# A Fast Algorithm for Adaptive Private Mean Estimation

John Duchi<sup>1,2</sup> Saminul Haque<sup>3</sup> Rohith Kuditipudi<sup>3</sup>

Departments of <sup>1</sup>Statistics, <sup>2</sup>Electrical Engineering, and <sup>3</sup>Computer Science Stanford University

January 2023

### Abstract

We design an  $(\varepsilon, \delta)$ -differentially private algorithm to estimate the mean of a d-variate distribution, with unknown covariance  $\Sigma$ , that is adaptive to  $\Sigma$ . To within polylogarithmic factors, the estimator achieves optimal rates of convergence with respect to the induced Mahalanobis norm  $\|\cdot\|_{\Sigma}$ , takes time  $\widetilde{O}(nd^2)$  to compute, has near linear sample complexity for sub-Gaussian distributions, allows  $\Sigma$  to be degenerate or low rank, and adaptively extends beyond sub-Gaussianity. Prior to this work, other methods required exponential computation time or the superlinear scaling  $n = \Omega(d^{3/2})$  to achieve non-trivial error with respect to the norm  $\|\cdot\|_{\Sigma}$ .

# 1 Introduction

We cannot consider the theory of differential privacy complete until we have—at least—a sample and computationally efficient estimator of the mean. To within logarithmic factors in the dimension d and sample size n, we achieve both.

To make this a bit more precise, let P be a distribution on  $\mathbb{R}^d$  with unknown mean  $\mu = \mathbb{E}_P[X]$  and unknown covariance  $\Sigma = \mathbb{E}_P[(X - \mu)(X - \mu)^T]$ , and let  $X_i \stackrel{\text{iid}}{\sim} P$ ,  $i \leq n$ . For an estimator  $\widehat{\mu}$ , consider the covariance-normalized error  $\text{err}_{\Sigma}(\widehat{\mu}, \mu) := (\widehat{\mu} - \mu)^T \Sigma^{-1}(\widehat{\mu} - \mu)$ . We give an  $(\varepsilon, \delta)$ -differentially private estimator  $\widetilde{\mu}$  of  $\mu$  such that, assuming the vectors  $\Sigma^{-1/2}X_i$  are sub-Gaussian and  $n = \widetilde{\Omega}(d/\varepsilon^2)$ ,

$$\operatorname{err}_{\Sigma}(\widetilde{\mu}, \mu) = (\widetilde{\mu} - \mu)^{T} \Sigma^{-1}(\widetilde{\mu} - \mu) \leq \widetilde{O}(1) \left[ \frac{d + \log \frac{1}{\delta}}{n} + \frac{d^{2} \log^{2} \frac{1}{\delta}}{n^{2} \varepsilon^{2}} \right]$$
(1)

with probability at least  $1 - \delta$ , where the  $\widetilde{O}(1)$  term hides dependence on the sub-Gaussian parameter of  $\Sigma^{-1/2}X$  and logarithmic factors in n. Except for a factor of  $\log \frac{1}{\delta}$  and the hidden logarithmic factors in n, this is optimal, and the method extends naturally to distributions with heavier tails for which we can provide similar near-optimal guarantees.

By measuring error with respect to the covariance  $\Sigma$  of the data itself, we adopt the familiar efficiency goals of classical theoretical statistics: that an estimator should be adaptive to structure in covariates and should have (near)-optimal covariance. Mean estimation is, of course, one of the most basic problems in statistics, and we have known for seventy-odd years that the sample mean  $\overline{X}_n := \frac{1}{n} \sum_{i=1}^n X_i$  is efficient [9, 30], achieving the optimal error  $\mathbb{E}[(\overline{X}_n - \mu)^T \Sigma^{-1}(\overline{X}_n - \mu)] = \frac{d}{n}$ , with high-probability guarantees under appropriate moment assumptions [39]. Perhaps stating the obvious, the sample mean is adaptive to the covariance of the distribution: no matter  $\Sigma$ , the sample mean is efficient.

When we require estimators to be private, however, the story is less clear. While differential privacy [14, 13] has become the *de facto* choice for protecting sensitive data in the sixteen or so years since its release—with substantial theoretical advances and successful applications [17,

3, 1, 28, 18, 12]—we know of no computationally efficient procedures that achieve order-optimal sample complexity with respect to the natural Mahalanobis norm  $||v||_{\Sigma} = \sqrt{v^T \Sigma^{-1} v}$  the population P induces via its covariance. Brown et al. [7] highlight this, developing sample efficient procedures that achieve small error in the Mahalanobis metric even when  $\Sigma$  is unknown. When the covariance  $\Sigma$  is known, estimators that truncate the data relative to  $\Sigma$  and add Gaussian noise to such a trimmed mean with covariance proportional to  $\Sigma$  suffice to privately estimate  $\mu$  (under approximate differential privacy) with the essentially optimal rate (1), so that  $n = \Omega(d)$  observations suffice to estimate  $\mu$  (see, e.g., [6, 7]). But in the more realistic setting that  $\Sigma$  is unknown, to the best of our knowledge all prior work either requires a sample of size  $n = \Omega(d^{3/2})$ ; is intractable, taking time exponential in n or d to compute; or assumes P is isotropic. Many of these further assume P is Gaussian, a stringent assumption that never obtains in practice. See Section 1.1 for more discussion.

Our contribution is a polynomial-time private estimator (Algorithm 4, PRIVMEAN) whose error matches the error achievable when the covariance is known (equivalently, the data is isotropic) to polylogarithmic factors. In essence, our estimator privatizes a stable estimate of the empirical mean by adding Gaussian noise with covariance proportional to a stable estimate of the empirical covariance; it takes time  $\tilde{O}(nd^2)$  to compute, has (nearly) linear sample complexity for sub-Gaussian distributions, allows  $\Sigma$  to be degenerate or low-rank, and naturally extends beyond sub-Gaussianity.

