

Strategyproof Classification Under Constant Hypotheses: A Tale of Two Functions

Reshef Meir, Ariel D. Procaccia, and Jeffrey S. Rosenschein

School of Engineering and Computer Science
The Hebrew University of Jerusalem
{reshef24,arielpro,jeff}@cs.huji.ac.il

Abstract

We consider the following setting: a decision maker must make a decision based on reported data points with binary labels. Subsets of data points are controlled by different selfish agents, which might misreport the labels in order to sway the decision in their favor. We design mechanisms (both deterministic and randomized) that reach an approximately optimal decision and are strategyproof, i.e., agents are best off when they tell the truth. We then recast our results into a classical machine learning classification framework, where the decision maker must make a decision (choose between the constant positive hypothesis and the constant negative hypothesis) based only on a sampled subset of the agents' points.

Introduction

In the design and analysis of multiagent systems, one often cannot assume that the agents are cooperative. Rather, the agents might be self-interested, seeking to maximize their own utility, possibly at the expense of the social good. With the growing awareness of this situation, game-theoretic notions and tools are increasingly brought into play.

One such setting, which we shall consider here, arises when a decision has to be made based on data points that are controlled by multiple (possibly) selfish agents, and the decision affects all the agents. The decision maker would like to make a decision which is consistent, as much as possible, with all the available data. However, the agents might misreport their data in an attempt to influence the final decision in their favor.

Motivating Examples

Consider, for instance, a spatial sensor array (represented as points in \mathbb{R}^3), and assume that each agent controls a subset of the sensors (such as the ones positioned in its own territory). A sensor's output is only available, as private information, to the controlling agent. One such scenario might be battlefield acoustic sensors (Lesser & Erman 1980): every agent controls a sector, and is charged with a specific mission in this sector. An agent might be interested in retreating (and thus failing to complete its mission) only if massive enemy movement is detected *in its own sector*. However, a global decision, to proceed or retreat, has to be made. An

agent may misreport its sensor readings in order to bring about a favorable allied decision.

A second example with an economic aspect might be a common central bank, such as the European Central Bank (ECB). The governing council makes decisions that are based on reports from the various national central banks (so one can think of the national central bankers as the agents). The national central bankers, in turn, collect private information, by means of their own institutions, regarding various economic indicators (these are the data points). Naturally, decisions taken at the European level (about, for instance, whether or not to support certain monetary policies) affect all national central banks. This strongly incentivizes the national central bankers to misreport their national statistics in a way that guarantees a decision they find desirable (though in this particular case, fear of discovery does incentivize truthfulness).

Overview of Models and Results

We present our results at two levels of generality. The more specific level is strongly motivated in its own right, but technically is also directly applied in order to obtain more general results in a machine learning framework.

Our specific model concerns n agents, each controlling a set of data points. Each point is labeled either as positive or negative; a positive label should be construed as implying that this data point supports some decision or proposition. Now, all agents report the labels of their points to some central authority, which in turn outputs a positive or negative decision. An agent's *risk* is (proportional to) the number of its points that the final decision mislabels, e.g., the number of negative points it controls in case of a positive decision. The decision maker is seeking to minimize the global risk, i.e., the total number of mislabeled points.

As noted above, an agent might find it advantageous to misreport the labels of its points. We are interested in designing decision-making mechanisms that are *strategyproof*: agents cannot benefit by lying. In return we only ask for approximate optimality. We put forward a simple deterministic decision-making mechanism which is group strategyproof (i.e., even coalitions of agents do not gain from lying) and gives a 3-approximation of the optimal global risk; in other words, the number of mislabeled points is at most 3 times the minimal number. Moreover, we show that no determinis-

tic strategyproof mechanism can do better. Interestingly, we circumvent this result by designing a strategyproof *randomized* mechanism which gives a 2-approximation, and further demonstrate that this is as far as randomization can take us.

The second part of the paper recasts the first into a more general model, which deals with the classical machine learning classification framework. It is often the case that the decision maker cannot query agents regarding all their points, due to, for example, communication or privacy constraints (think of the European bank example given above; in this example, both abovementioned constraints apply, as the number of economic indicators is enormous, and economic institutions are well-aware of privacy considerations). To complicate matters, in the general model each agent holds a different distribution over the input space, which reflects the relative importance it gives to different data points. So, we assume that a mechanism receives labels of points from agents, where each agent’s points are sampled from its individual distribution. The mechanism then outputs a decision: one of *two functions*, the *constant positive hypothesis* or the *constant negative hypothesis*. The goal is to guarantee that the concept returned by the algorithm gives a good approximation of the optimal risk, in expectation. Crucially, we demonstrate that the results of the previous, more specific, model can be leveraged to achieve this goal.

Related Work

Our work is closely related to the work of Dekel, Fischer and Procaccia (2008). They also investigated game-theoretic aspects of machine learning, albeit in a *regression learning* setting. Specifically, in their setting the label of each data point is a real number, and the risk of some hypothesis is the total *distance* to the correct labels. Dekel et al. put forward approximately optimal and strategyproof algorithms for some limited hypothesis classes. Our work should be seen as an extension of theirs to the world of classification, specifically under the very interesting (as will become apparent later) hypothesis class that contains only the two constant (positive and negative) functions.

Several existing works study issues on the border of machine learning and game theory (Balcan *et al.* 2005; Procaccia *et al.* 2007). For example, some papers on multiagent learning (see, e.g., Littman (1994), or Hu and Wellman (2004)) attempt to learn a Nash equilibrium in Markov games (which model multiagent interactions), usually via reinforcement learning. These works do not consider incentives in the learning process itself, but rather using learning to deal with strategic situations.

