
Distinguishing Distributions

When Samples Are Strategically Transformed

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Often, a principal must make a decision based on data provided by an agent.
2 Moreover, typically, that agent has an interest in the decision that is not perfectly
3 aligned with that of the principal. Thus, the agent may have an incentive to select
4 from or modify the samples he obtains before sending them to the principal. In
5 other settings, the principal may not even be able to observe samples directly;
6 instead, she must rely on signals that the agent is able to send based on the samples
7 that he obtains, and he will choose these signals strategically.
8 In this paper, we give necessary and sufficient conditions for when the principal can
9 distinguish between agents of “good” and “bad” types, when the type affects the
10 distribution of samples that the agent has access to. We also study the computational
11 complexity of checking these conditions. Finally, we study how many samples are
12 needed.

13 1 Introduction

14 Anyone can have a bad day. Or a lucky one. Thus, in general, to determine with reasonable confidence
15 who are the highly capable agents—whether they be people, companies, or anything else—we need
16 to observe their output over an extended period of time. Moreover, capability is generally not
17 one-dimensional, and who should be considered highly capable depends on what it is that we are
18 looking for. Finally, the policy that we set to evaluate agents’ output will in general affect how they
19 strategically try to shape that output. Thus, we must choose our policy to enable the agents that are
20 highly capable (according to our definition) to distinguish themselves from others.

Example. Suppose that there are researchers of different *types*. Specifically, suppose we have the following set of types:

$$\Theta = \{\text{TML-H, TML-L, AML-H, AML-L}\}$$

where “TML” stands for “theoretical machine learning,” “AML” for “applied machine learning,” and “L” and “H” for “low quality” and “high quality,” respectively. Each researcher generates high-quality *ideas* (which we will in this paper refer to as *samples*) according to some probabilistic process. Suppose here the sample space is

$$S = \{\text{T, A, B}\}$$

where “T” stands for a purely theoretical idea without immediate applied significance, “A” for an applied idea without immediate theoretical significance, and “B” for an idea that has both theoretical and applied significance. Finally, suppose there are only 3 conferences: COLT, KDD, and NeurIPS (we will in this paper refer to papers published in these conferences as “signals”).

$$\Sigma = \{\text{COLT, KDD, NeurIPS}\}$$

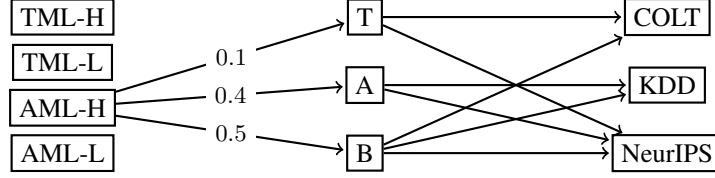


Figure 1: Illustration of the example.

21 A T or a B idea (sample) can be turned into a COLT paper (signal);¹ an A or a B idea can be turned
 22 into a KDD paper; and a T, A, or B idea can be turned into a NeurIPS paper.² Each idea, of course,
 23 can be published in only one conference.

24 Suppose a university would like to hire an AML-H researcher (but none of the other types). The
 25 faculty recruiting committee, unfortunately, is excessively lazy and only looks at the publication
 26 counts in the various venues. While the candidate researchers of course are committed to improving
 27 this terrible process once they get the job, for now their only concern is getting the job. In particular,
 28 everyone will attempt to pretend to be an AML-H researcher by sending their papers to the appropriate
 29 venues. But what exactly does this mean?

30 Suppose an AML-H researcher generates ideas at the following rates: 0.5 B, 0.4 A, 0.1 T. Moreover
 31 suppose that a TML-H researcher generates ideas at the following rates: 0.5 B, 0.1 A, 0.4 T. If the
 32 AML-H researcher sends all her papers to NeurIPS, then, even in the long run, she cannot distinguish
 33 herself from the TML-H researcher, who could do the same. On the other hand, if she sends strictly
 34 more than 0.6 of her ideas to KDD, then in the long run she will be able to distinguish herself from
 35 the TML-H researcher, because 0.4 of the latter’s ideas cannot go to KDD.

36 Now consider the AML-L researcher. First, an easy case: suppose he generates ideas at the following
 37 rates: 0.4 B, 0.3 A, 0 T. (These numbers do not sum to 1, but this is not necessary, since they are rates.
 38 Equivalently, we can suppose him to have “the empty idea” \emptyset with the remaining probability 0.3,
 39 which can be sent only to “the empty conference” where anything can be sent. This “empty signal”
 40 can also be used to model that the researchers sometimes only have ideas that they do not consider
 41 worth publishing, i.e., that they strategically select only a subset of their samples to pursue.) Clearly
 42 the AML-H researcher will in the long run distinguish herself from the AML-L researcher simply by
 43 the overall number of papers published (as long as the AML-H researcher does not unnecessarily send
 44 papers to the empty conference!). Alternatively, suppose the AML-L researcher generates ideas at
 45 the following rates: 0.4 B, 0.5 A, 0.1 T (so that the only weakness of the AML-L researcher relative
 46 to the AML-H researcher is that fewer of his ideas have both theoretical and applied significance). In
 47 this case, the AML-H researcher can, in the long run, distinguish herself from the AML-L researcher
 48 by sending strictly more than 0.5 of her ideas to COLT. Of course, this conflicts with what she needs
 49 to do distinguish herself from the TML-H researcher. Still, she can distinguish herself from both the
 50 TML-H and the AML-L researcher in the long run by, in odd-numbered years, sending strictly more
 51 than 0.6 of her ideas to KDD, and, in even-numbered years, sending strictly more than 0.5 of her
 52 ideas to COLT.

53 In the long run we are all dead. —*John Maynard Keynes*

54 In reality, the candidates will have only finite time to prove themselves. Still, the lazy committee may
 55 hope to distinguish them with high probability. How many years suffice for this (and, therefore, should
 56 be the length of a typical Ph.D. program, potentially extended with a postdoctoral appointment)?

57 While this is example is a bit tongue-in-cheek, it is not hard to see that this basic phenomenon
 58 frequently occurs in society. People select from their opportunities and craft them to fit what they

¹Of course, having the basic idea is generally only a small part of the work that needs to be done for a conference paper; but for our purposes here, we may imagine that the idea incorporates all the work that needs to be done.

²We use the names of actual conferences strictly for amusement value, and while we think our example roughly aligns with the focus of these conferences, we do not mean to imply anything about their selectivity (all these ideas are high-quality) or open-mindedness. We also do not mean to imply anything about other conferences—e.g., ICML could just as well have been used instead of NeurIPS—or (in what follows) about different types of researchers or the priorities and effort levels of actual hiring committees.

think will appeal to future employers. A start-up company may select from its opportunities and craft them to fit what they think will impress future backers. In this paper, we introduce a general model that captures all these and other cases. Within this model, we characterize conditions under which agents of certain types can distinguish themselves from others, as well as how many samples are needed for this.

1.1 Related Work

Zhang et al. [19] study a related problem in which an agent draws samples and has to submit a subset of size k of these samples to a principal, where k is exogenous. In that paper, the motivation is that the principal can inspect only so many samples. In contrast, in this paper there is no such constraint, but samples can be modified or turned into signals according to a given (arbitrary) graph. This paper also allows for uncertainty about how many samples an agent has available, via the “empty sample/signal” trick illustrated in the introductory example.

Our setting is related to mechanism design with *partial verification* [8, 18], where an agent’s type restricts which signals he can send. This can be thought of as corresponding to the special case of our model in which an agent only has a single sample which is fully determined by his type. More generally, our setting is related to the literature in economics on *signaling* (along the lines of [16]). However, our model does not involve the agents taking any costly actions. There is other work that generalizes the partial verification setting to allow costly signaling [12, 13], motivated in part by settings where agents are being classified but they can strategically change their features at some cost (as also studied in [9]).³ In contrast to this line of work, in this paper we consider settings where a single agent with a single type *repeatedly* generates samples according to a distribution (which are then strategically transformed into signals). This allows us to study the question of how many samples are needed to, with high confidence, distinguish types from each other.

Our results can be viewed as generalizations of classical results in *efficient statistics*, and in particular, results for learning and testing discrete distributions, to strategic settings. One of our main results, Theorem 6, relies on a subroutine which generalizes the folklore result that $\Theta(n/\epsilon^2)$ samples are needed to estimate a discrete distribution over support $[n]$ with error at most ϵ , with probability $2/3$. Another main result, Theorem 7, uses as a building block the identity testing algorithm by Valiant and Valiant [17]⁴, which given $O(\sqrt{n}/\epsilon^2)$ samples to a distribution y , with probability $2/3$, distinguishes between the following two cases: y is identical to a given distribution x over support $[n]$, or the distance between x and y is at least ϵ . Theorem 7 generalizes their algorithm into an environment where samples can be strategically modified according to a partial order.

