

# Private Linear Regression via a Down-Sensitivity to Privacy Reduction

**Ittai Rubinstein**

**Chris Ge**

**Samuel B. Hopkins**

*Electrical Engineering and Computer Science  
Massachusetts Institute of Technology  
Cambridge, MA 02139, USA*

ITTAIR@MIT.EDU

CGE7@MIT.EDU

SAMHOP@MIT.EDU

**Editors:** Steve Hanneke and Tor Lattimore

## Abstract

We present a sample- and time-efficient  $(\epsilon, \delta)$ -differentially private (DP) algorithm for  $d$ -dimensional linear regression with a sample complexity of

$$n_{\text{STAR}} = \tilde{O} \left( \frac{d}{\alpha^2} + \frac{d \log(1/\delta)}{\alpha \epsilon} + \frac{d \log(1/\delta)}{\epsilon} \right) + o(d).$$

This improves upon prior polynomial-time algorithms whose sample complexity either depends on the condition number of the design matrix  $\kappa$  (for DP-SGD with gradient clipping), scales quadratically with the dimension (for Sum-of-Squares algorithms) or with the inverse of the privacy parameter (for outlier removal algorithms such as insufficient statistics perturbation or ISSP),

$$n_{\text{SOS}} = \tilde{\Omega} \left( \frac{d^2}{\alpha^2} \right), \quad n_{\text{DP-SGD}} = \tilde{\Omega} \left( \frac{d\sqrt{\kappa}}{\epsilon} \right), \quad n_{\text{ISSP}} = \tilde{\Omega} \left( \frac{d}{\epsilon^2} \right).$$

Our algorithm is based on a novel *subsample-test-aggregate* (STA) approach for ensuring privacy given only bounded *down-sensitivity* – robustness to removal, but not addition, of a small number of samples. The intuition that down-sensitivity should be related to privacy is not new, but STA formalizes this by providing an *efficient black-box reduction from down-sensitivity to privacy* which we expect to be applicable beyond the setting of linear regression.

**Keywords:** differential privacy, linear regression, OLS, high-dimensional

## 1. Introduction

As the fields of machine learning, medicine and econometrics become ever more data-driven, the problem of ensuring the privacy of individuals also grows in importance. *Differential privacy* (DP) is the gold standard mathematical definition for an algorithm to protect privacy of individuals represented in its input [Dwork et al. \(2006b\)](#).

Despite two decades of intensive research, DP algorithms for foundational tasks in statistics and machine learning lag far behind their non-private counterparts in accuracy, running time, or both. Indeed, this is true even for basic tasks such as mean estimation, covariance estimation, and linear regression [Alabi et al. \(2020\)](#); [Barrientos et al. \(2024\)](#). One of the main challenges associated with privatizing learning algorithms for these problems is their sensitivity to extreme / adversarial samples. In this paper, we introduce STA (**S**ubsample-**T**est-**A**ggregate) – a novel approach to making learners private even if the underlying statistic may be highly sensitive to the addition of extreme samples.

As a concrete application of STA, we introduce STAR (STA **R**egression) – an  $(\epsilon, \delta)$ -DP algorithm for ordinary least squares (OLS) linear regression, with runtime comparable to non-private linear regression (up to logarithmic factors) and a significantly lower sample complexity than previous polynomial-time approaches. Measured in terms of the number of samples required to obtain generalization error  $\alpha$  when the covariates are drawn from an unknown Gaussian distribution, previous polynomial-time algorithms require a number of samples that depends on the condition number of the covariates  $\kappa$ , or scales quadratically either with the privacy parameter or with the dimension [Anderson et al. \(2025\)](#); [Liu et al. \(2023\)](#); [Brown et al. \(2024\)](#):

$$n_{\text{PRIOR ALGORITHMS}} \geq \tilde{\Omega} \left( \min \left\{ \frac{d^2}{\alpha^2}, \frac{d\sqrt{\kappa}}{\epsilon}, \frac{d}{\epsilon^2} \right\} \right).$$

By contrast, in the same Gaussian-design setting, STAR can be run with a sample complexity of

$$n_{\text{STAR}} = \tilde{O} \left( \frac{d}{\alpha^2} + \frac{d \log(1/\delta)}{\alpha \epsilon} + \frac{d \log(1/\delta)}{\epsilon} \right) + o(d).$$

The privacy guarantees of both STAR and these prior algorithms hold regardless of the dataset and although we expect their generalization error to behave similarly for a much wider class of data distributions, we restrict this part of our analysis to the Gaussian setting to keep our evaluation in line with the existing literature [Brown et al. \(2024\)](#); [Anderson et al. \(2025\)](#).

**Outliers Are the Main Challenge.** If our underlying learner is robust in the sense that the resulting model cannot be changed significantly by altering a single input, then it can be made private by adding appropriately calibrated noise to the learned model [Dwork and Roth \(2014\)](#).

The key challenge with privatizing linear regression is dealing with “outliers”: samples, or configurations of samples, whose presence can have a large effect on the learned model (see [Figure 1](#)). This motivates the question:

*How can we privatize an algorithm that is sensitive to outliers?*

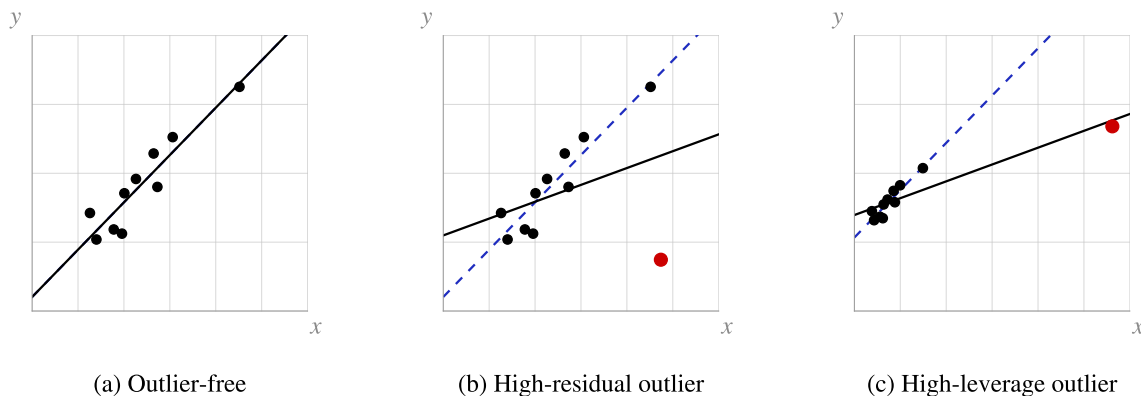


Figure 1: Sensitivity of linear regression to high-leverage and high-residual outliers.

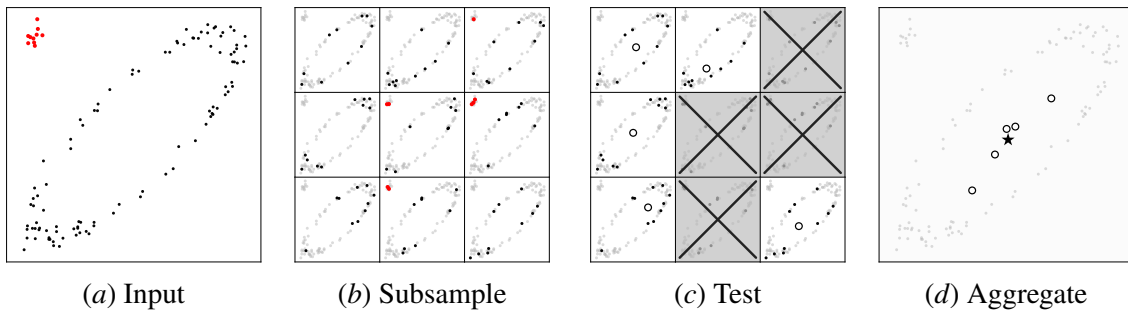


Figure 2: Illustration of our Subsample–Test–Aggregate Algorithm 1 for the case of mean estimation. We plot the input **inliers** in black and the **outliers** in red. Starting from the original “corrupted” dataset, the algorithm draws small random subsamples, tests each subsample for outliers, and aggregates only the estimates from subsamples that pass the test. Failed subsamples are greyed out and crossed out, white circles denote the empirical means of accepted subsamples, and the black star denotes the final aggregate. In the STA algorithm, we also add noise to the estimator before release; we omit that step here for simplicity.

Even though it is often relatively simple to test for the existence of outliers, it is much harder to identify *which samples are the outliers*. For instance, we might say that a linear regression is outlier-free if all of its samples have bounded leverages and residuals, but the outliers are not always the samples with the highest leverages or residuals [Huang et al. \(2024\)](#); [Hu et al. \(2024\)](#).

This is the main difficulty limiting previous private regression algorithms, which attempt to find the largest-possible outlier-free subsets of a dataset via iterative or convex-programming methods [Liu et al. \(2023\)](#); [Brown et al. \(2024\)](#). We introduce a much simpler solution: we subsample the dataset (i.e., draw a small random subset of the training set), then test each subsample to retain only outlier-free subsamples, and output an aggregate of models trained on the outlier-free subsamples (see Figure 2).

**Theorem 1 (Subsample-Test-Aggregate – Informal (see Theorem 30))** *Any learning algorithm  $A$ , outlier detection test  $\mathcal{T}$ , and aggregator  $\text{AGG}$ , can be converted into a private Subsample-Test-Aggregate estimator for  $A$ .*

The key conceptual difference between STA and approaches based on iterative outlier removal is that the test only needs to solve a *decision* problem (does the dataset contain outliers or not?), and not a *search* problem (which samples are the outliers?).

STAR shows that the simplest of all outlier removal schemes – subsampling and conditioning on the random subsample containing no outliers – suffices for private linear regression with  $\tilde{O}(d/\epsilon)$  samples. STA also bears some similarity to recent down-sensitivity based mechanisms (such as [Cummings and Durfee \(2020\)](#); [Linder et al. \(2025\)](#); [Steinke and Steinke \(2025\)](#)), but unlike these previous approaches, STA both supports learners with high-dimensional outputs and allows for aggregation of the output of the underlying learner across many subsamples which increases the accuracy of the private learner with the same privacy budget. See Section 1.3 for a more detailed comparison.

### 1.1. Results

We now describe our results for linear regression more formally; we discuss the STA framework in Section 1.2 after describing more background. Recall the definition of  $(\epsilon, \delta)$ -differential privacy:

**Definition 2** ( $(\epsilon, \delta)$ -Differential Privacy, **Dwork et al. (2006a)**) *A (randomized) algorithm  $\mathcal{A}$  with domain  $\mathcal{X}^n$  satisfies  $(\epsilon, \delta)$ -differential privacy if for all pairs of neighboring datasets  $X, X' \in \mathcal{X}^n$  that differ in at most one sample, and for all measurable subsets  $S$  of the output space of  $\mathcal{A}$ ,*

$$\Pr[\mathcal{A}(X) \in S] \leq e^\epsilon \Pr[\mathcal{A}(X') \in S] + \delta.$$

In keeping with much recent work on DP statistics, e.g. [Karwa and Vadhan \(2017\)](#); [Brown et al. \(2024, 2023, 2021\)](#); [Liu et al. \(2022\)](#); [Hopkins et al. \(2023\)](#); [Georgiev and Hopkins \(2022\)](#); [Narayanan et al. \(2022\)](#); [Asi et al. \(2023\)](#); [Kuditipudi et al. \(2023\)](#); [Cai et al. \(2021\)](#); [Kamath et al. \(2019\)](#); [Arbas et al. \(2023\)](#); [Biswas et al. \(2020\)](#); [Kamath et al. \(2020\)](#); [Ashtiani and Liaw \(2022\)](#); [Anderson et al. \(2025\)](#); [Dagan et al. \(2024\)](#), we study algorithms which satisfy DP with respect to all inputs and have good prediction/estimation error when the data is drawn from a “nice” distribution such as [Definition 3](#) (i.e., *worst-case* privacy guarantees, *average-case* utility guarantees).<sup>1</sup>

**Definition 3** *Let  $\Sigma \succ 0$  with  $\Sigma \in \mathbb{R}^{d \times d}$ ,  $\sigma^2 > 0$ , and  $\beta^* \in \mathbb{R}^d$ . The  $n$ -sample linear Gaussian model consists of  $n$  i.i.d. draws  $(x_1, y_1), \dots, (x_n, y_n)$  where*

$$x_i \sim \mathcal{N}(0, \Sigma) \text{ and } y_i \mid x_i \sim \mathcal{N}(x_i^\top \beta^*, \sigma^2).$$

*We measure the performance of a linear regression estimator  $\hat{\beta}$  by its (normalized) excess risk,*

$$\alpha := \frac{1}{\sigma} \sqrt{\left( \mathbb{E}_{(x,y)} \left[ (x^\top \hat{\beta} - y)^2 \right] - \mathbb{E}_{(x,y)} \left[ (x^\top \beta^* - y)^2 \right] \right)} = \frac{1}{\sigma} \left\| \Sigma^{1/2} (\hat{\beta} - \beta^*) \right\|.$$

We introduce a time and sample efficient  $(\epsilon, \delta)$ -DP regression algorithm, STAR (*Subsample-Test-Aggregate Regression*). We state its main guarantees here and describe the algorithm in [Section 1.2](#).

**Theorem 4 (Privacy and utility of STAR)** *STAR ([Algorithm 6](#)) is  $(\epsilon, \delta)$ -DP, runs in time  $\tilde{O}(nd^2(1 + \lambda))$ , has sample complexity*

$$n_{\text{STAR}} = \tilde{O} \left( \frac{d}{\alpha^2} + \frac{d(1 + \lambda) \log(1/\delta)}{\alpha \epsilon} + \frac{d(1 + \lambda) \log(1/\delta)}{\epsilon} \right) \quad \text{for } \lambda := \frac{\log(1/\delta)}{\sqrt{d}},$$

*and if the input data is drawn according to [Definition 3](#), then whp it achieves excess risk  $\alpha$ .*

---

1. A distinct line of work, including *sufficient statistics perturbation* (SSP) and variants like AdaSSP [Sheffet \(2017\)](#); [Wang \(2018\)](#); [Lev et al. \(2026\)](#) focuses on DP linear regression under stronger assumptions on the covariates – typically  $\ell_2$  boundedness and small condition number. The state-of-the-art of these methods also has a sample complexity of  $\frac{d^{3/2}}{\alpha \epsilon} \geq \min \left\{ \frac{d}{\epsilon^2}, \frac{d^2}{\alpha^2} \right\}$ , no better than the ISSP or SoS approaches [Anderson et al. \(2025\)](#); [Brown et al. \(2024\)](#).

The computational complexity of STAR is dominated by  $\tilde{O}(n/d)$  linear algebra operations over  $(1 + \lambda)d \times d$  matrices. As such, it can be implemented with fast matrix multiplication achieving an asymptotic runtime of  $\tilde{O}(nd^{\omega-1}(1 + \lambda))$  – matching the runtime of OLS up to logarithmic factors.

From a statistical standpoint, even without privacy, estimating a  $d$ -dimensional regression parameter to error  $\alpha$  requires on the order of  $d/\alpha^2$  samples, but privacy comes at a cost and information-theoretic lower bounds imply that any  $(\epsilon, \delta)$ -DP estimator achieving excess risk  $\alpha$  with constant success probability must use

$$n \geq \tilde{\Omega}\left(\frac{d}{\alpha^2} + \frac{d}{\alpha\epsilon} + \frac{\log(1/\delta)}{\epsilon}\right)$$

samples (see, e.g., [Karwa and Vadhan \(2017\)](#); [Cai et al. \(2023\)](#); [Brown et al. \(2024\)](#)).

These bounds are essentially tight for inefficient algorithms. Liu et al.’s high-dimensional propose-test-release (HPTR) framework gives an exponential-time estimator with sample complexity which nearly matches the lower bound up to logarithmic factors when  $\alpha$  is a constant [Liu et al. \(2022\)](#). This highlights that the bottleneck for private regression algorithms is computational rather than statistical.

The most notable gap between  $n_{\text{STAR}}$  and the optimal sample complexity is the leading-order dependence on  $\log(1/\delta)$ , which is  $\log(1/\delta)/\epsilon$  in the optimal sample complexity, but incurs an extra  $d$  factor in  $n_{\text{STAR}}$ . This is partially explained by statistical query lower bounds, which suggest that polynomial-time DP regression algorithms using  $o(d^2)$  samples must also use  $\omega(\log(1/\delta)/\epsilon)$  samples [Georgiev and Hopkins \(2022\)](#); [Anderson et al. \(2025\)](#); thus, we do not expect to eliminate this  $d$  factor gap entirely. The other gaps between  $n_{\text{STAR}}$  and the optimal sample complexity are  $o(d)$  additive terms, or logarithmic factors.

In [Table 1](#), we summarize the state of DP linear regression algorithms. All existing polynomial-time algorithms fall short of the optimal sample complexity in (at least) one of three ways: requiring  $\Omega\left(\frac{d^2}{\alpha^2}\right)$  samples [Anderson et al. \(2025\)](#), requiring  $\Omega\left(\frac{d}{\epsilon^2}\right)$  samples [Brown et al. \(2024\)](#), or requiring samples growing polynomially with the condition number of  $\Sigma$ .

## 1.2. Techniques

### 1.2.1. SETTING AND CHALLENGES

**The Gaussian Mechanism** To get started, it will be useful to review the Gaussian mechanism, one of the foundational algorithmic tools in DP [Dwork and Roth \(2014\)](#). Suppose  $f : \mathcal{X}^n \rightarrow \mathbb{R}^d$  is a function from  $n$ -element datasets to  $\mathbb{R}^d$  which satisfies the following *bounded sensitivity* property: for all datasets  $X, X' \in \mathcal{X}^n$  which differ on exactly one element,  $\|f(X) - f(X')\| \leq 1$ , where  $\|\cdot\|$  is Euclidean/ $\ell_2$  norm. Then the following algorithm is  $(\epsilon, \delta)$ -DP: on input  $X$ , output  $f(X) + Z$ , where  $Z \sim \mathcal{N}(0, O(\frac{\log(1/\delta)}{\epsilon^2}) \cdot I)$  is Gaussian.

However, many classical estimators, including the empirical mean and OLS, have unbounded sensitivity to the addition of a single influential sample, even when all other samples are drawn from a benign distribution. At first sight, this seems incompatible with DP, which requires privacy under arbitrary single-sample edits. Our subsample-test-aggregate framework resolves this tension by requiring only the weaker notion of “down-sensitivity”.

**Down-Sensitivity** We formalize this “inlier-like” condition using *bounded down-sensitivity*: we say that a function  $f$  has bounded down-sensitivity on the dataset  $\mathcal{D}$  if *removing* any single sample

Paper	Sample Complexity	Method	Poly-time?
Liu et al. (2022)	$\frac{d}{\alpha^2} + \frac{d+\log(1/\delta)}{\alpha\varepsilon}$	HPTR	No
Varshney et al. (2022)	$\frac{d}{\alpha^2} + \frac{\kappa\sqrt{L}Rd\log(1/\delta)}{\alpha\varepsilon}$	DP-SGD	Yes
Liu et al. (2023)	$\frac{d}{\alpha^2} + \frac{\sqrt{\kappa}d\log(1/\delta)}{\alpha\varepsilon}$	DP-SGD	Yes
Asi et al. (2023)	$\frac{d}{\alpha^2} + \frac{d\log(R+\sqrt{\kappa})}{\alpha\varepsilon}$	exp. mech.	No
Brown et al. (2024)	$\frac{d}{\alpha^2} + \frac{d\sqrt{\log(1/\delta)}}{\alpha\varepsilon} + \frac{d\log^2(1/\delta)}{\varepsilon^2}$	outlier filtering	Yes
Anderson et al. (2025)	$\frac{d^2}{\alpha^2} + \frac{d}{\alpha\varepsilon} + \frac{d\log RL}{\varepsilon}$	SoS	Yes
Anderson et al. (2025)	$\frac{d^2}{\alpha^2} + \frac{d}{\alpha\varepsilon} + \frac{\log(1/\delta)}{\varepsilon}$	SoS	Yes
STAR (Theorem 4)	$\frac{d}{\alpha^2} + \frac{d\log(1/\delta)}{\alpha\varepsilon} + \frac{d\log(1/\delta)}{\varepsilon} + o(d)$	STA	Yes

Table 1: Sample complexities of private linear regression algorithms for well-behaved data (Definition 3). As in Anderson et al. (2025), we set  $d$  to be the dimension,  $\varepsilon, \delta$  to be the privacy parameters,  $\alpha$  to be the excess risk,  $R \geq \|\beta^*\|_2$  to be an upper bound on the norm of the ground truth model,  $\kappa$  an upper bound on the condition number of  $\Sigma$  and  $L$  an upper bound on the operator norm of  $\Sigma$ .

from  $\mathcal{D}$  has only a bounded effect on the output of  $f$ . Mean and linear regression have unbounded sensitivity to the *addition* of outliers regardless of the initial dataset, but despite this, there are *efficiently verifiable* conditions that suffice for bounded down-sensitivity for OLS, and these conditions hold with high probability on data drawn from “nice” distributions (such as Definition 3). Verifiable conditions for bounded down-sensitivity have been at the heart of many modern algorithms for private and robust statistics in high dimensions.

For instance, it is a well-known fact in the statistics community that the OLS estimator is robust to removing any single sample if it has bounded leverages and residuals Allen (1974), and this is the core of the ISSP regression algorithm of Brown et al. (2024): starting from the entire regression, samples are carefully removed until this down-sensitivity condition is met. Similarly, the Sum of Squares (SoS) based regression algorithm of Anderson et al. (2025) utilizes a down-sensitivity condition introduced by Bakshi and Prasad (2021) that is particularly amenable to the SoS framework.

However, the route from bounded down-sensitivity to DP algorithms typically has many technical challenges, often overcome with problem-specific algorithmic strategies. For instance, the main challenge in the ISSP algorithm is deciding *which* samples to classify as “outliers”, leading to the  $\frac{d\log^2(1/\delta)}{\varepsilon^2}$  term in their sample complexity, and the SoS approach of Anderson et al. (2025) comes at a significant cost to the runtime of their algorithm and is limited to down-sensitivity bounds that fit the SoS framework resulting in the  $\frac{d^2}{\alpha^2}$  term of their sample complexity.

More generally, if we are given a new down-sensitivity bound, it might be possible to find a corresponding filtering method à la ISSP, and it might be possible to convert it to a (slow) polynomial-

time SoS algorithm, but we are not guaranteed to get an efficient private algorithm from either framework.

Our STA technique charts a new course from down-sensitivity to DP estimation. Unlike previous approaches, STA is agnostic to the algorithm used to check down-sensitivity – instead, STA makes *black box* use of an oracle which can certify down-sensitivity of a learning algorithm (in our case, OLS regression) on subsamples of an input dataset  $X$ . We turn next to a more detailed overview of STA.

**Subsample-and-Aggregate** Another important mechanism to consider in our comparison is *subsample-and-aggregate* [Dwork and Roth \(2014\)](#). In this framework, subsamples (i.e., small subsets of the dataset) are drawn at random, the (potentially sensitive) learner is applied to each of the subsets separately, and the outputs of the learners are aggregated using a robust aggregation algorithm. Stability is obtained from the fact that changing any single entry in the dataset can only affect a small fraction of the subsamples allowing us to rely on the robustness guarantee of the aggregator. Finally, a standard mechanism (e.g., the Gaussian mechanism) is applied to this now-stable statistic to ensure privacy.

As a concrete example, a natural way to apply the subsample-and-aggregate framework to the problem of private linear regression, could be to draw a large number  $N$  of random subsamples from our original regression, to apply the standard OLS formula to each of them, and then aggregate using the BHS private mean estimation algorithm of [Brown et al. \(2023\)](#). In order to apply the OLS formula on each of these subsamples, we would most likely need each subsample to contain at least  $t \geq d/n$  fraction of the samples (otherwise the design matrix would be singular).

Therefore, an adversary changing a single sample in the regression might corrupt up to  $\approx tN$  of the subsamples. Therefore, to ensure privacy of the overall subsample-and-aggregate algorithm to a single change in the dataset, we need to ensure privacy of the BHS aggregator to  $tN$  changes in its input. Composing the privacy guarantees of [Brown et al. \(2023\)](#) to the change in a single learner  $tN$  times yields a privacy-accuracy tradeoff of

$$N \geq \frac{d \log(1/\delta)}{\epsilon} \cdot tN \Rightarrow n \geq \frac{d}{t} \geq \frac{d^2 \log(1/\delta)}{\epsilon},$$

yielding a quadratic-in- $d$  sample complexity.<sup>2</sup>

**Subsample-Test-Aggregate** Motivated by these shortcomings of existing privacy mechanisms, we introduce *Subsample-Test-Aggregate* (STA) – a privacy framework combining elements of subsample-and-aggregate and propose-test-release. At a high level, the STA algorithm generates random subsamples of the dataset, uses a sensitivity test to reject sensitive subsamples of the dataset and then aggregates the results. STA differs from subsample-and-aggregate in that its stability now stems solely from the test (allowing us to use simple aggregators such as outputting the empirical mean) and from propose-test-release in that it suffices to test only for down-sensitivity.

Concretely, STA has three main building blocks: a *learner*  $A : \mathcal{X}^* \rightarrow \Omega$ , which takes a dataset  $X$  in  $\mathcal{X}^*$  and outputs a hypothesis in a hypothesis space  $\Omega$ , a *sensitivity test*  $\mathcal{T} : \mathcal{X}^* \rightarrow \{\text{ACCEPT}, \text{REJECT}\}$ , which tests whether  $A$  has bounded down-sensitivity on a given dataset, and an *aggregator*, taking several hypotheses and aggregating them into a single one. Later, we

2. This analysis is just one way to use subsample-and-aggregate for DP-OLS, and we do not claim to prove an impossibility result. However, the issues we raise seem to be structural and to the best of our knowledge, no known subsample-and-aggregate algorithm for high dimensional statistics like OLS or mean estimation avoids them.

will generalize the setting, to allow the tester  $\mathcal{T}(X)$  to output a quantitative bound on the degree of down-sensitivity that  $A$  has on dataset  $X$ ; we start with the simpler version for the sake of exposition.

**Algorithm 1: SUBSAMPLE-TEST-AGGREGATE**

**Input:** Dataset  $X_{1:n} \in \mathcal{X}^n$ ; learner  $A : \mathcal{X}^* \rightarrow \Omega_A$ ; stability test  $\mathcal{T} : \mathcal{X}^* \rightarrow \{\text{ACCEPT}, \text{REJECT}\}$ .  
**Output:** A model in  $\Omega$  or REJECT.  
**for**  $j \leftarrow 1$  **to**  $N$  **do** // draw  $N$  iid Poisson subsamples  
    | Draw  $S_j \subseteq [n]$  by including each  $i \in [n]$  independently with probability  $t$ .  
    |  $T_j \leftarrow \mathcal{T}(X_{S_j})$   
**end**  
**if**  $|\{j \in [N] \mid T_j = \text{ACCEPT}\}|/N$  *is not sufficiently high* **then**  
    | **return** REJECT // private acceptance-rate check; see Thm. 17  
**end**  
 $\mathcal{J} \leftarrow \{j \in [N] : T_j = \text{ACCEPT}\}$  **return** Aggregate( $\{A(X_{S_j}) : j \in \mathcal{J}\}$ )

The overall goal is to map an input dataset in  $\mathcal{X}^n$  to a single good hypothesis for that dataset, in  $\Omega$ , subject to DP. Crucially, the learner  $A$  may have large or unbounded down-sensitivity on the input dataset  $X$  – allowing such datasets is important to obtain a DP algorithm, since in our settings every dataset is adjacent to one where  $A$  has large down-sensitivity. The STA scheme leverages the existence of *subsets* of  $X$  on which  $A$  is down-stable.

