

On the power of adaptivity in statistical adversaries

Guy Blanc
Stanford University

GBLANC@STANFORD.EDU

Jane Lange
Massachusetts Institute of Technology

JLANGE@MIT.EDU

Ali Malik
Stanford University

MALIKALI@CS.STANFORD.EDU

Li-Yang Tan
Stanford University

LIYANG@CS.STANFORD.EDU

Editors: Po-Ling Loh and Maxim Raginsky

Abstract

We initiate the study of a fundamental question concerning adversarial noise models in statistical problems where the algorithm receives i.i.d. draws from a distribution \mathcal{D} . The definitions of these adversaries specify the *type* of allowable corruptions (noise model) as well as *when* these corruptions can be made (adaptivity); the latter differentiates between oblivious adversaries that can only corrupt the distribution \mathcal{D} and adaptive adversaries that can have their corruptions depend on the specific sample S that is drawn from \mathcal{D} .

We investigate whether oblivious adversaries are effectively equivalent to adaptive adversaries, across all noise models studied in the literature, under a unifying framework that we introduce. Specifically, can the behavior of an algorithm \mathcal{A} in the presence of oblivious adversaries always be well-approximated by that of an algorithm \mathcal{A}' in the presence of adaptive adversaries? Our first result shows that this is indeed the case for the broad class of *statistical query* algorithms, under all reasonable noise models. We then show that in the specific case of *additive noise*, this equivalence holds for *all* algorithms. Finally, we map out an approach towards proving this statement in its fullest generality, for all algorithms and under all reasonable noise models.

Keywords: Statistical problems, Adversary models, Robustness, Statistical query algorithms

1. Introduction

The possibility of noise pervades most problems in statistical estimation and learning. In this paper we will be concerned with *adversarial* noise models, as opposed to the class of more benign random noise models. Adversarial noise models are the subject of intensive study across statistics (Huber, 1964; Hampel, 1971; Tukey, 1975), learning theory (Valiant, 1985; Haussler, 1992; Kearns and Li, 1993; Kearns et al., 1994; Bshouty et al., 2002), and algorithms (Diakonikolas et al., 2019; Lai et al., 2016; Charikar et al., 2017; Diakonikolas and Kane, 2019). The definition of each model specifies:

1. The *type* of corruptions allowed. For example, the adversary may be allowed to add arbitrary points (additive noise (Huber, 1964; Valiant, 1985)), or in the context of supervised learning, allowed to change the labels in the data (agnostic noise (Haussler, 1992; Kearns et al., 1994)).
2. The *adaptivity* of the adversary.

The latter is the focus of our work. Consider any statistical problem where the algorithm is given i.i.d. draws from a distribution \mathcal{D} . On one hand we have *oblivious* adversaries: such an adversary corrupts \mathcal{D} to a different distribution $\widehat{\mathcal{D}}$, from which the algorithm then receives a sample. On the other hand we have *adaptive* adversaries: such an adversary first draws a sample S from \mathcal{D} , and upon seeing the specific outcomes in S , corrupts it to \widehat{S} which is then passed on to the algorithm. One can further consider adversaries with intermediate adaptive power, but we think of this as a dichotomy for now. A coupling argument shows that adaptive adversaries are at least as powerful as oblivious ones (Diakonikolas et al., 2019; Zhu et al., 2019). In this work we investigate whether they can be *strictly* more powerful.

Question 1 *Fix the type of corruptions allowed. Is it true that for any algorithm \mathcal{A} , there is an algorithm \mathcal{A}' whose behavior in the presence of adaptive adversaries well-approximates that of \mathcal{A} in the presence of oblivious adversaries?*

The distinction between oblivious and adaptive adversaries is frequently touched upon in works concerning statistical problems. Sometimes this distinction is brought up in service of emphasizing that the algorithms given in these works are robust against adaptive adversaries; other times it is brought up when the algorithms are shown to be robust against oblivious adversaries, and the viewpoint of an adaptive corruption process is provided as intuition for the noise model. However, the relative power of oblivious and adaptive adversaries in the statistical setting has not been systematically considered in the literature.

1.1. Our contributions

1.1.1. A UNIFIED FRAMEWORK FOR CHARACTERIZING DATA ADVERSARIES.

To reason generally about [Question 1](#), we associate every type of allowable corruptions with a cost function ρ between distributions. An oblivious ρ -adversary therefore corrupts \mathcal{D} to some $\widehat{\mathcal{D}}$ that is η -close with respect to ρ , meaning that $\rho(\mathcal{D}, \widehat{\mathcal{D}}) \leq \eta$. An adaptive ρ -adversary corrupts a sample S drawn from \mathcal{D} to some \widehat{S} such that the uniform distribution over \widehat{S} is η -close with respect to ρ to that over S . For example, when $\rho(\mathcal{D}, \widehat{\mathcal{D}})$ is the total variation distance between \mathcal{D} and $\widehat{\mathcal{D}}$, the resulting adaptive adversary represents *nasty noise* as defined in [Bshouty et al. \(2002\)](#), where the adversary is allowed to change an arbitrary η -fraction of the points in S . We discuss this framework in more detail in [Section 3](#).

1.1.2. YES TO [QUESTION 1](#) FOR ALL SQ ALGORITHMS

Our first result is an affirmative answer to [Question 1](#) for the broad class of *statistical query* (SQ) algorithms. Our proof will require a mild assumption on this cost function, the precise statement of which we defer to the body of the paper. For now, we simply refer to cost functions satisfying this assumption as “reasonable”, and mention that it is easily satisfied by all standard noise models, and can be seen to be necessary for our result to hold.

Theorem 1 (SQ algorithms are robust to adaptive adversaries) ¹ *For all reasonable cost functions ρ and SQ algorithms \mathcal{A} , the behavior of $\mathcal{A}' := \mathcal{A}$ in the presence of adaptive ρ -adversaries well-approximates that of \mathcal{A} in the presence of oblivious adversaries.*

1. See [Theorem 5](#) for the formal version of this theorem.

A key ingredient in our proof of [Theorem 1](#) is a novel reduction, using duality, from a k -query SQ algorithm to a single “representative” SQ.

In other words, the SQ framework neutralizes adaptive adversaries into oblivious ones, and in the context of [Question 1](#), we can take \mathcal{A}' to be \mathcal{A} itself. Looking ahead to our other results, we remark that such a statement cannot be true for all algorithms: there are trivial examples of (non-SQ) algorithms \mathcal{A} for which \mathcal{A}' has to be a modified version of \mathcal{A} .²

[Theorem 1](#) adds to the already-deep connections between the SQ framework and noise tolerance. The SQ framework was originally introduced in learning theory, where it continues to be influential in the design of learning algorithms that are resilient to *random classification noise* ([Kearns, 1998](#)). Most relevant to the topic of this paper, to our knowledge across all noise models, all existing statistical algorithms that have been shown to be robust to adversarial noise can be cast in the SQ framework.

1.1.3. YES TO [QUESTION 1](#) FOR ADDITIVE NOISE

Our second result is an affirmative answer to [Question 1](#) for one of the most natural types of corruptions:

Theorem 2³ *The answer to [Question 1](#) is “yes” for additive noise.*

In additive noise, an oblivious adversary corrupts \mathcal{D} to $\widehat{\mathcal{D}} = (1 - \eta)\mathcal{D} + \eta\mathcal{E}$ for an arbitrary distribution \mathcal{E} of their choosing. An adaptive adversary, on the other hand, gets to inspect the sample S drawn from \mathcal{D} , and adds $\frac{\eta}{1-\eta}|S|$ many arbitrary points of their choosing to S . The oblivious version of additive noise was introduced by Huber ([Huber, 1964](#)) and has become known as Huber’s contamination model; the adaptive version is commonly called “data poisoning” in the security and machine learning literature.

Additive noise also captures the well-studied *malicious noise* ([Definition 16](#)) from learning theory ([Valiant, 1985](#)) (see also ([Kearns and Li, 1993](#))). [Theorem 2](#) therefore shows that Huber’s contamination model, the malicious noise model, and the adaptive version of additive noise, are in fact all equivalent. Our proof of [Theorem 2](#) is constructive: we give an explicit description of how \mathcal{A}' can be obtained from \mathcal{A} , and \mathcal{A}' preserves the computational and sample efficiency of \mathcal{A} up to polynomial factors.

1.1.4. YES TO [QUESTION 1](#) IN ITS FULLEST GENERALITY? AN APPROACH VIA SUBSAMPLING

Our proof of [Theorem 2](#) is actually an instantiation of a broader approach towards answering [Question 1](#) affirmatively in its fullest generality: showing that the answer is “yes” for all (reasonable) types of allowable corruptions and all algorithms. We introduce the following definition:

Definition 1 (Neutralizing filter) *Let ρ be a cost function and \mathcal{A} be an algorithm for a statistical problem over a domain \mathcal{X} . We say that a randomized function $\Phi : \mathcal{X}^* \rightarrow \mathcal{X}^*$ is a neutralizing filter for \mathcal{A} with respect to ρ if the following holds. For all distributions \mathcal{D} , with high probability over the draw of S from \mathcal{D} , the behavior of*

2. One such example is $\mathcal{D} = \text{Bernoulli}(\frac{1}{2})$ and $\mathcal{A} = \mathbb{1}[\text{number of 1’s in the sample is } 0 \bmod 100]$. As the sample size grows, an oblivious adversary can barely change the acceptance probability of \mathcal{A} , whereas an adaptive one can change it completely.

3. See [Theorem 7](#) for the formal version

\mathcal{A} on $\Phi(\widehat{\mathcal{S}})$, where $\widehat{\mathcal{S}}$ is a corruption of \mathcal{S} by an adaptive ρ -adversary,

well-approximates the behavior of

\mathcal{A} on a sample from $\widehat{\mathcal{D}}$, where $\widehat{\mathcal{D}}$ is a corruption of \mathcal{D} by an oblivious ρ -adversary.

Perhaps the most natural filter in this context is the *subsampling* filter. For an n -sample algorithm \mathcal{A} , we request for a larger sample of size $m \geq n$, allow the adaptive adversary to corrupt it, and then run \mathcal{A} on a size- n subsample of the corrupted size- m sample. We call this the “ $m \rightarrow n$ subsampling filter” and denote it as $\Phi_{m \rightarrow n}$. The hope here is for the randomness of the subsampling step to neutralize the adaptivity of the adversary.

The efficiency of the subsampling filter is measured by the overhead in sample complexity that it incurs, i.e. how much larger m is relative to n . Subsampling from a sample that is roughly the size of the domain of course makes adaptive adversaries equivalent to oblivious ones, but this renders a sample-efficient algorithm inefficient. We are interested whether the subsampling filter can be effective while only incurring a mild overhead in sample complexity.

We propose the following conjecture as a general approach towards answering [Question 1](#):

Conjecture 1 *For all reasonable cost functions ρ and n -sample algorithms \mathcal{A} , the subsampling filter $\Phi_{m \rightarrow n}$ is a neutralizing filter for \mathcal{A} with respect to ρ with $m = \text{poly}(n, \log(|\mathcal{X}|))$.*

We obtain [Theorem 2](#) by proving [Conjecture 1](#) in the case of ρ being additive noise. While we have not been able to prove [Conjecture 1](#) for all ρ 's and all algorithms, we can show the following:

Theorem 3 (If it is possible, subsampling neutralizes adaptivity) ⁴ *Let ρ be a cost function and \mathcal{A} be an n -sample algorithm. Suppose there is an m -sample algorithm \mathcal{A}' whose behavior in the presence of adaptive ρ -adversaries well-approximates that of \mathcal{A} in the presence of oblivious ρ -adversaries. Then $\Phi_{M \rightarrow n}$ is a neutralizing filter for \mathcal{A} with respect to ρ with $M = O(m^2)$.*

We note that in the context of [Theorem 3](#), we do not require \mathcal{A}' to be computationally efficient: as long as \mathcal{A}' is sample efficient, then our resulting algorithm, the subsampling filter applied to \mathcal{A} , inherits the computational efficiency of \mathcal{A} .

Finally, we show that the bound on m in [Conjecture 1](#) cannot be further strengthened to be independent of $|\mathcal{X}|$, the size of the domain:

Theorem 4 (Subsampling lower bound) ⁵ *Let ρ be the cost function for additive noise and η be the corruption budget. There is an n -sample algorithm \mathcal{A} such that for $m \leq O_\eta(n \log(|\mathcal{X}|) / \log^2 n)$, $\Phi_{m \rightarrow n}$ is not a neutralizing filter for \mathcal{A} with respect to ρ .*

While [Theorem 4](#) shows that some dependence on $|\mathcal{X}|$ is necessary, we remark that the $\log(|\mathcal{X}|)$ dependence in [Conjecture 1](#) is fairly mild—this is the description length of a sample point $x \in \mathcal{X}$. [Theorem 4](#) also shows that the quantitative bounds that we achieve for [Theorem 2](#) has an optimal dependence on $|\mathcal{X}|$.

4. See [Theorem 6](#) for the formal version

5. See [Theorem 8](#) for the formal version

1.2. Other related work

A separation result Recently, Deng et al. gave a separation between the adaptive adversary, which they call “data-aware”, and the oblivious adversary (Deng et al., 2021). Specifically they showed that there are settings in which a natural Lasso-based algorithm for feature selection will often fail to select the correct features in the presence of an adaptive additive adversary, but would succeed in the presence of an oblivious additive adversary. Our Theorem 2 implies that this Lasso-based algorithm would succeed if it used the subsampling filter to preprocess its sample.

The online and dynamic setting While the focus of our work is on the statistical setting, the distinction between adaptive and oblivious adversaries has also been the subject of recent study in the *online* (Haghtalab et al., 2021; Alon et al., 2021) and *dynamic* (Beimel et al., 2022) setting, albeit with a notably different notion of adaptivity. In these settings, the adaptive adversaries can change the input distribution in response to the previous behavior of the algorithm, while oblivious adversaries must choose a fixed input distribution before the algorithms run.

Adaptive data analysis We emphasize the distinction between the focus of our work and the recent fruitful line of work on adaptive data analysis (Dwork et al. (2015); Hardt and Ullman (2014); Steinke and Ullman (2015); Bassily et al. (2021)). The focus of our work is on the adaptivity of the *adversary*, whereas the focus of this line of work is on the adaptivity of the *SQ algorithm*. Throughout this work, we reserve the use of “adaptive” to refer to the adversary, and all SQ algorithms will inherently be adaptive.

1.3. Discussion and future work

Implications of our results Theorem 1 says that for all reasonable cost functions, all SQ algorithms (existing and future ones) that are resilient to oblivious adversaries are “automatically” also resilient to adaptive adversaries. Likewise, lower bounds against adaptive adversaries immediately yield lower bounds against oblivious ones. The same remark further applies for all algorithms in the case of additive noise and malicious noise, by Theorem 2.

As a concrete example, we recall that the agnostic learning framework was originally defined with respect to oblivious adversaries (Haussler, 1992; Kearns et al., 1994). As in the PAC model there a concept class \mathcal{C} , but the target function f is no longer assumed to lie within \mathcal{C} —hence the name of the model. The learning algorithm is expected to achieve error close to opt , the distance from f to \mathcal{C} . However, many papers on agnostic learning provide the viewpoint of an adaptive corruption process as intuition for the model: the data *is* assumed to be labeled according to a function $f \in \mathcal{C}$, but an adversary corrupts an opt fraction of the labels given to the learning algorithm. This adaptive version was subsequently defined as a separate model called *nasty classification noise* (Bshouty et al., 2002) (as a special case of the nasty sample noise model introduced in that paper). Theorem 1 therefore shows that the agnostic learning model and the nasty classification noise model are in fact equivalent when it comes to SQ algorithms.