2 Preliminaries

For a set S , we use $\Delta(S)$ to denote the set of probability distributions over S . Given a distribution $x \in \Delta(S)$, we use $x(i)$ to denote the probability mass on the element $i \in S$, and $x(A)$ to denote the total probability mass on the set $A \subseteq S$. We are generally interested in distinguishing one or more *good* distributions from one or more *bad* distributions (where good and bad are determined by what we are looking for). We use g to denote the good distribution, and b to denote the bad distribution. (We use $(g_i)_i$ and $(b_i)_i$ when there are multiple good/bad distributions.) The agent, depending on his type being either good or bad, draws n samples i.i.d. from either g or b . How samples can be turned into signals is represented by a bipartite graph $G = (S \cup \Sigma, E)$ between the (discrete) sample space S and the (discrete) signal space Σ . An agent must convert each sample into a signal and then submit all n signals to the principal. E specifies which signals are valid for each sample: a sample $s \in S$ can be converted into a signal $\sigma \in \Sigma$ iff $(s, \sigma) \in E$.

Note that our model generalizes each of the following models:

1. The agent can choose to omit samples. We can add an “empty signal” to Σ , where converting a sample s to the empty signal corresponds to not reporting s .

³Other work that models strategic agents manipulating the data that they submit [5, 15] concerns *aggregating* the data of multiple agents into a *single* outcome that all these agents care about; as such, this is less related to our model here, as here we are interested in determining a given single agent’s type rather than choosing a single outcome that affects multiple agents.

⁴Diakonikolas et al. [6] give an algorithm of the same sample complexity.

- 106 2. The agent may or may not receive a sample in each round. E.g., in the example where samples
 107 correspond to ideas and signals correspond to papers, in some rounds the agent may not have any
 108 (worthwhile) idea. We can add an “empty sample” in S which can only be converted to the empty
 109 signal.
 110 3. The signal space is the same as the sample space: $S = \Sigma$. In this case it is more natural to replace
 111 the bipartite graph by one that has only one copy of each sample/signal, is no longer bipartite,
 112 and that represents the possibility of changing sample/signal u to sample/signal v by a directed
 113 edge (u, v) .

114 We will be interested in the probability of accepting good or bad types after T rounds (i.e., after the
 115 agent draws T samples). We call the T signals submitted a *report* $\mathcal{R} \in \Sigma^T$. The principal gets to
 116 choose an acceptance function (or policy, which could be randomized) $f : \mathcal{R} \rightarrow \{0, 1\}$ that maps
 117 the report into a binary decision. Her goal is to accept the good agent and reject the bad agent. The
 118 agent wants to be accepted regardless of his type. The principal can thus make two types of mistakes:
 119 false-positive (or type 1 error) when she accepts a bad agent, and false-negative (type 2 error) when
 120 she rejects a good agent. The principal wants to minimize the maximum probability of making either
 121 type of mistakes.

122 We recall the following definition of the total variation distance:

123 **Definition 1** (Total Variation Distance). The total variation distance between two distributions
 124 $x, y \in \Delta(\Sigma)$ over support Σ is defined to be

$$d_{\text{TV}}(x, y) = \frac{1}{2} \|x - y\|_1 = \frac{1}{2} \sum_{\sigma \in \Sigma} |x(\sigma) - y(\sigma)| = \max_{A \subseteq \Sigma} (x(A) - y(A)).$$

125 In our setting, the total variation distance provides a good way to measure the closeness between two
 126 *signal* distributions, which are observable by the principal. We will generalize this definition to our
 127 strategic setting, to measure how close two distributions over the *sample* space are to each other.

128 3 Basic Structural Results

129 In this section, we define a notion that we term “directed total variation distance” d_{DTV} . For two
 130 distributions x and y over samples, $d_{\text{DTV}}(x, y)$ measures how well x can distinguish itself from y
 131 in our strategic setting. As we will see in the later sections, d_{DTV} is a central notion in this paper,
 132 and often dictates the number of samples we need to distinguish the two distributions under strategic
 133 reporting.

134 In Section 3.1, we give the formal definitions of *reporting strategies* and the directed total variation
 135 distance $d_{\text{DTV}}(x, y)$. In Section 3.2, we define another notion $\text{MaxSep}(x, y)$ that measures how
 136 well x can distinguish itself from y from the principal’s perspective, using separating sets instead
 137 of reporting strategies. In Section 3.3, we present one of our key structural results (Proposition 1),
 138 which shows that the two notions are equivalent.

139 3.1 Directed Total Variation Distance

140 Before investigating distinguishing distributions under strategic reporting, we first generalize the
 141 classical measure of how close two distributions are, d_{TV} , to our strategic setting. We first give a
 142 formal definition the *reporting strategy* used by the agents.

143 **Definition 2** ((Single-Round) Reporting Strategy). Given $x \in \Delta(S)$, $\alpha \in \Delta(\Sigma)$, we say x can report
 144 α ($x \rightarrow \alpha$), if there exist a *reporting strategy* $R = \{r_{s,\sigma}\}_{(s,\sigma) \in E}$ satisfying:

- 145 • $r_{s,\sigma} \geq 0$ for all $(s, \sigma) \in E$.
- 146 • For each $s \in S$, $\sum_{\sigma: (s,\sigma) \in E} r_{s,\sigma} = 1$.
- 147 • For each $\sigma \in \Sigma$, $\sum_{s: (s,\sigma) \in E} x(s) \cdot r_{s,\sigma} = \alpha(\sigma)$.

148 We say x reports α by strategy R ($x \rightarrow_R \alpha$).

149 In other words, when each sample $s \in S$ is drawn from the distribution x and given this sample
 150 the agent is reporting $\sigma \in \Sigma$ with probability $r_{s,\sigma}$, the resulting distribution over the signal space is
 151 exactly α . For a fixed sample or a random variable s , we use $R(s) \in \Delta(\Sigma)$ to denote the random
 152 variable whose distribution over the signal space is induced by $\{R_{s,\sigma}\}_{\sigma \in \Sigma}$.

Given the definition of reporting strategies, we are ready to generalize d_{TV} to our setting. Intuitively, x chooses a report first, and then y chooses a report in response; they play a zero-sum game where x wants the reports to be as far away from each other as possible. $d_{\text{DTV}}(x, y)$ is the value of this two-player game when x must choose a report (i.e., a pure strategy) first, which measures how far x can stay away from y .

Definition 3 (Directed Total Variation Distance). Given (S, Σ, E) , the directed total variation distance between two distributions $x, y \in \Delta(S)$ over the sample space S is defined to be

$$d_{\text{DTV}}(x, y) = \max_{\alpha: x \rightarrow \alpha} \min_{\beta: y \rightarrow \beta} d_{\text{TV}}(\alpha, \beta).$$

3.2 (Maximum) Separation

Directed total variance distance nicely characterizes the distance between two distributions from the agent's perspective, but it is not immediately clear how that might help the principal. In particular, are two distributions easily separable by setting an appropriate policy if they have large directed total variation distance? To study this, we introduce several concepts to model the problem from the principal's perspective.

Definition 4 (Preimage of Signals). For any set of signals $A \subseteq \Sigma$, the preimage $\text{pre}(A)$ of A is defined to be the set of samples which can be mapped to a signal in A . That is

$$\text{pre}(A) = \{s \in S \mid \exists \sigma \in A, \text{ s.t. } (s, \sigma) \in E\}.$$

The principal could label a set A of signals as “good” signals and simply measure how many good signals the agent is able to send. Ideally, this A is chosen so that a good agent can send (significantly) more signals in A than a bad agent. This inspires the following definitions.

Definition 5 (Separation). For any $A \subseteq \Sigma$, if $x(\text{pre}(A)) - y(\text{pre}(A)) = \epsilon > 0$, then we say A separates x from y by a margin of ϵ .

Definition 6 (Max Separation). The max separation of $x \in \Delta(S)$ from $y \in \Delta(S)$ over the sample space S is defined to be $\text{MaxSep}(x, y) = \max_{A \subseteq \Sigma} (x(\text{pre}(A)) - y(\text{pre}(A)))$.

3.3 Equivalence of d_{DTV} and MaxSep

We now draw the connection between the agent's and the principal's perspectives. The following proposition can be viewed as a generalization of the classic Hall's Marriage Theorem. Proposition 1 states that g can distinguish itself from b under strategic reporting iff there exists a subset A^* of signals so that g can generate more signals in A^* than b . Equivalently, the best reporting strategy for g is to focus on a subset A^* of the signal space, and try to convert samples into signals in A^* whenever possible.

