

# Adaptive and oblivious statistical adversaries are equivalent

Guy Blanc

Gregory Valiant

*Stanford*

*Stanford*

September 3, 2025

## Abstract

We resolve a fundamental question about the ability to perform a statistical task, such as learning, when an adversary corrupts the sample. Such adversaries are specified by the types of corruption they can make and their level of knowledge about the sample. The latter distinguishes between sample-adaptive adversaries which know the contents of the sample when choosing the corruption, and sample-oblivious adversaries, which do not. We prove that for all types of corruptions, sample-adaptive and sample-oblivious adversaries are *equivalent* up to polynomial factors in the sample size. This resolves the main open question introduced by [BLMT22] and further explored in [CHL+23].

Specifically, consider any algorithm  $A$  that solves a statistical task even when a sample-oblivious adversary corrupts its input. We show that there is an algorithm  $A'$  that solves the same task when the corresponding sample-adaptive adversary corrupts its input. The construction of  $A'$  is simple and maintains the computational efficiency of  $A$ : It requests a polynomially larger sample than  $A$  uses and then runs  $A$  on a uniformly random subsample.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Our Results</b>	<b>2</b>
2.1	A unified framework to define statistical adversaries . . . . .	2
2.2	Our main result: Adaptive and oblivious adversaries are equivalent . . . . .	3
2.3	Lower bounds . . . . .	4
2.4	Relation to recent work . . . . .	6
<b>3</b>	<b>Instantiating common adversaries in our framework</b>	<b>6</b>
3.1	Partially-adaptive adversaries . . . . .	8
<b>4</b>	<b>Technical overview</b>	<b>8</b>
4.1	A first attempt and why it fails . . . . .	9
4.2	Our approach: A randomized simulation . . . . .	10
4.3	Bounding the distance from a product distribution . . . . .	11
4.4	Rounding to the nearest oblivious corruption . . . . .	12
<b>5</b>	<b>Preliminaries</b>	<b>12</b>
<b>6</b>	<b>Bounding the distance from a product distribution: Proof of Lemma 4.2</b>	<b>15</b>
6.1	Proof of Lemma 4.3 . . . . .	16
6.2	Bounding the mutual information for low-degree cost functions . . . . .	18
6.3	Proof of Lemma 4.2 . . . . .	19
<b>7</b>	<b>Rounding to a legal corruption: Proof of Lemma 4.4</b>	<b>20</b>
7.1	Proof of Lemma 7.1 . . . . .	21
7.2	Generalizing to when the adversary can corrupt: Proof of Claim 7.3 . . . . .	23
7.3	Proof of Lemma 4.4 . . . . .	24
<b>8</b>	<b>Adaptive adversaries are at least as strong as oblivious adversaries</b>	<b>26</b>
<b>9</b>	<b>Putting the pieces together: Proof of Theorem 3</b>	<b>29</b>
<b>10</b>	<b>Lower bounds</b>	<b>30</b>
10.1	Proof of Theorem 4 . . . . .	30
10.2	Proof of Theorem 5 . . . . .	32
<b>11</b>	<b>Acknowledgments</b>	<b>34</b>
<b>A</b>	<b>The subtractive and additive adversaries in our framework</b>	<b>37</b>
A.1	Subtractive adversaries . . . . .	37
A.2	Additive adversaries . . . . .	39
<b>B</b>	<b>Partially-adaptive adversaries</b>	<b>42</b>
<b>C</b>	<b>Brief overview of [BLMT22]’s approaches and their limitations</b>	<b>45</b>
C.1	The special case of additive adversaries . . . . .	45
C.2	The special case of statistical query algorithms . . . . .	45

# 1 Introduction

Classic models of data analysis assume that data is drawn independently from the distribution of interest, but the real world is rarely so kind. To be robust to the messiness of real-world data, we desire algorithms that succeed even in the presence of an adversary that corrupts the data. Such adversaries were first introduced in the seminal works of [Tuk60, Hub64, Ham71] and have since been the subject of intense study in a variety of settings [Val85, Hau92, KL93, KSS94, BEK02, DSFT<sup>+</sup>14, LRV16, CSV17, DKK<sup>+</sup>19, DKPZ21, HSSVG22, BLMT22, CHL<sup>+</sup>23, DK23].

By now, there are numerous models for how the adversary can corrupt the data, including additive, subtractive, “strong”/“nasty”, agnostic, and adaptive and non-adaptive variants of each of these. In many cases, the provable guarantees for our algorithms are only known for a subset of these models. We refer the interested reader to the excellent recent textbook [DK23] for a more complete background and survey of recent results and open directions. Our work focuses on a surprisingly under-explored question:

*What is the relationship between the various statistical adversaries?*

Specifically, we compare *adaptive* adversaries, which can look at the sample before deciding on a corruption, and *oblivious* adversaries, which must commit to their corruptions before the i.i.d. sample is drawn.

**Theorem 1** (Informal, see [Theorem 2](#) for the formal version). *Adaptive adversaries and their oblivious counterparts are equivalent up to scaling the sample size by a factor polynomial in the original sample size and polylogarithmic in the domain size.*

[Theorem 1](#) resolves the main question introduced by [BLMT22] and further explored in [CHL<sup>+</sup>23]. We defer its formal statement to [Section 2](#), but, for now, mention two points. First, it is a generic result that proves the equivalence between many distinct adaptive adversaries and their oblivious counterparts (e.g. the equivalence between “subtractive adaptive” and “subtractive oblivious” adversaries). Second, it is constructive. We give a simple transformation, the subsampling filter described in [Definition 5](#), which takes any algorithm that succeeds on a statistical task in the presence of the oblivious adversary and converts it to one that succeeds on the same task in the presence of the adaptive adversary. This transformation preserves the statistical and computational efficiency of the original algorithm up to polynomial factors.

In addition to answering a foundational question about the relative power of statistical adversaries, [Theorem 1](#) has several practical implications:

1. Given the many distinct definitions of robustness, it can be difficult for a practitioner to determine which definition is most appropriate for their setting and therefore which algorithm to utilize. [Theorem 1](#) partially alleviates this issue by greatly reducing the number of truly unique adversary models.
2. It shows that a single algorithmic idea, that of subsampling, *amplifies* robustness in many different models. Formally, it takes an algorithm that is only robust to the oblivious adversary and converts it to one robust to the adaptive counterpart. This suggests that, even if the practitioner cannot precisely determine the most appropriate model of robustness, they should try subsampling.

3. **Theorem 1** can be reformulated as an answer to an equivalent and independently interesting question: How useful is it to hide one’s dataset from the adversary? It shows that private data does *not* afford much more robustness than public data.

## 2 Our Results

Before formally describing our main result, we define a unified framework in which to express and analyze statistical adversaries. It may be instructive to view this framework with the following concrete adaptive adversary and its oblivious counterpart in mind:

**Example 1.** Consider the following adaptive and oblivious adversaries parameterized by  $\eta \in [0, 1]$ :

- Adaptive: When the algorithm requests  $n$  points, first an i.i.d. sample  $\mathbf{S} \sim \mathcal{D}^n$  is drawn. Then, may alter up to  $\lfloor \eta \cdot n \rfloor$  of them arbitrarily. The algorithm receives this corrupted sample.
- Oblivious: The adversary can choose any  $\mathcal{D}'$  that has a total variation distance to  $\mathcal{D}$  of at most  $\eta$ , and the algorithm receives  $n$  i.i.d. draws from  $\mathcal{D}'$ .

These two adversaries are well-studied, and are referred to by different names. The adaptive adversary is typically referred to as “strong contamination” in the statistical estimation literature [DK23] and “nasty noise” in the PAC learning literature [BEK02]. The oblivious adversary has been referred to as “general, non-adaptive, contamination” [DK23]. In our unified framework, these adversaries will be defined via the same “cost function,” and as a result, we prove them equivalent.

### 2.1 A unified framework to define statistical adversaries

Each adversary will be parameterized by a “cost” function  $\rho$  where  $\rho(x, y)$  specifies the cost the adversary pays to corrupt  $x$  to  $y$ , with a cost of  $\infty$  indicating that the adversary is not allowed to change  $x$  to  $y$ . The adversary can choose any corruptions subject to a budget constraint on the total cost incurred. This cost function is required to have two basic properties.

**Definition 1** (Cost function). *A function  $\rho : X \times X \rightarrow \mathbb{R}_{\geq 0} \cup \{\infty\}$  is said to be a “cost function” if it satisfies the following properties.*

1. For any  $x \in X$ ,  $\rho(x, x) = 0$ .
2. For any  $x, y \in X$ ,  $\rho(x, y) \geq 0$ .

The adversary is specified by both the cost function and whether it is adaptive or oblivious. Given the cost function,  $\rho$ , the corresponding adaptive adversary is defined as follows:

**Definition 2** (Adaptive adversary, corruptions to the sample). *For any cost function  $\rho$  and  $S \in X^n$ , we use  $\mathcal{C}_\rho(S)$  to denote all  $S' \in X^n$  for which*

$$\frac{1}{n} \sum_{i \in [n]} \rho(S_i, S'_i) \leq 1.$$

*The  $\rho$ -adaptive adversary is allowed to corrupt the clean sample  $S$  to any  $S' \in \mathcal{C}_\rho(S)$ . For any  $f : X^n \rightarrow \{0, 1\}$  and distribution  $\mathcal{D}$ , the max success probability of  $f$  in the presence of the  $\rho$ -adaptive adversary is denoted:*

$$\text{Adaptive-Max}_\rho(f, \mathcal{D}) := \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^n} \left[ \sup_{\mathbf{S}' \in \mathcal{C}_\rho(\mathbf{S})} \{f(\mathbf{S}')\} \right].$$

In the case of the adversaries of Example 1, their cost functions are simply  $\rho(x, y) = 1/\eta$  for all  $x \neq y$ . In that case, the budget constraint in the above definition ensures that, for this choice of cost function,  $\rho$ , the  $\rho$ -adaptive adversary can corrupt at most an  $\eta$  fraction of points in the sample, corresponding to the standard definition of the “strong contamination”/“nasty noise” models.

Given a cost function, the associated oblivious adversary replaces the budget constraint of the adaptive setting with a natural distributional analog. It is easy to see that the following definition of  $\rho$ -oblivious adversaries is equivalent to the “general, non-adaptive, contamination” model of Example 1 when the cost function is defined as  $\rho(x, y) = 1/\eta$  for all  $x \neq y$ .

**Definition 3** (Oblivious adversary, corruptions to a distribution). *For any cost function  $\rho$  and distribution  $\mathcal{D}$ , we overload  $\mathcal{C}_\rho(\mathcal{D})$  to refer to the set of all distributions  $\mathcal{D}'$  for which there exists a coupling of  $\mathbf{x} \sim \mathcal{D}$  and  $\mathbf{x}' \sim \mathcal{D}'$  satisfying*

$$\mathbb{E}[\rho(\mathbf{x}, \mathbf{x}')] \leq 1.$$

The  $\rho$ -oblivious adversary is allowed to corrupt the base distribution  $\mathcal{D}$  to any  $\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})$ . For any  $f : X^n \rightarrow \{0, 1\}$  and distribution  $\mathcal{D}$ , the max success probability of  $f$  in the presence of the  $\rho$ -oblivious adversary is denoted:

$$\text{Oblivious-Max}_\rho(f, \mathcal{D}) := \sup_{\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})} \left\{ \mathbb{E}_{\mathbf{S}' \sim (\mathcal{D}')^n} [f(\mathbf{S}')] \right\}.$$

In Definition 3, since  $\mathbf{x} \sim \mathcal{D}$  and  $\mathbf{x}' \sim \mathcal{D}'$  can be coupled so that the average cost to corrupt  $\mathbf{x}$  to  $\mathbf{x}'$  is at most 1, we can similarly couple  $\mathbf{S} \sim \mathcal{D}^n$  and  $\mathbf{S}' \sim (\mathcal{D}')^n$  so that the average cost to corrupt each point in  $\mathbf{S}$  to the corresponding point in  $\mathbf{S}'$  is at most 1. From this perspective, the crucial difference between the oblivious and adaptive adversary is that the oblivious adversary must commit to how it corrupts each  $\mathbf{x}$  without knowing the contents of the sample, whereas the adaptive adversary gets to view  $\mathbf{S}$  before deciding.

We show, in Section 3, that our framework can express many commonly studied statistical adversaries, including subtractive contamination, additive contamination, and agnostic noise.

**Remark 1** (Partially-adaptive statistical adversaries). Some statistical adversaries lie between their fully adaptive and fully oblivious counterparts. These include malicious noise [Val85] and the non-iid oblivious adversary defined in [CHL<sup>+</sup>23]. Our results readily extend to such adversaries (see Section 3.1 for details).

## 2.2 Our main result: Adaptive and oblivious adversaries are equivalent

Our main result is that for any algorithm  $A$  and cost function  $\rho$ , there exists an algorithm  $A'$  inheriting the efficiency of  $A$  for which the performance of  $A$  in the presence of the oblivious adversary is equivalent to the performance of  $A'$  in the presence of the adaptive adversary.

**Definition 4** ( $\varepsilon$ -equivalent algorithms). *For any algorithms  $A : X^n \rightarrow Y$  and  $A' : X^m \rightarrow Y$ , we say that  $A$  in the presence of the  $\rho$ -oblivious adversary is  $\varepsilon$ -equivalent to  $A'$  in the presence of the  $\rho$ -adaptive adversary if for any test function  $T : Y \rightarrow \{0, 1\}$  and distribution  $\mathcal{D}$  supported on  $X$ ,*

$$|\text{Oblivious-Max}_\rho(T \circ A, \mathcal{D}) - \text{Adaptive-Max}_\rho(T \circ A', \mathcal{D})| \leq \varepsilon.$$

Colloquially,  $A$  and  $A'$  are  $\varepsilon$ -equivalent if no test can distinguish their outputs with more than  $\varepsilon$  probability. Note that while the above definition is about the *maximum* acceptance probability of  $T$ , it also applies to the test  $\bar{T} := 1 - T$  and therefore the *minimum* acceptance probability of  $T$  also must be approximately the same for  $A$  and  $A'$ .

The algorithm  $A'$  will run  $A$  on a uniformly random subsample of its input.

**Definition 5** (Subsampling filter). *For any  $m \geq n$  we define the subsampling filter  $\Phi_{m \rightarrow n} : X^m \rightarrow X^n$  as the (randomized) algorithm that given  $S \in X^m$ , returns a sample of  $n$  points drawn uniformly without replacement from  $S$ .*

**Theorem 2** (Subsampling neutralizes the adaptivity in statistical adversaries). *For any algorithm  $A : X^n \rightarrow Y$ ,  $\varepsilon > 0$ , and cost function  $\rho$ , let  $m = \text{poly}(n, \ln |X|, 1/\varepsilon)$  and  $A' := A \circ \Phi_{m \rightarrow n}$ . Then,  $A$  in the presence of the  $\rho$ -oblivious adversary is  $\varepsilon$ -equivalent to  $A'$  in the presence of the  $\rho$ -adaptive adversary.*

For constant  $\varepsilon$ , **Theorem 2** says that if there is an algorithm  $A$  solving a statistical task with an oblivious adversary taking as input  $n \cdot \log |X|$  bits, there is an algorithm  $A'$  solving the same task with an adaptive adversary taking only polynomially more bits as input. Furthermore, if  $A$  is computationally efficient, then  $A'$  is too.

**Remark 2** (Continuous domains). In many statistical problems, the domain is  $\mathbb{R}^d$ . To apply **Theorem 2** to an algorithm  $A$  over continuous domains, we first discretize that domain to some  $X := \text{disc}(\mathbb{R})^d$  where the discretization depends on  $A$ . If  $A$  requires  $b$  bits of precision in each dimension, then  $\log_2 |X| = bd$ , which is typically polynomial in  $n$ . For example, under the mild assumption that  $A$  accesses the bits of each dimension sequentially, both  $b$  and  $d$  are upper bounded by the time complexity of  $A$ . In this setting, if the time complexity of  $A$  is polynomial in  $n$ , then so is  $\ln |X|$ .

**Theorem 2** is a special case of our main theorem in which the  $|X|$  is replaced with the *degree* of the cost function, a measure of the number of corruptions the adversary can make for each input.

**Definition 6** (Degree of a cost function). *For any cost function  $\rho : X \times X \rightarrow \mathbb{R}_{\geq 0} \cup \{\infty\}$ , the degree of  $\rho$  is defined as*

$$\text{deg}(\rho) := \sup_{x \in X} \{ \text{The number of distinct } y \in X \text{ for which } \rho(x, y) \neq \infty \}.$$

**Theorem 3** (Main result, generalization of **Theorem 2**). *For any algorithm  $A : X^n \rightarrow Y$ ,  $\varepsilon > 0$ , and cost function  $\rho$  with degree  $d \geq 2$ , let  $m = O\left(\frac{n^4 (\ln d)^2}{\varepsilon^4}\right)$  and  $A' := A \circ \Phi_{m \rightarrow n}$ . Then,  $A$  in the presence of the  $\rho$ -oblivious adversary is  $\varepsilon$ -equivalent to  $A'$  in the presence of the  $\rho$ -adaptive adversary.*

The degree is constant for many natural cost functions, such as the cost function corresponding to subtractive contamination. In these cases, **Theorem 3** has no dependence on the domain size.

### 2.3 Lower bounds

**Theorem 3** requires a polynomial increase of the sample size of  $A'$  relative to  $A$ . It is natural to wonder whether such an increase is necessary. The results of [CHL+23] show it is.

**Fact 2.1** ([CHL<sup>+</sup>23]). *For any  $n \in \mathbb{N}$ , the task of Gaussian mean testing with appropriate parameters (depending on  $n$ ) can be solved using  $n$  samples in the presence of the oblivious additive adversary, but requires  $\tilde{\Omega}(n^{4/3})$  in the presence of the adaptive additive adversary.*

In the setting of [Theorem 3](#), one difficulty in interpreting [Fact 2.1](#) is that the cost function corresponding to the additive adversary has a large degree<sup>1</sup>, and so it’s unclear if the increased sample size is due a dependence on the degree or an innate required polynomial increase.

For example, the subtractive adversary (formally defined in [Section 3](#)) has a degree of 2 because, for each point in the sample, it chooses between keeping that point or removing it. For this adversary, is a polynomial increase in sample size necessary? Our first lower bound gives a straightforward proof this is the case, even for a simple task.

**Theorem 4** (A polynomial increase in sample size is necessary). *Let  $\mathcal{D}$  be a distribution on  $X = [m] := \{1, \dots, m\}$  that is promised to be uniform on some  $X' \subseteq X$ . Then,*

1. *There is an algorithm that estimates  $|X'|$  using  $n := \tilde{O}(\sqrt{m})$  samples even with oblivious subtractive contamination.*
2. *Any algorithm that estimates  $|X'|$  to the same accuracy with adaptive subtractive contamination requires  $\tilde{\Omega}(m)$  samples.*

[Theorem 4](#) implies that in the statement of [Theorem 2](#), we must take  $m$  polynomial larger than  $n$ . Next, we show that this  $m$  must also depend polylogarithmically on a degree-like characteristic of the cost function.

**Definition 7** (Budget-bounded degree). *For any cost function  $\rho : X \times X \rightarrow \mathbb{R}_{\geq 0} \cup \{\infty\}$  and  $b \in \mathbb{R}_{\geq 0}$ , the  $b$ -bounded degree of  $\rho$  is defined as*

$$\text{deg}_b(\rho) := \sup_{x \in X} \{ \text{The number of distinct } y \in X \text{ for which } \rho(x, y) \leq b \}.$$

This lower bound will make one assumption on  $\rho$ : that  $\rho(x, y) \geq 1 + \delta$  whenever  $x \neq y$  for a small constant  $\delta$ . This corresponds to the adversary having a budget on how many points they can change and is satisfied by all of the well-studied models discussed in [Section 3](#).

**Theorem 5** (Dependence on  $\ln \text{deg}_b(\rho)$  is necessary). *For any constants  $b, \delta > 0$ , large enough  $n \in \mathbb{N}$ , and cost function  $\rho$  for which  $\rho(x, y) \geq 1 + \delta$  whenever  $x \neq y$ , there is an algorithm  $A : X^n \rightarrow \{0, 1\}$  for which the following holds. If  $A$  in the presence of the  $\rho$ -oblivious adversary is  $(\varepsilon = 0.9)$ -equivalent to  $A' := A \circ \Phi_{m \rightarrow n}$  in the presence of the  $\rho$ -adaptive adversary, then*

$$m \geq \tilde{\Omega}_{b, \delta}(n \cdot \ln \text{deg}_b(\rho)).$$

Comparing [Theorems 3](#) and [5](#), for “reasonable” cost functions in which  $\text{deg}(\rho) \approx \text{deg}_{1000}(\rho)$  and  $\rho(x, y) \geq 1.001$  for all  $x \neq y$ , a domain-size independent result is possible precisely when the degree does not grow with  $|X|$ .

---

<sup>1</sup>Since the domain for the task is over the continuous domain of  $\mathbb{R}^d$ , technically this cost function has infinite degree. However, as we discussed in [Remark 2](#), it makes more sense to think of the degree as  $\approx 2^d$  in this setting, which happens to be exponential in the  $n$  of [Fact 2.1](#).

## 2.4 Relation to recent work

Recent work of Blanc, Lange, Malik, and Tan initiated a formal study of the relationship between adaptive adversaries and their oblivious counterparts [BLMT22]. They conjectured the equivalence of adaptive and oblivious statistical adversaries but only proved it in two special cases.

1. They showed that *additive* oblivious and *additive* adaptive adversaries are equivalent. Our result, which applies to all statistical adversaries, requires an entirely different approach. This is because, in some sense, the adaptive additive adversary is *less adaptive* than other adaptive adversaries. We elaborate on this point in [Appendix C.1](#) and explain why [BLMT22]’s approach does not generalize to all adversaries.
2. They also showed that if a *statistical query* (SQ) algorithm is robust to an oblivious adversary, it can be upgraded to be robust to the corresponding adaptive adversary. The restriction to SQ algorithms greatly facilitates [BLMT22]’s analysis because we have a much better understanding of SQ algorithms than general algorithms. For example, the quality of the best SQ algorithm for a given task is captured by simple combinatorial measures [BFJ<sup>+</sup>94, Fel17]. In [Appendix C.2](#), we further describe [BLMT22]’s SQ result and give advantages of our result even for algorithms that can be cast in the SQ framework.

Other recent work of Canonne, Hopkins, Li, Liu, and Narayanan tackled the equivalence of statistical adversaries from the other direction [CHL<sup>+</sup>23]. While we aim to show that distinct statistical adversaries are equivalent, they showed a separation: For the well-studied problem of Gaussian mean testing, an adaptive adversary requires polynomial more samples than the corresponding oblivious adversary (see [Fact 2.1](#)).

## 3 Instantiating common adversaries in our framework

Here, we show how to express many common statistical adversaries within our framework. For completeness, we include the “strong contamination/nasty noise” adversary of [Example 1](#).