Another line of research attempts to learn in the face of noise (Littlestone 1991; Kearns & Li 1993; Goldman & Sloan 1995; Bshouty, Eiron, & Kushilevitz 2002). Perhaps closer to our work is the paper of Dalvi *et al.* (2004), who model classification as a game between a classifier and an adversary. Dalvi *et al.* examine the optimal strategies of the classifier and adversary, given their strategic considerations. In contrast (but similarly to Dekel *et al.* (2008)), our research concentrates on designing strategyproof algorithms, i.e., algorithms that preclude strategic behavior *in the first place*, rather than algorithms that work well *in spite of* strategic

behavior.

A Simple Setting

In this section we present a specific model, as described above: each agent controls a subset of data points; the decision maker has full information about the “identity” of the points controlled by the various agents, but does not know their labels. Rather, the labels are reported by the agents. This simple setting is strongly motivated in its own right (see the examples given above), but will also be leveraged later to obtain results in a learning-theoretic setting. In order to easily recast the results later, we introduce some learning theoretic notions already in this section.

Formally, let $I = \{1, \dots, n\}$ be the set of agents. Let \mathcal{X} be the input space, and $\{+, -\}$ be the set of labels to which points in \mathcal{X} can be mapped.

For each agent $i \in I$, let $X_i = \{x_{i,1}, \dots, x_{i,m_i}\} \subseteq \mathcal{X}^{m_i}$ be the set of points that agent i controls, and let $Y_i = \{y_{i,1}, \dots, y_{i,m_i}\} \subseteq \{+, -\}^{m_i}$ be the set of labels that are associated with these points. We refer to the pair $s_{i,j} = \langle x_{i,j}, y_{i,j} \rangle$ as an *example*. A positive label means, intuitively, that the example supports a decision, while a negative one means the example opposes it. We denote the subset of the dataset controlled by agent i with $S_i = \{s_{i,j}\}_{j=1}^{m_i}$ and the entire dataset, i.e., the multiset of all examples, by $S = \uplus_{i \in I} S_i$.

Let C be a class of functions from \mathcal{X} to $\{+, -\}$, i.e., each $c \in C$ is a classifier that maps all possible points to labels. In learning theory, C is referred to as the *concept class* of the problem. In this paper we will consider the special case where C contains only the two constant functions $\{c_+, c_-\}$ where $\forall x \in \mathcal{X}, c_+(x) = +; c_-(x) = -$, i.e., the classification mechanism may decide to classify *all* examples as positive, or all of them as negative. This should be interpreted as taking either a positive or a negative decision.

We evaluate each such classifier simply according to the number of errors it makes on the set of examples. Formally, we define the *subjective risk* associated by agent i with the classifier c as

$$R_i(c, S_i) = \frac{1}{m_i} \sum_{j=1}^{m_i} \ell(c(x_{i,j}), y_{i,j}),$$

where ℓ is the natural 0–1 loss function: $\ell(y, y')$ is 1 if $y \neq y'$ and 0 if $y = y'$. We define the *global risk* in a similar way to be the average risk with respect to all agents:

$$R(c, S) = \frac{\sum_{i \in I} m_i R_i(c, S_i)}{\sum_{i \in I} m_i} = \frac{1}{m} \sum_{\langle x, y \rangle \in S} \ell(c(x), y), \quad (1)$$

where $m = \sum_{i \in I} m_i$.

A *mechanism* receives as input a dataset S , and outputs one of the two concepts in C . Our goal is to design a mechanism that minimizes the global risk, i.e., a mechanism that chooses from $\{c_+, c_-\}$ the concept that makes fewer errors on S . Less formally, the decision maker would like to make either a positive or negative decision, in a way that is most consistent with the available data.

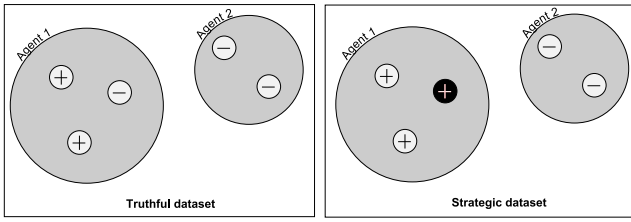


Figure 1: ERM is not strategyproof. Agent 1 changes one of its points from negative to positive, thus changing the risk minimizer from c_- to c_+ , to agent 1’s advantage. In this illustration, $\mathcal{X} = \mathbb{R}^2$.

In our model, agents report to the mechanism the labels of the points they control. If all agents report truthfully, the above problem is trivially solved by choosing c according to the majority of labels. This is a special case of the Empirical Risk Minimization (ERM) mechanism,¹ which by definition picks the concept in \mathcal{C} that minimizes the risk on the set of given examples. We denote by c^* the concept returned by ERM, formally:

$$c^* = \text{ERM}(S) = \text{argmin}_{c \in \mathcal{C}} R(c, S)$$

Unfortunately, if we choose ERM as our mechanism then agents may lie in order to decrease their subjective risk. Indeed, consider the following dataset (illustrated in Figure 1): agent 1 controls 3 examples, 2 positive and 1 negative. Agent 2 controls 2 examples, both negative. Since there is a majority of negative examples, ERM would return c_- ; agent 1 would suffer a subjective risk of $2/3$. On the other hand, if agent 1 reported his negative example to be positive as well, ERM would return c_+ , with a subjective risk of only $1/3$ for agent 1. Indeed, note that an agent’s utility is measured with respect to its real labels, rather than with respect to the reported labels.

Why is truth-telling important? Once we guarantee that agents are telling the truth, we may concentrate on minimizing the risk, knowing that this is equivalent to actually maximizing the social good (i.e., making the right decision). In other words, we would like our mechanism to be *strategyproof* (SP). By definition, a mechanism is SP (in dominant strategies) if no agent may gain (i.e., lower his subjective risk) by reporting labels that differ from his real labels.

Remark 1. If we allow payments to be transferred to and from the agents, ERM can be augmented with Vickrey-Clarke-Groves (VCG) payments to achieve strategyproofness (see, e.g., (Nisan 2007) for an overview of the VCG mechanism). However, in many multiagent systems, and in particular in internet settings, such payments are often not feasible. Therefore, we concentrate throughout the paper on achieving good mechanisms *without* payments. See Dekel et al. (2008) for a discussion of this point.