Proposition 1. For any $x, y \in \Delta(S)$, $d_{\text{DTV}}(x, y) = \text{MaxSep}(x, y)$.

The proof of the proposition, as well as all other proofs, is deferred to the appendix. This equivalence between d_{DTV} and MaxSep not only is a nice structural result; Proposition 1 plays a substantial part in our main algorithmic results.

It is worth noting that $d_{\text{DTV}}(x, y)$ in general is not equal to $d_{\text{DTV}}(y, x)$. However, the triangle inequality still holds for d_{DTV} , which also enables some of our main results.

Proposition 2. For any $x, y, z \in \Delta(S)$, $d_{\text{DTV}}(x, y) + d_{\text{DTV}}(y, z) \geq d_{\text{DTV}}(x, z)$.

4 Structural and Computational Results in the General Case

In this section, we define adaptive and non-adaptive reporting strategies (Definition 7), and the accepting probabilities of the optimal reporting strategies after T rounds (Definition 8). At a high level, we give a tight characterization result on when there exists a policy that can distinguish g from b under strategic reporting, and provide an asymptotically tight bound on the sample complexity of the optimal policy. Moreover, we show that while our structural result is clean and tight, it is computationally hard to check if the condition holds. That is, in the general case, it is NP-hard to determine whether there is a policy that can distinguish g from b .

More specifically, we first show that there exists a policy that can distinguish g from b in the limit (when $T \rightarrow \infty$) iff $d_{\text{DTV}}(g, b) > 0$ (Theorem 1). Next, we give an asymptotically tight sample complexity bound of $T = \Theta(1/\epsilon^2)$ when $d_{\text{DTV}}(g, b) = \epsilon$ and we want to distinguish g from b with high constant probability (Theorem 3). We then extend the existence result to more general settings when there are multiple good and bad distributions (Theorem 4). Finally, we show that it is NP-hard to decide if we are in the case where $d_{\text{DTV}}(g, b) = 0$ or $d_{\text{DTV}}(g, b) > \frac{1}{\text{poly}(m, n)}$ (Theorem 2).

We start with the definition of adaptive reporting strategies.

Definition 7 (Adaptive Reporting Strategy). An adaptive reporting strategy $\mathcal{R} = (R^1, \dots, R^T)$ is a sequence of (different) reporting strategies. The signal σ^i at time i is obtained by applying R^i to the sample s^i at time i . $R^i = R^i(\sigma^1, \dots, \sigma^{i-1})$ may depend on all past signals. A reporting strategy is non-adaptive if $R^i = R^1$ for any i and $(\sigma^1, \dots, \sigma^{i-1})$, and adaptive otherwise. For an adaptive policy $\mathcal{R} = (R^1, \dots, R^T)$, we interchangeably write $\sigma^i = R^i(s^i \mid \sigma^1, \dots, \sigma^{i-1})$ to indicate the dependence of R^i on $\sigma^1, \dots, \sigma^{i-1}$.

When we analyze the quality of a fixed T -round policy f , we are interested in the probability that f accepts g or b after T rounds, when the agent (of either type) best-responds to f .

Definition 8 (Acceptance Probabilities of the Best Reporting Strategies). Given $x \in \Delta(S)$, $T \in \mathbb{N}$, and the principal's policy f , let the acceptance rate under adaptive / non-adaptive reporting respectively be

$$p_{\text{ada}}(f, x, T) = \max_{\mathcal{R}=(R^1, \dots, R^T)} \mathbb{E}[f((R^i(s^i))_{i \in [T]})],$$

$$p_{\text{non}}(f, x, T) = \max_{\mathcal{R}=(R, \dots, R)} \mathbb{E}[f((R^i(s^i))_{i \in [T]})]$$

where the expectations are taken over T i.i.d. samples $(s^i)_i$ drawn from x . Observe that $p_{\text{ada}}(f, x, T) \geq p_{\text{non}}(f, x, T)$ for any f, x and T .

Intuitively, if $d_{\text{DTV}}(g, b) = 0$, then the bad distribution can mimic the good distribution perfectly in the signal space, no matter what reporting strategy g uses. Therefore, it is impossible to distinguish g from b . The next theorem formalizes this intuition. In particular, even if g reports adaptively, b can still mimic g 's conditional reporting strategy in every situation (i.e., for every combination of previously reported signals).

Theorem 1 (Separability in the Limit). *Given good and bad distributions g and b :*

(i) *If $d_{\text{DTV}}(g, b) > 0$, then there exists a policy f such that*

$$\lim_{T \rightarrow \infty} (p_{\text{non}}(f, g, T) - p_{\text{ada}}(f, b, T)) = 1.$$

That is, f accepts g and rejects b with probability 1 in the limit.

(ii) *If $d_{\text{DTV}}(g, b) = 0$, then for any policy f and any T ,*

$$p_{\text{ada}}(f, g, T) \leq p_{\text{ada}}(f, b, T), p_{\text{non}}(f, g, T) \leq p_{\text{non}}(f, b, T).$$

That is, no policy can separate g from b , regardless of whether the setting is adaptive.

The next theorem states that while our characterization result (Theorem 1) is clean and tight (we can distinguish iff $d_{\text{DTV}}(g, b) > 0$), it is in fact computationally hard to check if this condition holds. Intuitively, Theorem 2 constructs an instance where the good distribution needs to focus on as few signals as possible. The parameters are chosen carefully so that it is crucial that g finds a subset of signals $A \subseteq \Sigma$ with minimum cardinality that covers the support of g .

Theorem 2 (hardness of checking separability). *Given $x, y \in \Delta(S)$, it is NP-hard to distinguish between the following two cases: (1) $d_{\text{DTV}}(x, y) = 0$ and (2) $d_{\text{DTV}}(x, y) \geq \frac{1}{\text{poly}(m, n)}$, or equivalently, to determine the existence of a set $A \subseteq \Sigma$ such that $x(\text{pre}(A)) - y(\text{pre}(A)) \geq \frac{1}{\text{poly}(m, n)}$.*

Note that the hardness of checking the existence of separating sets implies the hardness of finding any separating set given that $d_{\text{DTV}}(x, y) > 0$. This is because given an algorithm for the latter problem, one could run that algorithm without knowing whether $d_{\text{DTV}}(x, y) > 0$ and see if it succeeds. Either the algorithm returns a separating set, or we know it must be the case that $d_{\text{DTV}}(x, y) = 0$ and no separating set exists.

Next, we focus on the case when there are finitely many samples. Theorem 3 is more refined than Theorem 1, in that it gives a tight sample complexity bound instead of only talking about distinguishing g and b in the limit.

Theorem 3 (Sample Complexity with Two Distributions). *For any g and b such that $d_{\text{DTV}}(g, b) \geq \epsilon$:*

- *There is a policy f such that for any $\delta > 0$ and $T \geq 2 \ln(1/\delta)/\epsilon^2$, $p_{\text{non}}(f, g, T) \geq 1 - \delta$ and $p_{\text{ada}}(f, b, T) \leq \delta$.*
- *When $d_{\text{DTV}}(g, b) = \epsilon$ and $T = o(1/\epsilon^2)$, for any f , $p_{\text{non}}(f, g, T) - p_{\text{non}}(f, b, T) < \frac{1}{3}$.*

Theorem 3 can be generalized to the case where there are multiple good and bad distributions. First, suppose there is one good distribution and multiple bad distributions. As long as $d_{\text{DTV}}(g, b_j) \geq \epsilon$ for every bad distribution b_j , we can use the testing algorithm in Theorem 3 to distinguish them in $T = O(1/\epsilon^2)$ rounds (with high constant probability). We potentially need to do so separately for every bad distribution, paying an extra factor of $\Omega(\ell)$ in the sample complexity if there are ℓ bad distributions. If there are k good distributions, then we can run the k testers in parallel, paying an additional factor of $\log(k)$ in the sample complexity to boost the success probability so that we can take a union bound.

Theorem 4 (Multiple Good and Bad Distributions, the General Case). *For any g_1, \dots, g_k and b_1, \dots, b_ℓ such that $d_{\text{DTV}}(g_i, b_j) \geq \epsilon$ for any $i \in [k]$ and $j \in [\ell]$, there is a policy f such that: For any $\delta > 0$ and $T \geq 2\ell \ln(k\ell/\delta)/\epsilon^2$, $p_{\text{ada}}(f, g_i, T) \geq 1 - \delta$ for any $i \in [k]$, and $p_{\text{ada}}(f, b_j, T) \leq \delta$ for any $j \in [\ell]$.*

We note that the policy in Theorem 4 requires the good distribution to report in different ways, which is not possible with a non-adaptive strategy according to our definition. In particular, the good distribution must know which bad distribution it is up against in each phase, and report accordingly. As our introductory example shows, this is in fact necessary when there are multiple bad distributions.