**Strong contamination/nasty noise:** As mentioned in [Example 1](#), both the “strong contamination/nasty noise” adversary, that can arbitrarily replace an  $\eta$  fraction of an i.i.d. sample, and the “general, non-adaptive, contamination” adversary, that can perturb the underlying distribution from which the sample is drawn by at most  $\eta$  in total variation distance, correspond to adaptive and oblivious adversaries with the following cost function:

$$\rho_{\text{strong}}(x, y) := \begin{cases} 0 & \text{if } x = y \\ \frac{1}{\eta} & \text{otherwise.} \end{cases}$$

**Agnostic learning** [[Hau92](#), [KSS94](#)]: Agnostic noise is a well-studied adversary [[KKMS08](#), [KK09](#), [Fel10](#), [DSFT<sup>+</sup>14](#), [DKPZ21](#)] specific to supervised learning problems, where each point in the sample is a pair  $(x, y)$  of the input and its label. This adversary is allowed to change  $\eta$  fraction of the labels but must keep the inputs unchanged. It corresponds to the cost function,

$$\rho_{\text{agnostic}}((x_1, y_1), (x_2, y_2)) := \begin{cases} 0 & \text{if } x_1 = x_2 \text{ and } y_1 = y_2 \\ \frac{1}{\eta} & \text{if } x_1 = x_2 \text{ and } y_1 \neq y_2 \\ \infty & \text{if } x_1 \neq x_2. \end{cases}$$

Agnostic learning typically refers to the  $\rho$ -oblivious adversary. It can be equivalently defined as the learner receiving an i.i.d. sample of points of the form  $(\mathbf{x}, g(\mathbf{x}))$  where  $g$  is close to the original target  $f$  in the sense that

$$\Pr_{\mathbf{x}}[f(\mathbf{x}) \neq g(\mathbf{x})] \leq \eta.$$

In the adaptive variant, first, a sample is drawn that is labeled by the true target function. Then, an adversary may corrupt  $\eta$ -fraction of the labels arbitrarily. This variant is sometimes referred to as *nasty classification noise* [BEK02].

Note that the cost function  $\rho_{\text{agnostic}}$  only has degree 2 in binary classification settings. **Theorem 3** therefore shows the equivalence between nasty classification noise and agnostic noise with no dependence on the domain size.

**Subtractive contamination:** In the adaptive variant of subtractive contamination, the adversary is allowed to remove  $\lfloor \eta n \rfloor$  points from a size- $n$  sample. The algorithm receives the remaining  $n - \lfloor \eta n \rfloor$  points.

In the oblivious variant [DK23], the algorithm receives i.i.d. samples from the distribution  $\mathcal{D}$  conditioned on some event  $E$  that occurs with probability  $1 - \eta$ . This can be thought of as the adversary removing  $\eta$ -fraction of the distribution corresponding to when the event  $E$  does not occur.

To fit subtractive contamination into our framework, we will augment the domain with a special element  $\emptyset$ , to indicate the adversary has removed this point. For the augmented domain  $X' := X \cup \{\emptyset\}$ , it uses the cost function,

$$\rho_{\text{sub}}(x, y) := \begin{cases} 0 & \text{if } x = y \\ \frac{1}{\eta} & \text{if } x \neq y \text{ and } y = \emptyset \\ \infty & \text{otherwise.} \end{cases}$$

Note that once again, this cost function has degree only 2, so by **Theorem 3**, the  $\rho$ -adaptive and  $\rho$ -oblivious adversaries are equivalent with no dependence on the domain size. In **Appendix A**, we give an easy reduction from the standard notions of subtractive noise (without the  $\emptyset$  element added to the domain) to the adversaries defined by  $\rho_{\text{sub}}$ . This reduction, combined with **Theorem 3**, shows that the standard oblivious and adaptive subtractive adversaries are equivalent.

**Additive contamination (Huber’s model [Hub64]):** In Huber’s original model [Hub64], rather than directly receive i.i.d. samples from the target distribution  $\mathcal{D}$ , the algorithm receives i.i.d. samples from  $\mathcal{D}'$ , the mixture distribution

$$\mathcal{D}' := (1 - \eta)\mathcal{D} + \eta\mathcal{E}$$

where the adversary chooses the outlier distribution  $\mathcal{E}$ . In the adaptive variant of this model, first a clean sample of  $\lfloor (1 - \eta)n \rfloor$  points are drawn i.i.d. from  $\mathcal{D}$ . Then, the adversary may add  $\lfloor \eta n \rfloor$  points arbitrarily. These  $n$  points are then randomly permuted so that the algorithm cannot trivially identify which points were added.

Similarly to subtractive contamination, we will use the augmented domain  $X' := X \cup \{\emptyset\}$ . For additive noise, we use the cost function

$$\rho_{\text{add}}(x, y) := \begin{cases} 0 & \text{if } x = y \\ \frac{1}{\eta} & \text{if } x \neq y \text{ and } x = \emptyset \\ \infty & \text{otherwise.} \end{cases}$$

In [Appendix A](#), we give an easy reduction showing how [Theorem 3](#) gives the equivalence between Huber’s contamination model and its adaptive variant. Note that this particular equivalence was already proven by [\[BLMT22\]](#), but for completeness, we show their result can be recovered using our framework.

### 3.1 Partially-adaptive adversaries

As alluded to in [Remark 1](#), some adversaries lie between the fully oblivious and fully adaptive adversaries. Our results show that such intermediate adversaries are equivalent to their fully oblivious and fully adaptive counterparts. The strategy for proving this equivalence is by showing the intermediate adversary is at least as strong as the oblivious adversary, and that it is no stronger than the adaptive adversary. Since [Theorem 2](#) implies the adaptive adversary is no more powerful than the oblivious adversary, we can conclude that all three adversaries are equivalent. We formalize this approach for the two adversaries described here in [Appendix B](#), showing both are equivalent to additive contamination.

**Malicious noise:** This model was first defined by [\[Val85\]](#). In it, the  $n$  samples are generated sequentially. For each point, independently with probability  $1 - \eta$ , that point is sampled from  $\mathcal{D}$ . Otherwise, the adversary chooses an arbitrary corrupted point with full knowledge of previous points generated but no knowledge of future points. Intuitively, this adversary is partially adaptive because when the adversary chooses how to corrupt a point, it has partial knowledge of the sample corresponding to the points generated previously.

**The non-independent additive adversary:** This model was recently studied in [\[CHL+23\]](#). In it, the adversary generates  $\lfloor \eta n \rfloor$  arbitrary points. The sample is then formed by combining  $\lfloor (1 - \eta)n \rfloor$  points drawn i.i.d. from  $\mathcal{D}$  with the adversary’s chosen points. Intuitively, this adversary is partially adaptive because the  $\eta n$  points it generates need not be i.i.d. from some distribution as they would for fully oblivious adversaries, but the adversary still does not know the sample when choosing corruptions as in a fully adaptive adversary.

## 4 Technical overview

To prove [Theorem 3](#), we begin with the observation that we can combine the algorithm  $A$  and test  $T$  into a single function  $f := T \circ A$ . Therefore, it suffices to prove the following.

**Theorem 6** ([Theorem 3](#) restated). *For any  $n, d \in \mathbb{N}$  where  $d \geq 2$ , domain  $X$ , and  $\varepsilon > 0$ , let  $m = O\left(\frac{n^4(\ln d)^2}{\varepsilon^4}\right)$ . Then, for any  $f : X^n \rightarrow \{0, 1\}$ , cost function  $\rho$  with degree  $d$ , and distribution  $\mathcal{D}$  supported on  $X$ ,*

$$|\text{Oblivious-Max}_\rho(f, \mathcal{D}) - \text{Adaptive-Max}_\rho(f \circ \Phi_{m \rightarrow n}, \mathcal{D})| \leq \varepsilon. \tag{1}$$

[Theorem 6](#) can be understood as a statement about the *indistinguishability* of the following two families of distributions, both over datasets in  $X^n$ .

1. The set of input distributions over  $n$  points the oblivious adversary can create,

$$\mathcal{D}_{\text{oblivious}}^{n, \rho, \mathcal{D}} := \{(\mathcal{D}')^n \mid \mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})\}.$$

2. For the adaptive adversary, we first define the set of input distributions over  $m$  points before subsampling: We say  $\mathcal{D}_{\text{adaptive}}$  is a valid adaptive corruption, denoted  $\mathcal{D}_{\text{adaptive}} \in \mathcal{D}_{\text{adaptive}}^{m,\rho,\mathcal{D}}$  if it is possible to couple  $\mathbf{S}' \sim \mathcal{D}_{\text{adaptive}}$  and a clean sample  $\mathbf{S} \sim \mathcal{D}^m$  so that  $\mathbf{S}' \in \mathcal{C}_\rho(\mathbf{S})$  with probability 1. Then, the distribution on  $n$  points is created via subsampling,

$$\Phi_{m \rightarrow n}(\mathcal{D}_{\text{adaptive}}^{m,\rho,\mathcal{D}}) := \left\{ \text{The distribution of } \Phi_{m \rightarrow n}(\mathbf{S}') \mid \mathbf{S}' \sim \mathcal{D}_{\text{adaptive}} \text{ for any } \mathcal{D}_{\text{adaptive}} \in \mathcal{D}_{\text{adaptive}}^{m,\rho,\mathcal{D}} \right\}.$$

With these definitions, [Theorem 6](#) can be recast as the following two statements:

1. **The oblivious adversary is no harder than the adaptive adversary:** For any distinguisher  $f : X^n \rightarrow \{0, 1\}$  and oblivious corruption  $\mathcal{D}_{\text{oblivious}} \in \mathcal{D}_{\text{oblivious}}^{n,\rho,\mathcal{D}}$ , there is a subsampled adaptive corruption  $\mathcal{D}_{\text{adaptive}} \in \Phi_{m \rightarrow n}(\mathcal{D}_{\text{adaptive}}^{m,\rho,\mathcal{D}})$  satisfying

$$\left| \mathbb{E}_{\mathbf{S} \sim \mathcal{D}_{\text{oblivious}}} [f(\mathbf{S})] - \mathbb{E}_{\mathbf{S} \sim \mathcal{D}_{\text{adaptive}}} [f(\mathbf{S})] \right| \leq \varepsilon.$$

This is the easy half of [Theorem 6](#) and is already known for some specific adversary models [[DKK<sup>+</sup>19](#), [ZJS19](#)]. We defer discussion of its proof to [Section 8](#).

2. **The adaptive adversary is no harder than the oblivious adversary:** For any distinguisher  $f : X^n \rightarrow \{0, 1\}$  and subsampled adaptive corruption  $\mathcal{D}_{\text{adaptive}} \in \Phi_{m \rightarrow n}(\mathcal{D}_{\text{adaptive}}^{m,\rho,\mathcal{D}})$ , there is an oblivious corruption  $\mathcal{D}_{\text{oblivious}} \in \mathcal{D}_{\text{oblivious}}^{n,\rho,\mathcal{D}}$  satisfying

$$\left| \mathbb{E}_{\mathbf{S} \sim \mathcal{D}_{\text{adaptive}}} [f(\mathbf{S})] - \mathbb{E}_{\mathbf{S} \sim \mathcal{D}_{\text{oblivious}}} [f(\mathbf{S})] \right| \leq \varepsilon. \quad (2)$$

The remainder of this overview is devoted to our proof of this harder half of [Theorem 6](#)

#### 4.1 A first attempt and why it fails

A natural approach towards proving this harder half of [Theorem 6](#) is to show that every adaptive adversary can be simulated by an oblivious adversary. This corresponds to switching the order of quantifiers in the desired statement: The goal of this approach is to show that for any adaptive corruption  $\mathcal{D}_{\text{adaptive}} \in \Phi_{m \rightarrow n}(\mathcal{D}_{\text{adaptive}}^{m,\rho,\mathcal{D}})$ , there is a single choice of oblivious corruption  $\mathcal{D}_{\text{oblivious}} \in \mathcal{D}_{\text{oblivious}}^{n,\rho,\mathcal{D}}$  satisfying

$$d_{\text{TV}}(\mathcal{D}_{\text{adaptive}}, \mathcal{D}_{\text{oblivious}}) \leq \varepsilon.$$

Indeed, as we discuss in [Section 8](#), such an approach works for the easier half of [Theorem 6](#). Here, we will explain why this approach fails for the harder half and latter use this counterexample to motivate our ultimately successful approach.

Our construction of this counterexample uses the adaptive and oblivious adversaries described in [Example 1](#) with a budget  $\eta = 1/2$ . Recall this means that,

1. For any  $S \in X^m$ , the adaptive adversary can change an arbitrary  $m/2$  points within  $S$ .
2. For any distribution  $\mathcal{D}$ , the oblivious adversary can choose any  $\mathcal{D}'$  with a total variation distance of at most  $1/2$  from  $\mathcal{D}$ .

Furthermore, we use perhaps the simplest possible base distribution,  $\mathcal{D} = \text{Unif}(\{0, 1\})$  and our counterexample works for *any*  $m \geq n := 2$ .

After receiving  $\mathbf{S} \sim \mathcal{D}^m$ , the adaptive adversary can choose a corruption so that  $\mathbf{S}'$  either contains only zeros or only ones, with both cases equally likely. They achieve this by flipping all the 0s or all the 1s in  $\mathbf{S}$ , whichever is less frequent (breaking ties uniformly). The result of this approach is that there is some  $\mathcal{D}_{\text{adaptive}} \in \mathcal{D}_{\text{adaptive}}$  for which

$$\mathcal{D}_{\text{adaptive}} = \text{Unif}([0, 0], [1, 1]).$$

The above distribution is far from any product distribution and, as a result, far from any possible  $\mathcal{D}_{\text{oblivious}}$ . Therefore, this naive simulation approach fails.

## 4.2 Our approach: A randomized simulation

A key observation about this counter example: Even though  $\mathcal{D}_{\text{adaptive}}$  is far from any single  $\mathcal{D}_{\text{oblivious}}$ , it is exactly equal to a mixture of oblivious adversaries. This is because the oblivious adversary can create the point-mass distribution that always outputs  $[0, 0]$  and that which always outputs  $[1, 1]$ . Our main lemma is that such a randomized simulation is always possible.

**Lemma 4.1** (With subsampling, the adaptive adversary can be simulated by a randomized oblivious adversary). *For any base distribution  $\mathcal{D}$ , sample size  $n$ , error parameter  $\varepsilon$ , and cost function  $\rho$  with degree  $d$ , set  $m = O\left(\frac{n^4(\ln d)^2}{\varepsilon^4}\right)$ . Then, for any subsampled adaptive corruption*

$\mathcal{D}_{\text{adaptive}} \in \Phi_{m \rightarrow n}\left(\mathcal{D}_{\text{adaptive}}^{m, \rho, \mathcal{D}}\right)$ , *there is a randomized oblivious corruption  $\mathcal{D}_{\text{oblivious}}$  supported on  $\mathcal{D}_{\text{oblivious}}^{n, \rho, \mathcal{D}}$  such that its mixture satisfies,*

$$d_{\text{TV}}(\mathcal{D}_{\text{adaptive}}, \mathbb{E}[\mathcal{D}_{\text{oblivious}}]) \leq \varepsilon.$$

At a high level, our approach to proving [Lemma 4.1](#) is to group the adaptively corrupted samples  $\mathbf{S} \sim \mathcal{D}_{\text{adaptive}}$  into a moderate number of groups that make “similar” corruptions. The goal of this grouping is that, if we look at the distribution  $\mathbf{S}$  conditioned on falling within a single group, the resulting distribution is close to a single oblivious adversary  $\mathcal{D}_{\text{oblivious}} \in \mathcal{D}_{\text{oblivious}}^{n, \rho, \mathcal{D}}$ .

Our method of determining how to group a  $\mathbf{S} \sim \mathcal{D}_{\text{adaptive}} \in \mathcal{D}_{\text{adaptive}}^{m, \rho, \mathcal{D}}$ . We specify each group by a “core”  $c \in X^k$  and group together all  $\mathbf{S}$  containing this core. Note this is a soft grouping in the sense that a single  $\mathbf{S}$  will be a member of many groups as it contains many distinct cores. This grouping strategy is summarized in the following definition.

**Definition 8** (Our grouping strategy). *For any distribution  $\mathcal{D}_{\text{adaptive}}$  over samples in  $X^m$  and parameters  $n + k \leq m$ , let  $\text{Grouped}_{n, k}(\mathcal{D}_{\text{adaptive}})$  be the joint distribution over  $(\mathbf{S}, \mathbf{c})$  formed by*

1. Drawing a  $\mathbf{S}_{\text{big}} \sim \mathcal{D}_{\text{adaptive}}$ .
2. Drawing  $\mathbf{S}$  and  $\mathbf{c}$  to be uniform size- $n$  and size- $k$  respectively disjoint subsamples of  $\mathbf{S}_{\text{big}}$ , meaning  $\mathbf{S} \sim \Phi_{m \rightarrow n}(\mathbf{S}_{\text{big}})$  and  $\mathbf{c} \sim \Phi_{(m-n) \rightarrow k}(\mathbf{S}_{\text{big}} \setminus \mathbf{S})$ .

*It will also be helpful to define  $\text{Group}_n(\mathcal{D}_{\text{adaptive}}, c)$  to be the distribution of  $\mathbf{S}$  conditioned on  $\mathbf{c} = c$ .*

To prove [Lemma 4.1](#), we will show that in expectation over the core  $\mathbf{c}$ , there is some  $\mathcal{D}_{\text{oblivious}}$  that is within  $\varepsilon$ -total variation distance to the distribution of  $\text{Group}_n(\mathcal{D}_{\text{adaptive}}, \mathbf{c})$ . We analyze this total variation distance in two steps which are described in the following two subsections.

### 4.3 Bounding the distance from a product distribution

Since every distribution  $\mathcal{D}_{\text{oblivious}} \in \mathcal{D}_{\text{oblivious}}^{n,\rho,\mathcal{D}}$  is a product distribution, the first step in our analysis is to show that the average group is close to some product distribution (though not necessarily one that is a valid oblivious corruption).

**Lemma 4.2** (The average group is close to a product distribution). *For any  $\mathcal{D}_{\text{adaptive}} \in \mathcal{D}_{\text{adaptive}}^{m,\rho,\mathcal{D}}$  where  $\rho$  is a degree- $d$  cost function, and any  $n + k_{\max} \leq m/2$ , there exists some  $k \leq k_{\max}$  for which*

$$\mathbb{E}_{\mathbf{S}, \mathbf{c} \sim \text{Grouped}_{n,k}(\mathcal{D}_{\text{adaptive}})} [d_{\text{TV}}(\text{Group}_n(\mathcal{D}_{\text{adaptive}}, \mathbf{c}), \mathcal{D}_{\text{goal}}(\mathbf{c})^n)] \leq \sqrt{\frac{n^2 \ln d}{2k_{\max}}}.$$

where  $\mathcal{D}_{\text{goal}}(\mathbf{c}) := \mathbb{E}_{\mathbf{S} \sim \text{Group}_n(\mathcal{D}_{\text{adaptive}}, \mathbf{c})}[\text{Unif}(\mathbf{S})]$  is the average data-point in the group corresponding to the core  $\mathbf{c}$ .

**Lemma 4.2** is closely related to the *correlation rounding* technique used to round semidefinite programs (see e.g. [BRS11, RT12]), also called the *pinning lemma* in the statistical physics community (see e.g. [Eld20]). Roughly speaking, these results say the following: For any (not necessarily independent) random variables  $\mathbf{x}_1, \dots, \mathbf{x}_m$ , by conditioning on a few of the  $\mathbf{x}_i$ , we can make the covariances between the remaining pairs close to independent. The version of this result we will need extends the concept of covariance from pairs to larger groups:

**Definition 9** (Multivariate total correlation). *The multivariate total correlation of random variables  $\mathbf{x}_1, \dots, \mathbf{x}_n$  is*

$$\text{Cor}(\mathbf{x}_1, \dots, \mathbf{x}_n) := d_{\text{KL}}(\mathcal{D} \parallel \mathcal{D}_1 \times \dots \times \mathcal{D}_n)$$

where  $\mathcal{D}$  is the distribution of  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$  and  $\mathcal{D}_i$  is the marginal distribution of  $\mathbf{x}_i$ . Similarly, for any  $\mathbf{y}$ , we define the conditional multivariate correlation as

$$\text{Cor}(\mathbf{x}_1, \dots, \mathbf{x}_n \mid \mathbf{y}) = \mathbb{E}_{\mathbf{y}}[\text{Cor}(\mathbf{x}_1 \mid \mathbf{y}, \dots, \mathbf{x}_n \mid \mathbf{y})].$$

We use the following to prove **Lemma 4.2**, where  $I(\mathbf{x}; \mathbf{y})$  denotes the mutual information between  $\mathbf{x}$  and  $\mathbf{y}$ .

**Lemma 4.3** (Correlation rounding). *For any random variable on  $\mathbf{S}$  on  $X^m$  and integers  $n + k_{\max} \leq m$ , there exists some  $k \leq k_{\max}$  for which,*

$$\mathbb{E}_{\substack{\mathbf{A} \sim \binom{[m]}{n} \\ \mathbf{B} \sim \binom{[m] \setminus \mathbf{A}}{k}}} [\text{Cor}(\mathbf{S}_{\mathbf{A}} \mid \mathbf{S}_{\mathbf{B}})] \leq \frac{n(n-1)}{2(k_{\max}+1)} \cdot \mathbb{E}_{\substack{\mathbf{i} \sim \text{Unif}(\binom{[m]}{n}) \\ \mathbf{B} \sim \binom{[m] \setminus \{\mathbf{i}\}}{n+k_{\max}-1}}} [I(\mathbf{S}_{\mathbf{i}}; \mathbf{S}_{\mathbf{B}})].$$

As far as we are aware, the closest results to **Lemma 4.3** already appearing in the literature are those of [MR17, JKR19]. They prove essentially the same result except the term  $\mathbb{E}[I(\mathbf{S}_{\mathbf{i}}; \mathbf{S}_{\mathbf{B}})]$  is replaced with  $\mathbb{E}[H(\mathbf{S}_{\mathbf{i}})]$ . Since entropy upper bounds mutual information, our result is always at least as strong as theirs.

The distinction between  $\mathbb{E}[I(\mathbf{S}_{\mathbf{i}}; \mathbf{S}_{\mathbf{B}})]$  and  $\mathbb{E}[H(\mathbf{S}_{\mathbf{i}})]$  ends up being crucial for our application. In **Lemma 6.5**, we are able to show that if  $\mathbf{S}$  is the result of an adaptive corruption with a degree- $d$  cost function, that  $\mathbb{E}[I(\mathbf{S}_{\mathbf{i}}; \mathbf{S}_{\mathbf{B}})] \leq O(\ln d)$ . In contrast, even when the adversary makes no corruptions, if the base (uncorrupted) distribution has high entropy, then  $\mathbb{E}[H(\mathbf{S}_{\mathbf{i}})]$  can be as large as  $\ln |X|$ . We furthermore view our version as having a natural interpretation: The term  $\mathbb{E}[I(\mathbf{S}_{\mathbf{i}}; \mathbf{S}_{\mathbf{B}})]$  measures some notion of how far  $\mathbf{S}_1, \dots, \mathbf{S}_m$  are from being independent. Thus, our result gives a correlation bound that improves when the starting distribution was already close to a product distribution.

## 4.4 Rounding to the nearest oblivious corruption

While [Lemma 4.2](#) guarantees that most groups are close to some product distribution, we require that product distribution to be a valid oblivious corruption. We show that this product distribution, on average, can be rounded to a nearby valid oblivious corruption.

**Lemma 4.4** (Error due to rounding for our grouping strategy). *Using the same notation as [Lemma 4.2](#), as long as  $k \leq m/2$*

$$\mathbb{E}_{\mathbf{c}} \left[ \inf_{\mathcal{D}_{\text{oblivious}} \in \mathcal{D}_{\text{oblivious}}^{n, \rho, \mathcal{D}}} \{d_{\text{TV}}(\mathcal{D}_{\text{oblivious}}, \mathcal{D}_{\text{goal}}(\mathbf{c})^n)\} \right] \leq 2n \cdot \sqrt{\frac{k \ln d}{m}}.$$

To prove [Lemma 4.4](#), we first show a more general but weaker result that holds for *any* grouping strategy. This result, formalized in [Lemma 7.1](#), gives a bound that scales with the number of groups. Since the number of groups we use is  $|X|^k$ , a direct application of [Lemma 7.1](#) would require the parameter  $m$  to scale with the domain size. Nonetheless, we show a more delicate application of [Lemma 7.1](#) allows for us to get a bound roughly as good as if there were only  $d^k$  groups.