### 1.1 Related work

There are many connections between differential privacy and robust statistics [11], in that the major focus of robust statistics is to develop estimators insensitive to outliers and corrupted data [36, 24, 25, 19], while differential privacy makes the output (distributions) of estimators similar even when individuals in the underlying data change [14, 13, 11]. While Tukey and Huber's initiation of robust statistics is more than sixty years old [36, 24], studying statistical limits of estimation and inference from corrupted data, computational tractability was elusive: only in the last decade have researchers developed computationally efficient methods for even robustly estimating a sample mean [10]. Similarly, only recently has the community elucidated trade-offs between statistical and computational considerations in robust estimation [10].

It is natural to wonder whether such trade-offs also arise with privacy. For example, classical procedures in private query evaluation require exponential time in natural problem parameters [20, 15]. Likewise, in estimation, following the "propose, test, release" framework of Dwork and Lei [11], a number of sample efficient private estimators [32, 7, 22] require testing whether a given statistic is robust to the removal of groups of data points, which can be computationally intractable in high-dimensions. In a number of these settings, computationally efficient estimators achieving comparable sample efficiency have emerged only within the last year or so [e.g. 22, 27, 4, 2]. Our mean estimation setting is a striking example of a seemingly simple problem for which no known sub-exponential time and sample efficient algorithm exists. In particular, to the best of our knowledge, all previous work has either (i) exponential runtime [7, 31]; (ii) is sample inefficient [26, 31], requiring sample size at least  $n = \Omega(d^{3/2})$ ; or (iii) otherwise essentially assumes the population covariance  $\Sigma$  is isotropic [26, 6, 23, 31] (nominally, the paper [23] allows arbitrary covariance, but the squared error of its estimator scales at least linearly with the condition number of the population covariance  $\Sigma$ , which is effectively equivalent to assuming isotropic covariance [6]). Here we have highlighted the most relevant (recent) examples; see in the paper [7] for coverage of earlier work.

The work most closely related to ours is that of Brown et al. [7], who also consider

covariance-adaptive mean estimation and also achieve (nearly) linear sample complexity. They give a roadmap to adaptive private mean estimation that circumvents private covariance estimation, a task whose sample complexity is necessarily  $\Omega(d^{3/2})$  (see the lower bound by Dwork et al. [16]), and are the first to achieve sample complexity  $o(d^{3/2})$ , let alone linear. However, their estimators take exponential time to compute; moreover, while their accuracy analysis is independent of the condition number of  $\Sigma$ , it assumes  $\Sigma$  is full rank. Finally, they only consider Gaussian and sub-Gaussian distributions.

In concurrent work, Hopkins et al. [22] give a generic reduction from private estimation to robust estimation and leverage this reduction to obtain private estimators with (near) optimal sample complexity. While their reduction is generic, the resulting estimators are efficient only in certain special cases, e.g., for Gaussian distributions whose algebraic moment relationships allow efficient formulation, and their results for mean estimation assume bounded covariance. They extend a line of work [27, 21] on obtaining efficient approximations of inefficient private mechanisms via sum-of-squares (SoS) relaxations. While technically efficient, SoS estimators typically incur large polynomial runtime and thus scale poorly to high-dimensional settings or large amounts of data. Unlike our estimator, however, they are robust to corruption of a constant fraction the data.

# 1.2 Organization

We provide a brief outline of the paper to come. Section 2 introduces notation and covers the preliminary privacy definitions we require for our development. Our main estimator, PRIVMEAN, consists of two main parts: stably estimating the covariance of the data to reasonable accuracy and then estimating a truncated mean to which we add noise. We present our algorithms in Section 3, where Section 3.1 gives the covariance estimator, Section 3.2 the mean estimator, and Section 3.3 presents the full procedure; we analyze PRIVMEAN's privacy in Section 4, deferring some of the requisite proofs to Sections 6 and 7. We provide accuracy analysis in Section 5, where we also present ADAMEAN (Algorithm 5), which allows PRIVMEAN to adapt to the scale of the observed data.

# 2 Preliminary definitions, privacy properties, and mechanisms

To make our coming development smoother and easier, here we introduce notation and recapitulate the privacy definitions we use throughout. We also review a few standard privacy mechanisms, providing guarantees on their behavior; for those results that are not completely standard, we include proofs in the appendices for completeness.

### 2.1 Notation

Semidefinite matrices and norms For a positive semidefinite (PSD) matrix  $A \in \mathbb{R}^{d \times d}$ , we let  $\operatorname{Col}(A)$  denote its columnspace and  $A^{\dagger}$  its pseudoinverse, while the square-root of the pseudoinverse is  $A^{\dagger/2}$ . We let  $\Pi_A := A^{\dagger}A = A^{\dagger/2}A^{1/2} \in \mathbb{R}^{d \times d}$  denote the orthogonal projector onto  $\operatorname{Col}(A)$ . Using the nuclear norm  $\|A\|_* = \sum_{i=1}^n \sigma_i(A)$  (the sum of A's singular values), we define the distance-like quantity for PSD matrices A, B as

$$d_{\mathrm{psd}}(A,B) := \begin{cases} \max \left\{ \left\| A^{\dagger/2}(B-A)A^{\dagger/2} \right\|_*, \left\| B^{\dagger/2}(A-B)B^{\dagger/2} \right\|_* \right\} & \text{if } \mathrm{Col}(A) = \mathrm{Col}(B) \\ \infty & \text{otherwise}, \end{cases}$$

setting  $d_{\mathrm{psd}}(A,B)=\infty$  if A or B are not PSD. When A and B are invertible,  $d_{\mathrm{psd}}(A,B)=\max\{\|A^{-1/2}BA^{-1/2}-I\|_*,\|B^{-1/2}AB^{-1/2}-I\|_*\}$ , though we note in passing that it is not a distance. The extended-value Mahalanobis norm  $\|\cdot\|_A$  corresponding to  $A\succeq 0$  is

$$||v||_A^2 := \lim_{t \downarrow 0} v^T (A + tI)^{-1} v = \begin{cases} v^T A^{\dagger} v & v \in \mathsf{Col}(A) \\ +\infty & \text{otherwise.} \end{cases}$$

When A is non-singular, this is the standard  $||v||_A = \sqrt{v^T A^{-1} v}$ , and the norm has the monotonicity property that if  $A \leq B$ , then  $||v||_A \geq ||v||_B$  for all  $v \in \mathbb{R}^d$ .