The STA scheme does this as follows. Given an  $n$ -element dataset  $X$ , pick  $X_{S_1}, \dots, X_{S_N}$  as random subsamples of  $X$ . For each of the subsamples  $X_{S_i}$ , run the tester  $\mathcal{T}(X_{S_i})$  to test whether  $A$  has bounded down-sensitivity on  $X_{S_i}$ . (We give the formal meaning of this below.) If sufficiently many of the tests pass, give the passing outputs of the learner  $\{A(X_{S_i}) : \mathcal{T}(X_{S_i}) = \text{ACCEPT}\}$  to the aggregator, and output its aggregated hypothesis. Unlike subsample-and-aggregate, the aggregator does not need to satisfy any robustness or stability properties – we show that testing down-sensitivity of the subsamples is instead enough to obtain strong privacy guarantees.

1.2.2. STAR-LITE: A CONCRETE EXAMPLE OF STA

To gain better intuition for the STA framework, we begin with a simplified version of the STAR algorithm – STAR-LITE (Algorithm 2).

STAR-LITE utilizes only the down-sensitivity implied by bounded leverages and residuals [Allen \(1974\)](#), and has similar privacy and runtime guarantees to those of STAR, with only a slightly worse sample complexity than the full STAR algorithm (though it still obtains a strictly better sample complexity than ISSP).

$$n_{\text{STAR-LITE}} = \tilde{O} \left( \frac{d}{\alpha^2} + \frac{d \log(1/\delta)}{\alpha \varepsilon} + \frac{d \log^2(1/\delta)}{\varepsilon} \right) < \tilde{O} \left( \frac{d}{\alpha^2} + \frac{d \log^2(1/\delta)}{\varepsilon^2} \right).$$

We will see that the test  $\mathcal{T}(X_{S_j})$  ensures that the matrix  $E_j$  bounds the magnitude and “shape” of the change to  $v_j$  which can be achieved by removing any single element of  $S_j$ . Concretely, if  $\Delta_j$  is this change, then we will show  $\Delta_j \Delta_j^\top \preceq E_j$ . Thus, the noise STAR-LITE adds to  $\frac{1}{k} \sum_{j=1}^k v_j$  adapts to the down-sensitivity of the subsampled OLS instances.

**Algorithm 2:** STAR-LITE: A simplified version of STAR.

**Input:** Data  $(X, y) \in \mathbb{R}^{n \times d} \times \mathbb{R}^n$

**Output:** Private linear model  $\hat{\beta} \in \mathbb{R}^d$

Set  $\mathcal{T}$  to be the *stability test* that accepts a subsample  $S \subseteq [n]$  if the corresponding regression has bounded leverages  $\lesssim L$  and residuals  $\lesssim R$ ; // where  $\lesssim$  denotes a randomized soft threshold

Define the Poisson-subsample process  $S \subseteq [n]$  to have  $i \in S \stackrel{i.i.d.}{\sim} \text{Bern}(t)$  with  $t \approx \frac{\log(1/\delta)d}{n}$

**if**  $\Pr[\mathcal{T}(S)] \ll 1$  // this can be checked privately - see Thm 17.

**then**

| **return** *REJECT*

**end**

Draw  $k \approx \frac{\log(n)}{\alpha^2 \log(1/\delta)}$  fresh Poisson-subsamples  $S_j$ , conditioned on  $\mathcal{T}$  accepting

Compute  $v_j := \text{OLS}(X_{S_j}, y_{S_j})$  and the down-sensitivity bounds  $E_j := R^2 L(X_{S_j}^\top X_{S_j})^{-1}$

**return** a sample from

$$\mathcal{N} \left( \frac{1}{k} \sum_{j=1}^k v_j, O \left( \frac{\log(1/\delta)}{k^2} \right) \sum_{j=1}^k E_j \right)$$

### 1.2.3. SUBSAMPLE-TEST-AGGREGATE: A NEW DP MECHANISM

**Down-Sensitivity of the Learning Algorithm** The most important relationship between the tester and learner is that if the tester outputs ACCEPT, it constitutes a promise that the learner is not down-sensitive. For purposes of this overview, we adopt the following definition, which assumes that  $\Omega = \mathbb{R}^d$ .

**Definition 5 (Down-Sensitivity Test (Simplified))** Let  $E \succ 0$ . We say that  $\mathcal{T}$  is an  $E \in \mathbb{R}^{d \times d}$  down-sensitivity test for  $A : \mathcal{X}^* \rightarrow \mathbb{R}^d$ , if

$$\forall S \subseteq [n] \forall i \in S \quad \Pr[\mathcal{T}(X_S) = \text{Accept}] > 0 \Rightarrow \left\| E^{-1/2} (A(X_S) - A(X_{S \setminus \{i\}})) \right\|_2 \leq 1.$$

The matrix  $E$  allows us to adapt down-sensitivity testing to a problem-appropriate norm – later we will return to the question of how to choose  $E$  for linear regression. Note that the definition of a down-sensitivity test is one-sided – it is possible for  $\mathcal{T}$  to REJECT on some  $X$  where  $A$  actually does have bounded down-sensitivity. This is important, since natural two-sided versions of this testing problem are computationally intractable [Moitra and Rohatgi \(2022\)](#). But it also means that when we use STA, we need to avoid trivial tests  $\mathcal{T}$  which always REJECT.

**Down-Stability of the Test** To ensure that STA is private, we will need a second down-stability property, this time *of the test itself* – the probability that the test accepts should not *decrease* by too much when a sample is removed. Notably, it could *grow* arbitrarily – this is crucial, since removing a single outlying sample can cause a dataset to go from highly down-sensitive to non-down-sensitive, meaning that the accept probability of the test should increase dramatically (perhaps from zero to nonzero).

**Definition 6 (Down-Stable Sensitivity Test (Simplified))** We say that the test  $\mathcal{T}$  is down-stable if for any dataset  $X$  and any index  $i$ , we have

$$\Pr[\mathcal{T}(X) = \text{Accept}] = O\left(\Pr[\mathcal{T}(X_{-i}) = \text{Accept}]\right)$$

where  $X_{-i}$  denotes  $X$  with the  $i$ -th sample removed. Here the probability is taken only with respect to internal randomness of the test  $\mathcal{T}$ .

Note that the existence of a down-stable sensitivity test  $\mathcal{T}$  for a learning algorithm  $A$  implies bounds on the down-sensitivity of  $A$  to removal of more than one sample – if  $\mathcal{T}(X)$  accepts with nonzero probability on  $X$ , then by down-stability it must accept with nonzero probability on all subsets  $X'$  of  $X$ , and hence  $A$  also has bounded down-sensitivity on all subsets  $X'$ . Later we will generalize the definition of down-stability of the test to allow for  $A$ 's down-sensitivity to increase as samples are removed – this is necessary for OLS, which becomes increasingly down-sensitive when the number of samples removed is so great that the design matrix becomes (near) singular.

**Stability and Privacy of STA** The following theorem captures the first key intuition about STA: the *average* of  $\mathcal{T}$ -accepted subsets is stable, not just to sample removals but also to arbitrary changes to a single sample.

**Theorem 7 (Stability of Mean Aggregation STA (Theorem 20 with  $\eta = 0$ ))** Let  $\mathcal{T}$  be a down-stable  $E \in \mathbb{R}^{d \times d}$  down-sensitivity test for  $A : \mathcal{X}^* \rightarrow \mathbb{R}^d$ . Let  $S \subseteq [n]$  be drawn by including each index independently with sufficiently small probability  $t$ . Let  $\mu(X) = \mathbb{E}[A(X_S) | \mathcal{T}(X_S) = \text{ACCEPT}]$  and  $M(X) = E + \text{Cov}(A(X_S) | \mathcal{T}(X_S) = \text{ACCEPT})$ . Then, for any neighboring datasets  $X, X'$  that differ on exactly one point, we have

$$\left\| M(X)^{-1/2} (\mu(X) - \mu(X')) \right\|_2 = O(t),$$

and

$$(1 - O(t))M(X) \preceq M(X') \preceq (1 + O(t))M(X).$$

The stability condition  $\|M(X)^{-1/2}(\mu(X) - \mu(X'))\| \leq O(t)$  suggests that if we add noise with covariance  $O(t^2) \cdot M(X)$  to  $\mu(X)$ , the result would be safe to release privately. However, this plan faces several complications. First,  $M(X)$  itself depends on  $X$ , which means that the *shape* of the noise we add might leak private information – although since  $M(X)$  itself is stable in spectral order we can hope to avoid too much privacy loss through this channel. Second, the privacy cost of adding Gaussian noise with covariance  $O(t^2) \cdot M(X)$  is difficult to analyze.<sup>3</sup> Third, STA doesn't have access directly to  $\mu(X)$ ; instead it can only draw samples from the distribution  $A(X_S)$  for a random subset  $S$ .

We address the challenges by using a different noise distribution with covariance proportional to  $M(X)$ . To sample  $\mu(X) + (\text{noise})$ , first, for an integer  $k$ , draw  $k$  samples from  $A(X_S)$  conditioned on  $\mathcal{T}(X_S) = \text{ACCEPT}$ . Then, average the samples and add  $Z \sim \mathcal{N}(0, (\log(1/\delta)/k)E)$ . The covariance of this distribution is similar to  $M(X)$ . Note that  $\mu(X)$  never needs to be computed directly. To prove privacy of the resulting mechanism, we combine the stability arguments underlying Theorem 7 with a generalization of the ‘‘advanced composition’’ method from DP. This argument is captured in the following theorem, which we can use to analyze the ‘‘Aggregate’’ step of STA when the aggregator is the empirical average plus Gaussian noise.

3. To the best of our knowledge, it may even violate privacy.

**Theorem 8 (Informal, see Theorem 30)** *The mechanism*

$$X \rightarrow \text{empirical mean over } k \text{ samples from } \text{Law} \left( \mathcal{N}(A(X_S), O(\log(1/\delta)E)) \mid S \sim \text{Acc}_X \right)$$

is  $(O(\sqrt{k \log(1/\delta)t}, \delta)$ -private.

The freedom to average over  $k$  draws from  $A(X_S)$  is crucial. In our application to OLS, each subset  $S$  will be too small to get a good estimator of  $\beta^*$  using only  $X_S$ . By averaging over many draws of  $A(X_S)$ , we get accuracy comparable to (non-private) OLS.

To complete the privacy analysis of Algorithm 1, we also need to analyze the privacy of the “acceptance rate check” step, which makes sure that  $\mathcal{T}(X_S)$  accepts with high-enough probability. Down-stability of  $\mathcal{T}$  turns out to be enough to privately check the acceptance rate – see Theorem 17.

**Generalizations Needed for OLS** We will need two key relaxations of the assumptions above to use STA for OLS regression. First, as we discussed above, we will need to allow for down-sensitivity of  $A$  to sometimes *increase* when a sample is removed. We can achieve this by modifying the definition of down-stability for  $\mathcal{T}$  to read:

$$\Pr[\mathcal{T}(X_S) = \text{Accept}] = O\left(\Pr[\mathcal{T}(X_{S \setminus \{i\}}) = \text{Accept}]\right) + \eta,$$

where we will eventually use  $\eta \approx \delta^{O(1)}$ .

The second relaxation we make is to Definition 5, where instead of assuming that  $\mathcal{T}(X) = \text{ACCEPT}$  implies a bound on the down-sensitivity of  $A$  in some fixed  $E$  geometry, we allow  $\mathcal{T}(X)$  to return a matrix  $E(X)$ , and an ACCEPT output guarantees bounded down-sensitivity for  $A$  in the  $E(X)$ -norm. We prove an analogue of Theorem 8 with varying  $E(X)$ , using a guarantee that  $E(X)$  itself is monotonically non-increasing and changes smoothly

$$E(X) \preceq E(X_{-i}) \quad \text{and} \quad \Omega(|E(X_{-i})|) \leq |E(X)| \leq |E(X_{-i})|.$$

#### 1.2.4. STAR: SUBSAMPLE-TEST-AGGREGATE REGRESSION

We now turn to our instantiation of STA for linear regression. Our learner will be OLS regression itself, and our aggregator will be the mean-plus-Gaussian described in Theorem 8 (generalized to allow data-dependent  $E(X_S)$  as described above). All that remains is to define our down-sensitivity test  $\mathcal{T}$  and its corresponding leave-one-out bound  $E(X_S)$ . Motivated by the Algorithm for Certifying Robustness Efficiently (ACRE) of Rubinfeld and Hopkins (2025), we propose the following test for robustness:

**Definition 9 ((Simplified) ACRE Robustness)** *Given a regression problem  $X = (X_1, \dots, X_n) \in \mathbb{R}^{n \times d}$ ,  $y = (Y_1, \dots, Y_n) \in \mathbb{R}^n$ , define its hat matrix to be  $H = X(X^\top X)^{-1}X^\top \in \mathbb{R}^{n \times n}$ . We say that the regression is  $(L, R, \ell)$ -ACRE if:*

1. All samples have bounded leverages – i.e., the diagonal elements of the hat matrix satisfy  $H_{i,i} \leq L$ .
2. All samples have bounded residuals. In terms of the hat matrix, this can be written as  $\|y - Hy\|_\infty \leq R$ .

3. All sample pairs have bounded cross-leverages – i.e., the off-diagonal elements of the hat matrix are also bounded  $|H_{i,j}| \leq \ell$  for all  $i \neq j$ .

Using leverage and residual scores to test the robustness of a linear regression is a common technique [Weisberg \(2005\)](#), but by also testing for cross-leverages, we can effectively test the robustness of the regression to a larger number of sample removals. This is a very simplified version of the tests performed by the ACRE algorithm. We capture in a later Lemma ([Lemma 12](#)) that the ACRE properties imply bounded down-sensitivity, and move on now to ACRE’s stability.

The key property of ACRE that we will use to ensure stability is that if a set of samples satisfies the ACRE conditions, any leave-one-out satisfies the ACRE condition with slightly worse parameters. This follows almost immediately from the Sherman-Morrison formula.

**Lemma 10 (Stability of the ACRE Condition)** *If  $X, y$  is  $(L, R, \ell)$ -ACRE for  $\ell \leq L$ , then any dataset obtained by dropping a single sample from  $X, y$  is  $((1 + c)L, (1 + c)R, (1 + c)\ell)$ -ACRE, for  $c = \frac{\ell}{1-L}$ .*

A careful reader might note two remaining discrepancies in our definition of the ACRE test. First, the current test is deterministic, and thus cannot achieve our definition of down-stability. Second, the test has several hyperparameters that need to be set.

**Soft-ACRE** We use [Lemma 10](#) to construct a randomized soft-threshold version of the ACRE test.

**Algorithm 3:** SOFT-ACRE

**Input:** Dataset  $X, y \in \mathbb{R}^{n \times d} \times \mathbb{R}^n$  and target thresholds  $(L, R, \ell) \in \mathbb{R}^3$

**Output:** ACCEPT or REJECT

Compute the hat matrix  $H \leftarrow X(X^\top X)^{-1}X^\top$  // if  $X^\top X$  is singular, we reject.

Find the worst violator of the ACRE conditions

$$\text{SCORE} = \max \left\{ \max_{i \in [n]} \left\{ \frac{H_{i,i}}{L} \right\}, \max_{i \neq j \in [n]} \left\{ \frac{|H_{i,j}|}{\ell} \right\}, \frac{\|y - Hy\|_\infty}{R} \right\}$$

Draw  $U \sim \text{Unif}[0, 1]$

**if** SCORE  $\leq 1$  **or** (SCORE  $\leq 2$  **and**  $U \leq \text{SCORE}^{-\frac{2L}{1-2\ell}}$ ) **then**

    | **return** ACCEPT

**else**

    | **return** REJECT

**end**

**Lemma 11 (SOFT-ACRE is Down-Stable)** *The randomized test  $\mathcal{T} = \text{SOFT-ACRE}$  is  $(2, \eta)$  down-stable (see [Definition 13](#)) for*

$$\eta = \exp \left( -\ln(2) \left( 1 - \frac{2L}{1-2\ell} \right) \frac{1-2L}{2\ell} \right).$$

Finally, we need to show that SOFT-ACRE also implies a bound on the leave-one-out effect of a regression.

**Lemma 12 (SOFT-ACRE implies bounded down-sensitivity)** *If  $X, y$  is a regression that is accepted by SOFT-ACRE with non-zero probability, then for  $E(X) := 8R^2 \frac{L}{(1-2L)^2} (X^\top X)^{-1}$  and any index  $i \in [n]$ , we have*

$$\left\| E(X)^{-1/2} (\text{OLS}(X, y) - \text{OLS}(X_{-i}, y_{-i})) \right\|_2 \leq 1.$$

Lemma 12 follows immediately from the fact that SOFT-ACRE never accepts regressions that are not  $(2L, 2R, 2\ell)$ -ACRE and the Sherman-Morrison formula.

**Selecting the Hyperparameters of ACRE** Given a subsample of size  $|S|$ , it is a well-known analysis that the average leverage scales like

$$\mathbb{E}_{i \in S} \left[ X_i^\top (X_S^\top X_S)^{-1} X_i \right] = \frac{1}{|S|} \text{tr}[H_S] = \frac{1}{|S|} \text{tr}[I_d] = \frac{d}{|S|}.$$

Therefore, for the STAR algorithm we will set the leverage threshold to  $L = \Theta\left(\frac{d}{tn}\right)$ . To get the STAR-LITE algorithm, we fix  $\ell = L$  (from the Cauchy-Schwarz inequality, this makes the cross-leverage test null, since we always have  $|H_{i,j}| \leq \sqrt{H_{i,i}H_{j,j}} \leq L$ ). However, we can often gain a “blessing-of-dimensionality” effect for the cross-leverages, which allows us to use  $\ell = \Theta\left(\frac{\sqrt{d \log(n)}}{tn}\right)$ .

Finally, for the residual bound  $R$ , if we are told the variance of the underlying sample distribution, we may set  $R = \Theta(\sigma \sqrt{\log(n)})$ . To select  $R$  when  $\sigma$  is not known, we first estimate  $\sigma$  using Algorithm 6 of Brown et al. (2024) (a subsample-and-aggregate algorithm for estimating  $\sigma$  up to a constant factor using  $\frac{d \log(1/\delta)}{\epsilon}$  samples).

### 1.3. Related Work

Related work on DP linear regression is in Section 1.1.

**Down-Sensitivity, Lipschitz Extensions, and Privacy** Down-sensitivity (defined slightly differently than we do here) as a tool for DP emerged in Chen and Zhou (2013), and has seen use especially in private analysis of graphs – see Raskhodnikova and Smith (2016) for a survey. These ideas are extended somewhat beyond graphs, though still in settings where the output is one-dimensional, by Fang et al. (2022). And they are extended to a graph-theoretic statistical inference setting, graphon estimation, in Borgs et al. (2015, 2018). More recent works relating down-sensitivity and privacy in the context of privately releasing real-valued (one dimensional) functions include Linder et al. (2025); Cummings and Durfee (2020); Steinke and Steinke (2025). And in the statistical estimation setting, Ashtiani and Liaw (2022); Tsfadia et al. (2022) also use approaches which combine sensitivity considerations with subsampling.

**Reductions from Private Algorithms to Robust Ones in High-Dimensional Statistics** A number of recent works offer frameworks to convert learning algorithms satisfying robustness guarantees to differentially private ones. Hopkins et al. (2023); Asi et al. (2023) describe a black-box reduction based on the *inverse sensitivity mechanism* Asi and Duchi (2020a); Asi et al. (2023), which can

obtain optimal privacy-accuracy tradeoffs, but the black-box reduction is not computationally efficient, even starting with a computationally efficient robust learning algorithm, and its computationally efficient variant in Hopkins et al. (2023) is non-black-box. These works were preceded by the foundational *high dimensional propose-test-release* framework of Liu et al. (2022), whose reduction is non-black-box and also not computationally efficient. Kothari et al. (2022) gives a framework for converting certain robust learning algorithms based on convex programming to private ones; their framework allows for computationally efficient private algorithms but is non-black-box. So far, none of these approaches offers a combination of computational efficiency and privacy-accuracy tradeoffs which achieves an algorithm with comparable guarantees to ours for linear regression.

**Beyond Worst-Case Analysis in DP** DP algorithms which add noise scaled to global sensitivity often do achieve optimal privacy-accuracy tradeoffs when accuracy is measured in a worst-case sense Dwork and Roth (2014), and yet the need for DP algorithms which adapt the amount of noise they add to the “shape” or “niceness” of each particular input has long been clear. How, then, do we measure the accuracy or quality of a DP algorithm in a way which differentiates algorithms which add less noise on “nice” inputs? The statistical estimation setting we study here is one way to formulate an average-case accuracy measure for a DP algorithm which distinguishes such algorithms. Other approaches include the *instance optimality* framework Asi and Duchi (2020a,b); Huang et al. (2021); McMillan et al. (2022); Feldman et al. (2024) and *subset-based instance optimality* Dick et al. (2023).

#### 1.4. Organization

In Appendix A, we introduce the subsample-test-aggregate in more detail and prove some of its basic results. Appendix B is devoted to our privacy mechanism based on adding Gaussian noise to an empirical average over subsamples with bounded down-sensitivity. In Appendix C, we give more details on our STAR algorithm and prove Theorem 4. Finally, in Appendix D we introduce basic DP primitives used throughout our construction.

#### Acknowledgments

We used large language models (LLMs) to assist with writing initial drafts of some of the proofs and for latex formatting. All LLM-generated text was reviewed and edited by human authors, and the authors take full responsibility for the content and originality of this submission.

We thank Gautam Kamath and Stefan Tiegel for helpful comments and feedback on earlier versions of this paper.

This work was supported by NSF Award No. 2238080 and CSAIL Alliances. Ittai Rubinstein was additionally supported by the MIT EECS MathWorks Fellowship.

#### References

Daniel Alabi, Audra McMillan, Jayshree Sarathy, Adam Smith, and Salil Vadhan. Differentially private simple linear regression. *arXiv preprint arXiv:2007.05157*, 2020.

David M Allen. The relationship between variable selection and data augmentation and a method for prediction. *technometrics*, 16(1):125–127, 1974.

- Prashanti Anderson, Ainesh Bakshi, Mahbod Majid, and Stefan Tiegel. Sample-optimal private regression in polynomial time. In *Proceedings of the 57th Annual ACM Symposium on Theory of Computing*, pages 2341–2349, 2025.
- Jamil Arbas, Hassan Ashtiani, and Christopher Liaw. Polynomial time and private learning of unbounded gaussian mixture models. In *International Conference on Machine Learning*, pages 1018–1040. PMLR, 2023.
- Hassan Ashtiani and Christopher Liaw. Private and polynomial time algorithms for learning gaussians and beyond. In *Conference on Learning Theory*, pages 1075–1076. PMLR, 2022.
- Hilal Asi and John C Duchi. Instance-optimality in differential privacy via approximate inverse sensitivity mechanisms. *Advances in neural information processing systems*, 33:14106–14117, 2020a.
- Hilal Asi and John C Duchi. Near instance-optimality in differential privacy. *arXiv preprint arXiv:2005.10630*, 2020b.
- Hilal Asi, Jonathan Ullman, and Lydia Zakyntinou. From robustness to privacy and back. In *International Conference on Machine Learning*, pages 1121–1146. PMLR, 2023.
- Ainesh Bakshi and Adarsh Prasad. Robust linear regression: Optimal rates in polynomial time. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, pages 102–115, 2021.
- Andrés F Barrientos, Aaron R Williams, Joshua Snoke, and Claire McKay Bowen. A feasibility study of differentially private summary statistics and regression analyses with evaluations on administrative and survey data. *Journal of the American Statistical Association*, 119(545):52–65, 2024.
- Sourav Biswas, Yihe Dong, Gautam Kamath, and Jonathan Ullman. Coinpress: Practical private mean and covariance estimation. *Advances in Neural Information Processing Systems*, 33:14475–14485, 2020.
- Christian Borgs, Jennifer Chayes, and Adam Smith. Private graphon estimation for sparse graphs. *Advances in Neural Information Processing Systems*, 28, 2015.
- Christian Borgs, Jennifer Chayes, Adam Smith, and Ilias Zadik. Revealing network structure, confidentially: Improved rates for node-private graphon estimation. In *2018 IEEE 59th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 533–543. IEEE, 2018.
- Gavin Brown, Marco Gaboardi, Adam Smith, Jonathan Ullman, and Lydia Zakyntinou. Covariance-aware private mean estimation without private covariance estimation. *Advances in Neural Information Processing Systems*, 34:7950–7964, 2021.
- Gavin Brown, Samuel B. Hopkins, and Adam D. Smith. Fast, sample-efficient, affine-invariant private mean and covariance estimation for subgaussian distributions. *CoRR*, abs/2301.12250, 2023. doi: 10.48550/ARXIV.2301.12250.