**Setting parameters.** We briefly sketch how to recover [Lemma 4.1](#) using [Lemmas 4.2](#) and [4.4](#). For any adaptive strategy  $\mathcal{D}_{\text{adaptive}}$ , [Lemma 4.2](#) gives there is core size  $k \leq k_{\text{max}} = O(n^2 \ln(d)/\varepsilon^2)$  for which

$$d_{\text{TV}}(\mathcal{D}_{\text{adaptive}}, \mathbb{E}_{\mathbf{c}}[\mathcal{D}_{\text{goal}}(\mathbf{c})^n]) \leq \varepsilon/2.$$

We will then apply [Lemma 4.4](#). Setting  $m = O(n^2 k (\ln d)/\varepsilon^2) = O(n^4 (\ln d)^2/\varepsilon^4)$ , we have that for each core  $c$ , there is some  $\mathcal{D}_{\text{oblivious}}(c) \in \mathcal{D}_{\text{oblivious}}^{n, \rho, \mathcal{D}}$  for which

$$\mathbb{E}_{\mathbf{c}}[d_{\text{TV}}(\mathcal{D}_{\text{oblivious}}(\mathbf{c}), \mathcal{D}_{\text{goal}}(\mathbf{c})^n)] \leq \varepsilon/2.$$

Combining the above two bounds recovers [Lemma 4.1](#).

## 5 Preliminaries

**Indexing.** For any  $n \in \mathbb{N}$ , we use  $[n]$  as shorthand for  $\{1, 2, \dots, n\}$ . Similarly, for  $n \leq m \in \mathbb{N}$ , we use  $[n, m]$  as shorthand for  $\{n, n+1, \dots, m\}$ . For any multiset  $S \in X^m$ , we use  $S_i$  to denote the  $i^{\text{th}}$  element of  $S$ . For any  $I \subseteq [m]$ , we use  $S_I$  to denote the multiset containing  $(S_{I_1}, S_{I_2}, \dots)$ . We'll also use  $S_{<j}$  and  $S_{-j}$  as shorthand for  $S_{[j-1]}$  and  $S_{[m] \setminus \{j\}}$  respectively. For any permutation  $\sigma : [m] \rightarrow [m]$  and  $S \in X^m$ , we'll use  $\sigma(S)$  as shorthand for the multiset in  $X^m$  satisfying  $\sigma(S)_i = S_{\sigma(i)}$ .

**Random variables and distributions.** We use **boldfont** to denote random variables and calligraphic font to denote distributions (e.g.  $\mathbf{x} \sim \mathcal{D}$ ). For a multiset  $S$ , we use  $\text{Unif}(S)$  to denote the uniform distribution over elements of  $S$ . For a distribution  $\mathcal{D}$ , we will use  $\mathbf{x}_1, \dots, \mathbf{x}_n \stackrel{\text{iid}}{\sim} \mathcal{D}$  and  $\mathbf{x} \sim \mathcal{D}^n$  interchangeably to denote that  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are independent and identically distributed according to  $\mathcal{D}$ . For any distributions  $\mathcal{D}_1, \mathcal{D}_2$ , we use  $\mathcal{D}_1 \times \mathcal{D}_2$  to denote the product distribution of  $\mathcal{D}_1$  and  $\mathcal{D}_2$ . We denote mixture distributions as convex combinations (e.g.  $\mathcal{D}_{\text{mix}} = 1/3 \cdot \mathcal{D}_1 + 2/3 \cdot \mathcal{D}_2$ ).

We use the following standard concentration inequality.

**Fact 5.1** (Chernoff bound). *Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be independent random variables on  $\{0, 1\}$ , and  $\mathbf{X}$  their sum. For  $\mu := \mathbb{E}[\mathbf{X}]$ ,*

$$\Pr[\mathbf{X} \geq 2\mu] \leq e^{-\mu/3} \quad \text{and} \quad \Pr[\mathbf{X} \leq \mu/2] \leq e^{-\mu/8}.$$

We will also use two commonly studied families of random variables. For any  $p \in [0, 1]$ , we use  $\text{Ber}(p)$  to denote the distribution that takes on value 1 with probability  $p$  and takes on 0 otherwise. Furthermore, for any  $n \in \mathbb{N}$ , we use  $\text{Bin}(n, p)$  to denote the sum of  $n$  independent random variables each distributed according to  $\text{Ber}(p)$ .

**Formalizing the corruption models.** We recap the notation used to formalize our corruption models. Beginning with the adaptive for any cost function  $\rho : X \times X \rightarrow \mathbb{R}_{\geq 0} \cup \infty$  and sample  $S \in X^m$ , we use  $\mathcal{C}_\rho(S)$  to denote legal adaptive corruptions of  $S$  under cost function  $\rho$ ,

$$\mathcal{C}_\rho(S) := \left\{ S' \in X^m \mid \frac{1}{m} \sum_{i \in [m]} \rho(S_i, S'_i) \leq 1 \right\}.$$

For a base distribution  $\mathcal{D}$ , the set of input distributions on size- $m$  data sets the adaptive adversary can create is denoted:

$$\mathcal{D}_{\text{adaptive}}^{m, \rho, \mathcal{D}} := \{\text{Distributions } \mathcal{D}' \text{ over } X^m \mid \text{Can couple } \mathbf{S}' \sim \mathcal{D}' \text{ and } \mathbf{S} \sim \mathcal{D}^m \text{ so } \mathbf{S}' \in \mathcal{C}_\rho(\mathbf{S}) \text{ w.p. } 1\}.$$

For the oblivious adversary, we overload  $\mathcal{C}_\rho(\mathcal{D})$  to denote all distributions the oblivious adversary can create,

$$\mathcal{C}_\rho(\mathcal{D}) := \{\text{Distributions } \mathcal{D}' \text{ over } X \mid \text{Can couple } \mathbf{x}' \sim \mathcal{D}' \text{ and } \mathbf{x} \sim \mathcal{D} \text{ so } \mathbb{E}[\rho(\mathbf{x}, \mathbf{x}') \leq 1]\}.$$

The set of input distributions on size- $n$  datasets the oblivious adversary can create is denoted:

$$\mathcal{D}_{\text{oblivious}}^{n, \rho, \mathcal{D}} := \{(\mathcal{D}')^n \mid \mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})\}.$$

**Subsampling filter.** Recall, in [Definition 5](#), we defined  $\Phi_{m \rightarrow n} : X^m \rightarrow X^n$  to be the (randomized) algorithm that given  $S \in X^m$  returns a sample of  $n$  points drawn uniformly without replacement from  $S$ . We will generalize this in three ways: First, we'll use  $\Phi_{\star \rightarrow n}$  the filter that takes in a sample  $S \in X^\star$  of at least  $n$  points and then subsamples it down to  $n$  points. Second, if  $\mathcal{D}$  is a distribution over  $S \in X^m$ , we'll use  $\Phi_{m \rightarrow n}(\mathcal{D})$  to denote the distribution of  $\Phi_{m \rightarrow n}(\mathbf{S})$ . Lastly, if  $\mathcal{D}$  is a family of distributions, we'll use  $\Phi_{m \rightarrow n}(\mathcal{D})$  to denote  $\{\Phi_{m \rightarrow n}(\mathcal{D}) \mid \mathcal{D} \in \mathcal{D}\}$ .

**TV distance and KL divergence.** We use two measures of statistical distance/divergence.

**Definition 10** (Total variation distance). *Let  $\mathcal{D}$  and  $\mathcal{D}'$  be any two distributions over the same domain  $X$ . The total variation distance between  $\mathcal{D}$  and  $\mathcal{D}'$ , is defined as*

$$d_{\text{TV}}(\mathcal{D}, \mathcal{D}') := \sup_{T: X \rightarrow [0, 1]} \left\{ \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[T(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \mathcal{D}'}[T(\mathbf{x})] \right\}.$$

*This quantity can be equivalently defined as the infimum over all couplings of  $\mathbf{x} \sim \mathcal{D}$  and  $\mathbf{x}' \sim \mathcal{D}'$  of  $\Pr[\mathbf{x} \neq \mathbf{x}']$ .*

TV distance is a true distance in the sense that it satisfies the triangle inequality.

**Fact 5.2** (The triangle inequality for TV distance). *For any distributions  $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3$ ,*

$$d_{\text{TV}}(\mathcal{D}_1, \mathcal{D}_3) \leq d_{\text{TV}}(\mathcal{D}_1, \mathcal{D}_2) + d_{\text{TV}}(\mathcal{D}_2, \mathcal{D}_3).$$

We can also give a (sometimes coarse) upper bound on the TV distance of product distribution.

**Fact 5.3** (Total variation distance of a product). *For any distributions  $\mathcal{D}_1, \mathcal{D}_2$  and  $n \in \mathbb{N}$ ,*

$$d_{\text{TV}}(\mathcal{D}_1^n, \mathcal{D}_2^n) \leq n \cdot d_{\text{TV}}(\mathcal{D}_1, \mathcal{D}_2).$$

Straight from the definition, we see that TV distance is convex in one argument.

**Fact 5.4** (Convexity of TV distance). *For any distributions  $\mathcal{D}_1, \mathcal{D}_2, \mathcal{E}_1$ , and  $\mathcal{E}_2$ , and mixture weight  $\lambda \in [0, 1]$ ,*

$$d_{\text{TV}}(\mathcal{D}_\lambda, \mathcal{E}_\lambda) \leq \lambda \cdot d_{\text{TV}}(\mathcal{D}_1, \mathcal{E}_1) + (1 - \lambda) \cdot d_{\text{TV}}(\mathcal{D}_2, \mathcal{E}_2),$$

where  $\mathcal{D}_\lambda$  is the mixture  $\lambda\mathcal{D}_1 + (1 - \lambda)\mathcal{D}_2$  and similarly  $\mathcal{E}_\lambda := \lambda\mathcal{E}_1 + (1 - \lambda)\mathcal{E}_2$ .

The other measure of statistical distance/divergence that plays a key role in our results in KL divergence.

**Definition 11** (Kullback-Leibler (KL) Divergence). *For distributions  $\mathcal{D}, \mathcal{E}$  supported on the same domain  $X$ , the KL divergence between  $\mathcal{D}$  and  $\mathcal{E}$  as defined as,*

$$d_{\text{KL}}(\mathcal{D} \parallel \mathcal{E}) := \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \ln \left( \frac{\mathcal{D}(\mathbf{x})}{\mathcal{E}(\mathbf{x})} \right) \right],$$

where  $\mathcal{D}(x)$  and  $\mathcal{E}(x)$  denote the probability mass or density functions of  $\mathcal{D}$  and  $\mathcal{E}$  respectively at the point  $x$  (or more generally,  $\mathcal{D}(x)/\mathcal{E}(x)$  is the Radon-Nikodym derivative of  $\mathcal{D}$  with respect to  $\mathcal{E}$ ).

Unlike TV distance, KL divergence is not a true distance in the sense that it does not satisfy triangle inequality. For us, it will suffice that it upper bounds TV distance via Pinsker's inequality [Pin64]. [Can22] has a nice summary of different forms and proofs of the below inequality and appropriate references for each.

**Fact 5.5** (Pinsker's inequality [Pin64, Can22]). *For any distributions  $\mathcal{D}, \mathcal{E}$ ,*

$$d_{\text{TV}}(\mathcal{D}, \mathcal{E}) \leq \sqrt{\frac{d_{\text{KL}}(\mathcal{D} \parallel \mathcal{E})}{2}}.$$

## Mutual information.

**Definition 12** (Mutual information). *For random variables  $\mathbf{x}, \mathbf{y}$  jointly distribution according to a distribution  $\mathcal{D}$ , let  $\mathcal{D}_x$  and  $\mathcal{D}_y$  be the marginal distributions of  $\mathbf{x}$  and  $\mathbf{y}$  respectively, and  $\mathcal{D}_{x|y}$  be the marginal distribution of  $\mathbf{x}$  conditioned on  $\mathbf{y} = y$ . The mutual information between  $\mathbf{x}$  and  $\mathbf{y}$  is defined as*

$$I(\mathbf{x}; \mathbf{y}) = d_{\text{KL}}(\mathcal{D} \parallel \mathcal{D}_x \times \mathcal{D}_y) = \mathbb{E}_{\mathbf{y}} [d_{\text{KL}}(\mathcal{D}_{x|\mathbf{y}} \parallel \mathcal{D}_x)]$$

**Definition 13** (Conditional mutual information). *For random variables  $\mathbf{x}, \mathbf{y}, \mathbf{z}$  jointly distributed, the mutual information of  $\mathbf{x}$  and  $\mathbf{y}$  conditioned on  $\mathbf{z}$  is*

$$I(\mathbf{x}; \mathbf{y} \mid \mathbf{z}) := \mathbb{E}_{\mathbf{z}' \sim \mathcal{D}_{\mathbf{z}}} [I((\mathbf{x} \mid \mathbf{z} = \mathbf{z}'); (\mathbf{y} \mid \mathbf{z} = \mathbf{z}'))]$$

where  $\mathcal{D}_{\mathbf{z}}$  is the marginal distribution of  $\mathbf{z}$ .

The *chain rule* connects mutual information and conditional mutual information.

**Fact 5.6** (Chain rule for mutual information). *For any  $\mathbf{x}, \mathbf{y}, \mathbf{z}$ ,*

$$I(\mathbf{x}; (\mathbf{y}, \mathbf{z})) = I(\mathbf{x}; \mathbf{z}) + I(\mathbf{x}; \mathbf{y} \mid \mathbf{z}).$$

*This is sometimes rewritten as*

$$I(\mathbf{x}; \mathbf{y} \mid \mathbf{z}) = I(\mathbf{x}; (\mathbf{y}, \mathbf{z})) - I(\mathbf{x}; \mathbf{z}).$$

Mutual information is always nonnegative

**Fact 5.7** (Nonnegativity of mutual information). *For any random variables  $\mathbf{x}, \mathbf{y}$ ,*

$$I(\mathbf{x}; \mathbf{y}) \geq 0.$$

As an easy consequence of the chain rule and mutual information being nonnegative, we have that mutual information can only increase if we consider more information.

**Fact 5.8.** *For any random variables  $\mathbf{x}, \mathbf{y}, \mathbf{z}$ ,*

$$I(\mathbf{x}; (\mathbf{y}, \mathbf{z})) \geq I(\mathbf{x}; \mathbf{y}).$$

Mutual information is also symmetric.

**Fact 5.9** (Symmetry of mutual information). *For any random variables  $\mathbf{x}, \mathbf{y}$ ,*

$$I(\mathbf{x}; \mathbf{y}) = I(\mathbf{y}; \mathbf{x}).$$

Another nice property of mutual information is that it is bounded by the support size of each variable.

**Fact 5.10** (Mutual information with a finite support). *For any random variables  $\mathbf{x}, \mathbf{y}$ , if one of  $\mathbf{x}$  or  $\mathbf{y}$  has a finite support of size  $d$ , then,*

$$I(\mathbf{x}; \mathbf{y}) \leq \ln d.$$

## 6 Bounding the distance from a product distribution: Proof of Lemma 4.2

In this section, we prove the following, restated for convenience.

**Lemma 4.2** (The average group is close to a product distribution). *For any  $\mathcal{D}_{\text{adaptive}} \in \mathcal{D}_{\text{adaptive}}^{m, \rho, \mathcal{D}}$  where  $\rho$  is a degree- $d$  cost function, and any  $n + k_{\max} \leq m/2$ , there exists some  $k \leq k_{\max}$  for which*

$$\mathbb{E}_{\mathbf{S}, \mathbf{c} \sim \text{Grouped}_{n, k}(\mathcal{D}_{\text{adaptive}})} [d_{\text{TV}}(\text{Group}_n(\mathcal{D}_{\text{adaptive}}, \mathbf{c}), \mathcal{D}_{\text{goal}}(\mathbf{c})^n)] \leq \sqrt{\frac{n^2 \ln d}{2k_{\max}}}.$$

where  $\mathcal{D}_{\text{goal}}(\mathbf{c}) := \mathbb{E}_{\mathbf{S} \sim \text{Group}_n(\mathcal{D}_{\text{adaptive}}, \mathbf{c})}[\text{Unif}(\mathbf{S})]$  is the average data-point in the group corresponding to the core  $\mathbf{c}$ .

We structure this proof into three steps:

1. In [Section 6.1](#), we prove our version of correlation rounding. As discussed in [Section 4.3](#), this is improvement of similar results from [\[MR17, JKR19\]](#) that replaces an entropy term with a mutual information term.
2. In [Section 6.2](#), we prove that this mutual information term is small when the adaptive adversary has low degree.
3. In [Section 6.3](#), we combine these two results to prove [Lemma 4.2](#)

## 6.1 Proof of [Lemma 4.3](#)

In this proof, we will often be reasoning about mutual information ([Definition 12](#)) and multivariate total correlation ([Definition 9](#)) of subsets of the random variable  $\mathbf{S}$  supported on  $X^m$ . It will be convenient to have the following concise notation: For any  $a + b + c \leq m$ , we define,

$$\begin{aligned} \text{Cor}_{\mathbf{S}}(a) &:= \mathbb{E}_{\mathbf{A} \sim \binom{[m]}{a}} [\text{Cor}(\mathbf{S}_{\mathbf{A}})], \\ \text{Cor}_{\mathbf{S}}(a \mid b) &:= \mathbb{E}_{\mathbf{A} \sim \binom{[m]}{a}, \mathbf{B} \sim \binom{[m] \setminus \mathbf{A}}{b}} [\text{Cor}(\mathbf{S}_{\mathbf{A}} \mid \mathbf{S}_{\mathbf{B}})], \\ I_{\mathbf{S}}(a; b) &:= \mathbb{E}_{\mathbf{A} \sim \binom{[m]}{a}, \mathbf{B} \sim \binom{[m] \setminus \mathbf{A}}{b}} [I(\mathbf{S}_{\mathbf{A}}; \mathbf{S}_{\mathbf{B}})], \\ I_{\mathbf{S}}(a; b \mid c) &:= \mathbb{E}_{\mathbf{A} \sim \binom{[m]}{a}, \mathbf{B} \sim \binom{[m] \setminus \mathbf{A}}{b}, \mathbf{C} \sim \binom{[m] \setminus (\mathbf{A} \cup \mathbf{B})}{c}} [I(\mathbf{S}_{\mathbf{A}}; \mathbf{S}_{\mathbf{B}} \mid \mathbf{S}_{\mathbf{C}})]. \end{aligned}$$

With this notation, we can succinctly restated [Lemma 4.3](#).

**Lemma 6.1** (Restatement of [Lemma 4.3](#)). *For any random variable on  $\mathbf{S}$  on  $X^m$  and integers  $n + k_{\max} \leq m$ , there exists some  $k \leq k_{\max}$  for which,*

$$\text{Cor}_{\mathbf{S}}(n \mid k) \leq \frac{n(n-1)}{2(k_{\max} + 1)} \cdot I_{\mathbf{S}}(1; n + k_{\max} - 1).$$

We assemble the ingredients used in the proof of [Lemma 6.1](#). First, is a simple application of the chain rule.

**Proposition 6.2.** *For any  $a + b + c \leq m$  and  $\mathbf{S}$  on  $X^m$ ,*

$$I_{\mathbf{S}}(a; b \mid c) = I_{\mathbf{S}}(a; b + c) - I_{\mathbf{S}}(a; c).$$

*Proof.* For any disjoint  $A, B, C \subseteq [m]$ , we apply [Fact 5.6](#) which gives that

$$I(\mathbf{S}_A; \mathbf{S}_B \mid \mathbf{S}_C) = I(\mathbf{S}_A; \mathbf{S}_{B \cup C}) - I(\mathbf{S}_A; \mathbf{S}_C).$$

Averaging over  $\mathbf{A}, \mathbf{B}$ , and  $\mathbf{C}$  gives the desired result.  $\square$

Second, we show the following.

**Proposition 6.3.** *For any random variable  $\mathbf{S}$  supported on  $X^m$  and  $a + b \leq m$ ,*

$$I_{\mathbf{S}}(a; b) \leq a \cdot I_{\mathbf{S}}(1, a + b - 1)$$

*Proof.* It suffices to show, for any random variables  $\mathbf{x}_1, \dots, \mathbf{x}_a$  and  $\mathbf{y}$ , that

$$I(\mathbf{x}_1, \dots, \mathbf{x}_a; \mathbf{y}) \leq \sum_{i \in [a]} I(\mathbf{x}_i; \mathbf{y}, \mathbf{x}_{\neq i}),$$

as the desired result then follows by averaging over all  $\mathbf{A}, \mathbf{B}$  and setting  $\mathbf{x} = \mathbf{S}_A$  and  $\mathbf{y} = \mathbf{S}_B$ . We bound,

$$\begin{aligned} I(\mathbf{x}_1, \dots, \mathbf{x}_a; \mathbf{y}) &= \sum_{i \in [a]} I(\mathbf{x}_i; \mathbf{y} \mid \mathbf{x}_{<i}) && \text{(Fact 5.6)} \\ &= \sum_{i \in [a]} I(\mathbf{x}_i; \mathbf{y}, \mathbf{x}_{<i}) - I(\mathbf{x}_i; \mathbf{x}_{<i}) && \text{(Fact 5.6 again)} \\ &\leq \sum_{i \in [a]} I(\mathbf{x}_i; \mathbf{y}, \mathbf{x}_{<i}) && \text{(Fact 5.7)} \\ &\leq \sum_{i \in [a]} I(\mathbf{x}_i; \mathbf{y}, \mathbf{x}_{<i}, \mathbf{x}_{>i}), && \text{(Fact 5.8)} \end{aligned}$$

which is exactly the desired bound.  $\square$

Third, we give an alternative form of multivariate total correlation.

**Proposition 6.4** (Multivariate total correlation in terms of mutual information). *For any random variables  $\mathbf{x}_1, \dots, \mathbf{x}_n$ ,*

$$\text{Cor}(\mathbf{x}_1, \dots, \mathbf{x}_n) = \sum_{i \in [n-1]} I(\mathbf{x}_{\leq i}; \mathbf{x}_{i+1}).$$

*Proof.* Throughout this proof, we use  $\mathcal{D}$  to denote the distribution of  $\mathbf{x}$ ,  $\mathcal{D}_i$  to denote the marginal distribution of  $\mathbf{x}_i$ , and  $\mathcal{D}_{\leq i}$  to denote the marginal distribution of  $\mathbf{x}_{\leq i}$ . Expanding the right-hand side,

$$\begin{aligned} \sum_{i \in [n-1]} I(\mathbf{x}_{\leq i}; \mathbf{x}_{i+1}) &= \sum_{i \in [n-1]} d_{\text{KL}}(\mathcal{D}_{\leq i+1} \parallel \mathcal{D}_{\leq i} \times \mathcal{D}_{i+1}) && \text{(Definition of mutual information)} \\ &= \sum_{i \in [n-1]} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\leq i+1}} \left[ \ln \left( \frac{\mathcal{D}_{\leq i+1}(\mathbf{x})}{\mathcal{D}_{\leq i}(\mathbf{x}_{\leq i}) \mathcal{D}_{i+1}(\mathbf{x}_{i+1})} \right) \right] && \text{(Definition of KL divergence)} \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \ln \left( \prod_{i=1}^{n-1} \frac{\mathcal{D}_{\leq i+1}(\mathbf{x}_{\leq i+1})}{\mathcal{D}_{\leq i}(\mathbf{x}_{\leq i}) \mathcal{D}_{i+1}(\mathbf{x}_{i+1})} \right) \right] && \text{(Linearity of expectation)} \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \ln \left( \frac{\mathcal{D}(\mathbf{x})}{\prod_{i=1}^n \mathcal{D}_i(\mathbf{x}_i)} \right) \right] && \text{(Cancellation of terms)} \\ &= d_{\text{KL}}(\mathcal{D} \parallel \mathcal{D}_1 \times \mathcal{D}_2 \times \dots \times \mathcal{D}_n) \end{aligned}$$

which is exactly  $\text{Cor}(\mathbf{x})$ .  $\square$

We are now ready to prove the main result of this subsection.