5 When Signals Are Partially Ordered

In many real-world situations, the sample and signal spaces are structured. For example, when a band is recruiting new members, applicants may be asked to submit video recordings of themselves playing. An applicant would probably videotape herself playing for an entire event as a sample, and then crop the recording to create a signal that demonstrates only her best performance. This cropping procedure is irreversible: the complete recording may be cropped to keep a part, but from a part, it is impossible to recover the full recording. The signal space in this scenario is partially ordered by the cropping procedure—the samples/signals can be transformed in one direction (shortening), but never the other. Also, there is a “default” signal for each sample, which is simply to submit the complete recording without cropping. The default signal can be transformed into any signal that can be reported from this sample. In this section, we consider the following abstraction of such scenarios:

- $S = \Sigma$,
- $(s, s) \in E$ for any $s \in S$,
- $(s, t) \in E$ and $(t, u) \in E \implies (s, u) \in E$, and
- E is acyclic except for self-cycles.

This abstraction also covers, for example, scenarios where the agent can choose to hide certain samples—any sample can be transformed into a non-sample, but not reversely. Note that given the above conditions, the sample/signal space is essentially a partially ordered set, where a sample can only be transformed according to this partial order. Let $n = |S|$ be the cardinality of the sample/signal space.

We first show some useful structural results in the partially ordered case. The following proposition demonstrates that the revelation principle holds in this case.

Proposition 3 (Revelation Principle). *For any policy f :*

- *There exists a policy f' such that for any $x \in \Delta(S)$, $T \in \mathbb{N}$,*

$$p_{\text{non}}(f, x, T) = p_{\text{non}}(f', x, T) = \mathbb{E}[f'((s^i)_i)].$$
- *There exists a policy f'' such that for any $x \in \Delta(S)$, $T \in \mathbb{N}$,*

$$p_{\text{ada}}(f, x, T) = p_{\text{ada}}(f'', x, T) = p_{\text{non}}(f'', x, T) = \mathbb{E}[f''((s^i)_i)].$$

The next proposition simplifies the definition of d_{DTV} in the partially ordered case, based on the insight that, per the revelation principle, the best way for x to avoid being mimicked by y is to always report the unmodified samples.

Proposition 4 (d_{DTV} Simplified). *In the transitive case, $d_{\text{DTV}}(x, y) = \min_{y \rightarrow y'} d_{\text{TV}}(x, y')$.*

This also gives us an efficient algorithm for finding the set that supports the max separation $\text{MaxSep}(x, y)$ of x from y :

Corollary 1 (Efficient Computation of Max Separation). *Given any $x, y \in \Delta(S)$, there is a poly-time algorithm which computes a set A^* satisfying $x(\text{pre}(A^*)) - y(\text{pre}(A^*)) = \text{MaxSep}(x, y)$.*

We show in Theorem 5 that in the partially ordered case we can separate multiple good distributions from multiple bad ones with much smaller overhead. The proof of Theorem 5 is similar to that of Theorem 4. The only difference is that, because of the revelation principle, we no longer require good distributions to report adaptively.

Theorem 5 (Multiple Good and Bad Distributions: The Partially Ordered Case). *For any g_1, \dots, g_k and b_1, \dots, b_ℓ where $d_{\text{DTV}}(g_i, b_j) \geq \epsilon$ for any $i \in [k], j \in [\ell]$, there is a policy f such that: For any $\delta > 0$ and $T \geq 2 \ln(k\ell/\delta)/\epsilon^2$, $p_{\text{non}}(f, g_i, T) \geq 1 - \delta$ for any $i \in [k]$, and $p_{\text{ada}}(f, b_j, T) \leq \delta$ for any $j \in [\ell]$.*

In the partially ordered case, we cannot only deal with multiple good and bad distributions much more efficiently, but also deal with any bad distribution using a single sample-efficient policy. Before stating the result, recall the following definition of the *width* of a partially ordered set.

Definition 9 (Width of Partially Ordered Sets). The width $\rho(G)$ of a partially ordered set represented as graph $G = (S, E)$ is defined to be $\rho(G) = \max\{|A| \mid A \subseteq S, \forall s_1, s_2 \in A, (s_1, s_2) \notin E\}$. In other words, the width is the maximum size of a set $A \subseteq S$ where any two elements in A are not comparable. Such a set A is called an *anti-chain*.

We now provide our generic policy, whose sample complexity, quite surprisingly, depends roughly linearly on the width of the sample space.

Theorem 6 (Efficient Policy against Any Bad Distribution). *For any $g \in \Delta(S)$, there is a policy f such that for any $\delta > 0$, and $T \geq \frac{2\rho \ln(1+n/\rho) \ln(1/\delta)}{\epsilon^2}$: (1) $p_{\text{non}}(f, g, T) \geq 1 - \delta$, and (2) for any b such that $d_{\text{DTV}}(g, b) \geq \epsilon$, $p_{\text{ada}}(f, b, T) \leq \delta$. Moreover, the outcome of the policy can be computed in polynomial time.*

The above policy is able to detect any bad distribution with adaptive reporting. For bad distributions without adaptive reporting, when $\rho = \Omega(\sqrt{n}/\log n)$, the following policy achieves even better sample complexity.

Theorem 7 (Efficient Policy against Non-adaptive Bad Distributions). *For any $g \in \Delta(S)$, there is a policy f such that for any $\delta > 0$, with $T = O\left(\frac{\sqrt{n} \ln(1/\delta)}{\epsilon^2}\right)$ samples: (1) $p_{\text{non}}(f, g, T) \geq 1 - \delta$, and (2) for any b such that $d_{\text{DTV}}(g, b) \geq \epsilon$, $p_{\text{non}}(f, b, T) \leq \delta$. Moreover, the outcome of the policy can be computed in polynomial time.*

6 Future research

In this paper, we have focused on distinguishing good and bad types with near certainty. In reality, the number of available samples may not always be sufficient for this. If so, it may be worthwhile to move beyond simple acceptance and rejection decisions to a more general mechanism design setup. For example, when the signals we receive from an agent are not decisive one way or another, perhaps an intermediate outcome between rejection and acceptance allows us to improve our objective, by avoiding the damage of either accepting a bad type or rejecting a good type. One may also consider settings in which signaling is costly (or at least sending high-quality signals comes at an effort cost, in line with traditional signaling models [16]) or in which agents can in fact improve their actual types via some investment cost. Any of these directions would further enrich the specific connections between mechanism design and learning theory that we have begun to explore in this paper (and that in turn complement other fascinating connections between these topics that have earlier been established by others [1, 11, 2, 10, 3, 4, 14, 7]).