When can quantifiers be swapped? The distinction between oblivious and adaptive adversaries can be viewed as a difference in the order of “for all” and “with high probability” quantifiers in the performance guarantees of statistical algorithms. An algorithm succeeds in the presence of oblivious adversaries if *for all* distributions $\widehat{\mathcal{D}}$ that are close to \mathcal{D} , the algorithm succeeds *with high probability* over a sample S drawn from $\widehat{\mathcal{D}}$. On the other hand, an algorithm succeeds in the presence of adaptive

adversaries if *with high probability* over a sample S drawn from \mathcal{D} , the algorithm succeeds *for all* corruptions \widehat{S} that are close to S . Our work formalizes the question of when these quantifiers can be swapped, and our results provide several answers.

Future work In this work we initiate the systematic study of the power of adaptivity in statistical adversaries. A concrete direction for future work is to answer [Question 1](#) for other broad classes of algorithms and natural noise models, either via the subsampling filter ([Conjecture 1](#)) or otherwise. Here we highlight the specific case of subtractive noise: having resolved the case of additive noise in this work, doing so for subtractive noise as well would be a significant step towards resolving [Question 1](#) for all three generic noise models described in [Section 3](#).

2. Preliminaries

We use **boldface** (e.g. $\mathbf{x} \sim \mathcal{D}$) to denote random variables. Throughout this paper \mathcal{X} denotes an arbitrary finite domain, and we write $S \in \mathcal{X}^*$ to represent a multiset of elements in \mathcal{X} , meaning $S \in \mathcal{X}^0 \cup \mathcal{X}^1 \cup \mathcal{X}^2 \dots$. We use the notation $a = b \pm \varepsilon$ to indicate that $|a - b| < \varepsilon$. For any $m \in \mathbb{N}$, the notation $[m]$ indicates the set $\{1, 2, \dots, m\}$.

Distributions. For any $S \in \mathcal{X}^*$, we use $\mathcal{U}(S)$ to refer to the uniform distribution over S . For simplicity, we enforce that all distributions only have rational probabilities, meaning $\Pr_{\mathbf{x} \sim \mathcal{D}}[\mathbf{x} = x]$ is rational for any distribution \mathcal{D} and element $x \in \mathcal{X}$.⁶ For any distributions $\mathcal{D}_1, \mathcal{D}_2$ and parameter $\theta \in [0, 1]$, we use $\theta\mathcal{D}_1 + (1 - \theta)\mathcal{D}_2$ to refer to the mixture distribution which samples from \mathcal{D}_1 with probability θ and from \mathcal{D}_2 with probability $(1 - \theta)$.

Definition 2 (Total variation distance) *Let \mathcal{D} and \mathcal{D}' be any two distributions over the same domain, \mathcal{X} . It is well known that the following are equivalent definitions for the total variation distance between \mathcal{D} and \mathcal{D}' , denoted $\text{dist}_{\text{TV}}(\mathcal{D}, \mathcal{D}')$:*

1. *It is characterized by the best test distinguishing the two distributions:*

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

2. *It is characterized by the coupling which makes the two random variables different with the smallest probability.*

$$\text{dist}_{\text{TV}}(\mathcal{D}, \mathcal{D}') := \inf_{(\mathbf{x}, \mathbf{x}') \text{ a coupling of } \mathcal{D}, \mathcal{D}'} \left\{ \Pr[\mathbf{x} \neq \mathbf{x}'] \right\}.$$

3. Adversarial noise models

To reason generally about various noise adversarial models, we represent the types of allowable corruptions by a budget η and cost function ρ , which maps each ordered pair of distributions to some non-negative cost (or infinity). The cost need not be symmetrical.

6. Alternatively, one could enforce that the cost function smoothly interpolate to irrational probabilities, which is the case for all standard noise models.

Definition 3 ((ρ, η) -oblivious adversary) *Given some cost function ρ and budget η , an algorithm operating in the oblivious adversary model will receive the following input: If the true data distribution is \mathcal{D} , then the algorithm will receive iid samples from an adversarial chosen $\widehat{\mathcal{D}}$ satisfying $\rho(\mathcal{D}, \widehat{\mathcal{D}}) \leq \eta$.*

Definition 4 ((ρ, η) -adaptive adversary) *Given some cost function ρ and budget η , an n -sample algorithm operating in adaptive adversary model will receive the following input: If the true data distribution is \mathcal{D} , first a clean sample $\mathcal{S} \sim \mathcal{D}^n$ is generated, and then the algorithm will get an adversarially chosen $\widehat{\mathcal{S}}$ satisfying $\rho(\mathcal{U}(\mathcal{S}), \mathcal{U}(\widehat{\mathcal{S}})) \leq \eta$.*

Throughout this paper, we use $\widehat{\square}$ to denote the corrupted version of a set or distribution.

We remark that the adaptive adversary as defined in [Definition 4](#) is slightly stronger than the definitions usually considered, in the sense that there is usually a bound on the size of $\widehat{\mathcal{S}}$ the adversary is allowed to produce. For example, in the nasty noise model ([Definition 7](#)), the adversary is only allowed to change points in the sample, and so $|\widehat{\mathcal{S}}| = |\mathcal{S}|$. All of our results apply regardless of the size of $\widehat{\mathcal{S}}$.

3.1. Standard noise models from the literature

In this subsection, we present three generic adversary models and show how they are special cases of our framework with an appropriate choice of cost function. Other standard models, and how they fit within our framework, are given in [Appendix A](#).

Definition 5 (Additive noise) *Given a size- n sample $S \in \mathcal{X}^n$ and a corruption budget η , the adaptive additive noise adversary is allowed to add $\lfloor n \cdot \eta / (1 - \eta) \rfloor$ points to S arbitrarily.*

The additive noise model is captured by the cost function:

$$\text{cost}_{\text{add}}(\mathcal{D}, \widehat{\mathcal{D}}) := \inf_{\eta \in \mathbb{R}_{\geq 0}} \left\{ \widehat{\mathcal{D}} = (1 - \eta)\mathcal{D} + \eta\mathcal{E} \right\} \quad \text{for some distribution } \mathcal{E}.$$

The oblivious version of additive noise is the well-known Huber contamination model ([Huber, 1964](#)).

Definition 6 (Subtractive noise) *Given a size- n sample $S \in \mathcal{X}^n$ and a corruption budget η , the adaptive subtractive noise adversary is allowed to remove $\lfloor \eta n \rfloor$ points from S arbitrarily.*

The subtractive noise adversary is captured by the cost function:

$$\text{cost}_{\text{sub}}(\mathcal{D}, \widehat{\mathcal{D}}) := \text{cost}_{\text{add}}(\widehat{\mathcal{D}}, \mathcal{D}).$$

Definition 7 (Nasty noise) *Given a size- n sample $S \in \mathcal{X}^n$ and a corruption budget η , the adaptive nasty noise adversary is allowed to change up to $\lfloor \eta n \rfloor$ of the points in S arbitrarily.*

This noise model is also known as *strong contamination*, and as *nasty sample noise* (or simply nasty noise) in the context of supervised learning ([Bshouty et al., 2002](#)). It is captured by the cost function $\rho = \text{dist}_{\text{TV}}$.

4. Overview of [Theorem 1](#): The SQ framework neutralizes adaptive adversaries

In this section we formally state and prove [Theorem 1](#)—namely that the behaviour of a Statistical Query (SQ) algorithm in the presence of adaptive adversaries is equivalent to its behaviour in the presence of oblivious adversaries.

Basics of the SQ framework. Let \mathcal{D} be a distribution over some domain \mathcal{X} . A statistical query is a pair (ϕ, τ) where $\phi : \mathcal{X} \rightarrow [-1, 1]$ is the query and $\tau > 0$ is a tolerance parameter. These queries can be answered by a statistical query oracle $\text{STAT}_{\mathcal{D}}$ which, given an SQ (ϕ, τ) , returns a value v equal to $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\phi(\mathbf{x})]$ up to an additive error of τ i.e. $v \in \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\phi(\mathbf{x})] \pm \tau$.

Throughout this section, we will use some convenient shorthand. For any distribution \mathcal{D} and multiset $S \in \mathcal{X}^*$, we write

$$\phi(\mathcal{D}) := \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\phi(\mathbf{x})] \quad \text{and} \quad \phi(S) := \phi(\mathcal{U}(S)) = \frac{1}{|S|} \sum_{x \in S} \phi(x).$$

A k -query statistical query algorithm, \mathcal{A} , is an algorithm that makes a sequence of statistical queries to the oracle $\text{STAT}_{\mathcal{D}}$ one by one, using the result of the previous queries to decide which statistical query to make next.

Definition 8 (k -query SQ Algorithm) A k -query SQ algorithm \mathcal{A} is a sequence of k SQs

$$(\phi^{(1)}, \tau^{(1)}), (\phi_{v_1}^{(2)}, \tau^{(2)}), (\phi_{v_1, v_2}^{(3)}, \tau^{(3)}), \dots, (\phi_{v_1, \dots, v_{k-1}}^{(k)}, \tau^{(k)})$$

to $\text{STAT}_{\mathcal{D}}$, where \mathcal{A} 's choice of the $(i+1)$ -st SQ can depend on v_1, \dots, v_i which are the answers of $\text{STAT}_{\mathcal{D}}$ to the previous i SQs. For notational simplicity, we make the standard assumption that all the τ 's are the same.

Using mechanisms to implement $\text{STAT}_{\mathcal{D}}$. The SQ framework is a stylized model that cleanly facilitates theoretical analyses; It allows the algorithm designer to abstract away an algorithm's interaction with a random sample and instead *assume* $\phi(\mathcal{D})$ can be accessed up to $\pm\tau$ accuracy.

For an SQ algorithm to be useful, the $\text{STAT}_{\mathcal{D}}$ oracle must be implemented. This is done by a *mechanism* which uses a random sample $\mathcal{S} \sim \mathcal{D}^n$ to simulate $\text{STAT}_{\mathcal{D}}$ with high probability. The interaction between a mechanism and SQ algorithm is depicted in [Figure 1](#).

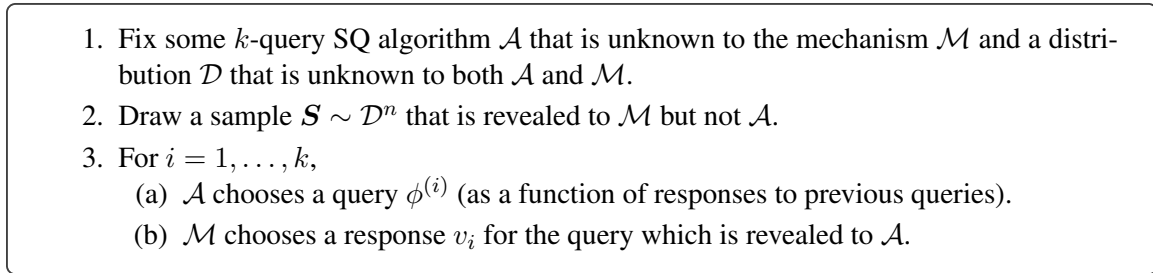


Figure 1: The interaction between a mechanism \mathcal{M} and SQ algorithm \mathcal{A} .

Definition 9 ((τ, δ) -accurate mechanisms) A mechanism \mathcal{M} is (τ, δ) -accurate for k -query SQ algorithms, if for any distribution \mathcal{D} and SQ algorithm \mathcal{A} , with probability at least $(1 - \delta)$ over the

randomness of \mathcal{S} and \mathcal{M} ,

$$v_i = \phi^{(i)}(\mathcal{D}) \pm \tau \quad \text{for all } i = 1, \dots, k$$

where v_i and $\phi^{(i)}$ are defined as in [Figure 1](#).

In this work, we focus on the τ -rounding mechanism.

Definition 10 (τ -rounding mechanism) *Given a sample $S \in \mathcal{X}^n$ and query ϕ , the τ -rounding mechanism, denoted \mathcal{M}_τ , returns the answer $v = \text{round}(\phi(S), \tau)$ where $\text{round}(x, \tau)$ refers to x rounded to the nearest integer multiple of τ .*

Fact 1 (The τ -rounding mechanism is accurate) *For any $k \in \mathbb{N}$ and $\delta, \tau > 0$, the τ -rounding mechanism with a sample size of*

$$n = O\left(\frac{k \log(1/\tau) + \log(1/\delta)}{\tau^2}\right)$$

is (τ, δ) accurate for k -query SQ algorithms.

Proof For each query, the \mathcal{M}_τ can return one of only $O(1/\tau)$ possible values (after rounding). The i^{th} query is chosen as a function of v_1, \dots, v_{i-1} , so there are at most $O(1/\tau)^{i-1}$ possible choices for the i^{th} query, and only $O(1/\tau)^k$ total unique queries \mathcal{A} could choose. Using a Chernoff bound and union bound over all possible queries, the probability \mathcal{A} asks a query, ϕ , where $|\phi(\mathcal{S}) - \phi(\mathcal{D})| \geq \tau/2$ is at most

$$\exp(-\Omega(\tau^2 n) + O(k \log(1/\tau)))$$

For the n given in [Fact 1](#), that probability is at most δ . The desired result follows from triangle inequality and $|\text{round}(\phi(\mathcal{S}), \tau) - \phi(\mathcal{S})| \leq \tau/2$. \blacksquare

The existence of accurate mechanisms (as in [Fact 1](#)) is the key to the SQ framework: SQ algorithms can assume that they have access to a $\text{STAT}_{\mathcal{D}}$ oracle, because for modest sample sizes and tiny failure probabilities, mechanisms are a $\text{STAT}_{\mathcal{D}}$ oracle.

4.1. The SQ framework in the presence of adversarial noise.

The SQ framework naturally extends to oblivious adversaries.

Definition 11 (k -query SQ Algorithm with an oblivious adversary) *Fix a cost function ρ and budget η . An k -query SQ algorithm, \mathcal{A} , in the presence of a (ρ, η) -oblivious adversary is a sequence of k SQs*

$$(\phi^{(1)}, \tau), (\phi_{v_1}^{(2)}, \tau), (\phi_{v_1, v_2}^{(3)}, \tau), \dots, (\phi_{v_1, \dots, v_{k-1}}^{(k)}, \tau)$$

each of which are answered according to $\text{STAT}_{\widehat{\mathcal{D}}}$ for some $\widehat{\mathcal{D}}$ satisfying $\rho(\mathcal{D}, \widehat{\mathcal{D}}) \leq \eta$. This $\widehat{\mathcal{D}}$ is the same for all queries, but adversarially chosen. \mathcal{A} 's choice of the $(i+1)$ -st SQ can depend on v_1, \dots, v_i which are the answers of $\text{STAT}_{\widehat{\mathcal{D}}}$ to the previous i SQs.

Once again, to run an SQ algorithm, the $\text{STAT}_{\widehat{\mathcal{D}}}$ oracle is implemented by a mechanism. Given a sample $\mathcal{S} \sim \widehat{\mathcal{D}}$, any mechanism satisfying [Definition 9](#) will be able to simulate the $\text{STAT}_{\widehat{\mathcal{D}}}$ oracle.

In the presence of (ρ, η) -adaptive adversaries, first a clean sample $\mathcal{S} \sim \mathcal{D}^n$ is drawn, and then the adversary chooses an η -corruption $\widehat{\mathcal{S}}$ that is passed to the mechanism, as shown in [Figure 2](#).

1. Fix some k -query SQ algorithm \mathcal{A} that is unknown to the mechanism \mathcal{M} and a distribution \mathcal{D} that is unknown to both \mathcal{A} and \mathcal{M} .
2. Draw a sample $\mathcal{S} \sim \mathcal{D}^n$ that is revealed to neither \mathcal{A} or \mathcal{M} .
3. An adversary chooses an $\widehat{\mathcal{S}}$ that is η -close to \mathcal{S} which is revealed to \mathcal{M} but not \mathcal{A} .
4. For $i = 1, \dots, k$,
 - (a) \mathcal{A} chooses a query $\phi^{(i)}$ (as a function of responses to previous queries).
 - (b) \mathcal{M} chooses a response v_i for the query which is revealed to \mathcal{A} .

Figure 2: The interaction between a mechanism \mathcal{M} and SQ algorithm \mathcal{A} in the presence of an adaptive adversary.

