without the constraint that \mathbf{q} is a probability distribution, the best response is to set all entries of \mathbf{q} as high as possible.) So instead we need to identify the report \mathbf{q} that, among all probability distributions, maximizes (5). We claim that \mathbf{q} is optimal only if all of the partial derivatives $\frac{d}{dq_i}$ are equal. For suppose one (that of i) was bigger than another (that of j) — then shifting a little bit of probability mass from j to i would yield a new distribution \mathbf{q}' with a strictly larger expected score (5).⁴ Concavity of (5) implies that the converse is also true — if \mathbf{q} equalizes the derivatives (6), then \mathbf{q} must be optimal (further details omitted). The only way to equalize the partial derivatives in (6) is to set \mathbf{q} proportional to \mathbf{p} . The unique probability distribution with this property is \mathbf{p} . ■

⁴By inspection of (6), this can only happen when $q_j > 0$ and hence there is mass available to move.

1.6 Final Comments

There are many strictly proper scoring rules beyond the quadratic and logarithmic rules; you'll encounter one of them on Exercise Set #9.

Which is better, the quadratic or logarithmic scoring rule? In general, there is no clear answer. If you're bothered by the fact that the quadratic scoring rule makes use of probability reports on the unrealized outcomes, you might prefer the logarithmic rule. If you're bothered by the fact that the logarithmic rule is very sensitive to changes of small probabilities (and equals $-\infty$ for "impossible" events), then you might prefer the quadratic rule. Both rules have been implemented in practice, though the logarithmic rule has been more widely used (especially in the context of prediction markets, discussed in the next lecture). In experiments, both seem to do a fine job of eliciting truthful predictions (see e.g. [1]).

Scoring rules are a neat idea, but are they actually useful for anything? We'll see two applications: in this lecture, to the problem of incentivizing honest feedback; and next lecture, to the design of prediction markets.

2 Incentivizing Honest Feedback

2.1 Motivating Examples

Suppose you ask someone to rate a movie, on a scale of 1 to 5. Can we use a scoring rule to incentivize the reviewer to state their true opinion? Not immediately — the issue is that scoring rules rely on the realization of some verifiable "ground truth" outcome. This assumption can fail for two conceptually different reasons: first, if there simply is no ground truth; second, if there is a ground truth but it is too costly to determine. The assumption of a verifiable outcome is fine for, say, weather forecasts, but it doesn't seem appropriate for scoring subjective opinions like movie ratings.

A second relevant example is peer grading (where students grade the assignments of other students), especially at large scale, for example in a MOOC ("massive open online course"). One can imagine a seasoned instructor supplying the "ground truth grade" of an assignment, and then using a scoring rule to score the grades assigned by students to the assignment. But this defeats the whole point of peer grading in a massive course (where the instructor cannot possibly evaluate all of the assignments).

So how is peer grading done in MOOCs now? Without any consideration of incentives. For example, a student's assignment might be graded by five peers, with the student receiving the median of these five grades. Sometimes grades get reweighted according to the apparent accuracy of the grader.

Unsurprisingly, in practice there's quite a bit of variance in peer grades (though the median grade is often surprisingly accurate [9]). What if we wanted to explicitly incentivize accurate grading, without direct verification? For example, one could make grading quality part of a student's grade, or introduce a reputation system to publicly recognize good (and maybe bad) graders.

2.2 The Model

Here's our setup for the problem of incentivizing honest feedback when there is no verifiable ground truth:

- There are n players (e.g., graders of an assignment in a MOOC).
- Player i has a “signal” s_i . This signal is supposed to represent the useful information possessed by player i , such as a grader's true opinion of the quality of an assignment (after a close inspection).⁵
- Each player i submits a report r_i (like a grade) to a mechanism (like a MOOC platform). The report r_i can be equal to the player's signal s_i , or not.
- The mechanism pays player i $\pi_i(r_1, \dots, r_n)$.