References

- [1] Pranjal Awasthi, Avrim Blum, Nika Haghtalab, and Yishay Mansour. Efficient PAC Learning from the Crowd. In *Conference on Learning Theory*, pages 127–150, 2017.
- [2] Avrim Blum, Nika Haghtalab, Ariel D Procaccia, and Mingda Qiao. Collaborative PAC Learning. In *Advances in Neural Information Processing Systems*, pages 2392–2401, 2017.
- [3] Yiling Chen, Nicole Immorlica, Brendan Lucier, Vasilis Syrgkanis, and Juba Ziani. Optimal data acquisition for statistical estimation. In *Proceedings of the 2018 ACM Conference on Economics and Computation*, pages 27–44. ACM, 2018.
- [4] Yiling Chen, Chara Podimata, Ariel D Procaccia, and Nisarg Shah. Strategyproof linear regression in high dimensions. In *Proceedings of the 2018 ACM Conference on Economics and Computation*, pages 9–26. ACM, 2018.
- [5] Ofer Dekel, Felix Fischer, and Ariel D. Procaccia. Incentive compatible regression learning. In *Proceedings of the Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 884–893, Philadelphia, PA, USA, 2008. Society for Industrial and Applied Mathematics.
- [6] Ilias Diakonikolas, Daniel M Kane, and Vladimir Nikishkin. Testing identity of structured distributions. In *Proceedings of the twenty-sixth annual ACM-SIAM symposium on Discrete algorithms*, pages 1841–1854. Society for Industrial and Applied Mathematics, 2015.
- [7] Jinshuo Dong, Aaron Roth, Zachary Schutzman, Bo Waggoner, and Zhiwei Steven Wu. Strategic classification from revealed preferences. In *Proceedings of the 2018 ACM Conference on Economics and Computation*, pages 55–70. ACM, 2018.
- [8] Jerry Green and Jean-Jacques Laffont. Partially verifiable information and mechanism design. *Review of Economic Studies*, 53:447–456, 1986.
- [9] Moritz Hardt, Nimrod Megiddo, Christos Papadimitriou, and Mary Wootters. Strategic classification. In *Innovations in Theoretical Computer Science (ITCS)*, Cambridge, MA, USA, 2016.
- [10] Lily Hu, Nicole Immorlica, and Jennifer Wortman Vaughan. The disparate effects of strategic manipulation. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 259–268. ACM, 2019.
- [11] Shahin Jabbari, Ryan M Rogers, Aaron Roth, and Steven Z Wu. Learning from rational behavior: Predicting solutions to unknown linear programs. In *Advances in Neural Information Processing Systems*, pages 1570–1578, 2016.
- [12] Andrew Kephart and Vincent Conitzer. Complexity of mechanism design with signaling costs. In *Proceedings of the Fourteenth International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*, pages 357–365, Istanbul, Turkey, 2015.
- [13] Andrew Kephart and Vincent Conitzer. The revelation principle for mechanism design with reporting costs. In *Proceedings of the Seventeenth ACM Conference on Economics and Computation (EC)*, pages 85–102, Maastricht, the Netherlands, 2016.
- [14] Annie Liang, Xiaosheng Mu, and Vasilis Syrgkanis. Optimal and myopic information acquisition. In *Proceedings of the 2018 ACM Conference on Economics and Computation*, pages 45–46. ACM, 2018.
- [15] Reshef Meir, Ariel D. Procaccia, and Jeffrey S. Rosenschein. Algorithms for strategyproof classification. *Artificial Intelligence*, 186:123–156, 2012.
- [16] Michael Spence. Job market signaling. *Quarterly Journal of Economics*, 87(3):355–374, 1973.
- [17] Gregory Valiant and Paul Valiant. An automatic inequality prover and instance optimal identity testing. *SIAM Journal on Computing*, 46(1):429–455, 2017.
- [18] Lan Yu. Mechanism design with partial verification and revelation principle. *Autonomous Agents and Multi-Agent Systems*, 22(1):217–223, 2011.
- [19] Hanrui Zhang, Yu Cheng, and Vincent Conitzer. When samples are strategically selected. In *Thirty-sixth International Conference on Machine Learning*, 2019.

387 A Omitted Proofs From Section 3

388 We need the following fact:

389 **Proposition 5** (Saturation). *If $x \rightarrow \alpha$, then for any $A \subseteq \Sigma$,*

$$x(\text{pre}(A)) \geq \alpha(A).$$

390 *Moreover, there exists α_A where $x \rightarrow \alpha_A$, such that*

$$x(\text{pre}(A)) = \alpha_A(A).$$

391 *We call the corresponding reporting strategy that achieves $x \rightarrow \alpha_A$ “saturating” for A .*

392 *Proof of Proposition 5.* Let $R = \{r_{s,\sigma}\}_{(s,\sigma) \in E}$ be the reporting strategy by which x reports α .

$$\begin{aligned} x(\text{pre}(A)) &= \sum_{s \in \text{pre}(A)} x(s) \\ &\geq \sum_{s \in \text{pre}(A)} \sum_{\sigma \in A} r_{s,\sigma} x(s) && (\sum_{\sigma \in A} r_{s,\sigma} \leq 1) \\ &= \sum_{\sigma \in A} \sum_{s: (s,\sigma) \in E} r_{s,\sigma} x(s) \\ &= \sum_{\sigma \in A} \alpha(\sigma) && (\text{definition of } R) \\ &= \alpha(A). \end{aligned}$$

393 Now we show α_A exists by constructing the corresponding reporting strategy. Let $R' = \{r'_{s,\sigma}\}$ be
 394 any reporting strategy satisfying: if $s \in \text{pre}(A)$, $r'_{s,\sigma} = 0$ for all $\sigma \notin A$. Such an R' exists because
 395 by the definition of $\text{pre}(A)$, for every $s \in \text{pre}(A)$, there is at least one $\sigma \in A$ that connects to s .

396 Now for any $s \in \text{pre}(A)$,

$$\sum_{\sigma \in A} r'_{s,\sigma} = 1.$$

397 Hence, for this reporting strategy, the single inequality in the derivation above becomes an equality,
 398 allowing us to conclude $x(\text{pre}(A)) = \alpha_A(A)$. \square

399 *Proof of Proposition 1.* We first show $\text{MaxSep}(x, y) \leq d_{\text{DTV}}(x, y)$. Let $A^* =$
 400 $\text{argmax}_A (x(\text{pre}(A)) - y(\text{pre}(A)))$.

$$\begin{aligned} d_{\text{DTV}}(x, y) &= \max_{\alpha: x \rightarrow \alpha} \min_{\beta: y \rightarrow \beta} d_{\text{TV}}(\alpha, \beta) \\ &\geq \max_{\alpha: x \rightarrow \alpha} \min_{\beta: y \rightarrow \beta} \sum_{\sigma \in A^*} \max\{\alpha(\sigma) - \beta(\sigma), 0\} && (\text{Definition 1 of } d_{\text{TV}}) \\ &\geq \max_{\alpha: x \rightarrow \alpha} \min_{\beta: y \rightarrow \beta} (\alpha(A^*) - \beta(A^*)) \\ &\geq \max_{\alpha: x \rightarrow \alpha} (\alpha(A^*) - y(\text{pre}(A^*))) && (\text{Proposition 5}) \\ &= x(\text{pre}(A^*)) - y(\text{pre}(A^*)) && (\text{Proposition 5, existence of saturating distribution}) \\ &= \text{MaxSep}(x, y). \end{aligned}$$

401 Now we show $\text{MaxSep}(x, y) \geq d_{\text{DTV}}(x, y)$. Let α^* be a signal distribution reported by x that
 402 achieves $d_{\text{DTV}}(x, y)$. Let β^* be a signal distribution reported by y that best-response to α^* , where
 403 we require as a tie-breaker that β^* minimizes the number of signals σ with $\alpha^*(\sigma) \geq \beta^*(\sigma)$.

404 Let $A^* = \{\sigma \mid \alpha^*(\sigma) \geq \beta^*(\sigma)\}$. We will show that A^* separates x from y by a margin of
 405 $d_{\text{DTV}}(x, y)$.

406 We first show that $\beta^*(A^*) = y(\text{pre}(A^*))$. Suppose otherwise $\beta^*(A^*) < y(\text{pre}(A^*))$. Let $R =$
 407 $\{r_{s,\sigma}\}$ be the reporting strategy that gives $y \rightarrow \beta^*$. We know that there exists some $s_0 \in \text{pre}(A^*)$

408 with $y(s_0) > 0$ where R does not convert all probability mass on s_0 into signals in A^* . Formally,
 409 we have $\sum_{\sigma \in A^*: (s_0, \sigma) \in E} r_{s_0, \sigma} < 1$. Consider any $\sigma_1, \sigma_2 \in \Sigma$ satisfying: $\sigma_1 \notin A^*$, $(s_0, \sigma_1) \in E$,
 410 $r_{s_0, \sigma_1} > 0$, $\sigma_2 \in A^*$, and $(s_0, \sigma_2) \in E$. We have $\alpha^*(\sigma_1) < \beta^*(\sigma_1)$ and $\alpha^*(\sigma_2) \geq \beta^*(\sigma_2)$. Now we
 411 discuss the following two cases and show there is a contradiction in both cases.

412 • If $\alpha^*(\sigma_2) > \beta^*(\sigma_2)$, then by moving

$$\min\{r_{s_0, \sigma_1} y(s_0), \beta^*(\sigma_1) - \alpha^*(\sigma_1), \alpha^*(\sigma_2) - \beta^*(\sigma_2)\} > 0$$

413 mass from σ_1 to σ_2 , y can report β' such that $d_{TV}(\alpha^*, \beta') < d_{TV}(\alpha^*, \beta^*)$, a contradiction.

414 • If $\alpha^*(\sigma_2) = \beta^*(\sigma_2)$, then by moving

$$\min\{r_{s_0, \sigma_1} y(s_0), (\beta^*(\sigma_1) - \alpha^*(\sigma_1))/2\} > 0$$

415 mass from σ_1 to σ_2 , y can report β' such that $d_{TV}(\alpha^*, \beta^*) = d_{TV}(\alpha^*, \beta')$. But now $\alpha^*(\sigma_2) -$
 416 $\beta'(\sigma_2) < 0$, and for any $\sigma \neq \sigma_2$, the sign of $\alpha^*(\sigma) - \beta'(\sigma)$ is the same as that of $\alpha^*(\sigma) - \beta^*(\sigma)$.
 417 So we have

$$|\{\sigma \mid \alpha^*(\sigma) \geq \beta^*(\sigma)\}| > |\{\sigma \mid \alpha^*(\sigma) \geq \beta'(\sigma)\}|,$$

418 which contradicts the choice of β^* .