- Gavin Brown, Jonathan Hayase, Samuel Hopkins, Weihao Kong, Xiyang Liu, Sewoong Oh, Juan C Perdomo, and Adam Smith. Insufficient statistics perturbation: Stable estimators for private least squares. *arXiv preprint arXiv:2404.15409*, 2024.
- T Tony Cai, Yichen Wang, and Linjun Zhang. The cost of privacy: Optimal rates of convergence for parameter estimation with differential privacy. *The Annals of Statistics*, 49(5):2825–2850, 2021.
- T Tony Cai, Yichen Wang, and Linjun Zhang. Score attack: A lower bound technique for optimal differentially private learning. *arXiv preprint arXiv:2303.07152*, 2023.
- Shixi Chen and Shuigeng Zhou. Recursive mechanism: towards node differential privacy and unrestricted joins. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, pages 653–664, 2013.
- Rachel Cummings and David Durfee. Individual sensitivity preprocessing for data privacy. In *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 528–547. SIAM, 2020.
- Yuval Dagan, Michael Jordan, Xuelin Yang, Lydia Zakyntinou, and Nikita Zhivotovskiy. Dimension-free private mean estimation for anisotropic distributions. *Advances in Neural Information Processing Systems*, 37:120667–120698, 2024.
- Travis Dick, Alex Kulesza, Ziteng Sun, and Ananda Theertha Suresh. Subset-based instance optimality in private estimation. In *International Conference on Machine Learning*, pages 7992–8014. PMLR, 2023.
- Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. *Foundations and trends® in theoretical computer science*, 9(3-4):211–487, 2014.
- Cynthia Dwork, Krishnaram Kenthapadi, Frank McSherry, Ilya Mironov, and Moni Naor. Our data, ourselves: Privacy via distributed noise generation. In *Annual international conference on the theory and applications of cryptographic techniques*, pages 486–503. Springer, 2006a.
- Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006b.
- Juanru Fang, Wei Dong, and Ke Yi. Shifted inverse: A general mechanism for monotonic functions under user differential privacy. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, pages 1009–1022, 2022.
- Vitaly Feldman, Audra McMillan, Satchit Sivakumar, and Kunal Talwar. Instance-optimal private density estimation in the wasserstein distance. *Advances in Neural Information Processing Systems*, 37:90061–90131, 2024.
- Kristian Georgiev and Samuel B. Hopkins. Privacy induces robustness: Information-computation gaps and sparse mean estimation. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9*,

- 2022, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/2d76b6a9f96181ab717c1a15ab9302e1-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/2d76b6a9f96181ab717c1a15ab9302e1-Abstract-Conference.html).
- Samuel B Hopkins, Gautam Kamath, Mahbod Majid, and Shyam Narayanan. Robustness implies privacy in statistical estimation. In *Proceedings of the 55th Annual ACM Symposium on Theory of Computing*, pages 497–506, 2023.
- Yuzheng Hu, Pingbang Hu, Han Zhao, and Jiaqi Ma. Most influential subset selection: Challenges, promises, and beyond. *Advances in Neural Information Processing Systems*, 37:119778–119810, 2024.
- Jenny Y Huang, David R Burt, Tin D Nguyen, Yunyi Shen, and Tamara Broderick. Approximations to worst-case data dropping: unmasking failure modes. *arXiv preprint arXiv:2408.09008*, 2024.
- Ziyue Huang, Yuting Liang, and Ke Yi. Instance-optimal mean estimation under differential privacy. *Advances in Neural Information Processing Systems*, 34:25993–26004, 2021.
- Gautam Kamath, Jerry Li, Vikrant Singhal, and Jonathan Ullman. Privately learning high-dimensional distributions. In *Conference on Learning Theory*, pages 1853–1902. PMLR, 2019.
- Gautam Kamath, Vikrant Singhal, and Jonathan Ullman. Private mean estimation of heavy-tailed distributions. In *Conference on Learning Theory*, pages 2204–2235. PMLR, 2020.
- Vishesh Karwa and Salil Vadhan. Finite sample differentially private confidence intervals. *arXiv preprint arXiv:1711.03908*, 2017.
- Pravesh Kothari, Pasin Manurangsi, and Ameya Velingker. Private robust estimation by stabilizing convex relaxations. In *Conference on Learning Theory*, pages 723–777. PMLR, 2022.
- Rohith Kuditipudi, John Duchi, and Saminul Haque. A pretty fast algorithm for adaptive private mean estimation. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 2511–2551. PMLR, 2023.
- Omri Lev, Moshe Shenfeld, Vishwak Srinivasan, Katrina Ligett, and Ashia C Wilson. Near-optimal private linear regression via iterative hessian mixing. *arXiv preprint arXiv:2601.07545*, 2026.
- Ephraim Linder, Sofya Raskhodnikova, Adam Smith, and Thomas Steinke. Privately evaluating untrusted black-box functions. In *Proceedings of the 57th Annual ACM Symposium on Theory of Computing*, pages 2350–2361, 2025.
- Xiyang Liu, Weihao Kong, and Sewoong Oh. Differential privacy and robust statistics in high dimensions. In *Conference on Learning Theory*, pages 1167–1246. PMLR, 2022.
- Xiyang Liu, Prateek Jain, Weihao Kong, Sewoong Oh, and Arun Sai Suggala. Near optimal private and robust linear regression. *arXiv preprint arXiv:2301.13273*, 2023.
- Audra McMillan, Adam Smith, and Jon Ullman. Instance-optimal differentially private estimation. *arXiv preprint arXiv:2210.15819*, 2022.
- Ankur Moitra and Dhruv Rohatgi. Provably auditing ordinary least squares in low dimensions. *arXiv preprint arXiv:2205.14284*, 2022.

- Shyam Narayanan, Vahab Mirrokni, and Hossein Esfandiari. Tight and robust private mean estimation with few users. In *International Conference on Machine Learning*, pages 16383–16412. PMLR, 2022.
- Sofya Raskhodnikova and Adam Smith. Differentially private analysis of graphs. In *Encyclopedia of Algorithms*, pages 543–547. Springer, 2016.
- Ittai Rubinstein and Samuel B. Hopkins. Robustness auditing for linear regression: To singularity and beyond. In *Proceedings of the Thirteenth International Conference on Learning Representations (ICLR 2025)*, 2025. URL <https://openreview.net/forum?id=V5ns6uvRZ9>.
- Or Sheffet. Differentially private ordinary least squares. In *International Conference on Machine Learning*, pages 3105–3114. PMLR, 2017.
- Günter F Steinke and Thomas Steinke. Privately estimating black-box statistics. *arXiv preprint arXiv:2510.00322*, 2025.
- Eliad Tsfadia, Edith Cohen, Haim Kaplan, Yishay Mansour, and Uri Stemmer. Friendlycore: Practical differentially private aggregation. In *International Conference on Machine Learning*, pages 21828–21863. PMLR, 2022.
- Prateek Varshney, Abhradeep Thakurta, and Prateek Jain. (nearly) optimal private linear regression via adaptive clipping. *arXiv preprint arXiv:2207.04686*, 2022.
- Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- Yu-Xiang Wang. Revisiting differentially private linear regression: optimal and adaptive prediction & estimation in unbounded domain. *arXiv preprint arXiv:1803.02596*, 2018.
- Sanford Weisberg. *Applied Linear Regression*. Wiley Series in Probability and Statistics. John Wiley & Sons, Hoboken, NJ, 3 edition, 2005. ISBN 9780471663799.

## Appendix A. The Subsample-Test-Aggregate Privacy Mechanism

In this section we formalize the subsample-test-aggregate (STA) mechanism abstractly. We focus on the minimal properties we need from a Test procedure in order to obtain strong privacy guarantees by aggregating over many test-passing subsamples.

### A.1. Notations

We fix notation for the *subsample-test* stage of the STA mechanism. Throughout, we will be careful about the two sources of randomness: the subsampling step and the randomized Test.

**Subsampling randomness.** Fix a dataset  $X = (x_1, \dots, x_n) \in \mathcal{X}^n$  and a subsampling rate  $t \in (0, 1)$ . As previously described, STA first draws a **random** subset  $S \subseteq [n]$  by *Poisson subsampling* at rate  $t$ , which we write as  $S \sim \text{PoisSub}_t$ . We write  $X_S \stackrel{\text{def}}{=} (x_i)_{i \in S}$ . When we need to represent a fixed subset, we will frequently use  $T \subseteq [n]$ .

**Test randomness.** The Test procedure  $\mathcal{T}$  is randomized. Given a subsample  $S$ , we run the test on  $X_S$  and write its accept/reject bit as

$$b_S(X) \in \{\text{ACCEPT}, \text{REJECT}\}.$$

(When the test produces additional outputs—e.g., a certificate matrix—we will denote them explicitly later; for this notation section we only need the accept/reject bit.)

**Acceptance event and pass probability.** We define the acceptance event and its marginal probability by

$$\text{Acc}_X \stackrel{\text{def}}{=} \{b_S(X) = \text{ACCEPT}\}, \quad p_X \stackrel{\text{def}}{=} \Pr[\text{Acc}_X].$$

Unless otherwise specified, all probabilities and expectations involving  $\text{Acc}_X$  (e.g.,  $\Pr[\text{Acc}_X]$ ,  $\Pr[\cdot | \text{Acc}_X]$ ,  $\mathbb{E}[\cdot | \text{Acc}_X]$ ) are taken over the *joint* randomness of: (i)  $S \sim \text{PoisSub}_t$  and (ii) the internal randomness of  $\mathcal{T}$ .

When we condition on  $S$  (e.g., in  $\Pr[\text{Acc}_X | S]$ ), the probability is only over the internal randomness of  $\mathcal{T}$ .

**The post-test subsample distribution.** When  $p_X > 0$ , the *post-test subsample distribution* is simply the conditional distribution of  $S$  given acceptance,  $S | \text{Acc}_X$ .

We will commonly switch between conditional and unconditional probabilities, as well as conditional and unconditional expectations, using the following identities from elementary probability.

$$\Pr[S = T | \text{Acc}_X] = \frac{\Pr[S = T \wedge \text{Acc}_X]}{p_X}.$$

and for any function  $f : 2^{[n]} \rightarrow \mathbb{R}$ ,

$$\mathbb{E}[f(S) | \text{Acc}_X] = \frac{\mathbb{E}[f(S) \cdot \mathbf{1}\{\text{Acc}_X\}]}{p_X}.$$

**General Notation** We write  $\mathbb{S}_+^d$  for the cone of  $d \times d$  symmetric positive semidefinite matrices,

$$\mathbb{S}_+^d \stackrel{\text{def}}{=} \left\{ M \in \mathbb{R}^{d \times d} : M = M^\top \text{ and } M \succeq 0 \right\}.$$

## A.2. Down-Stable Test

### A.2.1. DEFINITION

**Definition 13 (Down-Stable Test)** We say that the test  $\mathcal{T} : \mathcal{X}^* \rightarrow \{\text{ACCEPT}, \text{REJECT}\}$  is  $C_p, \eta$ -down-stable, if for any set of samples  $X \in \mathcal{X}^*$ , and any leave-one-out  $X^{-i}$  of  $X$ , we have

$$\Pr[\mathcal{T}(X) = \text{ACCEPT}] \leq C_p \Pr[\mathcal{T}(X^{-i}) = \text{ACCEPT}] + \eta.$$

Even before tying  $\mathcal{T}$  to a learning algorithm, we can derive some properties just from the fact that it is down-stable.

A.2.2. KEY PROPERTY: NO SINGLE SAMPLE CAN BE TOO OVER-REPRESENTED BY  
CONDITIONING ON A DOWN-STABLE TEST

**Theorem 14 (No index is over-represented)**

If  $\mathcal{T}$  is a  $(C_p, \eta)$ -down-stable test,  $i \in [n]$  is any index and  $T \subseteq [n]$  is any subset of the training set containing  $i \in T$ , then

$$\Pr[S = T \mid \text{Acc}_X] \leq \frac{C_p t}{1-t} \cdot \left( \Pr[S = T \setminus \{i\} \mid \text{Acc}_X] + \frac{\eta}{C_p p_X} \cdot \Pr[S = T \setminus \{i\}] \right),$$

where all probabilities are taken over the subsampling (and, for the event  $\text{Acc}_X$ , also over the randomness of  $\mathcal{T}$ ).

**Remark 15 (Interpretation of Theorem 14)** *Theorem 14 says that conditioning on acceptance cannot dramatically increase the relative likelihood of subsamples that contain any fixed point  $i$ . Up to the expected inclusion rate (the factor  $\frac{t}{1-t}$ ) and a constant factor  $C_p$ , sets containing  $i$  behave similarly to the corresponding sets without  $i$ , plus an additive slack controlled by  $\eta/p$ . In applications we will choose parameters so that  $\eta \ll p$ , making the slack term negligible.*

Theorem 14 is central to the STA framework. First, we can use it to certify that the acceptance rate of  $\mathcal{T}$  on random subsamples is stable (Corollary 16), which in turn allows us to privately check that this acceptance probability is not too low (Theorem 17). Privately checking that the acceptance rate  $p_X$  is not too low is crucial, since otherwise our STA algorithm could have an unbounded runtime or condition on very low probability events.

This no over-representation property is also what powers our privacy guarantee for the single-draw aggregation mechanism (Lemma 28), and we re-use the same proof technique to show that mean-aggregation STA is stable (Theorem 20).

**Proof of Theorem 14** **Proof** [Proof of Theorem 14] Fix  $i \in [n]$  and  $T \ni i$ , and let  $R \stackrel{\text{def}}{=} T \setminus \{i\}$ . Write

$$q_T \stackrel{\text{def}}{=} \Pr [\mathcal{T}(X_T) = \text{ACCEPT}] \quad \text{and} \quad q_R \stackrel{\text{def}}{=} \Pr [\mathcal{T}(X_R) = \text{ACCEPT}],$$

where the probability is over the internal randomness of  $\mathcal{T}$ .

By our down-stability assumption (Definition 13) applied to  $(X_T, X_R)$  which differ by one deletion, we have

$$q_T \leq C_p q_R + \eta.$$

Multiplying by  $\Pr[T]$  and using  $\Pr[T \wedge \text{Acc}_X] = \Pr[T] \cdot q_T$ , we get

$$\Pr[T \wedge \text{Acc}_X] \leq C_p \Pr[T] \cdot q_R + \eta \Pr[T].$$

Under Poisson subsampling with rate  $t$ , inclusion of  $i$  is independent of the other indices, so for  $T \ni i$  we have

$$\Pr[T] = \frac{t}{1-t} \Pr[R].$$

Substituting this identity and using  $\Pr[R \wedge \text{Acc}_X] = \Pr[R] \cdot q_R$ , we obtain

$$\Pr[T \wedge \text{Acc}_X] \leq \frac{tC_p}{1-t} \Pr[R \wedge \text{Acc}_X] + \frac{t\eta}{1-t} \Pr[R].$$

Finally, dividing by  $p_X = \Pr[\text{Acc}_X]$  gives

$$\Pr[T \mid \text{Acc}_X] = \frac{\Pr[T \wedge \text{Acc}_X]}{p_X} \leq \frac{tC_p}{1-t} \Pr[R \mid \text{Acc}_X] + \frac{t\eta}{(1-t)p_X} \Pr[R],$$

which is the desired inequality. ■

**Corollary: The acceptance rate of  $\mathcal{T}$  is stable**

**Corollary 16 (Stability of STA Acceptance Probability)** *Assume  $\mathcal{T}$  is a  $C_p, \eta$ -down-stable test, run under independent Poisson subsampling at rate*

$$t \in \left(0, \frac{1}{5(1+C_p)}\right),$$

and let  $X, X' \in \mathcal{X}^n$  be neighboring datasets that differ in exactly one coordinate.

Then, if

$$p_X \geq 10\eta,$$

we have

$$(1 - 3C_p t) p_{X'} \leq p_X \leq (1 + 3C_p t) p_{X'}. \quad (1)$$

### A.2.3. PRIVATELY TESTING $p_X$

The goal of this section is to develop a differentially private check that the probability of a subsample passing the test,  $p_X$ , is not too small. This is necessary because many of our stability results like Corollary 16 require this passing probability to be much larger than the  $\eta$  error terms. Furthermore, having too small of a  $p_X$  would increase the number of subsamples we have to draw in the Monte Carlo estimates, increasing runtime.

We provide an algorithm for privately testing that  $p_X$  is sufficiently large, and prove a theorem about its guarantees. This algorithm calls `PrivateBernoulliThreshold` (Algorithm 7) as a subroutine.

**Algorithm 4:** `PRIVATEPCHECK $_{\varepsilon, \delta, \gamma, p_{\text{acc}}}$ ( $X$ )`

**Input:** Dataset  $X$ ; subsampling rate  $t$  (used by  $\mathcal{T}$ ); test  $\mathcal{T}$  with stability parameter  $C_p$ ; privacy parameters  $(\varepsilon, \delta) \in (0, 1)^2$ ; slack  $\gamma \in (0, 1/2)$ ; threshold  $p_{\text{acc}} \in (0, 1)$ ; number of runs  $N \in \mathbb{N}$ .

**Output:** `ACCEPT` or `REJECT`.

$\tau \leftarrow -\log(1 - 3C_p t)$ ; // well-defined since  $3C_p t < 1$  under Cor. 16

$\varepsilon_0 \leftarrow \varepsilon/3$   $\delta_0 \leftarrow \delta/12$   $p_{\text{th}} \leftarrow p_{\text{acc}} e^{-\gamma}$

**for**  $j = 1$  **to**  $N$  **do**

    | Draw  $S_j$  by independent Poisson subsampling at rate  $t$  Run  $\mathcal{T}$  on  $X_{S_j}$  and set  $b_j \leftarrow$   
     |  $\mathbb{1}\{\mathcal{T}(X_{S_j}) = \text{ACCEPT}\}$

**end**

$y \leftarrow \text{PBT}_{\tau, \varepsilon_0, \delta_0}(b_1, \dots, b_N; p_{\text{th}})$ ; // Algorithm 7

**return** `ACCEPT` if  $y = \text{REJECT}$ , else `REJECT`

**Theorem 17 (PRIVATEPCHECK: privacy, soundness, and completeness)** *Assume  $\mathcal{T}$  is a  $C_p, \eta$ -down-stable test run under independent Poisson subsampling at rate  $t$  satisfying the conditions of Corollary 16. For each dataset  $X$ , define the acceptance probability*

$$p_X := \Pr[\mathcal{T}(X_S) = \text{ACCEPT}],$$

where the probability is over the Poisson subsample  $S$  and the internal randomness of  $\mathcal{T}$ .

Fix  $\varepsilon \in (0, 1)$ ,  $\delta \in (0, 1/10)$ , a slack parameter  $\gamma \in (0, 1/2)$ , and an acceptance threshold  $p_{\text{acc}} \in (0, 1)$ . Define

$$\tau := -\log(1 - 3C_p t), \quad \text{and assume} \quad \gamma \leq \tau/2.$$

Define the low/high acceptance thresholds

$$p_{\text{high}} := p_{\text{acc}}, \quad p_{\text{low}} := p_{\text{acc}} \cdot \exp\left(-\frac{6\tau}{\varepsilon} \log \frac{24}{\delta}\right) \cdot e^{-2\gamma}.$$

Assume also  $\eta \leq p_{\text{low}}/20$  and

$$N \geq \frac{12}{\gamma^2 p_{\text{low}}} \cdot \log \frac{48}{\delta}. \quad (2)$$

Run Algorithm 4 with parameters  $(\varepsilon, \delta, \gamma, p_{\text{acc}}, N)$ . Then it satisfies:

- **Completeness.** If  $p_X \geq p_{\text{high}}$  then  $\Pr[\text{PRIVATEPCHECK}(X) = \text{ACCEPT}] \geq 1 - \delta$ .
- **Soundness.** If  $p_X \leq p_{\text{low}}$  then  $\Pr[\text{PRIVATEPCHECK}(X) = \text{REJECT}] \geq 1 - \delta$ .
- **Privacy.** PRIVATEPCHECK is  $(\varepsilon, \delta)$ -differentially private.

**Proof** [Proof Sketch of Theorem 17] From Corollary 16, we know that  $z = \log(p_X)$  is stable. Moreover, standard concentration bounds can be used to show that when  $p_X > p_{\text{th}}$  is not too small, the empirical estimate  $p_X^*$  of  $p_X$  from sufficiently many samples ( $N \gtrsim 1/p_{\text{th}}$ ) suffices to estimate  $z \approx z^* := \log(1/p_X^*)$ .

Therefore, so long  $p_X$  is not too small, adding Laplace noise suffices to privatize  $z^*$ , and by post-processing, we may also reveal

$$\text{DECISION} = \mathbb{1}\{z^* + \text{LAPLACE} > \text{THRESHOLD}\}$$

for any fixed THRESHOLD.

Finally, to deal with cases where  $p_X$  is very small, we simply note that these are deep into the REJECT range for  $z$ , so even a bad approximation of  $p_X$  will suffice for them. ■

For a more detailed proof of Theorem 17, see Section A.5.

### A.3. Mean-Aggregation STA

Once we have a framework for privately drawing test-passing subsamples of our training set, we want to start relating this test to the down-sensitivity of a learning algorithm, which will require some new definitions.

**Definition 18 (Down-Sensitivity Bound)** Let  $A : \mathcal{X}^* \rightarrow \mathbb{R}^d$  be a deterministic learning algorithm that maps a dataset  $X \in \mathcal{X}^*$  to a model  $A(X) \in \mathbb{R}^d$ . We say that the mapping  $E : \mathcal{X}^* \rightarrow \mathbb{S}_+^d$  is a valid down-sensitivity bound corresponding to the test  $\mathcal{T} : \mathcal{X}^* \rightarrow \{\text{ACCEPT}, \text{REJECT}\}$ , if for any dataset  $X = (X_1, \dots, X_n)$  and for any  $X^{-i}$  obtained from  $X$  by dropping a single sample  $i$ ,

$$\Pr[\mathcal{T}(X) = \text{ACCEPT}] > 0 \implies \forall i \in [n] \Delta_i \Delta_i^\top \preceq E(X) \quad \text{where} \quad \Delta_i = A(X) - A(X^{-i}).$$

We further say that  $E$  is  $C_E$ -stable, for  $C_E > 0$ , if

$$\Pr[\mathcal{T}(X) = \text{ACCEPT}] > 0 \implies E(X) \preceq C_E E(X^{-i})$$

**Definition 19 (Envelope Condition)** Let  $A : \mathcal{X}^* \rightarrow \mathbb{R}^d$  be a deterministic learning algorithm, and let  $\mathcal{T} : \mathcal{X}^* \rightarrow \{\text{ACCEPT}, \text{REJECT}\}$  be a randomized test, with down-sensitivity bound  $E : \mathcal{X}^* \rightarrow \mathbb{S}_+^d$ .

Fix PSD matrices  $\mathcal{E}, \mathcal{V} \in \mathbb{S}_+^d$ . We say that  $\mathcal{T}$  implies a  $\mathcal{E}, \mathcal{V}$  envelope on  $E$ , if

$$\Pr[\mathcal{T}(X) = \text{ACCEPT}] > 0 \implies E(X) \preceq \mathcal{E}, \quad A(X)A(X)^\top \preceq \mathcal{V}$$

**Theorem 20 (Stability of Mean-Aggregation STA)** Assume  $\mathcal{T}$  is a  $C_p, \eta$ -down-stable test for  $A$ , run under independent Poisson subsampling at rate

$$t \leq \frac{1}{c_0(1 + C_p(2 + 3C_E))}$$

for a sufficiently large universal constant  $c_0$ . Assume  $\mathcal{T}$  has corresponding  $C_E$ -stable down-sensitivity bound  $E$ , and that  $\eta = 0$  or  $\mathcal{T}$  implies a  $\mathcal{E}, \mathcal{V}$  envelope on  $E$ .

Let  $X, X' \in \mathcal{X}^n$  be datasets that differ in exactly one coordinate, such that  $p_X \stackrel{\text{def}}{=} \Pr[\text{Acc}_X] > 10\eta$ .

For a (deterministic) learning algorithm  $A : \mathcal{X}^* \rightarrow \mathbb{R}^d$ , define the subsample output,

$$v_S(X) \stackrel{\text{def}}{=} A(X_S).$$

the conditional mean

$$\mu(X) \stackrel{\text{def}}{=} \mathbb{E}[v_S(X) | \text{Acc}_X]$$

and the proxy matrix (using the returned certificate  $E_S(X)$  from running the test  $\mathcal{T}$  on subset  $S$ )

$$M(X) \stackrel{\text{def}}{=} \text{Cov}(v_S(X)) + \mathbb{E}[E_S(X) | \text{Acc}_X] \tag{3}$$

$$= \mathbb{E} \left[ (v_S(X) - \mu(X))(v_S(X) - \mu(X))^\top + E_S(X) \middle| \text{Acc}_X \right] \tag{4}$$

Define  $\Delta \stackrel{\text{def}}{=} \mu(X) - \mu(X')$ . Then there are universal constants  $c_1, \dots, c_6$  such that:

1. *Conditional mean stability.*

$$\Delta\Delta^\top \preceq c_1(C_p(1+C_E)t)^2 \cdot M(X) + c_2 \frac{C_p t^2}{p_X} \eta \mathcal{E} + c_3 \frac{C_p t^2}{p_X} \eta \mathcal{V}.$$

2. *Proxy stability (multiplicative up to additive error).*

$$M(X) \preceq (1 + c_4((C_p + C_p C_E)t)) M(X') + \frac{c_5 t}{p_X} \eta \mathcal{E} + \frac{c_6 t}{p_X} \eta \mathcal{V},$$

and

$$M(X) \succeq (1 - c_4((C_p + C_p C_E)t)) M(X') - \frac{c_5 t}{p_X} \eta \mathcal{E} - \frac{c_6 t}{p_X} \eta \mathcal{V}.$$

#### A.4. Proof of Corollary 16

**Proof** [Proof of Corollary 16] Let  $i$  be the (unique) index on which  $X$  and  $X'$  differ. If  $i \notin S$  then  $X_S = X'_S$ , so  $\mathcal{T}$  sees identical input and hence

$$\Pr[\text{Acc}_X \wedge (i \notin S)] = \Pr[\text{Acc}_{X'} \wedge (i \notin S)].$$

Therefore the entire difference comes from subsamples containing  $i$ :

$$|p_X - p_{X'}| = \left| \Pr[\text{Acc}_X \wedge (i \in S)] - \Pr[\text{Acc}_{X'} \wedge (i \in S)] \right| \leq \Pr[\text{Acc}_X \wedge (i \in S)] + \Pr[\text{Acc}_{X'} \wedge (i \in S)].$$

Writing

$$\delta_i(X) \stackrel{\text{def}}{=} \Pr[i \in S \mid \text{Acc}_X], \quad \delta_i(X') \stackrel{\text{def}}{=} \Pr[i \in S \mid \text{Acc}_{X'}],$$

this becomes

$$|p_X - p_{X'}| \leq p_X \delta_i(X) + p_{X'} \delta_i(X').$$

We now bound  $\delta_i(X)$  via Theorem 14. Let

$$\alpha \stackrel{\text{def}}{=} \frac{C_p t}{1 - t}.$$

For every  $T \subseteq [n]$  with  $i \in T$ , Theorem 14 gives

$$\Pr[S = T \mid \text{Acc}_X] \leq \alpha \left( \Pr[S = T \setminus \{i\} \mid \text{Acc}_X] + \frac{\eta}{C_p p_X} \Pr[S = T \setminus \{i\}] \right).$$

Summing over all  $T \ni i$  yields

$$\delta_i(X) = \sum_{T \ni i} \Pr[S = T \mid \text{Acc}_X] \leq \alpha \sum_{T \ni i} \Pr[S = T \setminus \{i\} \mid \text{Acc}_X] + \alpha \cdot \frac{\eta}{C_p p_X} \sum_{T \ni i} \Pr[S = T \setminus \{i\}].$$

The change of variables  $U = T \setminus \{i\}$  gives

$$\sum_{T \ni i} \Pr[S = T \setminus \{i\} \mid \text{Acc}_X] = \Pr[i \notin S \mid \text{Acc}_X] = 1 - \delta_i(X),$$

and under Poisson subsampling,

$$\sum_{T \ni i} \Pr[S = T \setminus \{i\}] = \Pr[i \notin S] = 1 - t.$$

Hence

$$\delta_i(X) \leq \alpha(1 - \delta_i(X)) + t \cdot \frac{\eta}{p_X}.$$

Rearranging,

$$(1 + \alpha)\delta_i(X) \leq \alpha + t \cdot \frac{\eta}{p_X} \implies \delta_i(X) \leq \frac{\alpha + t \cdot \eta/p_X}{1 + \alpha}.$$

The same argument gives

$$\delta_i(X') \leq \frac{\alpha + t \cdot \eta/p_{X'}}{1 + \alpha}.$$

Plugging into  $|p_X - p_{X'}| \leq p_X \delta_i(X) + p_{X'} \delta_i(X')$ ,

$$|p_X - p_{X'}| \leq \frac{\alpha p_X + t\eta}{1 + \alpha} + \frac{\alpha p_{X'} + t\eta}{1 + \alpha} = \frac{\alpha(p_X + p_{X'}) + 2t\eta}{1 + \alpha}.$$

Let  $\beta \stackrel{\text{def}}{=} \frac{\alpha}{1 + \alpha} = \frac{C_p t}{1 + (C_p - 1)t}$ , so

$$|p_X - p_{X'}| \leq \beta(p_X + p_{X'}) + \frac{2t\eta}{1 + \alpha} \leq \beta(p_X + p_{X'}) + 2t\eta.$$

Using the assumption  $p_X \geq 10\eta$ , we have  $2t\eta \leq \frac{t}{5}p_X$ , and thus

$$|p_X - p_{X'}| \leq \beta(p_X + p_{X'}) + \frac{t}{5}p_X.$$

This implies both

$$p_X \leq p_{X'} + \beta(p_X + p_{X'}) + \frac{t}{5}p_X \implies p_X \left(1 - \beta - \frac{t}{5}\right) \leq (1 + \beta)p_{X'},$$

and

$$p_X \geq p_{X'} - \beta(p_X + p_{X'}) - \frac{t}{5}p_X \implies p_X \left(1 + \beta + \frac{t}{5}\right) \geq (1 - \beta)p_{X'}.$$

Therefore

$$\frac{1 - \beta}{1 + \beta + \frac{t}{5}} p_{X'} \leq p_X \leq \frac{1 + \beta}{1 - \beta - \frac{t}{5}} p_{X'}.$$

Finally, since  $t \leq \frac{1}{5(1 + C_p)}$ , we have  $C_p t \leq \frac{1}{5}$  and also  $\beta \leq C_p t$  (and  $t/5 \leq (C_p t)/5$ ). Writing  $u \stackrel{\text{def}}{=} C_p t$ , we get

$$\frac{1 + \beta}{1 - \beta - \frac{t}{5}} = 1 + \frac{2\beta + \frac{t}{5}}{1 - \beta - \frac{t}{5}} \leq 1 + \frac{\frac{11}{5}u}{1 - \frac{6}{5}u} \leq 1 + \frac{\frac{11}{5}u}{\frac{19}{25}} < 1 + 3u = 1 + 3C_p t,$$

and similarly,

$$\frac{1 - \beta}{1 + \beta + \frac{t}{5}} = 1 - \frac{2\beta + \frac{t}{5}}{1 + \beta + \frac{t}{5}} \geq 1 - \left(2\beta + \frac{t}{5}\right) \geq 1 - \frac{11}{5}u \geq 1 - 3u = 1 - 3C_p t.$$

Plugging these into the previous display yields (1). ■

### A.5. Proof of Theorem 17

**Proof** Fix a dataset  $X$ . In Algorithm 4, each run draws an independent  $\text{Poisson}(t)$  subsample  $S_j$  and runs  $\mathcal{T}$  with fresh internal randomness. Therefore the bits

$$b_j := \mathbf{1}\{\mathcal{T}(X_{S_j}) = \text{ACCEPT}\} \quad (j = 1, \dots, N)$$

are i.i.d.  $\text{Ber}(p_X)$ , where  $p_X = \Pr[\mathcal{T}(X_S) = \text{ACCEPT}]$ . Let  $y$  denote the output of the call to PBT inside Algorithm 4, and recall that  $\text{PRIVATEPCHECK}(X)$  is the deterministic post-processing that flips  $y$ .