Several comments. First, the “payment” π_i might be in money, or it might be in some other currency that the player cares about (like extra credit points or reputation). Currently in most MOOCs, π_i is just zero. In any case, we assume that player i wants to maximize $\pi_i(r_1, \dots, r_n)$. Finally, observe that the decisions made by the mechanism (i.e., the payments π_i) depend on the only information that it has, the players' reports r_1, \dots, r_n (and not on any “ground truth”).

A player i wants to choose her report r_i to maximize her payment $\pi_i(r_1, \dots, r_n)$ from the mechanism. But this payment depends on the (presumably unknown) reports of the other players. So how should player i compare the relative benefits of different reports? Analogously to last lecture (for revenue-maximizing auctions), we proceed by assuming a prior distribution over signal profiles.

Precisely, we will assume that the signal profile (s_1, \dots, s_n) is drawn from a distribution D , and that D is known to all of the players. For example, with two players and a binary signal space, such a distribution might look like

	$s_2 = 0$	$s_2 = 1$
$s_1 = 0$.3	.1
$s_1 = 1$.1	.5

We also assume that D is symmetric, meaning that (i) zooming in on two players i and j , the joint distribution of their two signals can be represented by a symmetric matrix, as above; and (ii) this symmetric matrix is the same for every pair i, j of distinct players.

In our example distribution above, the signals of the two players are *not* independent. (This stands in contrast to the auction model discussed last lecture, where we assumed that bidders had independent valuations.) For example, knowing nothing about s_1 , the probabilities that s_2 is 0 or 1 are $\frac{2}{5}$ and $\frac{3}{5}$, respectively. But if we condition on the event that $s_1 = 0$ (e.g., that one TA gives a bad grade to an assignment), then the (conditional)

⁵The simple case where $s_i \in \{0, 1\}$ is already interesting, but the results we'll cover only assume that s_i belongs to some finite set (like $\{A, B, C, D, F\}$).

probabilities that s_2 is 0 or 1 become $\frac{3}{4}$ and $\frac{1}{4}$, respectively. That is, it becomes much more likely that the other TA will give the assignment a bad grade. In our motivating applications, it makes sense that players would have correlated signals. For example, players' signals could be perturbed or noisy versions of some unobserved ground truth, like the true quality of an assignment.

2.3 Output Agreement

So how should we choose the payment functions π_1, \dots, π_n to incentivize truthful reporting by players? The first idea is to reward agreement between players. This makes intuitive sense if we're hoping to extract a rough consensus from players' reports. Formally, here's how the mechanism works:

Output Agreement

1. For each player i :
 - (a) Pick a random player $j \neq i$.
 - (b) Set i 's payoff π_i equal to 1 if $r_i = r_j$, and to 0 if $r_i \neq r_j$.

The output agreement mechanism will be familiar if you've ever played the ESP game (the canonical "game with a purpose"). The point of the ESP game is to make it fun for humans to annotate images, thereby generating a nice labeled data set for supervised machine learning algorithms. The ESP game works by taking two random people, showing them the same image, asking them to type in descriptive words for the image, and rewarding them (in a fictitious currency) whenever they type in the same word.⁶ This is exactly the output agreement mechanism.

Is the output agreement mechanism truthful? That is, does a player maximize her expected reward by reporting her true signal? Does this guarantee hold at least in the case where all other players are reporting truthfully?

The answer depends on the prior distribution D . In the example above, the output agreement mechanism is truthful, meaning that if all other players are reporting truthfully, then truthful reporting is the unique best response for a player. For example, if player 1 receives a signal of 0, then conditionally, player 2 is more likely to have signal 0 ($\frac{3}{4}$ probability) than 1 ($\frac{1}{4}$ probability). This means it's a best response for the first player to report a 0, which is a truthful report. Similarly, if $s_1 = 1$, then the second player is more likely to have signal 1 ($\frac{5}{6}$ probability) than 0 ($\frac{1}{6}$), so it is again a best response to (truthfully) report 1.