Sets and Partitions For sets S, S', define the distance  $d_{\text{sym}}(S, S') := \max\{|S \setminus S'|, |S' \setminus S|\}$ . Given integers n and b, where we assume b divides n for simplicity, we let  $\mathcal{P}_{n,b}$  be the set of all partitions of [n] such that each subset constituting the partition has b elements. We represent a given partition in  $\mathcal{P}_{n,b}$  as a tuple of subsets  $S = (S_1, \ldots, S_{n/b})$ , where each  $S_j \subset [n]$  has b elements and are pairwise disjoint.

**Distributions** We let  $W \sim \mathsf{Lap}(\sigma)$  denote that W has Laplace distribution with scale  $\sigma$ , with density  $p(w) = \frac{1}{2\sigma} \exp(-|w|/\sigma)$ .  $X \sim \mathsf{N}(\mu, \Sigma)$  indicates that X is normal with mean  $\mu$  and covariance  $\Sigma \succeq 0$ , where if  $\Sigma$  is not full rank we mean that X has support  $\mu + \mathsf{Col}(\Sigma)$ .

## 2.2 Privacy definition and basic properties

It will be convenient for us to use closeness of distributions in our derivations (cf. [12, Ch. 3.5]), so we frame differential privacy as a type of closeness in distribution.

**Definition 1**  $((\varepsilon, \delta)$ -closeness). Probability distributions P and Q are  $(\varepsilon, \delta)$ -close in distribution, denoted  $P \stackrel{d}{=}_{\varepsilon, \delta} Q$ , if for all measurable sets  $A \subset \mathcal{X}$ ,

$$P(A) \le e^{\varepsilon} Q(A) + \delta$$
 and  $Q(A) \le e^{\varepsilon} P(A) + \delta$ .

Similarly, random variables X and Y are  $(\varepsilon, \delta)$ -close,  $X \stackrel{d}{=}_{\varepsilon, \delta} Y$ , if their induced distributions are:  $\mathbb{P}(X \in \cdot) \stackrel{d}{=}_{\varepsilon, \delta} \mathbb{P}(Y \in \cdot)$ 

Differential privacy [14, 13] is then equivalent to this notion of closenss: a randomized function (or mechanism) M from an input space  $\mathcal{X}^n$  to  $\mathcal{Y}$  is then  $(\varepsilon, \delta)$ -differentially private if and only if for any vectors  $x, x' \in \mathcal{X}^n$  differing in only a single element,

$$M(x) \stackrel{d}{=}_{\varepsilon,\delta} M(x').$$

The following results on closeness are standard [12, Ch. 3].

**Lemma 2.1** (Basic composition). Let X, X', Y, Y' be random variables satisfying  $X \stackrel{d}{=}_{\varepsilon_X, \delta_X} X'$ , and  $Y \stackrel{d}{=}_{\varepsilon_Y, \delta_Y} Y'$ . Then  $(X, Y) \stackrel{d}{=}_{\varepsilon_X + \varepsilon_Y, \delta_X + \delta_Y} (X', Y')$ .

**Lemma 2.2** (Group composition). Let  $X_1, \ldots, X_k$  be random variables with  $X_i \stackrel{d}{=}_{\varepsilon_i, \delta_i} X_{i+1}$  for each i. Let  $\varepsilon_{>i} := \sum_{j=i+1}^{k-1} \varepsilon_j$ ,  $\varepsilon = \sum_{i=1}^{k-1} \varepsilon_i$ , and  $\delta = \sum_{i=1}^k e^{\varepsilon_{>i}} \delta_i$ . Then  $X_1 \stackrel{d}{=}_{\varepsilon, \delta} X_k$ .

**Lemma 2.3** (Post-Processing). Let X, Y, W be random variables. Then for any function f, if  $X \stackrel{d}{=}_{\varepsilon, \delta} Y$ , then  $f(X, W) \stackrel{d}{=}_{\varepsilon, \delta} f(Y, W)$ .

### 2.3 Mechanisms

We use several known mechanisms, and our procedures rely on their distributional closeness properties. The first is the TOPk mechanism, which (approximately) returns the largest k elements of a sample. In our analysis, it will be convenient to call the procedures we develop with noise as an argument to allow easier tracking of distributional closeness.

```
Algorithm 1: Top-k DP (TOPk)

Input: data x \in \mathbb{R}^p, threshold k

Noise: \xi_1, \xi_2 \in \mathbb{R}^p

Output: R \subseteq [n] such that |R| = k, \tilde{x} \in (\mathbb{R} \cup \{\bot\})^p

1 y_1 \leftarrow x + \xi_1

2 y_2 \leftarrow x + \xi_2

3 R \leftarrow index set comprising the k largest y_{1,j}'s

4 for j \in [p] do

5 | if j \in R then

6 | \tilde{x}_j \leftarrow y_{2,j}

7 | else

8 | \tilde{x}_j \leftarrow 1

9 return \tilde{x}
```

**Lemma 2.4** (TOPk mechanism, [34], Theorem 2.1). Let  $\gamma, \varepsilon \in \mathbb{R}_+$ . Let  $x, x' \in \mathbb{R}^p$  be such that  $||x - x'||_{\infty} \leq \gamma$ . Then for  $\xi_1, \xi_2 \sim \mathsf{Lap}(\frac{2k\gamma}{\varepsilon})^p$ ,

$$\mathtt{TOPk}(x,k;\xi_1,\xi_2) \stackrel{d}{=}_{\varepsilon,0} \mathtt{TOPk}(x',k;\xi_1,\xi_2).$$

As our procedures rely on adding Gaussian noise, we require two distributional closeness results for normal distributions. See Appendices A.1 and A.2 for proofs, which we include for completeness, as they are both tweaks of existing results [12, 33].