*Proof of Lemma 6.1.* We wish to show there is some  $k \leq k_{\max}$  for which  $\text{Cor}_{\mathcal{S}}(n \mid k)$  is small. For any such  $k$ , we have that

$$\begin{aligned} \text{Cor}_{\mathcal{S}}(n \mid k) &= \sum_{i \in [n-1]} I_{\mathcal{S}}(i; 1 \mid k) && \text{(Proposition 6.4)} \\ &= \sum_{i \in [n-1]} I_{\mathcal{S}}(i; k+1) - I_{\mathcal{S}}(i; k) && \text{(Proposition 6.2.)} \end{aligned}$$

Summing up the above for all  $k = 0, \dots, k_{\max}$ , we obtain

$$\begin{aligned} \sum_{k \in [0, k_{\max}]} \text{Cor}_{\mathcal{S}}(n \mid k) &= \sum_{k \in [0, k_{\max}]} \sum_{i \in [n-1]} I_{\mathcal{S}}(i; k+1) - I_{\mathcal{S}}(i; k) \\ &= \sum_{i \in [n-1]} I_{\mathcal{S}}(i, k_{\max} + 1) \quad (\text{Cancel telescoping terms and } I_{\mathcal{S}}(i, 0) = 0) \\ &\leq \sum_{i \in [n-1]} i \cdot I_{\mathcal{S}}(1, i + k_{\max}) && \text{(Proposition 6.3)} \\ &\leq \sum_{i \in [n-1]} i \cdot I_{\mathcal{S}}(1, n-1 + k_{\max}) && \text{(Fact 5.8)} \\ &= \frac{n(n-1)}{2} \cdot I_{\mathcal{S}}(1, n-1 + k_{\max}) \end{aligned}$$

Therefore, using the fact that minimum over all  $k \in [0, k_{\max}]$  is at most the mean, there exists one choice of  $k$  for which

$$\text{Cor}_{\mathcal{S}}(n \mid k) \leq \frac{n(n-1)}{2(k_{\max} + 1)} \cdot I_{\mathcal{S}}(1, n-1 + k_{\max}). \quad \square$$

## 6.2 Bounding the mutual information for low-degree cost functions

The main result of this subsection is the following.

**Lemma 6.5** (Bounding mutual information for low-degree corruptions). *For any  $\mathcal{S}' \sim \mathcal{D}_{\text{adaptive}} \in \mathcal{D}_{\text{adaptive}}^{m, \rho, \mathcal{D}}$  where  $\rho$  is a degree- $d$  cost function and  $r < m$ ,*

$$\mathbb{E}_{i \sim \text{Unif}([m]), \mathcal{B} \sim \binom{[m] \setminus \{i\}}{r}} [I(\mathcal{S}'_i; \mathcal{S}'_{\mathcal{B}})] \leq \frac{m}{m-r} \cdot \ln d.$$

We will use the following.

**Proposition 6.6.** *Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be independent random variables and  $\mathbf{y}$  be any (not necessarily independent) random variable. Then,*

$$\sum_{i \in [n]} I(\mathbf{x}_i; \mathbf{y}) \leq I((\mathbf{x}_1, \dots, \mathbf{x}_n); \mathbf{y}).$$

*Proof.* We bound,

$$\begin{aligned}
I((\mathbf{x}_1, \dots, \mathbf{x}_n); \mathbf{y}) &= \sum_{i \in [n]} I(\mathbf{x}_i; \mathbf{y} \mid \mathbf{x}_{<i}) && \text{(Fact 5.6)} \\
&= \sum_{i \in [n]} I(\mathbf{x}_i; (\mathbf{y}, \mathbf{x}_{<i})) - I(\mathbf{x}_i; \mathbf{x}_{<i}) && \text{(Fact 5.6 again)} \\
&= \sum_{i \in [n]} I(\mathbf{x}_i; (\mathbf{y}, \mathbf{x}_{<i})) && (\mathbf{x}_i \text{ and } \mathbf{x}_{<i} \text{ are independent)} \\
&\geq \sum_{i \in [n]} I(\mathbf{x}_i; \mathbf{y}) && \text{(Fact 5.8)}
\end{aligned}$$

□

*Proof of Lemma 6.5.* Since  $\mathcal{D}_{\text{adaptive}}$  is a legal adaptive corruption, there is a coupling of  $\mathbf{S} \sim \mathcal{D}^m$  and  $\mathbf{S}'$  for which  $\mathbf{S}' \in \mathcal{C}_\rho(\mathbf{S})$  with probability 1. In particular, this implies that that once we condition on  $\mathbf{S}_i$ , there are at most  $d$  choices for  $\mathbf{S}'_i$ .

For any fixed choice of  $\mathbf{i} = i, \mathbf{B} = B$ , we have that

$$\begin{aligned}
I(\mathbf{S}'_i; \mathbf{S}'_B) &\leq I((\mathbf{S}_i, \mathbf{S}'_i); (\mathbf{S}_B, \mathbf{S}'_B)) && \text{(Fact 5.8)} \\
&= I(\mathbf{S}_i; (\mathbf{S}_B, \mathbf{S}'_B)) + I(\mathbf{S}'_i; (\mathbf{S}_B, \mathbf{S}'_B) \mid \mathbf{S}_i) && \text{(Fact 5.6)} \\
&= I(\mathbf{S}_i; \mathbf{S}_B) + I(\mathbf{S}_i; \mathbf{S}'_B \mid \mathbf{S}_B) + I(\mathbf{S}'_i; (\mathbf{S}_B, \mathbf{S}'_B) \mid \mathbf{S}_i) && \text{(Fact 5.6 again.)}
\end{aligned}$$

The first term,  $I(\mathbf{S}_i; \mathbf{S}_B)$ , is zero because  $\mathbf{S}_i$  and  $\mathbf{S}_B$  are independent. The third term,  $I(\mathbf{S}'_i; (\mathbf{S}_B, \mathbf{S}'_B) \mid \mathbf{S}_i)$ , is at most  $\ln d$  by Fact 5.10 and the fact that conditioned on  $\mathbf{S}_i$  there are only  $d$  possible values for  $\mathbf{S}'_i$ . For the remaining term, we bound it in expectation over  $\mathbf{i}$ ,

$$\begin{aligned}
\mathbb{E}_{\mathbf{i} \sim \text{Unif}([m] \setminus B)} [I(\mathbf{S}_i; \mathbf{S}'_B \mid \mathbf{S}_B)] &= \frac{1}{m-r} \cdot \sum_{i \in ([m] \setminus B)} I(\mathbf{S}_i; \mathbf{S}'_B \mid \mathbf{S}_B) \\
&\leq \frac{1}{m-r} I(\mathbf{S}_{[m] \setminus B}; \mathbf{S}'_B \mid \mathbf{S}_B) && \text{(Proposition 6.6)} \\
&\leq \frac{r \ln d}{m-r}. && \text{(Fact 5.10 and } \mathbf{S}'_B \text{ has } \leq d^r \text{ options given } \mathbf{S}_B)
\end{aligned}$$

Combining these bounds, we have that

$$\mathbb{E}_{\mathbf{i} \sim \text{Unif}([m]), \mathbf{B} \sim \binom{[m] \setminus \{\mathbf{i}\}}{r}} [I(\mathbf{S}'_i; \mathbf{S}'_B)] \leq \ln d + \frac{r \ln d}{m-r} = \frac{m \ln d}{m-r}. \quad \square$$

### 6.3 Proof of Lemma 4.2

*Proof.* Since  $\mathcal{D}_{\text{adaptive}} \in \mathcal{D}_{\text{adaptive}}^{m, \rho, \mathcal{D}}$ , it is possible to couple  $\mathbf{S} \sim \mathcal{D}^m$  and  $\mathbf{S}' \sim \mathcal{D}_{\text{adaptive}}$  so that  $\mathbf{S}' \in \mathcal{C}_\rho(\mathbf{S})$  with probability 1. Next, draw a uniform permutation  $\sigma : [m] \rightarrow [m]$  and define the random variables  $\mathbf{T}$  and  $\mathbf{T}'$  each over  $X^m$  as

$$\mathbf{T}_i = \mathbf{S}_{\sigma(i)} \quad \text{and} \quad \mathbf{T}'_i = \mathbf{S}'_{\sigma(i)}.$$

Note that the marginal distribution of  $\mathbf{T}$  is still  $\mathcal{D}^m$ , and it still holds that  $\mathbf{T}' \in \mathcal{C}_\rho(\mathbf{T})$  with probability 1. Therefore, the distribution of  $\mathbf{T}'$  is still in  $\mathcal{D}_{\text{adaptive}}^{m,\rho,\mathcal{D}}$ .

The upshot is now the distribution of all permutations of  $\mathbf{T}'$  are identical. Therefore, setting  $\mathbf{c} = \mathbf{T}'_{[\leq k]}$  and  $\mathbf{R} := \mathbf{T}'_{[k+1, k+n]}$  we have that the distribution of  $(\mathbf{R}, \mathbf{c})$  is exactly  $\text{Grouped}_{n,k}(\mathcal{D}_{\text{adaptive}})$ . Therefore,

$$\mathcal{D}_{\text{goal}}(\mathbf{c}) = \mathbb{E}[\text{Unif}(\mathbf{R}) \mid \mathbf{c} = \mathbf{c}].$$

Reusing permutation invariance of  $\mathbf{T}$  (and therefore  $\mathbf{R}$ ), we have that the distribution of  $\mathbf{R}_i$  is the same for all  $i \in [n]$  (even after conditioning on  $\mathbf{c} = \mathbf{c}$ ), and therefore all equal to  $\mathcal{D}_{\text{goal}}(\mathbf{c})$ . As a result, we have that

$$\mathbb{E}_{\mathbf{c}}[d_{\text{KL}}(\text{Group}_n(\mathcal{D}_{\text{adaptive}}, \mathbf{c}) \parallel \mathcal{D}_{\text{goal}}(\mathbf{c})^n)] = \text{Cor}\left(\mathbf{T}'_{[k+1, k+n]} \mid \mathbf{T}'_{\leq k}\right).$$

We then choose  $k \leq k_{\max}$  to be the parameter selected by [Lemma 4.3](#) and bound

$$\begin{aligned} \text{Cor}\left(\mathbf{T}'_{[k+1, k+n]} \mid \mathbf{T}'_{\leq k}\right) &= \mathbb{E}_{\substack{\mathbf{A} \sim \binom{[m]}{n} \\ \mathbf{B} \sim \binom{[m] \setminus \mathbf{A}}{k}}} [\text{Cor}(\mathbf{T}'_{\mathbf{A}} \mid \mathbf{T}'_{\mathbf{B}})] && \text{(Permutation invariance of } \mathbf{T}') \\ &\leq \frac{n(n-1)}{2(k_{\max} + 1)} \cdot \mathbb{E}_{\substack{i \sim \text{Unif}([m]) \\ \mathbf{B} \sim \binom{[m] \setminus \{i\}}{n+k-1}}} [I(\mathbf{T}'_i; \mathbf{T}'_{\mathbf{B}})] && \text{(Lemma 4.3)} \\ &\leq \frac{n(n-1)}{2(k_{\max} + 1)} \cdot 2 \ln d && \text{(Lemma 6.5 and } n + k_{\max} \leq m/2) \end{aligned}$$

We are ready to give the desired bound:

$$\begin{aligned} &\mathbb{E}_{\mathbf{c}}[d_{\text{TV}}(\text{Group}_n(\mathcal{D}_{\text{adaptive}}, \mathbf{c}), \mathcal{D}_{\text{goal}}(\mathbf{c})^n)] \\ &\leq \mathbb{E}_{\mathbf{c}} \left[ \sqrt{\frac{d_{\text{KL}}(\text{Group}_n(\mathcal{D}_{\text{adaptive}}, \mathbf{c}) \parallel \mathcal{D}_{\text{goal}}(\mathbf{c})^n)}{2}} \right] && \text{(Fact 5.5)} \\ &\leq \sqrt{\mathbb{E}_{\mathbf{c}} \left[ \frac{d_{\text{KL}}(\text{Group}_n(\mathcal{D}_{\text{adaptive}}, \mathbf{c}) \parallel \mathcal{D}_{\text{goal}}(\mathbf{c})^n)}{2} \right]} && \text{(Jensen's inequality)} \\ &\leq \sqrt{\frac{n(n-1) \ln d}{2(k_{\max} + 1)}}, \end{aligned}$$

which is easily upper bounded by  $\sqrt{\frac{n^2 \ln d}{2k_{\max}}}$ . □

## 7 Rounding to a legal corruption: Proof of [Lemma 4.4](#)

In this section, we prove the following, restated for convenience.

**Lemma 4.4** (Error due to rounding for our grouping strategy). *Using the same notation as [Lemma 4.2](#), as long as  $k \leq m/2$*

$$\mathbb{E}_{\mathbf{c}} \left[ \inf_{\mathcal{D}_{\text{oblivious}} \in \mathcal{D}_{\text{oblivious}}^{n,\rho,\mathcal{D}}} \{d_{\text{TV}}(\mathcal{D}_{\text{oblivious}}, \mathcal{D}_{\text{goal}}(\mathbf{c})^n)\} \right] \leq 2n \cdot \sqrt{\frac{k \ln d}{m}}.$$

Our proof of [Lemma 4.4](#) will first prove the following statement, which applies to any grouping strategy.

**Lemma 7.1** (Rounding to few groups incurs small error). *For any base distribution  $\mathcal{D}$ , sample size  $m$ , cost function  $\rho$ , adaptive corruption  $\mathcal{D}_{\text{adaptive}} \in \mathcal{D}_{\text{adaptive}}^{m, \rho, \mathcal{D}}$ , draw  $\mathbf{S}' \sim \mathcal{D}_{\text{adaptive}}$  and let  $\mathbf{g}$  be any random variable supported on a set  $G$  specifying which group  $\mathbf{S}'$  falls in. Then,*

$$\mathbb{E}_{\mathbf{g}} \left[ \inf_{\mathcal{D}' \in \mathcal{C}_{\rho}(\mathcal{D})} \{d_{\text{TV}}(\mathcal{D}', \mathcal{D}'_{\text{goal}}(\mathbf{g}))\} \right] \leq \sqrt{\frac{\ln(|G|)}{2m}}.$$

where  $\mathcal{D}'_{\text{goal}}(g) := \mathbb{E}[\text{Unif}(\mathbf{S}') \mid \mathbf{g} = g]$ .

We begin by proving [Lemma 7.1](#). Later, in [Section 7.3](#), we show how to use it to derive [Lemma 4.4](#).

## 7.1 Proof of [Lemma 7.1](#)

The proof of [Lemma 7.1](#) is broken into two pieces. First, we prove it when the adversary cannot make any corruptions:

**Claim 7.2.** *For any distribution  $\mathcal{D}$ , draw  $\mathbf{S} \sim \mathcal{D}^m$  and  $\mathbf{g}$  be any random variable on the same probability space as  $\mathbf{S}$  supported on a set  $G$ . Then,*

$$\mathbb{E}_{\mathbf{g}}[d_{\text{TV}}(\mathcal{D}, \mathcal{D}_{\text{goal}}(\mathbf{g}))] \leq \sqrt{\frac{\ln(|G|)}{2m}} \quad \text{where } \mathcal{D}_{\text{goal}}(g) := \mathbb{E}[\text{Unif}(\mathbf{S}) \mid \mathbf{g} = g].$$

This is a special case of [Lemma 7.1](#) corresponding to when the cost function is simply

$$\rho(x, y) = \begin{cases} 0 & \text{if } x = y \\ \infty & \text{if } x \neq y. \end{cases}$$

Then, we will use the following to show that the special case in [Claim 7.2](#) is sufficient to prove the more general case.

**Claim 7.3** (Lipschitzness of corruptions). *For any cost function  $\rho$ , distributions  $\mathcal{D}_1$  and  $\mathcal{D}_2$ , and any  $\mathcal{D}'_1 \in \mathcal{C}_{\rho}(\mathcal{D}_1)$ , there is some  $\mathcal{D}'_2 \in \mathcal{C}_{\rho}(\mathcal{D}_2)$  for which*

$$d_{\text{TV}}(\mathcal{D}'_1, \mathcal{D}'_2) \leq d_{\text{TV}}(\mathcal{D}_1, \mathcal{D}_2).$$

*Proof of [Lemma 7.1](#) using [Claims 7.2](#) and [7.3](#).* Since  $\mathcal{D}_{\text{adaptive}}$  is a valid adaptive corruption, it is possible to couple a sample  $\mathbf{S} \sim \mathcal{D}^m$  with  $\mathbf{S}' \sim \mathcal{D}_{\text{adaptive}}$  so that  $\mathbf{S}' \in \mathcal{C}_{\rho}(\mathbf{S})$  with probability 1. Then, define  $\mathcal{D}_{\text{goal}}(g) := \mathbb{E}[\text{Unif}(\mathbf{S}) \mid \mathbf{g} = g]$ . Then, [Claim 7.2](#) gives that

$$\mathbb{E}_{\mathbf{g}}[d_{\text{TV}}(\mathcal{D}, \mathcal{D}_{\text{goal}}(\mathbf{g}))] \leq \sqrt{\frac{\ln(|G|)}{2m}} \quad \text{where } \mathcal{D}_{\text{goal}}(g) := \mathbb{E}[\text{Unif}(\mathbf{S}) \mid \mathbf{g} = g].$$

It therefore suffices to show that, for every  $g \in G$ , there is some  $\mathcal{D}' \in \mathcal{C}_{\rho}(\mathcal{D})$  for which

$$d_{\text{TV}}(\mathcal{D}', \mathcal{D}'_{\text{goal}}(g)) \leq d_{\text{TV}}(\mathcal{D}, \mathcal{D}_{\text{goal}}(g)).$$

Using [Claim 7.3](#), this is implied by showing that  $\mathcal{D}'_{\text{goal}}(g) \in \mathcal{C}_\rho(\mathcal{D}_{\text{goal}}(g))$ . This means showing a coupling of  $\mathbf{x} \sim \mathcal{D}_{\text{goal}}(g)$  and  $\mathbf{x}' \sim \mathcal{D}'_{\text{goal}}(g)$  for which  $\mathbb{E}[\rho(\mathbf{x}, \mathbf{x}')] \leq 1$ .

For this coupling, first condition on  $\mathbf{g} = g$  and then take a uniform index  $i \sim \text{Unif}([m])$ . We will set  $\mathbf{x} := \mathbf{S}_i$  and  $\mathbf{x}' := \mathbf{S}'_i$ . By definition, the marginal distribution of  $\mathbf{x}$  is  $\mathcal{D}_{\text{goal}}(g)$  and of  $\mathbf{x}'$  is  $\mathcal{D}'_{\text{goal}}(g)$ . Finally,

$$\mathbb{E}[\rho(\mathbf{x}, \mathbf{x}')] = \mathbb{E}\left[\frac{1}{m} \sum_{i \in [m]} \rho(\mathbf{S}_i, \mathbf{S}'_i)\right] \leq 1. \quad \square$$

### 7.1.1 The special case when the adversary cannot corrupt: Proof of [Claim 7.2](#)

This proof will use a few standard facts about sub-Gaussian random variables. We refer the reader to [\[Ver18\]](#) for a more thorough treatment of sub-Gaussian random variables.

**Definition 14** (Sub-Gaussian random variables). *A random variable  $z$  is said to be sub-Gaussian with variance-proxy  $\sigma^2$  if, for all  $\lambda \in \mathbb{R}$ ,*

$$\mathbb{E}[e^{\lambda z}] \leq \exp\left(\frac{\lambda^2 \sigma^2}{2}\right).$$

**Fact 7.4** (Mean of bounded and independent random variables is sub-Gaussian). *Let  $z_1, \dots, z_m$  each be independent mean-0 random variables bounded on  $[a, a+1]$  for some  $a \in \mathbb{R}$ . Then  $\mathbf{Z} := \frac{z_1 + \dots + z_m}{m}$  is sub-Gaussian with variance-proxy  $\frac{1}{4m}$ .*

**Fact 7.5** (Expected maximum of sub-Gaussian random variables). *Let  $\mathbf{Z}_1, \dots, \mathbf{Z}_n$  be (not necessarily independent) sub-Gaussian random variables each with variance proxy  $\sigma^2$ . Then,*

$$\mathbb{E}\left[\max_{i \in [n]} \mathbf{Z}_i\right] \leq \sqrt{2\sigma^2 \ln(n)}.$$

*Proof of [Claim 7.2](#).* Our goal is to bound

$$\begin{aligned} \mathbb{E}_{\mathbf{g}}[d_{\text{TV}}(\mathcal{D}_{\text{goal}}(\mathbf{g}), \mathcal{D})] &= \mathbb{E}_{\mathbf{g}}\left[\sup_{T: X \rightarrow [0,1]} \left\{ \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{goal}}(\mathbf{g})}[T(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[T(\mathbf{x})] \right\}\right] \\ &= \sup_{T_g: X \rightarrow [0,1] \text{ for all } g \in G} \left\{ \mathbb{E}_{\mathbf{g}} \left[ \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{goal}}(\mathbf{g})}[T_g(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[T_g(\mathbf{x})] \right] \right\}. \end{aligned}$$

For the remainder of this proof we will fix an arbitrary choice of  $\{T_g\}_{g \in G}$  and upper bound the above quantity. First, we shift the  $T_g$  so that the second term is 0 by defining

$$T'_g(x) = T_g(x) - \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[T_g(\mathbf{x})].$$

The result is that the range of  $T'_g(x)$  is of the form  $[a, a+1]$  for some  $a \in \mathbb{R}$  and that  $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[T'_g(\mathbf{x})] = 0$ .

Our goal is to upper bound the quantity  $\mathbb{E}_{\mathbf{g}} \left[ \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{goal}}(\mathbf{g})}[T'_g(\mathbf{x})] \right]$  over all such  $\{T'_g\}_{g \in G}$ .

Since  $\mathcal{D}_{\text{goal}}(g) := \mathbb{E}[\text{Unif}(\mathbf{S}) \mid \mathbf{g} = g]$ , for any  $T'_g$  we have that

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{goal}}(g)}[T'_g(\mathbf{x})] = \mathbb{E}_{\mathbf{S} | \mathbf{g} = g} \left[ \frac{1}{m} \cdot \sum_{i \in [m]} T'_g(\mathbf{S}_i) \right].$$

Therefore, we wish to upper bound

$$\mathbb{E}_{\mathbf{S} \sim \mathcal{D}^m} \left[ \mathbb{E}_{\mathbf{g} | \mathbf{S}} \left[ \frac{1}{m} \cdot \sum_{i \in [m]} T'_g(\mathbf{S}_i) \right] \right] \leq \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^m} \left[ \max_{g \in G} \left( \frac{1}{m} \cdot \sum_{i \in [m]} T'_g(\mathbf{S}_i) \right) \right].$$

For each  $g$ , let us define the random variable  $\mathbf{Z}_g$  to be  $\frac{1}{m} \cdot \sum_{i \in [m]} T'_g(\mathbf{S}_i)$ . Then, the above is simply  $\mathbb{E}[\max_{g \in G} \mathbf{Z}_g]$ . Furthermore, by [Fact 7.4](#), each  $\mathbf{Z}_g$  is sub-Gaussian with variance-proxy  $\frac{1}{4m}$ . Finally, we apply [Fact 7.5](#) to give that

$$\mathbb{E}_{\mathbf{g}} [d_{\text{TV}}(\mathcal{D}_{\text{goal}}(\mathbf{g}), \mathcal{D})] \leq \mathbb{E}[\max_{g \in G} \mathbf{Z}_g] \leq \sqrt{\frac{\ln |G|}{2m}}. \quad \square$$

## 7.2 Generalizing to when the adversary can corrupt: Proof of [Claim 7.3](#)

*Proof.* By [Definition 10](#), it suffices to show that for any coupling of  $\mathbf{x}_1 \sim \mathcal{D}_1$  and  $\mathbf{x}_2 \sim \mathcal{D}_2$ , there is a coupling of  $\mathbf{y}_1 \sim \mathcal{D}'_1$  and  $\mathbf{y}_2 \sim \mathcal{D}'_2$  for which

$$\Pr[\mathbf{y}_1 \neq \mathbf{y}_2] \leq \Pr[\mathbf{x}_1 \neq \mathbf{x}_2].$$

Fix any such coupling of  $\mathbf{x}_1$  and  $\mathbf{x}_2$  and let  $\mathcal{D}_{\mathbf{x}_2 | \mathbf{x}_1}$  denote the distribution of  $\mathbf{x}_2$  conditioned on  $\mathbf{x}_1 = x_1$  under this coupling.