419 Now given that $y(\text{pre}(A^*)) = \beta^*(A^*)$, we have

$$\begin{aligned} \text{MaxSep}(x, y) &= \max_A (x(\text{pre}(A)) - y(\text{pre}(A))) \\ &\geq x(\text{pre}(A^*)) - y(\text{pre}(A^*)) \\ &\geq \alpha^*(A^*) - y(\text{pre}(A^*)) && \text{(Proposition 5)} \\ &= \alpha^*(A^*) - \beta^*(A^*) \\ &= d_{TV}(\alpha, \beta) \\ &= d_{DTV}(x, y). \end{aligned} \quad \square$$

420 *Proof of Proposition 2.* Let $A^* = \arg\max_A (x(\text{pre}(A)) - z(\text{pre}(A)))$. We have

$$\begin{aligned} d_{DTV}(x, y) + d_{DTV}(y, z) &= \text{MaxSep}(x, y) + \text{MaxSep}(y, z) \\ &= \max_A (x(\text{pre}(A)) - y(\text{pre}(A))) + \max_A (y(\text{pre}(A)) - z(\text{pre}(A))) \\ &\geq (x(\text{pre}(A^*)) - y(\text{pre}(A^*))) + (y(\text{pre}(A^*)) - z(\text{pre}(A^*))) \\ &= x(\text{pre}(A^*)) - z(\text{pre}(A^*)) \\ &= \text{MaxSep}(x, z) \\ &= d_{DTV}(x, z). \end{aligned} \quad \square$$

421 B Omitted Proofs From Section 4

422 *Proof of Theorem 1.* Part (i) follows from Theorem 3.

423 For part (ii), suppose $d_{DTV}(g, b) = 0$. Let s_g^i (resp. s_b^i) be a random variable that denotes the sample
 424 drawn from g (resp. b) at time i . Abusing notation, for two random variables X and Y , we write
 425 $d_{TV}(X, Y)$ for the d_{TV} between the underlying distributions of X and Y .

426 We show that given an adaptive / non-adaptive \mathcal{R}_g , there is an adaptive / non-adaptive \mathcal{R}_b , such that

$$d_{TV}((R_g^i(s_g^i))_{i \in [T]}, (R_b^i(s_b^i))_{i \in [T]}) = 0. \quad (1)$$

427 Because the good and bad distributions have identical distributions over the signal space, and this
 428 holds for all possible reporting strategies \mathcal{R}_g , part (ii) follows immediately.

429 Consider first non-adaptive reporting. Fix $\mathcal{R}_g = (R_g^1, \dots, R_g^T)$ where $R_g^i = R_g$ for all i , let
 430 $\mathcal{R}_b = (R_b^1, \dots, R_b^T)$, where

$$d_{TV}(R_g^i(s_g^i), R_b^i(s_b^i)) = 0.$$

431 The existence of such an \mathcal{R}_b follows from the fact that $d_{DTV}(g, b) = 0$. Now since $R_g^i(s_g^i)$ and
 432 $R_b^i(s_b^i)$ are i.i.d., Equation (1) holds.

Now consider adaptive reporting. For any adaptive reporting strategy \mathcal{R}_g , we will construct an adaptive \mathcal{R}_b inductively, such that for any k ,

$$d_{\text{TV}}((R_g^i(s_g^i))_{i \in [k]}, (R_b^i(s_b^i))_{i \in [k]}) = 0.$$

For the base case when $k = 1$, observe that since $d_{\text{DTV}}(g, b) = 0$, for any R_g^1 , there exists R_b^1 such that

$$d_{\text{TV}}(R_g^1(s_g^1), R_b^1(s_b^1)) = 0.$$

For the inductive case, suppose that $d_{\text{TV}}((R_g^i(s_g^i))_{i \in [k]}, (R_b^i(s_b^i))_{i \in [k]}) = 0$. Given (R_b^1, \dots, R_b^k) , we construct R_b^{k+1} in the following way. Let R_b^{k+1} be such that

$$R_b^{k+1}(s_b^{k+1} \mid \sigma^1, \dots, \sigma^k) = R_g^{k+1}(s_g^{k+1} \mid \sigma^1, \dots, \sigma^k),$$

for any $(\sigma^1, \dots, \sigma^k)$. Now for any $(\sigma^1, \dots, \sigma^{k+1})$,

$$\begin{aligned} & \Pr[(R_b^1(s_b^1), \dots, R_b^{k+1}(s_b^{k+1})) = (\sigma^1, \dots, \sigma^{k+1})] \\ &= \Pr[(R_b^1(s_b^1), \dots, R_b^k(s_b^k)) = (\sigma^1, \dots, \sigma^k)] \cdot \Pr[R_b^{k+1}(s_b^{k+1} \mid \sigma^1, \dots, \sigma^k) = \sigma^{k+1}] \\ &= \Pr[(R_g^1(s_g^1), \dots, R_g^k(s_g^k)) = (\sigma^1, \dots, \sigma^k)] \cdot \Pr[R_b^{k+1}(s_b^{k+1} \mid \sigma^1, \dots, \sigma^k) = \sigma^{k+1}] \\ & \quad \text{(induction hypothesis)} \\ &= \Pr[(R_g^1(s_g^1), \dots, R_g^k(s_g^k)) = (\sigma^1, \dots, \sigma^k)] \cdot \Pr[R_g^{k+1}(s_g^{k+1} \mid \sigma^1, \dots, \sigma^k) = \sigma^{k+1}] \\ & \quad \text{(construction of } R_b^{k+1}) \\ &= \Pr[(R_g^1(s_g^1), \dots, R_g^{k+1}(s_g^{k+1})) = (\sigma^1, \dots, \sigma^{k+1})]. \end{aligned}$$

In other words, we have

$$d_{\text{TV}}((R_g^i(s_g^i))_{i \in [k+1]}, (R_b^i(s_b^i))_{i \in [k+1]}) = 0,$$

which concludes the inductive proof for Equation (1) in the adaptive case. \square

Proof of Theorem 2. We reduce from Set Cover. More specifically, we use the following decision version of Set Cover: given ground set $X = [n]$, family of sets $\mathcal{F} = \{F_1, \dots, F_m\}$ where $F_i \subseteq X$, and integer $k = m/2$, determine whether there are k sets in \mathcal{F} whose union is X . Note that it is without generality to set $k = m/2$, since given any Set Cover instance with an arbitrary k , we could always pad the instance by adding at most m elements into X and m sets into \mathcal{F} , to obtain an equivalent new instance with $k' = m'/2$. Fixing a Set Cover instance, we construct S, Σ, E, x and y in the following way.

- $S = U \cup V$, where $U = \{u_1, \dots, u_{n+1}\}$, $V = \{v_1, \dots, v_{m+2}\}$, and $U \cap V = \emptyset$.
- $\Sigma = \{\sigma_1, \dots, \sigma_{m+1}\}$.
- $x(u_i) = \frac{1}{2n}$ for $i \in [n]$, and $x(u_{n+1}) = \frac{1}{2}$.
- $y(v_i) = (1/2 + 1/(2n) - t)/m$ for $i \in [m]$, $y(v_{m+1}) = \frac{1}{2} - \frac{1}{2n}$, and $y(v_{m+2}) = t$, where $t \in [0, 1/2 + 1/(2n)]$ is a constant to be determined later.
- For $i \in [m]$ and $j \in F_i$, let $(u_j, \sigma_i) \in E$. Let $(u_{n+1}, \sigma_{m+1}) \in E$.
- For any $i \in [m]$, let $(v_i, \sigma_i) \in E$. For any $i \in [m+2]$, let $(v_i, \sigma_{m+1}) \in E$. For any $i \in [m+1]$, let $(v_{m+1}, \sigma_i) \in E$.
- E contains only edges that mentioned above.