Our goal is to show that there are mechanisms that can simulate $\text{STAT}_{\widehat{\mathcal{D}}}$ for some $\widehat{\mathcal{D}}$ η -close to \mathcal{D} given just the corrupted sample $\widehat{\mathcal{S}}$.

Definition 12 ((τ, δ) -accurate in the presence of adaptive noise) *Fix a cost function ρ , budget η . A mechanism \mathcal{M} is said to be (τ, δ) -accurate for k -query SQ algorithms in the presence of adaptive noise, if for any distribution \mathcal{D} and SQ algorithm \mathcal{A} , the following holds. With probability at least $1 - \delta$ over the randomness of $\mathcal{S} \sim \mathcal{D}^n$,*

$$v_i = \phi^{(i)}(\widehat{\mathcal{D}}) \pm \tau \quad \text{for all } i = 1, \dots, k$$

for some $\widehat{\mathcal{D}}$ η -close to \mathcal{D} , where v_i and $\phi^{(i)}$ are defined as in [Figure 2](#). In particular, this holds regardless of how the adversary chooses $\widehat{\mathcal{S}}$.

If \mathcal{A} succeeds given $\text{STAT}_{\widehat{\mathcal{D}}}$ for every distribution $\widehat{\mathcal{D}}$ that is η -close to \mathcal{D} (i.e. \mathcal{A} is resilient to oblivious adversaries), then, with high probability, \mathcal{A} also succeeds in the presence of an adaptive adversary when using a mechanism satisfying [Definition 12](#). We will show that the rounding mechanism meets [Definition 12](#) whenever ρ is “reasonable” in the sense of the following definition; this property is easily satisfied by all standard cost functions.

Definition 13 (Closed under mixtures) *We say that ρ is closed under mixtures if for any distributions $\mathcal{D}_1, \mathcal{D}_2, \widehat{\mathcal{D}}_1, \widehat{\mathcal{D}}_2$ and $\theta \in (0, 1)$,*

$$\rho(\theta\mathcal{D}_1 + (1 - \theta)\mathcal{D}_2, \theta\widehat{\mathcal{D}}_1 + (1 - \theta)\widehat{\mathcal{D}}_2) \leq \max(\rho(\mathcal{D}_1, \widehat{\mathcal{D}}_1), \rho(\mathcal{D}_2, \widehat{\mathcal{D}}_2)).$$

Requiring that ρ is closed under mixtures enforces that the adaptive and oblivious adversaries “match up” in the sense of making the same types of changes. This is formalized in the following fact, for which we provide a short proof in [Appendix F](#).

Fact 2 Let ρ be closed under mixtures, \mathcal{D} be a distribution over \mathcal{X} , η be a corruption budget, and $n \in \mathbb{N}$. Suppose that for all $S \in \mathcal{X}^n$, there is a corresponding \widehat{S} satisfying $\rho(\mathcal{U}(S), \mathcal{U}(\widehat{S})) \leq \eta$. Let $\widehat{\mathcal{D}}$ be the distribution where $\mathbf{x} \sim \widehat{\mathcal{D}}$ is generated by: 1) Drawing $\mathbf{S} \sim \mathcal{D}^n$ and 2), drawing $\mathbf{x} \sim \mathcal{U}(\widehat{\mathbf{S}})$. Then, $\rho(\mathcal{D}, \widehat{\mathcal{D}}) \leq \eta$.

Our quantitative bounds will depend on a parameter that is related to the types of corruption the adversaries can make. Suppose that the adaptive adversary is required to keep the size of the corrupted sample the same as the clean sample ($|\widehat{S}| = |S|$). In this case, we require that if two samples S_1 and S_2 differ in only one point, then for any \widehat{S}_1 that is η -close to S_1 , there is some \widehat{S}_2 that is η -close to S_2 where \widehat{S}_1 and \widehat{S}_2 differ in only a small number of points. The following definition generalizes that notion to the case where the adversary can also change the number of points in the sample.

Definition 14 (ℓ -local) For any $\ell > 0$, a cost function ρ with budget η is ℓ -local if for any distributions $\mathcal{D}_1, \mathcal{D}_2$ and η -corruption $\widehat{\mathcal{D}}_1$ of \mathcal{D}_1 , there is some η -corruption $\widehat{\mathcal{D}}_2$ of \mathcal{D}_2 satisfying $\text{dist}_{\text{TV}}(\widehat{\mathcal{D}}_1, \widehat{\mathcal{D}}_2) \leq \ell \cdot \text{dist}_{\text{TV}}(\mathcal{D}_1, \mathcal{D}_2)$

All of the adversary models in [Section 3](#) are 1-local with the exception of the η -subtractive noise, which is $\frac{1}{1-\eta}$ -local. We encourage the reader to think of ℓ as a constant.

We are now ready to state the formal version of [Theorem 1](#), which generalizes [Fact 1](#) to the setting of adversarial noise.

Theorem 5 (Formal version of [Theorem 1](#)) For any ℓ -local cost function, adversary budget η , $\delta, \tau > 0$, and $k \in \mathbb{N}$, the τ -rounding mechanism with a sample size of

$$n = O\left(\frac{\ell^2(k \log(1/\tau) + \log(1/\delta))}{\tau^2}\right)$$

is (τ, δ) accurate for k -query SQ algorithm in the presence of adaptive noise.

Proof sketch. Here, we prove [Theorem 5](#) contingent on [Lemmas 1](#) and [2](#), which we prove in [Appendix B](#).

First, we will prove the special case where \mathcal{A} makes only a single query.⁷

Lemma 1 (Theorem 5 in the case of a single SQ) Let $\Psi : \mathcal{X} \rightarrow [-1, 1]$ be a statistical query, $T \in [-1, 1]$, and suppose:

$$\Psi(\widehat{\mathcal{D}}) \leq T \quad \text{for all } \widehat{\mathcal{D}} \text{ that are } \eta\text{-close to } \mathcal{D}.$$

Then, for any $\tau > 0$ and sample size $n \in \mathbb{N}$, the probability over $\mathbf{S} \sim \mathcal{D}^n$ that there is some $\widehat{\mathbf{S}}$ that is η -close to \mathbf{S} satisfying

$$\Psi(\widehat{\mathbf{S}}) \geq T + \frac{\tau}{2}$$

is at most $\exp\left(-\frac{\tau^2 n}{8\ell^2}\right)$.

7. To prove [Theorem 5](#) in the case of a single statistical query, we could apply [Lemma 1](#) twice: Once to bound how large the adversary can make $\Psi(\widehat{\mathbf{S}})$ and once to bound how small it can make $\Psi(\widehat{\mathbf{S}})$. Instead, we will directly apply [Lemma 1](#) to prove the multi-query case of [Theorem 5](#).

In order to prove [Theorem 5](#), we want to bound the probability that for a random $\mathcal{S} \sim \mathcal{D}^n$, there is some corruption $\widehat{\mathcal{S}}$ for which \mathcal{M}_τ is not τ -accurate. In detail, that means, for

$$\begin{aligned} v_1 &:= \text{round}(\phi^{(1)}(\widehat{\mathcal{S}}), \tau) \\ v_{i+1} &:= \text{round}(\phi_{v_1, \dots, v_i}^{(i+1)}(\widehat{\mathcal{S}}), \tau) \quad \text{for } i \in \{0, 1, \dots, k-1\}, \end{aligned} \tag{1}$$

and, for every $\widehat{\mathcal{D}}$ that is η -close to \mathcal{D} , there is some i for which $|v_i - \phi_{v_1, \dots, v_{i-1}}^{(i)}(\widehat{\mathcal{D}})| > \tau$. Consider a single possible choice for v_1, \dots, v_k (which also fixes the k -statistical queries, $\phi^{(1)}, \dots, \phi^{(k)}$). We use the separating hyperplane theorem to reduce to the case of a single statistical query.

Lemma 2 (*k*-query SQ algorithm to a single SQ) *Fix any k statistical queries $\phi^{(1)}, \dots, \phi^{(k)} : \mathcal{X} \rightarrow [-1, 1]$ and k values v_1, \dots, v_k . Suppose that there is no $\widehat{\mathcal{D}}$ that is η -close to \mathcal{D} satisfying*

$$\phi^{(i)}(\widehat{\mathcal{D}}) \in v_i \pm \tau \quad \text{for every } i \in [k].$$

Then there exists a single statistical query $\Psi : \mathcal{X} \rightarrow [-1, 1]$ and threshold T with the following properties.

1. $\Psi(\widehat{\mathcal{D}}) \leq T$ for every $\widehat{\mathcal{D}}$ that is η -close to \mathcal{D} .
2. For any sample $\widehat{\mathcal{S}}$ satisfying $\phi^{(i)}(\widehat{\mathcal{S}}) \in v_i \pm \frac{\tau}{2}$ for each $i \in [k]$, it is also true that $\Psi(\widehat{\mathcal{S}}) \geq T + \frac{\tau}{2}$.

Applying [Lemmas 1](#) and [2](#), for any fixed choice of v_1, \dots, v_k , the probability there is some $\widehat{\mathcal{S}}$ η -close to a sample $\mathcal{S} \sim \mathcal{D}^n$ satisfying [Equation \(1\)](#), and for every $\widehat{\mathcal{D}}$ that is η -close to \mathcal{D} , there is some i for which $|v_i - \phi_{v_1, \dots, v_{i-1}}^{(i)}(\widehat{\mathcal{D}})| > \tau$ is only $\exp\left(-\frac{\tau^2 n}{8\ell^2}\right)$. Since each $v_i \in [-1, 1]$ is an integer multiple of τ , there are at most $(\frac{2}{\tau} + 1)^k$ many choices for (v_1, \dots, v_k) . A union bound over all these choices completes the proof of [Theorem 5](#).

Acknowledgments

We thank Adam Klivans and Greg Valiant for helpful conversations. We are grateful to Greg for allowing us to include [Theorem 4](#) which was proved jointly with him.

Guy and Li-Yang are supported by NSF CAREER Award 1942123. Jane is supported by NSF Award CCF-2006664. Ali is supported by a graduate fellowship award from Knight-Hennessy Scholars at Stanford University.

References

- Noga Alon, Omri Ben-Eliezer, Yuval Dagan, Shay Moran, Moni Naor, and Eylon Yogev. Adversarial laws of large numbers and optimal regret in online classification. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, pages 447–455, 2021.
- Raef Bassily, Kobbi Nissim, Adam Smith, Thomas Steinke, Uri Stemmer, and Jonathan Ullman. Algorithmic stability for adaptive data analysis. *SIAM Journal on Computing*, 2021.

- Amos Beimel, Haim Kaplan, Yishay Mansour, Kobbi Nissim, Thatchaphol Saranurak, and Uri Stemmer. Dynamic algorithms against an adaptive adversary: Generic constructions and lower bounds. In *Proceedings of the 54rd Annual ACM SIGACT Symposium on Theory of Computing*, 2022.
- Nader H Bshouty, Nadav Eiron, and Eyal Kushilevitz. PAC learning with nasty noise. *Theoretical Computer Science*, 288(2):255–275, 2002.
- 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.
- Samuel Deng, Sanjam Garg, Somesh Jha, Saeed Mahloujifar, Mohammad Mahmoody, and Abhradeep Guha Thakurta. A separation result between data-oblivious and data-aware poisoning attacks. *Advances in Neural Information Processing Systems*, 34, 2021.
- Ilias Diakonikolas and Daniel M Kane. Recent advances in algorithmic high-dimensional robust statistics. *arXiv preprint arXiv:1911.05911*, 2019.
- 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.
- Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Leon Roth. Preserving statistical validity in adaptive data analysis. In *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, pages 117–126, 2015.
- Nika Haghtalab, Tim Roughgarden, and Abhishek Shetty. Smoothed analysis with adaptive adversaries. In *Proceedings of the 62nd Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, 2021.
- Frank R. Hampel. A General Qualitative Definition of Robustness. *The Annals of Mathematical Statistics*, 42(6):1887 – 1896, 1971. doi: 10.1214/aoms/1177693054. URL <https://doi.org/10.1214/aoms/1177693054>.
- Moritz Hardt and Jonathan Ullman. Preventing false discovery in interactive data analysis is hard. In *Proceedings of the 55th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 454–463, 2014.
- David Haussler. Decision theoretic generalizations of the pac model for neural net and other learning applications. *Information and Computation*, 100(1):78–150, 1992.
- Peter Huber. Robust estimation of a location parameter. *The Annals of Mathematical Statistics*, 35(1), 1964.
- Michael Kearns. Efficient noise-tolerant learning from statistical queries. *Journal of the ACM (JACM)*, 45(6):983–1006, 1998.
- Michael Kearns and Ming Li. Learning in the presence of malicious errors. *SIAM Journal on Computing*, 22(4):807–837, 1993.

- Michael Kearns, Robert Schapire, and Linda Sellie. Toward efficient agnostic learning. *Machine Learning*, 17(2/3):115–141, 1994.
- 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.
- Thomas Steinke and Jonathan Ullman. Interactive fingerprinting codes and the hardness of preventing false discovery. In *Proceedings of the 28th Conference on Learning Theory (COLT)*, pages 1588–1628, 2015.
- John W Tukey. Mathematics and the picturing of data. In *Proceedings of the International Congress of Mathematicians, Vancouver*, volume 2, pages 523–531, 1975.
- Leslie G. Valiant. Learning disjunction of conjunctions. In *Proceedings of the 9th International Joint Conference on Artificial Intelligence (IJCAI)*, pages 560–566, 1985.
- Banghua Zhu, Jiantao Jiao, and Jacob Steinhardt. Generalized resilience and robust statistics. *arXiv*, abs/1909.08755, 2019. URL <http://arxiv.org/abs/1909.08755>.

Appendix A. Other standard noise models

Here, we list other standard noise models and show they fall under our framework.

Definition 15 (Nasty classification noise (Bshouty et al., 2002)) *Let the domain be $\mathcal{X} = X \times Y$. Given a size- n sample $S \in \mathcal{X}^n$ and a corruption budget η , the adaptive nasty classification noise adversary is allowed to choose $\lfloor \eta n \rfloor$ points and for each one, change it from (x, y) to (x, \hat{y}) for arbitrary $\hat{y} \in Y$.*

This model is captured by the cost function:

$$\text{cost}_{\text{agn}}(\mathcal{D}, \hat{\mathcal{D}}) = \begin{cases} \infty & \text{if } \mathcal{D} \text{ and } \mathcal{D}' \text{ do not have the same marginal distribution over } \mathcal{X} \\ \text{dist}_{\text{TV}}(\mathcal{D}, \hat{\mathcal{D}}) & \text{otherwise} \end{cases}$$

The $(\text{cost}_{\text{agn}}, \eta)$ -oblivious adversary corresponds exactly to the well-studied agnostic learning model (Haussler, 1992; Kearns et al., 1994). Hence, Theorem 1 implies that nasty classification noise and the agnostic learning model are identical for SQ algorithms.

The final noise model that we discuss is defined with respect to an adversary that has intermediate adaptive power:

Definition 16 (Malicious noise (Valiant, 1985)) *In the malicious noise model where the adversary has corruption budget η , a sample is generated point-by-point. For each point, independently with probability $1 - \eta$, that point is $x \sim \mathcal{D}$. Otherwise, the adversary is allowed to make that point an arbitrary $x \in \mathcal{X}$ with knowledge of the previous points sampled but not the future points.*

On the relationship between malicious noise and additive noise. The malicious noise adversary does not have full adaptivity, as when they decide what point to add, they only have knowledge of previous points sampled and not future ones. We now show how to encode the fully adaptive version of malicious noise, in which the adversary knows all points in the sample when deciding corruptions, in our framework.