Since  $\mathcal{D}'_1 \in \mathcal{C}_\rho(\mathcal{D}_1)$ , there is a coupling of  $\mathbf{x}_1 \sim \mathcal{D}$  and  $\mathbf{y}_1 \sim \mathcal{D}'_1$  for which  $\mathbb{E}[\rho(\mathbf{x}_1, \mathbf{y}_1)] \leq 1$ . Let  $\mathcal{D}_{\mathbf{y}_1 | \mathbf{x}_1}$  be the distribution of  $\mathbf{y}_1$  conditioned on  $\mathbf{x}_1 = x_1$  in this coupling.

We will now specify a joint distribution over all of  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2$ .

1. Draw  $\mathbf{x}_1 \sim \mathcal{D}$ .
2. Draw  $\mathbf{x}_2 \sim \mathcal{D}_{\mathbf{x}_2 | \mathbf{x}_1}$ . This will result in the marginal distribution of  $\mathbf{x}_2$  being  $\mathcal{D}'_2$ .
3. Draw  $\mathbf{y}_1 \sim \mathcal{D}_{\mathbf{y}_1 | \mathbf{x}_1}$ . This will result in the marginal distribution of  $\mathbf{y}'_1$  being  $\mathcal{D}'_1$ .
4. Set  $\mathbf{y}_2$  to

$$\mathbf{y}_2 = \begin{cases} \mathbf{y}_1 & \text{if } \mathbf{x}_1 = \mathbf{x}_2 \\ \mathbf{x}_2 & \text{otherwise,} \end{cases} \quad (3)$$

and define  $\mathcal{D}'_2$  to be the distribution over  $\mathbf{y}_2$ .

The desired result follows from the following two claims:

**Claim 1,**  $\Pr[\mathbf{y}_1 \neq \mathbf{y}_2] \leq \Pr[\mathbf{x}_1 \neq \mathbf{x}_2]$ : For this, we simply observe from [Equation \(3\)](#) that if  $\mathbf{x}_1 = \mathbf{x}_2$  then  $\mathbf{y}_1 = \mathbf{y}_2$ . Therefore,  $\Pr[\mathbf{y}_1 = \mathbf{y}_2] \geq \Pr[\mathbf{x}_1 = \mathbf{x}_2]$ , and negating this gives the desired result.

**Claim 2,**  $\mathcal{D}'_2 \in \mathcal{C}_\rho(\mathcal{D}_2)$ : For this, it suffices to show that  $\mathbb{E}[\rho(\mathbf{x}_2, \mathbf{y}_2)] \leq 1$ . We bound,

$$\begin{aligned} \mathbb{E}[\rho(\mathbf{x}_2, \mathbf{y}_2)] &= \mathbb{E}[\rho(\mathbf{x}_1, \mathbf{y}_1) \cdot \mathbf{1}[\mathbf{x}_1 = \mathbf{x}_2]] && \text{(Equation (3) and } \rho(x_2, x_2) = 0 \text{ for any } x_2) \\ &\leq \mathbb{E}[\rho(\mathbf{x}_1, \mathbf{y}_1)] \leq 1. \end{aligned}$$

□

### 7.3 Proof of Lemma 4.4

We conclude this section with a better error bound fine tuned for the particular grouping strategy we utilize. The proof of this improved bound will use the more general and weaker bound of Lemma 7.1 as a black-box. Roughly speaking, Lemma 4.4 obtains the bound that Lemma 7.1 would obtain if there were  $d^k$  possible groups, even though there are actually  $|X|^k$ , which can be substantially larger. This improvement comes the fact that, in Claim 7.6 below, after conditioning on  $\mathbf{c}$ , there are only  $d^k$  choices remaining for  $\mathbf{c}'$ .

**Claim 7.6.** *For any  $\mathcal{D}_{\text{adaptive}} \in \mathcal{D}_{\text{adaptive}}^{m,\rho,\mathcal{D}}$  where  $\rho$  is a degree- $d$  cost function, draw  $\mathbf{T} \sim \mathcal{D}_{\text{adaptive}}$  and  $\mathbf{I} \subseteq [m]$  a uniform random subset of  $k$  indices from  $[m]$  and define*

$$\begin{aligned} \mathbf{c} &= \mathbf{T}_{\mathbf{I}} & \mathbf{R} &= \mathbf{T}_{[m]\setminus\mathbf{I}} \\ \mathbf{c}' &= \mathbf{T}'_{\mathbf{I}} & \mathbf{R}' &= \mathbf{T}'_{[m]\setminus\mathbf{I}}. \end{aligned}$$

Then, for  $\mathcal{D}_{\text{goal}}(\mathbf{c}, \mathbf{c}') := \mathbb{E}[\text{Unif}(\mathbf{R}') \mid \mathbf{c} = \mathbf{c} \text{ and } \mathbf{c}' = \mathbf{c}']$  and any fixed value of  $\mathbf{c}$ , we have that

$$\mathbb{E}_{\mathbf{c}'|\mathbf{c}=\mathbf{c}} \left[ \inf_{\mathcal{D}' \in \mathcal{C}_{\rho}(\mathcal{D})} \{d_{\text{TV}}(\mathcal{D}', \mathcal{D}_{\text{goal}}(\mathbf{c}, \mathbf{c}'))\} \right] \leq \sqrt{\frac{k \ln d}{2(m-k)}} + \frac{k}{m}$$

*Proof.* The key observation underlying this proof is that  $\mathbf{R}'$  can be understood as the corruption of an appropriately defined adaptive adversary, even after conditioning on  $\mathbf{c} = \mathbf{c}$ . Since  $\mathbf{S}' \in \mathcal{C}_{\rho}(\mathbf{S})$  with probability 1, we have that

$$\sum_{i \in [m-k]} \rho(\mathbf{R}_i, \mathbf{R}'_i) \leq \sum_{i \in [m]} \rho(\mathbf{S}_i, \mathbf{S}'_i) \leq m.$$

Note that this does not necessarily mean that  $\mathbf{R}' \in \mathcal{C}_{\rho}(\mathbf{R})$ , as that would require  $\sum_{i \in [m-k]} \rho(\mathbf{R}_i, \mathbf{R}'_i) \leq m - k$ . Instead, we have with probability 1 that  $\mathbf{R}'$  is a valid corruption of  $\mathbf{R}$  for a slightly more permissive cost function,

$$\mathbf{R}' \in \mathcal{C}_{\bar{\rho}}(\mathbf{R}) \quad \text{for} \quad \bar{\rho}(x, y) = \frac{m-k}{m} \cdot \rho(x, y).$$

We also observe that even if we condition on  $\mathbf{c} = \mathbf{c}$  the distribution of  $\mathbf{R}$  is still simply  $\mathcal{D}^{m-k}$ . This is because we had originally drawn  $\mathbf{S} \sim \mathcal{D}^m$  and so  $\mathbf{c}$  and  $\mathbf{R}$  are independent. Combining these two observations, we have that the distribution of  $\mathbf{R}'$  conditioned on  $\mathbf{c} = \mathbf{c}$  is in the class  $\mathcal{D}_{\text{adaptive}}^{m-k, \bar{\rho}, \mathcal{D}}$ . Then, since there are at most  $d^k$  choices for  $\mathbf{c}'$  after conditioning on  $\mathbf{c} = \mathbf{c}$ , we can apply Lemma 7.1 to obtain,

$$\mathbb{E}_{\mathbf{c}'|\mathbf{c}=\mathbf{c}} \left[ \inf_{\bar{\mathcal{D}}' \in \mathcal{C}_{\bar{\rho}}(\mathcal{D})} \{d_{\text{TV}}(\mathcal{D}', \mathcal{D}_{\text{goal}}(\mathbf{c}, \mathbf{c}'))\} \right] \leq \sqrt{\frac{\ln d^k}{2(m-k)}} = \sqrt{\frac{k \ln d}{2(m-k)}}.$$

All that remains is to show that for every  $\bar{\mathcal{D}}' \in \mathcal{C}_{\bar{\rho}}(\mathcal{D})$ , there is some  $\mathcal{D}' \in \mathcal{C}_{\rho}(\mathcal{D})$  satisfying  $d_{\text{TV}}(\mathcal{D}', \bar{\mathcal{D}}') \leq k/m$ , as the desired result will then follow by the triangle inequality (Fact 5.2).

Since  $\bar{\mathcal{D}}' \in \mathcal{C}_{\bar{\rho}}(\mathcal{D})$ , there is a coupling of  $\mathbf{x} \sim \mathcal{D}$  and  $\bar{\mathbf{x}}' \sim \bar{\mathcal{D}}'$  for which  $\mathbb{E}[\bar{\rho}(\mathbf{x}, \bar{\mathbf{x}}')] \leq 1$ . Let  $\mathcal{D}'$  be the distribution of

$$\mathbf{x}' = \begin{cases} \bar{\mathbf{x}}' & \text{with probability } \frac{m-k}{m} \\ \mathbf{x} & \text{with probability } \frac{k}{m}. \end{cases}$$

Then,  $d_{\text{TV}}(\mathcal{D}', \bar{\mathcal{D}}') \leq \Pr[\mathbf{x}' \neq \bar{\mathbf{x}}'] \leq \frac{k}{m}$ , and  $\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})$  because

$$\mathbb{E}[\rho(\mathbf{x}, \mathbf{x}')] = \frac{m-k}{m} \mathbb{E}[\rho(\mathbf{x}, \bar{\mathbf{x}}')] = \frac{m-k}{m} \cdot \frac{m}{m-k} \mathbb{E}[\bar{\rho}(\mathbf{x}, \bar{\mathbf{x}}')] \leq 1. \quad \square$$

We also use the following simple ingredient to prove [Lemma 4.4](#)

**Proposition 7.7.** *Let  $\{\mathcal{D}(g)\}_{g \in G}$  be a family of distributions. Then for any base distribution  $\mathcal{D}$ , cost function  $\rho$ , and random variable  $\mathbf{g}$  supported on  $G$ ,*

$$\inf_{\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})} \left\{ d_{\text{TV}}\left(\mathcal{D}', \mathbb{E}_{\mathbf{g}}[\mathcal{D}(\mathbf{g})]\right) \right\} \leq \mathbb{E}_{\mathbf{g}} \left[ \inf_{\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})} \{d_{\text{TV}}(\mathcal{D}', \mathcal{D}(\mathbf{g}))\} \right].$$

*Proof.* Fix any  $\mathcal{D}'(g) \in \mathcal{C}_\rho(\mathcal{D})$  for all  $g \in G$ . We will show there exists a choice of  $\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})$  for which

$$d_{\text{TV}}\left(\mathcal{D}', \mathbb{E}_{\mathbf{g}}[\mathcal{D}(\mathbf{g})]\right) \leq \mathbb{E}_{\mathbf{g}}[d_{\text{TV}}(\mathcal{D}'(\mathbf{g}), \mathcal{D}(\mathbf{g}))], \quad (4)$$

which implies the desired inequality. Define,

$$\mathcal{D}' := \mathbb{E}_{\mathbf{g}}[\mathcal{D}'(\mathbf{g})].$$

Then, [Equation \(4\)](#) holds by [Fact 5.4](#). All that remains is to show that  $\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})$ . For all  $g \in G$ , there must be a coupling of  $\mathbf{x}' \sim \mathcal{D}'(g)$  and  $\mathbf{x} \sim \mathcal{D}$  for which  $\mathbb{E}[\rho(\mathbf{x}, \mathbf{x}')] \leq 1$ . Therefore, if we first draw  $\mathbf{g}$  and then  $\mathbf{x}, \mathbf{x}'$  from this coupling for  $\mathcal{D}'(\mathbf{g})$  and  $\mathcal{D}$ , then

1. The marginal distribution of  $\mathbf{x}$  will be that of  $\mathcal{D}$ . This even holds conditioning on any value of  $\mathbf{g} = g$ .
2. The marginal distribution of  $\mathbf{x}'$  will be  $\mathcal{D}'$ . This is by the definition of mixture distributions.
3. We can bound,

$$\mathbb{E}[\rho(\mathbf{x}, \mathbf{x}')] \leq \max_g \mathbb{E}[\rho(\mathbf{x}, \mathbf{x}') \mid \mathbf{g} = g] \leq 1.$$

Therefore,  $\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})$ . □

We are now ready to prove the main result of this subsection.

*Proof of [Lemma 4.4](#).* First, if  $d = 1$ , the only legal cost function is

$$\rho(x, y) = \begin{cases} 0 & \text{if } x = y \\ \infty & \text{otherwise.} \end{cases}$$

In this case, the adversaries can make no corruptions implying that  $\mathcal{D}_{\text{adaptive}}$  must be  $\mathcal{D}^m$  and  $\mathcal{D}_{\text{goal}}(c) = \mathcal{D}$  for all choices of  $c$ . The desired error bound of 0 easily holds. Therefore, we may assume  $d \geq 2$ .

Let  $\mathbf{c}, \mathbf{c}', \mathbf{R}, \mathbf{R}'$  be as defined in [Claim 7.6](#). We first observe that by taking  $\mathbf{S}' \sim \Phi_{m-k \rightarrow n}(\mathbf{R}')$ , we have that  $(\mathbf{S}', \mathbf{c}')$  are distributed according to  $\text{Grouped}_{n,k}(\mathcal{D}_{\text{adaptive}})$ . Therefore,

$$\mathcal{D}_{\text{goal}}(\mathbf{c}') = \mathbb{E}[\text{Unif}(\mathbf{R}') \mid \mathbf{c}' = \mathbf{c}'] = \mathbb{E}_{\mathbf{c} \mid \mathbf{c}' = \mathbf{c}'}[\mathcal{D}_{\text{goal}}(\mathbf{c}, \mathbf{c}')],$$

where  $\mathcal{D}_{\text{goal}}(\mathbf{c}, \mathbf{c}')$  is as defined in [Claim 7.6](#). We proceed to bound

$$\begin{aligned}
\mathbb{E}_{\mathbf{c}'} \left[ \inf_{\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})} \{d_{\text{TV}}(\mathcal{D}', \mathcal{D}_{\text{goal}}(\mathbf{c}'))\} \right] &= \mathbb{E}_{\mathbf{c}'} \left[ \inf_{\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})} \left\{ d_{\text{TV}} \left( \mathcal{D}', \mathbb{E}_{\mathbf{c}|\mathbf{c}'} [\mathcal{D}_{\text{goal}}(\mathbf{c}, \mathbf{c}')] \right) \right\} \right] \\
&\leq \mathbb{E}_{\mathbf{c}'} \left[ \mathbb{E}_{\mathbf{c}|\mathbf{c}'} \left[ \inf_{\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})} \{d_{\text{TV}}(\mathcal{D}', \mathcal{D}_{\text{goal}}(\mathbf{c}, \mathbf{c}'))\} \right] \right] \quad (\text{Proposition 7.7}) \\
&\leq \mathbb{E}_{\mathbf{c}} \left[ \mathbb{E}_{\mathbf{c}'|\mathbf{c}} \left[ \inf_{\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})} \{d_{\text{TV}}(\mathcal{D}', \mathcal{D}_{\text{goal}}(\mathbf{c}, \mathbf{c}'))\} \right] \right] \\
&\leq \sqrt{\frac{k \ln d}{2(m-k)}} + \frac{k}{m} \quad (\text{Claim 7.6}) \\
&\leq 2\sqrt{\frac{k \ln d}{m}} \quad (k \leq m/2 \text{ and } d \geq 2.)
\end{aligned}$$

Recall that  $\mathcal{D}_{\text{oblivious}}^{n,\rho,\mathcal{D}} := \{(\mathcal{D}')^n \mid \mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})\}$ , and so the desired bound follows from [Fact 5.3](#).  $\square$

## 8 Adaptive adversaries are at least as strong as oblivious adversaries

In this section, we prove the easy direction of [Theorem 6](#).

**Claim 8.1** (The adaptive adversary can simulate the oblivious adversary). *For any  $m \geq 25n^2/\varepsilon^2$ , distribution  $\mathcal{D}$ , cost function  $\rho$ , and oblivious adversary  $\mathcal{D}_{\text{oblivious}} \in \mathcal{D}_{\text{oblivious}}^{n,\rho,\mathcal{D}}$ , there is a corresponding adaptive adversary  $\mathcal{D}_{\text{adaptive}} \in \Phi_{m \rightarrow n}(\mathcal{D}_{\text{adaptive}}^{m,\rho,\mathcal{D}})$  satisfying*

$$d_{\text{TV}}(\mathcal{D}_{\text{adaptive}}, \mathcal{D}_{\text{oblivious}}) \leq \varepsilon.$$

Such a statement is well-known to hold for some specific [\[DKK<sup>+</sup>19, ZJS19\]](#) adversary models. Here, we show it holds with any cost function.

The high-level idea in the proof of [Claim 8.1](#) is simple: The oblivious adversary will have chosen some  $\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})$  (in which case  $\mathcal{D}_{\text{oblivious}} = (\mathcal{D}')^n$ ) for which it is possible to couple samples  $\mathbf{S} \sim \mathcal{D}^n$  and  $\mathbf{S}' \sim (\mathcal{D}')$  so that the *average* cost of corrupting a point in  $\mathbf{S}_i$  to the corresponding point in  $\mathbf{S}'_i$  is at most 1. If it were always the case that  $\mathbf{S}' \in \mathcal{C}_\rho(\mathbf{S})$ , we would be done, as the adaptive adversary could then always corrupt  $\mathbf{S}$  to  $\mathbf{S}'$ . However, even though the corruption of  $\mathbf{S}$  to  $\mathbf{S}'$  has an average cost of 1, it can exceed 1 for some draws of  $\mathbf{S}, \mathbf{S}'$ , in which case the adaptive adversary can not exactly simulate the oblivious adversary.

Therefore, our strategy will be to round  $\mathbf{S}'$  to a valid corruption. The quantity we need to bound is how many points of  $\mathbf{S}'$  we need to change to make the corruption valid, which we do in the following lemma.

**Claim 8.2.** *For any distribution  $\mathcal{D}_{\text{cost}}$  supported on  $\mathbb{R}_{\geq 0}$  with mean 1, draw  $\mathbf{x}_1, \dots, \mathbf{x}_n \stackrel{\text{iid}}{\sim} \mathcal{D}_{\text{cost}}$  and define  $\Delta(\mathbf{x}_1, \dots, \mathbf{x}_n)$  to be the minimum number of  $x_i$  that must be removed so that the sum of the remaining elements is at most  $n$ . Then,*

$$\mathbb{E}[\Delta(\mathbf{x}_1, \dots, \mathbf{x}_n)] \leq 5\sqrt{n}.$$

The main ingredient in the proof of [Claim 8.2](#) is an upper bound on the probability that  $\Delta(\mathbf{x}_1, \dots, \mathbf{x}_n)$  exceeds a value.

**Proposition 8.3.** *In the setting of [Claim 8.2](#), for any  $v \geq 0$*

$$\Pr[\Delta(\mathbf{x}_1, \dots, \mathbf{x}_n) \geq v] \leq \frac{2}{v} + \frac{4n}{v^2}.$$

*Proof.* The main idea in this proof is, for any  $r \in [0, 1]$ , to exhibit a strategy with the following properties.

1. The probability the strategy removes more than  $2nr$  elements is at most  $\frac{1}{nr}$ .
2. The probability that, after this strategy removes elements, the remaining sum is more than  $n$  is at most  $\frac{1}{nr^2}$ .

Combining the above with a union bound gives that

$$\Pr[\Delta(\mathbf{x}) \geq 2nr] \leq \frac{1}{nr} + \frac{1}{nr^2}.$$

The above is equivalent to the desired result as we can set  $r = \frac{v}{2n}$ .

For any  $\tau \in \mathbb{R}$ ,  $p \in [0, 1]$  consider the strategy that, for each  $i \in [m]$  keeps  $\mathbf{x}_i$  with probability  $f(\mathbf{x}_i)$  and otherwise removes it for

$$f(x) := \begin{cases} 1 & \text{if } x < \tau \\ p & \text{if } x = \tau \\ 0 & \text{if } x > \tau. \end{cases}$$

It is always possible to choose  $\tau$  and  $p$  so that the probability  $\mathbf{x}_i$  is removed is any desired value. We set them so that the probability  $\mathbf{x}_i$  is removed is exactly  $r$ .

**It is unlikely many elements are removed.** Let  $\mathbf{R}$  be the random variable representing the number of removed elements. Then, the distribution of  $\mathbf{R}$  is simply  $\text{Bin}(n, r)$  and so it satisfies  $\mathbb{E}[\mathbf{R}] = nr$  and  $\text{Var}[\mathbf{R}] \leq nr$ . By Chebyshev's inequality:

$$\Pr[\mathbf{R} \geq 2nr] \leq \frac{\text{Var}[\mathbf{R}]}{(2nr - \mathbb{E}[\mathbf{R}])^2} \leq \frac{nr}{(nr)^2} = \frac{1}{nr}.$$

**It is unlikely the sum of the remaining elements is more than  $n$ .** Let  $\mathbf{X}$  be the sum of the remaining elements. Then,

$$\mathbf{x} := \sum_{i \in [n]} z_i \cdot \mathbf{x}_i \quad \text{for } z_i \sim \text{Ber}(f(\mathbf{x}_i)).$$

We will analyze the mean and variance of this  $\mathbf{X}$  and then use Chebyshev's inequality to bound the probability it is more than 1. For the mean,

$$\begin{aligned} \mathbb{E}[\mathbf{X}] &= n \cdot \mathbb{E}[z_i \cdot \mathbf{x}_i] && \text{(Linearity of expectation)} \\ &= n \cdot (\mathbb{E}[\mathbf{x}_i] - \Pr[z_i = 0] \cdot \mathbb{E}[\mathbf{x}_i \mid z_i = 0]) && (z_i \text{ supported on } \{0, 1\}) \\ &\leq n \cdot (1 - r \cdot \tau). \end{aligned}$$

For the variance of  $\mathbf{X}$ , since  $\mathbf{x}_1, \dots, \mathbf{x}_m$  are independent, the variances sum. Therefore,

$$\begin{aligned} \text{Var}[\mathbf{X}] &= n \cdot \text{Var}[\mathbf{z}_i \cdot \mathbf{x}_i] \\ &\leq n \cdot \mathbb{E}[(\mathbf{z}_i \cdot \mathbf{x}_i) \cdot (\mathbf{z}_i \cdot \mathbf{x}_i)] \\ &\leq n \cdot \max(\mathbf{z}_i \cdot \mathbf{x}_i) \cdot \mathbb{E}[\mathbf{z}_i \cdot \mathbf{x}_i] \\ &\leq n\tau. \end{aligned}$$

Therefore, by applying Chebyshev's inequality, we see that

$$\Pr[\mathbf{X} \geq n] \leq \frac{n\tau}{n^2\tau^2r^2} = \frac{1}{n\tau^2r^2}.$$

Note furthermore that if  $\tau < 1$ , then every term in the sum is less than 1, in which case  $\Pr[\mathbf{P} \geq 1] = 0$ . Therefore, the worst-case choice of  $\tau$  for our bound is  $\tau = 1$  in which case

$$\Pr[\mathbf{P} > n] \leq \frac{1}{nr^2}. \quad \square$$

*Proof of Claim 8.2.* We write,

$$\begin{aligned} \mathbb{E}[\Delta(\mathbf{z})] &= \int_0^n \Pr[\Delta(\mathbf{z}) \geq v] dv \\ &\leq \int_0^n \max\left(1, \frac{2}{v} + \frac{4n}{v^2}\right) dv \\ &\leq 2\sqrt{n} + \int_{2\sqrt{n}}^n \frac{4n}{v^2} dv + \int_{2\sqrt{n}}^n \frac{2}{v} dv \\ &\leq 2\sqrt{n} + \int_{2\sqrt{n}}^\infty \frac{4n}{v^2} dv + \int_{2\sqrt{n}}^n \frac{2}{2\sqrt{n}} dv \\ &= 2\sqrt{n} + 2\sqrt{n} + \sqrt{n} = 5\sqrt{n}. \quad \square \end{aligned}$$

We are now ready to prove the main result of this section.