**Lemma 2.5** (Gaussians, distinct means). Let  $\mu_1, \mu_2 \in \mathbb{R}^d$  and let  $\Sigma \in \mathbb{R}^{d \times d}$  be PSD. Suppose  $\|\mu_1 - \mu_2\|_{\Sigma} \leq \rho$  and define

$$\tau = \begin{cases} \frac{\rho}{\varepsilon} \sqrt{2 \log \frac{5}{4\delta}} & \text{if } 0 < \varepsilon \le 1\\ \rho / \left( \sqrt{2 \log \frac{1}{\delta} + 2\varepsilon} - \sqrt{2 \log \frac{1}{\delta}} \right) & \text{otherwise.} \end{cases}$$

Then  $N(\mu_1, \tau^2 \Sigma) \stackrel{d}{=}_{\varepsilon, \delta} N(\mu_2, \tau^2 \Sigma)$ .

Brown et al. [7, Lemma 4.15] essentially give the next result, but we allow low rank covariance matrices.

**Lemma 2.6** (Gaussians, distinct covariances). Let  $\mu \in \mathbb{R}^d$  and  $\Sigma_1, \Sigma_2 \in \mathbb{R}^{d \times d}$  be PSD and satisfy  $d_{psd}(\Sigma_1, \Sigma_2) \leq \gamma < \infty$ . Then  $\mathsf{N}(\mu, \Sigma_1) \stackrel{d}{=}_{\varepsilon, \delta} \mathsf{N}(\mu, \Sigma_2)$  for  $\varepsilon \geq 6\gamma \log(2/\delta)$ .

We conclude with a standard guarantee for Laplacian random vectors [e.g. 12, Thm. 3.6].

**Lemma 2.7** (Laplace mechanism). Let  $\alpha, \beta > 0$  and  $Z \stackrel{\text{iid}}{\sim} \mathsf{Lap}(\beta/\alpha)$ . Then for any  $A \subseteq \mathbb{R}^m$  and  $\eta \in \mathbb{R}^m$  such that  $\|\eta\|_1 \leq \beta$ ,

$$\mathbb{P}(Z \in A) \le \exp(\alpha)\mathbb{P}(Z \in A + \eta).$$

# 3 Algorithms

As our estimator and its full analysis are fairly involved, we provide a broad overview of our procedures here. We compute the estimator, whose full treatment we give in Algorithm 4 (PRIVMEAN) in section 3.3, in two phases, consisting of a stable covariance estimate and a stable mean estimate. Each carefully prunes outliers from the data, using plug-in quantities from the remaining observations as substitutes for the usual plug-in mean and covariance.

In the first phase (Algorithm 2, COVSAFE), we obtain a robust but non-private estimate  $\widehat{\Sigma}$  of the covariance. Assuming for convenience n is even, we pair observations and initially let

$$\widehat{\Sigma} := \frac{1}{n} \sum_{i=1}^{n/2} (x_i - x_{n/2+i}) (x_i - x_{n/2+i})^T.$$

As  $x_i - x_{n/2+i}$  is symmetric, we can prune pairs of observations for which  $||x_i - x_{n/2+i}||_{\widehat{\Sigma}}$  is large (regardless of the population mean  $\mu$ ), recompute  $\widehat{\Sigma}$  on the remaining observations, and repeat until convergence. The key is that this pruning, while it provides no formal robustness guarantees, is stable to changes of a single example  $x_i$ , ensuring  $\widehat{\Sigma}$  itself is stable.

In the second phase (Algorithm 3, MEANSAFE), we first obtain a robust estimate  $\widehat{\mu}$  of the empirical mean by trimming outliers with respect to  $\|\cdot\|_{\widehat{\Sigma}}$ . Using  $\|x_i\|_{\widehat{\Sigma}}$  to determine whether  $x_i$  is influential for a mean estimate is unreliable, as the quantity may be arbitrarily large even for non-outliers if  $\|\mu\|_{\widehat{\Sigma}}$  itself is large; unfortunately, paired observations (as in the stable covariance estimation phase) are similarly unhelpful, as  $\|x_i - x_j\|_{\widehat{\Sigma}}$  could be small if both  $x_i, x_j$  are "outlying" in the same way. Instead, we randomly partition the n observations into groups S of size  $O(\log \frac{n}{\delta})$  and prune all observations in a group S if any two observations in S are far with respect to  $\|\cdot\|_{\widehat{\Sigma}}$ , so that there is at least a pair of outlying observations in the group. Assuming the total number of pruned observations across both phases is not too large—and much of our analysis shows how to make the pruned observations stable across different samples x, x'—we let  $\widehat{\mu}$  be the empirical mean of the un-pruned observations, then release  $\widehat{\mu} \sim N(\widehat{\mu}, \sigma^2(\varepsilon, \delta)\widehat{\Sigma})$ , where the privacy budget determines  $\sigma^2(\varepsilon, \delta)$ .

### 3.1 Stable covariance estimation

The first component of the private mean estimation algorithm is the covariance estimation procedure COVSAFE in Alg. 2, which removes suitably unusual pairs of data points from the sample  $x \in (\mathbb{R}^d)^n$ , then uses the remaining pairs to actually construct the covariance. The procedure maintains an empirical covariance  $\Sigma_t$  of the remaining data at each iteration t, so that  $\{\Sigma_t\}$  is a non-increasing (in the semidefinite order) sequence of matrices, and stores removed indices in an iteratively growing collection  $R_t$  for  $t = 1, 2, \ldots$ ; the procedure thus necessarily terminates after at most n/2 rounds of index removal. For convenience of our analysis, COVSAFE also returns a transcript  $\Gamma$  of the removed indices and iteratively constructed covariances, returning  $\bot$  if the data is so unstable that it removes too many indices.