Despite the fact that ERM is not SP, the concept that minimizes the global risk is clearly optimal. Thus we would like

¹In the current setting, there is no distinction between empirical risk and “real” risk. This distinction will become apparent in the next section.

to use it to evaluate other concepts and mechanisms. Formally, define the optimal risk to be

$$r^* = R(c^*, S) = \min\{R(c_+, S), R(c_-, S)\}.$$

As is common in computer science, we will be satisfied with only approximate optimality (if this guarantees strategyproofness). Indeed:

Definition 2. A mechanism M is an α -approximation mechanism if for any dataset S it holds that $R(M(S), S) \leq \alpha \cdot r^*$.

ERM, for example, is a 1-approximation mechanism, but is not SP. On the other hand, a mechanism that always returns c_- is SP but does not give any finite approximation ratio (it is sufficient to consider a dataset with one positive example).

Remark 3. Informally we state that in our current setting, we can obtain similar approximation results even under mechanisms that are not SP, assuming agents lie only when this is beneficial to them. Nevertheless, strategyproofness gives us a very clean framework to analyze mechanisms in the face of strategic behavior. When we discuss our learning theoretic framework, where obtaining strategyproofness is next to impossible, we shall apply the former, less elegant, type of analysis.

Deterministic Mechanisms

We start with some observations. Note that the identity of each sampled point is not important, only the *number* of positive and negative points each agent controls. Thus we denote by $P_i = |\{(x, y) \in S_i : y = +\}|$, $N_i = m_i - P_i = |\{(x, y) \in S_i : y = -\}|$. For convenience we also let $P = \sum_{i \in I} P_i$, $N = \sum_{i \in I} N_i$. We emphasize that $\{P_i, N_i\}_{i \in I}$ contain all the information relevant for our problem and can thus replace S .

Now, denote by c_i the ERM on S_i , i.e., $c_i = c_+$ if $P_i \geq N_i$ and c_- otherwise. Clearly c_i is the best classifier agent i can hope for. Consider the following mechanism

Mechanism 1

1. Based on the labels of each agent P_i, N_i , calculate c_i . Define each agent as a *negative agent* if $c_i = c_-$, and as a *positive agent* if $c_i = c_+$.
2. Denote by $P' = \sum_{i: c_i = c_+} m_i$ the number of examples that belong to positive agents, and similarly $N' = \sum_{i: c_i = c_-} m_i = m - P'$.
3. If $P' \geq N'$ return c_+ , otherwise return c_- .

Remark 4. Mechanism 1 can be thought of as a specialized, imported version of the Project-and-Fit mechanism of Dekel et al. (Dekel, Fischer, & Procaccia 2008). However, the results regarding Mechanism 1’s guarantees do not follow from their results, since the setting is different (regression vs. classification).

We will show that this mechanism has the excellent game-theoretic property of being *group strategyproof*: no coalition of players can gain by lying. In other words, if some agent in the coalition strictly gains from the joint lie, some other agent in the coalition must strictly lose.

Theorem 5. *Mechanism 1 is a 3-approximation group strategyproof mechanism.*

Proof. We first show group strategyproofness. Let $B \subseteq I$. We can assume without loss of generality that either all agents in B are positive or all of them are negative, since a positive (resp., negative) agent cannot gain from lying if the mechanism returns c_+ (resp., c_-). Again w.l.o.g., the agents are all positive. Therefore, if some agent is to benefit from lying, the mechanism has to return c_- on the truthful dataset. However, since the mechanism considers all agents in B to be positive agents when the truthful dataset is given, an agent in B can only hope to influence the outcome by reporting a majority of negative examples. However, this only increases N' , reinforcing the mechanism's decision to return c_- .

It remains to demonstrate that the approximation ratio is as claimed. We assume without loss of generality that the mechanism returned c_+ , i.e., $P' \geq N'$. We first prove that if the mechanism returned the positive concept, at least $1/4$ of the examples are indeed positive.

Lemma 6. $P \geq \frac{1}{4}m$.

Proof. Clearly $P' \geq \frac{m}{2} \geq N'$ otherwise we would get $c = c_-$. Now, if an agent is *positive* ($c_i = c_+$), at least half of its examples are also positive. Thus

$$P = \sum_{i \in I} P_i \geq \sum_{i: c_i = c_+} P_i \geq \sum_{i: c_i = c_+} \frac{m_i}{2} = \frac{P'}{2},$$

and so:

$$P \geq \frac{P'}{2} \geq \frac{m}{4}$$

□

Now, we know that $P + N = m$, so:

$$N = m - P \leq m - \left(\frac{m}{4}\right) = \frac{3m}{4} \leq 3P$$

Clearly if the mechanism decided “correctly”, i.e., $P \geq m/2$, then

$$R(c, S) = R(c_+, S) = \frac{N}{m} = r^*.$$

Otherwise, if $P < m/2$, then

$$R(c, S) = R(c_+, S) = \frac{N}{m} \leq 3 \frac{P}{m} = 3R(c_-, S) = 3r^*.$$

In any case we have that $R(c, S) \leq 3r^*$, proving that Mechanism 1 is indeed a 3-approximation mechanism. □

As 3-approximation is achieved by such a trivial mechanism, we would naturally like to know whether it is possible to get a better approximation ratio, without waiving the SP property. We show that this is *not* the case by proving a matching lower bound on the best possible approximation ratio achievable by an SP mechanism. Note that the lower bound only requires strategyproofness, not group strategyproofness.