*Proof of Claim 8.1.* Consider any  $\mathcal{D}_{\text{oblivious}} \in \mathcal{D}_{\text{oblivious}}^{n,\rho,\mathcal{D}}$ . Then, there is some  $\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})$  for which  $\mathcal{D}_{\text{oblivious}} = (\mathcal{D}')^n$ . By definition, there is a coupling between  $\mathbf{x} \sim \mathcal{D}$  and  $\mathbf{y} \sim \mathcal{D}'$  so that  $\mathbb{E}[\rho(\mathbf{x}, \mathbf{y})] \leq 1$ . We can extend this to a coupling of  $\mathbf{S} \sim \mathcal{D}^m$  and  $\mathbf{S}' \sim (\mathcal{D}')^m$  so that

$$\mathbb{E}\left[\frac{1}{m} \cdot \sum_{i \in [m]} \rho(\mathbf{S}_i, \mathbf{S}'_i)\right] \leq 1.$$

Let  $\mathcal{D}_{\text{cost}}$  be the distribution of  $\rho(\mathbf{x}, \mathbf{y})$  in this coupling. By Claim 8.2, we can construct  $\mathbf{S}''$  for which  $\frac{1}{m} \cdot \sum_{i \in [m]} \rho(\mathbf{S}_i, \mathbf{S}''_i) \leq 1$  with probability 1 and for which the expected number of coordinates on which  $\mathbf{S}''$  and  $\mathbf{S}'$  differ is at most  $5\sqrt{m}$ . This is because, whenever Claim 8.2 asks to “remove” some  $\mathbf{x}_i$ , we can simply set  $\mathbf{S}''_i = \mathbf{S}_i$  in which case  $\rho(\mathbf{S}_i, \mathbf{S}''_i) = 0$ .

The adaptive adversary, given a sample  $S \in X^m$  corrupts it to the distribution of  $\mathbf{S}'' \mid \mathbf{S} = S$ . The result is that the distribution of  $\mathcal{D}_{\text{adaptive}}$  is equivalent to the distribution of  $\mathcal{D}_{\text{adaptive}}$ . Then,

for any test function  $f : X^n \rightarrow [0, 1]$

$$\begin{aligned} \mathbb{E}_{\mathbf{T} \sim \mathcal{D}_{\text{adaptive}}} [f(\mathbf{T})] &= \mathbb{E}[f \circ \Phi_{m \rightarrow n}(\mathbf{S}'')] \\ &\leq \mathbb{E}[f \circ \Phi_{m \rightarrow n}(\mathbf{S}')] + \Pr[\mathbf{S}'' \text{ differs from } \mathbf{S}' \text{ on one of } n \text{ sampled points}] \\ &\leq \mathbb{E}[f \circ \Phi_{m \rightarrow n}(\mathbf{S}')] + \frac{5n}{\sqrt{m}}. \end{aligned}$$

The distribution of  $\Phi_{m \rightarrow n}(\mathbf{S}')$  for  $\mathbf{S}' \sim (\mathcal{D}')^m$  is simply  $(\mathcal{D}')^n$ , and so the first term is simply  $\mathbb{E}_{\mathbf{T} \sim \mathcal{D}_{\text{oblivious}}} [f(\mathbf{T})]$ . By our choice of  $m$ , the second term's magnitude is at most  $\varepsilon$ . The desired result follows from the definition of total variation distance.  $\square$

## 9 Putting the pieces together: Proof of Theorem 3

In this section, we combine all the previous ingredients to prove the below theorem, restated for convenience, which also easily implies Theorems 2 and 3.

**Theorem 6** (Theorem 3 restated). *For any  $n, d \in \mathbb{N}$  where  $d \geq 2$ , domain  $X$ , and  $\varepsilon > 0$ , let  $m = O\left(\frac{n^4(\ln d)^2}{\varepsilon^4}\right)$ . Then, for any  $f : X^n \rightarrow \{0, 1\}$ , cost function  $\rho$  with degree  $d$ , and distribution  $\mathcal{D}$  supported on  $X$ ,*

$$|\text{Oblivious-Max}_\rho(f, \mathcal{D}) - \text{Adaptive-Max}_\rho(f \circ \Phi_{m \rightarrow n}, \mathcal{D})| \leq \varepsilon. \quad (1)$$

The above is a direct consequence of the following simple result combined with Lemma 4.1 and Claim 8.1

**Proposition 9.1** (Indistinguishability from simulations). *Let  $\mathcal{D}_1$  and  $\mathcal{D}_2$  be two families of distributions on the domain  $X$  satisfying,*

1. *For all  $\mathcal{D}_1 \in \mathcal{D}_1$ , there is a random variable  $\mathcal{D}_2$  supported on  $\mathcal{D}_2$  such that the mixture satisfies*

$$d_{\text{TV}}(\mathcal{D}_1, \mathbb{E}[\mathcal{D}_2]) \leq \varepsilon.$$

2. *For all  $\mathcal{D}_2 \in \mathcal{D}_2$ , there is a random variable  $\mathcal{D}_1$  supported on  $\mathcal{D}_1$  such that the mixture satisfies*

$$d_{\text{TV}}(\mathcal{D}_2, \mathbb{E}[\mathcal{D}_1]) \leq \varepsilon.$$

Then, for any  $f : X \rightarrow \{0, 1\}$ ,

$$\left| \sup_{\mathcal{D}_1 \in \mathcal{D}_1} \left\{ \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_1} [f(\mathbf{x})] \right\} - \sup_{\mathcal{D}_2 \in \mathcal{D}_2} \left\{ \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_2} [f(\mathbf{x})] \right\} \right| \leq \varepsilon.$$

*Proof.* By symmetry, it suffices to show that

$$\sup_{\mathcal{D}_1 \in \mathcal{D}_1} \left\{ \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_1} [f(\mathbf{x})] \right\} \leq \sup_{\mathcal{D}_2 \in \mathcal{D}_2} \left\{ \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_2} [f(\mathbf{x})] \right\} + \varepsilon. \quad (5)$$

Fix any  $\mathcal{D}_1 \in \mathcal{D}_1$ . Then, there is some random variable  $\mathcal{D}_2$  supported on  $\mathcal{D}_2$  such that the mixture satisfies

$$d_{\text{TV}}(\mathcal{D}_1, \mathbb{E}[\mathcal{D}_2]) \leq \varepsilon.$$

The above implies that,

$$\mathbb{E}_{\mathcal{D}_2} [\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_2} [f(\mathbf{x})]] \geq \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_1} [f(\mathbf{x})] - \varepsilon.$$

Therefore, there must be some fixed choice of  $\mathcal{D}_2 \in \mathcal{D}_2$  for which

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_2} [f(\mathbf{x})] \leq \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_1} [f(\mathbf{x})] - \varepsilon.$$

Hence [Equation \(5\)](#) holds. □

## 10 Lower bounds

### 10.1 Proof of [Theorem 4](#)

Here, we prove [Theorem 4](#). Given a distribution  $\mathcal{D}$  promised to be uniform on some  $X' \subseteq X := [n]$ , we show that approximating the cardinality of  $X'$  is harder in the presence of adaptive subtractive contamination than it is in the presence of oblivious subtractive contamination. To begin, we formalize both models. For ease of notation, we stick with the setting where the adversary can remove half of the sample or distribution, though all the conclusions would remain the same if this 1/2 were replaced with any other constant.

**Definition 15** (1/2-Subtractive contamination, special case of [Definition 16](#)). *For any distribution  $\mathcal{D}$ , we say that  $\mathcal{D}' \in \text{sub}(\mathcal{D})$  if a sample of  $\mathbf{x}' \sim \mathcal{D}'$  is equivalent to a sample  $\mathbf{x} \sim \mathcal{D}$  conditioned on an event that occurs with probability 1/2.*

*Similarly, for any sample  $S \in X^{2m}$ , we say that  $S' \in \text{sub}(S)$  if, for  $m$  unique indices  $i_1, \dots, i_m \in [2m]$ ,*

$$(S')_j = S_{i_j} \quad \text{for all } j \in [m].$$

We prove the following.

**Theorem 7** (A polynomial increase in sample size is necessary, formal version [Theorem 4](#)). *Let  $\mathcal{D}$  be a distribution on  $X = [m] := \{1, \dots, m\}$  that is promised to be uniform on some  $X' \subseteq X$ . Then for some absolute constant  $c < 1$ ,*

1. *For  $n_{\text{small}} := O(\sqrt{m})$  and any  $k \leq m$ , there is an algorithm  $f_{\text{oblivious}} : X^{n_{\text{small}}} \rightarrow \{0, 1\}$  that distinguishes between the cases where  $|X'| \geq k$  vs  $|X'| \leq ck$  with high probability even in the presence of the oblivious adversary,*

$$\begin{aligned} |X'| \geq k &\implies \inf_{\mathcal{D}' \in \text{sub}(\text{Unif}(X'))} \left( \Pr_{\mathcal{S} \sim (\mathcal{D}')^{n_{\text{small}}}} [f_{\text{oblivious}}(\mathcal{S})] \right) \geq 0.99 \\ |X'| \leq ck &\implies \sup_{\mathcal{D}' \in \text{sub}(\text{Unif}(X'))} \left( \Pr_{\mathcal{S} \sim (\mathcal{D}')^{n_{\text{small}}}} [f_{\text{oblivious}}(\mathcal{S})] \right) \leq 0.01. \end{aligned}$$

2. *For  $n_{\text{large}} = \Omega(m)$  and  $k = m$ , there is no algorithm with the same guarantees: Formally, for any  $f_{\text{adaptive}} : X^{n_{\text{large}}} \rightarrow \{0, 1\}$ , either there is an  $X'$  containing  $k$  elements for which*

$$\mathbb{E}_{\mathcal{S} \sim \text{Unif}(X')^{2n_{\text{large}}}} \left[ \inf_{\mathcal{S}' \in \text{sub}(\mathcal{S})} \mathbb{1}[f_{\text{adaptive}}(\mathcal{S}')] \right] \leq 0.51$$

or there is an  $X'$  containing at most  $ck$  elements for which

$$\mathbb{E}_{\mathbf{S} \sim \text{Unif}(X')^{2n_{\text{large}}}} \left[ \sup_{\mathbf{S}' \in \text{sub}(\mathbf{S})} \mathbb{1}[f_{\text{adaptive}}(\mathbf{S}')] \right] \geq 0.49$$

Note that by standard techniques, the above can be converted to a separation in the search version of the problem, where the goal is to approximate  $|X'|$  to multiplicative accuracy, at the cost of a  $\text{polylog}(m)$  dependence.

The construction of  $f_{\text{oblivious}}$  follows from standard results on the probability of a collision in a sample (see e.g. [GR11]).

**Fact 10.1** (The probability of a collision in a sample). *Let  $\mathcal{D}'$  be any distribution supported on  $k$  points for which  $\Pr_{\mathbf{x} \sim \mathcal{D}'}[\mathbf{x} = x] \leq \frac{2}{k}$  for all possible  $x$ . Then, for  $\mathbf{x}_1, \dots, \mathbf{x}_n \stackrel{\text{iid}}{\sim} \mathcal{D}'$ , let  $\mathbf{E}$  indicate whether there is  $i \neq j$  for which  $\mathbf{x}_i = \mathbf{x}_j$ . There are absolute constants  $c_1, c_2$  for which,*

$$c_1 n^2 \leq k \implies \Pr[\mathbf{E}] \leq 0.01 \quad \text{and} \quad c_2 n^2 \geq k \implies \Pr[\mathbf{E}] \geq 0.99.$$

Given the above fact, we just have  $f_{\text{oblivious}}$  accept if and only if it finds a collision in the sample.

The lower bound against the adaptive adversary is an easy consequence of the following simple proposition which shows the adaptive adversary can remove all duplicates from the sample.

**Proposition 10.2.** *Let  $\mathcal{D}$  be uniform on a distribution supported on  $k$  points and*

$$n := \lfloor \frac{\varepsilon k}{2} \rfloor.$$

*For  $\mathbf{x}_1, \dots, \mathbf{x}_{2n} \stackrel{\text{iid}}{\sim} \mathcal{D}$ , with probability at least  $1 - \varepsilon$ , the number  $i \in [2n]$  for which there exists  $j \neq i$  satisfying  $\mathbf{x}_i = \mathbf{x}_j$  is at most  $n$ .*

*Proof.* Let  $\mathbf{z}_i$  be the indicator that there is some  $j \neq i$  for which  $\mathbf{x}_j = \mathbf{x}_i$ . By union bound,

$$\mathbb{E}[\mathbf{z}_i] \leq \frac{2n-1}{k} \leq \frac{2n}{k}.$$

Applying Markov's inequality to  $\mathbf{Z} = \sum_{i \in [2n]} \mathbf{z}_i$ ,

$$\Pr[\mathbf{Z} \geq n] \leq \frac{\mathbb{E}[\mathbf{Z}]}{n} \leq \frac{\frac{2n^2}{k}}{n} \leq \frac{2n}{k},$$

which is at most  $\varepsilon$  for our choice of  $n$ . □

*Proof of Theorem 7.* For small enough constant  $c$ , **Fact 10.1** implies the following. Picking  $n(k) = \Theta(\sqrt{k})$  appropriately, the algorithm  $f_{\text{oblivious}}(x_1, \dots, x_n)$  which accepts iff it has a collision in its first  $n(k) \leq n_{\text{small}}$  samples has the desired behavior.

For the lower bound against adaptive adversaries, consider the adaptive adversary that given a sample  $x_1, \dots, x_{2n_{\text{large}}}$  does the following:

1. If there are less than  $n_{\text{large}}$  choices for  $i \in [2n_{\text{large}}]$  such that there is  $j \neq i$  for which  $x_i \neq x_j$ , the adaptive adversary takes any size- $n$  subset of  $x_1, \dots, x_{2n_{\text{large}}}$  not containing these collisions.
2. Otherwise, it just returns the first  $n$  elements  $x_1, \dots, x_n$ .

**Proposition 10.2** implies that, for any  $X' \subseteq X$  and  $\mathbf{x}_1, \dots, \mathbf{x}_{2n_{\text{large}}}$ , the probability the second case occurs is at most 0.01.

Now, for any  $f_{\text{adaptive}} : X^{n_{\text{large}}} \rightarrow \{0, 1\}$ , define,

$$p := \mathbb{E}_{\mathbf{S} \sim \binom{X}{n_{\text{large}}}} [f_{\text{adaptive}}(\mathbf{S})].$$

If  $p \geq 0.5$ , suppose we draw  $\mathbf{X}'$  uniformly among all size  $ck$  subsets of  $X$ . Then, conditioned on the second case of the adaptive adversary's strategy not occurring,  $f_{\text{adaptive}}$  be equally likely to receive any size  $n_{\text{large}}$  subset of  $X$ . Therefore, it accepts with probability at least  $0.5 - 0.01 = 0.49$ . There must hence exist at least one  $X'$  of size  $ck$  for which, with this strategy for the adaptive adversary, the probability that  $f_{\text{adaptive}} : X^{n_{\text{large}}}$  accepts is at least 0.49.

On the other hand, if  $p \leq 0.5$ , we make a similar argument: For  $X' = X$ , conditioned on the second case of the adaptive adversary's strategy not occurring, the sample  $f_{\text{adaptive}}$  sees is equally likely to be any size- $n_{\text{large}}$  subset of  $X$ . Therefore, it can accept with probability at most  $0.5 + 0.01 = 0.51$  for this strategy of the adaptive adversary.

In both cases,  $f_{\text{adaptive}}$  fails, completing this lower bound.  $\square$

## 10.2 Proof of **Theorem 5**

**Theorem 8** (Dependence on degree is necessary, formal version of **Theorem 5**). *For any  $b, \delta > 0$ , large enough  $n \in \mathbb{N}$  (as a function of  $b$  and  $\delta$ ), and cost function  $\rho : X \times X \rightarrow \mathbb{R}_{\geq 0} \cup \{\infty\}$  for which  $\rho(x, y) \geq 1 + \delta$  whenever  $x \neq y$ , let*

$$m = \Omega_{b, \delta} \left( \frac{n}{(\ln n)^2} \cdot \ln \deg_b(\rho) \right).$$

*If  $m > n$ , there is a function  $f : X^n \rightarrow \{0, 1\}$  and distribution  $\mathcal{D}$  over  $X$  for which*

$$\text{Adaptive-Max}_\rho(f \circ \Phi_{m \rightarrow n}, \mathcal{D}) \geq 1 - O(1/n) \quad \text{and} \quad \text{Oblivious-Max}_\rho(f, \mathcal{D}) \leq O(1/n). \quad (6)$$

We use similar ideas as [BLMT22, Theorem 8], which shows that **Theorem 8** holds in the specific case where  $\rho$  corresponds to additive noise, though we do need to generalize those ideas to this more general setting.

Let  $d$  be the largest integer such that  $2^d \leq \deg_b(\rho)$ . By **Definition 7**, we can choose a subset  $X' \subseteq X$  of cardinality  $2^d$  and point  $x^* \in X'$  for which  $\rho(x^*, y) \leq b$  for all  $y \in X'$ . Let  $M : X \rightarrow \{\pm 1\}^d \cup \{\vec{0}\}$  be any mapping that takes every element of  $X'$  to a unique element of  $\{\pm 1\}^d$  and all other elements to  $\vec{0}$ . For an appropriate threshold  $\tau > 0$ , we'll define

$$f(x_1, \dots, x_n) := \begin{cases} 1 & \text{if for every } x_i, \text{ there is an } x_j \text{ with } j \neq i \text{ s.t. } \langle M(x_i), M(x_j) \rangle \geq \tau \\ 0 & \text{otherwise.} \end{cases} \quad (7)$$

This choice of  $f$  was analyzed by [BLMT22].

**Fact 10.3** (Choosing the threshold  $\tau$ , [BLMT22]). *For any  $p \in (0, 1)$  and*

$$m \geq \Omega_p \left( \frac{nd}{(\ln n)^2} \right),$$

*there is a choice of threshold  $\tau$  for which both of the following hold:*

1. Lemma 7.1 of [BLMT22], a uniform point is hard to correlate with: For any  $x_1, \dots, x_{n-1} \in X$ ,

$$\Pr_{\mathbf{u} \sim \text{Unif}(X')} [\text{There is an } i \in [n-1] \text{ for which } \langle \mathbf{u}, x_i \rangle \geq \tau] \leq \frac{1}{n}.$$

2. Lemma 7.2 of [BLMT22], an adaptive adversary make all points correlated: Take any  $n_{\text{small}} \in \mathbb{N}$  for which  $pm \leq n_{\text{small}} \leq m$ . For  $\mathbf{S} \sim \text{Unif}(X')^{m_{\text{small}}}$ , there is a strategy for adding  $m - m_{\text{small}}$  points to  $\mathbf{S}$  to form  $\mathbf{S}'$  for which

$$\mathbb{E}[f \circ \Phi_{m \rightarrow n}(\mathbf{S}')] \geq 1 - \frac{1}{n}.$$

*Proof of Theorem 8.* Define

$$c := \min\left(\frac{1}{b}, 1 - \frac{1}{1+\delta}\right),$$

and set  $\mathcal{D}$  to the distribution that is equal to  $x^*$  with probability  $\frac{c}{2}$  and otherwise uniform over  $X'$ ,

$$\mathcal{D} := \frac{c}{2} \cdot \{x^*\} + \left(1 - \frac{c}{2}\right) \cdot \text{Unif}(X'). \quad (8)$$

Also, set

$$p := 1 - \frac{c}{4}$$

and let  $\tau$  be the threshold in Fact 10.3. We will show that both Theorem 8 holds with this choice of  $\tau$ ,  $f$  as in Equation (7), and  $\mathcal{D}$  as in Equation (8).

We begin by analyzing the oblivious adversary. First, for any  $\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})$ , since  $\rho(x, x') \geq 1 + \delta$  for each  $x \neq x'$ , there must be a coupling of  $\mathbf{x} \sim \mathcal{D}$  and  $\mathbf{x}' \sim \mathcal{D}'$  for which

$$\Pr[\mathbf{x} \neq \mathbf{x}'] \leq \frac{1}{1+\delta}$$

Furthermore, based on Equation (8), there is a coupling of  $\mathbf{x} \sim \mathcal{D}$  and  $\mathbf{u} \sim \text{Unif}(X')$  for which

$$\Pr[\mathbf{x} \neq \mathbf{u}] \leq \frac{c}{2}.$$

Combining the above, there is a coupling of  $\mathbf{x}' \sim \mathcal{D}'$  and  $\mathbf{u} \sim \text{Unif}(X')$  for which

$$\Pr[\mathbf{x}' \neq \mathbf{u}] \leq \frac{c}{2} + \frac{1}{1+\delta}.$$

In particular,

$$\Pr[\mathbf{x}' = \mathbf{u}] \geq \left(1 - \frac{1}{1+\delta}\right) - \frac{c}{2} \geq \frac{c}{2}.$$

This means there is a coupling of  $\mathbf{S} \sim \mathcal{D}'$  and  $\mathbf{U} \sim \text{Unif}(X')$  for which, independently for each  $i \in [n]$ , with probability at least  $c/2$ ,  $\mathbf{S}_i = \mathbf{U}_i$ . Because we assumed  $n$  is sufficiently large as a function of  $b$  and  $\delta$ , we are free to assume that  $c \geq \Omega((\ln n)/n)$ . As a result, with probability at least  $1 - 1/n$ , there is some  $i \in [n]$  for which  $\mathbf{S}_i = \mathbf{U}_i$ . Then,

$$\mathbb{E}_{\mathbf{S} \sim (\mathcal{D}')^n} [f(\mathbf{S})] \leq \frac{1}{n} + \mathbb{E}_{\mathbf{S}} [f(\mathbf{S}) \mid \mathbf{S}_i = \mathbf{U}_i \text{ for some } i \in [n]] \leq \frac{2}{n},$$

where the second inequality is by the first part of [Fact 10.3](#).

We proceed to analyze the adaptive adversary. To draw  $\mathbf{S} \sim \mathcal{D}^m$ , we can first draw  $\mathbf{E} \sim \text{Ber}(c/2)^m$ . Then, for each  $i \in [m]$ , if  $\mathbf{E}_i = 1$  we set  $\mathbf{S} = x^*$  and otherwise draw it uniformly from  $X'$ .