In general, the output agreement mechanism is not truthful. For example, suppose we modify the prior distribution above as follows:

⁶A nice twist is to list already-discovered descriptors as "taboo words," thereby forcing the two people to come up with a new and useful descriptor for the image.

	$s_2 = 0$	$s_2 = 1$
$s_1 = 0$.1	.2
$s_1 = 1$.2	.5

Suppose the first player receives a signal of 0. The conditional probability that the second player also received a signal of 0 is $\frac{1}{3}$. This is higher than the unconditional probability ($\frac{3}{10}$), but it's still more likely that the second player received a signal of 1. Thus the best response of the first player is to report a 1, even though her signal was a 0.

In general, the output agreement mechanism is truthful if and only if the most likely outcome is a matching of signals—if and only if $\Pr[s_2 = x | s_1 = x] > \Pr[s_2 = y | s_1 = x]$ for every x and $y \neq x$. In terms of the defining matrix of the distribution D , this means that each diagonal entry should be larger than the sum of all of the other entries in the same row (or column). This condition holds in our first example, but not in our second example.

Is there a mechanism that incentivizes truthful reporting for *all* symmetric priors?

2.4 The Peer Prediction Mechanism

We next make an additional assumption, that the distribution D over signals is known to the mechanism. (By contrast, the output agreement mechanism is well defined even without knowledge of the prior.) This assumption may or may not be reasonable, depending on the setting. In the peer grading example, for CS161 (undergrad algorithms), which I've taught many times and with hundreds of students, I have a pretty good sense of the distribution (of the “ground truth” grades, and of TA's grades given the ground truth grade). In CS269I, at least in this first offering, I'm not so confident about the distribution D of signals.

We next present the elegant *peer prediction* mechanism, due to Miller et al. [8]. The high-level idea is to treat a player's report as a prediction of other players' reports, and then judge this prediction using a strictly proper scoring rule S . Formally, fix such an S (quadratic, or logarithmic, etc.). The mechanism is:

Peer Prediction

1. For each player i :
 - (a) Pick a random player $j \neq i$.
 - (b) Let $D_j(r_i)$ denote the distribution of s_j , conditioned on the event that $s_i = r_i$.⁷
 - (c) Set i 's payoff π_i to

$$S(\underbrace{D_j(r_i)}_{\text{prediction}}, \underbrace{r_j}_{\text{realization}}).$$

⁷In our first example of this section, $D_j(0)$ would be the $\frac{3}{4}$ - $\frac{1}{4}$ distribution, while $D_j(1)$ would be the $\frac{1}{6}$ - $\frac{5}{6}$ distribution.

In effect, the mechanism takes player i 's report r_i at face value (using r_i as a proxy for s_i), which induces a conditional probability distribution on s_j (and r_j , if j is reporting truthfully).

The Peer Prediction mechanism works (i.e., is truthful) under a weak condition: distinct reports (r_i 's) should induce distinct conditional distributions ($D_j(r_i)$'s). For example, my conditional distribution on the grade that will be assigned by the second TA is different for each grade (in $\{A,B,C,D,F\}$) that could be assigned by the first TA. This assumption is more or less without loss of generality: if two signals induce the same conditional distribution over others' signals, you may as well combine the two signals into one.

The next proposition states that truthful reporting is an equilibrium of the Peer Prediction mechanism.

Proposition 2.1 *In the Peer Prediction mechanism, if every player other than i reports truthfully, then the unique best response of player i is to report truthfully (i.e., set $r_i = s_i$).*

Proof: Since every player $j \neq i$ is assumed truthful, $r_j = s_j$ for every $j \neq i$. Thus, given that i has the signal s_i , the conditional distribution on r_j (and s_j) is exactly $D_j(s_i)$.