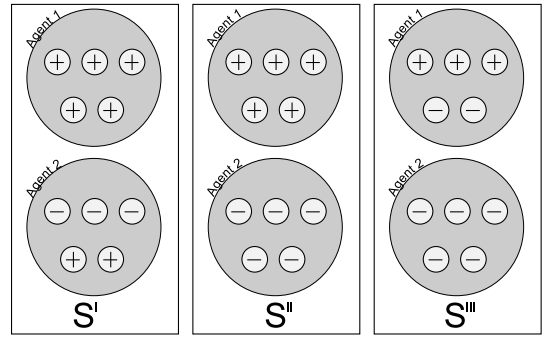


Figure 2: The examples of each agent in the three datasets are shown (for $k = 2$). Agent 1 can make dataset II look like dataset III and vice versa by reporting false labels. The same goes for agent 2 regarding datasets I and II.

Theorem 7. *Let $\epsilon > 0$. There is no $(3 - \epsilon)$ -approximation strategyproof mechanism.*

Proof. To prove the bound, we present 3 different datasets. We show that any SP mechanism must return the same result on all of them, while neither concept in C yields an approximation ratio of $(3 - \epsilon)$ in all three.

Let $\epsilon > 0$. We will use $I = \{1, 2\}$, and an integer $k = k(\epsilon)$ to be defined later. Note that in all 3 datasets $m_1 = m_2 = 2k + 1$. We define the three datasets as follows (see Figure 2 for an illustration):

- S^I : $P_1 = 2k + 1, N_1 = 0$; $P_2 = k, N_2 = k + 1$
- S^{II} : $P_1 = 2k + 1, N_1 = 0$; $P_2 = 0, N_2 = 2k + 1$
- S^{III} : $P_1 = k + 1, N_1 = k$; $P_2 = 0, N_2 = 2k + 1$

Let M be some strategyproof mechanism. Then it must hold that $M(S^I) = M(S^{III})$. Indeed, otherwise assume first that $M(S^I) = c_+$ and $M(S^{III}) = c_-$. Notice that the only difference between the two settings is agent 2's labels. If agent 2's truthful labels are as in S^I , his subjective ERM is c_- . Therefore, he can report his labels to be as in S^{III} (i.e., all negative) and obtain c_- . Now, if $M(S^I) = c_-$ and $M(S^{III}) = c_+$, agent 2 can gain by deviating from S^{III} to S^I . A symmetric argument, with respect to agent 1 (that in all settings prefers c_+) shows that $M(S^{II}) = M(S^{III})$.

So, without loss of generality assume that $c = M(S^I) = M(S^{II}) = M(S^{III}) = c_+$ (otherwise, symmetric arguments yield the same result). Therefore:

$$R(c, S^{III}) = R(c_+, S^{III}) = \frac{N_1 + N_2}{m} = \frac{3k + 1}{4k + 2} \quad (2)$$

On the other hand, the negative concept is much better:

$$r^* = R(c_-, S^{III}) = \frac{k + 1}{4k + 2}$$

By combining the last two equations:

$$\frac{R(c, S^{III})}{r^*} = \frac{\frac{3k+1}{4k+2}}{\frac{k+1}{4k+2}} = \frac{3k+1}{k+1}$$

Let us set $k > \frac{3}{\epsilon}$; then the last expression is strictly greater than $3 - \epsilon$, and thus $R(c, S^{III}) > (3 - \epsilon)r^*$. We conclude that any SP mechanism cannot have an approximation ratio of $3 - \epsilon$. \square

Randomized mechanisms

What if we let our mechanism flip coins? Can we find an SP randomized mechanism that beats (in expectation) the 3-approximation deterministic lower bound? To answer the question we first need to formally define the risk of such a mechanism, since it may return different concepts on the same dataset. We do this by simply by taking the *expected* risk over all possible outcomes.

Definition 8. Let M be a randomized mechanism, which returns each concept $c \in C$ with probability $p_M(c|S)$.

$$R(M(S), S) = \mathbb{E}[R(c, S)] = \sum_{c \in C} p_M(c|S) \cdot R(c, S)$$

For our simple concept class $C = \{c_+, c_-\}$, a randomized mechanism is defined only by the probability of returning a positive or negative concept, given S . Accordingly, the risk is

$$R(M(S), S) = p_M(c_+|S)R(c_+, S) + p_M(c_-|S)R(c_-, S)$$

We start our investigation of SP randomized mechanisms by establishing a lower bound of 2 on their approximation ratio.

Theorem 9. *Let $\epsilon > 0$. There is no $(2 - \epsilon)$ -approximation strategyproof randomized mechanism.*

The proof of the theorem is given in the appendix that was submitted with this paper as supplementary material.

We presently put forward a randomized SP 2-approximation mechanism, thereby matching the lower bound with an upper bound. We will calculate P' and N' as in our deterministic Mechanism 1. The natural thing to do would be simply to select c_+ with probability P'/m and c_- with probability N'/m . Unfortunately, this simple randomization (which is clearly SP) cannot even beat the deterministic bound of $3 - \epsilon$.²

Crucially, a more sophisticated (and less intuitive) randomization can do the trick.

Mechanism 2

1. Compute P' and N' as in Mechanism 1.
2. If $P' \geq N'$, set $t = \frac{N'}{m}$; return c_+ with probability $\frac{2-3t}{2-2t}$, and c_- with probability $\frac{t}{2-2t}$.
3. Else if $N' > P'$, set $t = \frac{P'}{m}$; return c_- with probability $\frac{2-3t}{2-2t}$, and c_+ with probability $\frac{t}{2-2t}$.

Theorem 10. *Mechanism 2 is a group strategyproof 2-approximation randomized mechanism.*

The proof of the theorem is given in the appendix that was submitted with this paper as supplementary material.

²We will not prove this formally, but shortly consider $P_1 = k + 1$, $N_1 = k$; $N_2 = m_2 = k(2k + 1)$ as k increases.

A Learning Theoretic Setting

In this section we extend our simple setting to a more general machine learning framework. Our previous results will be leveraged to obtain powerful learning theoretic results.