**Internal parameters used by PBT.** By inspection of Algorithm 4, the call to PBT uses

$$\kappa = \tau, \quad \varepsilon_0 = \varepsilon/3, \quad \delta_0 = \delta/12, \quad p_{\text{th}} = p_{\text{acc}}e^{-\gamma}$$

Define

$$\beta_0 := \delta/12, \quad C_0 := \frac{2\log(2/\delta_0)}{\varepsilon_0} = \frac{6}{\varepsilon} \log \frac{24}{\delta}.$$

Then the thresholds appearing in Theorem 46 (applied to  $\text{Ber}(p_X)$ ) are

$$p_{\text{high}}^{\text{PBT}} = p_{\text{th}}e^{\gamma} = p_{\text{acc}} = p_{\text{high}}, \quad p_{\text{low}}^{\text{PBT}} = p_{\text{th}}e^{-\kappa C_0}e^{-\gamma} = p_{\text{acc}}e^{-\tau C_0}e^{-2\gamma} = p_{\text{low}}.$$

Moreover, the sample size condition in Theorem 46 with failure probability  $\beta_0$  is exactly (2), since  $\log(4/\beta_0) = \log(48/\delta)$ .

**Completeness and soundness.** Apply Theorem 46 to the i.i.d. bits  $b_1, \dots, b_N \sim \text{Ber}(p_X)$  with parameters  $(\kappa, \varepsilon_0, \delta_0, p_{\text{th}}, \gamma, \beta_0)$ . If  $p_X \geq p_{\text{high}}$  then Theorem 46(2) implies  $y = \text{REJECT}$  with probability at least  $1 - \beta_0$ , hence  $\text{PRIVATEPCHECK}(X) = \text{ACCEPT}$  with probability at least  $1 - \beta_0$ . If  $p_X \leq p_{\text{low}}$  then Theorem 46(1) implies  $y = \text{ACCEPT}$  with probability at least  $1 - \beta_0$ , hence  $\text{PRIVATEPCHECK}(X) = \text{REJECT}$  with probability at least  $1 - \beta_0$ . Since  $\beta_0 = \delta/12 \leq \delta$ , both statements hold with probability at least  $1 - \delta$  as claimed.

**Privacy.** Fix neighboring datasets  $X \sim X'$  and write  $p := p_X$  and  $p' := p_{X'}$ . Let  $\mathcal{D}_p$  and  $\mathcal{D}_{p'}$  denote the output distributions of  $y$  under  $\text{Ber}(p)^{\otimes N}$  and  $\text{Ber}(p')^{\otimes N}$ , respectively. Since  $\text{PRIVATEPCHECK}$  is deterministic post-processing of  $y$ , it suffices to show  $\mathcal{D}_p \approx_{\varepsilon, \delta} \mathcal{D}_{p'}$ .

We consider two cases.

*Case 1:*  $\min\{p, p'\} \geq p_{\text{low}}/2$ . By  $\eta \leq p_{\text{low}}/20$ , this implies  $\min\{p, p'\} \geq 10\eta$ , so Corollary 16 applies (to both orientations) and yields

$$(1 - 3C_p t)p' \leq p \leq (1 + 3C_p t)p'.$$

Let  $r := 1 - 3C_p t$ , so  $\tau = -\log r$  and  $e^{-\tau} = r$ . From the above equation,

$$p' \leq \frac{p}{r} = e^{\tau}p, \quad p' \geq \frac{p}{1 + 3C_p t} \geq (1 - 3C_p t)p = rp = e^{-\tau}p,$$

where we used  $\frac{1}{1+a} \geq 1-a$  for  $a \geq 0$ . Hence  $e^{-\kappa}p \leq p' \leq e^{\kappa}p$  with  $\kappa = \tau$ , and also  $p, p' \geq p_{\text{low}}/2$ . Therefore Theorem 46(3) applies to  $\mathcal{D}_p$  and  $\mathcal{D}_{p'}$  and gives

$$\mathcal{D}_p \approx_{3\varepsilon_0, \delta_{\text{priv}}} \mathcal{D}_{p'}, \quad \delta_{\text{priv}} = (1 + e^{\varepsilon_0} + e^{2\varepsilon_0})\delta_0 + (1 + e^{2\varepsilon_0})\beta_0.$$

Since  $\varepsilon \in (0, 1)$  we have  $e^{\varepsilon_0} = e^{\varepsilon/3} < 2$  and  $e^{2\varepsilon_0} < 2$ , so

$$\delta_{\text{priv}} \leq (1 + 2 + 2) \frac{\delta}{12} + (1 + 2) \frac{\delta}{12} = \frac{8}{12} \delta \leq \delta.$$

Also  $3\varepsilon_0 = \varepsilon$ . Thus  $\mathcal{D}_p \approx_{\varepsilon, \delta} \mathcal{D}_{p'}$ .

*Case 2:*  $\min\{p, p'\} < p_{\text{low}}/2$ . Without loss of generality assume  $p < p_{\text{low}}/2$ . We claim that then  $p' \leq p_{\text{low}}$ .

If  $p' < 10\eta$ , then  $p' \leq 10\eta \leq p_{\text{low}}/2 < p_{\text{low}}$ . Otherwise  $p' \geq 10\eta$ , so Corollary 16 applies to  $X'$  (with roles swapped) and yields  $p' \leq (1 + 3C_p t) p$ . Since  $t$  satisfies the conditions of Corollary 16, we have  $3C_p t < 1$ ; thus using  $p < p_{\text{low}}/2$  gives  $p' < p_{\text{low}}$ . Hence  $p, p' \leq p_{\text{low}}$ .

Now apply Theorem 46(1): under both  $p$  and  $p'$ , PBT outputs **ACCEPT** with probability at least  $1 - \beta_0$ . Equivalently,

$$\mathcal{D}_p(\text{REJECT}) \leq \beta_0, \quad \mathcal{D}_{p'}(\text{REJECT}) \leq \beta_0.$$

Since the output space is  $\{\text{ACCEPT}, \text{REJECT}\}$ , for every event  $\mathcal{S} \subseteq \{\text{ACCEPT}, \text{REJECT}\}$  we have

$$|\mathcal{D}_p(\mathcal{S}) - \mathcal{D}_{p'}(\mathcal{S})| \leq \beta_0 \leq \delta.$$

In particular, for every  $\mathcal{S}$ ,  $\mathcal{D}_p(\mathcal{S}) \leq \mathcal{D}_{p'}(\mathcal{S}) + \delta \leq e^\varepsilon \mathcal{D}_{p'}(\mathcal{S}) + \delta$ , and symmetrically with  $p$  and  $p'$  swapped. Hence  $\mathcal{D}_p \approx_{\varepsilon, \delta} \mathcal{D}_{p'}$  in this case as well.

Combining the two cases proves  $\mathcal{D}_p \approx_{\varepsilon, \delta} \mathcal{D}_{p'}$  for all neighboring  $X \sim X'$ . Therefore the internal output  $y$  is  $(\varepsilon, \delta)$ -DP, and since **PRIVATEPCHECK** is a deterministic post-processing of  $y$ , it is also  $(\varepsilon, \delta)$ -DP.  $\blacksquare$

## A.6. Proof of Theorem 20

Our privacy analysis studies distributions and expectations *conditioned on* **Acc**, and the following lemmas show that this conditioning does not allow any single example to have outsized influence.

### Lemma 21 (No single sample dominates the expected certificate)

Assume  $\mathcal{T}$  is a  $C_p, \eta$ -down-stable test for  $A$ , with corresponding  $C_E$ -stable down-sensitivity bound  $E$ , and that  $\mathcal{T}$  implies a  $\mathcal{E}, \mathcal{V}$  envelope on  $E$ . Then for every  $i \in [n]$ ,

$$\mathbb{E}[E_S \cdot \mathbf{1}\{i \in S\} | \text{Acc}_X] \preceq \frac{tC_p C_E}{1-t} \cdot \mathbb{E}[E_S | \text{Acc}_X] + t \cdot \frac{\eta \mathcal{E}}{p_X}.$$

**Remark 22 (Interpretation of Lemma 21)** Lemma 21 formalizes the intuition that, provided the marginal acceptance probability  $p$  is not too small, conditioning on **Acc** cannot cause any single point to dominate the typical sensitivity certificate  $E_S$ . The second term is an additive error that becomes negligible when  $\eta/p$  is small (since  $\mathcal{E}$  upper bounds  $E_S$  on acceptance).

**Proof** [Proof of Lemma 21] Fix  $i \in [n]$ . Expanding the conditional expectation,

$$\mathbb{E}[E_S \cdot \mathbf{1}\{i \in S\} | \text{Acc}_X] = \sum_{S \ni i} \Pr[S | \text{Acc}_X] E_S.$$

Apply Theorem 14 (with  $R = S \setminus \{i\}$ ) to each summand:

$$\sum_{S \ni i} \Pr[S \mid \text{Acc}] E_S \preceq \frac{tC_p}{1-t} \sum_{S \ni i} \Pr[R \mid \text{Acc}_X] E_S + \frac{t\eta}{(1-t)p_X} \sum_{S \ni i} \Pr[R] E_S.$$

For the first term, note that under the conditioning on  $\text{Acc}$ , every set  $S$  in the sum satisfies  $b(X_S) = \text{ACCEPT}$ , and hence  $C_E$  stability (Definition 18) gives

$$E_S \preceq C_E E_{S \setminus \{i\}} = C_E E_R.$$

Therefore,

$$\frac{tC_p}{1-t} \sum_{S \ni i} \Pr[R \mid \text{Acc}_X] E_S \preceq \frac{tC_p C_E}{1-t} \sum_{S \ni i} \Pr[R \mid \text{Acc}_X] E_R.$$

Changing variables from  $S \ni i$  to  $R = S \setminus \{i\}$  (which ranges over sets  $R \subseteq [n] \setminus \{i\}$ ), the term on the right above becomes

$$\sum_{S \ni i} \Pr[R \mid \text{Acc}_X] E_R = \sum_{R: i \notin R} \Pr[R \mid \text{Acc}_X] E_R \preceq \sum_R \Pr[R \mid \text{Acc}_X] E_R = \mathbb{E}[E_S \mid \text{Acc}_X],$$

since all matrices are PSD.

For the second term, again every  $S$  in the sum satisfies  $b(X_S) = \text{ACCEPT}$ , so the envelope condition (Definition 19) implies  $E_S \preceq \mathcal{E}$ . Hence

$$\frac{t\eta}{(1-t)p_X} \sum_{S \ni i} \Pr[R] E_S \preceq \frac{t\eta}{(1-t)p_X} \sum_{S \ni i} \Pr[R] \cdot \mathcal{E} = \frac{t\eta}{(1-t)p_X} \left( \sum_{S \ni i} \Pr[S \setminus \{i\}] \right) \mathcal{E}.$$

Finally,  $\sum_{S \ni i} \Pr[S \setminus \{i\}] = \Pr[i \notin S] = 1 - t$  under Poisson subsampling, so this term equals  $\frac{t\eta}{p_X} \mathcal{E}$ .

Combining the bounds proves the lemma.  $\blacksquare$

#### A.6.1. AUXILIARY INEQUALITIES

##### Lemma 23 (Bounded drop-set second moment implies bounded mean shift)

Let  $\mathbf{v}$  be an  $\mathbb{R}^d$ -valued random vector and let  $\Sigma \in \mathbb{S}_+^d$  be any PSD matrix. Let  $D$  be an event with  $\Pr[D] \geq 1 - \delta$ . Suppose there is a parameter  $\gamma \geq 0$  such that

$$\mathbb{E} \left[ \mathbb{1}\{D^c\} \mathbf{v} \mathbf{v}^\top \right] \preceq \gamma \Sigma.$$

Define the (truncation) mean shift

$$\mathbf{m} \stackrel{\text{def}}{=} \mathbb{E}[\mathbf{v}] - \mathbb{E}[\mathbb{1}\{D\} \mathbf{v}] = \mathbb{E}[\mathbb{1}\{D^c\} \mathbf{v}].$$

Then

$$\mathbf{m} \mathbf{m}^\top \preceq \gamma \delta \Sigma.$$

In particular, if  $\Sigma \succ 0$ , then  $\|\mathbf{m}\|_{\Sigma^{-1}} \leq \sqrt{\gamma \delta}$ .

**Proof** Fix any  $\mathbf{u} \in \mathbb{R}^d$ . Since  $\mathbf{m} = \mathbb{E}[\mathbf{1}\{D^c\}\mathbf{v}]$ , we have

$$\langle \mathbf{u}, \mathbf{m} \rangle = \mathbb{E}[\mathbf{1}\{D^c\} \langle \mathbf{u}, \mathbf{v} \rangle].$$

By Cauchy–Schwarz,

$$\langle \mathbf{u}, \mathbf{m} \rangle^2 \leq \mathbb{E}[\mathbf{1}\{D^c\}] \cdot \mathbb{E}[\mathbf{1}\{D^c\} \langle \mathbf{u}, \mathbf{v} \rangle^2] \leq \delta \cdot \mathbf{u}^\top \mathbb{E}[\mathbf{1}\{D^c\} \mathbf{v} \mathbf{v}^\top] \mathbf{u} \leq \gamma \delta \cdot \mathbf{u}^\top \Sigma \mathbf{u}.$$

Equivalently,

$$\mathbf{u}^\top (\mathbf{m} \mathbf{m}^\top) \mathbf{u} \leq \gamma \delta \cdot \mathbf{u}^\top \Sigma \mathbf{u} \quad \text{for all } \mathbf{u} \in \mathbb{R}^d.$$

This implies  $\mathbf{m} \mathbf{m}^\top \preceq \gamma \delta \Sigma$ .

If  $\Sigma \succ 0$ , taking the supremum of  $\langle \mathbf{u}, \mathbf{m} \rangle^2 \leq \gamma \delta \|\mathbf{u}\|_\Sigma^2$  over  $\mathbf{u} \neq 0$  yields  $\|\mathbf{m}\|_{\Sigma^{-1}}^2 \leq \gamma \delta$ . ■

**Lemma 24 (A useful PSD inequality)** For any  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ ,

$$(\mathbf{a} + \mathbf{b})(\mathbf{a} + \mathbf{b})^\top \preceq 2 \mathbf{a} \mathbf{a}^\top + 2 \mathbf{b} \mathbf{b}^\top.$$

**Proof** Rearranging,

$$2 \mathbf{a} \mathbf{a}^\top + 2 \mathbf{b} \mathbf{b}^\top - (\mathbf{a} + \mathbf{b})(\mathbf{a} + \mathbf{b})^\top = (\mathbf{a} - \mathbf{b})(\mathbf{a} - \mathbf{b})^\top \succeq 0,$$

which proves the claim. ■

**Corollary 25 (Propagating a leave-one-out certificate to a second-moment bound)** Let  $\mathbf{o} \in \mathbb{R}^d$  be any fixed offset. Let  $S \subseteq [n]$  and suppose  $\mathcal{T}(X_S) = (\text{ACCEPT}, E_S)$ . Fix any  $i \in S$ , and define

$$\mathbf{v} \stackrel{\text{def}}{=} A(X_S) + \mathbf{o}, \quad \mathbf{u} \stackrel{\text{def}}{=} A(X_{S \setminus \{i\}}) + \mathbf{o}.$$

Then

$$\mathbf{v} \mathbf{v}^\top \preceq 2(\mathbf{u} \mathbf{u}^\top + E_S).$$

**Proof** Let  $\Delta \stackrel{\text{def}}{=} \mathbf{v} - \mathbf{u} = A(X_S) - A(X_{S \setminus \{i\}})$ . Since  $\mathcal{T}(X_S)$  accepts, the down-sensitivity certificate condition (Definition 18(1)) implies  $\Delta \Delta^\top \preceq E_S$ . Applying Lemma 24 with  $\mathbf{a} = \mathbf{u}$  and  $\mathbf{b} = \Delta$  gives

$$\mathbf{v} \mathbf{v}^\top = (\mathbf{u} + \Delta)(\mathbf{u} + \Delta)^\top \preceq 2 \mathbf{u} \mathbf{u}^\top + 2 \Delta \Delta^\top \preceq 2 \mathbf{u} \mathbf{u}^\top + 2 E_S,$$

which is equivalent to the stated bound. ■

## A.6.2. STABILITY UNDER NEIGHBORING DATASETS

We now formalize the key structural fact needed for privacy: under a stable down-sensitivity test, the *conditional* distribution over accepted subsamples does not change too much when we modify a single datapoint. Concretely, this implies that both (i) the conditional mean of the accepted estimator and (ii) the associated “proxy covariance” matrix are stable under neighboring datasets.

A.6.3. CONDITIONING ON  $j \notin S$  DOES NOT SHIFT THE MEAN TOO MUCH

**Lemma 26 (Mean shift from excluding one index)** *Assume  $\mathcal{T}$  is a  $C_p, \eta$ -down-stable test for  $A$ , with corresponding  $C_E$ -stable down-sensitivity bound  $E$ , and that  $\mathcal{T}$  implies a  $\mathcal{E}, \mathcal{V}$  envelope on  $E$ . Let*

$$\mu^{-j}(X) \stackrel{\text{def}}{=} \mathbb{E}[v_S(X) | \text{Acc}_X \wedge (j \notin S)] \quad \text{and} \quad \Delta_j(X) \stackrel{\text{def}}{=} \mu(X) - \mu^{-j}(X).$$

Fix  $X \in \mathcal{X}^n$  and  $j \in [n]$ . Let  $\delta_j(X) \stackrel{\text{def}}{=} \Pr[j \in S | \text{Acc}_X]$ . Then

$$\Delta_j(X) \Delta_j(X)^\top \preceq \frac{\delta_j(X)}{(1 - \delta_j(X))^2} \cdot \left( \frac{2tC_p(1 + C_E)}{1 - t} M(X) + \frac{2t}{p_X} \eta \mathcal{E} + \frac{8t}{(1 - t)p_X} \eta \mathcal{V} \right). \quad (5)$$

**Proof** Let  $\tilde{v}_S(X) \stackrel{\text{def}}{=} v_S(X) - \mu(X)$ . Then  $\mathbb{E}[\tilde{v}_S(X) | \text{Acc}_X] = 0$  and

$$\Delta_j(X) = \mu(X) - \mu^{-j}(X) = -\mathbb{E}[\tilde{v}_S(X) | \text{Acc}_X \wedge (j \notin S)].$$

Let  $D$  be the event  $\{j \notin S\}$ , so  $\Pr[D | \text{Acc}_X] = 1 - \delta_j(X)$ . Applying Lemma 23 to the random vector  $\tilde{v}_S(X)$  under the conditional distribution  $\Pr[\cdot | \text{Acc}_X]$  and the event  $D$  gives (after dividing by  $\Pr[D | \text{Acc}_X]^2$  from both sides and using  $\mathbb{E}[\tilde{v}_S(X) | \text{Acc}_X] = 0$ )

$$(\mathbb{E}[\tilde{v}_S(X) | \text{Acc}_X \wedge D]) (\mathbb{E}[\tilde{v}_S(X) | \text{Acc}_X \wedge D])^\top \preceq \frac{\delta_j(X)}{(1 - \delta_j(X))^2} \cdot \mathbb{E}[\tilde{v}_S(X) \tilde{v}_S(X)^\top \cdot \mathbf{1}\{j \in S\} | \text{Acc}_X].$$

Since  $\Delta_j(X) = -\mathbb{E}[\tilde{v}_S(X) | \text{Acc}_X \wedge D]$ , it remains to upper bound the RHS.

On the event  $\text{Acc}_X$  and  $j \in S$ , the down-sensitivity bound (Definition 18) with  $X^{-j}$  obtained by dropping  $j$  implies

$$(v_S(X) - v_{S \setminus \{j\}}(X)) (v_S(X) - v_{S \setminus \{j\}}(X))^\top \preceq E_S(X).$$

By the corollary of Lemma 24, applied to the centered vectors

$$\tilde{v}_S(X) = \tilde{v}_{S \setminus \{j\}}(X) + (v_S(X) - v_{S \setminus \{j\}}(X)),$$

we get

$$\tilde{v}_S(X) \tilde{v}_S(X)^\top \preceq 2\tilde{v}_{S \setminus \{j\}}(X) \tilde{v}_{S \setminus \{j\}}(X)^\top + 2E_S(X).$$

Multiplying by  $\mathbf{1}\{j \in S\}$  and taking  $\mathbb{E}[\cdot | \text{Acc}_X]$  yields

$$\mathbb{E}[\tilde{v}_S \tilde{v}_S^\top \cdot \mathbf{1}\{j \in S\} | \text{Acc}_X] \preceq 2\mathbb{E}[\tilde{v}_{S \setminus \{j\}} \tilde{v}_{S \setminus \{j\}}^\top \cdot \mathbf{1}\{j \in S\} | \text{Acc}_X] + 2\mathbb{E}[E_S \cdot \mathbf{1}\{j \in S\} | \text{Acc}_X]. \quad (6)$$

We bound the second term using Lemma 21, giving

$$\mathbb{E}[E_S \cdot \mathbf{1}\{j \in S\} | \text{Acc}_X] \preceq \frac{tC_p C_E}{1 - t} \cdot \mathbb{E}[E_S | \text{Acc}_X] + \frac{t}{p_X} \eta \mathcal{E}. \quad (7)$$

For the first term in (6), apply Theorem 14 to reweight from sets containing  $j$  to sets with  $j$  removed. Concretely, Theorem 14 implies (after the  $S \mapsto S \setminus \{j\}$  change of variables) that

$$\mathbb{E}[\tilde{v}_{S \setminus \{j\}} \tilde{v}_{S \setminus \{j\}}^\top \cdot \mathbf{1}\{j \in S\} | \text{Acc}_X] \preceq \frac{tC_p}{1 - t} \cdot \mathbb{E}[\tilde{v}_S \tilde{v}_S^\top | \text{Acc}_X] + \frac{t\eta}{(1 - t)p_X} \cdot \mathbb{E}_S[\tilde{v}_S \tilde{v}_S^\top].$$

By Definition 19 and Lemma 24,  $\tilde{v}_S \tilde{v}_S^\top \preceq 2v_S v_S^\top + 2\mu(X)\mu(X)^\top \preceq 4\mathcal{V}$ , hence  $\mathbb{E}_S [\tilde{v}_S \tilde{v}_S^\top] \preceq 4\mathcal{V}$ . Therefore,

$$\mathbb{E} \left[ \tilde{v}_{S \setminus \{j\}} \tilde{v}_{S \setminus \{j\}}^\top \cdot \mathbf{1}\{j \in S\} \middle| \text{Acc}_X \right] \preceq \frac{tC_p}{1-t} \cdot \mathbb{E} \left[ \tilde{v}_S \tilde{v}_S^\top \middle| \text{Acc}_X \right] + \frac{4t\eta}{(1-t)p_X} \mathcal{V} \quad (8)$$

Plugging (7) and (8) into (6), and using the definition of  $M(X)$  in (3) (so that  $\mathbb{E} [\tilde{v}_S \tilde{v}_S^\top \middle| \text{Acc}_X] \preceq M(X)$  and  $\mathbb{E} [E_S \middle| \text{Acc}_X] \preceq M(X)$ ) yields

$$\mathbb{E} \left[ \tilde{v}_S \tilde{v}_S^\top \cdot \mathbf{1}\{j \in S\} \middle| \text{Acc}_X \right] \preceq \frac{2tC_p(1+C_E)}{1-t} M(X) + \frac{2t\eta}{p_X} \mathcal{E} + \frac{8t\eta}{(1-t)p_X} \mathcal{V}.$$

Substituting this bound into the initial application of Lemma 23 completes the proof of (5).  $\blacksquare$

#### A.6.4. STABILITY OF $M$ UNDER DELETIONS

Let  $X^{-j}$  denote the dataset obtained by dropping the  $j$ th sample from  $X$ . Under Poisson sampling, sampling from  $X^{-j}$  is equivalent to sampling  $S$  from  $X$  and conditioning on  $j \notin S$ ; moreover, on this event the subsample itself is identical.