We first augment the domain to $\mathcal{X}' = \mathcal{X} \cup \{\emptyset\}$ where \emptyset will be used to indicate the adversary can change this point arbitrarily. We then let \mathcal{D}' be the distribution satisfying, for each $x \in \mathcal{X}$:

$$\mathcal{D}' = (1 - \eta)\mathcal{D} + \eta\mathcal{D}_{\emptyset}$$

where \mathcal{D}_{\emptyset} is the distribution that always outputs \emptyset . We define the cost function to be

$$\text{cost}_{\text{mal}}(\mathcal{D}', \widehat{\mathcal{D}}') = \begin{cases} \infty & \text{if } \Pr_{\mathbf{x} \sim \widehat{\mathcal{D}}'}[\mathbf{x} = \emptyset] > 0 \\ \infty & \text{if } \Pr_{\mathbf{x} \sim \widehat{\mathcal{D}}'}[\mathbf{x} = x] < \Pr_{\mathbf{x} \sim \widehat{\mathcal{D}}}[\mathbf{x} = x] \text{ for any } x \in \mathcal{X} \\ 0 & \text{otherwise.} \end{cases}$$

Note that in order for $\widehat{\mathcal{D}}'$ to be a valid corruption of \mathcal{D}' (i.e. $\text{cost}_{\text{mal}}(\mathcal{D}', \widehat{\mathcal{D}}') \neq \infty$), all of the probability mass of $\widehat{\mathcal{D}}'$ must be over \mathcal{X} , with none on \emptyset . The above provides an encoding of an adaptive adversary that is at least as powerful as malicious noise. To further understand the corresponding oblivious adversary, we note that, after fixing η ,

$$\text{cost}_{\text{mal}}(\mathcal{D}', \widehat{\mathcal{D}}') \neq \infty \quad \text{if and only if} \quad \widehat{\mathcal{D}}' = (1 - \eta)\mathcal{D} + \eta\mathcal{E} \quad \text{for some distribution } \mathcal{E}.$$

Hence, the oblivious adversary corresponding to malicious noise is the same as the oblivious adversary for additive noise. [Theorem 2](#) implies that the adaptive and oblivious versions of additive noise, as well as malicious noise, are all equivalent.

A.1. Technical remarks

Fixed budget vs. variable budget. Consider the nasty noise model ([Definition 7](#)), corresponding to $\rho = \text{dist}_{\text{TV}}$ in our framework. Given a size- n clean sample \mathcal{S} , the adaptive adversary can choose any η -fraction of the points to change arbitrarily to create the corrupted sample $\widehat{\mathcal{S}}$. Often an alternative definition is used where the adversary is allowed to arbitrarily change m points in \mathcal{S} , where m can vary based on the specific sample \mathcal{S} , as long as the marginal distribution of m over samples $\mathcal{S} \sim \mathcal{D}^n$ is $\text{Bin}(n, \eta)$. This definition is used by ([Diakonikolas et al., 2019](#); [Zhu et al., 2019](#)) to show that the adaptive adversary can simulate any oblivious adversary. Technically, for our definition of an adaptive adversary with fixed budget η , this fact is not strictly true.⁸ However, all our results are readily extendable to these slightly stronger adaptive adversaries with random-budgets.

We define an adversary model which encompasses those adversaries with variable budgets. First, the adversary will generate a corrupted data set $\widehat{\mathcal{S}}$. Then, it is allowed to change some points in that set to create $\widehat{\mathcal{S}}'$ as long as, most of the time, $\widehat{\mathcal{S}}'$ and $\widehat{\mathcal{S}}$ are close.

Definition 17 (Strong (ρ, η) -adaptive adversary) *Given a cost function ρ and budget η , an n -sample algorithm operating in the strong (ρ, η) -adaptive adversary model will receive as input the sample $\widehat{\mathcal{S}}'$ where*

8. For example, suppose $\mathcal{X} = \{0, 1\}$, \mathcal{D} is the identically 0 distribution, and $\eta = 0.1$. If a size-10 sample is taken, under our adaptive definition, $\widehat{\mathcal{S}}$ will never have more than a single 1. However, the oblivious adversary can choose $\widehat{\mathcal{D}}$ that returns 1 with probability 0.1 and as a result is a non-zero chance that two 1's appear in a size-10 sample.

1. If the true data distribution is \mathcal{D} , first a clean sample $\mathbf{S} \sim \mathcal{D}^n$ is generated.
2. The adversary chooses a $\widehat{\mathbf{S}}$ satisfying $\rho(\mathcal{U}(\mathbf{S}), \mathcal{U}(\widehat{\mathbf{S}})) \leq \eta$.
3. The adversary chooses some $\widehat{\mathbf{S}}'$ where, over the randomness of the original sample and the adversary's decisions, the following holds.

$$\Pr_{\widehat{\mathbf{S}}, \widehat{\mathbf{S}}'} \left[\text{dist}_{\text{TV}}(\mathcal{U}(\widehat{\mathbf{S}}), \mathcal{U}(\widehat{\mathbf{S}}')) \geq t \right] \leq \exp(-O(nt^2)) \quad \text{for all } t \in (0, 1). \quad (2)$$

All of our upper bounds on the strength of adaptive adversaries also apply to the strong adaptive adversary. See [Remark 1](#) for changes in the proof needed for [Theorem 1](#). For [Theorems 2](#) and [3](#), we can use an even weaker restriction on the adversary and only require that

$$\mathbb{E}_{\widehat{\mathbf{S}}, \widehat{\mathbf{S}}'} \left[\text{dist}_{\text{TV}}(\mathcal{U}(\widehat{\mathbf{S}}), \mathcal{U}(\widehat{\mathbf{S}}')) \right] = O(1/\sqrt{n}) \quad (3)$$

in place of [Equation \(2\)](#). See [Remark 4](#) for how [Equation \(3\)](#) can be used to make [Theorems 2](#) and [3](#) work with strong adaptive adversaries.

Finite vs infinite domains. For our analyses, we assume that the domain, \mathcal{X} , is finite. Our goal is to understand when an algorithm that succeeds in the presence of an oblivious adversary implies an algorithm that succeeds in the presence of an adaptive adversary. Any algorithm that succeeds in the presence of an oblivious adversary can only read finitely many bits of each data point, effectively discretizing the domain.

That said, [Theorem 1](#) also applies to infinite domains. If the domain is infinite, a more general definition of closed under mixtures is required in place of the simpler [Definition 13](#)

Definition 18 (Closed under mixtures, infinite domain) *We say that ρ is closed under mixtures if for any distributions $\mathcal{D}, \mathcal{D}'$, and coupling of $\mathbf{x} \sim \mathcal{D}$, $\mathbf{x}' \sim \mathcal{D}'$ and a latent variable \mathbf{z} (over any domain),*

$$\rho(\mathcal{D}, \mathcal{D}') \leq \sup_{\mathbf{z}} \left\{ \rho((\mathbf{x} \mid \mathbf{z}), (\mathbf{x}' \mid \mathbf{z})) \right\}.$$

When the domain is finite, [Definitions 13](#) and [18](#) are equivalent. When it is infinite, [Definition 18](#) is needed to prove [Fact 2](#). The remainder of the proof of [Theorem 1](#) is identical.

Appendix B. Missing Lemmas from the proof of [Theorem 5](#)

We will use a few standard technical tools:

Fact 3 (Separating hyperplane theorem) *Let $A, B \in \mathbb{R}^k$ be disjoint, nonempty, and convex. There exists a nonzero vector $w \in \mathbb{R}^k$ and $T \in \mathbb{R}$ such that $a \cdot w \leq T$ and $b \cdot w \geq T$ for all $a \in A$ and $b \in B$.*

Fact 4 (McDiarmid's inequality) *Suppose that $f : \mathcal{X}^n \rightarrow \mathbb{R}$ satisfies the c -bounded difference property: for any $(x_1, \dots, x_n), (x'_1, \dots, x'_n) \in \mathcal{X}^n$ that differ on only on a single coordinate, f satisfies*

$$|f(x_1, \dots, x_n) - f(x'_1, \dots, x'_n)| \leq c.$$

Then, for any $\tau > 0$ and any distribution \mathcal{D} over \mathcal{X} ,

$$\Pr_{\mathbf{S} \sim \mathcal{D}^n} [f(\mathbf{S}) - \mu \geq \tau] \leq \exp\left(-\frac{2\tau^2}{c^2 n}\right) \quad \text{where } \mu := \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^n} [f(\mathbf{S})].$$

We prove the following two Lemmas, restated for convenience.

Lemma 1 (Theorem 5 in the case of a single SQ) *Let $\Psi : \mathcal{X} \rightarrow [-1, 1]$ be a statistical query, $T \in [-1, 1]$, and suppose:*

$$\Psi(\widehat{\mathcal{D}}) \leq T \quad \text{for all } \widehat{\mathcal{D}} \text{ that are } \eta\text{-close to } \mathcal{D}.$$

Then, for any $\tau > 0$ and sample size $n \in \mathbb{N}$, the probability over $\mathbf{S} \sim \mathcal{D}^n$ that there is some $\widehat{\mathbf{S}}$ that is η -close to \mathbf{S} satisfying

$$\Psi(\widehat{\mathbf{S}}) \geq T + \frac{\tau}{2}$$

is at most $\exp\left(-\frac{\tau^2 n}{8\ell^2}\right)$.

Proof For any sample $S \in \mathcal{X}^n$, define

$$f(S) := \sup_{\widehat{\mathcal{U}}(S) \text{ is } \eta\text{-close to } \mathcal{U}(S)} \left\{ \Psi(\widehat{\mathcal{U}}(S)) \right\}.$$

We will show that f satisfies the $c := \frac{2\ell}{n}$ -bounded difference property. Consider any samples $S, S' \in \mathcal{X}^n$ that differ in only a single point. We will show that for every point in the set

$$\{\Psi(\widehat{\mathcal{U}}(S)) \mid \widehat{\mathcal{U}}(S) \text{ is } \eta\text{-close to } \mathcal{U}(S)\},$$

there is some point in the set

$$\{\Psi(\widehat{\mathcal{U}}(S')) \mid \widehat{\mathcal{U}}(S') \text{ is } \eta\text{-close to } \mathcal{U}(S')\}$$

that differs from it by at most $\pm \frac{2\ell}{n}$, and vice versa. This implies that f satisfies the $c := \frac{2\ell}{n}$ -bounded difference property.

As S and S' only differ in a single point,

$$\text{dist}_{\text{TV}}(\mathcal{U}(S), \mathcal{U}(S')) \leq \frac{1}{n}.$$

Furthermore by [Definition 14](#), for any $\widehat{\mathcal{U}}(S)$ that is η -close to $\mathcal{U}(S)$, there is some $\widehat{\mathcal{U}}(S')$ that is η -close to $\mathcal{U}(S')$ satisfying

$$\text{dist}_{\text{TV}}(\widehat{\mathcal{U}}(S), \widehat{\mathcal{U}}(S')) \leq \frac{\ell}{n}.$$

Using the [Definition 2](#) and the fact that the range of Ψ is a length-2 interval, the above implies that

$$\left| \Psi(\widehat{\mathcal{U}}(S)) - \Psi(\widehat{\mathcal{U}}(S')) \right| \leq \frac{2\ell}{n}$$

proving that f satisfies the $\frac{2\ell}{n}$ -bounded difference property. By McDiarmid's inequality, $f(\mathbf{S})$ concentrates around its mean. Lastly, we show that

$$\mathbb{E}_{\mathbf{S} \sim \mathcal{D}^n} [f(\mathbf{S})] \leq T. \quad (4)$$

Suppose for the sake of contradiction that $\mathbb{E}[f(\mathbf{S})] > T + \varepsilon$ for some $\varepsilon > 0$. For each $S \in \mathcal{X}^n$, let $\widehat{\mathcal{U}}(S)$ be η -close to $\mathcal{U}(S)$ and satisfy $\Psi(\widehat{\mathcal{U}}(S)) \geq f(S) - \varepsilon$ (which exists by the definition of f). We can define $\widehat{\mathcal{D}}$ as the distribution where, to sample $x \sim \widehat{\mathcal{D}}$, we

1. Draw an i.i.d. sample $S \sim \mathcal{D}^n$.
2. Draw $x \sim \widehat{\mathcal{U}}(S)$ uniformly.

By [Fact 2](#), $\widehat{\mathcal{D}}$ is η -close to \mathcal{D} . Then,

$$T \geq \Psi(\widehat{\mathcal{D}}) = \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^n} [\Psi(\widehat{\mathcal{U}}(\mathbf{S}))] \geq \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^n} [f(\mathbf{S}) - \varepsilon] = \mathbb{E}[f(\mathbf{S})] - \varepsilon > T.$$

This is a contradiction, so [Equation \(4\)](#) holds. [Lemma 1](#) follows from McDiarmid's inequality applied to f . ■

Lemma 2 (*k*-query SQ algorithm to a single SQ) Fix any k statistical queries $\phi^{(1)}, \dots, \phi^{(k)} : \mathcal{X} \rightarrow [-1, 1]$ and k values v_1, \dots, v_k . Suppose that there is no $\widehat{\mathcal{D}}$ that is η -close to \mathcal{D} satisfying

$$\phi^{(i)}(\widehat{\mathcal{D}}) \in v_i \pm \tau \quad \text{for every } i \in [k].$$

Then there exists a single statistical query $\Psi : \mathcal{X} \rightarrow [-1, 1]$ and threshold T with the following properties.

1. $\Psi(\widehat{\mathcal{D}}) \leq T$ for every $\widehat{\mathcal{D}}$ that is η -close to \mathcal{D} .
2. For any sample \widehat{S} satisfying $\phi^{(i)}(\widehat{S}) \in v_i \pm \frac{\tau}{2}$ for each $i \in [k]$, it is also true that $\Psi(\widehat{S}) \geq T + \frac{\tau}{2}$.

Proof We'll actually prove a slightly more general result. We'll show that for any any *distribution* \mathcal{E} satisfying $\phi^{(i)}(\mathcal{E}) \in v_i \pm \frac{\tau}{2}$ for each $i \in [k]$, it is also true that $\Psi(\mathcal{E}) \geq T + \frac{\tau}{2}$. [Lemma 2](#) follows by setting $\mathcal{E} = \mathcal{U}(\widehat{S})$.

We define $A \in \mathbb{R}^d$ to be

$$A := \left\{ \left(\phi^{(1)}(\widehat{\mathcal{D}}), \dots, \phi^{(k)}(\widehat{\mathcal{D}}) \right) \mid \widehat{\mathcal{D}} \text{ is } \eta\text{-close to } \mathcal{D} \right\},$$

which is convex since the cost function is closed under mixtures ([Definition 13](#)). We define B to be

$$B := \left\{ b \in \mathbb{R}^k \mid b_i \in v_i \pm \tau \text{ for all } i \in [k] \right\},$$

which is convex since it is the intersection of halfspaces. By the assumptions of [Lemma 2](#), A and B are disjoint. Let $w \in \mathbb{R}^k$ and $T \in \mathbb{R}$ be the vector and threshold respectively guaranteed to exist by the separating hyperplane theorem, normalized so that $\|w\|_1 = 1$. We define

$$\Psi(x) := \sum_{i \in [k]} w_i \cdot \phi^{(i)}(x).$$

As $\|w\|_1 = 1$ and the range of each $\phi^{(i)}(x) \in [-1, 1]$ for each $x \in \mathcal{X}$, it is also true that $\Psi(x) \in [-1, 1]$. We show that Ψ, T meet the two criteria of [Lemma 2](#). The [first criteria](#) holds by the separating hyperplane theorem. The [second criteria](#) is equivalent to showing that any $b \in B_{\text{inner}}$ satisfies $b \cdot w \geq T + \frac{\tau}{2}$ where

$$B_{\text{inner}} := \left\{ b \in \mathbb{R}^k \mid b_i \in v_i \pm \frac{\tau}{2} \text{ for all } i \in [k] \right\}.$$

There must be a minimal point for $b \cdot w$ at a ‘‘corner’’ of B . Let b^* be such a minimal point (i.e. $(b^*)_i = v_i + c_i \cdot \tau$ for $c_i \in \{\pm 1\}$). For any $b \in B_{\text{inner}}$

$$\begin{aligned} b \cdot w &= b^* \cdot w + \sum_{i \in [k]} w_i (b_i - b_i^*) \\ &= b^* \cdot w + \sum_{i \in [k]} |w_i| \cdot |b_i - b_i^*| && (b^* \text{ minimal for } b \cdot w \text{ over } b \in B) \\ &\geq b^* \cdot w + \sum_{i \in [k]} |w_i| \cdot \frac{\tau}{2} && (|b_i - b_i^*| \geq \tau/2 \text{ for } b \in B_{\text{inner}}) \\ &\geq T + \frac{\tau}{2}. && (\text{Separating hyperplane theorem, } \|w\|_1 = 1) \end{aligned}$$

This implies the [second criteria](#) of [Lemma 2](#). ■

Remark 1 *Theorem 5 also holds with strong adaptive adversaries ([Definition 17](#) in [Appendix A.1](#)) rather than just adaptive adversaries, with the appropriate changes in constants. The proof only differs in [Lemma 1](#). We then wish to bound the probability that the adversary can make $\Psi(\widehat{\mathcal{S}}') \geq T + \tau/2$, where $\widehat{\mathcal{S}}$ and $\widehat{\mathcal{S}}'$ are as in [Definition 17](#). For constant ℓ ,*

$$\begin{aligned} \Pr[\Psi(\widehat{\mathcal{S}}') \geq T + \tau/2] &\leq \Pr_{\widehat{\mathcal{S}}, \widehat{\mathcal{S}}'} \left[|\Psi(\widehat{\mathcal{S}}) - \Psi(\widehat{\mathcal{S}}')| \geq \tau/4 \right] + \Pr_{\widehat{\mathcal{S}}} [\Psi(\widehat{\mathcal{S}}) \geq T + \tau/4] \\ &\leq \Pr_{\widehat{\mathcal{S}}, \widehat{\mathcal{S}}'} \left[\text{dist}_{\text{TV}}(\mathcal{U}(\widehat{\mathcal{S}}), \mathcal{U}(\widehat{\mathcal{S}}')) \geq \tau/8 \right] + \exp(-O(\tau^2 n)) \quad (\text{Lemma 1}) \\ &\leq \exp(-O(\tau^2 n)). \quad (\text{Definition 17}) \end{aligned}$$

Once [Lemma 1](#) is modified to handle strong adaptive adversaries, the remainder of the proof of [Theorem 5](#) applies unchanged.