Instead of looking at a fixed set of examples and selecting the concept that fits them best, each agent $i \in I$ now has a private function $Y_i : \mathcal{X} \rightarrow \{+, -\}$, which assigns a label to every point in the input space. In addition, every agent holds a (known) distribution ρ_i over the input space, which reflects the relative importance it attributes to each point. The new definition of the subjective risk naturally extends the previous setting by expressing the errors a concept makes when compared to Y_i , given the distribution ρ_i :

$$R_i(c) = \mathbb{E}_{x \sim \rho_i}[\ell(c(x), Y_i(x))]$$

The global risk is calculated similarly to the way it was before. For ease of exposition, we will assume in this section that all agents have equal weight.³

$$R(c) = \sum_{i \in I} \frac{1}{n} \cdot R_i(c)$$

Since we cannot directly evaluate the risk in this learning theoretic framework, we may only sample points from the agents' distributions and ask the agents to label them. We then try to minimize the *real* global risk, using the *empirical risk* as a proxy. The empirical risk is the risk on the sampled dataset, as defined in the previous section.

Mechanism 3

1. for each agent $i \in I$, sample m points i.i.d. from ρ_i . Denote i 's set of points as $X_i = \{x_{i1}, \dots, x_{im}\}$.
2. For every $i \in I$, $j = 1, \dots, m$, ask agent i to label x_{ij} . Denote $\bar{S}_i = \{(x_{i,j}, y_{i,j})\}_{j=1}^m$.
3. Use Mechanism 2 on $\bar{S} = \{\bar{S}_1, \dots, \bar{S}_n\}$, and return the result.

We presently establish a theorem that explicitly states the number of examples we need to sample in order to properly estimate the real risk. We will get that in expectation (taken over the randomness of the sampling procedure and Mechanism 2's randomization) Mechanism 3 yields close to a 2-approximation with relatively few examples, even in the face of strategic behavior. The subtle point here is that Mechanism 3 is not strategyproof. Indeed, even if an agent gives greater weight to negative points (according to Y_i and ρ_i), it might be the case that (by miserable chance) the agent's sampled dataset only contains positive points.

However, since Mechanism 2 is SP in the previous section's setting, if an agent's sampled dataset faithfully represents its true distribution, and the agent is strongly inclined towards c_+ or c_- , the agent still cannot benefit by lying. If an agent is almost indifferent between c_+ and c_- , it might wish to lie—but crucially, such an agent contributes little to the global risk.

³The results can be generalized to varying weights by sampling for each agent a number of points proportional to its weight, yet still large enough.

Our game theoretic assumption in the theorem is that agents that cannot gain by lying will tell the truth (so under this assumption, some agents may tell the truth even if they gain by lying). This is a weaker assumption than the common assumption that all agents are utility maximizing (i.e., simply wish to minimize their subjective risk). It is useful to employ the weaker version, as in many settings it might be the case that some of the agents are centrally designed, and so are bound to tell the truth regardless (even if they can gain by lying).

Remark 11. Consider the following simple mechanism: sample one point per agent, and let the agent label this single point. If the agent labels the point positively, the agent is positive; otherwise it is negative. Now apply Mechanism 2. Under the latter (strong) assumption this mechanism provides good guarantees, but under the former (weak) assumption it provides bad guarantees (since truthful agents might be assigned datasets that do not reflect their risk)—unlike Mechanism 3, as will become apparent momentarily.

One can also consider a mechanism that just asks each agent to report whether it prefers c_+ or c_- . Such a mechanism, though, is not consistent with our learning theoretic framework, and so is outside the scope of this paper.

Theorem 12. *Given sampled datasets, assume that agents are truthful if they cannot gain by lying. Let $R(M_3)$ denote the expected risk of Mechanism 3, where the expectation is taken over the randomness of the sampling and Mechanism 2. For any $\epsilon' > 0$, there is m (polynomial in n and $\frac{1}{\epsilon'}$) such that by sampling m points for each agent, it holds that*

$$R(M_3) \leq 2r^* + \epsilon'.$$

The proof of the theorem, which is technically the most involved, is given in the appendix that was submitted with this paper as supplementary material.

Remark 13. In our current learning theoretic setting there are no reasonable SP mechanisms. Indeed, even dictatorship, i.e., choosing some fixed agent's best classifier given its reported examples, is not SP, as one can sample a majority of positive examples when the agent in fact prefers c_- . In their Theorem 5.1, Dekel et al. (2008) do not obtain SP in the (regression) learning theoretic setting, but rather ϵ -SP: agents cannot gain more than ϵ by lying—with high probability, given enough examples. We circumvent the strategyproofness issue with a more complicated assumption, and thereby obtain a far stronger result (which is not true in the Dekel et al. setting). On the other hand, our result only holds for the very small hypothesis class $\{c_+, c_-\}$, while theirs is more general.

Conclusions

We explored the problem of making a decision based on labeled data, under the assumption that the labels are not directly accessible. Rather, they are reported by agents that may lie in order to bias the final decision in their favor.

Using the classic definition of optimal risk as the minimal number of mislabeled data points, we presented a very simple deterministic strategyproof mechanism whose risk is at

most three times optimal. Moreover, we demonstrated that no deterministic mechanism can do better while maintaining the strategyproofness property. We further showed that the deterministic 3-approximation bound can be improved to 2-approximation using the notion of expected risk and a nonintuitive randomized mechanism. Finally, in the last section we demonstrated how to reformulate this mechanism in a learning theoretic setting, where the mechanism essentially learns a constant concept based on sampled data that is controlled by selfish agents.

Our mechanisms can serve human and automated decision makers that wish to maximize social welfare in the face of data that is biased by conflicting interests. Crucially, our results in the learning theoretic setting constitute first steps in designing classifiers that can function well in non-cooperative environments; in the future we intend to extend the results to richer concept classes.