**Lemma 27 (Proxy stability under dropping one sample)** *Assume  $\mathcal{T}$  is a  $C_p, \eta$ -down-stable test for  $A$ , with corresponding  $C_E$ -stable down-sensitivity bound  $E$ , and that  $\mathcal{T}$  implies a  $\mathcal{E}, \mathcal{V}$  envelope on  $E$ . Fix  $X \in \mathcal{X}^n$  and  $j \in [n]$ .*

*Let  $\delta_j(X), \Delta_j(X)$  be as in Lemma 26 and define*

$$a \stackrel{\text{def}}{=} \frac{tC_p}{1-t} (2 + 3C_E).$$

*Then:*

$$(1 - \delta_j(X)) \cdot M(X^{-j}) \preceq M(X), \quad (9)$$

$$M(X) \preceq \frac{1}{1-a} \cdot \left( M(X^{-j}) + \Delta_j(X) \Delta_j(X)^\top + \frac{3t\eta}{p_X} \mathcal{E} + \frac{8t\eta}{(1-t)p_X} \mathcal{V} \right), \quad (10)$$

*whenever  $a < 1$ .*

**Proof** Let  $\tilde{v}_S(X) = v_S(X) - \mu(X)$  and define the PSD random matrix

$$W_S(X) \stackrel{\text{def}}{=} \tilde{v}_S(X) \tilde{v}_S(X)^\top + E_S(X).$$

By definition,  $M(X) = \mathbb{E} [W_S(X) \middle| \text{Acc}_X]$ .

Let  $D$  be the event  $j \notin S$ . As noted above, the conditional distribution of  $X_S$  given  $\text{Acc}_X \wedge D$  is the same as the accepted-subsample distribution for  $X^{-j}$ . A direct expansion gives

$$M(X^{-j}) = \mathbb{E} [W_S(X) \middle| \text{Acc}_X \wedge D] - \Delta_j(X) \Delta_j(X)^\top. \quad (11)$$

Next, observe that

$$\begin{aligned}\mathbb{E}[W_S(X)|\text{Acc}_X \wedge D] &= \frac{1}{1 - \delta_j(X)} \cdot \mathbb{E}[W_S(X) \cdot \mathbf{1}\{D\}|\text{Acc}_X] \\ &= \frac{1}{1 - \delta_j(X)} \cdot (M(X) - \mathbb{E}[W_S(X) \cdot \mathbf{1}\{j \in S\}|\text{Acc}_X]).\end{aligned}$$

Plugging into (11) yields

$$M(X^{-j}) = \frac{1}{1 - \delta_j(X)} \cdot (M(X) - \mathbb{E}[W_S(X) \cdot \mathbf{1}\{j \in S\}|\text{Acc}_X]) - \Delta_j(X)\Delta_j(X)^\top. \quad (12)$$

Since  $\mathbb{E}[W_S(X) \cdot \mathbf{1}\{j \in S\}|\text{Acc}_X] \geq 0$  and  $\Delta_j(X)\Delta_j(X)^\top \geq 0$ , we immediately get  $M(X^{-j}) \leq \frac{1}{1 - \delta_j(X)}M(X)$ , which is equivalent to (9).

For the other direction, rearrange (12) as

$$M(X) = (1 - \delta_j(X)) \left( M(X^{-j}) + \Delta_j(X)\Delta_j(X)^\top \right) + \mathbb{E}[W_S(X) \cdot \mathbf{1}\{j \in S\}|\text{Acc}_X]. \quad (13)$$

We upper bound the last term. By Lemma 26's intermediate bounds (equations 6, 7, 8), we have

$$\begin{aligned}\mathbb{E}\left[\tilde{v}_S\tilde{v}_S^\top \cdot \mathbf{1}\{j \in S\}|\text{Acc}_X\right] &\leq \frac{2tC_p}{1-t} \cdot \mathbb{E}\left[\tilde{v}_S\tilde{v}_S^\top|\text{Acc}_X\right] \\ &\quad + \frac{2tC_pC_E}{1-t} \cdot \mathbb{E}[E_S|\text{Acc}_X] \\ &\quad + \frac{2t\eta}{p_X}\mathcal{E} + \frac{8t\eta}{(1-t)p_X}\mathcal{V}.\end{aligned}$$

and by Lemma 21,

$$\mathbb{E}[E_S \cdot \mathbf{1}\{j \in S\}|\text{Acc}_X] \leq \frac{tC_pC_E}{1-t} \cdot \mathbb{E}[E_S|\text{Acc}_X] + \frac{t\eta}{p_X}\mathcal{E}.$$

Adding these gives

$$\begin{aligned}\mathbb{E}[W_S(X) \cdot \mathbf{1}\{j \in S\}|\text{Acc}_X] &\leq \frac{tC_p}{1-t} (2 + 3C_E) M(X) + \frac{3t\eta}{p_X}\mathcal{E} + \frac{8t\eta}{(1-t)p_X}\mathcal{V} \\ &= a M(X) + \frac{3t\eta}{p_X}\mathcal{E} + \frac{8t\eta}{(1-t)p_X}\mathcal{V}.\end{aligned}$$

Plugging into (13) and using  $1 - \delta_j(X) \leq 1$  yields

$$M(X) \leq M(X^{-j}) + \Delta_j(X)\Delta_j(X)^\top + a M(X) + \frac{3t\eta}{p_X}\mathcal{E} + \frac{8t\eta}{(1-t)p_X}\mathcal{V}.$$

Rearranging completes the proof of (10). ■

### A.6.5. CONCLUDING THEOREM 20

We now upgrade Lemmas 26 and 27 into the main theorem of this subsection for *substitution* neighbors. We'll use the fact from Theorem 14 that  $\delta_j(X) = O(t)$  when  $p_X \gtrsim \eta$ .

**Proof**

For (i), note that on the event  $j \notin S$  we have  $X_S = X'_S$ , hence  $v_S(X) = v_S(X')$  and the test output is identical; therefore  $\mu^{-j}(X) = \mu^{-j}(X')$ . Thus

$$\mu(X) - \mu(X') = (\mu(X) - \mu^{-j}(X)) - (\mu(X') - \mu^{-j}(X')).$$

Applying Lemma 24 gives  $\Delta\Delta^\top \preceq 2\Delta_j(X)\Delta_j(X)^\top + 2\Delta_j(X')\Delta_j(X')^\top$ . Lemma 26 bounds each term.

For (ii), let  $Y = X^{-j} = (X')^{-j}$ . By Lemma 27 applied to  $X$ , we have  $M(Y) \preceq \frac{1}{1-\delta_j(X)}M(X)$ . Applying Lemma 27 to  $X'$  and substituting this bound on  $M(Y)$ , along with Lemma 26 on  $\Delta_j(X')\Delta_j(X')$  yields the displayed inequality. The reverse direction follows symmetrically.  $\blacksquare$

## Appendix B. Black-Box Aggregation via Advanced Composition

### B.1. Advanced Composition

In this section, we will analyze the privacy of our main STA mechanism – Algorithm 5. The key idea is that every draw from

$$\left( \mathcal{N}(A(X_S), O(\log(1/\delta))E(X_S)) \mid S \sim \text{Acc}_X \right)$$

is  $(O(t), \delta)$ -private (Lemma 28).

Using just a single draw from this distribution would not give us good utility, since it would only utilize a small fraction of our dataset. Therefore, we would want to average over  $k \gg 1$  draws from the mechanism above. Utilizing the advanced composition theorem (Theorem 29), we can achieve this for a very small cost to the privacy of our algorithm, since advanced composition yields a  $(O(t\sqrt{k\log(1/\delta)} + kt^2), O(k\delta))$ -privacy guarantee.

**Lemma 28 (Privacy of a Single Draw)** *Let  $\mathcal{T}$  be a  $(C_p, \eta)$ -down-stable test for a down-stability bound  $E(\cdot)$ , and assume that*

$$E(X_S) \preceq E(X_{S \setminus \{i\}}) \quad \text{and} \quad C_E |E(X_{S \setminus \{i\}})| \leq |E(X_S)| \leq |E(X_{S \setminus \{i\}})|.$$

Assume  $t \leq \frac{1}{5(C_p+1)}$  and fix any  $\delta \in (0, 1/10)$ .

Let  $X, X' \in \mathcal{X}^n$  be neighboring datasets such that  $p_X \geq \frac{2\eta}{\delta_{\text{draw}}}$ . Then

$$\left( \mathcal{N}(A(X_S), \sigma^2 E(X_S)) \mid S \sim \text{Acc}_X \right) \approx_{\varepsilon_{\text{draw}}, \delta_{\text{draw}}} \left( \mathcal{N}(A(X'_S), \sigma^2 E(X'_S)) \mid S \sim \text{Acc}_{X'} \right),$$

for  $\sigma^2 \stackrel{\text{def}}{=} 2 \log(5/\delta_{\text{draw}})$  and  $\varepsilon_{\text{draw}} = O(t)$ .

**Theorem 29 (Advanced composition (Dwork and Roth, 2014, Thm. 3.20))** *Let  $k \in \mathbb{N}$ . For each  $i \in [k]$ , let  $\mathcal{M}_i$  be an (interactive) mechanism that is  $(\varepsilon, \delta)$ -differentially private. Then for every  $\delta' \in (0, 1)$ , the adaptive  $k$ -fold composition  $\mathcal{M}_{1:k}$  is*

$$\left( \varepsilon\sqrt{2k\log(1/\delta')} + k\varepsilon(e^\varepsilon - 1), k\delta + \delta' \right)\text{-DP}.$$

**Algorithm 5:** AC-STA: STA with a Black-Box Aggregator

**Input:**  $X \in \mathcal{X}^*$ ,  $k \in \mathbb{N}$ ,  $\text{AGG} : \mathbb{R}^{d \times k} \rightarrow \mathbb{R}^d$

**Output:** a model in  $\mathbb{R}^d$  or REJECT

Draw  $S_1, \dots, S_k \stackrel{\text{iid}}{\sim} \text{Acc}_X$

**for**  $j \leftarrow 1$  **to**  $k$  **do**

    | Draw  $v_j \leftarrow \mathcal{N}(A(X_{S_j}), 2 \log(5/\delta) E(X_{S_j}))$

**end**

**return**  $\text{AGG}(v_1, \dots, v_k)$

Averaging over  $k \approx \frac{\log(1/\delta)}{\alpha^2}$  draws from this mechanism suffices for the utility of our private regression algorithm. The advanced composition theorem allows us to ensure privacy for any aggregation function so long as  $t \leq \frac{\varepsilon \alpha}{\log(1/\delta)}$ .

**Theorem 30 (Privacy of AC-STA)** *Under the assumptions of Lemma 28, the output of Algorithm 5 is private, in the sense that for any  $\delta > 0$  and for any datasets  $X, X'$  as in Lemma 28, it satisfies*

$$\text{AC-STA}(X) \approx \sqrt{k \log(1/\delta) \varepsilon_{\text{draw}} + 2k \varepsilon_{\text{draw}}^2, k \delta_{\text{draw}} + \delta} \text{AC-STA}(X').$$

Theorem 30 follows immediately from Lemma 28 and Theorem 29. The rest of this section will be devoted to proving Lemma 28. Throughout this section, let  $X, X'$  be neighboring datasets and let  $i \in [n]$  be the index on which they differ.

## B.2. Proof of Lemma 28 – Privacy of a Single Draw

We now proceed to prove that a single draw from  $(\mathcal{N}(A(X_S), \log(1/\delta)E(X_S)) \mid S \sim \text{Acc}_X)$  is private. The key idea we will use is that Theorem 14 will allow us to show that  $S \setminus \{i\}$  behaves the same for the datasets  $X$  and  $X'$  (Lemma 31).

**Lemma 31 ( $S \setminus \{i\}$  is Private)** *Under the assumptions of Lemma 28,*

$$(S \setminus \{i\} \mid S \sim \text{Acc}_X) \approx_{\varepsilon_S, \delta_S} (S \setminus \{i\} \mid S \sim \text{Acc}_{X'}),$$

where

$$\varepsilon_S := \log \left( \frac{1 + \frac{C_p t}{1-t}}{1 - 3C_p t} \right) = O(t) \quad \text{and} \quad \delta_S := \frac{t}{1-t} (1 + 3C_p t) \cdot \frac{\eta}{p_X} = O \left( \frac{\eta}{p_X} \right) \quad (14)$$

It is important to note that Lemma 31 *does not* mean that we can reveal  $S \setminus \{i\}$  while maintaining privacy, since the choice of  $i$  depended on  $X, X'$ , while privacy of a statistic would require that it is private for any neighboring datasets.

We will then combine Theorem 14 with our assumptions on the properties of  $E$  to show that the  $\mathcal{N}(0, \log(1/\delta)E(X_S))$  normal noise is itself private (Lemma 32). The idea is that if we are given  $T = S \setminus \{i\}$ , but not told if  $S = T$  or  $S = T \cup \{i\}$ , then  $S = T$  is much more likely (from Theorem 14). If we are then also told that the output of the mechanism  $v$  is very far from 0, then the property  $E(X_T) \succeq E(X_{T \cup \{i\}})$  tells us that  $v$  is more likely to be generated from

$S = T$  than from  $S = T \cup \{i\}$ . Alternatively, if  $v$  is close to 0, then the second assumption that  $\det(E(X_T)) = O(\det(E(X_{T \cup \{i\}})))$  tells us that the likelihood of drawing  $v$  from  $T \cup \{i\}$  also cannot be too high. Combining the two tells us that regardless of the value of  $v$ , its likelihood of being drawn from  $T$  or a mixture of  $T$  and  $T \cup \{i\}$  are similar.

**Lemma 32** ( $\mathcal{N}(0, E(X_S))$  is Private) *Under the assumptions of Lemma 28,*

$$\left( \mathcal{N}(0, E(X_S)) \mid S \sim \text{Acc}_X \right) \approx_{\varepsilon_E, \delta_E} \left( \mathcal{N}(0, E(X_{S \setminus \{i\}})) \mid S \sim \text{Acc}_X \right),$$

where

$$\varepsilon_E := \log \left( 1 + \sqrt{C_E} \cdot \frac{C_p t}{1-t} \right) \quad \text{and} \quad \delta_E := \sqrt{C_E} \cdot \frac{t}{1-t} \cdot \frac{\eta}{p_X}. \quad (15)$$

We then utilize a similar argument, but this time combined with a standard argument for the privacy of the Gaussian mechanism to shifts to prove Lemma 33.

**Lemma 33 (The Gaussian Mechanism for Spikes)** *Fix  $\tau \in [0, 1]$  and  $\delta \in (0, 1/2)$ . Let  $\Sigma \in \mathbb{S}_+^d$  be deterministic, and let  $v \sim \mathcal{V} \in \mathbb{R}^d$  be random such that*

$$\Pr_{v \sim \mathcal{V}} [v = 0] \geq 1 - \tau \quad \text{and} \quad \Sigma \succeq 2 \log \left( \frac{4}{\delta} \right) v v^\top \text{ almost surely.}$$

*Interpreting  $\mathcal{N}(v, \Sigma)$  as the mixture distribution obtained by first drawing  $v \sim \mathcal{V}$  and then drawing  $z \sim \mathcal{N}(v, \Sigma)$ , we have*

$$\mathcal{N}(0, \Sigma) \approx_{\varepsilon, \delta'} \mathcal{N}(v, \Sigma),$$

where one may take

$$\varepsilon := \log((1 - \tau) + \tau e) = \log(1 + \tau(e - 1)) \quad \text{and} \quad \delta' := \tau \cdot \frac{\delta}{2}. \quad (16)$$

*In particular,  $\varepsilon \leq (e - 1)\tau = O(\tau)$  and  $\delta' \leq \delta$ .*

### B.2.1. PROOF OF LEMMA 31

**Proof** [Proof of Lemma 31] Fix any set  $R \subseteq [n] \setminus \{i\}$ . By a simple partition,

$$\Pr[S \setminus \{i\} = R \mid S \sim \text{Acc}_X] = \Pr[S = R \mid S \sim \text{Acc}_X] + \Pr[S = R \cup \{i\} \mid S \sim \text{Acc}_X]. \quad (17)$$

Applying Theorem 14 to the set  $T := R \cup \{i\}$  (so  $T \setminus \{i\} = R$ ), we get

$$\Pr[S = R \cup \{i\} \mid S \sim \text{Acc}_X] \leq \frac{C_p t}{1-t} \cdot \left( \Pr[S = R \mid S \sim \text{Acc}_X] + \frac{\eta}{C_p p_X} \cdot \Pr[S = R] \right).$$

Plugging this into (17) yields the pointwise bound

$$\begin{aligned} \Pr[S \setminus \{i\} = R \mid S \sim \text{Acc}_X] &\leq \left( 1 + \frac{C_p t}{1-t} \right) \Pr[S = R \mid S \sim \text{Acc}_X] \\ &\quad + \frac{t}{1-t} \cdot \frac{\eta}{p_X} \Pr[S = R]. \end{aligned} \quad (18)$$

Since  $i \notin R$ , we have  $X_R = X'_R$ , and therefore the conditional acceptance behavior is identical under  $X$  and  $X'$  given the event  $S = R$ :

$$\Pr[\text{Acc}_X \mid S = R] = \Pr[\text{Acc}_{X'} \mid S = R].$$

Hence,

$$\Pr[S = R \mid S \sim \text{Acc}_X] = \frac{\Pr[S = R] \Pr[\text{Acc}_X \mid S = R]}{p_X} = \frac{p_{X'}}{p_X} \Pr[S = R \mid S \sim \text{Acc}_{X'}]. \quad (19)$$

By Corollary 16, we have  $p_X \geq (1 - 3C_p t)p_{X'}$ , i.e.

$$\frac{p_{X'}}{p_X} \leq \frac{1}{1 - 3C_p t}. \quad (20)$$

Moreover, for every  $R$ ,

$$\Pr[S = R \mid S \sim \text{Acc}_{X'}] \leq \Pr[S \setminus \{i\} = R \mid S \sim \text{Acc}_{X'}], \quad (21)$$

since the event  $\{S = R\}$  implies  $\{S \setminus \{i\} = R\}$ .

Combining (18), (19), (20), and (21), we obtain

$$\begin{aligned} \Pr[S \setminus \{i\} = R \mid S \sim \text{Acc}_X] &\leq \frac{1 + \frac{C_p t}{1-t}}{1 - 3C_p t} \cdot \Pr[S \setminus \{i\} = R \mid S \sim \text{Acc}_{X'}] \\ &\quad + \frac{t}{1-t} \cdot \frac{\eta}{p_X} \Pr[S = R]. \end{aligned} \quad (22)$$

Let  $\mathcal{E} \subseteq 2^{[n] \setminus \{i\}}$  be any event. Summing (22) over  $R \in \mathcal{E}$  gives

$$\begin{aligned} \Pr[S \setminus \{i\} \in \mathcal{E} \mid S \sim \text{Acc}_X] &\leq \frac{1 + \frac{C_p t}{1-t}}{1 - 3C_p t} \cdot \Pr[S \setminus \{i\} \in \mathcal{E} \mid S \sim \text{Acc}_{X'}] \\ &\quad + \frac{t}{1-t} \cdot \frac{\eta}{p_X} \sum_{R \in \mathcal{E}} \Pr[S = R] \\ &\leq \exp(\varepsilon_S) \cdot \Pr[S \setminus \{i\} \in \mathcal{E} \mid S \sim \text{Acc}_{X'}] \\ &\quad + \frac{t}{1-t} \cdot \frac{\eta}{p_X}, \end{aligned}$$

where  $\varepsilon_S$  is as in (14) and we used  $\sum_{R \in \mathcal{E}} \Pr[S = R] \leq 1$ .

The reverse inequality (with  $X$  and  $X'$  swapped) is identical, except the additive term becomes  $\frac{t}{1-t} \cdot \frac{\eta}{p_{X'}}$ ; using Corollary 16 again, we have  $p_{X'} \geq p_X / (1 + 3C_p t)$  and hence

$$\frac{t}{1-t} \cdot \frac{\eta}{p_{X'}} \leq \frac{t}{1-t} (1 + 3C_p t) \cdot \frac{\eta}{p_X} = \delta_S.$$

Thus the two conditional laws are  $(\varepsilon_S, \delta_S)$ -close, proving the lemma.

Finally, to see the advertised scaling, note that for  $t \in (0, \frac{1}{5(1+C_p)})$  we have  $t \leq 1/5$  and  $3C_p t \leq 3/5$ , so

$$\varepsilon_S = \log \left( 1 + \frac{C_p t}{1-t} \right) + \log \left( \frac{1}{1 - 3C_p t} \right) \leq \frac{C_p t}{1-t} + \frac{3C_p t}{1 - 3C_p t} \leq \frac{5}{4} C_p t + \frac{15}{2} C_p t \leq 9C_p t. \quad \blacksquare$$

## B.2.2. PROOF OF LEMMA 32

**Proof** [Proof of Lemma 32] Fix an arbitrary measurable event  $\mathcal{A} \subseteq \mathbb{R}^d$  (where  $d$  is the ambient dimension of the output). For each  $T \subseteq [n] \setminus \{i\}$ , abbreviate

$$\Sigma_T := E(X_T) \quad \text{and} \quad \Sigma_{T \cup \{i\}} := E(X_{T \cup \{i\}}),$$

and let

$$P_T(\mathcal{A}) := \Pr_{z \sim \mathcal{N}(0, \Sigma_T)} [z \in \mathcal{A}] \quad \text{and} \quad Q_T(\mathcal{A}) := \Pr_{z \sim \mathcal{N}(0, \Sigma_{T \cup \{i\}})} [z \in \mathcal{A}].$$

**Comparing the two Gaussians for a fixed  $T$ .** Let  $p_T(\cdot)$  and  $q_T(\cdot)$  denote the PDFs of  $\mathcal{N}(0, \Sigma_T)$  and  $\mathcal{N}(0, \Sigma_{T \cup \{i\}})$ , respectively. Since  $\Sigma_T \succeq \Sigma_{T \cup \{i\}}$ , we have  $\Sigma_T^{-1} \preceq \Sigma_{T \cup \{i\}}^{-1}$ , and hence for every  $v \in \mathbb{R}^d$ ,

$$\frac{q_T(v)}{p_T(v)} = \sqrt{\frac{\det(\Sigma_T)}{\det(\Sigma_{T \cup \{i\}})}} \cdot \exp\left(-\frac{1}{2}v^\top (\Sigma_{T \cup \{i\}}^{-1} - \Sigma_T^{-1})v\right) \leq \sqrt{\frac{\det(\Sigma_T)}{\det(\Sigma_{T \cup \{i\}})}} \leq \sqrt{C_E}.$$

Integrating over  $\mathcal{A}$  gives the uniform comparison

$$Q_T(\mathcal{A}) \leq \sqrt{C_E} P_T(\mathcal{A}). \quad (23)$$

**Expanding the two mixture distributions.** Write  $S \sim \text{Acc}_X$  for Poisson subsampling (and test randomness) conditioned on  $\text{Acc}_X$ . Grouping by  $T := S \setminus \{i\}$ , we may write

$$\Pr_{S \sim \text{Acc}_X} [\mathcal{N}(0, E(X_S)) \in \mathcal{A}] = \sum_{T \subseteq [n] \setminus \{i\}} \Pr_{S \sim \text{Acc}_X} [S = T] \cdot P_T(\mathcal{A}) \quad (24)$$

$$+ \Pr_{S \sim \text{Acc}_X} [S = T \cup \{i\}] \cdot Q_T(\mathcal{A}), \quad (25)$$

$$\Pr_{S \sim \text{Acc}_X} [\mathcal{N}(0, E(X_{S \setminus \{i\}})) \in \mathcal{A}] = \sum_{T \subseteq [n] \setminus \{i\}} \Pr_{S \sim \text{Acc}_X} [S = T] \cdot P_T(\mathcal{A}) \quad (26)$$

$$+ \Pr_{S \sim \text{Acc}_X} [S = T \cup \{i\}] \cdot P_T(\mathcal{A}). \quad (27)$$

**Applying Theorem 14** to the set  $T \cup \{i\}$  yields, for every  $T \subseteq [n] \setminus \{i\}$ ,

$$\Pr_{S \sim \text{Acc}_X} [S = T \cup \{i\}] \leq \frac{C_p t}{1-t} \cdot \Pr_{S \sim \text{Acc}_X} [S = T] + \frac{t}{1-t} \cdot \frac{\eta}{p_X} \cdot \Pr_{S \sim \text{Acc}_X} [S = T]. \quad (28)$$

Define  $\rho := \frac{C_p t}{1-t}$  and  $\xi := \frac{t}{1-t} \cdot \frac{\eta}{p_X}$  for brevity.

**The first DP inequality.** Using (23) and then (28) in (24) gives

$$\begin{aligned} \Pr_{S \sim \text{Acc}_X} [\mathcal{N}(0, E(X_S)) \in \mathcal{A}] &\leq \sum_T \Pr_{S \sim \text{Acc}_X} [S = T] \cdot P_T(\mathcal{A}) \\ &\quad + \Pr_{S \sim \text{Acc}_X} [S = T \cup \{i\}] \cdot \sqrt{C_E} P_T(\mathcal{A}) \\ &\leq \sum_T \left(1 + \sqrt{C_E} \rho\right) \Pr_{S \sim \text{Acc}_X} [S = T] \cdot P_T(\mathcal{A}) \\ &\quad + \sqrt{C_E} \xi \sum_T \Pr_{S \sim \text{Acc}_X} [S = T] \cdot P_T(\mathcal{A}). \end{aligned}$$

Since  $\sum_T \Pr_{S \sim \text{Acc}_X} [S = T] \cdot P_T(\mathcal{A}) \leq \sum_T \Pr_{S \sim \text{Acc}_X} [S = T] \leq 1$ , and since

$$\begin{aligned} \sum_T \Pr_{S \sim \text{Acc}_X} [S = T] \cdot P_T(\mathcal{A}) &\leq \sum_T \left( \Pr_{S \sim \text{Acc}_X} [S = T] + \Pr_{S \sim \text{Acc}_X} [S = T \cup \{i\}] \right) \cdot P_T(\mathcal{A}) \\ &= \Pr_{S \sim \text{Acc}_X} [\mathcal{N}(0, E(X_{S \setminus \{i\}})) \in \mathcal{A}], \end{aligned}$$

we conclude

$$\Pr_{S \sim \text{Acc}_X} [\mathcal{N}(0, E(X_S)) \in \mathcal{A}] \leq \left(1 + \sqrt{C_E \rho}\right) \cdot \Pr_{S \sim \text{Acc}_X} [\mathcal{N}(0, E(X_{S \setminus \{i\}})) \in \mathcal{A}] + \sqrt{C_E} \xi.$$

This is exactly the first  $(\varepsilon_E, \delta_E)$  inequality with parameters (15).