Appendix C. Proof of [Theorem 3](#): If adaptivity can be neutralized, subsampling does it

Towards tackling [Question 1](#) for all algorithms, we define the subsampling filter, a natural ‘‘wrapper algorithm’’ that operates only on samples and can be applied to any existing algorithm:

Definition 19 (Subsampling Filter) *Define the subsampling filter $\Phi_{m \rightarrow n} : \mathcal{X}^m \rightarrow \mathcal{X}^n$ that, given a set $S \in \mathcal{X}^m$, subsamples n elements $S' \sim \mathcal{U}(S)^n$ and returns them. We will also write $\Phi_{* \rightarrow n}$ when the size of S is variable.*

Intuitively, by requesting a large number of points m and randomly subsampling points for the original algorithm, the filter should be able to neutralise some of the power of the adaptive adversary, since the adversary cannot know which subsample the algorithm will receive. In this section we will prove [Theorem 3](#) which, informally speaking, states that *if* the noise model is such that adaptivity can be neutralised, then subsampling does it. In the next section, we carry out this proof strategy for the specific case of additive noise and establish [Theorem 2](#).

C.1. Definitions and the formal statement of [Theorem 3](#)

We begin by formalizing what it means for the behavior of an algorithm \mathcal{A} in the presence of an oblivious adversary to be equivalent to that of an algorithm \mathcal{A}' in the presence of an adaptive adversary. Roughly speaking, that corresponds to the range of acceptance probabilities \mathcal{A} can have with all possible oblivious adversaries being close to the range of acceptance probabilities that \mathcal{A}' can have with all possible adaptive adversaries.

Definition 20 Fix a cost function ρ , budget $\eta \geq 0$, and distribution \mathcal{D} over \mathcal{X} . For an algorithm $\mathcal{A} : \mathcal{X}^* \rightarrow \{0, 1\}$, we define:

$$\begin{aligned} \text{Oblivious-Max}_{\rho, \eta}(\mathcal{A}, \mathcal{D}, n) &:= \sup_{\rho(\mathcal{D}, \widehat{\mathcal{D}}) \leq \eta} \left\{ \mathbb{E}_{\mathcal{S} \sim \widehat{\mathcal{D}}^n} [\mathcal{A}(\mathcal{S})] \right\}, \\ \text{Oblivious-Min}_{\rho, \eta}(\mathcal{A}, \mathcal{D}, n) &:= \inf_{\rho(\mathcal{D}, \widehat{\mathcal{D}}) \leq \eta} \left\{ \mathbb{E}_{\mathcal{S} \sim \widehat{\mathcal{D}}^n} [\mathcal{A}(\mathcal{S})] \right\}, \end{aligned}$$

the maximum and minimum acceptance probabilities of \mathcal{A} given an obviously corrupted \mathcal{D} . We similarly define the adaptive versions:

$$\begin{aligned} \text{Adaptive-Max}_{\rho, \eta}(\mathcal{A}, \mathcal{D}, m) &:= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \left[\sup_{\rho(\mathcal{U}(\mathcal{S}), \mathcal{U}(\widehat{\mathcal{S}})) \leq \eta} \left\{ \mathcal{A}(\widehat{\mathcal{S}}) \right\} \right], \\ \text{Adaptive-Min}_{\rho, \eta}(\mathcal{A}, \mathcal{D}, m) &:= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \left[\inf_{\rho(\mathcal{U}(\mathcal{S}), \mathcal{U}(\widehat{\mathcal{S}})) \leq \eta} \left\{ \mathcal{A}(\widehat{\mathcal{S}}) \right\} \right]. \end{aligned}$$

Definition 21 (ε -equivalent) Fix a cost function ρ and a budget $\eta \geq 0$. Let $\mathcal{A}, \mathcal{A}' : \mathcal{X}^* \rightarrow \{0, 1\}$ be two algorithms. We say that \mathcal{A} in the presence of (ρ, η) -oblivious adversaries is (n, m, ε) -equivalent to \mathcal{A}' in the presence of (ρ, η) -adaptive adversaries if the following holds for all distributions \mathcal{D} over \mathcal{X} :

$$\text{Adaptive-Max}_{\rho, \eta}(\mathcal{A}', \mathcal{D}, m) = \text{Oblivious-Max}_{\rho, \eta}(\mathcal{A}, \mathcal{D}, n) \pm \varepsilon,$$

and likewise for Min instead of Max. If the algorithms \mathcal{A} or \mathcal{A}' are randomized, then these expectations are also over the randomness of the algorithms.

We now state the formal version of [Theorem 3](#):

Theorem 6 (Formal version of [Theorem 3](#): If it is possible, subsampling does it) Fix a cost function ρ and budget $\eta \geq 0$. Suppose that \mathcal{A} and \mathcal{A}' are algorithms where \mathcal{A} in the presence of (ρ, η) -oblivious adversaries is (n, m, ε) -equivalent to \mathcal{A}' in the presence of (ρ, η) -adaptive adversaries.

Consider the subsampling algorithm $\mathcal{A}_{\text{sub}} := \mathcal{A} \circ \Phi_{* \rightarrow n}$ which, given a sample $S \in \mathcal{X}^*$, subsamples n elements $\mathbf{S}' \sim \mathcal{U}(S)^n$ and returns $\mathcal{A}(\mathbf{S}')$. For

$$M := O\left(\frac{m^2 \log(1/\varepsilon)^2}{\varepsilon^5}\right),$$

we have that \mathcal{A} in the presence of (ρ, η) -oblivious adversaries is $(n, M, 9\varepsilon)$ -equivalent to \mathcal{A}_{sub} in the presence of (ρ, η) -adaptive adversaries.

Remark 2 (Search to decision reduction) In both [Theorems 2](#) and [3](#), we focus on decision algorithms that output a single bit $\{0, 1\}$, rather than the more general setting of search algorithms that output an answer from some set \mathcal{Y} . This is without loss of generality. Given a search algorithm $\mathcal{A} : \mathcal{X}^* \rightarrow \mathcal{Y}$, and any set of “good outputs” $Y \subseteq \mathcal{Y}$, we could define an algorithm $\mathcal{B} := \mathbb{1}_Y \circ \mathcal{A}$ where $\mathcal{B}(S) = 1$ iff $\mathcal{A}(S) \in Y$. Then, we can directly apply [Theorems 2](#) and [3](#) to \mathcal{B} . Hence, [Theorems 2](#) and [3](#) hold for search algorithms with the appropriate definition of “equivalence” for search algorithms: We say $\mathcal{A} : \mathcal{X}^* \rightarrow \mathcal{Y}$ and $\mathcal{A}' : \mathcal{X}^* \rightarrow \mathcal{Y}$ are ε -equivalent if for every $Y \subseteq \mathcal{Y}$, $\mathbb{1}_Y \circ \mathcal{A}$ and $\mathbb{1}_Y \circ \mathcal{A}'$ are ε -equivalent (according to [Definition 21](#)).

C.2. Proof of [Theorem 6](#)

Our proof of [Theorem 6](#) relies on the following simple lemma. Roughly speaking, it states that sampling with replacement and sampling without replacement are nearly indistinguishable when the population is a quadratic factor larger than the number of samples.

Lemma 3 For any distribution \mathcal{D} and integers $m, M \in \mathbb{N}$, let $\Phi_{M \rightarrow m} \circ \mathcal{D}^M$ be the distribution with the following generative process: first draw a size- M sample $\mathbf{S} \sim \mathcal{D}^M$, and then subsample, with replacement, m points from \mathbf{S} . Then

$$\text{dist}_{\text{TV}}(\mathcal{D}^m, \Phi_{M \rightarrow m} \circ \mathcal{D}^M) \leq \frac{\binom{m}{2}}{M}.$$

Proof We describe a coupling of $\mathbf{S} \sim \mathcal{D}^m$ and $\mathbf{S}' \sim \Phi_{M \rightarrow m} \circ \mathcal{D}^M$ such that $\Pr[\mathbf{S} \neq \mathbf{S}'] \leq \binom{m}{2}/M$.

1. Initialize S and S' to be empty sets, and y_1, \dots, y_M to be unset variables.
2. Repeat m times:
 - (a) Draw $i \sim [M]$ uniformly.
 - (b) If y_i is unset, draw $x \sim \mathcal{D}$ and set $y_i \leftarrow x$. Then, add x to both \mathbf{S} and \mathbf{S}' .
 - (c) Otherwise, add y_i to \mathbf{S}' and sample $x \sim \mathcal{D}$ to add to \mathbf{S} .

It is straightforward to verify that the above generative process leads to the distribution of \mathbf{S} and \mathbf{S}' being that of \mathcal{D}^m and $\Phi_{M \rightarrow m} \circ \mathcal{D}^M$ respectively. Furthermore, if $\mathbf{S} \neq \mathbf{S}'$, that means there is some index $i \in [M]$ that was sampled at least twice. If we fix $j_1 \neq j_2 \in [m]$ and some $i \in [M]$, the probability i is the index chosen at steps j_1 and j_2 is $1/M^2$. Union bounding over the M choices for i and $\binom{m}{2}$ for j_1, j_2 gives that

$$\Pr[\mathbf{S} \neq \mathbf{S}'] \leq \frac{\binom{m}{2}}{M}.$$

The desired result follows from the definition of total variation cost, [Definition 2](#). ■

[Lemma 3](#) will be used in conjunction with the following fact:

Fact 5 Suppose there exists a test which, given c samples from a distribution \mathcal{E} that is either \mathcal{D}_0 or \mathcal{D}_1 , returns 0 if $\mathcal{E} = \mathcal{D}_0$ with probability at least $\frac{3}{4}$, and returns 1 if $\mathcal{E} = \mathcal{D}_1$ with probability at least $\frac{3}{4}$. Then,

$$\text{dist}_{\text{TV}}(\mathcal{D}_1, \mathcal{D}_2) \geq \frac{1}{2c}.$$

Together, [Lemma 3](#) and [Fact 5](#) imply that for appropriately chosen m and M , there is no sample-efficient test distinguishing \mathcal{D}^m from $\Phi_{M \rightarrow m} \circ \mathcal{D}_M$ with high probability. We will use this to prove [Theorem 6](#) by contradiction. Using the assumption that \mathcal{A}' is equivalent to \mathcal{A} , we will design a sample-efficient test that approximates Oblivious- $\text{Max}_{\rho, \eta}$ and Oblivious- $\text{Min}_{\rho, \eta}$ for \mathcal{A} with respect to both \mathcal{D}^m and $\Phi_{M \rightarrow m} \circ \mathcal{D}^M$. We then show that if $\mathcal{A}_{\text{sub}} := \Phi_{* \rightarrow n} \circ \mathcal{A}$ is not equivalent to \mathcal{A} , then these values will distinguish the two distributions.

The following lemma carries out the first part of this plan:

Lemma 4 Let \mathcal{A} and \mathcal{A}' be as in [Theorem 6](#). There is an estimator $\text{Est-}\text{Max}_{\rho, \eta}$ that uses

$$m' := \frac{m \log(2/\varepsilon)}{2\varepsilon^2}$$

samples from a distribution \mathcal{D} and returns an estimate of $\text{Oblivious-}\text{Max}_{\rho, \eta}(\mathcal{A}, \mathcal{D}, n)$ that is accurate to $\pm 2\varepsilon$ with probability at least $1 - \varepsilon$, and likewise an estimator $\text{Est-}\text{Min}_{\rho, \eta}$ for $\text{Oblivious-}\text{Min}_{\rho, \eta}$. Formally, for all distributions \mathcal{D} over \mathcal{X} ,

$$\Pr_{\mathcal{S} \sim \mathcal{D}^{m'}} \left[\text{Est-}\text{Max}_{\rho, \eta}(\mathcal{S}) = \text{Oblivious-}\text{Max}_{\rho, \eta}(\mathcal{A}, \mathcal{D}, n) \pm 2\varepsilon \right] \geq 1 - \varepsilon,$$

and likewise for $\text{Est-}\text{Min}_{\rho, \eta}$ and $\text{Oblivious-}\text{Min}_{\rho, \eta}$.

Proof We will prove the lemma for $\text{Est-}\text{Max}_{\rho, \eta}$ and $\text{Oblivious-}\text{Max}_{\rho, \eta}$; the proof for Min instead of Max is identical. $\text{Est-}\text{Max}_{\rho, \eta}$ computes an estimate satisfying:

$$\Pr_{\mathcal{S} \sim \mathcal{D}^{m'}} \left[\text{Est-}\text{Max}_{\rho, \eta}(\mathcal{S}) = \text{Adaptive-}\text{Max}_{\rho, \eta}(\mathcal{A}', \mathcal{D}, m) \pm \varepsilon \right] \leq \varepsilon. \quad (5)$$

This is sufficient to guarantee that $\text{Est-}\text{Max}_{\rho, \eta}$'s estimate is within $\pm 2\varepsilon$ of $\text{Oblivious-}\text{Max}_{\rho, \eta}(\mathcal{A}, \mathcal{D}, n)$ by our assumption that \mathcal{A} is (n, m, ε) -equivalent to \mathcal{A}' . To provide such an estimate, $\text{Est-}\text{Max}_{\rho, \eta}(\mathcal{S})$ draws $\log(2/\varepsilon)/(2\varepsilon^2)$ many size- m samples $\mathcal{S}' \sim \mathcal{D}^m$. For each, it computes:

$$\sup_{\rho(\widehat{\mathcal{S}}', \mathcal{S}') \leq \eta} \{ \mathcal{A}'(\widehat{\mathcal{S}}') \} \quad (6)$$

and returns the average of these supremums. By the Chernoff bound, this average satisfies [Equation \(5\)](#). ■

Remark 3 We are only concerned with the sample efficiency of these estimators, not their time efficiency or even whether they are computable. Indeed, as stated, an algorithm computing the estimators would need to loop or infinitely many $\widehat{\mathcal{S}}' \in \mathcal{X}^*$ to compute [Equation \(6\)](#). For us they are just an analytical tool used to prove [Theorem 6](#), the conclusion of which gives an algorithm $\mathcal{A}_{\text{sub}} := \mathcal{A} \circ \Phi_{* \rightarrow n}$ that inherits the time efficiency of \mathcal{A} .