References

- Balcan, M.-F.; Blum, A.; Hartline, J. D.; and Mansour, Y. 2005. Mechanism design via machine learning. In *The 46th Annual IEEE Symposium on Foundations of Computer Science (FOCS 2005)*, 605–614.
- Bshouty, N. H.; Eiron, N.; and Kushilevitz, E. 2002. PAC learning with nasty noise. *Theoretical Computer Science* 288(2):255–275.
- Dalvi, N.; Domingos, P.; Mausam; Sanghai, S.; and Verma, D. 2004. Adversarial classification. In *Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD 2004)*, 99–108.
- Dekel, O.; Fischer, F.; and Procaccia, A. D. 2008. Incentive compatible regression learning. In *The ACM-SIAM Symposium on Discrete Algorithms (SODA 2008)*, 277–286.
- Goldman, S. A., and Sloan, R. H. 1995. Can PAC learning algorithms tolerate random attribute noise? *Algorithmica* 14(1):70–84.
- Hu, J., and Wellman, M. 2004. Nash Q-learning for general-sum stochastic games. *Journal of Machine Learning Research* 4:1039–1069.
- Kearns, M., and Li, M. 1993. Learning in the presence of malicious errors. *SIAM J. on Computing* 22(4):807–837.
- Lesser, V. R., and Erman, L. D. 1980. Distributed interpretation: A model and experiment. *IEEE Transactions on Computers* 29(12):1144–1163.
- Littlestone, N. 1991. Redundant noisy attributes, attribute errors, and linear-threshold learning using Winnow. In *COLT*, 147–156.
- Littman, M. L. 1994. Markov games as a framework for multi-agent reinforcement learning. In *ICML*, 157–163.
- Nisan, N. 2007. Introduction to mechanism design (for computer scientists). In Nisan, N.; Roughgarden, T.; Tardos, E.; and Vazirani, V., eds., *Algorithmic Game Theory*. Cambridge University Press. chapter 9.
- Procaccia, A. D.; Zohar, A.; Peleg, Y.; and Rosenschein, J. S. 2007. Learning voting trees. In *The National Conference on Artificial Intelligence (AAAI 2007)*, 110–115.

Proof of Theorem 9

We will use the same datasets used in the proof of Theorem 7, and illustrated in Figure 2. Let M be a SP randomized mechanism, and denote by $p_M(c|S)$ its probability of outputting c given S .

We first show that the mechanism chooses the positive hypothesis with the same probability in all three datasets.

Lemma 14. $p_M(c_+|S^I) = p_M(c_+|S^{II}) = p_M(c_+|S^{III})$.

Proof (sketch). If $p_M(c_+|S^I) \neq p_M(c_+|S^{II})$ then agent 2 will report its labels in a way that guarantees a higher probability of c_- . Similarly, $p_M(c_+|S^{II}) \neq p_M(c_+|S^{III})$ implies that agent 1 can increase the probability of c_+ by lying. \square

Denote

$$p_+ = p_M(c_+|S^I) = p_M(c_+|S^{II}) = p_M(c_+|S^{III}),$$

and

$$p_- = p_M(c_-|S^I) = p_M(c_-|S^{II}) = p_M(c_-|S^{III}).$$

Without loss of generality $p_+ \geq \frac{1}{2} \geq p_-$. Then:

$$\begin{aligned} R(M(S^{III}), S^{III}) &= p_+ R(c_+, S^{III}) + p_- R(c_-, S^{III}) \\ &= p_+ \cdot \frac{3k+1}{4k+2} + p_- \cdot \frac{k+1}{4k+2} \\ &\geq \frac{1}{2} \cdot \frac{3k+1}{4k+2} + \frac{1}{2} \cdot \frac{k+1}{4k+2} = \frac{1}{2}, \end{aligned}$$

whereas

$$r^* = R(c_-, S^{III}) = \frac{k+1}{4k+2}$$

For $k > \frac{1}{\epsilon}$ it holds that

$$\frac{R(M(S^{III}), S^{III})}{r^*} = \frac{4k+2}{2(k+1)} = 2 - \frac{1}{k+1} > 2 - \epsilon$$

As before, if $p_- > p_+$, a symmetric argument shows that $R(M(S^I), S^I) > (2 - \epsilon)r^*$. Therefore no SP mechanism can achieve a $(2 - \epsilon)$ -approximation, even through randomization. \square

Proof of Theorem 10

Similarly to Mechanism 1, Mechanism 2 is clearly group SP, since declaring a false label may only increase the probability of obtaining a classifier that labels correctly less than half of the agent's examples, thus increasing the subjective expected risk.

Assume without loss of generality that $P' \geq N'$, so $t = \frac{N'}{m}$, and c_+ is returned with probability $\frac{2-3t}{2-2t}$, c_- with probability $\frac{t}{2-2t}$. Recall that N, P denote the number of negative and positive examples, respectively. First notice that

$$N' \leq 2N. \quad (3)$$

This is trivially true since every negative agent has a majority of negative examples.

Case 1: $P \geq N$. From (3), the real risk of the best classifier satisfies:

$$r^* = R(c_+, S) = \frac{N}{m} \geq \frac{N'}{2m} = \frac{t}{2}, \quad (4)$$

whereas Mechanism 2 satisfies:

$$\begin{aligned} R(M(S), S) &= \frac{2-3t}{2-2t} r^* + \frac{t}{2-2t} (1-r^*) \\ &= \frac{2-3t}{2-2t} r^* + \frac{tr^*}{2-2t} \left(\frac{1}{r^*} - 1\right) \\ &= \frac{r^*}{2-2t} (2-3t + t(\frac{1}{r^*} - 1)) \\ &= \frac{r^*}{2-2t} (2 + \frac{t}{r^*} - 4t) \\ &\leq \frac{r^*}{2-2t} (2 + 2 - 4t) \\ &= 2r^*, \end{aligned}$$

where the inequality follows from (4).

Case 2: $N > P$. In this case, the optimal risk is $r^* = R(c_-, S) = \frac{P}{m}$.