**The second DP inequality.** Starting from (26) and applying (28) directly (without needing (23)), we get

$$\begin{aligned} \Pr_{S \sim \text{Acc}_X} [\mathcal{N}(0, E(X_{S \setminus \{i\}})) \in \mathcal{A}] &= \sum_T \Pr_{S \sim \text{Acc}_X} [S = T] \cdot P_T(\mathcal{A}) \\ &\quad + \Pr_{S \sim \text{Acc}_X} [S = T \cup \{i\}] \cdot P_T(\mathcal{A}) \\ &\leq \sum_T (1 + \rho) \Pr_{S \sim \text{Acc}_X} [S = T] \cdot P_T(\mathcal{A}) \\ &\quad + \xi \sum_T \Pr_{S \sim \text{Acc}_X} [S = T] \cdot P_T(\mathcal{A}) \\ &\leq (1 + \rho) \cdot \sum_T \Pr_{S \sim \text{Acc}_X} [S = T] \cdot P_T(\mathcal{A}) + \xi. \end{aligned}$$

Finally, by (24) we have

$$\sum_T \Pr_{S \sim \text{Acc}_X} [S = T] \cdot P_T(\mathcal{A}) \leq \Pr_{S \sim \text{Acc}_X} [\mathcal{N}(0, E(X_S)) \in \mathcal{A}],$$

and since  $1 + \rho \leq 1 + \sqrt{C_E \rho}$  and  $\xi \leq \sqrt{C_E} \xi$ , this yields

$$\Pr_{S \sim \text{Acc}_X} [\mathcal{N}(0, E(X_{S \setminus \{i\}})) \in \mathcal{A}] \leq \left(1 + \sqrt{C_E \rho}\right) \cdot \Pr_{S \sim \text{Acc}_X} [\mathcal{N}(0, E(X_S)) \in \mathcal{A}] + \sqrt{C_E} \xi,$$

which is the second  $(\varepsilon_E, \delta_E)$  inequality with the same parameters (15). This completes the proof.  $\blacksquare$

### B.2.3. PROOF OF LEMMA 33

**Proof** [Proof of Lemma 33] Fix any measurable event  $\mathcal{A} \subseteq \mathbb{R}^d$ .

**Warm-up: analysis of the standard Gaussian mechanism with  $\varepsilon = 1$ .** Fix any deterministic vector  $v \in \mathbb{R}^d$  satisfying  $\Sigma \succeq 2 \log\left(\frac{4}{\delta}\right) vv^\top$ . Let  $p_0$  and  $p_v$  denote the PDFs of  $\mathcal{N}(0, \Sigma)$  and  $\mathcal{N}(v, \Sigma)$  (with respect to Lebesgue measure on the support subspace; equivalently one may interpret  $\Sigma^{-1}$  below as the Moore–Penrose pseudoinverse  $\Sigma^\dagger$ ). A direct calculation gives the privacy-loss random variable

$$L(u) := \log\left(\frac{p_0(u)}{p_v(u)}\right) = \frac{1}{2} v^\top \Sigma^\dagger v - u^\top \Sigma^\dagger v.$$

Let  $\alpha := v^\top \Sigma^\dagger v$ . The matrix domination assumption implies (e.g. by the rank-one PSD ordering characterization) that

$$\alpha \leq \frac{1}{2 \log\left(\frac{4}{\delta}\right)}. \quad (29)$$

Now, if  $u \sim \mathcal{N}(0, \Sigma)$  then  $u^\top \Sigma^\dagger v$  is a one-dimensional Gaussian with mean 0 and variance  $\alpha$ , and therefore  $L(u) \sim \mathcal{N}(\alpha/2, \alpha)$ . Hence, using the standard Gaussian tail bound  $\Pr_{Z \sim \mathcal{N}(0,1)} [Z \geq t] \leq e^{-t^2/2}$ ,

$$\begin{aligned} \Pr_{u \sim \mathcal{N}(0, \Sigma)} [L(u) > 1] &= \Pr_{Z \sim \mathcal{N}(0,1)} \left[ Z > \frac{1 - \alpha/2}{\sqrt{\alpha}} \right] \\ &\leq \exp\left(-\frac{(1 - \alpha/2)^2}{2\alpha}\right) \\ &= \exp\left(-\left(\frac{1}{2\alpha} - \frac{1}{2} + \frac{\alpha}{8}\right)\right) \\ &\leq \exp\left(-\left(\log\left(\frac{4}{\delta}\right) - \frac{1}{2}\right)\right) \\ &= e^{1/2} \cdot \frac{\delta}{4} \leq \frac{\delta}{2}, \end{aligned}$$

where in the last line we used (29) to lower bound  $\frac{1}{2\alpha} \geq \log(4/\delta)$ .

Likewise, if  $u \sim \mathcal{N}(v, \Sigma)$  then  $u = v + g$  with  $g \sim \mathcal{N}(0, \Sigma)$ , so  $u^\top \Sigma^\dagger v = \alpha + g^\top \Sigma^\dagger v$  and thus  $L(u) \sim \mathcal{N}(-\alpha/2, \alpha)$ . By symmetry, the same tail bound yields

$$\Pr_{u \sim \mathcal{N}(v, \Sigma)} [L(u) < -1] \leq \frac{\delta}{2}.$$

Define the good set  $\mathcal{G} := \{u : |L(u)| \leq 1\}$ . On  $\mathcal{G}$  we have both  $p_0(u) \leq e \cdot p_v(u)$  and  $p_v(u) \leq e \cdot p_0(u)$ . Therefore, for every event  $\mathcal{A}$ ,

$$\Pr_{u \sim \mathcal{N}(0, \Sigma)} [u \in \mathcal{A}] \leq e \Pr_{u \sim \mathcal{N}(v, \Sigma)} [u \in \mathcal{A}] + \Pr_{u \sim \mathcal{N}(0, \Sigma)} [u \notin \mathcal{G}] \leq e \Pr_{u \sim \mathcal{N}(v, \Sigma)} [u \in \mathcal{A}] + \frac{\delta}{2},$$

and similarly

$$\Pr_{u \sim \mathcal{N}(v, \Sigma)} [u \in \mathcal{A}] \leq e \Pr_{u \sim \mathcal{N}(0, \Sigma)} [u \in \mathcal{A}] + \frac{\delta}{2}.$$

Thus, for every fixed  $v$  satisfying the matrix domination assumption,

$$\mathcal{N}(0, \Sigma) \approx_{1, \delta/2} \mathcal{N}(v, \Sigma). \quad (30)$$

**Mixing over a spike at  $v = 0$ .** Let  $P := \mathcal{N}(0, \Sigma)$ . Write  $\mathcal{V}_{\neq 0}$  for the conditional law of  $v$  given  $v \neq 0$ , and define the conditional mixture

$$Q := \mathbb{E}_{v \sim \mathcal{V}_{\neq 0}} [\mathcal{N}(v, \Sigma)].$$

By (30) and linearity of expectation, we have  $Q \approx_{1, \delta/2} P$ . Moreover, using  $\Pr[v = 0] \geq 1 - \tau$ , the full mixture distribution is

$$\mathcal{N}(v, \Sigma) = (1 - \tau)P + \tau Q.$$

We now compare  $P$  and  $(1 - \tau)P + \tau Q$ . First,

$$\begin{aligned} \Pr_{\mathcal{N}(v, \Sigma)} [u \in \mathcal{A}] &= (1 - \tau)\Pr_P [u \in \mathcal{A}] + \tau\Pr_Q [u \in \mathcal{A}] \\ &\leq (1 - \tau)\Pr_P [u \in \mathcal{A}] + \tau \left( e\Pr_P [u \in \mathcal{A}] + \delta/2 \right) \\ &= ((1 - \tau) + \tau e)\Pr_P [u \in \mathcal{A}] + \tau \cdot \frac{\delta}{2} = e^\varepsilon \Pr_P [u \in \mathcal{A}] + \delta', \end{aligned}$$

with  $\varepsilon, \delta'$  as in (16).

For the reverse inequality, using  $P \approx_{1, \delta/2} Q$  we have  $\Pr_P [u \in \mathcal{A}] \leq e\Pr_Q [u \in \mathcal{A}] + \delta/2$ . Substituting  $\Pr_Q [u \in \mathcal{A}] = \frac{1}{\tau} \left( \Pr_{\mathcal{N}(v, \Sigma)} [u \in \mathcal{A}] - (1 - \tau)\Pr_P [u \in \mathcal{A}] \right)$  and rearranging gives

$$\Pr_P [u \in \mathcal{A}] \leq \frac{e}{\tau + e(1 - \tau)} \Pr_{\mathcal{N}(v, \Sigma)} [u \in \mathcal{A}] + \frac{\tau}{\tau + e(1 - \tau)} \cdot \frac{\delta}{2}.$$

Since  $\frac{e}{\tau + e(1 - \tau)} = \frac{1}{(1 - \tau) + \tau/e} \leq (1 - \tau) + \tau e = e^\varepsilon$  and  $\frac{\tau}{\tau + e(1 - \tau)} \cdot \frac{\delta}{2} \leq \tau \cdot \frac{\delta}{2} = \delta'$ , we obtain

$$\Pr_P [u \in \mathcal{A}] \leq e^\varepsilon \Pr_{\mathcal{N}(v, \Sigma)} [u \in \mathcal{A}] + \delta'.$$

Thus  $\mathcal{N}(0, \Sigma) \approx_{\varepsilon, \delta'} \mathcal{N}(v, \Sigma)$ , completing the proof. ■

#### B.2.4. CONCLUDING LEMMA 28

**Proof** [Proof of Lemma 28] Let  $i$  be the unique coordinate on which  $X$  and  $X'$  differ, and set

$$\sigma^2 \stackrel{\text{def}}{=} 2 \log \left( \frac{5}{\delta} \right), \quad \alpha \stackrel{\text{def}}{=} \frac{C_p t}{1 - t}.$$

Write

$$\tau_X \stackrel{\text{def}}{=} \Pr[i \in S \mid \text{Acc}_X], \quad \tau_{X'} \stackrel{\text{def}}{=} \Pr[i \in S \mid \text{Acc}_{X'}].$$

Summing Theorem 14 over all  $T \ni i$  gives

$$\tau_X \leq \alpha(1 - \tau_X) + t \cdot \frac{\eta}{p_X} \quad \implies \quad \tau_X \leq \frac{\alpha + t \cdot \eta/p_X}{1 + \alpha}.$$

The same argument yields

$$\tau_{X'} \leq \frac{\alpha + t \cdot \eta/p_{X'}}{1 + \alpha}.$$

Since  $\delta < 1/5$ , the assumption  $p_X \geq 2\eta/\delta$  implies  $p_X \geq 10\eta$ ; hence Corollary 16 applies and gives  $p_{X'} \geq p_X/(1 + 3C_p t)$ , so

$$\frac{\eta}{p_{X'}} \leq (1 + 3C_p t) \cdot \frac{\eta}{p_X}.$$

Therefore,

$$\tau_X, \tau_{X'} \leq \tau \stackrel{\text{def}}{=} \frac{\alpha + t(1 + 3C_p t) \cdot \eta/p_X}{1 + \alpha}. \quad (31)$$

**Applying Lemma 33 to remove  $i$  from the mean.** Define the spike

$$v \stackrel{\text{def}}{=} A(X_S) - A(X_{S \setminus \{i\}}).$$

By definition,  $v = 0$  whenever  $i \notin S$ , and from the definition of  $E$ ,

$$E(X_S) \succeq vv^\top \quad \text{almost surely.}$$

Condition on  $S$  (so  $\Sigma \stackrel{\text{def}}{=} \sigma^2 E(X_S)$  is fixed). Since  $\sigma^2 = 2 \log(5/\delta) \geq 2 \log(4/\delta)$ , we have

$$\Sigma \succeq 2 \log\left(\frac{4}{\delta}\right) vv^\top.$$

Applying Lemma 33 with  $\tau$  from (31) and then shifting by the common mean  $A(X_{S \setminus \{i\}})$  gives

$$\left( \mathcal{N}(A(X_S), \sigma^2 E(X_S)) \mid S \sim \text{Acc}_X \right) \approx_{\varepsilon_G, \tau\delta/2} \left( \mathcal{N}(A(X_{S \setminus \{i\}}), \sigma^2 E(X_S)) \mid S \sim \text{Acc}_X \right),$$

where  $\varepsilon_G \stackrel{\text{def}}{=} \log(1 + \tau(e - 1))$ .

**Applying Lemma 32 to remove  $i$  from the covariance.** By translation invariance and post-processing, Lemma 32 applies with the same  $\varepsilon_E$  and  $\delta_E$  to Gaussians whose mean is  $A(X_{S \setminus \{i\}})$ , and also with both covariances scaled by the common factor  $\sigma^2$ . Hence

$$\left( \mathcal{N}(A(X_{S \setminus \{i\}}), \sigma^2 E(X_S)) \mid S \sim \text{Acc}_X \right) \approx_{\varepsilon_E, \delta_E} \left( \mathcal{N}(A(X_{S \setminus \{i\}}), \sigma^2 E(X_{S \setminus \{i\}})) \mid S \sim \text{Acc}_X \right),$$

where  $\varepsilon_E$  is as in (15) and

$$\delta_E = \sqrt{C_E} \cdot \frac{t}{1-t} \cdot \frac{\eta}{p_X}.$$

**Applying Lemma 31 to replace  $S \sim \text{Acc}_X$  with  $S \sim \text{Acc}_{X'}$ .** Since  $X$  and  $X'$  differ only at coordinate  $i$ , for every set  $U \subseteq [n]$  we have

$$X_U = X'_U \quad \text{whenever } i \notin U.$$

In particular,  $X_{S \setminus \{i\}} = X'_{S \setminus \{i\}}$  and hence also  $A(X_{S \setminus \{i\}}) = A(X'_{S \setminus \{i\}})$  and  $E(X_{S \setminus \{i\}}) = E(X'_{S \setminus \{i\}})$ . Thus, Lemma 31 and post-processing imply

$$\left( \mathcal{N}(A(X_{S \setminus \{i\}}), \sigma^2 E(X_{S \setminus \{i\}})) \mid S \sim \text{Acc}_X \right) \approx_{\varepsilon_S, \delta_S} \left( \mathcal{N}(A(X'_{S \setminus \{i\}}), \sigma^2 E(X'_{S \setminus \{i\}})) \mid S \sim \text{Acc}_{X'} \right),$$

where  $\varepsilon_S$  and  $\delta_S$  are as in (14).

**Re-introducing  $i$  into the covariance using Lemma 32 (again).** Applying Lemma 32 to  $X'$  (and using  $\eta/p_{X'} \leq (1 + 3C_p t)\eta/p_X$  from above) yields

$$\left(\mathcal{N}\left(A(X'_{S \setminus \{i\}}), \sigma^2 E(X'_{S \setminus \{i\}})\right) \mid S \sim \text{Acc}_{X'}\right) \approx_{\varepsilon_E, \delta'_E} \left(\mathcal{N}\left(A(X'_{S \setminus \{i\}}), \sigma^2 E(X'_S)\right) \mid S \sim \text{Acc}_{X'}\right),$$

where

$$\delta'_E = \sqrt{C_E} \cdot \frac{t}{1-t} \cdot \frac{\eta}{p_{X'}} \leq \sqrt{C_E} \cdot \frac{t}{1-t} (1 + 3C_p t) \cdot \frac{\eta}{p_X}.$$

**Re-introducing  $i$  into the mean using Lemma 33 (again).** Exactly as in Step 1 (now for  $X'$ ), we obtain

$$\left(\mathcal{N}\left(A(X'_{S \setminus \{i\}}), \sigma^2 E(X'_S)\right) \mid S \sim \text{Acc}_{X'}\right) \approx_{\varepsilon_G, \tau\delta/2} \left(\mathcal{N}\left(A(X'_S), \sigma^2 E(X'_S)\right) \mid S \sim \text{Acc}_{X'}\right).$$

**Composition and the final  $\delta$  bookkeeping.** Composing the five steps, the total  $\varepsilon$  is

$$\varepsilon_{\text{draw}} = \varepsilon_S + 2\varepsilon_E + 2\varepsilon_G = O(t),$$

as claimed, while the total  $\delta$  is at most

$$\delta_{\text{tot}} \leq \delta_S + \delta_E + \delta'_E + \tau\delta.$$

Using  $\eta/p_X \leq \delta/2$ ,  $t \leq 1/5$ , and  $\alpha \leq 1/4$  (from  $t \leq 1/(5(C_p + 1))$ ), one checks that  $\delta_S \leq \delta/5$ ,  $\delta_E \leq \delta/8$ ,  $\delta'_E \leq \delta/5$ , and  $\tau \leq 3/10$ , hence  $\delta_{\text{tot}} \leq \delta$ . This completes the proof.  $\blacksquare$

## Appendix C. STAR: Subsample-Test-Aggregate for Linear Regression

In this section, we will instantiate our subsample-test-aggregate framework for linear regressions.

### C.1. The STAR Algorithm

**Algorithm 6:** SUBSAMPLE-TEST-AGGREGATE REGRESSION (STAR)

**Input:** Data  $(X, y) \in \mathbb{R}^{n \times d} \times \mathbb{R}^n$

**Output:** Private linear model  $\hat{\beta} \in \mathbb{R}^d$

Define the Poisson-subsample process  $i \in S \stackrel{i.i.d.}{\sim} \text{Bern}(t)$  with  $t = \Theta\left(\frac{d + \sqrt{d \log(n) \log(1/\delta)}}{n}\right)$ .

Set  $\mathcal{T}$  to be the SOFT-ACRE stability test with parameters as defined in Section 1.2.4.

Use the private  $p$ -check algorithm given in Theorem 17 to privately check if  $\Pr[\mathcal{T}(S)] \ll 1$ .

**if**  $\Pr[\mathcal{T}(S)] \ll 1$  **then**

**return** REJECT

**end**

Draw  $k = \frac{\log(n) \log(1/\delta)}{\alpha^2}$  fresh Poisson-subsamples  $S_j$ , conditioned on  $\mathcal{T}$  accepting.

For each  $S_j$ , draw a sample

$$v_j \leftarrow \mathcal{N}\left(\text{OLS}(X_{S_j}, y_{S_j}), \frac{16R^2L}{(1-2L)^2} \log(5/\delta) (X_{S_j}^\top X_{S_j})^{-1}\right).$$

**return** the empirical mean  $\frac{1}{k} \sum_{j=1}^k v_j$

## C.2. Utility Analysis

In this section we assemble our utility analysis for STAR.

### C.2.1. ACCURACY BOUNDS FOR SUBSAMPLED OLS

For the duration of this section, let  $(X_1, y_1), \dots, (X_n, y_n)$  be iid draws from the Gaussian linear model where  $X \sim \mathcal{N}(0, I)$  and  $y = \langle X, \beta_0 \rangle + \epsilon$ , where  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  and  $\beta_0 \in \mathbb{R}^d$ . Let  $t \in [0, 1]$  be a subsampling rate. For  $S \subseteq [n]$ , let  $\hat{\beta}_S = (X_S^\top X_S)^{-1} X_S^\top y_S$  be the (minimum norm) OLS regressor using only the samples contained in  $S$ . Suppose we draw random subsets  $S_1, \dots, S_k$ , where each subset is defined by including each sample  $i$  independently with probability  $t$ , and let  $\bar{\beta} = \frac{1}{k} \sum_{i \leq k} \hat{\beta}_{S_i}$ . Our main result in this section is:

**Theorem 34** *Let  $E$  be event that  $|S_i| \geq tn/4$  for all  $i \leq k$ , and assume  $tn \geq Cd$  for a big-enough constant  $C$ . Then*

$$\mathbb{E}[\|\bar{\beta} - \beta_0\|^2 \mid E] \leq O\left(\frac{d}{tnk} + \frac{d}{n}\right) \cdot \sigma^2,$$

where the expectation is taken with respect to  $(X_1, y_1), \dots, (X_n, y_n)$  and  $S_1, \dots, S_k$ .

To prove the theorem we assemble a few lemmas.

**Lemma 35** *Let  $S, T \subseteq [n]$  be two independent Bernoulli( $t$ ) subsamples and let  $I = S \cap T$ . Let  $W_S \in \mathbb{R}^{n \times n}$  be the diagonal indicator matrix for membership in  $S$  and similarly for  $W_T$ . Let  $X$  have rows  $X_1, \dots, X_n$ . Write*

$$M_S := X^\top W_S X = \sum_{i \in S} x_i x_i^\top, \quad M_T := X^\top W_T X = \sum_{i \in T} x_i x_i^\top.$$

Then

$$\mathbb{E}_\epsilon \left[ \left\langle \hat{\beta}_S - \beta_0, \hat{\beta}_T - \beta_0 \right\rangle \mid X, S, T \right] = \sigma^2 \sum_{i \in I} x_i^\top M_S^{-1} M_T^{-1} x_i,$$

where the expectation is taken only with respect to  $\epsilon_1, \dots, \epsilon_n$ .

### Proof

Expanding the formula for OLS,

$$\hat{\beta}_S - \beta_0 = M_S^{-1} \sum_{i \in S} x_i \epsilon_i, \quad \hat{\beta}_T - \beta_0 = M_T^{-1} \sum_{i \in T} x_i \epsilon_i.$$

Now condition on  $(X, S, T)$  and take expectation over the Gaussian noise  $\epsilon$ : only the indices in the overlap  $I := S \cap T$  contribute, because  $\epsilon_i$ 's are independent. A short calculation gives

$$\mathbb{E}_\epsilon \left[ \left\langle \hat{\beta}_S - \beta_0, \hat{\beta}_T - \beta_0 \right\rangle \mid X, S, T \right] = \sigma^2 \sum_{i \in I} x_i^\top M_S^{-1} M_T^{-1} x_i$$

which proves the lemma. ■

**Lemma 36** Let  $M_S = X_S X_S^T$  and  $M_T = X_T X_T^T$  be Wishart matrices, where  $X_S \in \mathbb{R}^{d \times m_1}$  has independent  $\mathcal{N}(0, 1)$  entries and similarly for  $X_T \in \mathbb{R}^{d \times m_2}$ , with  $m_1, m_2 \geq Cd$  for a sufficiently large constant  $C$ . Let  $x \in \mathbb{R}^d$  be  $\mathcal{N}(0, I)$ -distributed. ( $X_S, X_T$ , and  $x$  may not be independent of each other.) Then

$$\mathbb{E} x^\top M_S^{-1} M_T^{-1} x \leq O\left(\frac{d}{m_1 m_2}\right)$$

**Proof** Using Cauchy-Schwarz, we have

$$\mathbb{E} x^\top M_S^{-1} M_T^{-1} x \leq (\mathbb{E} \|x\|^4)^{1/2} (\mathbb{E} \|M_S^{-1}\|^4)^{1/4} (\mathbb{E} \|M_T^{-1}\|^4)^{1/4}$$

By Vershynin (2010)(Corollary 5.35),  $\mathbb{E} \|M_S^{-1}\|^4 \leq O(m_1^{-4})$ . The lemma follows.  $\blacksquare$

Now we can prove Theorem 34.

**Proof** [Proof of Theorem 34] We expand the error of the subsampled estimator, eliding the conditioning on  $E$  for simplicity:

$$\mathbb{E} \|\bar{\beta} - \beta_0\|^2 = \frac{1}{k^2} \sum_{i \leq k} \mathbb{E} \|\hat{\beta}_{S_i} - \beta_0\|^2 + \frac{1}{k^2} \sum_{i \neq j \leq k} \mathbb{E} \langle \hat{\beta}_{S_i} - \beta_0, \hat{\beta}_{S_j} - \beta_0 \rangle$$

In the first term,  $\mathbb{E} \|\hat{\beta}_{S_i} - \beta_0\|^2$  is the squared error of an OLS estimator using  $|S_i|$  samples, hence is at most  $O(d/|S_i|)$ . Since  $|S_i| \geq tn/4$ , this is at most  $O(d/(tn))$ , and so the entire first term is at most  $O(\sigma^2 d/(tnk))$ . For the second term, combining Lemmas 35 and 36, we get  $O(d/(tn)^2) \cdot \sigma^2 \cdot \mathbb{E} |S \cap T|$ , where  $S$  and  $T$  are independent subsamples. This is  $O(d/(tn)^2) \cdot \sigma^2 \cdot t^2 n = O(\sigma^2 d/n)$ .  $\blacksquare$

### C.2.2. ANALYSIS OF ACRE TEST ON GAUSSIAN DATA

Our goal in this section is to prove Theorem 37. For this section, let  $(x_1, y_1), \dots, (x_n, y_n)$  be draws from a linear Gaussian model with  $n \geq Cd$  for a big-enough constant  $C$ .

**Theorem 37 (ACRE passes with high probability on well-behaved data)** *If  $X, y$  is a regression drawn according to Definition 3, and  $S \subseteq [n]$  is drawn from a Poisson- $t$  subsampling procedure with  $tn = C(d + \sqrt{d \log(n)} \log(1/\delta))$  for a sufficiently large constant  $C$ , then  $X_S, y_S$  is  $(L, \ell, R)$ -ACRE (for  $L, \ell, R$  as defined in Section 1.2.4) with probability  $1 - O(n^{-100}) - \exp(-\Omega(tn))$ .*

Theorem 37 follows from the following lemmas:

**Conditioning on the subsample size.** Let  $S \subseteq [n]$  be drawn by Poisson- $t$  subsampling and write  $m := |S|$ . (Equivalently,  $m \sim \text{Binomial}(n, t)$  in the usual Poisson-subsampling sense.) A standard Chernoff bound gives

$$\Pr[m \notin [\frac{1}{2}tn, 2tn]] \leq 2 \exp(-\Omega(tn)). \quad (32)$$

In the rest of this subsection we write  $X := X_S$  and  $y := y_S$ , and we work on the event  $m \in [\frac{1}{2}tn, 2tn]$ .

**Lemma 38** *With probability at least  $1 - n^{-100} - \exp(-\Omega(tn))$ , every leverage score satisfies*

$$\max_{i \in S} x_i^\top (X^\top X)^{-1} x_i \leq \frac{C(d + \sqrt{d \log n} + \log n)}{tn},$$

for a sufficiently large universal constant  $C$ .