The next lemma notes that the Adaptive- $\text{Max}_{\rho,\eta}$ of $\mathcal{A}_{\text{sub}} := \mathcal{A} \circ \Phi_{* \rightarrow n}$ can be expressed in terms of the Oblivious- $\text{Max}_{\rho,\eta}$ of \mathcal{A} . Formally:

Lemma 5 Adaptive- $\text{Max}_{\rho,\eta}(\mathcal{A}_{\text{sub}}, \mathcal{D}, M) = \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^M} [\text{Oblivious-}\text{Max}_{\rho,\eta}(\mathcal{A}, \mathcal{U}(\mathcal{S}), n)]$, and likewise for Min instead of Max.

Proof The lemma follows from this series of identities:

$$\begin{aligned}
 \text{Adaptive-}\text{Max}_{\rho,\eta}(\mathcal{A}_{\text{sub}}, \mathcal{D}, M) &= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^M} \left[\sup_{\rho(\mathcal{U}(\mathcal{S}), \mathcal{U}(\widehat{\mathcal{S}})) \leq \eta} \left\{ \mathcal{A}_{\text{sub}}(\widehat{\mathcal{S}}) \right\} \right] \\
 &\quad \text{(Definition of Adaptive-}\text{Max}_{\rho,\eta}) \\
 &= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^M} \left[\sup_{\rho(\mathcal{U}(\mathcal{S}), \mathcal{U}(\widehat{\mathcal{S}})) \leq \eta} \left\{ (\mathcal{A} \circ \Phi_{* \rightarrow n})(\widehat{\mathcal{S}}) \right\} \right] \\
 &\quad \text{(Definition of } \mathcal{A}_{\text{sub}}) \\
 &= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^M} \left[\sup_{\rho(\mathcal{U}(\mathcal{S}), \mathcal{E}) \leq \eta} \left\{ \mathbb{E}_{\mathcal{S}' \sim \mathcal{E}^n} [\mathcal{A}(\mathcal{S}')] \right\} \right] \\
 &\quad \text{(Definition of the subsampling filter } \Phi_{* \rightarrow n}) \\
 &= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^M} [\text{Oblivious-}\text{Max}_{\rho,\eta}(\mathcal{A}, \mathcal{U}(\mathcal{S}), n)], \\
 &\quad \text{(Definition of Oblivious-}\text{Max}_{\rho,\eta})
 \end{aligned}$$

where the penultimate identity also uses our convention that distributions have rational weights, and therefore can be expressed as the uniform distribution over a sufficiently large multiset of elements. ■

The following lemma completes the proof of [Theorem 6](#).

Lemma 6 Let m' be as in [Lemma 4](#) and define $M := 14(m')^2/\varepsilon$. Then

$$\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^M} [\text{Oblivious-}\text{Max}_{\rho,\eta}(\mathcal{A}, \mathcal{U}(\mathcal{S}), n)] = \text{Oblivious-}\text{Max}_{\rho,\eta}(\mathcal{A}, \mathcal{D}, n) \pm 9\varepsilon,$$

and likewise for Min instead of Max.

Proof Our proof proceeds by contradiction: assuming that $\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^M} [\text{Oblivious-}\text{Max}_{\rho,\eta}(\mathcal{A}, \mathcal{U}(\mathcal{S}), n)]$ is more than 9ε far from $\text{Oblivious-}\text{Max}_{\rho,\eta}(\mathcal{A}, \mathcal{D}, n)$, we will prove that

$$\text{dist}_{\text{TV}}(\mathcal{D}^{m'}, \Phi_{M \rightarrow m'} \circ \mathcal{D}^M) \geq \frac{\varepsilon}{14}. \quad (7)$$

From [Lemma 3](#), we know that

$$\text{dist}_{\text{TV}}(\mathcal{D}^{m'}, \Phi_{M \rightarrow m'} \circ \mathcal{D}^M) < \frac{(m')^2}{M} = \frac{\varepsilon}{14},$$

which yields the desired contradiction. To establish [Equation \(7\)](#), we design an algorithm that given $\lceil \frac{6}{\varepsilon} \rceil \leq \frac{7}{\varepsilon}$ samples from either $\mathcal{D}^{m'}$ from $\Phi_{M \rightarrow m'} \circ \mathcal{D}^M$ is able to distinguish them with probability $\frac{3}{4}$. Once we do, the desired result follows from [Fact 5](#).

Let $\mu := \text{Oblivious-Max}_{\rho,\eta}(\mathcal{A}, \mathcal{D}, n)$. First, we define

$$\delta := \Pr_{\mathcal{S} \sim \mathcal{D}^M} [|\text{Oblivious-Max}_{\rho,\eta}(\mathcal{A}, \mathcal{U}(\mathcal{S}), n) - \mu| > 4\varepsilon]$$

and we bound

$$\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^M} [\text{Oblivious-Max}_{\rho,\eta}(\mathcal{A}, \mathcal{U}(\mathcal{S}), n)] \leq (1 - \delta) \cdot (\mu + 4\varepsilon) + \delta \cdot 1 \leq \mu + 4\varepsilon + \delta.$$

Similarly,

$$\begin{aligned} \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^M} [\text{Oblivious-Max}_{\rho,\eta}(\mathcal{A}, \mathcal{U}(\mathcal{S}), n)] &\geq (1 - \delta) \cdot (\mu - 4\varepsilon) + \delta \cdot 0 \\ &\geq \mu - 4\varepsilon - \delta. \end{aligned} \quad (\mu - 4\varepsilon \leq 1)$$

Hence by our assumption that

$$\left| \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^M} [\text{Oblivious-Max}_{\rho,\eta}(\mathcal{A}, \mathcal{U}(\mathcal{S}), n)] - \text{Oblivious-Max}_{\rho,\eta}(\mathcal{A}, \mathcal{D}, n) \right| > 9\varepsilon,$$

we can conclude that $\delta > 5\varepsilon$.

Now let \mathcal{E} be either $\mathcal{D}^{m'}$ or $\Phi_{M \rightarrow m'} \circ \mathcal{D}^M$. Our test to determine which \mathcal{E} is will do the following: draw $\lceil \frac{6}{\varepsilon} \rceil$ samples from \mathcal{E} , $\mathcal{S} \sim \mathcal{E}$, and run $\text{Est-Max}_{\rho,\eta}(\mathcal{S})$ on each. If less than a 2ε fraction of the estimates returned by $\text{Est-Max}_{\rho,\eta}(\mathcal{S})$ differ from μ in more than $\pm 2\varepsilon$, return that $\mathcal{E} = \mathcal{D}^{m'}$. Otherwise, return that $\mathcal{E} = \Phi_{M \rightarrow m'} \circ \mathcal{D}^M$. We will prove this test succeeds with probability at least $1 - e^{-2} \geq \frac{3}{4}$. We consider the two possible cases:

1. Case 1: $\mathcal{E} = \mathcal{D}^{m'}$. In this case, given a sample from \mathcal{E} , $\text{Est-Max}_{\rho,\eta}(\mathcal{S})$ returns an estimate that is within $\pm 2\varepsilon$ of μ with probability at least $1 - \varepsilon$. By the Chernoff bound, given $\frac{6}{\varepsilon}$ such samples, the probability more than 2ε fraction deviate from μ by more than $\pm \varepsilon$ is at most $\exp(-\frac{1}{3} \cdot \frac{6}{\varepsilon} \cdot \varepsilon) = e^{-2}$. Therefore, the test succeeds with probability at least $1 - e^{-2}$.
2. Case 2: $\mathcal{E} = \Phi_{M \rightarrow m'} \circ \mathcal{D}^M$. We showed above that with probability at least 5ε over a sample $\mathcal{S} \sim \mathcal{D}^M$ we have $|\text{Oblivious-Max}_{\rho,\eta}(\mathcal{A}, \mathcal{U}(\mathcal{S}), n) - \mu| > 4\varepsilon$. When that's the case, $\text{Est-Max}_{\rho,\eta}(\mathcal{S})$ returns an estimate that is further than $\pm 2\varepsilon$ from μ with probability at least $1 - \varepsilon$. Therefore, on a single sample, the probability that the estimate of \mathcal{A} deviates from μ by more than $\pm 2\varepsilon$ is at least $5\varepsilon(1 - \varepsilon) \geq 4\varepsilon$. By the Chernoff bound, given $\frac{6}{\varepsilon}$ samples, the probability that at most 2ε fraction deviate μ by at most $\pm 2\varepsilon$ is at most $\exp(-\frac{1}{8} \cdot \frac{6}{\varepsilon} \cdot 4\varepsilon) = e^{-3}$. Therefore, the test succeeds with probability at least $1 - e^{-2}$.

Hence, given $\lceil \frac{6}{\varepsilon} \rceil$ samples, it is possible to distinguish $\mathcal{D}^{m'}$ from $\Phi_{M \rightarrow m'} \circ \mathcal{D}^M$ with a success probability of at least $\frac{3}{4}$. Equation (7) follows from Fact 5, completing the proof by contradiction. \blacksquare

Remark 4 (Strong adaptive adversaries) *If $\mathcal{A} : \mathcal{X}^n \rightarrow \{0, 1\}$ in the presence of oblivious adversaries is (n, M, ε) -equivalent to $\mathcal{A}_{\text{sub}} := \mathcal{A} \circ \Phi_{* \rightarrow n}$ in the presence of adaptive adversaries, then it is also $(n, M, 2\varepsilon)$ -equivalent to \mathcal{A}_{sub} in the presence of strong adaptive adversaries (Definition 17 in Appendix A.1) as long as $M = \Omega(n^2/\varepsilon^2)$. This applies to both Theorem 6 in this section and Theorem 7 in the next.*

Let $\widehat{\mathbf{S}}$ and $\widehat{\mathbf{S}}'$ be defined as in [Definition 17](#). Then, for any strong adaptive adversary supplying the sample $\widehat{\mathbf{S}}'$, there is some adaptive adversary supplying the sample $\widehat{\mathbf{S}}$, so we wish to compare $\mathbb{E}[\mathcal{A}_{\text{sub}}(\widehat{\mathbf{S}}')]$ to $\mathbb{E}[\mathcal{A}_{\text{sub}}(\widehat{\mathbf{S}})]$. Recall that \mathcal{A}_{sub} first subsamples to n points, so

$$\begin{aligned} \left| \mathbb{E}[\mathcal{A}_{\text{sub}}(\widehat{\mathbf{S}}')] - \mathbb{E}[\mathcal{A}_{\text{sub}}(\widehat{\mathbf{S}})] \right| &\leq n \mathbb{E}_{\widehat{\mathbf{S}}, \widehat{\mathbf{S}}'} [\text{dist}_{\text{TV}}(\mathcal{U}(\widehat{\mathbf{S}}), \mathcal{U}(\widehat{\mathbf{S}}'))] \\ &\leq O(n/\sqrt{M}) = \varepsilon. \end{aligned} \quad (\text{by Equation (3)})$$

Appendix D. Proof of [Theorem 2](#): The subsampling filter neutralizes adaptive additive noise

In this section, we prove the following theorem:

Theorem 7 (Formal version of [Theorem 2](#))

Fix a budget $\eta \geq 0$, distribution \mathcal{D} over \mathcal{X} , and algorithm $\mathcal{A} : \mathcal{X}^n \rightarrow \{0, 1\}$. Consider the subsampling algorithm $\mathcal{A}_{\text{sub}} := \mathcal{A} \circ \Phi_{* \rightarrow n}$ which, given a sample $S \in \mathcal{X}^*$, subsamples n elements $S' \sim \mathcal{U}(S)^n$ and returns $\mathcal{A}(S')$. For

$$M := O\left(\frac{n^4 \log(|\mathcal{X}|)}{\varepsilon^2}\right),$$

we have that \mathcal{A} in the presence of $(\text{cost}_{\text{add}}, \eta)$ -oblivious adversaries is (n, M, ε) -equivalent to \mathcal{A}_{sub} in the presence of $(\text{cost}_{\text{add}}, \eta)$ -adaptive adversaries.

Implicit in the proof of [Theorem 6](#), we proved the following.

Lemma 7 Fix a cost function ρ , budget $\eta \geq 0$, and $\varepsilon > 0$. Suppose that for an algorithm $\mathcal{A} : \mathcal{X}^n \rightarrow \{0, 1\}$, there are estimators $\text{Est-Max}_{\rho, \eta}, \text{Est-Min}_{\rho, \eta} : \mathcal{X}^m \rightarrow \{0, 1\}$ that use m samples from a distribution \mathcal{D} and returns estimates satisfying

$$\Pr_{S \sim \mathcal{D}^m} \left[\text{Est-Max}_{\rho, \eta}(S) = \text{Oblivious-Max}_{\rho, \eta}(\mathcal{A}, \mathcal{D}, n) \pm 2\varepsilon \right] \geq 1 - \varepsilon,$$

and likewise for Min instead of Max. Then, for

$$M = \frac{14m^2}{\varepsilon},$$

we have that \mathcal{A} in the presence of (ρ, η) -oblivious adversaries is $(n, M, 9\varepsilon)$ -equivalent to $\mathcal{A}_{\text{sub}} := \mathcal{A} \circ \Phi_{* \rightarrow n}$ in the presence of (ρ, η) -adaptive adversaries.

In order to prove [Theorem 7](#), we'll construct the estimator $\text{Est-Max}_{\text{cost}_{\text{add}}, \eta}$ (and likewise for Min). The goal is to estimate

$$\text{Oblivious-Max}_{\text{cost}_{\text{add}}, \eta}(\mathcal{A}, \mathcal{D}, n) := \sup_{\widehat{\mathcal{D}} = (1-\eta) \cdot \mathcal{D} + \eta \cdot \mathcal{E}} \left\{ \mathbb{E}_{S \sim \widehat{\mathcal{D}}^n} [\mathcal{A}(S)] \right\}.$$

The key insight is rather than trying all possible distributions \mathcal{E} , to compute a $\pm\varepsilon$ approximation, it suffices to consider those \mathcal{E} that are equal to $\mathcal{U}(T)$ for some $T \in \mathcal{X}^{n^2/\varepsilon}$. Our final result will have a logarithmic dependence on the number of \mathcal{E} we need to try, which results in just a logarithmic dependence on $|\mathcal{X}|$.

The following definition and accompany fact will be useful.

Definition 22 (Stochastic function) For any sets \mathcal{X}, \mathcal{Y} , a stochastic function $\mathbf{f} : \mathcal{X} \rightarrow \mathcal{Y}$ is a collection of distributions $\{\mathcal{D}_x \mid x \in \mathcal{X}\}$ each supported on \mathcal{Y} where the notation $\mathbf{f}(x)$ indicates an independent draw from \mathcal{D}_x .