The key is that the covariance estimates are appropriately stable (see Conditions (C.i) and (C.ii) to come in Section 4), and with high probability on any given input x, the algorithm guarantees that its output changes little when we remove index i or, if the data has too much variance relative to itself, that the procedure simply returns  $\hat{\Sigma} = \bot$ . To allow cleaner description of the precise results we require in our main privacy result in Section 4, for a putative bound B on  $||x_i - x_j||_{\Sigma}^2$ , acceptable number of outliers m, and privacy random

## Algorithm 2: Robust Covariance Estimation (COVSAFE)

```
Input : data x_{1:n}
      Params: threshold B, threshold m
      Noise : z \in \mathbb{R}^{n/2+1}. w \in \mathbb{R}
  \tilde{x} \leftarrow x_{1:n/2} - x_{n/2+1:n}
  2 R_0 \leftarrow \emptyset, \Sigma_0 \leftarrow \frac{1}{n} \sum_{i=1}^{n/2} \tilde{x}_i \tilde{x}_i^T
  3 converged \leftarrow false, t \leftarrow 0
  4 while not converged do
             t \leftarrow t+1, R_t \leftarrow R_{t-1}, \, \texttt{truncated} \leftarrow 0
  \mathbf{5}
             for i \in [n/2] \setminus R_{t-1} do
  6
                   if 2\log \|\tilde{x}_i\|_{\Sigma_{t-1}} + z_i + z_{n/2+1} > \log(B) then  R_t \leftarrow R_t \cup \{i\}  truncated \leftarrow truncated +1
  7
  8
   9
             if truncated = 0 then
               converged \leftarrow true, T \leftarrow t
12 \sum_{t} \leftarrow \frac{1}{n} \sum_{i \in [n/2] \setminus R_t} \tilde{x}_i \tilde{x}_i^T
13 \Gamma \leftarrow ([R_i]_{t=0}^T, [\Sigma_t]_{t=0}^T, T)
14 if |R_T| > m + w then
           return \perp, \Gamma
16 return \Sigma_T, \Gamma
```

variables Z and W to be specified, let

$$(\widehat{\Sigma}, \Gamma) := \mathtt{COVSAFE}_{B,m}(x; Z, W),$$
 (2a)

where  $\Gamma = ([\Sigma_t]_{t \leq T}, [R_t]_{t \leq T}, T)$  is the transcript of intermediate covariances and removed indices, and for  $\tilde{x} = x_{1:n/2} - x_{n/2+1:n}$  (as in Line 1) define the leave-one-out covariance

$$\widehat{\Sigma}_{-i} := \begin{cases} \widehat{\Sigma} - \frac{1}{n} 1 \{ i \in R_T \} \, \widetilde{x}_i \widetilde{x}_i^T & \text{if } \widehat{\Sigma} \neq \bot \\ \bot & \text{otherwise,} \end{cases}$$
 (2b)

which is  $\widehat{\Sigma}$  whenever COVSAFE does not remove index pair  $(i, n/2 + i) \in [n]^2$ .

### 3.2 Stable mean estimation

The second component of the private mean estimation algorithm is a sample mean estimator, adding noise commensurate with an estimated (positive semidefinite) noise covariance that we abstractly call  $A \in \mathbb{R}^{d \times d}$ . The procedure MEANSAFE removes elements  $x_i$  of the data x that are "too far" from the bulk of the data, measured by  $||x_i - x_{i'}||_A$ , using randomization to be sure that the removed indices are appropriately private. The algorithm uses TOPk to select groups of indices that contain too many outlying datapoints, then removes all data associated with these groups. By evaluating (random) groups of data, the procedure enforces privacy in that if the majority of the data are appropriately close to a center point as measured by covariance, then few groups have large diameter, and adding or removing a single datapoint  $x_i$  can only effect the removal of one group and the method may privately return a noisy empirical mean. When many datapoints are outliers, the method is likely to return  $\bot$  regardless of the behavior of any individual datapoint.

## Algorithm 3: Robust Mean Estimation (MEANSAFE)

```
Input: data x_{1:n}, PSD matrix A \in \mathbb{R}^{d \times d}
Params: threshold B, batchsize b, threshold k
Noise: S = (S_1, \dots, S_{n/b}) \in \mathcal{P}_{n,b}, z, z' \in \mathbb{R}^{n/b}, w \in \mathbb{R}, z^{\mathsf{N}} \in \mathbb{R}^d
Output: mean estimate \widetilde{\mu}

1 for j \in [n/b] do
2 \Big| D_j \leftarrow \log(\operatorname{diam}_A(x_{S_j}))
3 \widetilde{D} \leftarrow \operatorname{TOPk}(D, k; z, z')
4 R \leftarrow \emptyset, t \leftarrow 0 /* initialize removed indices to empty */
5 for j \in [n/b] do
6 \Big| \text{ if } \widetilde{D}_j \neq \bot \text{ and } \widetilde{D}_j > \log(\sqrt{B}/4) \text{ then}
7 \Big| R \leftarrow R \cup S_j
8 \Big| t \leftarrow t + 1
9 \widehat{\mu} \leftarrow \frac{1}{n - |R|} \sum_{i \notin R} x_i
10 if t > \frac{2k}{3} + w then
11 \Big| \widetilde{\mu} \leftarrow \bot
12 else
13 \Big| \widetilde{\mu} \leftarrow \widehat{\mu} + A^{1/2} z^{\mathsf{N}}
14 \Gamma \leftarrow (D, \widetilde{D}, R, t, \widehat{\mu})
15 return \widetilde{\mu}, \Gamma
```

For use in Section 4, as with COVSAFE, we assign notation to the outputs of MEANSAFE. Let  $x \in \mathbb{R}^{n \times d}$  be an arbitrary sample and A an arbitrary positive semidefinite matrix. For parameters defining the supposed bound B on  $||x_i - x_j||_{\Sigma}^2$ , group size b, acceptable outlier count k, and privacy random variables  $(\mathcal{S}, Z, Z', W, Z^{\mathsf{N}})$ , all to be specified later, define

$$(\widetilde{\mu}(x,A),\Gamma(x,A)) := \mathtt{MEANSAFE}_{B,b,k}(x,A;\mathcal{S},Z,Z',W,Z^{\mathsf{N}}). \tag{3}$$

### 3.3 The private mean estimation algorithm

Given COVSAFE and MEANSAFE, Algorithm 4 (PRIVMEAN) combines the two (with appropriate parameter settings) to perform private mean estimation. First, PRIVMEAN computes a stable covariance estimate via COVSAFE, and assuming the returned covariance estimate  $\widehat{\Sigma} \neq \perp$ , then computes a trimmed mean to which it adds Gaussian noise with covariance proportional to  $\widehat{\Sigma}$  using MEANSAFE. Theorem 2 in Section 4 shows that the parameter choices guarantee privacy.