**Proof** Condition on  $m = |S|$  and define the whitened design

$$Z := X \Sigma^{-1/2} \in \mathbb{R}^{m \times d}.$$

The rows of  $Z$  are i.i.d.  $N(0, I)$ . By Vershynin (2010)[Corollary 5.35], for  $m \geq C_0 d$ ,

$$\|Z^\top Z - mI\| \leq \frac{1}{2}m$$

with probability at least  $1 - 2 \exp(-\Omega(m))$ . Equivalently,

$$\frac{1}{2}mI \preceq Z^\top Z \preceq \frac{3}{2}mI.$$

Multiplying on both sides by  $\Sigma^{1/2}$  gives

$$\frac{1}{2}m \Sigma \preceq X^\top X = \Sigma^{1/2} Z^\top Z \Sigma^{1/2} \preceq \frac{3}{2}m \Sigma,$$

which implies the displayed constant-factor sandwich (after loosening  $\frac{3}{2}$  to 2).

On this event,

$$(X^\top X)^{-1} \preceq \frac{2}{m} \Sigma^{-1}.$$

Hence for any  $i \in S$ ,

$$x_i^\top (X^\top X)^{-1} x_i \leq \frac{2}{m} x_i^\top \Sigma^{-1} x_i = \frac{2}{m} \|\Sigma^{-1/2} x_i\|^2.$$

Since  $\Sigma^{-1/2} x_i \sim N(0, I)$ ,  $\|\Sigma^{-1/2} x_i\|^2 \sim \chi_d^2$ . By the Laurent–Massart tail bound, for any  $u > 0$ ,

$$\Pr[\chi_d^2 \geq d + 2\sqrt{du} + 2u] \leq e^{-u}.$$

Taking  $u = 110 \log n$  and union bounding over  $i \in [n]$  gives, with probability at least  $1 - n^{-100}$ ,

$$\max_{i \in [n]} x_i^\top \Sigma^{-1} x_i \leq d + 2\sqrt{110 d \log n} + 220 \log n.$$

Finally, on the size event  $m \geq \frac{1}{2}tn$  from (32), we obtain

$$\max_{i \in S} x_i^\top (X^\top X)^{-1} x_i \leq \frac{4}{tn} \left( d + 2\sqrt{110 d \log n} + 220 \log n \right) = O\left( \frac{d + \sqrt{d \log n} + \log n}{tn} \right),$$

as claimed. Combining failure probabilities yields  $n^{-100} + \exp(-\Omega(tn))$ .  $\blacksquare$

**Lemma 39** *With probability at least  $1 - n^{-100} - \exp(-\Omega(tn))$ , every cross-leverage satisfies*

$$\max_{\substack{i,j \in S \\ i \neq j}} \left| x_i^\top (X^\top X)^{-1} x_j \right| \leq \frac{C(\sqrt{d \log n} + \log n)}{tn},$$

for a sufficiently large universal constant  $C$ .

**Proof** Fix distinct  $i, j \in S$  and let  $M_{-i} := X_{-i}^\top X_{-i}$  denote the Gram matrix with row  $i$  removed. By Sherman–Morrison for adding  $x_i x_i^\top$ ,

$$(X^\top X)^{-1} = M_{-i}^{-1} - \frac{M_{-i}^{-1} x_i x_i^\top M_{-i}^{-1}}{1 + x_i^\top M_{-i}^{-1} x_i},$$

hence

$$x_i^\top (X^\top X)^{-1} x_j = \frac{x_i^\top M_{-i}^{-1} x_j}{1 + x_i^\top M_{-i}^{-1} x_i}, \quad \text{so} \quad \left| x_i^\top (X^\top X)^{-1} x_j \right| \leq \left| x_i^\top M_{-i}^{-1} x_j \right|.$$

Now condition on  $X_{-i}$  (equivalently on  $M_{-i}$ ) and on  $x_j$ . Let  $v := M_{-i}^{-1} x_j$ . Since  $x_i$  is independent of  $(X_{-i}, x_j)$  and  $x_i \sim N(0, \Sigma)$ , we have

$$x_i^\top v \mid (X_{-i}, x_j) \sim N\left(0, v^\top \Sigma v\right).$$

We next upper bound  $v^\top \Sigma v$ . For each fixed  $i$ ,  $M_{-i}$  is the sum of  $(m-1)$  independent  $N(0, \Sigma)$  outer products. By [Vershynin \(2010\)](#)[Corollary 5.35] (applied to  $(m-1)$  samples),

$$\Pr[M_{-i} \not\geq \frac{1}{2}(m-1)\Sigma] \leq 2 \exp(-\Omega(m)).$$

Union bounding over  $i \in S$  (at most  $n$  choices) keeps this at  $\exp(-\Omega(m))$ . On the event  $M_{-i} \succeq \frac{1}{2}(m-1)\Sigma$ , we have

$$\Sigma^{1/2} M_{-i}^{-1} \Sigma^{1/2} \preceq \frac{2}{m-1} I \preceq \frac{4}{m} I,$$

and thus

$$v^\top \Sigma v = \|\Sigma^{1/2} M_{-i}^{-1} x_j\|^2 = \|\Sigma^{1/2} M_{-i}^{-1} \Sigma^{1/2} \Sigma^{-1/2} x_j\|^2 \leq \left(\frac{4}{m}\right)^2 \|\Sigma^{-1/2} x_j\|^2 = \frac{16}{m^2} x_j^\top \Sigma^{-1} x_j.$$

From the same  $\chi^2$  union bound as in [Lemma 38](#), with probability at least  $1 - n^{-100}$  we have simultaneously for all  $j \in [n]$  that

$$x_j^\top \Sigma^{-1} x_j \leq B \quad \text{where} \quad B := d + 2\sqrt{110 d \log n} + 220 \log n.$$

Therefore, on the intersection of these events,

$$\text{Var}\left(x_i^\top M_{-i}^{-1} x_j \mid X_{-i}, x_j\right) \leq \frac{16B}{m^2}.$$

A Gaussian tail bound then gives

$$\Pr \left[ \left| x_i^\top M_{-i}^{-1} x_j \right| \geq \frac{4\sqrt{B}}{m} \sqrt{220 \log n} \mid X_{-i}, x_j \right] \leq 2n^{-110}.$$

Union bounding this tail over all ordered pairs  $(i, j)$  with  $i \neq j$  (at most  $n^2$  pairs) gives overall failure probability  $\leq n^{-100}$  for the cross-term bound.

Finally, on the size event  $m \geq \frac{1}{2}tn$  from (32), the right-hand side is

$$O\left(\frac{\sqrt{B \log n}}{m}\right) = O\left(\frac{\sqrt{d \log n} + \log n}{m}\right) \leq O\left(\frac{\sqrt{d \log n} + \log n}{tn}\right),$$

proving the lemma. Combining failure probabilities yields  $n^{-100} + \exp(-\Omega(tn))$ .  $\blacksquare$

**Lemma 40** *With probability at least  $1 - n^{-100} - \exp(-\Omega(tn))$ , every residual satisfies*

$$\max_{i \in S} |y_i - (Hy)_i| \leq C \sigma \sqrt{\log n},$$

where  $H := X(X^\top X)^{-1}X^\top$  is the hat matrix for  $(X, y) = (X_S, y_S)$ .

**Proof** Work on the event from (32) that  $m \in [\frac{1}{2}tn, 2tn]$ , and on the conditioning event from Lemma 38 that  $X^\top X \succeq \frac{1}{2}m\Sigma$ . This intersection fails with probability  $\exp(-\Omega(tn))$ , and on it  $X^\top X$  is invertible so  $H$  is well-defined.

Write  $y = X\beta^* + \epsilon$  with  $\epsilon \sim \mathcal{N}(0, \sigma^2 I_m)$  independent of  $X$ . Since  $HX = X$ , the residual vector is

$$r := y - Hy = (I - H)\epsilon.$$

Condition on  $X$  so  $H$  is fixed. Then  $r$  is mean-zero Gaussian with covariance  $\sigma^2(I - H)$ . For each coordinate  $k \in [m]$ ,

$$\text{Var}(r_k \mid X) = \sigma^2(I - H)_{kk} = \sigma^2(1 - H_{kk}) \leq \sigma^2,$$

because  $H$  is an orthogonal projector and thus  $0 \leq H_{kk} \leq 1$ . Therefore, for any  $u > 0$ ,

$$\Pr \left[ |r_k| \geq \sigma \sqrt{2u} \mid X \right] \leq 2e^{-u}.$$

Taking  $u = 110 \log n$  and union bounding over  $k \in [m]$  (and using  $m \leq n$ ) gives

$$\Pr \left[ \|r\|_\infty \geq \sigma \sqrt{220 \log n} \mid X \right] \leq 2me^{-110 \log n} \leq 2n \cdot n^{-110} \leq n^{-100}.$$

Combining with the  $\exp(-\Omega(tn))$  failure probability of the size/conditioning event completes the proof.  $\blacksquare$

## C.2.3. PROOFS OF LEMMAS FROM SECTION 1.2.4

**Proof** [Proof of Lemma 10]

Fix  $i \in [n]$  and let  $X_{-i}, y_{-i}$  be the dataset with sample  $i$  removed. The new covariance matrix  $X_{-i}^\top X_{-i} = X^\top X - X_i X_i^\top$  is a rank-1-update of the original covariance matrix. Therefore, updating the hat matrix using the Sherman-Morrison formula, we have

$$\begin{aligned} \forall j, k \neq i \quad (H_{-i})_{j,k} &= X_j^\top (X^\top X - X_i X_i^\top)^{-1} X_k = \\ &= X_j^\top \left( (X^\top X)^{-1} + \frac{1}{1 - H_{i,i}} (X^\top X)^{-1} X_i X_i^\top (X^\top X)^{-1} \right) X_k = \\ &= H_{j,k} + \frac{H_{i,j} H_{i,k}}{1 - H_{i,i}}. \end{aligned}$$

Immediately, we see that the right-hand-side is at most  $\ell + \frac{\ell^2}{1-L}$  for  $j \neq k$  and  $L + \frac{\ell^2}{1-L} \leq L + \frac{\ell L}{1-L}$  for  $j = k$ , yielding the stability of the leverages and cross-leverages.

To see the stability of the residuals, fix  $j \neq i$  and write

$$(H_{-i}y)_j = X_j^\top (X^\top X - X_i X_i^\top)^{-1} (X^\top y - X_i y_i).$$

Applying Sherman–Morrison and rearranging gives the standard leave-one-out identity

$$(H_{-i}y)_j = (Hy)_j - \frac{H_{j,i}}{1 - H_{i,i}} (y_i - (Hy)_i).$$

Therefore,

$$|y_j - (H_{-i}y)_j| \leq |y_j - (Hy)_j| + \frac{|H_{j,i}|}{1 - H_{i,i}} |y_i - (Hy)_i| \leq R + \frac{\ell}{1-L} R = (1 + c)R. \quad \blacksquare$$

**Proof** [Proof of Lemma 11] Lemma 11 follows from Lemma 10.

If  $X, y$  is a regression that has probability  $> \eta$  of being accepted by soft-ACRE, then it must have  $\text{SCORE}(X, y) \leq 2 \left(1 - \frac{2\ell}{1-2L}\right) \leq 2$ , so  $X, y$  must be  $(2L, 2R, 2\ell)$ -ACRE. Therefore, from Lemma 10, for any leave-one-out  $X_{-i}, y_{-i}$ , we know that

$$\text{SCORE}(X_{-i}, y_{-i}) \leq \left(1 + \frac{2\ell}{1-2L}\right) \text{SCORE}(X, y),$$

so it must be accepted by SOFT-ACRE with probability at least  $\frac{1}{2}$  that of  $X, y$ . \blacksquare

### C.3. Proof of Theorem 4

**Proof** [Proof of Theorem 4]

To prove Theorem 4, we need to show that the STAR algorithm is private and efficient and that it has high utility for well-behaved data.

**Privacy** The STAR algorithm utilizes the AC-STA framework. It relies on a down-sensitivity test  $\mathcal{T}$ , which we proved to be  $(O(1), \exp(-\Omega(1/\ell)))$ -down-stable (Lemma 11). We set the sub-sampling rate  $t$  to be of order  $\frac{d + \sqrt{d \log(n) \log(1/\delta)}}{n}$ , and this allows us to set

$$L = \Theta\left(\frac{d}{tn}\right) < \frac{1}{10} \quad \text{and} \quad \ell = \Theta\left(\frac{\sqrt{d \log(n)}}{tn}\right) = O\left(\frac{1}{\log(1/\delta)}\right).$$

We select the constants to ensure that  $\eta = \exp(-\Omega(1/\ell)) < \delta^2$ .

From our assumption that  $n \geq \frac{(1+\lambda)d \log(1/\delta)}{\varepsilon}$ , we have  $t = O\left(\frac{\varepsilon}{\log(1/\delta)}\right)$ . From Theorem 17, this allows us to privately verify that  $\Pr[\mathcal{T}(S) = \text{ACCEPT}] > p_0$ , but because of our regime of  $t$ , we have  $p_0 = \Omega(1)$ , and in particular  $p_0/\delta \gg \delta^2 \geq \eta$ .

From Lemma 12, we know that passing SOFT-ACRE implies bounded down-sensitivity in the geometry of

$$E(X_S) = \frac{8LR^2}{(1-2L)^2} (X_S^\top X_S)^{-1}.$$

By its definition, removing a sample can only cause  $(X_S^\top X_S)^{-1}$  to grow in the PSD sense

$$E(X_S) \preceq E(X_{S \setminus \{i\}}),$$

and because passing SOFT-ACRE with non-zero probability implies that  $X_S, y_S$  is  $(2L, 2R, 2\ell)$ -ACRE, and in particular, that implies that  $X_S$  has bounded leverages  $\leq 2L$ , a standard identity (see e.g., Brown et al. (2023)) gives us

$$(1-2L)|E(X_{S \setminus \{i\}})| \leq |E(X_S)|.$$

Therefore, the assumptions of Lemma 28 hold with  $C_E = 1 - 2L = \Theta(1)$ , implying that the output of the STAR algorithm is  $(O(\sqrt{k \log(1/\delta)t + kt^2}), O(k \frac{\eta}{p_0} + \delta))$ -DP. From our assumption that

$$n \geq \frac{(1+\lambda)d \log(1/\delta) \sqrt{\log(n)}}{\alpha \varepsilon_{\text{target}}}, \quad t = \Theta\left(\frac{d + \sqrt{d \log(n) \log(1/\delta)}}{n}\right), \quad k = \frac{\log(n) \log(1/\delta)}{\alpha^2},$$

we get that the desired privacy of the algorithm

$$\varepsilon = \Theta(\sqrt{k \log(1/\delta)t + kt^2}) = \Theta\left(\frac{\log(1/\delta) \sqrt{\log(n)}}{\alpha} \times \frac{\alpha \varepsilon_{\text{target}}}{\log(1/\delta) \sqrt{\log(n)}}\right) = \Theta(\varepsilon_{\text{target}}).$$

**Utility** The utility of the STAR algorithm for well-behaved data follows directly from a combination of Theorem 34 (which says that the empirical mean over the OLS fits of  $k$  subsamples not conditioned on being ACRE is more than needed to get the desired excess risk  $\alpha$ ), and Theorem 37 (which tells us that SOFT-ACRE will accept all  $O(k/p_0) \ll n^{100}$  tested subsamples with high probability). From a union bound on the two failure events (that even a single ACRE test failed or that the mean over the empirical  $k$  subsamples was far from  $\beta^*$ ).

Finally, note that the STAR algorithm also added Gaussian noise. The covariance of this noise was

$$\text{Covariance} = \frac{\log(1/\delta)}{k^2} \sum_{j=1}^k E(X_{S_j}) = \frac{16LR^2 \log(5/\delta)}{(1-2L)^2 k^2} \sum_{j=1}^k (X_{S_j}^\top X_{S_j})^{-1}.$$

From Lemma 38 and a union bound on the  $k$  subsamples, we know that

$$\text{Covariance} \preceq \Theta \left( \frac{d \log(n) \log(1/\delta) \sigma^2}{(tn)^2} \frac{\alpha^2}{\log(d) \log(1/\delta)} \right) \Sigma^{-1} = O \left( \frac{\alpha^2 \sigma^2}{d} \right) \Sigma^{-1},$$

so with high probability

$$\|\text{Added Noise}\|_{\Sigma} = \Theta(\alpha\sigma).$$

Therefore, by the triangle inequality, with high probability, the Mahalanobis-norm distance between the output of our algorithm  $\beta^{\text{STAR}}$ , and the ground-truth answer  $\beta^*$ , is at most  $O(\alpha\sigma)$ , as desired.

**Time Complexity** Finally, we consider the time complexity of STAR. Because  $p_0 = \Omega(1)$  from our assumption that  $t = O\left(\frac{\varepsilon}{\log(1/\delta)}\right)$ , the private  $p$ -check could be run for the cost of  $\log(1/\delta) \ll n/d$  calls to the SOFT-ACRE oracle, and each such call was dominated by the cost of computing the hat matrix which was in turn dominated by a matrix multiplication between two  $d \times tn = d \times \tilde{O}(d + \lambda d)$  matrices.

Similarly, to get our  $k = \tilde{O}(n/d)$  fresh draws of a  $\mathcal{T}$ -passing subsample  $S_j$ , it suffices to run  $\mathcal{T}$  on  $O(k/p_0) = O(k)$  hypothesis subsamples.  $\blacksquare$

## Appendix D. Basic Mechanisms for Privacy

In this appendix, we will explore some basic primitives we use to ensure the privacy of certain intermediate values throughout the paper.

### D.1. Soft-Truncation Primitives for Enforcing Boundedness

#### D.1.1. DIFFERENTIALLY-PRIVATE SOFT-THRESHOLD TEST

**Lemma 41 (Differentially-Private soft-Threshold Test)** *For any  $\varepsilon, \delta \in (0, 1)$ , there exists a randomized mechanism  $\varphi_{\varepsilon, \delta}$  that takes as input a real number  $z$  and returns ACCEPT or REJECT such that:*

1. **Privacy:** for any  $|z - z'| < 1$ ,

$$\varphi_{\varepsilon, \delta}(z) \approx_{\varepsilon, \delta} \varphi_{\varepsilon, \delta}(z')$$

2. **Correctness:** almost surely:

$$\forall z \leq 0 \varphi_{\varepsilon, \delta}(z) = \text{ACCEPT} \quad \text{and} \quad \forall z \geq \frac{2 \log(2/\delta)}{\varepsilon} \varphi_{\varepsilon, \delta}(z) = \text{REJECT}$$

**Proof** [Proof of Lemma 41] It is a classic result in DP that for a 1-sensitive variable  $z$ , adding truncated Laplace noise suffices to ensure its privacy. In particular, for  $z_{\text{private}} = z + L$ , where  $L$  is drawn from the truncated Laplace law

$$\Pr[L = \lambda] \propto \mathbf{1}_{|\lambda| \leq \tau} \exp(-\varepsilon|\lambda|),$$

is  $\varepsilon, \delta$  approx-DP when  $\tau = \frac{\log(2/\delta)}{\varepsilon}$ .

Therefore, from the post-processing theorem, we know that

$$\varphi_{\varepsilon, \delta}(z) := \mathbb{1}_{z_{\text{private}} \leq \tau}$$

is also private.

Finally, to ensure the correctness of  $\varphi$ , we note that  $|L| < \tau$  almost surely, so if  $z \leq 0$ , then  $z_{\text{private}} \leq \tau$  almost surely, and if  $z \geq 2\tau$ , then  $z_{\text{private}} > \tau$  almost surely, completing our proof. ■

#### D.1.2. SOFT-TRUNCATION KERNELS

We use a soft-threshold acceptance function  $\text{Acc}(\cdot) \in [0, 1]$ . Formally, each  $\text{Acc}(s)$  is the acceptance probability of an underlying randomized accept/reject rule; when writing  $\text{Acc}(s) \approx_{\varepsilon, \delta} \text{Acc}(s')$  we mean that the induced Bernoulli laws on  $\{\text{ACCEPT}, \text{REJECT}\}$  are  $(\varepsilon, \delta)$ -approximately indistinguishable.

**Definition 42 (Power-decay and exponential-decay kernels)** Fix threshold  $\tau > 0$  and floor  $\zeta \in (0, 1)$ .

1. (Power) For exponent  $\alpha > 0$  and  $s > 0$ ,

$$\text{Acc}_{\tau, \alpha, \zeta}^{\text{pow}}(s) := \Pr[\varphi_{\alpha, \zeta}(\log(s/\tau)) = \text{ACCEPT}].$$

2. (Exponential) For slope  $\gamma > 0$  and  $s \geq 0$ ,

$$\text{Acc}_{\tau, \gamma, \zeta}^{\text{exp}}(s) := \Pr[\varphi_{\gamma, \zeta}(s - \tau) = \text{ACCEPT}].$$

**Lemma 43 (Truncation implies global bounds)** Fix  $\tau > 0$  and  $\zeta \in (0, 1)$ .

1. (Power) For any  $\alpha > 0$  and any  $s > 0$ ,

$$s \leq \tau \implies \text{Acc}_{\tau, \alpha, \zeta}^{\text{pow}}(s) = 1, \quad s \geq \tau \exp\left(2 \frac{\log(2/\zeta)}{\alpha}\right) \implies \text{Acc}_{\tau, \alpha, \zeta}^{\text{pow}}(s) = 0.$$

2. (Exponential) For any  $\gamma > 0$  and any  $s \geq 0$ ,

$$s \leq \tau \implies \text{Acc}_{\tau, \gamma, \zeta}^{\text{exp}}(s) = 1, \quad s \geq \tau + 2 \frac{\log(2/\zeta)}{\gamma} \implies \text{Acc}_{\tau, \gamma, \zeta}^{\text{exp}}(s) = 0.$$

**Proof** For the power kernel, set  $z := \log(s/\tau)$ . If  $s \leq \tau$  then  $z \leq 0$  and by Lemma 41 we have  $\varphi_{\alpha, \zeta}(z) = \text{ACCEPT}$  almost surely, hence  $\text{Acc}_{\tau, \alpha, \zeta}^{\text{pow}}(s) = 1$ . If instead  $s \geq \tau \exp\left(2 \frac{\log(2/\zeta)}{\alpha}\right)$ , then  $z \geq 2 \frac{\log(2/\zeta)}{\alpha}$  and Lemma 41 gives  $\varphi_{\alpha, \zeta}(z) = \text{REJECT}$  almost surely, hence  $\text{Acc}_{\tau, \alpha, \zeta}^{\text{pow}}(s) = 0$ .

For the exponential kernel, set  $z := s - \tau$  and apply Lemma 41 identically: if  $s \leq \tau$  then  $z \leq 0$  so acceptance is almost sure, while if  $s \geq \tau + 2 \frac{\log(2/\zeta)}{\gamma}$  then  $z \geq 2 \frac{\log(2/\zeta)}{\gamma}$  so rejection is almost sure. ■

**Lemma 44 (Smoothness of Soft-Truncation)** Fix  $\tau > 0$  and  $\zeta \in (0, 1)$ , and define  $k := \lceil L \rceil$ .

1. (Power; multiplicative smoothness) For any  $s, s' > 0$  such that  $e^{-L}s' \leq s \leq e^L s'$ ,

$$\text{Acc}_{\tau, \alpha, \zeta}^{\text{pow}}(s) \approx_{\varepsilon_L, \delta_L} \text{Acc}_{\tau, \alpha, \zeta}^{\text{pow}}(s') \quad \text{with} \quad \varepsilon_L := k\alpha, \quad \delta_L := \zeta \sum_{i=0}^{k-1} e^{i\alpha} \leq ke^{(k-1)\alpha} \zeta.$$

2. (Exponential; additive smoothness) For any  $s, s' \geq 0$  such that  $s' - L \leq s \leq s' + L$ ,

$$\text{Acc}_{\tau, \gamma, \zeta}^{\text{exp}}(s) \approx_{\varepsilon'_L, \delta'_L} \text{Acc}_{\tau, \gamma, \zeta}^{\text{exp}}(s') \quad \text{with} \quad \varepsilon'_L := k\gamma, \quad \delta'_L := \zeta \sum_{i=0}^{k-1} e^{i\gamma} \leq ke^{(k-1)\gamma} \zeta.$$

**Proof** We first note the following chaining bound for  $\varphi_{\varepsilon, \delta}$ . Let  $z_0, \dots, z_k$  satisfy  $|z_{i+1} - z_i| < 1$  for all  $i$ . By Lemma 41(1), for every measurable event  $S$  and each  $i$ ,

$$\Pr[\varphi_{\varepsilon, \delta}(z_i) \in S] \leq e^\varepsilon \Pr[\varphi_{\varepsilon, \delta}(z_{i+1}) \in S] + \delta.$$

Iterating this inequality yields

$$\Pr[\varphi_{\varepsilon, \delta}(z_0) \in S] \leq e^{k\varepsilon} \Pr[\varphi_{\varepsilon, \delta}(z_k) \in S] + \delta \sum_{i=0}^{k-1} e^{i\varepsilon},$$

and the reverse direction follows by swapping the roles of  $z_0$  and  $z_k$ . Hence  $\varphi_{\varepsilon, \delta}(z_0) \approx_{k\varepsilon, \delta \sum_{i=0}^{k-1} e^{i\varepsilon}} \varphi_{\varepsilon, \delta}(z_k)$ .

For the power kernel, let  $z := \log(s/\tau)$  and  $z' := \log(s'/\tau)$ . The assumption  $e^{-L}s' \leq s \leq e^L s'$  implies  $|z - z'| = |\log(s/s')| \leq L$ . Define  $k := \lceil L \rceil + 1$  and interpolate  $z_i := z + \frac{i}{k}(z' - z)$  so that  $|z_{i+1} - z_i| = |z' - z|/k < 1$ . Applying the chaining bound with  $(\varepsilon, \delta) = (\alpha, \zeta)$  gives

$$\varphi_{\alpha, \zeta}(z) \approx_{k\alpha, \zeta \sum_{i=0}^{k-1} e^{i\alpha}} \varphi_{\alpha, \zeta}(z').$$

Finally, by definition,  $\text{Acc}_{\tau, \alpha, \zeta}^{\text{pow}}(s)$  and  $\text{Acc}_{\tau, \alpha, \zeta}^{\text{pow}}(s')$  are the acceptance probabilities of these two accept/reject distributions, so the same  $(\varepsilon_L, \delta_L)$ -approximate indistinguishability holds for the induced Bernoulli laws.

The exponential case is identical with  $z := s - \tau$  and  $z' := s' - \tau$ : the condition  $|s - s'| \leq L$  implies  $|z - z'| \leq L$ , and chaining with  $(\varepsilon, \delta) = (\gamma, \zeta)$  yields the stated parameters.  $\blacksquare$

## D.2. PRIVATEBERNOULLITHRESHOLD

Let  $\varphi_{\varepsilon, \delta}$  denote the mechanism from Lemma 41.

**Algorithm 7:** PRIVATEBERNOULLITHRESHOLD $_{\kappa, \varepsilon, \delta}$  (PBT)

**Input:** Bits  $b_1, \dots, b_N \in \{0, 1\}$ ; sensitivity parameter  $\kappa > 0$ ; privacy parameters  $(\varepsilon, \delta) \in (0, 1)^2$ ; threshold  $p_{\text{th}} \in (0, 1)$ .