Fact 6 For any distribution $\mathcal{D}_1, \mathcal{D}_2$ supported on a domain \mathcal{X} , and any stochastic function $\mathbf{f} : \mathcal{X} \rightarrow \mathcal{Y}$, let $\mathbf{f} \circ \mathcal{D}_i$ be the distribution where to sample $\mathbf{y} \sim \mathbf{f} \circ \mathcal{D}_i$, we first sample $\mathbf{x} \sim \mathcal{D}_i$ and then sample $\mathbf{y} = \mathbf{f}(\mathbf{x})$. Then,

$$\text{dist}_{\text{TV}}(\mathbf{f} \circ \mathcal{D}_1, \mathbf{f} \circ \mathcal{D}_2) \leq \text{dist}_{\text{TV}}(\mathcal{D}_1, \mathcal{D}_2).$$

Proof Given a coupling of $\mathbf{x}_1 \sim \mathcal{D}_1$ and $\mathbf{x}_2 \sim \mathcal{D}_2$, consider the coupling of $\mathbf{y}_1 \sim \mathbf{f} \circ \mathcal{D}_1$ and $\mathbf{y}_2 \sim \mathbf{f} \circ \mathcal{D}_2$ where if $\mathbf{x}_1 = \mathbf{x}_2$ then $\mathbf{y}_1 = \mathbf{y}_2 = \mathbf{f}(\mathbf{x}_1)$ and otherwise $\mathbf{y}_1 = \mathbf{f}(\mathbf{x}_1)$ and $\mathbf{y}_2 = \mathbf{f}(\mathbf{x}_2)$ independently. Then,

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

implying the desired result by [Definition 2](#). ■

To prove [Theorem 7](#), we will design the estimators $\text{Est-Max}_{\text{cost}_{\text{add}}, \eta}$ and $\text{Est-Min}_{\text{cost}_{\text{add}}, \eta}$. Plugging the sample complexity of those estimators into [Lemma 7](#) would be sufficient for [Theorem 7](#) with the slightly worse $M = n^6 \log(|\mathcal{X}|)^2 / \varepsilon^7$. To get the optimal $\log(|\mathcal{X}|)$ dependence on the domain size (and an improved dependence on n and ε), we need a more refined version of [Lemma 7](#) that takes advantage of the structure of the particular estimators we derive. The following Lemma applies for any cost function ρ , but we will then design \mathcal{F} satisfying [Equation \(8\)](#) specifically for $\rho = \text{cost}_{\text{add}}$ in [Lemma 9](#).

Lemma 8 Fix a cost function ρ , budget $\eta \geq 0$, $\varepsilon \in (0, 1]$, and algorithm $\mathcal{A} : \mathcal{X}^n \rightarrow \{0, 1\}$. Suppose there is a set of stochastic functions \mathcal{F} each $\mathcal{X} \rightarrow \mathcal{X}$, that satisfy, for any distribution \mathcal{D} over \mathcal{X} ,

$$\max_{\mathbf{f} \in \mathcal{F}} \left\{ \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^n} [\mathcal{A}(\mathbf{f}(\mathbf{S}))] \right\} = \text{Oblivious-Max}_{\rho, \eta}(\mathcal{A}, \mathcal{D}, n) \pm \varepsilon, \quad (8)$$

and likewise for Min instead of Max, where $\mathbf{f}(S)$ is shorthand for applying \mathbf{f} element wise and independently to S . Then:

1. There are estimators $\text{Est-Max}_{\text{cost}_{\text{add}}, \eta}$ and $\text{Est-Min}_{\text{cost}_{\text{add}}, \eta}$ meeting the requirements of [Lemma 7](#) for

$$m' = O\left(\frac{n \log(|\mathcal{F}|/\varepsilon)}{\varepsilon^2}\right).$$

In particular, this implies that for

$$M = O\left(\frac{n^2 \log(|\mathcal{F}|/\varepsilon)^2}{\varepsilon^5}\right),$$

we have that \mathcal{A} in the presence of (ρ, η) -oblivious adversaries is $(n, M, 9\varepsilon)$ -equivalent to \mathcal{A}_{sub} in the presence of (ρ, η) -adaptive adversaries.

2. More directly, for

$$M = O\left(\frac{n^2 \log(|\mathcal{F}|/\varepsilon)}{\varepsilon}\right) \quad (9)$$

we have that \mathcal{A} in the presence of (ρ, η) -oblivious adversaries is $(n, M, 5\varepsilon)$ -equivalent to \mathcal{A}_{sub} in the presence of (ρ, η) -adaptive adversaries.

Proof We first give $\text{Est-Max}_{\text{cost}_{\text{add}}, \eta}$ satisfying the first item. For $r := O\left(\frac{\log(|\mathcal{F}|/\varepsilon)}{\varepsilon^2}\right)$, let

$$\text{Est-Max}_{\text{cost}_{\text{add}}, \eta} = \max_{\mathbf{f} \in \mathcal{F}} \left\{ \underbrace{\mathbb{E}_{\text{trial} \in [r]} \left[\mathbb{E}_{\mathbf{S} \sim \mathcal{D}^n} [\mathcal{A}(\mathbf{f}(\mathbf{S}))] \right]}_{:= \text{Est}(\mathbf{f})} \right\}.$$

We note that as long as the samples are reused across different $\mathbf{f} \in \mathcal{F}$ that $m' = rn$ samples from \mathcal{D} suffices to compute the above expression. By Hoeffding's inequality, for any fixed $\mathbf{f} \in \mathcal{F}$, with probability at least $1 - \frac{\varepsilon}{T}$

$$\text{Est}(\mathbf{f}) = \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^n} [\mathcal{A}(\mathbf{f}(\mathbf{S}))] \pm \varepsilon.$$

By union bound, with probability at least $1 - \varepsilon$, the above holds for every $\mathbf{f} \in \mathcal{F}$. If so, $\text{Est-Max}_{\text{cost}_{\text{add}}, \eta}$ has the desired accuracy by [Equation \(8\)](#).

Next, we prove the M from [Equation \(9\)](#) suffices. By [Lemma 5](#), it suffices to prove that

$$\mathbb{E}_{\mathbf{S} \sim \mathcal{D}^M} [\text{Oblivious-Max}_{\rho, \eta}(\mathcal{A}, \mathcal{U}(\mathbf{S}), n)] = \text{Oblivious-Max}_{\rho, \eta}(\mathcal{A}, \mathcal{D}, n) \pm 5\varepsilon,$$

and likewise for Min instead of Max. Applying [Equation \(8\)](#) to both sides of the above equation, it is sufficient to prove that

$$\mathbb{E}_{\mathbf{S} \sim \mathcal{D}^M} \left[\max_{\mathbf{f} \in \mathcal{F}} \left\{ \underbrace{\mathbb{E}_{\mathbf{S}_n \sim \mathcal{U}(\mathbf{S})^n} [\mathcal{A}(\mathbf{f}(\mathbf{S}_n))] }_{:= g^{(\mathbf{f})}(\mathbf{S})} \right\} \right] = \max_{\mathbf{f} \in \mathcal{F}} \left\{ \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^n} [\mathcal{A}(\mathbf{f}(\mathbf{S}))] \right\} \pm 3\varepsilon. \quad (10)$$

Fix a single $\mathbf{f} \in \mathcal{F}$. We'll use McDiarmid's inequality ([Fact 4](#)) to say that $g^{(\mathbf{f})}(\mathbf{S})$ concentrates around its mean. Take any $S \in \mathcal{X}^M$ and suppose we change one point in it to create S' . Then $|g^{(\mathbf{f})}(S) - g^{(\mathbf{f})}(S')|$ is at most the probability that the changed point appears in \mathbf{S}_n , which is at most $\frac{n}{M}$. Applying McDiarmid's inequality,

$$\begin{aligned} \Pr_{\mathbf{S} \sim \mathcal{D}^M} \left[g^{(\mathbf{f})}(\mathbf{S}) = \mu^{(\mathbf{f})} \pm \varepsilon \right] &\geq 1 - 2 \exp\left(-\frac{2\varepsilon^2}{(n/M)^2 M}\right) \quad \text{where } \mu^{(\mathbf{f})} := \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^M} [g^{(\mathbf{f})}(\mathbf{S})] \\ &\geq 1 - \frac{\varepsilon}{|\mathcal{F}|}. \quad \text{(using } M = O\left(\frac{n^2 \log(|\mathcal{F}|/\varepsilon)}{\varepsilon}\right)) \end{aligned}$$

By union bound, with probability at least $1 - \varepsilon$, for every $\mathbf{f} \in \mathcal{F}$ we have that $g^{(\mathbf{f})}(\mathbf{S}) = \mu^{(\mathbf{f})} \pm \varepsilon$ allowing us to bound the left hand side of [Equation \(10\)](#),

$$\begin{aligned} \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^M} \left[\max_{\mathbf{f} \in \mathcal{F}} \left\{ g^{(\mathbf{f})}(\mathbf{S}) \right\} \right] &= \max_{\mathbf{f} \in \mathcal{F}} \left\{ \mu^{(\mathbf{f})} \right\} \pm \varepsilon \pm \Pr_{\mathbf{S} \sim \mathcal{D}^M} \left[g^{(\mathbf{f})}(\mathbf{S}) \neq \mu^{(\mathbf{f})} \pm \varepsilon \text{ for some } \mathbf{f} \in \mathcal{F} \right] \\ &= \max_{\mathbf{f} \in \mathcal{F}} \left\{ \mu^{(\mathbf{f})} \right\} \pm 2\varepsilon. \end{aligned}$$

Lastly, we want to compare the above to the right hand side of Equation (10) to $\max_{\mathbf{f} \in \mathcal{F}} \mu^{(\mathbf{f})}$. Fix any $\mathbf{f} \in \mathcal{F}$. Using the notation of Lemma 3,

$$\mu^{(\mathbf{f})} = \mathbb{E}_{\mathbf{S} \sim \Phi_{M \rightarrow n} \circ \mathcal{D}^M} [\mathcal{A}(\mathbf{f}(\mathbf{S}))].$$

Therefore,

$$\begin{aligned} \left| \max_{\mathbf{f} \in \mathcal{F}} \left\{ \mu^{(\mathbf{f})} \right\} - \max_{\mathbf{f} \in \mathcal{F}} \left\{ \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^n} [\mathcal{A}(\mathbf{f}(\mathbf{S}))] \right\} \right| &\leq \max_{\mathbf{f} \in [\mathcal{F}]} \left| \mu^{(\mathbf{f})} - \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^n} [\mathcal{A}(\mathbf{f}(\mathbf{S}))] \right| \\ &\leq \text{dist}_{\text{TV}}(\mathcal{D}^n, \Phi_{M \rightarrow n} \circ \mathcal{D}^M) \quad (\text{Definition 2}) \\ &\leq \frac{\binom{n}{2}}{M} \leq \varepsilon. \quad (\text{Lemma 3}) \end{aligned}$$

Therefore, the left hand term of Equation (10) is within $\pm 2\varepsilon$ of $\max_{\mathbf{f} \in \mathcal{F}} \mu^{(\mathbf{f})}$ and the right hand term is within $\pm \varepsilon$ of $\max_{\mathbf{f} \in \mathcal{F}} \mu^{(\mathbf{f})}$. Hence, Equation (10) holds, completing this proof. \blacksquare

In order to use Lemma 8, we need to design \mathcal{F} and prove that it satisfies Equation (8) when $\rho = \text{cost}_{\text{add}}$. The below lemma completes the proof of Theorem 7.

Lemma 9 Fix a budget $\eta \geq 0$, $\varepsilon \in (0, 1]$, and algorithm $\mathcal{A} : \mathcal{X}^n \rightarrow \{0, 1\}$. Let \mathcal{F} be the set of $|\mathcal{X}|^{n^2/\varepsilon}$ stochastic functions defined below.

$$\mathcal{F} := \{\mathbf{f}^{(T)} \mid T \in \mathcal{X}^{n^2/\varepsilon}\} \quad \text{where} \quad \mathbf{f}^{(T)}(x) := \begin{cases} x & \text{with probability } 1 - \eta \\ \mathbf{y} \text{ where } \mathbf{y} \sim \mathcal{U}(T) & \text{with probability } \eta. \end{cases}$$

Then for any \mathcal{D} over \mathcal{X} ,

$$\begin{aligned} \max_{\mathbf{f} \in \mathcal{F}} \left\{ \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^n} [\mathcal{A}(\mathbf{f}(\mathbf{S}))] \right\} &= \text{Oblivious-Max}_{\text{cost}_{\text{add}}, \eta}(\mathcal{A}, \mathcal{D}, n) \pm \varepsilon \\ &= \sup_{\widehat{\mathcal{D}} = (1-\eta) \cdot \mathcal{D} + \eta \cdot \mathcal{E}} \left\{ \mathbb{E}_{\mathbf{S} \sim \widehat{\mathcal{D}}^n} [\mathcal{A}(\mathbf{S})] \right\} \pm \varepsilon, \end{aligned}$$

and likewise for Min instead of Max.

Proof Fix any distribution \mathcal{D} . First, we note that the distribution of $\mathbf{f}^{(T)}(x)$ where $x \sim \mathcal{D}$ is simply that of $(1 - \eta) \cdot \mathcal{D} + \eta \cdot \mathcal{U}(T)$. Therefore, one direction of the desired result is easy:

$$\max_{\mathbf{f} \in \mathcal{F}} \left\{ \mathbb{E}_{\mathbf{S} \sim \mathcal{D}^n} [\mathcal{A}(\mathbf{f}(\mathbf{S}))] \right\} \leq \text{Oblivious-Max}_{\text{cost}_{\text{add}}, \eta}(\mathcal{A}, \mathcal{D}, n).$$

The remainder of this proof is devoted to proving the left hand side of the above equation is at most ε smaller than the right hand side. Fix any distribution \mathcal{E} and consider $\widehat{\mathcal{D}} = (1 - \eta) \cdot \mathcal{D} + \eta \cdot \mathcal{E}$. We'll show that

$$\mathbb{E}_{\mathbf{T} \sim \mathcal{E}^{n^2/\varepsilon}} \left[\mathbb{E}_{\mathbf{S} \sim \mathcal{D}^n} [\mathcal{A}(\mathbf{f}^{\mathbf{T}}(\mathbf{S}))] \right] \geq \mathbb{E}_{\mathbf{S} \sim \widehat{\mathcal{D}}^n} [\mathcal{A}(\mathbf{S})] - \varepsilon. \quad (11)$$

In particular, the above implies there is a single choice for T that is within ε of $\mathbb{E}_{\mathcal{S} \sim \widehat{\mathcal{D}}^n}[\mathcal{A}(\mathcal{S})]$. As this holds for all corruptions $\widehat{\mathcal{D}}$, it implies the desired result.

To sample from $\widehat{\mathcal{D}}^n$ where $\widehat{\mathcal{D}} = (1 - \eta) \cdot \mathcal{D} + \eta \cdot \mathcal{E}$ we can first draw $\mathcal{S} \sim \mathcal{E}^n$ and then return $\mathbf{h}(\mathcal{S})$, applied element wise, where:

$$\mathbf{h}(x) := \begin{cases} \mathbf{y} \text{ where } \mathbf{y} \sim \mathcal{D} & \text{with probability } 1 - \eta \\ x & \text{with probability } \eta. \end{cases}$$

Therefore, the distribution of $\mathbf{f}^T(\mathcal{S})$ on the left hand side of Equation (11) is $\mathbf{h} \circ \Phi_{n^2/\varepsilon \rightarrow n} \circ \mathcal{E}^{n^2/\varepsilon}$ (using the notation of Lemma 3 and Fact 6). Finally, we prove Equation (11) (recalling that $\widehat{\mathcal{D}} = (1 - \eta) \cdot \mathcal{D} + \eta \cdot \mathcal{E}$)

$$\begin{aligned} \left| \mathbb{E}_{\mathcal{T} \sim \mathcal{E}^{n^2/\varepsilon}} \left[\mathbb{E}_{\mathcal{S} \sim \widehat{\mathcal{D}}^n} [\mathcal{A}(\mathbf{f}^T(\mathcal{S}))] \right] - \mathbb{E}_{\mathcal{S} \sim \widehat{\mathcal{D}}^n} [\mathcal{A}(\mathcal{S})] \right| &\leq \text{dist}_{\text{TV}}(\mathbf{h} \circ \Phi_{n^2/\varepsilon \rightarrow n} \circ \mathcal{E}^{n^2/\varepsilon}, \mathbf{h} \circ \mathcal{E}^n) \\ &\leq \text{dist}_{\text{TV}}(\Phi_{n^2/\varepsilon \rightarrow n} \circ \mathcal{E}^{n^2/\varepsilon}, \mathcal{E}^n) \\ &\leq \varepsilon. \end{aligned} \quad \begin{array}{l} \text{(Definition 2)} \\ \text{(Fact 6)} \\ \text{(Lemma 3)} \end{array}$$

■

Appendix E. Proof of Theorem 4: Lower bounds against the subsampling filter

In this section, we show a lower bound on m needed for the subsampling filter to work. Our lower bound holds in the setting of additive noise and therefore also shows that the dependence on $|\mathcal{X}|$ in Theorem 7 is optimal.