We remark briefly on the runtime of PRIVMEAN. Each iteration of the **while** loop (beginning in Line 4) of COVSAFE involves a  $d \times d$  matrix inversion followed by taking (at most)  $n \geq d$  matrix-vector products, requiring  $O(nd^2)$ . We may modify COVSAFE without changing its behavior to terminate after  $m + W_{\rm cov}$  iterations, as rejecting more than  $m + W_{\rm cov}$  indices guarantees that COVSAFE (and hence PRIVMEAN) returns  $\bot$  (see Line 14). With high probability, we have  $m+W_{\rm cov}=O(\frac{1}{\varepsilon}\log\frac{1}{\delta})$ , and giving runtime  $O(nd^2\min\{n,\frac{1}{\varepsilon}\log\frac{1}{\delta}\})$ . COVSAFE's runtime dominates MEANSAFE's, giving total (high probability) runtime  $O(nd^2\min\{n,\frac{1}{\varepsilon}\log\frac{1}{\delta}\})$ . As an aside, we may convert this expected runtime into a worst-case runtime of the same order without effecting the privacy of PRIVMEAN by truncating  $W_{\rm cov}$  to scale  $\frac{1}{\varepsilon}\log\frac{1}{\delta}$ .

# Algorithm 4: Covariance Adaptive Private Mean Estimation (PRIVMEAN)

Input : data  $x_{1:n}$ 

**Params:** threshold B, privacy budget  $(\varepsilon, \delta)$ 

Output: mean estimate  $\widetilde{\mu}$ 

### Robust Covariance Estimation

```
\begin{array}{c} \mathbf{1} & \overline{m \leftarrow \frac{16}{\varepsilon} \log \frac{1}{\delta}, \ m_{\max} \leftarrow m + \frac{16}{\varepsilon} \log \frac{1 + e^{\varepsilon/4}}{\delta}} \\ \mathbf{2} & \sigma_Z \leftarrow \frac{32\sqrt{\varepsilon}B(m_{\max}+1)}{n\varepsilon}, \ \sigma_{W_{\mathrm{cov}}} \leftarrow \frac{16}{\varepsilon} \\ \mathbf{3} & Z_{\mathrm{cov}} \sim \mathsf{Lap}(\sigma_Z)^{n/2+1}, \ W_{\mathrm{cov}} \sim \mathsf{Lap}(\sigma_{W_{\mathrm{cov}}}) \\ \mathbf{4} & \widehat{\Sigma}, \Gamma_{\mathrm{cov}} \leftarrow \mathsf{COVSAFE}_{B,m}(x; Z_{\mathrm{cov}}, W_{\mathrm{cov}}) \\ \mathbf{5} & \mathbf{if} \ \widehat{\Sigma} = \bot \ \mathbf{then} \\ \mathbf{6} & \mathbf{return} \ \bot \end{array}
```

### Private Mean Estimation

```
7 b \leftarrow 1 + \log_2 \frac{6n^2}{\delta}, k \leftarrow \frac{24}{\varepsilon} \log \frac{3}{\delta} - 3

8 \sigma_{\text{top}} \leftarrow \frac{8k}{n\varepsilon} \frac{B\sqrt{e}}{1 - B\sqrt{e}/n}, \sigma_{\text{N}} \leftarrow \frac{20b\sqrt{B}}{n\varepsilon} \exp(3\sigma_{\text{top}} \log \frac{12n}{b\delta}), \sigma_{W_{\text{mean}}} \leftarrow \frac{8}{\varepsilon}

9 S \sim \text{Uni}(\mathcal{P}_{n,b}), Z_{\text{top}}, Z'_{\text{top}} \stackrel{\text{iid}}{\sim} \text{Lap}(\sigma_{\text{top}})^{n/b}, W \sim \text{Lap}(\sigma_{W_{\text{mean}}}), Z^{\text{N}} \sim \text{N}(0, \sigma_{\text{N}}^2 I_{d \times d})

10 \widetilde{\mu}, \Gamma_{\text{mean}} \leftarrow \text{MEANSAFE}_{B,b,k}(x, \widehat{\Sigma}; S, Z_{\text{top}}, Z'_{\text{top}}, W, Z^{\text{N}})

11 \text{return } \widetilde{\mu}
```

# 4 Main privacy result

The analysis of PRIVMEAN is fairly involved, though there are four key building blocks. The first two conditions involve what we term internal and external leave-one-out stability of the covariance estimates (2a) and (2b) COVSAFE returns. These conditions require that the covariance estimates (2) are appropriately stable, both in terms of removing a single element contributing to the covariance estimate  $\hat{\Sigma}$  on input x and in terms of stability across two inputs x, x' whose transformations in Line 1 of COVSAFE, i.e.,  $\tilde{x} = x_{1:n/2} - x_{n/2+1:n}$  and  $\tilde{x}' = x'_{1:n/2} - x'_{n/2+1:n}$ , differ only in a single element. Letting  $0 \le a < \infty$  be a constant to be determined later and  $\gamma \in (0,1)$  be a probability, consider the conditions

(C.i) Internal leave-one-out stability. Let  $\widehat{\Sigma}$  and  $\widehat{\Sigma}_{-i}$  be the outputs (2) of COVSAFE on an arbitrary input x of size n. Then for any index  $i \in [n/2]$ , with probability at least  $1 - \gamma$ ,

$$d_{\mathrm{psd}}(\widehat{\Sigma}, \widehat{\Sigma}_{-i}) \leq \frac{a}{n} \text{ or } \widehat{\Sigma} = \perp.$$

(C.ii) External leave-one-out stability. Let  $\widehat{\Sigma}$  and  $\widehat{\Sigma}'$  be the outputs of COVSAFE on inputs x, x' of size n such that  $\widetilde{x}$  and  $\widetilde{x}'$  differ only in index  $i \in [n/2]$ , where  $\widehat{\Sigma}_{-i}$  and  $\widehat{\Sigma}'_{-i}$  are defined as in (2b). Then

$$\widehat{\Sigma}_{-i} \stackrel{d}{=}_{\varepsilon,\delta} \widehat{\Sigma}'_{-i}$$

The second two conditions involve the noisy truncated mean estimate (3) MEANSAFE outputs. The first of these conditions (C.iii) essentially states MEANSAFE is stable over inputs x and x' differing in a single element, while the second states that MEANSAFE applied with

identical input samples x, x' but different covariance estimates A, A' is stable so long as A, A' are close in the same sense as in Condition (C.i).