Since S is a strictly proper scoring rule, the (unique) best-case scenario for the payoff $S(D_j(r_i), r_j)$ is when $D_j(r_i)$ is exactly the distribution from which r_j is drawn (i.e., $D_j(s_i)$). Because distinct reports induce distinct conditional distributions on others' signals, the unique way to achieve this best-case scenario is to set $r_i = s_i$. ■

2.5 Implementation Challenges

The Peer Prediction mechanism is a neat idea, but there are some obstacles to implementing it directly in practice (which have motivated much follow-up work in the last 10 years). First, we already mentioned the drawback that the mechanism requires advance knowledge of the prior distribution D over signal profiles. We already noted that this assumption is sometimes palatable, and sometimes not. A number of subsequently proposed mechanisms address this issue by, in effect, learning a good approximation of D from players' reports (at the cost of additional mechanism complexity and/or additional assumptions on the number of players).

A second and quite serious issue is the presence of additional (non-truthful and “bad”) equilibria. This issue is easiest to see with the output agreement mechanism—if everybody always reports “1” (for example), then everybody gets their best-case payoff of 1 with certainty. (This would correspond to all graders agreeing to give everybody the maximum-possible score.) There is a clear incentive for players to coordinate on this equilibrium of the output agreement mechanism: not only does everybody receive their maximum-possible payoff, but no effort is required to formulate the report (unlike, say, accurate grading). But this equilibrium is a disaster for whoever is running the mechanism — the reports are independent of players' signals, and provide no information whatsoever. Such “uninformative” equilibria also plague the Peer Prediction mechanism.

This problem is not merely theoretical. In experiments, there have been cases where participants really do seem to coordinate on high-payoff but uninformative equilibria [3, 5].

To add insult to injury, paying participants a fixed reward (independent of their report) can lead to more truthful reports than a Peer Prediction-style mechanism [5]! Using a non-trivial mechanism to elicit feedback seems to nudge participants into thinking strategically, which in turn can lead to less informative behavior.

Lots of ongoing research is attempting to mitigate these problems. For example, one approach to escape the obstacles above is to do a limited amount of costly verification to produce a ground truth for a small number of outcomes (e.g., a few assignments graded by the instructor rather than the students) [3, 4]. Another approach is to tweak the Peer Prediction mechanism so that all of the uninformative (i.e., signal-independent) equilibria provide lower payoffs to all players than the truthful equilibrium, which should encourage players to coordinate on the truthful equilibrium [7].

References

- [1] J. E. Bickel. Some comparisons among quadratic, spherical, and logarithmic scoring rules. *Decision Analysis*, 4(2):49–65, 2007.
- [2] G. W. Brier. Verification of forecasts expressed in terms of probability. *Monthly Weather Review*, 78(1):1–3, 1950.
- [3] L. de Alfaro, M. Shavlovsky, and V. Polychronopoulos. Incentives for truthful peer grading. Working paper, 2016.
- [4] A. Gao, J. R. Wright, and K. Leyton-Brown. Incentivizing evaluation via limited access to ground truth: Peer-prediction makes things worse. Working paper, 2016.
- [5] X. A. Gao, J. Zhang, and Y. Chen. Trick or treat: putting peer prediction to the test. In *Proceedings of the 15th ACM Conference on Economics and Computation (EC)*, pages 507–524, 2014.
- [6] I. J. Good. Rational decisions. *Journal of the Royal Statistical Society, Series B*, 14(1):107–114, 1952.
- [7] Y. Kong, G. Schoenebeck, and K. Ligett. Putting peer prediction under the micro(economic)scope and making truth-telling focal. Working paper, 2016.
- [8] N. Miller, P. Resnick, and R. Zeckhauser. Eliciting informative feedback: The peer prediction method. *Management Science*, 51(9):1359–1373, 2005.
- [9] C. Piech, J. Huang, Z. Chen, C. Do, A. Ng, and D. Koller. Tuned models of peer assessment in MOOCs. In *Proceedings of the 6th International Conference on Educational Data Mining*, 2013.