Lemma 15. $\frac{1}{r^*} - 1 \leq \frac{1+t}{1-t}$

Proof. The largest possible number of negative examples is achieved when all the negative agents control only negative examples, and all the positive agents control only a slight majority of positive labels. Formally, we have that $N \leq N' + \frac{P'}{2}$, and thus:

$$\frac{N}{m} \leq \frac{N'}{m} + \frac{P'}{2m} = \frac{N'}{m} + \frac{m-N'}{2m} = \frac{N'}{2m} + \frac{1}{2}.$$

It follows that $1 - r^* \leq \frac{t}{2} + \frac{1}{2}$; therefore $r^* \geq \frac{1-t}{2}$. We now conclude that

$$\frac{1}{r^*} - 1 \leq \frac{2}{1-t} - 1 = \frac{1+t}{1-t}$$

\square

Now, we have:

$$\begin{aligned} R(M(S), S) &= \frac{t}{2-2t} r^* + \left(\frac{2-3t}{2-2t}\right) (1-r^*) \\ &= \left(\frac{t}{2-2t} + \left(1 - \frac{t}{2-2t}\right) \left(\frac{1}{r^*} - 1\right)\right) r^* \\ &\leq \left(\frac{t}{2-2t} + \left(1 - \frac{t}{2-2t}\right) \frac{1+t}{1-t}\right) r^* \\ &= \frac{1-2t^2}{(1-t)^2} r^* = f(t)r^*, \end{aligned}$$

where the inequality is due to Lemma 15.

It is now sufficient to show that $f(t) \leq 2$. By taking the derivative of $f(t)$ we find that

$$f'(t) = \frac{2-4t}{(1-t)^3}.$$

Note that both numerator and denominator are nonnegative in the range $t \in [0, 1/2]$, thus $f'(t)$ is nonnegative and $f(t)$ is monotonically nondecreasing:

$$\forall t \in [0, 1/2], f(t) \leq f\left(\frac{1}{2}\right) = 2$$

As in the deterministic proof, we have that in any case $R(M(S), S) \leq 2r^*$, thus 2-approximation is assured. \square

Proof of Theorem 12

In this proof we will differentiate the real risk, as defined for the learning-theoretic setting, from the *empirical* risk on a given sample, as defined in the simple setting. The empirical risk will be denoted by

$$\hat{R}(c, S) = \frac{1}{m} \sum_{(x,y) \in S} \ell(c(x), y).$$

Without loss of generality we assume that $r^* = R(c_-) < R(c_+)$. Notice that if $r^* = R(c^*) = R(c_-) > \frac{1}{2} - 3\epsilon$ then any concept our mechanism returns will trivially attain a risk of at most $\frac{1}{2} + 3\epsilon \leq r^* + 6\epsilon$. Therefore, we can assume for the rest of this proof that

$$R(c_-) + 3\epsilon \leq \frac{1}{2} \leq R(c_+) - 3\epsilon. \quad (5)$$

Let us introduce some new notations and definitions. Denote the data set with the real labels by $S_i = \{(x_{i,j}, Y_i(x_{i,j}))\}_{j \leq m}$; $S = \{S_1, \dots, S_n\}$. Note that the mechanism has no direct access to S , but only to the reported labels as they appear in \tilde{S} .

Define G as the event “the empirical and real risk differ by at most ϵ for all agents”; formally:

$$\forall c \in \{c_+, c_-\}, \forall i \in I, |\hat{R}_i(c, S_i) - R_i(c)| < \epsilon. \quad (6)$$

Lemma 16. *Let $\delta > 0$. If $m > \frac{1}{2\epsilon^2} \ln(\frac{2n}{\delta})$, then with probability of at least $1 - \delta$, G occurs.*

Proof. Fix $i \in I$. Let $e(x)$ be the indicator random variable of the event $Y_i(x) = +$. We can now rewrite the empirical and real risk as the sum and the expectation of $e(x)$:

$$R_i(c_-) = \mathbb{E}_{x \sim \rho_i}[e(x)]$$

$$\hat{R}_i(c_-, S_i) = \frac{1}{m} \sum_{x \in S_i} e(x)$$

Since S_i is sampled i.i.d. from ρ_i , the empirical risk is the sum of independent Bernoulli random variables with expectation $R_i(c_-)$. We derive from the Chernoff bound that for any data set of size $|S_i| = m$:

$$\Pr[|\hat{R}_i(c_-, S_i) - R_i(c_-)| > \epsilon] < 2e^{-2\epsilon^2 m}$$

Taking $m > \frac{1}{2\epsilon^2} \ln(\frac{2n}{\delta})$, we get:

$$\begin{aligned} \Pr[\neg G] &= \Pr[\exists i \in I, |\hat{R}_i(c_-, S_i) - R_i(c_-)| > \epsilon] \\ &\leq \sum_{i \in I} \Pr[|\hat{R}_i(c_-, S_i) - R_i(c_-)| > \epsilon] \\ &\leq |I| 2e^{-2\epsilon^2 m} < n \frac{\delta}{n} = \delta, \end{aligned}$$

where the first inequality is due to the union bound.

Note that it is enough to show the above for c_- since

$$|\hat{R}_i(c_-, S_i) - R_i(c_-)| = |\hat{R}_i(c_+, S_i) - R_i(c_+)|. \quad \square$$

If G occurs, then from (6) and the triangle inequality it holds that for all $c \in \{c_+, c_-\}$ and $i \in I$,

$$|R(c) - \hat{R}(c, S)| \leq \sum_{i \in I} \frac{1}{n} |R_i(c) - \hat{R}_i(c, S)| \leq \epsilon. \quad (7)$$

Using (7) we could have bounded the risk of $M(S)$, but unfortunately this would not do as the mechanism may only access \tilde{S} and not S . In order to bound $R(M(\tilde{S}))$, we need to know, or estimate, how the agents label their examples. To handle this problem, we will first analyze which agents may gain by lying, and then define a new data set \tilde{S} with the following two properties: no agent has motivation to lie (thus we can assess the result of running M on \tilde{S}), and \tilde{S}, S are very similar.