Conditioned on  $\sum_i \mathbf{E}_i = E$ , we have that  $E$  of the elements in  $\mathbf{S}$  are set to  $x^*$  and the other  $m - E$  are drawn independently and uniformly from  $X'$ . By the second part of [Fact 10.3](#), whenever  $m - E \geq pm$  (or equivalently,  $E \leq \frac{mc}{4}$ ), there is a way to modify the  $E$  many elements for which  $\mathbf{E}_i = 1$  to form a corrupted sample  $\mathbf{S}'$  satisfying

$$\mathbb{E}[f \circ \Phi_{m \rightarrow n}(\mathbf{S}')] \geq 1 - \frac{1}{n}.$$

Furthermore, if  $E \leq mc \leq \frac{m}{b}$ , the adversary has enough budget to modify all of the indices for which  $\mathbf{E}_i = 1$  to arbitrary elements of  $X'$ . Therefore, there is a strategy for the adaptive adversary so that

$$\mathbb{E} \left[ f \circ \Phi_{m \rightarrow n}(\mathbf{S}') \mid \frac{mc}{4} \leq \sum_{i \in [m]} \mathbf{E}_i \leq mc \right] \geq 1 - \frac{1}{n}.$$

We once again assume that  $c \geq \Omega((\ln n)/n)$ , which implies that  $c \geq \Omega((\ln m)/m)$ . By a Chernoff bound, this gives that

$$\Pr \left[ \frac{mc}{4} \leq \sum_{i \in [m]} \mathbf{E}_i \leq mc \right] \geq 1 - \frac{1}{n}.$$

So by union bound,

$$\mathbb{E}[f \circ \Phi_{m \rightarrow n}(\mathbf{S}')] \geq 1 - \frac{2}{n}. \quad \square$$

## 11 Acknowledgments

The authors thank Li-Yang Tan, Abhishek Shetty, and the anonymous STOC reviewers for their helpful discussions and feedback. Gregory is supported by a Simons Foundation Investigator Award, NSF award AF-2341890 and UT Austin's Foundations of ML NSF AI Institute. Guy is supported by NSF awards 1942123, 2211237, and 2224246 and a Jane Street Graduate Research Fellowship.

## References

- [BEK02] Nader H Bshouty, Nadav Eiron, and Eyal Kushilevitz. PAC learning with nasty noise. *Theoretical Computer Science*, 288(2):255–275, 2002. [1](#), [1](#), [3](#)
- [BFJ<sup>+</sup>94] Avrim Blum, Merrick Furst, Jeffrey Jackson, Michael Kearns, Yishay Mansour, and Steven Rudich. Weakly learning dnf and characterizing statistical query learning using fourier analysis. In *Proceedings of the twenty-sixth annual ACM symposium on Theory of computing*, pages 253–262, 1994. [2](#), [C.2](#)
- [BLMT22] Guy Blanc, Jane Lange, Ali Malik, and Li-Yang Tan. On the power of adaptivity in statistical adversaries. In *Conference on Learning Theory*, pages 5030–5061. PMLR, 2022. [\(document\)](#), [1](#), [1](#), [2.4](#), [1](#), [2](#), [3](#), [10.2](#), [10.2](#), [10.3](#), [1](#), [2](#), [A.2](#), [C](#), [C.1](#), [C.1](#), [C.1](#), [C.2](#)

- [BRS11] Boaz Barak, Prasad Raghavendra, and David Steurer. Rounding semidefinite programming hierarchies via global correlation. In *2011 IEEE 52nd Annual Symposium on Foundations of Computer Science*, pages 472–481. IEEE, 2011. [4.3](#)
- [Can22] Clément L Canonne. A short note on an inequality between kl and tv. *arXiv preprint arXiv:2202.07198*, 2022. [5](#), [5.5](#)
- [CHL<sup>+</sup>23] Clément Canonne, Samuel B Hopkins, Jerry Li, Allen Liu, and Shyam Narayanan. The full landscape of robust mean testing: Sharp separations between oblivious and adaptive contamination. In *2023 IEEE 64th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 2159–2168. IEEE, 2023. ([document](#)), [1](#), [1](#), [1](#), [2.3](#), [2.1](#), [2.4](#), [3.1](#), [20](#)
- [CSV17] Moses Charikar, Jacob Steinhardt, and Gregory Valiant. Learning from untrusted data. In *Proceedings of the 49th Annual Symposium on Theory of Computing (STOC)*, pages 47–60, 2017. [1](#)
- [DK23] Ilias Diakonikolas and Daniel M Kane. *Algorithmic high-dimensional robust statistics*. Cambridge university press, 2023. [1](#), [1](#), [3](#)
- [DKK<sup>+</sup>19] Ilias Diakonikolas, Gautam Kamath, Daniel Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. Robust estimators in high-dimensions without the computational intractability. *SIAM Journal on Computing*, 48(2):742–864, 2019. [1](#), [1](#), [8](#)
- [DKPZ21] Ilias Diakonikolas, Daniel M Kane, Thanasis Pittas, and Nikos Zarifis. The optimality of polynomial regression for agnostic learning under gaussian marginals in the sq model. In *Conference on Learning Theory*, pages 1552–1584. PMLR, 2021. [1](#), [3](#)
- [DSFT<sup>+</sup>14] Dana Dachman-Soled, Vitaly Feldman, Li-Yang Tan, Andrew Wan, and Karl Wimmer. Approximate resilience, monotonicity, and the complexity of agnostic learning. In *Proceedings of the twenty-sixth annual ACM-SIAM symposium on Discrete algorithms*, pages 498–511. SIAM, 2014. [1](#), [3](#)
- [Eld20] Ronen Eldan. Taming correlations through entropy-efficient measure decompositions with applications to mean-field approximation. *Probability Theory and Related Fields*, 176(3):737–755, 2020. [4.3](#)
- [Fel10] Vitaly Feldman. Distribution-specific agnostic boosting. *Innovations in Computer Science*, 2010. [3](#)
- [Fel17] Vitaly Feldman. A general characterization of the statistical query complexity. In *Conference on learning theory*, pages 785–830. PMLR, 2017. [2](#)
- [GR11] Oded Goldreich and Dana Ron. On testing expansion in bounded-degree graphs. *Studies in Complexity and Cryptography. Miscellanea on the Interplay between Randomness and Computation: In Collaboration with Lidor Avigad, Mihir Bellare, Zvika Brakerski, Shafi Goldwasser, Shai Halevi, Tali Kaufman, Leonid Levin, Noam Nisan, Dana Ron, Madhu Sudan, Luca Trevisan, Salil Vadhan, Avi Wigderson, David Zuckerman*, pages 68–75, 2011. [10.1](#)

- [Ham71] Frank R. Hampel. A General Qualitative Definition of Robustness. *The Annals of Mathematical Statistics*, 42(6):1887 – 1896, 1971. [1](#)
- [Hau92] David Haussler. Decision theoretic generalizations of the pac model for neural net and other learning applications. *Information and Computation*, 100(1):78–150, 1992. [1](#), [3](#)
- [HSSVG22] Daniel J Hsu, Clayton H Sanford, Rocco Servedio, and Emmanouil Vasileios Vlatakis-Gkaragkounis. Near-optimal statistical query lower bounds for agnostically learning intersections of halfspaces with gaussian marginals. In *Conference on Learning Theory*, pages 283–312. PMLR, 2022. [1](#)
- [Hub64] Peter Huber. Robust estimation of a location parameter. *The Annals of Mathematical Statistics*, 35(1), 1964. [1](#), [3](#), [A.2](#)
- [JKR19] Vishesh Jain, Frederic Koehler, and Andrej Risteski. Mean-field approximation, convex hierarchies, and the optimality of correlation rounding: a unified perspective. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pages 1226–1236, 2019. [4.3](#), [1](#)
- [Kea98] Michael Kearns. Efficient noise-tolerant learning from statistical queries. *Journal of the ACM (JACM)*, 45(6):983–1006, 1998. [C.2](#)
- [KK09] Varun Kanade and Adam Kalai. Potential-based agnostic boosting. *Advances in neural information processing systems*, 22, 2009. [3](#)
- [KKMS08] Adam Tauman Kalai, Adam R Klivans, Yishay Mansour, and Rocco A Servedio. Agnostically learning halfspaces. *SIAM Journal on Computing*, 37(6):1777–1805, 2008. [3](#)
- [KL93] Michael Kearns and Ming Li. Learning in the presence of malicious errors. *SIAM Journal on Computing*, 22(4):807–837, 1993. [1](#)
- [KSS94] Michael Kearns, Robert Schapire, and Linda Sellie. Toward efficient agnostic learning. *Machine Learning*, 17(2/3):115–141, 1994. [1](#), [3](#)
- [LRV16] Kevin A Lai, Anup B Rao, and Santosh Vempala. Agnostic estimation of mean and covariance. In *Proceedings of the 57th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 665–674, 2016. [1](#)
- [MR17] Pasin Manurangsi and Prasad Raghavendra. A birthday repetition theorem and complexity of approximating dense csps. In *44th International Colloquium on Automata, Languages, and Programming (ICALP 2017)*, pages 78–1. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2017. [4.3](#), [1](#)
- [Pin64] Mark S Pinsker. Information and information stability of random variables and processes. *Holden-Day*, 1964. [5](#), [5.5](#)
- [RT12] Prasad Raghavendra and Ning Tan. Approximating csps with global cardinality constraints using sdp hierarchies. In *Proceedings of the twenty-third annual ACM-SIAM symposium on Discrete Algorithms*, pages 373–387. SIAM, 2012. [4.3](#)

- [Tuk60] John Wilder Tukey. A survey of sampling from contaminated distributions. *Contributions to probability and statistics*, pages 448–485, 1960. [1](#)
- [Val85] Leslie G. Valiant. Learning disjunction of conjunctions. In *Proceedings of the 9th International Joint Conference on Artificial Intelligence (IJCAI)*, pages 560–566, 1985. [1](#), [1](#), [3.1](#), [19](#)
- [Ver18] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018. [7.1.1](#)
- [ZJS19] Banghua Zhu, Jiantao Jiao, and Jacob Steinhardt. Generalized resilience and robust statistics. *arXiv*, abs/1909.08755, 2019. [1](#), [8](#)

## A The subtractive and additive adversaries in our framework

We show how to fit subtractive and additive contamination into our framework.

### A.1 Subtractive adversaries

First, we formally define  $\eta$ -subtractive corruptions

**Definition 16.** For any distribution  $\mathcal{D}$  and  $\eta \in (0, 1)$ , we say that  $\mathcal{D}'$  is an  $\eta$ -subtractive contamination of  $\mathcal{D}$  if a sample from  $\mathcal{D}'$  is equivalent to a sample from  $\mathcal{D}$  conditioned on an event that occurs with probability at least  $1 - \eta$ . We use  $\text{sub}_\eta(\mathcal{D})$  to denote the set of all such  $\mathcal{D}'$ .

Similarly, for any  $S \in X^m$ , we say that  $S'$  is an  $\eta$ -subtractive contamination of  $S$  if it is formed by removing at most  $\lceil \eta \cdot m \rceil$  arbitrary points from  $S$ . In a slight overload of notation, we use  $\text{sub}_\eta(S)$  to denote the set of all such  $S'$ .

We will show, as an easy consequence of [Theorem 3](#), that the oblivious and adaptive variants of subtractive contamination are equivalent.

**Theorem 9** (Oblivious and adaptive subtractive contamination are equivalent). For any  $\eta, \varepsilon \in (0, 1)$ ,  $f : X^n \rightarrow \{0, 1\}$ , and distribution  $\mathcal{D}$  over  $X$ , let  $M = \text{poly}(n, 1/\varepsilon, 1/(1 - \eta))$ . Then,

$$\left| \sup_{\mathcal{D}' \in \text{sub}_\eta(\mathcal{D})} \left( \mathbb{E}_{\mathbf{S} \sim (\mathcal{D}')^n} [f(\mathbf{S})] \right) - \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^M} \left[ \sup_{\mathbf{S}' \in \text{sub}_\eta(\mathbf{S})} (\mathbb{E}[f \circ \Phi_{\star \rightarrow n}(\mathbf{S}')]) \right] \right| \leq \varepsilon.$$

[Theorem 9](#) is an easy consequence of [Theorem 6](#) as well as the below lemma.

**Lemma A.1** (Converting the subtractive adversary to our framework). For any  $\eta, \varepsilon \in (0, 1)$ ,  $f : X^n \rightarrow \{0, 1\}$ , and distribution  $\mathcal{D}$  over  $X$ , let

$$m := \left\lceil \frac{\max(2n, 8 \ln(1/\varepsilon))}{1 - \eta} \right\rceil, \tag{9}$$

and  $X' := X \cup \{\emptyset\}$ . There is a degree-2 cost function  $\rho$  and  $f' : (X')^m \rightarrow \{0, 1\}$  for which

$$\left| \sup_{\mathcal{D}' \in \text{sub}_\eta(\mathcal{D})} \left( \mathbb{E}_{\mathbf{S} \sim (\mathcal{D}')^n} [f(\mathbf{S})] \right) - \text{Oblivious-Max}_\rho(f', \mathcal{D}) \right| \leq \varepsilon,$$

and, for all  $M \geq m$ ,

$$\mathbb{E}_{\mathbf{S} \sim \mathcal{D}^M} \left[ \sup_{\mathbf{S}' \in \text{sub}_\eta(\mathbf{S})} (\mathbb{E}[f \circ \Phi_{\star \rightarrow n}(\mathbf{S}')]]) \right] = \text{Adaptive-Max}_\rho(f' \circ \Phi_{M \rightarrow m}, \mathcal{D})$$

*Proof of Theorem 9 assuming Lemma A.1.* Let  $f'$ ,  $\rho$ , and  $\mathcal{D}$  be as in Lemma A.1. By Theorem 6, for  $M = \text{poly}(m, 1/\varepsilon) = \text{poly}(n, 1/(1-\eta), \varepsilon)$ ,

$$|\text{Oblivious-Max}_\rho(f', \mathcal{D}') - \text{Adaptive-Max}_\rho(f' \circ \Phi_{M \rightarrow m}, \mathcal{D}')| \leq \varepsilon.$$

By the first part of Lemma A.1 and triangle inequality, we have that

$$\left| \sup_{\mathcal{D}' \in \text{sub}_\eta(\mathcal{D})} \left( \mathbb{E}_{\mathbf{S} \sim (\mathcal{D}')^n} [f(\mathbf{S})] \right) - \text{Adaptive-Max}_\rho(f' \circ \Phi_{M \rightarrow m}, \mathcal{D}') \right| \leq 2\varepsilon.$$

By the first part of Lemma A.1 this implies that

$$\left| \sup_{\mathcal{D}' \in \text{sub}_\eta(\mathcal{D})} \left( \mathbb{E}_{\mathbf{S} \sim (\mathcal{D}')^n} [f(\mathbf{S})] \right) - \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^M} \left[ \sup_{\mathbf{S}' \in \text{sub}_\eta(\mathbf{S})} (\mathbb{E}[f \circ \Phi_{\star \rightarrow n}(\mathbf{S}')]]) \right] \right| \leq 2\varepsilon.$$

The desired result holds by renaming  $\varepsilon' = \varepsilon/2$ . □

The proof of Lemma A.1 will use the following basic concentration inequality

**Proposition A.2.** *Let  $\mathcal{D}$  be a distribution on  $X' := X \cup \{\emptyset\}$  for which*

$$\Pr_{\mathbf{x} \sim \mathcal{D}'}[\mathbf{x} = \emptyset] \leq \eta.$$

For  $m$  as in Equation (9),

$$\Pr_{\mathbf{S} \sim (\mathcal{D}')^m} \left[ \sum_{x \in \mathbf{S}} \mathbf{1}[x \neq \emptyset] \leq n \right] \leq \varepsilon.$$

*Proof.* Let  $\mathbf{z}$  be the random variable that counts the number of entries  $\mathbf{S}'$  has that are not equal to  $\emptyset$ . Then,

$$\mu := \mathbb{E}[\mathbf{z}] \geq \frac{m}{1-\eta} \geq \max(2n, 8 \ln(1/\varepsilon)).$$

Furthermore,  $\mathbf{z}$  is the sum of independent random variables taking on values in  $\{0, 1\}$  (each indicating whether  $\mathbf{S}_i \neq \emptyset$  for some index  $i$ ). By a standard Chernoff bound (Fact 5.1),

$$\Pr[\mathbf{z} \leq n] \leq e^{-\mu/8} \leq \varepsilon. \quad \square$$

*Proof of Lemma A.1.* We begin by constructing the cost function. In the original subtractive adversary, for each  $x$ , the adversary can either choose to keep  $x$  in the sample or remove it at a cost of  $1/\eta$ . We will construct  $\rho$  so that the adversary has the same options and represent this “removal” option as converting an input  $x$  to  $\emptyset$ :

$$\rho(x, y) := \begin{cases} 0 & \text{if } x = y \\ \frac{1}{\eta} & \text{if } x \neq y \text{ and } y = \emptyset \\ \infty & \text{otherwise.} \end{cases} \quad (10)$$

The function  $f'$  simply runs  $f$  on a random subset of its non-null input. For any  $S \in X^m$ , let  $S_{\neq \emptyset}$  denote the subset of  $S$  consisting of all points not equal to  $\emptyset$ . Then,

$$f'(S) := \begin{cases} f(\Phi_{\star \rightarrow n}(S_{\neq \emptyset})) & \text{if } |S_{\neq \emptyset}| \geq n \\ 0 & \text{otherwise.} \end{cases}$$

Next, we analyze the oblivious adversaries. Consider a draw  $\mathbf{x} \sim \mathcal{D}$  coupled to an event  $\mathbf{E}$  occurring with probability at least  $1 - \delta$ . Then,  $\text{sub}_\eta(\mathcal{D})$  consists of all possible distributions of  $\mathbf{x}$  conditioned on  $\mathbf{E}$ , whereas, based on [Equation \(10\)](#),  $\mathcal{C}_\rho(\mathcal{D})$  consists of all the possible distributions of  $\mathbf{y}$  where

$$\mathbf{y} := \begin{cases} \mathbf{x} & \text{if } \mathbf{E} \\ \emptyset & \text{otherwise.} \end{cases}$$

For any such event  $\mathbf{E}$ , let  $\mathcal{D}'_1 \in \text{sub}_\eta(\mathcal{D})$  and  $\mathcal{D}'_2 \in \mathcal{C}_\rho(\mathcal{D})$  be the corresponding distribution. We will show that

$$\left| \mathbb{E}_{\mathbf{S}_1 \sim (\mathcal{D}'_1)^n} [f(\mathbf{S}_1)] - \mathbb{E}_{\mathbf{S}_2 \sim (\mathcal{D}'_2)^m} [f(\mathbf{S}_2)] \right| \leq \varepsilon.$$

Each element of  $\mathbf{S}_2$  is set to  $\emptyset$  with a probability that is at most  $\eta$  and otherwise has the same distribution as an element of  $\mathbf{S}_1$ . Therefore, the above difference is bounded by the probability that  $\mathbf{S}_2$  less than  $n$  non-null elements. This is at most  $\varepsilon$  by [Proposition A.2](#).

For the adaptive equivalence, consider any  $S \in X^M$ . Then any  $S_1 \in \text{sub}_\eta(S)$  is formed by removing at most  $\lfloor \eta M \rfloor$  of the points in  $S$ , whereas  $S_2 \in \mathcal{C}_\rho(S)$  is formed by setting  $\lfloor \eta M \rfloor$  of the points to  $\emptyset$ . Suppose we remove the same set of points to form  $S_1$  as we set to  $\emptyset$  to form  $S_2$ . Then, using the fact that at least  $n$  points must remain unchanged since  $M \geq m \geq n/(1 - \eta)$  and that the subsampling filter composes

$$\mathbb{E}[f \circ \Phi_{\star \rightarrow n}(S_1)] = \mathbb{E}[f' \circ \Phi_{M \rightarrow m}(S_2)].$$

Hence, for any choice of  $S \in X^M$ ,

$$\sup_{S_1 \in \text{sub}_\eta(S)} (\mathbb{E}[f \circ \Phi_{\star \rightarrow n}(S_1)]) = \sup_{S_2 \in \mathcal{C}_\rho(S)} (\mathbb{E}[f' \circ \Phi_{M \rightarrow m}(S_2)]).$$

This implies the desired result. □

## A.2 Additive adversaries

First, we formally define  $\eta$ -additive corruptions. Note that the oblivious adversary below exactly corresponds to Huber's original contamination model [\[Hub64\]](#).

**Definition 17** (Standard additive adversaries). *For any distribution  $\mathcal{D}$  and  $\eta \in (0, 1)$ , we say that  $\mathcal{D}'$  is an  $\eta$ -additive contamination of  $\mathcal{D}$  if, for some distribution  $\mathcal{E}$ ,*

$$\mathcal{D}' := (1 - \eta)\mathcal{D} + \eta\mathcal{E}$$

We use  $\text{add}_\eta(\mathcal{D})$  to denote the set of all such  $\mathcal{D}'$ , and for any function  $f : X^n \rightarrow \{0, 1\}$ , define

$$\text{Oblivious-Add-Max}_\eta(f, \mathcal{D}) := \sup_{\mathcal{D}' \in \text{add}_\eta(\mathcal{D})} \left\{ \mathbb{E}_{\mathbf{S}' \sim (\mathcal{D}')^n} (f(\mathbf{S}')) \right\}.$$

Similarly, for any  $S \in X^{\lceil(1-\eta)m\rceil}$ , we say  $S'$  is an  $\eta$ -additive contamination of  $S$  if it is formed by adding  $\lfloor\eta m\rfloor$  points to  $S$  and then arbitrarily permuting it. In a slight overload of notation, we use  $\text{add}_\eta(S)$  to denote the set of all such  $S'$ , and for any  $f : X^m \rightarrow \{0, 1\}$ , write

$$\text{Adaptive-Add-Max}_\eta(f, \mathcal{D}) := \mathbb{E}_{S \sim \mathcal{D}^{\lceil(1-\eta)m\rceil}} \left[ \sup_{S' \in \text{add}_\eta(S)} (\mathbb{E}[f(S')]) \right].$$

The equivalence between these two adversaries was already shown by [BLMT22], but we also prove it here to show our framework can recover their result.

**Theorem 10** (Oblivious and adaptive additive contamination are equivalent). *For any  $\eta, \varepsilon \in (0, 1)$ ,  $f : X^n \rightarrow \{0, 1\}$ , and distribution  $\mathcal{D}$  over  $X$ , let  $m = \text{poly}(n, 1/\varepsilon, \ln |X|)$ . Then,*

$$|\text{Oblivious-Add-Max}_\eta(f, \mathcal{D}) - \text{Adaptive-Add-Max}_\eta(f \circ \Phi_{m \rightarrow n}, \mathcal{D})| \leq \varepsilon.$$

To prove this equivalence, we will introduce a variant of the adaptive adversary, the *binomial* adversary. In this variant, rather than drawing *exactly*  $\lceil(1-\eta)m\rceil$  clean points and the adversary being able to add  $\lfloor\eta m\rfloor$  corrupted points, the number of clean points is itself drawn randomly from a binomial distribution.

**Definition 18** (Binomial adversary). *For any distribution  $\mathcal{D}$  and sample size  $m$ , the binomial adversary first draws  $z \sim \text{Bin}(m, (1-\eta))$  clean points from  $\mathcal{D}$ , then adds  $m - z$  arbitrary points to this clean sample, and finally permutes all  $m$  points arbitrarily. For any  $f : X^m \rightarrow \{0, 1\}$ , we define*

$$\text{Binomial-Max}_\eta(f, \mathcal{D}) := \mathbb{E}_{z \sim \text{Bin}(m, 1-\eta)} \left[ \mathbb{E}_{S \sim \mathcal{D}^z} \left[ \sup_{S' \in \text{complete}_m(S)} (\mathbb{E}[f(S')]) \right] \right]$$

where  $\text{complete}_m(S)$  to denote all samples that can be created by adding  $m - |S|$  points to  $S$ .

Our proof of [Theorem 10](#) proceeds in two steps. We first use [Theorem 2](#) to show that the oblivious additive adversary is equivalent to the binomial adversary, and then show equivalence between the binomial adversary and adaptive additive adversary.