Now consider the problem of finding a set that separates x from y with a positive margin. First observe that such a set A would never include σ_{m+1} , since $y(\text{pre}(\{\sigma_{m+1}\})) = 1$. Our goal is to set t , such that iff $|A| \leq k$ and $\text{pre}(A) = \{u_1, \dots, u_n\}$, A separates x from y with a positive margin. Such an A in the Set Cover instance would correspond to at most k sets in \mathcal{F} whose union cover X . Note that if $\sigma_{m+1} \notin A$,

$$y(\text{pre}(A)) = \frac{1/2 + 1/(2n) - t}{m} \cdot |A| + \frac{1}{2} - \frac{1}{2n} \geq \frac{1}{2} - \frac{1}{2n}.$$

If $\text{pre}(A)$ covers $\{u_1, \dots, u_n\}$, then $x(\text{pre}(A)) = \frac{1}{2}$. Otherwise, $x(\text{pre}(A)) \leq \frac{1}{2} - \frac{1}{2n} \leq y(\text{pre}(A))$. So if $\text{pre}(A)$ does not cover $\{u_1, \dots, u_n\}$, A cannot be a separating set. We set t such that $y(\text{pre}(A)) = \frac{1}{2}$ if $|A| = k+1 = (m+2)/2$. Such a t always exists. Moreover, observe that such a value of t guarantees that whenever $|A| \leq k$, $y(\text{pre}(A)) \leq \frac{1}{2} - \frac{1}{\text{poly}(m, n)}$. Now iff $|A| \leq k$

467 and A covers $\{u_1, \dots, u_n\}$, A separates x from y with a margin of $\frac{1}{\text{poly}(m, n)}$. In other words, there
 468 is a separating set with a positive margin iff there are at most k sets that cover X in the Set Cover
 469 instance. Our NP-hardness result follows. \square

470 *Proof of Theorem 3.* For the first bullet point, let A^* be a set which separates g from b by a margin
 471 of ϵ . Consider the following policy: accept $(\sigma^1, \dots, \sigma^T)$ iff

$$\frac{1}{T} \sum_{i \in [T]} \mathbb{I}[\sigma^i \in A^*] \geq g(\text{pre}(A^*)) - \frac{1}{2}\epsilon.$$

472 That is, the policy accepts the distribution iff $\bar{\alpha}(A^*) \geq g(\text{pre}(A^*)) - \frac{1}{2}\epsilon$, where $\bar{\alpha}$ is the empirical
 473 distribution of the reported signals. We now bound the probability of g being accepted. Using some
 474 saturating reporting strategy R_{A^*} for A^* (Proposition 5), we have

$$s_g^i \in \text{pre}(A^*) \iff R_{A^*}(s_g^i) \in A^*.$$

475 So by the Chernoff-Hoeffding bound, f rejects g with probability

$$\Pr \left[\frac{1}{T} \sum_i \mathbb{I}[s_g^i \in \text{pre}(A^*)] - g(\text{pre}(A^*)) < -\frac{1}{2}\epsilon \right] \leq \exp(-T\epsilon^2/2) \leq \delta.$$

476 On the other hand, by Proposition 5 for any reporting strategy R_b of b ,

$$\Pr[R_b(s_b^i) \in A^*] \leq b(\text{pre}(A^*)) \leq g(\text{pre}(A^*)) - \epsilon.$$

477 So f accepts b with probability at most

$$\Pr \left[\frac{1}{T} \sum_i \mathbb{I}[s_b^i \in \text{pre}(A^*)] - b(\text{pre}(A^*)) \geq \frac{1}{2}\epsilon \right] \geq \exp(-T\epsilon^2/2) \leq \delta.$$

478 For the second bullet point, consider the following instance: $S = \Sigma = (s_1, s_2)$, $g(s_1) = \frac{1}{2} + \epsilon$,
 479 $g(s_2) = \frac{1}{2} - \epsilon$, $b(s_1) = b(s_2) = \frac{1}{2}$, and $E = \{(s_1, s_1), (s_2, s_2)\}$. In words, s_1 is a good sample/signal,
 480 and s_2 is a bad one. Agents must report the sample drawn as is. The good distribution draws good
 481 samples with slightly higher probability than the bad distribution. For this instance, distinguishing
 482 between g and b is exactly equivalent to distinguishing a coin with bias ϵ with a fair coin. In the latter
 483 problem, it is well-known that $\Omega(1/\epsilon^2)$ samples are required. \square

484 *Proof of Theorem 4.* Consider the following policy which uses the policy in Theorem 3 as a building
 485 block. Let the policy in Theorem 3 be $f_{g,b}$ for good distribution g and bad distribution b . Let
 486 $T_0 = 2 \ln(k\ell/\delta)/\epsilon^2$, where $\ell T_0 = T$. Given the T reported signals (σ^i) , our policy f proceeds in
 487 the following way:

- 488 • For each $i \in [k]$, $j \in [\ell]$, feed the T_0 signals

$$\sigma^{(j-1)T_0+1}, \dots, \sigma^{jT_0}$$

489 to policy f_{g_i, b_j} , and let the output be $o_{i,j} = f_{g_i, b_j}(\sigma^{(j-1)T_0+1}, \dots, \sigma^{jT_0})$.

- 490 • f outputs 1 iff

$$\bigvee_{i \in [k]} \bigwedge_{j \in [\ell]} o_{i,j} = 1.$$

491 To see the correctness of the policy, observe that for each any i, j , with probability $1 - \frac{\delta}{k\ell}$, f_{g_i, b_j}
 492 accepts g_i and rejects b_j given the signals fed in. Taking a union bound over all such (i, j) , with
 493 probability at least $1 - \delta$, all these policies succeed simultaneously. Now for some good distribution
 494 g_{i^*} , as long as the above event happens, we have $o_{i^*, j} = 1$ for all $j \in [\ell]$, so

$$\bigvee_{i \in [k]} \bigwedge_{j \in [\ell]} o_{i,j} \geq \prod_{j \in [\ell]} o_{i^*, j} = 1.$$

495 On the other hand, for some bad distribution b_{j^*} , we have $o_{i, j^*} = 0$ for any $i \in [k]$, and therefore

$$\bigvee_{i \in [k]} \bigwedge_{j \in [\ell]} o_{i,j} \leq \sum_i \prod_j o_{i,j} = 0. \quad \square$$

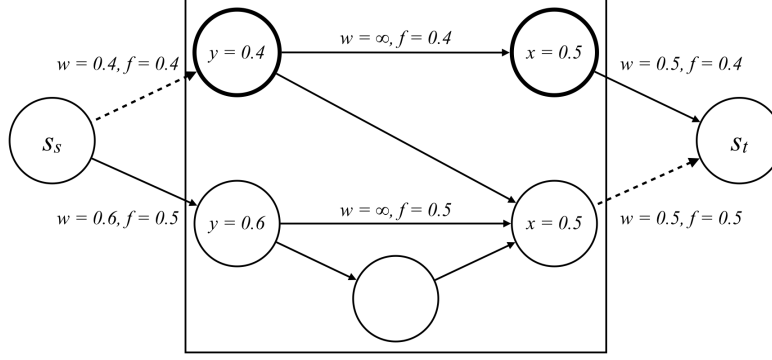


Figure 2: Illustration of Proposition 4. Vertices in the frame are from S , and the rest of the network is constructed as described in the proof. The dashed edges are saturated in the max flow. The boldface vertices are cut to s_t , and therefore constitute the prefix supporting the max separation.

C Omitted Proofs From Section 5

Proof of Proposition 3. Consider f' (resp. f'') which first applies the optimal non-adaptive (resp. adaptive) reporting strategy for x to the original samples, and then applies f to the transformed samples. Now the optimal reporting strategy for x given policy f' (or f'') is simply reporting the original sample received from x . The proposition follows. \square

Proof of Proposition 4. We show that $\text{MaxSep}(x, y) = \min_{y \rightarrow y'} d_{\text{TV}}(x, y')$, which implies the proposition given Proposition 1.

Consider the following flow network $G = (V, E', w)$:

- $V = S \cup \{s_s, s_t\}$, where s_s is the source and s_t is the sink.
- $E' = E \cup \{(s_s, s)\}_{s \in S} \cup \{(s, s_t)\}_{s \in S}$.
- $w(s_1, s_2) = \infty$ for any $(s_1, s_2) \in E$, $w(s_s, s) = b(s)$ for $s \in S$, and $w(s, s_t) = g(s)$ for $s \in S$.

See Figure 2 for illustration of an example network. Now observe that

- $1 - \text{MaxSep}(x, y)$ is the s_s - s_t min-cut of this network. This is because every set $A \subseteq S$ corresponds to a cut, where $S \setminus \text{pre}(A)$ is cut to s_s and $\text{pre}(A)$ is cut to s_t . The value of $1 - (x(\text{pre}(A)) - y(\text{pre}(A)))$ is exactly the value of the cut. Similarly, any cut corresponds to a separating set. It follows that $\text{MaxSep}(x, y)$ corresponds to the min-cut.
- $1 - \min_{y \rightarrow y'} d_{\text{TV}}(x, y')$ is the s_s - s_t max-flow of the network. This is because every y' corresponds to a feasible flow in the network, whose capacity is

$$\sum_s \min\{x(s), y'(s)\} = 1 - d_{\text{TV}}(x, y').$$

Taking max over y' , we see that the max-flow has capacity

$$\max_{y \rightarrow y'} (1 - d_{\text{TV}}(x, y')) = 1 - \min_{y \rightarrow y'} d_{\text{TV}}(x, y').$$

Strong duality immediately gives the desired statement. \square

Proof of Corollary 1. Run max-flow on the flow network constructed in the proof of Proposition 4, compute the min-cut on the residual network, and return the subset of S on the same side as s_s . \square

Proof of Theorem 5. Let the policy in Theorem 3 be the *truthful version* of $f_{g,b}$ for good distribution g and bad distribution b .⁵ Given the T reported signals (σ^i) , our policy f proceeds in the following way:

⁵The policy in Theorem 3 is itself truthful, but the construction here works even if it is not.