(C.iii) Mean sample stability. Let  $\widetilde{\mu}(A, x)$  be the mean MEANSAFE outputs (3) on input covariance A and data x, and let x, x' differ only in one element. Then

$$\widetilde{\mu}(x,A) \stackrel{d}{=}_{\varepsilon,\delta} \widetilde{\mu}(x',A).$$

(C.iv) Mean covariance stability. If  $d_{psd}(A, A') \leq \frac{a}{n}$ , then  $\widetilde{\mu}(x, A) \stackrel{d}{=}_{\varepsilon, \delta} \widetilde{\mu}(x, A')$ .

Conditions (C.i)–(C.iv) form the basic privacy building blocks to show that PRIVMEAN is differentially private, and the following proposition—a warm-up for the full Theorem 2 to come—shows how we may relatively easily synthesize the conditions to achieve privacy.

**Proposition 1.** Let samples x, x' differ in a single element, and let  $\widehat{\Sigma}$  and  $\widetilde{\mu}(x, \widehat{\Sigma})$  and  $\widehat{\Sigma}'$  and  $\widetilde{\mu}(x', \widehat{\Sigma}')$  be the covariance and mean estimates (2) and (3) for inputs x and x', respectively. Let Conditions (C.i)–(C.iv) hold. Then

$$\widetilde{\mu}(x,\widehat{\Sigma}) \stackrel{d}{=}_{4\varepsilon,(e^{3\varepsilon}+e^{\varepsilon})\delta+(e^{2\varepsilon}+1)(\delta+\gamma)} \widetilde{\mu}(x',\widehat{\Sigma}').$$

*Proof.* As x and x' are adjacent, there exists  $i \in [n/2]$  such that  $\tilde{x}_{-i} = \tilde{x}'_{-i}$ . We have a string of approximate distributional equalities that, together with the transitivity of distributional closeness implied by group privacy (Lemma 2.2), make the proposition immediate. First, we show that conditions (C.i) and (C.iv) imply

$$\widetilde{\mu}(x,\widehat{\Sigma}) \stackrel{d}{=}_{\varepsilon,\delta+\gamma} \widetilde{\mu}(x,\widehat{\Sigma}_{-i})$$
 and  $\widetilde{\mu}(x',\widehat{\Sigma}') \stackrel{d}{=}_{\varepsilon,\delta+\gamma} \widetilde{\mu}(x',\widehat{\Sigma}'_{-i})$ .

We prove the first equality as the argument for the second is identical. Treating x as fixed, let  $\mathcal{E}$  be the event that  $d_{\text{psd}}(\widehat{\Sigma}, \widehat{\Sigma}_{-i}) \leq \frac{a}{n}$  or  $\widehat{\Sigma} = \perp$ . Then for any measurable set O we have

$$\mathbb{P}(\widetilde{\mu}(x,\widehat{\Sigma}) \in O) = \mathbb{E}\left[\mathbb{P}(\widetilde{\mu}(x,\widehat{\Sigma}) \in O \mid \widehat{\Sigma})1\{\mathcal{E}\}\right] + \mathbb{E}\left[\mathbb{P}(\widetilde{\mu}(x,\widehat{\Sigma}) \in O \mid \widehat{\Sigma})1\{\mathcal{E}^c\}\right]$$

$$\stackrel{(i)}{\leq} \mathbb{E}\left[\left(e^{\varepsilon}\mathbb{P}(\widetilde{\mu}(x,\widehat{\Sigma}_{-i}) \in O \mid \widehat{\Sigma}_{-i}) + \delta\right)1\{\mathcal{E}\}\right] + \mathbb{P}(\mathcal{E}^c)$$

$$\leq e^{\varepsilon}\mathbb{P}(\widetilde{\mu}(x,\widehat{\Sigma}_{-i}) \in O) + \delta + \gamma,$$

where inequality (i) is Condition (C.iv) and the final inequality follows from the  $\gamma$  probability bound in Condition (C.i). Second, we have the distributional approximations

$$\widetilde{\mu}(x,\widehat{\Sigma}_{-i}) \stackrel{d}{=}_{\varepsilon,\delta} \widetilde{\mu}(x,\widehat{\Sigma}'_{-i})$$

by Condition (C.ii), because post-processing preserves distributional closeness (Lemma 2.3). Finally, we observe from the mean sample stability condition (C.iii) that

$$\widetilde{\mu}(x,\widehat{\Sigma}') \stackrel{d}{=}_{\varepsilon,\delta} \widetilde{\mu}(x',\widehat{\Sigma}').$$

Combining each distributional equality, we have

$$\widetilde{\mu}(x,\widehat{\Sigma}) \stackrel{d}{=}_{\varepsilon,\delta+\gamma} \widetilde{\mu}(x,\widehat{\Sigma}_{-i}) \stackrel{d}{=}_{\varepsilon,\delta} \widetilde{\mu}(x,\widehat{\Sigma}'_{-i}) \stackrel{d}{=}_{\varepsilon,\delta+\gamma} \widetilde{\mu}(x,\widehat{\Sigma}') \stackrel{d}{=}_{\varepsilon,\delta} \widetilde{\mu}(x',\widehat{\Sigma}').$$

Apply Lemma 2.2.

Finally, then, we come to our main privacy theorem, which verifies that the procedures making up PRIVMEAN indeed satisfy Conditions (C.i)–(C.iv) with appropriate constants. We state the theorem here, giving a proof that consists of lemmas making precise the constants that appear in the conditions and whose proofs we defer.

**Theorem 2.** Let  $B < \infty$ ,  $\delta \in (0,1)$ , and let  $x, x' \in (\mathbb{R}^d)^n$  be adjacent samples, and let  $\varepsilon \leq 8$ . Define  $\delta' = (e^{3\varepsilon/4} + e^{\varepsilon/4})\delta + 2(e^{\varepsilon/2} + 1)\delta$  and let  $m \in \mathbb{N}$  be as in line 1 of PRIVMEAN. Assume that  $\delta \leq \frac{1}{n}$  and n is large enough that

$$n \geq \frac{128\sqrt{e}B\log\frac{n(1+e^{\varepsilon/4})}{\delta}}{\varepsilon}\left(m+1+\frac{16}{\varepsilon}\log\frac{1+e^{\varepsilon/4}}{\delta}\right) = O(1)\frac{B\log^2\frac{1}{\delta}}{\varepsilon^2}.$$

Then  $\text{PRIVMEAN}_{B,(\varepsilon,\delta)}(x)$  is  $(\varepsilon,\delta')$ -differentially private.