We now divide I into two types of agents:

$$I' = \{i \in I : |R_i(c_-) - \frac{1}{2}| < \epsilon\},$$

and $I'' = I \setminus I'$. For each agent $i \in I$, we denote by P_i, N_i the number of positive/negative examples the agent controls in S_i . Note that $P_i = m\hat{R}_i(c_-, S_i)$. Since $R(c_-) < R(c_+)$ we may assume without loss of generality that all agents $i \in I'$ prefer c_+ (otherwise lying only lowers the expected risk of our mechanism). Agents in I'' , on the other hand, cannot benefit by lying, since S_i must reflect i 's truthful preferences, and Mechanism 2 (which is used by Mechanism 3 in step 3) is SP.

For each agent i define a new set of examples \tilde{S}_i as follows:

- If $i \in I''$, $\tilde{S}_i = S_i$.
- If $i \in I'$, define $\tilde{P}_i = P_i + \lceil \epsilon m \rceil$ and let \tilde{S}_i contain \tilde{P}_i positive examples and $m - \tilde{P}_i$ negative ones.

Lemma 17. *If G occurs, then for all agents in I*

$$\tilde{N}_i \leq \tilde{P}_i \iff R_i(c_-) \geq R_i(c_+)$$

Proof. If $i \in I''$ then w.l.o.g. $R_i(c_-) \leq R_i(c_+) - 2\epsilon$, thus from (6)

$$\begin{aligned} \tilde{P}_i &= P_i = m\hat{R}_i(c_-, S_i) \leq m(R_i(c_-) + \epsilon) \\ &\leq m(R_i(c_+) - \epsilon) \leq m\hat{R}_i(c_+, S_i) = N_i = \tilde{N}_i \end{aligned}$$

If $i \in I'$ then according to our assumption

$$R_i(c_+) \leq R_i(c_-) \leq R_i(c_+) + 2\epsilon.$$

Moreover, by the definition of \tilde{P}_i ,

$$\tilde{P}_i \geq P_i + m\epsilon; \tilde{N}_i \leq N_i - m\epsilon.$$

Thus

$$\begin{aligned} \tilde{P}_i &\geq P_i + m\epsilon = m\hat{R}_i(c_-, S_i) + m\epsilon \geq mR_i(c_-) \\ &\geq mR_i(c_+) \geq m(\hat{R}_i(c_+, S_i) - \epsilon) \geq N_i - m\epsilon \geq \tilde{N}_i \end{aligned} \quad \square$$

Lemma 17 implies that, if G occurs, agents cannot do better than report \tilde{S} under Mechanism 3, since \tilde{S}_i reflects the real preferences of agent i . Now, if agent i reports truthfully, then $\bar{P}_i = P_i$. If i decides to lie, it may report more positive labels, but cannot gain from reporting more than \bar{P}_i such labels, and, crucially, the Mechanism's outcome will not change in this case. The immediate result is that we can assume:

$$P \leq \bar{P} = \sum_{i \in I} \frac{1}{n} \bar{P}_i \leq \sum_{i \in I} \frac{1}{n} \tilde{P}_i = \tilde{P},$$

and, since the expected risk of M only increases with the number of positive examples (the probability of Mechanism 3 choosing the positive classifier increases),

$$R(M(S)) \leq R(M(\bar{S})) \leq R(M(\tilde{S})). \quad (8)$$

We can now concentrate on bounding the empirical risk on \tilde{S} .

Lemma 18. *If G occurs,*

$$\forall c \in \{c_+, c_-\}, |R(c) - \hat{R}(c, \tilde{S})| \leq 3\epsilon. \quad (9)$$

As in Lemma 16, it will suffice to show this only for c_- .

Proof. From (6), for $m > \frac{1}{\epsilon}$,

$$\begin{aligned} \hat{R}(c_-, \tilde{S}) &= \frac{\tilde{P}_i}{m} = \frac{P_i + \lceil m\epsilon \rceil}{m} \leq \frac{P_i + m\epsilon + 1}{m} \\ &\leq \frac{P_i}{m} + 2\epsilon = \hat{R}(c_-, S) + 2\epsilon \\ &\leq R(c_-) + \epsilon + 2\epsilon = R(c_-) + 3\epsilon \end{aligned}$$

□

From (5) and (9)

$$\hat{R}(c_-, \tilde{S}) \leq R(c_-) + 3\epsilon \leq R(c_+) - 3\epsilon \leq \hat{R}(c_+, \tilde{S}) \quad (10)$$

So c_- is also empirically the best concept for \tilde{S} ; Mechanism 2 guarantees:

$$\hat{R}(M(\tilde{S}), \tilde{S}) \leq 2\hat{R}(c_-, \tilde{S}) \quad (11)$$

Furthermore, since the risk of Mechanism 3 is a convex combination of the risk of c_+ , c_- , we get from (9),

$$R(M(\tilde{S})) \leq \hat{R}(M(\tilde{S}), \tilde{S}) + 3\epsilon \quad (12)$$

Finally, by using (8), (12), (11) and (10) in this order, we get that if G occurs:

$$\begin{aligned} R(M(\bar{S})) &\leq R(M(\tilde{S})) \leq \hat{R}(M(\tilde{S}), \tilde{S}) + 3\epsilon \\ &\leq 2\hat{R}(c_-, \tilde{S}) + 3\epsilon \\ &\leq 2(R(c_-) + 3\epsilon) + 3\epsilon = 2r^* + 9\epsilon \end{aligned} \quad (13)$$

If G does not occur, the risk cannot exceed 1. Thus by applying Lemma 16 with $\delta = \epsilon = \frac{\epsilon'}{10}$ we find that for $m > 50 \frac{1}{\epsilon'^2} \ln(\frac{10n}{\epsilon'})$:

$$\begin{aligned} R(M) &\leq \Pr[G](2r^* + 9\epsilon) + \Pr[\neg G]1 \\ &\leq 2r^* + 9\epsilon + \epsilon \\ &\leq 2r^* + \epsilon' \end{aligned}$$

□