**Output:** ACCEPT or REJECT.

$$\hat{p} \leftarrow \frac{1 + \sum_{j=1}^N b_j}{N + 1}; \quad // \text{ Laplace smoothing; } \hat{p} \in (0, 1] \text{ always}$$

$$\hat{z} \leftarrow \frac{1}{\kappa} \log\left(\frac{\hat{p}}{p_{\text{th}}}\right) + (2 \log(2/\delta)/\varepsilon) \text{ return } \varphi_{\varepsilon, \delta}(\hat{z})$$

**Deterministic accept/reject regions.** By Lemma 41 applied to  $\hat{z}$ :

- If  $\hat{z} \leq 0$  then PBT outputs **ACCEPT** almost surely.
- If  $\hat{z} \geq (2 \log(2/\delta)/\varepsilon)$  then PBT outputs **REJECT** almost surely.

Equivalently,

$$\hat{p} \leq p_{\text{th}} e^{-\kappa(2 \log(2/\delta)/\varepsilon)} \implies \text{ACCEPT a.s.}, \quad \hat{p} \geq p_{\text{th}} \implies \text{REJECT a.s.}$$

This is the “gap” that permits privacy: the mechanism can be randomized only when  $\hat{p} \in (p_{\text{th}} e^{-\kappa(2 \log(2/\delta)/\varepsilon)}, p_{\text{th}})$ .

First, we will prove a useful monotonicity lemma that the accepted probability of PBT increases monotonically with  $p$ , when the bits are sampled from  $\text{Ber}(p)$

**Lemma 45 (PBT is monotone)** Fix parameters  $N \in \mathbb{N}$ ,  $\kappa > 0$ ,  $(\varepsilon, \delta) \in (0, 1)^2$ , and  $p_{\text{th}} \in (0, 1)$ , and let  $\text{PBT}_{\kappa, \varepsilon, \delta}$  be Algorithm 7. Define, for  $z \in \mathbb{R}$ ,

$$f(z) := \Pr [\varphi_{\varepsilon, \delta}(z) = \text{REJECT}],$$

where the probability is over the internal randomness of  $\varphi_{\varepsilon, \delta}$ .

Assume that  $f$  is nondecreasing, i.e.,

$$z \leq z' \implies f(z) \leq f(z'). \quad (\heartsuit)$$

(For the threshold-type construction in Lemma 41, this monotonicity holds.)

Then the following monotonicity properties hold.

1. **Monotonicity in the dataset.** For any two bit-vectors  $b, b' \in \{0, 1\}^N$  with  $\sum_{j=1}^N b_j \leq \sum_{j=1}^N b'_j$  (in particular, if  $b_j \leq b'_j$  for all  $j$ ),

$$\Pr [\text{PBT}(b) = \text{REJECT}] \leq \Pr [\text{PBT}(b') = \text{REJECT}],$$

and hence  $\Pr [\text{PBT}(b) = \text{ACCEPT}] \geq \Pr [\text{PBT}(b') = \text{ACCEPT}]$ .

2. **Monotonicity in the Bernoulli mean.** Let  $\mathcal{D}_p$  denote the output distribution of PBT when  $b_1, \dots, b_N \stackrel{\text{i.i.d.}}{\sim} \text{Ber}(p)$ . If  $p \leq p'$ , then

$$\mathcal{D}_p(\text{REJECT}) \leq \mathcal{D}_{p'}(\text{REJECT}), \quad \text{equivalently} \quad \mathcal{D}_p(\text{ACCEPT}) \geq \mathcal{D}_{p'}(\text{ACCEPT}).$$

**Proof** Let  $X(b) := \sum_{j=1}^N b_j$ . Algorithm 7 forms

$$\hat{p}(b) = \frac{1 + X(b)}{N + 1}, \quad \hat{z}(b) = \frac{1}{\kappa} \log \left( \frac{\hat{p}(b)}{p_{\text{th}}} \right) + (2 \log(2/\delta)/\varepsilon).$$

By inspecting the expressions, we have:

- $X(b) \leq X(b') \implies \hat{p}(b) \leq \hat{p}(b')$  (because  $\hat{p}$  is affine in  $X$  with positive slope);
- $\hat{p}(b) \leq \hat{p}(b') \implies \hat{z}(b) \leq \hat{z}(b')$  (because  $\log(\cdot)$  is increasing and  $\kappa > 0$ ).

Hence, whenever  $X(b) \leq X(b')$ , we have  $\hat{z}(b) \leq \hat{z}(b')$ .

Now condition on the (fixed) dataset  $b$ . The only remaining randomness in  $\text{PBT}(b)$  comes from the call to  $\varphi_{\varepsilon, \delta}$ , so

$$\Pr [\text{PBT}(b) = \text{REJECT}] = \Pr [\varphi_{\varepsilon, \delta}(\hat{z}(b)) = \text{REJECT}] = f(\hat{z}(b)).$$

Therefore, if  $X(b) \leq X(b')$ , then  $\hat{z}(b) \leq \hat{z}(b')$ , and by the monotonicity assumption ( $\heartsuit$ ),

$$\Pr [\text{PBT}(b) = \text{REJECT}] = f(\hat{z}(b)) \leq f(\hat{z}(b')) = \Pr [\text{PBT}(b') = \text{REJECT}],$$

which proves part (1). Since the output is binary, the ACCEPT inequality follows immediately.

For part (2), fix  $p \leq p'$  and couple the two product measures  $\text{Ber}(p)^{\otimes N}$  and  $\text{Ber}(p')^{\otimes N}$  as follows: let  $U_1, \dots, U_N \stackrel{\text{i.i.d.}}{\sim} \text{Unif}[0, 1]$  and set

$$b_j := \mathbf{1}\{U_j \leq p\}, \quad b'_j := \mathbf{1}\{U_j \leq p'\}.$$

Then  $b_j \leq b'_j$  for all  $j$ , hence  $X(b) \leq X(b')$  almost surely. Applying part (1) pointwise (for each realization of  $(U_1, \dots, U_N)$ ) gives

$$\Pr [\text{PBT}(b) = \text{REJECT} \mid U_1, \dots, U_N] \leq \Pr [\text{PBT}(b') = \text{REJECT} \mid U_1, \dots, U_N].$$

Taking expectations over the uniforms yields

$$\Pr_{\text{Ber}(p)^{\otimes N}} [\text{PBT} = \text{REJECT}] \leq \Pr_{\text{Ber}(p')^{\otimes N}} [\text{PBT} = \text{REJECT}],$$

i.e.  $\mathcal{D}_p(\text{REJECT}) \leq \mathcal{D}_{p'}(\text{REJECT})$ . Again, the ACCEPT inequality follows since the output is binary.  $\blacksquare$

**Theorem 46 (PBT: privacy, soundness, and completeness)** *Let  $b_1, \dots, b_N \stackrel{\text{i.i.d.}}{\sim} \text{Ber}(p)$  and let  $\text{PBT}_{\kappa, \varepsilon, \delta}$  be Algorithm 7 with threshold  $p_{\text{th}}$ . Fix parameters  $\beta \in (0, 1)$  and  $\gamma \in (0, 1/2)$  with  $\gamma \leq \frac{\kappa}{2}$ , and define*

$$p_{\text{low}} := p_{\text{th}} e^{-\kappa(2 \log(2/\delta)/\varepsilon)} e^{-\gamma}, \quad p_{\text{high}} := p_{\text{th}} e^{\gamma}.$$

Assume

$$N \geq \frac{12}{\gamma^2 p_{\text{low}}} \cdot \log \frac{4}{\beta}. \quad (\star)$$

1. **Completeness (low mean).** *If  $p \leq p_{\text{low}}$ , then PBT outputs ACCEPT with probability at least  $1 - \beta$ .*
2. **Soundness (high mean).** *If  $p \geq p_{\text{high}}$ , then PBT outputs REJECT with probability at least  $1 - \beta$ .*
3. **Privacy-style stability (log-sensitivity  $\kappa$ ).** *Let  $p, p'$  satisfy  $e^{-\kappa} p \leq p' \leq e^{\kappa} p$  and also  $p, p' \geq p_{\text{low}}/2$ . Let  $\mathcal{D}_p$  denote the output distribution of PBT under  $\text{Ber}(p)^{\otimes N}$ . Then*

$$\mathcal{D}_p \approx_{3\varepsilon, \delta_{\text{priv}}} \mathcal{D}_{p'}, \quad \text{where} \quad \delta_{\text{priv}} := (1 + e^{\varepsilon} + e^{2\varepsilon}) \delta + (1 + e^{2\varepsilon}) \beta.$$

Moreover, the sample size requirement  $(\star)$  can be rewritten as

$$N = O \left( \frac{\exp(2\kappa \log(2/\delta)/\varepsilon)}{\gamma^2 p_{\text{th}}} \cdot \log \frac{1}{\beta} \right).$$

**Proof** Write  $\bar{p} := \frac{1}{N} \sum_{j=1}^N b_j$  and recall  $\hat{p} = \frac{1+N\bar{p}}{N+1}$ .

**Step 1: Useful facts about  $\hat{p}, \bar{p}$ , and  $p$**  For  $p \in (0, 1)$ , write  $\Pr_p[\cdot]$  for the probability over  $b_1, \dots, b_N$   $\overset{\text{i.i.d.}}{\sim} \text{Ber}(p)$  and the internal randomness of  $\varphi_{\varepsilon, \delta}$ . Let

$$C := (2 \log(2/\delta)/\varepsilon), \quad t_{\text{acc}} := p_{\text{th}} e^{-\kappa C}.$$

By the deterministic regions of Lemma 41 applied to  $\hat{z}$ ,

$$\hat{p} \leq t_{\text{acc}} \implies \text{PBT OUTPUTS ACCEPT a.s.}, \quad \hat{p} \geq p_{\text{th}} \implies \text{PBT OUTPUTS REJECT a.s.} \quad (33)$$

Moreover,

$$\hat{p} = \frac{1 + N\bar{p}}{N+1} = \bar{p} + \frac{1 - \bar{p}}{N+1} \implies \bar{p} \leq \hat{p} \leq \bar{p} + \frac{1}{N+1}. \quad (34)$$

In particular, for any threshold  $t \in (0, 1)$ ,

$$\Pr_p[\hat{p} > t] \leq \Pr_p \left[ \bar{p} > t - \frac{1}{N+1} \right], \quad \Pr_p[\hat{p} < t] \leq \Pr_p[\bar{p} < t]. \quad (35)$$

Finally, by Lemma 45,

$$p \leq p_0 \implies \Pr_p[\text{PBT} = \text{ACCEPT}] \geq \Pr_{p_0}[\text{PBT} = \text{ACCEPT}], \quad (36)$$

$$p \geq p_0 \implies \Pr_p[\text{PBT} = \text{REJECT}] \geq \Pr_{p_0}[\text{PBT} = \text{REJECT}] \quad (37)$$

**Step 2: completeness.** Fix any  $p \leq p_{\text{low}}$ . By (36) with  $p_0 = p_{\text{low}}$ , it suffices to analyze  $p = p_{\text{low}}$ . Using (33) and (35) (with  $t = t_{\text{acc}}$ ),

$$\begin{aligned} \Pr_{p_{\text{low}}}[\text{PBT OUTPUTS ACCEPT}] &\geq \Pr_{p_{\text{low}}}[\hat{p} \leq t_{\text{acc}}] = 1 - \Pr_{p_{\text{low}}}[\hat{p} > t_{\text{acc}}] \\ &\geq 1 - \Pr_{p_{\text{low}}} \left[ \bar{p} > t_{\text{acc}} - \frac{1}{N+1} \right]. \end{aligned} \quad (38)$$

Since  $p_{\text{low}} = t_{\text{acc}} e^{-\gamma}$ , we have  $t_{\text{acc}} = e^{\gamma} p_{\text{low}} \geq (1 + \gamma) p_{\text{low}}$ . Also,  $(\star)$  implies  $N p_{\text{low}} \geq \frac{12}{\gamma^2} \log \frac{4}{\beta} \geq \frac{2}{\gamma}$ , hence  $\frac{1}{N+1} \leq \frac{1}{N} \leq \frac{\gamma}{2} p_{\text{low}}$ . Therefore

$$t_{\text{acc}} - \frac{1}{N+1} \geq \left(1 + \frac{\gamma}{2}\right) p_{\text{low}}.$$

Applying the multiplicative Chernoff upper-tail bound (with  $\gamma/2 \in (0, 1)$ ) gives

$$\Pr_{p_{\text{low}}} \left[ \bar{p} > \left(1 + \frac{\gamma}{2}\right) p_{\text{low}} \right] \leq \exp \left( -\frac{(\gamma/2)^2}{3} p_{\text{low}} N \right) = \exp \left( -\frac{\gamma^2}{12} p_{\text{low}} N \right) \leq \frac{\beta}{4} \leq \beta,$$

where the penultimate inequality uses  $(\star)$ . Plugging this into (38) yields  $\Pr_{p_{\text{low}}}[\text{PBT OUTPUTS ACCEPT}] \geq 1 - \beta$ , and (36) implies the same for all  $p \leq p_{\text{low}}$ .

**Step 3: soundness.** Fix any  $p \geq p_{\text{high}}$ . By (37) with  $p_0 = p_{\text{high}}$ , it suffices to analyze  $p = p_{\text{high}}$ . Using (33) and (35) (with  $t = p_{\text{th}}$ ),

$$\begin{aligned} \Pr_{p_{\text{high}}} [\text{PBT OUTPUTS REJECT}] &\geq \Pr_{p_{\text{high}}} [\hat{p} \geq p_{\text{th}}] = 1 - \Pr_{p_{\text{high}}} [\hat{p} < p_{\text{th}}] \\ &\geq 1 - \Pr_{p_{\text{high}}} [\bar{p} < p_{\text{th}}]. \end{aligned} \quad (39)$$

Since  $p_{\text{high}} = p_{\text{th}}e^\gamma$ , we have  $p_{\text{th}} = e^{-\gamma}p_{\text{high}}$ . Thus the multiplicative Chernoff lower-tail bound implies

$$\Pr_{p_{\text{high}}} [\bar{p} < p_{\text{th}}] = \Pr_{p_{\text{high}}} [\bar{p} < e^{-\gamma}p_{\text{high}}] \leq \exp\left(-\frac{(1-e^{-\gamma})^2}{2}p_{\text{high}}N\right) \leq \exp\left(-\frac{\gamma^2}{8}p_{\text{high}}N\right),$$

where we used  $1 - e^{-\gamma} \geq \gamma/2$  for  $\gamma \in (0, 1/2)$ . Since  $p_{\text{high}} \geq p_{\text{low}}$ , assumption  $(\star)$  yields

$$\exp\left(-\frac{\gamma^2}{8}p_{\text{high}}N\right) \leq \exp\left(-\frac{\gamma^2}{8}p_{\text{low}}N\right) \leq \exp\left(-\frac{12}{8}\log\frac{4}{\beta}\right) = \left(\frac{\beta}{4}\right)^{3/2} \leq \beta.$$

Plugging into (39) yields  $\Pr_{p_{\text{high}}} [\text{PBT OUTPUTS REJECT}] \geq 1 - \beta$ , and (37) implies the same for all  $p \geq p_{\text{high}}$ .

**Step 4: privacy-style stability.** Recall our assumptions that  $e^{-\kappa}p \leq p' \leq e^\kappa p$  and  $p, p' \geq p_{\text{low}}/2$ . Let

$$\hat{z} = \frac{1}{\kappa} \log\left(\frac{\hat{p}}{p_{\text{th}}}\right) + C, \quad z(p) := \frac{1}{\kappa} \log\left(\frac{p}{p_{\text{th}}}\right) + C, \quad C := (2\log(2/\delta)/\varepsilon).$$

Then

$$|z(p) - z(p')| = \frac{1}{\kappa} |\log(p/p')| \leq 1. \quad (3)$$

**Step 4(a): multiplicative concentration implies log concentration.** Define the good event

$$\mathcal{G}(p) := \{e^{-\gamma}p \leq \bar{p} \leq e^\gamma p\}.$$

For  $\gamma \in (0, 1/2)$  we have  $e^\gamma - 1 \geq \gamma$  and also  $1 - e^{-\gamma} \geq \frac{2}{3}\gamma$ .

By standard multiplicative Chernoff bounds for  $\bar{p}$ :

$$\begin{aligned} \Pr_p [\bar{p} > e^\gamma p] &\leq \exp\left(-\frac{(e^\gamma - 1)^2}{3}pN\right) \leq \exp\left(-\frac{\gamma^2}{3}pN\right), \\ \Pr_p [\bar{p} < e^{-\gamma}p] &\leq \exp\left(-\frac{(1 - e^{-\gamma})^2}{2}pN\right) \leq \exp\left(-\frac{2\gamma^2}{9}pN\right), \end{aligned}$$

Therefore, for all  $p$ ,

$$\Pr_p [\mathcal{G}(p)^c] \leq \exp\left(-\frac{\gamma^2}{3}pN\right) + \exp\left(-\frac{2\gamma^2}{9}pN\right) \leq 2\exp\left(-\frac{2\gamma^2}{9}pN\right) \quad (40)$$

On  $\mathcal{G}(p)$  we have  $|\log(\bar{p}/p)| \leq \gamma$ .

**Step 4(b): a good event implies  $|\hat{z} - z(p)| \leq 1$ .** Fix  $p$  with  $p \geq p_{\text{low}}/2$  and suppose we are in the good event  $\mathcal{G}(p)$  from the previous step. We always have  $\bar{p} \leq \hat{p} \leq \bar{p} + \frac{1}{N+1}$  from Step 1. Under  $(\star)$  we showed in Step 2 that

$$\frac{1}{N+1} \leq \frac{\gamma}{2} p_{\text{low}} \leq \gamma p. \quad (41)$$

On  $\mathcal{G}(p)$  we thus have

$$\hat{p} \leq \bar{p} + \frac{1}{N+1} \leq e^\gamma p + \gamma p = p(e^\gamma + \gamma) \leq p e^{2\gamma},$$

where the last inequality holds for all  $\gamma \in (0, 1/2)$  since  $e^{2\gamma} - e^\gamma - \gamma \geq 0$ . Also  $\hat{p} \geq \bar{p} \geq e^{-\gamma} p$  on  $\mathcal{G}(p)$ . Therefore, on  $\mathcal{G}(p)$ ,

$$e^{-\gamma} \leq \frac{\hat{p}}{p} \leq e^{2\gamma} \quad \implies \quad \left| \log \left( \frac{\hat{p}}{p} \right) \right| \leq 2\gamma.$$

Equivalently,

$$|\hat{z} - z(p)| = \frac{1}{\kappa} \left| \log \left( \frac{\hat{p}}{p} \right) \right| \leq \frac{2\gamma}{\kappa} \leq 1, \quad (42)$$

where we used the assumption  $\gamma \leq \kappa/2$ .

By (40) and  $p \geq p_{\text{low}}/2$ ,

$$\Pr_p[\mathcal{G}(p)^c] \leq 2 \exp\left(-\frac{2\gamma^2}{9} pN\right) \leq 2 \exp\left(-\frac{\gamma^2}{9} p_{\text{low}}N\right) \leq 2 \exp\left(-\frac{4}{3} \log \frac{4}{\beta}\right) = 2 \left(\frac{\beta}{4}\right)^{4/3} \leq \beta,$$

so we may summarize:

$$\Pr_p[|\hat{z} - z(p)| \leq 1] \geq 1 - \beta, \quad \Pr_{p'}[|\hat{z}' - z(p')| \leq 1] \geq 1 - \beta. \quad (43)$$

**Step 4(c): a ‘‘DP under a good event’’ claim.** For  $z \in \mathbb{R}$  and an event  $\mathcal{S} \subseteq \{\text{ACCEPT}, \text{REJECT}\}$ , write

$$q(z) := \Pr[\varphi_{\varepsilon, \delta}(z) \in \mathcal{S}].$$

By Lemma 41, if  $|z - z'| \leq 1$  then

$$q(z) \leq e^\varepsilon q(z') + \delta \quad \text{and} \quad q(z') \leq e^\varepsilon q(z) + \delta. \quad (44)$$

*Claim.* Let  $Z$  be any real-valued random variable and  $z_0 \in \mathbb{R}$ . If  $\Pr[|Z - z_0| \leq 1] \geq 1 - \beta$ , then for every event  $\mathcal{S}$ ,

$$\Pr[\varphi_{\varepsilon, \delta}(Z) \in \mathcal{S}] \leq e^\varepsilon \Pr[\varphi_{\varepsilon, \delta}(z_0) \in \mathcal{S}] + \delta + \beta, \quad (45)$$

and also

$$\Pr[\varphi_{\varepsilon, \delta}(z_0) \in \mathcal{S}] \leq e^\varepsilon \Pr[\varphi_{\varepsilon, \delta}(Z) \in \mathcal{S}] + \delta + \beta. \quad (46)$$

*Proof of claim.* Let  $E = \{|Z - z_0| \leq 1\}$ . Then, using (44) pointwise on  $E$  and the trivial bound  $q(Z) \leq 1$ ,

$$\Pr[\varphi(Z) \in \mathcal{S}] = \mathbb{E}[q(Z)\mathbf{1}_E] + \mathbb{E}[q(Z)\mathbf{1}_{E^c}] \leq \mathbb{E}[(e^\varepsilon q(z_0) + \delta)\mathbf{1}_E] + \Pr[E^c] \leq e^\varepsilon q(z_0) + \delta + \beta,$$

which is (45). For (46), apply the reverse inequality in (44) pointwise on  $E$ :  $q(z_0) \leq e^\varepsilon q(Z) + \delta$ , multiply by  $\mathbf{1}_E$ , take expectations, and add the term  $q(z_0) \Pr[E^c] \leq \Pr[E^c] \leq \beta$ .

**Step 4(d): chaining the three comparisons.** Let  $\mathcal{D}_p$  denote the output distribution of PBT under  $\text{Ber}(p)^{\otimes N}$ . For an event  $\mathcal{S} \subseteq \{\text{ACCEPT}, \text{REJECT}\}$ , we have

$$\mathcal{D}_p(\mathcal{S}) = \Pr_p[\varphi_{\varepsilon, \delta}(\hat{z}) \in \mathcal{S}].$$

Applying the claim with  $Z = \hat{z}$  and  $z_0 = z(p)$ , and using (43), gives

$$\mathcal{D}_p(\mathcal{S}) \leq e^\varepsilon \Pr[\varphi_{\varepsilon, \delta}(z(p)) \in \mathcal{S}] + \delta + \beta. \quad (47)$$

Next, since  $|z(p) - z(p')| \leq 1$  by (3), Lemma 41 implies

$$\Pr[\varphi_{\varepsilon, \delta}(z(p)) \in \mathcal{S}] \leq e^\varepsilon \Pr[\varphi_{\varepsilon, \delta}(z(p')) \in \mathcal{S}] + \delta. \quad (48)$$

Finally, applying the claim again (now in the reverse direction (46)) with  $Z = \hat{z}'$  and  $z_0 = z(p')$  yields

$$\Pr[\varphi_{\varepsilon, \delta}(z(p')) \in \mathcal{S}] \leq e^\varepsilon \mathcal{D}_{p'}(\mathcal{S}) + \delta + \beta. \quad (49)$$

Combining (47), (48), and (49) gives

$$\mathcal{D}_p(\mathcal{S}) \leq e^{3\varepsilon} \mathcal{D}_{p'}(\mathcal{S}) + (1 + e^\varepsilon + e^{2\varepsilon})\delta + (1 + e^{2\varepsilon})\beta.$$

Swapping the roles of  $p$  and  $p'$  yields the reverse inequality, hence  $\mathcal{D}_p \approx_{3\varepsilon, \delta_{\text{priv}}} \mathcal{D}_{p'}$  with  $\delta_{\text{priv}} = (1 + e^\varepsilon + e^{2\varepsilon})\delta + (1 + e^{2\varepsilon})\beta$ .  $\blacksquare$

### D.3. Privacy of Poisson

**Lemma 47 (Poisson rate change is  $(\varepsilon, \delta)$ -close)** *Let  $\delta \in (0, 1/4]$ ,  $\varepsilon \in (0, 1]$ , and  $\lambda \geq 8 \log(4/\delta)$ . Fix  $\alpha \in (0, 1/4]$  and  $\lambda' := (1 - 2\alpha)\lambda$ . If*

$$\alpha \leq \frac{\varepsilon}{16\sqrt{3\lambda \log(4/\delta)}},$$

then

$$\text{Pois}(\lambda) \approx_{\varepsilon, \delta} \text{Pois}(\lambda').$$

**Proof** [Proof of Lemma 47] Set  $q := 1 - 2\alpha$  and write  $P$  for  $\text{Pois}(\lambda)$  and  $Q$  for  $\text{Pois}(\lambda')$ . For  $k \in \mathbb{N}$ ,

$$\log \frac{P(k)}{Q(k)} = (\lambda' - \lambda) + k \log(\lambda/\lambda') = -2\alpha\lambda + k \log(1/q).$$

Let  $r := \sqrt{3\lambda \log(4/\delta)}$  and define the good sets

$$G := \{k : |k - \lambda| \leq r\}, \quad G' := \{k : |k - \lambda'| \leq r\}.$$

A standard two-sided Poisson tail bound implies

$$P(G^c) \leq \delta/2 \quad \text{and} \quad Q((G')^c) \leq \delta/2,$$

using  $\lambda \geq 8 \log(4/\delta)$  and  $\lambda' \geq \lambda/2$  (since  $\alpha \leq 1/4$ ).

Now take any event  $T \subseteq \mathbb{N}$ . Then

$$P(T) \leq P(T \cap G) + \delta/2 \leq \exp\left(\sup_{k \in G} \log \frac{P(k)}{Q(k)}\right) Q(T) + \delta/2.$$

For  $k \in G$ , write  $k = \lambda + u$  with  $|u| \leq r$ . Then

$$\log \frac{P(k)}{Q(k)} = \lambda(\log(1/q) - 2\alpha) + u \log(1/q).$$

For  $\alpha \leq 1/4$  we have  $\log(1/q) = \log\left(\frac{1}{1-2\alpha}\right) \leq 4\alpha$  and  $\log(1/q) - 2\alpha \leq 4\alpha^2$ , hence

$$\sup_{k \in G} \log \frac{P(k)}{Q(k)} \leq 4\alpha^2 \lambda + 4\alpha r.$$

Under the stated bound on  $\alpha$  and  $\varepsilon \leq 1$  (so that  $4\alpha^2 \lambda \leq 4\alpha r \leq \varepsilon/2$ ), we get  $\sup_{k \in G} \log \frac{P(k)}{Q(k)} \leq \varepsilon$ , and therefore

$$P(T) \leq e^\varepsilon Q(T) + \delta/2.$$

The reverse inequality  $Q(T) \leq e^\varepsilon P(T) + \delta/2$  is identical by swapping  $(\lambda, P, G)$  with  $(\lambda', Q, G')$  and using the same bound on the likelihood ratio. Combining the two directions gives  $P \approx_{\varepsilon, \delta} Q$ . ■