**Proposition A.3** (The oblivious adversary and binomial adversary are equivalent). *For any  $\eta, \varepsilon \in (0, 1)$ ,  $f : X^n \rightarrow \{0, 1\}$ , and distribution  $\mathcal{D}$  over  $X$ , let  $m = \text{poly}(n, 1/\varepsilon, \ln |X|)$ . Then,*

$$|\text{Oblivious-Add-Max}_\eta(f, \mathcal{D}) - \text{Binomial-Max}_\eta(f \circ \Phi_{m \rightarrow n}, \mathcal{D})| \leq \varepsilon.$$

*Proof.* This will be a fairly straightforward application of [Theorem 2](#). Let  $X' := X \cup \{\emptyset\}$  and  $\mathcal{D}_\eta$  be the distribution that outputs  $\emptyset$  with probability  $\eta$  and otherwise outputs  $\mathcal{D}$ ,

$$\mathcal{D}_\eta := (1 - \eta) \cdot \mathcal{D} + \eta \cdot \{\emptyset\}.$$

Then, we'll define an adversary that can send  $\emptyset$  to any element of  $X'$  but otherwise cannot change its input.

$$\rho(x, y) := \begin{cases} 0 & \text{if } x = y \text{ or } x = \emptyset \\ \infty & \text{otherwise.} \end{cases}$$

Finally, let  $f' : X^n \rightarrow \{0, 1\}$  be defined as

$$f'(S) := \begin{cases} 0 & \text{if } \emptyset \in S \\ f(S) & \text{otherwise.} \end{cases}$$

We observe that,

$$\begin{aligned} \text{Oblivious-Add-Max}_\eta(f, \mathcal{D}) &= \text{Oblivious-Max}_\rho(f', \mathcal{D}_\eta), \quad \text{and} \\ \text{Binomial-Max}_\eta(f \circ \Phi_{m \rightarrow n}, \mathcal{D}) &= \text{Adaptive-Max}_\rho(f' \circ \Phi_{m \rightarrow n}, \mathcal{D}_\eta). \end{aligned}$$

because in order to maximize the success probability of  $f'$ , the adversaries should send every  $\emptyset$  they see to their adversarial choice of an element in  $X$ . The desired result follows from [Theorem 6](#).  $\square$

Next, we show that the binomial adversary and adaptive additive adversary are equivalent.

**Proposition A.4** (The binomial and adaptive additive adversaries are equivalent). *For any  $f : X^n \rightarrow \{0, 1\}$ ,  $\varepsilon, \eta \in (0, 1)$ , and distribution  $\mathcal{D}$ , let*

$$m = O\left(\frac{n^2}{\varepsilon^2}\right).$$

Then, for  $f' = f \circ \Phi_{m \rightarrow n}$

$$|\text{Adaptive-Add-Max}_\eta(f', \mathcal{D}) - \text{Binomial-Max}_\eta(f', \mathcal{D})| \leq \varepsilon.$$

*Proof.* Expanding the definitions, we wish to show that,

$$\left| \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^{\lceil (1-\eta)m \rceil}} \left[ \sup_{\mathbf{S}'_1 \in \text{add}_\eta(\mathbf{S})} (\mathbb{E}[f'(\mathbf{S}'_1)]) \right] - \mathbb{E}_{\mathbf{z} \sim \text{Bin}(m, 1-\eta)} \left[ \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^{\mathbf{z}}} \left[ \sup_{\mathbf{S}'_2 \in \text{complete}_m(\mathbf{S})} (\mathbb{E}[f'(\mathbf{S}'_2)]) \right] \right] \right| \leq \varepsilon.$$

Regardless of the strategy of one adversary, it is possible to choose a strategy for the other adversary so that the following holds: There is a coupling of  $\mathbf{S}'_1$  and  $\mathbf{S}'_2$  for which the expected number of differences between  $\mathbf{S}'_1$  and  $\mathbf{S}'_2$  is at most  $\mathbb{E}[|\mathbf{z} - \lceil (1-\eta)m \rceil|]$ . Furthermore, for any  $S_1, S_2$  differing in at most  $\Delta$  points and function  $f : X^n \rightarrow \{0, 1\}$ ,

$$|\mathbb{E}[f \circ \Phi_{m \rightarrow n}(S_1)] - \mathbb{E}[f \circ \Phi_{m \rightarrow n}(S_2)]| \leq \frac{n\Delta}{m} = \varepsilon/2,$$

because, in order for the two above quantities to differ, the subsample must select one of the  $\Delta$  differences. Therefore,

$$|\text{Adaptive-Add-Max}_\eta(f', \mathcal{D}) - \text{Binomial-Max}_\eta(f', \mathcal{D})| \leq \frac{n}{m} \cdot \mathbb{E}[|\mathbf{z} - \lceil (1-\eta)m \rceil|] \leq O\left(\frac{n}{\sqrt{m}}\right) \leq \varepsilon. \quad \square$$

Finally, we prove the main result of this section.

*Proof of Theorem 10.* Let  $f' = f \circ \Phi_{m \rightarrow n}$ . Then, by triangle inequality

$$\begin{aligned} & |\text{Oblivious-Add-Max}_\eta(f, \mathcal{D}) - \text{Adaptive-Add-Max}_\eta(f', \mathcal{D})| \\ & \leq |\text{Oblivious-Add-Max}_\eta(f, \mathcal{D}) - \text{Binomial-Max}_\eta(f', \mathcal{D})| \\ & \quad + |\text{Adaptive-Add-Max}_\eta(f', \mathcal{D}) - \text{Binomial-Max}_\eta(f', \mathcal{D})|. \end{aligned}$$

Each of the above terms is at most  $\varepsilon/2$  by [Propositions A.3](#) and [A.4](#).  $\square$

## B Partially-adaptive adversaries

In this section, we introduce two partially adaptive adversaries, *malicious noise* and the *non-iid* adversary, and prove that both are equivalent to fully adaptive and fully oblivious adversaries.

**Definition 19** (Malicious noise, [Val85]). *In malicious noise with base distribution  $\mathcal{D}$  and noise rate  $\eta$ , the sample  $\mathbf{S} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$  is generated sequentially. For each  $i \in [n]$ , an  $\eta$ -coin is flipped and then,*

1. *If the coin is tails, a clean point is sampled  $\mathbf{x}_i \sim \mathcal{D}$ .*
2. *If the coin is heads, the adversary gets to choose  $\mathbf{x}_i$  arbitrarily with full knowledge of  $\mathbf{x}_1, \dots, \mathbf{x}_{i-1}$  but no knowledge of the future points  $(\mathbf{x}_{i+1}, \dots, \mathbf{x}_n)$ .*

*For any  $f : X^n \rightarrow \{0, 1\}$  and distribution  $\mathcal{D}$ , we'll use  $\text{Mal-Max}_\eta(f, \mathcal{D})$  to denote the maximum expected value of  $f(\mathbf{S})$  over any  $\mathbf{S}$  generated by a malicious adversary with noise rate  $\eta$ .*

The malicious adversary is partially adaptive in the sense that, when it chooses how to corrupt  $\mathbf{x}_i$ , it only knows the points generated in the past.

**Definition 20** (The non-iid adversary, [CHL+23]). *In the non-iid adversary model with base distribution  $\mathcal{D}$  and noise rate  $\eta$ , to generate  $n$  samples, first the adversary arbitrarily chooses  $\lfloor \eta n \rfloor$  points, and then  $\lfloor (1 - \eta)n \rfloor$  points are generated iid from  $\mathcal{D}$ , added to the generated points, and permuted arbitrarily. For any  $f : X^n \rightarrow \{0, 1\}$  and distribution  $\mathcal{D}$ , we'll use  $\text{Non-iid-Max}_\eta(f, \mathcal{D})$  to denote the maximum expected value of  $f(\mathbf{S})$  over any  $\mathbf{S}$  generated by a non-iid adversary with noise rate  $\eta$ .*

This adversary is referred to as *non-iid* because the  $\lfloor \eta n \rfloor$  points can be generated arbitrarily and need not be iid. If they were, this adversary would be extremely similar to the oblivious additive adversary, with the only difference being that the non-iid adversary generates exactly  $\lfloor \eta n \rfloor$  corruptions, whereas the oblivious adversary generates  $\text{Bin}(n, \eta)$  corruptions.

**Theorem 11** (Equivalence of all additive adversaries). *For any  $\eta, \varepsilon \in (0, 1)$ ,  $f : X^n \rightarrow \{0, 1\}$ , and distribution  $\mathcal{D}$  over  $X$ , let  $m = \text{poly}(n, 1/\varepsilon, \ln |X|)$  and  $f' = f \circ \Phi_{m \rightarrow n}$ . The following are all within  $\pm \varepsilon$  of one another.*

1. *The maximum success probability of the oblivious additive adversary,*

$$\text{Oblivious-Add-Max}_\eta(f, \mathcal{D}).$$

2. *The maximum success probability of the adaptive additive adversary,*

$$\text{Adaptive-Add-Max}_\eta(f', \mathcal{D}).$$

3. *The maximum success probability of the malicious adversary,*

$$\text{Mal-Max}_\eta(f', \mathcal{D}).$$

4. *The maximum success probability of the non-iid adversary,*

$$\text{Non-iid-Max}_\eta(f', \mathcal{D}).$$

We already proved that Oblivious-Add-Max $_{\eta}(f, \mathcal{D})$  and Adaptive-Add-Max $_{\eta}(f', \mathcal{D})$  are within  $\pm\varepsilon$  of one another. To prove the same for malicious noise, we will show that malicious noise is no more powerful than the adaptive adversary, and at least as powerful as the oblivious adversary.

**Proposition B.1.** *In the setting [Theorem 11](#),*

$$\text{Mal-Max}_{\eta}(f', \mathcal{D}) \leq \text{Adaptive-Add-Max}_{\eta}(f', \mathcal{D}) + \varepsilon.$$

*Proof.* Consider an arbitrary malicious adversary. This adversary can make  $z$  many corruptions, where  $z \sim \text{Bin}(m, (1 - \eta))$ . Therefore, if the malicious adversary knew the full sample, it would be the same adversary as the binomial adversary. As a result, for any choices of the malicious adversary, there is a binomial adversary simulating it (generating the same distribution over corrupted samples). This gives that

$$\text{Mal-Max}_{\eta}(f', \mathcal{D}) \leq \text{Binomial-Max}_{\eta}(f', \mathcal{D}).$$

The desired result then follows from [Proposition A.4](#).  $\square$

**Proposition B.2.** *In the setting of [Theorem 11](#),*

$$\text{Oblivious-Add-Max}_{\eta}(f, \mathcal{D}) \leq \text{Mal-Max}_{\eta}(f', \mathcal{D}).$$

*Proof.* Consider any strategy for the oblivious adversary. It chooses an arbitrary distribution  $\mathcal{E}$  and sets

$$\mathcal{D}' = (1 - \eta) \cdot \mathcal{D} + \eta \cdot \mathcal{E}.$$

Now, consider the malicious adversary that, whenever it can corrupt a point, it draws a point from  $\mathcal{E}$  as its corruption. Then, each of the  $m$  points in this malicious adversary's sample is independent and drawn from  $\mathcal{D}'$ . After subsampling uniformly without replacement, the  $n$  points will still be independent and drawn from  $\mathcal{E}$ . This means that for any choices of the oblivious adversary, the malicious adversary can simulate them, giving the desired inequality.  $\square$

We execute the same two steps for the non-iid adversary.

**Proposition B.3.** *In the setting [Theorem 11](#),*

$$\text{Non-iid-Max}_{\eta}(f', \mathcal{D}) \leq \text{Adaptive-Add-Max}_{\eta}(f', \mathcal{D}).$$

*Proof.* Consider any strategy for the non-iid adversary. This is a set of points  $x_1, \dots, x_{\lfloor \eta m \rfloor}$  it will add to the sample. Now, consider the adaptive adversary that adds these same points regardless of what the clean points are. It's straightforward to see this adaptive adversary simulates the non-iid adversary, giving the desired inequality.  $\square$

**Proposition B.4.** *In the setting [Theorem 11](#),*

$$\text{Oblivious-Add-Max}_{\eta}(f, \mathcal{D}) \leq \text{Non-iid-Max}_{\eta}(f', \mathcal{D}) + \varepsilon.$$

*Proof.* Consider any strategy for the oblivious adversary. It chooses an arbitrary distribution  $\mathcal{E}$  and sets

$$\mathcal{D}' = (1 - \eta) \cdot \mathcal{D} + \eta \cdot \mathcal{E}.$$

To draw a sample  $\mathbf{S} \sim (\mathcal{D}')^n$ , we can first independently draw indicators  $\mathbf{a}_1, \dots, \mathbf{a}_n \stackrel{\text{iid}}{\sim} \text{Ber}(\eta)$ . For every  $i$  in which  $\mathbf{a}_i = 1$ ,  $\mathbf{S}_i$  is sampled from  $\mathcal{E}$ . In contrast, if  $\mathbf{a}_i = 0$ , then  $\mathbf{S}_i$  is sampled from  $\mathcal{D}$ .

Now consider the non-iid adversary that draws  $\mathbf{x}_1, \dots, \mathbf{x}_{\lfloor \eta m \rfloor} \stackrel{\text{iid}}{\sim} \mathcal{E}$  and adds them to the sample, and let  $\mathbf{T}$  be a size- $n$  subsample from the resulting non-iid adversaries subsample. Let  $\mathbf{b}_i$  be the indicator for whether the  $i^{\text{th}}$  element of  $\mathbf{T}$  comes from one of these  $\lfloor \eta m \rfloor$  points added. Then, we observe that, after conditioning on the values of  $\mathbf{b}_1, \dots, \mathbf{b}_n$ , each element of  $\mathbf{T}$  is independently drawn from  $\mathcal{E}$  if  $\mathbf{b}_i = 1$  and  $\mathcal{D}$  otherwise. We observe that the distribution of  $(\mathbf{b}_1, \dots, \mathbf{b}_n)$  is equivalent to the distribution obtained by first drawing  $\mathbf{i}_1, \dots, \mathbf{i}_n$  uniformly without replacement from  $[m]$  and then setting  $\mathbf{b}_i = \mathbb{1}[\mathbf{i}_i \leq \lfloor \eta m \rfloor]$ .

Therefore, the desired result follows from showing that the TV distance of the distributions of  $\mathbf{a}$  and  $\mathbf{b}$  is at most  $\varepsilon$ . We prove by exhibiting a coupling of  $\mathbf{a}$  and  $\mathbf{b}$  for which they differ with probability at most  $\varepsilon$ .

1. Draw  $\mathbf{z}_1, \dots, \mathbf{z}_n$  uniformly and independently  $[0, 1]$ .
2. Set  $\mathbf{a}_j = \mathbb{1}[\mathbf{z}_j \leq \eta]$  for each  $j \in [n]$ .
3. For each  $j \in [n]$ , let  $\mathbf{i}_n = \lfloor \mathbf{z}_j \cdot m \rfloor + 1$ . Note this gives that  $\mathbf{i}_1, \dots, \mathbf{i}_n$  are each uniform on  $[m]$  and they are independent.
4. If  $\mathbf{i}_1, \dots, \mathbf{i}_n$  are not unique, resample them by drawing them uniformly from  $[m]$  without replacement.
5. Set  $\mathbf{b}_j = \mathbb{1}[\mathbf{i}_j \leq \lfloor \eta m \rfloor]$ .

First, we confirm that the marginal distributions are correct. Each  $\mathbf{a}_j$  is independent and drawn from  $\text{Ber}(\eta)$ , so  $\mathbf{a}$  has the correct marginal distribution.

If we resample, then  $\mathbf{i}_1, \dots, \mathbf{i}_n$  are a uniform set of  $n$  distinct indices from  $[m]$ . If we don't resample, then they are also a uniform set of  $n$  distinct indices from  $[m]$ , because before resampling, they are independent and uniform from  $[m]$ , and we only don't resample if they are distinct. Therefore, the marginal distribution of  $\mathbf{b}$  is correct.

Finally, we bound the probability  $\mathbf{a} \neq \mathbf{b}$ . There are two ways that  $\mathbf{a}$  and  $\mathbf{b}$  could be different.

1. One the  $\mathbf{z}_i$  is between  $\frac{\lfloor \eta m \rfloor}{m}$  and  $\eta$ . This occurs with probability at most  $\frac{1}{m}$ .
2. We needed to resample  $\mathbf{i}_1, \dots, \mathbf{i}_n$  because they were not unique. This occurs if  $\mathbf{i}_j = \mathbf{i}_k$  for  $j \neq k$ . By union bound, it occurs with probability at most  $\frac{\binom{n}{2}}{m} \leq n^2/m$ .

Union bounding over the above two, we have that

$$\text{Oblivious-Add-Max}_\eta(f, \mathcal{D}) \leq \text{Non-iid-Max}_\eta(f', \mathcal{D}) + \frac{1}{m} + \frac{n^2}{m}. \quad \square$$

**Theorem 11** is immediate from **Propositions B.1** to **B.4** and **Theorem 10**.

## C Brief overview of [BLMT22]’s approaches and their limitations

### C.1 The special case of additive adversaries

[BLMT22] proved that oblivious and adaptive *additive* adversaries are equivalent (corresponding to [Theorem 10](#)). Here we briefly describe their proof strategy, and why it does not generalize to other adversary models. For simplicity, we set  $\eta = 1/2$  in the below exposition.

Recall that the adaptive additive adversary, given a sample  $S \in X^{M/2}$  can construct  $S \cup T$  for arbitrary  $T \in X^{M/2}$ . Using a standard concentration inequality, for any  $f : X^n \rightarrow \{0, 1\}$  fixed choice of the corruption  $T$ ,

$$\Pr_{S \sim \mathcal{D}^{(1-\eta) \cdot M}} [f \circ \Phi_{M \rightarrow n}(S \cup T) > \text{Oblivious-Max} + \varepsilon] \leq 2^{-\Omega_{n,\varepsilon}(M)} \quad (11)$$

where Oblivious-Max is appropriately defined for the setting. At first glance, the number of choices for  $T$  is  $|X|^{M/2}$ , which also grows exponentially in  $M$ . This makes it impossible to union bound over all choices of  $T$ . [BLMT22]’s key observation is that the space of possible corruptions (choices of  $T$ ) can be easily discretized to one that is much smaller.

In particular, given any  $T \in X^{M/2}$ , let  $\mathbf{T}$  be formed by,

1. First taking  $m \leq M$  samples  $\mathbf{x}_1, \dots, \mathbf{x}_m \stackrel{\text{iid}}{\sim} \text{Unif}(T)$ .
2. Then, for  $k := \frac{M}{2m}$ , constructing  $\mathbf{T}$  by taking  $k$  copies of each of  $\mathbf{x}_1, \dots, \mathbf{x}_m$ .

Then,  $\Phi_{M \rightarrow n}(T)$  and  $\Phi_{M \rightarrow n}(\mathbf{T})$  look identical unless a collision occurs (i.e. the same  $\mathbf{x}_i$  is sampled twice). If  $m = n^2/\varepsilon$ , that collision occurs with probability at most  $\varepsilon$ , which is negligible. Furthermore, there are only  $|X|^m$  possible choices for  $\mathbf{T}$ , which does not grow exponentially with  $M$ . Therefore, the desired result can proven using [Equation \(11\)](#) and an appropriate concentration inequality.

This strategy crucially relies on the fact that, while the adaptive additive adversary can choose its corruption (choice of  $T$ ) as a function of the clean sample  $S$ , the set of possible corruptions does not depend on  $S$ . Hence, adaptivity is inherently weaker for the additive adversary than for other models where the space of possible corruptions depends on the clean sample.

For example, consider the case of subtractive adversaries, where the adaptive adversary can remove half the points in the sample. For a sample  $S \in X^M$ , there are  $\approx 2^M$  ways to remove half the points, each parameterized by a bit string  $b \in \{0, 1\}^M$  where  $b_i$  indicates whether the  $i^{\text{th}}$  point is removed. Crucially, the “effect” of a bit string  $b$  depends on the clean sample  $S$  — in the sense that for the adversary to determine whether removing  $S_i$  is a good idea, it must know the value of  $S_i$ . In particular, if the adversary is only allowed to choose  $b$  from a subset  $B \subseteq \{0, 1\}^M$  of size much smaller than  $2^M$  that is fixed before seeing the clean sample  $S$ , the power of the adversary is greatly diminished. This makes it not clear how a similar discretization argument as [BLMT22] used for additive adversaries would work.

### C.2 The special case of statistical query algorithms

[BLMT22] also approved the equivalence between oblivious and adaptive adversaries for algorithms that never directly examine their dataset and only access it through *statistical queries* (SQ) [Kea98].

**Basics of the SQ framework.** A SQ is a pair  $(\varphi, \tau)$  where  $\varphi : X \rightarrow [0, 1]$  is the query and  $\tau > 0$  is the tolerance. For any distribution  $\mathcal{D}$ , a valid response to the query  $(\varphi, \tau)$  is any value that is

within  $\pm\tau$  of  $\mathbb{E}_{\mathbf{x}\sim\mathcal{D}}[\varphi(\mathbf{x})]$ . An SQ algorithm  $A$  using  $k$  queries of tolerance  $\tau$  specifies a sequence of  $k$  adaptively chosen queries,  $(\varphi_1, \tau), \dots, (\varphi_k, \tau)$ . For each  $t \in [k]$ , it receives a response  $v_t$  which is within  $\pm\tau$  of  $\mathbb{E}_{\mathbf{x}\sim\mathcal{D}}[\varphi_t(\mathbf{x})]$ , and the identity of  $\varphi_{t+1}$  is allowed to depend on the prior response  $v_1, \dots, v_t$ . After receiving all responses  $v_1, \dots, v_k$ , the  $A$  chooses an output  $y \in Y$ .

We say  $y \in Y$  is a valid output of  $A$  on distribution  $\mathcal{D}$  if it is a response that  $A$  can generate given valid responses  $v_1, \dots, v_k$  which are each within  $\pm\tau$  of  $\mathbb{E}_{\mathbf{x}\sim\mathcal{D}}[\varphi_t(\mathbf{x})]$ . We can now state [BLMT22]’s main result for SQ algorithms.

**Fact C.1** (Oblivious and adaptive adversaries are equivalent for the SQ framework). *Let  $A$  be any SQ algorithm making  $k$  queries of tolerance  $\tau$ , and  $\rho$  be any cost function. For  $m = \text{poly}(k, \tau)$ , there is an algorithm  $A' : X^m \rightarrow Y$  with the following guarantee: For any distribution  $\mathcal{D}$  over  $X^m$ , draw  $\mathbf{S} \sim \mathcal{D}^m$  and let an adversary choose  $\mathbf{S}' \in \mathcal{C}_\rho(\mathbf{S})$ . Then,  $A'(\mathbf{S}')$  is a valid output of  $A$  on distribution  $\mathcal{D}'$  for some  $\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})$  with high probability over the randomness of  $\mathbf{S}$ .*

To understand the utility of **Fact C.1**, suppose we have an SQ algorithm  $A$  that solves some task in the presence of an oblivious adversary. This means that, for all  $\mathcal{D}' \in \mathcal{C}_\rho(\mathcal{D})$ , any valid output of  $A$  on  $\mathcal{D}'$  is a good answer for this task. Then, **Fact C.1** says that  $A'$  given an adaptively corrupted sample will also provide a good answer for this task with high probability.

One straightforward weakness of this result compared to ours is not every task that admits an efficient solution also admits an efficient solution by an SQ algorithm [BFJ<sup>+</sup>94]. Even for tasks that can be cast into the SQ framework, our result has advantages.

1. The SQ equivalence in **Fact C.1** is not black box. Given an algorithm  $A$  not already in the SQ framework, in order to design an  $A'$  that defeats adaptive adversaries, first the algorithm designer must find an SQ algorithm that is “equivalent” to  $A$  in order to apply **Fact C.1**, a task that is not always trivial. In contrast, our result gives a black-box technique, via subsampling, to construct  $A'$ .
2. The SQ equivalence in **Fact C.1** does not have a well-defined sample overhead. Even if an algorithm  $A : X^n \rightarrow Y$  can be cast into some  $A_{\text{SQ}}$  operating in the SQ framework, the number of queries and tolerance  $A_{\text{SQ}}$  needs is not a predictable function of  $n$ . Therefore, it’s unclear how much larger the  $m$  in **Fact C.1** will be than  $n$ . In contrast, **Theorem 3** gives a simple expression for what  $m$  is needed as a function of  $n$ .