Theorem 8 (Formal version of Theorem 4) *For any sample size n , domain \mathcal{X} with $|\mathcal{X}| = 2^d$ for an integer d , adversary budget η , and $\varepsilon > 0$, there exists an algorithm $\mathcal{A} : \mathcal{X}^n \rightarrow \{0, 1\}$ and corresponding subsampled algorithm $\mathcal{A}_{\text{sub}} := \mathcal{A} \circ \Phi_{* \rightarrow n}$ for which \mathcal{A} in the presence of $(\text{cost}_{\text{add}}, \eta)$ -oblivious adversaries is not $(n, m, 1 - \varepsilon)$ -equivalent to \mathcal{A}_{sub} in the presence of $(\text{cost}_{\text{add}}, \eta)$ -adaptive adversaries for any $m = O_\eta(n \log |\mathcal{X}| / \log^2 n)$.*

Proof overview Without loss of generality, we can consider the domain to be the Boolean hypercube, $\mathcal{X} = \{\pm 1\}^d$. Otherwise, we could map the domain to the hypercube. For an appropriate threshold t , we'll define

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

Let \mathcal{D} be uniform over \mathcal{X} . For any $\varepsilon > 0$, n , and d , we'll show that there is a choice of t such that:

1. **Lemma 10:** Oblivious- $\text{Max}_{\text{cost}_{\text{add}}, \eta}(\mathcal{A}, \mathcal{D}, n) \leq \varepsilon/2$, meaning, for any $\widehat{\mathcal{D}} = (1 - \eta)\mathcal{D} + \eta\mathcal{E}$, it is the case that $\mathbb{E}_{\mathcal{S} \sim \widehat{\mathcal{D}}^n}[\mathcal{A}(\mathcal{S})] \leq \varepsilon/2$.

2. **Lemma 11:** Adaptive- $\text{Max}_{\text{cost}_{\text{add}}, \eta}(\mathcal{A}, \mathcal{D}, m) \geq 1 - \frac{\varepsilon}{2}$ whenever $m = O_\eta(n \log |\mathcal{X}| / \log^2 n)$, meaning that for $\mathbf{S} \sim \mathcal{D}^m$, the adaptive adversary can choose $\lfloor m \cdot \eta / (1 - \eta) \rfloor$ points \mathbf{T} to add to the sample for which

$$\mathbb{E}[\mathcal{A}_{\text{sub}}(\widehat{\mathbf{S}})] \geq 1 - \frac{\varepsilon}{2} \quad \text{where } \widehat{\mathbf{S}} = \mathbf{S} \cup \mathbf{T}.$$

Together, these prove that \mathcal{A} in the presence of $(\text{cost}_{\text{add}}, \eta)$ -oblivious adversaries is not $(n, m, 1 - \varepsilon)$ -equivalent to \mathcal{A}_{sub} in the presence of $(\text{cost}_{\text{add}}, \eta)$ -adaptive adversaries.

Lemma 10 For any distribution \mathcal{E} and $\widehat{\mathcal{D}} = (1 - \eta)\mathcal{D} + \eta\mathcal{E}$,

$$\mathbb{E}_{\mathbf{S} \sim \widehat{\mathcal{D}}^n} [\mathcal{A}(\mathbf{S})] < \varepsilon/2.$$

Proof First, we note that for any $x' \in \mathcal{X}$ and a clean sample $\mathbf{x} \sim \mathcal{D}$, the probability that $\langle \mathbf{x}, \mathbf{x}' \rangle \geq t$ is small: by Hoeffding's inequality, $\Pr_{\mathbf{x} \sim \mathcal{D}}[\langle \mathbf{x}, \mathbf{x}' \rangle \geq t] \leq \exp(-t^2/2d)$. Combining this with a simple union bound, we can show that the probability of even a single clean point forming a correlated pair in the sample is small:

$$\begin{aligned} \mathbb{E}_{\mathbf{S} \sim \widehat{\mathcal{D}}^n} [\mathcal{A}(\mathbf{S})] &\leq \eta^n + \mathbb{E}_{\mathbf{S} \sim \widehat{\mathcal{D}}^n} [\mathcal{A}(\mathbf{S}) \mid \text{at least one clean point in } \mathbf{S}] \\ &= \eta^n + \mathbb{E}_{\substack{\mathbf{x} \sim \mathcal{D} \\ \mathbf{S} \sim \widehat{\mathcal{D}}^{n-1}}} [\mathcal{A}(\mathbf{S} \cup \{\mathbf{x}\})] \\ &\leq \eta^n + \Pr_{\substack{\mathbf{x} \sim \mathcal{D} \\ \mathbf{S} \sim \widehat{\mathcal{D}}^{n-1}}} [\exists \mathbf{x}' \in \mathbf{S} \text{ with } \langle \mathbf{x}, \mathbf{x}' \rangle \geq t] \quad (\text{weaken to just one clean point}) \\ &\leq \eta^n + (n-1) \Pr_{\substack{\mathbf{x} \sim \mathcal{D} \\ \mathbf{x}' \sim \widehat{\mathcal{D}}}} [\langle \mathbf{x}, \mathbf{x}' \rangle \geq t] \quad (\text{union bound over } \mathbf{S}) \\ &\leq \eta^n + n \exp\left(-\frac{t^2}{2d}\right). \quad (\text{Hoeffding's, over the randomness of } \mathbf{x}) \end{aligned}$$

This will be vanishingly small for the particular choice of t determined in [Lemma 11](#). ■

Lemma 11 For any $m = O_\eta(nd / \log^2 n)$, there exists an adversarial strategy that, given $\mathbf{S} \in \mathcal{D}^m$, chooses $\lfloor m \cdot \eta / (1 - \eta) \rfloor$ points \mathbf{T} to add to the sample for which

$$\mathbb{E}[\mathcal{A}_{\text{sub}}(\widehat{\mathbf{S}})] \geq 1 - \varepsilon/2 \quad \text{where } \widehat{\mathbf{S}} = \mathbf{S} \cup \mathbf{T}$$

Proof Let $C = \lfloor m \cdot \eta / (1 - \eta) \rfloor$ be the number of points the adversary can add and denote the sample $\mathbf{S} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$. The adversary constructs $\mathbf{T} = \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(C)}\}$ by setting each $\mathbf{y}^{(j)}$ to be the elementwise majority of k chosen points from \mathbf{S} (where k will be determined later). The idea is that $\mathbf{y}^{(j)}$ is a cluster center that will form a correlated pair with every one of these k points, with high probability.

More formally, for each $j = 1, \dots, C$, define $\mathbf{S}^{(j)} = \{x^{(1+(j-1)k \pmod{|\mathbf{S}|})}, \dots, x^{(jk \pmod{|\mathbf{S}|})}\}$ to be the j -th chunk of k points from \mathbf{S} , with the indices wrapping around to the start of \mathbf{S} as necessary. We take $\mathbf{y}^{(j)}$ to be the elementwise majority⁹ of the points in $\mathbf{S}^{(j)}$:

⁹ We can assume k is odd for simplicity.

$$\mathbf{y}_\ell^{(j)} := \operatorname{Maj}_{\mathbf{x} \in \mathcal{S}^{(j)}} \{\mathbf{x}_\ell\} \quad \text{for } \ell = 1, \dots, d. \quad (12)$$

First, we note that for a given $\mathbf{y}^{(j)}$, with high probability, any point in $\mathbf{x} \in \mathcal{S}^{(j)}$ will have a large dot product with $\mathbf{y}^{(j)}$:

$$\begin{aligned} \mathbb{E}_{\mathcal{S}^{(j)} \sim \mathcal{D}^k} [\langle \mathbf{x}, \mathbf{y}^{(j)} \rangle] &= \sum_{\ell=1}^d \mathbb{E}_{\mathcal{S}^{(j)} \sim \mathcal{D}^k} [\mathbf{x}_\ell \mathbf{y}_\ell^{(j)}] \\ &= \sum_{\ell=1}^d \left(\Pr_{\mathcal{S}^{(j)} \sim \mathcal{D}^k} [\mathbf{x}_\ell = \mathbf{y}_\ell^{(j)}] - \Pr_{\mathcal{S}^{(j)} \sim \mathcal{D}^k} [\mathbf{x}_\ell \neq \mathbf{y}_\ell^{(j)}] \right) \\ &= \sum_{\ell=1}^d \left(\Pr_{\mathbf{u} \sim \operatorname{Bin}(k-1, \frac{1}{2})} \left[\mathbf{u} \geq \frac{(k-1)}{2} \right] - \Pr_{\mathbf{u} \sim \operatorname{Bin}(k-1, \frac{1}{2})} \left[\mathbf{u} < \frac{(k-1)}{2} \right] \right) \\ &= d \Pr_{\mathbf{u} \sim \operatorname{Bin}(k-1, \frac{1}{2})} \left[\mathbf{u} = \frac{(k-1)}{2} \right] \\ &= \sqrt{\frac{2}{\pi}} \frac{d}{\sqrt{k}} (1 \pm o(1)). \end{aligned}$$

Define $\mu := (\sqrt{2/\pi})(d/\sqrt{k})$. We will take $t = \mu/2$ as the threshold for \mathcal{A} , both here and in [Lemma 10](#) as well. For $\mathcal{S}^{(j)} \sim \mathcal{D}^k$, with $\mathbf{y}^{(j)}$ the elementwise majority of $\mathcal{S}^{(j)}$, and any $\mathbf{x} \in \mathcal{S}^{(j)}$, this gives:

$$\begin{aligned} \Pr_{\mathcal{S}^{(j)} \sim \mathcal{D}^k} [\langle \mathbf{x}, \mathbf{y}^{(j)} \rangle < t] &= \Pr_{\mathcal{S}^{(j)} \sim \mathcal{D}^k} [\langle \mathbf{x}, \mathbf{y}^{(j)} \rangle < \frac{\mu}{2}] \\ &\leq \exp \left[-\Theta \left(\frac{\mu^2}{d} \right) \right] \quad (\text{Hoeffding's inequality}) \\ &= \exp \left[-\Theta \left(\frac{d}{k} \right) \right]. \quad (13) \end{aligned}$$

The subsampling filter $\Phi_{* \rightarrow n}$ takes a random subsample of size n from $\widehat{\mathcal{S}} = \mathcal{S} \cup \mathcal{T}$ (with replacement). We want to show, with high probability over size n subsamples $\mathcal{S}' \sim \mathcal{U}(\widehat{\mathcal{S}})^n$, that $\mathcal{A}(\mathcal{S}') = 1$. For any point $\mathbf{x} \in \mathcal{S}$, we say $\mathbf{y}^{(j)} \in \mathcal{T}$ is “good” for \mathbf{x} if \mathbf{x} was in the cluster used to compute $\mathbf{y}^{(j)}$, meaning $\mathbf{x} \in \mathcal{S}^{(j)}$. Similarly, for any $\mathbf{y}^{(j)} \in \mathcal{T}$, we say that $\mathbf{x} \in \mathcal{S}$ is “good” for $\mathbf{y}^{(j)}$ if $\mathbf{x} \in \mathcal{S}^{(j)}$. By construction, a given $\mathbf{x} \in \mathcal{S}$ participates in the computation of at least $\lfloor Ck/m \rfloor = \Theta_\eta(k)$ many cluster centers $\mathbf{y}^{(j)}$'s, and so for each $\mathbf{x} \in \mathcal{S}$, there are $\Theta_\eta(k)$ good \mathbf{y} 's $\in \mathcal{T}$. Similarly, there are exactly k good \mathbf{x} 's $\in \mathcal{S}$ for each $\mathbf{y} \in \mathcal{T}$. As $\widehat{\mathcal{S}}$ has size $\Theta_\eta(m)$, for any $\mathbf{x} \in \mathcal{S}'$, using c_x to denote the number of good points for \mathbf{x} that show up in \mathcal{S}' , we have that

\mathbf{c}_x is distributed as $\text{Bin}(n, \Omega_\eta(k/m))$. This gives:

$$\begin{aligned}
 \Pr_{\mathbf{S}' \sim \mathcal{U}(\widehat{\mathcal{S}})^n} [\mathcal{A}(\mathbf{S}') = 0] &= \Pr_{\mathbf{S}' \sim \mathcal{U}(\widehat{\mathcal{S}})^n} [\exists \mathbf{x} \in \mathbf{S}' \text{ s.t. } \forall \mathbf{x}' \in \mathbf{S}', \langle \mathbf{x}, \mathbf{x}' \rangle < t] \\
 &\leq n \Pr_{\mathbf{x}, \mathbf{S}' \sim \mathcal{U}(\widehat{\mathcal{S}})^n} [\forall \mathbf{x}' \in \mathbf{S}', \langle \mathbf{x}, \mathbf{x}' \rangle < t] \quad (\text{union bound}) \\
 &\leq n \left[\Pr[\mathbf{c}_x = 0] + \Pr_{\mathbf{x}, \mathbf{S}' \sim \mathcal{U}(\widehat{\mathcal{S}})^n} [\forall \mathbf{x}' \in \mathbf{S}', \langle \mathbf{x}, \mathbf{x}' \rangle < t \mid \mathbf{c}_x \geq 1] \right] \\
 &\leq n \left[\Pr[\mathbf{c}_x = 0] + \Pr_{\mathbf{x}, \mathbf{S}' \sim \mathcal{U}(\widehat{\mathcal{S}})^n} [\langle \mathbf{x}, \mathbf{y}_x \rangle < t] \right] \quad (\text{weaken to good point}) \\
 &\leq n \left(1 - \Theta_\eta \left(\frac{k}{m} \right) \right)^n + n \exp \left[-\Theta_\eta \left(\frac{d}{k} \right) \right]. \quad (\text{from Equation (13)})
 \end{aligned}$$

■

Setting of parameters. To complete the proof of [Theorem 8](#), we need to choose m and k so that, all the terms in [Lemma 10](#) and [Lemma 11](#) are vanishingly small. In particular, we need $n(1 - \Theta_\eta(k/m))^n \rightarrow 0$ and $n \exp(-\Theta_\eta(d/k)) \rightarrow 0$ as $n \rightarrow \infty$. If $k = \Omega_\eta(m \log n/n)$ and $d = \Omega_\eta(k \log n) = \Omega_\eta(m \log^2 n/n)$ with sufficiently large constant factors, we get the desired result. In other words, as long as $m \leq O_\eta(nd/\log^2 n)$, the adaptive adversary is stronger than the oblivious adversary.

Appendix F. Proof of [Fact 2](#)

Proof By a standard inductive argument, [Definition 13](#) implies that for any $m \in \mathbb{N}$, weights $\theta_1, \dots, \theta_m \geq 0$ summing to 1 and distributions $\mathcal{D}_1, \dots, \mathcal{D}_m, \widehat{\mathcal{D}}_1, \dots, \widehat{\mathcal{D}}_m$, that

$$\rho \left(\sum_{i \in [m]} \theta_i \mathcal{D}_i, \sum_{i \in [m]} \theta_i \widehat{\mathcal{D}}_i \right) \leq \max_{i \in [m]} \left\{ \rho(\mathcal{D}_i, \widehat{\mathcal{D}}_i) \right\}.$$

The distribution \mathcal{D} and $\widehat{\mathcal{D}}$ can be written as the mixtures

$$\begin{aligned}
 \mathcal{D} &= \sum_{S \in \mathcal{X}^n} \Pr_{\mathbf{S} \sim \mathcal{D}^n} [\mathbf{S} = S] \mathcal{U}(S), \\
 \widehat{\mathcal{D}} &= \sum_{S \in \mathcal{X}^n} \Pr_{\mathbf{S} \sim \widehat{\mathcal{D}}^n} [\mathbf{S} = S] \mathcal{U}(\widehat{S}).
 \end{aligned}$$

Since $\rho(\mathcal{U}(S), \mathcal{U}(\widehat{S})) \leq \eta$ for every $S \in \mathcal{X}^n$, we can conclude $\rho(\mathcal{D}, \widehat{\mathcal{D}}) \leq \eta$. ■