- For each $i \in [k]$, $j \in [\ell]$, feed all T signals reported to policy f_{g_i, b_j} , and let the output be $o_{i,j} = f_{g_i, b_j}(\sigma^1, \dots, \sigma^T)$.
- f outputs 1 iff

$$\bigvee_{i \in [k]} \bigwedge_{j \in [\ell]} o_{i,j} = 1.$$

The rest of the proof is essentially the same as that of Theorem 4. \square

Our policy against any adaptive bad distribution in Theorem 6 uses an efficient learner as a building block, which generalizes classical results for learning discrete distributions.

Theorem 8 (Efficient Learner). *Let $\rho = \rho(G)$ be the width of graph $G = (S, E)$. For any $x \in \Delta(S)$, $\epsilon > 0$, $\delta > 0$, and $T = \frac{\rho \ln(1+n/\rho) \ln(1/\delta)}{2\epsilon^2}$, for any valid reporting strategy that satisfies $(s^i, \sigma^i) \in E$, with probability at least $1 - \delta$, $d_{\text{DTV}}(\bar{\alpha}, x) \leq \epsilon$, where $\bar{\alpha}$ is the empirical distribution given by the reports $(\sigma^i)_i$, i.e., $\bar{\alpha}(s) = \frac{\sum_i \mathbb{I}[\sigma^i = s]}{T}$.*

The following well-known fact about the width is used in the analysis of our learner:

Theorem 9 (Dilworth's Theorem). *A chain in a partially ordered set $G = (S, E)$ is an ordered set $C = (c_1, \dots, c_\ell)$, where $c_i \in S$ for $i \in [\ell]$ and $(c_i, c_{i+1}) \in E$ for any $i \in [\ell - 1]$. Dilworth's Theorem states that for any partially ordered set $G = (S, E)$, the width of $\rho(G)$ is equal to the minimum number of chains whose union covers S .*

Proof of Theorem 8. We show that $\text{MaxSep}(\bar{\alpha}, x) \leq \epsilon$ w.p. $1 - \delta$. More specifically, if for all A where $A = \text{pre}(A)$, $\bar{\alpha}(A) - x(A) \leq \epsilon$, then duality gives immediately that $d_{\text{DTV}}(\bar{\alpha}, x) \leq \epsilon$. We will show that this happens with probability $1 - \delta$.

Let \bar{x} be the empirical distribution of $(s^i)_i$. Fix $A \subseteq S$ where $A = \text{pre}(A)$. Observe that $\bar{x}(A) \geq \bar{\alpha}(A)$, so $x(A) = \mathbb{E}[\bar{x}(A)] \geq \mathbb{E}[\bar{\alpha}(A)]$. The Chernoff bound gives

$$\Pr[\bar{\alpha}(A) \geq x(A) + \epsilon] \leq \exp(-2T\epsilon^2) \leq \frac{\delta}{(1+n/\rho)^\rho}.$$

We only need to show that the number of different sets A where $A = \text{pre}(A)$ is at most $(1+n/\rho)^\rho$. We call such sets prefixes of graph (S, E) . Dilworth's Theorem (Theorem 9) states that the width ρ of (S, E) is equal to the minimum number of chains whose union covers S . Let $\mathcal{C} = \{C_k\}_{k \in [\rho]}$ be such a covering family, where for any k , $C_k = (s_{k,1}, \dots, s_{k,\ell_k})$ is a chain (i.e., $(s_{k,i}, s_{k,i+1}) \in E$ for $i \in [\ell_k - 1]$). For any prefix A , let $p_k(A) = |A \cap C_k|$. Observe that if two prefixes A_1 and A_2 are distinct, then there is some $k \in [\rho]$ such that $p_k(A_1) \neq p_k(A_2)$. On the other hand, consider vector $(p_1(A), \dots, p_\rho(A))$. The number of possible values of this vector is $\prod_k (\ell_k + 1) \leq (1+n/\rho)^\rho$, which is an upper bound of the number of different prefixes. Taking union bound over all these prefixes, we have

$$\Pr[\forall A \text{ where } A = \text{pre}(A), \bar{\alpha}(A) \geq x(A) + \epsilon] \leq \frac{\delta}{(1+n/\rho)^\rho} \cdot (1+n/\rho)^\rho = \delta.$$

The theorem follows. \square

Given the efficient learner constructed above, we are ready to prove Theorem 6.

Proof of Theorem 6. Consider the following policy: compute the empirical distribution $\bar{\alpha}$ of the reported signals. Accept iff $d_{\text{DTV}}(g, \bar{\alpha}) < \frac{1}{2}\epsilon$. Note that since g is known, $d_{\text{DTV}}(g, \bar{\alpha})$ can be computed in polynomial time using the algorithm in Corollary 1.

We first show that $p_{\text{non}}(f, g, T) \geq 1 - \delta$. In particular, we show that if g reports truthfully, then with probability $1 - \delta$, $d_{\text{DTV}}(g, \bar{g}) < 1 - \frac{1}{2}\epsilon$. The argument is similar to that in the proof of Theorem 8. For any $A \subseteq S$ where $A = \text{pre}(A)$, the Chernoff bound implies

$$\Pr[g(A) - \bar{g}(A) \geq \epsilon/2] \leq \frac{\delta}{(1+n/\rho)^\rho}.$$

Since there are at most $(1+n/\rho)^\rho$ such sets, from a simple union bound, with probability $1 - \delta$, $d_{\text{DTV}}(g, \bar{g}) = \text{MaxSep}(g, \bar{g}) \leq \frac{1}{2}\epsilon$.

560 Now we show that $p_{\text{ada}}(f, b, T) \leq \delta$ for any b where $d_{\text{DTV}}(g, b) \geq \epsilon$. No matter what adaptive
 561 reporting strategy b uses, the signals reported by b must satisfy $(s_b^i, \sigma_b^i) \in E$ for all i . By Theorem 8,
 562 with probability $1 - \delta$, the empirical distribution $\bar{\alpha}$ satisfies $d_{\text{DTV}}(\bar{\alpha}, b) \leq \frac{1}{2}\epsilon$. Now since d_{DTV}
 563 satisfies the triangle inequality (Proposition 2),

$$d_{\text{DTV}}(g, \bar{\alpha}) \geq d_{\text{DTV}}(g, b) - d_{\text{DTV}}(\bar{\alpha}, b) \geq \epsilon - \frac{1}{2}\epsilon = \frac{1}{2}\epsilon.$$

564 Whenever this happens, b is rejected by f , which means $p_{\text{ada}}(f, b, T) \leq \delta$. □

565 *Proof of Theorem 7.* We use the algorithm by Valiant and Valiant [17] for testing identity of discrete
 566 distributions as a building block. Given a distribution $x \in \Delta([n])$, with $T = O\left(\frac{\sqrt{n} \ln(1/\delta)}{\epsilon^2}\right)$ samples
 567 to an unknown distribution y , their algorithm distinguishes between the following two cases: (1)
 568 $y = x$ and (2) $d_{\text{TV}}(x, y) \geq \epsilon$. Our policy for non-adaptive reporting is simply running the algorithm
 569 by Valiant and Valiant on the good distribution g and the signals reported $(\sigma^i)_i$.

570 The good distribution g , in order to be accepted with high probability, simply reports truthfully. The
 571 distribution of signals of g is therefore exactly g , which with probability $1 - \delta$ passes the test.

572 As for the bad distribution, observe that any non-adaptive reporting strategy $\mathcal{R}_b = (R_b, \dots, R_b)$
 573 induces a distribution α_b of signals reported, where $b \rightarrow_{R_b} \alpha_b$. No matter how b reports, because
 574 $d_{\text{DTV}}(g, b) \geq \epsilon$, we always have $d_{\text{TV}}(g, \alpha_b) \geq \epsilon$, in which case α_b fails the test with probability at
 575 least $1 - \delta$. □