As a brief remark, the condition  $\varepsilon \leq 8$  is only for convenience; a minor modification of the proof allows arbitrary  $\varepsilon$  at the expense of a more convoluted theorem statement but in which n remains of the same order.

*Proof.* By Proposition 1, it suffices to verify Conditions (C.i)–(C.iv), where we demonstrate each holding with parameters  $(\varepsilon/4, \delta)$ . Throughout the proof, the value  $m \in \mathbb{N}$  (line 1 in PRIVMEAN) and parameter  $B < \infty$  remain tacit, as the privacy guarantee holds regardless.

First, we consider Conditions (C.i) and (C.ii) on the covariance estimates. We prove the coming two lemmas in Section 6, which begins with preliminaries that we require for their proofs before giving the proofs proper. Our first lemma provides sufficient conditions to verify Condition (C.i), internal stability. Let  $z \in \mathbb{R}^{n/2+1}$  and  $w \in \mathbb{R}$  be variables—these will be random to allow privacy presently, but we use them for the definition—and let

$$\widehat{\Sigma}(x,z,w), \left( [R_t]_{t=0}^T, [\Sigma_t]_{t=0}^T, T \right) := \texttt{COVSAFE}_{B,m}(x;z,w), \tag{4a}$$

where we leave the dependence of the transcript ( $[R_t], [\Sigma_t], T$ ) on (x, z, w) implicit, and redefine  $\widehat{\Sigma}_{-i}$  as in the definitions (2):

$$\widehat{\Sigma}_{-i}(x, z, w) := \widehat{\Sigma}(x, z, w) - \frac{1}{n} \mathbb{1}\{i \in R_T\} \, \widetilde{x}_i \widetilde{x}_i^T$$
(4b)

whenever  $\widehat{\Sigma}(x,z,w) \neq \bot$ , and  $\bot$  otherwise. Then we have

**Lemma 4.1** (Internal stability). Let  $Z_j \stackrel{\text{iid}}{\sim} \mathsf{Lap}(\sigma)$  and  $i \in [n/2]$ . Then with probability at least  $1 - \exp(-\frac{1}{4\sigma})$ , either  $\widehat{\Sigma}(x, Z, w) = \bot$  or

$$d_{\mathrm{psd}}(\widehat{\Sigma}(x,Z,w),\widehat{\Sigma}_{-i}(x,Z,w)) \le \frac{1}{1 - B\sqrt{e}/n} \frac{B\sqrt{e}}{n}.$$

See Section 6.2 for a proof of Lemma 4.1. Turning to the condition (C.ii) on external stability of COVSAFE, we compare the leave-one-out covariances  $\widehat{\Sigma}_{-i}(x,z,w)$  and  $\widehat{\Sigma}_{-i}(x',z,w)$  with input samples x and x', respectively, with identical (randomization) parameters z,w. Recalling  $\widetilde{x} = x_{1:n/2} - x_{n/2+1:n}$  and  $\widetilde{x}' = x'_{1:n/2} - x'_{n/2+1:n}$ , we have the following guarantee:

**Lemma 4.2** (External stability). Let  $\gamma \in (0,1)$ ,  $Z_j \stackrel{\text{iid}}{\sim} \mathsf{Lap}(\sigma_Z)$ ,  $W \sim \mathsf{Lap}(\sigma_{W_{\text{cov}}})$  and  $k \in \mathbb{N}$ . Define  $m_{\max} = m + \sigma_{W_{\text{cov}}} \log \frac{1}{\gamma}$  and  $\alpha = \frac{1}{\sigma_{W_{\text{cov}}}} + \frac{2\sqrt{eB}(m_{\max}+1)}{n\sigma_Z}$  and  $\beta = \frac{\gamma}{2} + \frac{n}{2} \exp(-\frac{1}{4\sigma_Z})$ . For all  $i \in [n/2]$ , if  $\tilde{x}_{-i} = \tilde{x}'_{-i}$  then

$$\widehat{\Sigma}_{-i}(x, Z, W) \stackrel{d}{=}_{2\alpha, (1+e^{\alpha})\beta} \widehat{\Sigma}_{-i}(x', Z, W).$$

The lemma effectively shows that the sets of removed indices R and R' are stable, and as they determine  $\widehat{\Sigma}_{-i}$  and  $\widehat{\Sigma}'_{-i}$ , this yields their closeness. See Section 6.3 for a proof of Lemma 4.2.

We turn to the guarantees of MEANSAFE, realizing Conditions (C.iii) and (C.iv). Recall the definition (3) of  $\widetilde{\mu}(x,A)$  as the output of MEANSAFE on input  $x \in (\mathbb{R}^d)^n$  with positive semidefinite  $A \in \mathbb{R}^{d \times d}$ , with parameters bound B, batchsize b, and threshold  $k \in \mathbb{N}$ , and  $\mathcal{S}$ , Z, Z', and W as noise. We take  $Z, Z' \in \mathbb{R}^{n/b}, W \in \mathbb{R}$  to be Laplacian random variables,  $Z^{\mathbb{N}} \in \mathbb{R}^d$  to be Gaussian, and  $\mathcal{S}$  to be a uniformly random partition of [n] into blocks of size n/b; we track their scales in giving our distributional approximation guarantees.

To more cleanly state a general sample stability guarantee, which we may use to verify Condition (C.iii), we define a number of additional parameters whose values we can determine. Let the batchsize  $b \in \mathbb{N}$  and threshold k > 0 satisfy  $b \geq 4$  and  $2b(k+1) \leq n$ . Let  $\beta_1, \gamma \in (0,1)$ , let  $\alpha \geq 0$ , and let  $\sigma_{\text{top}} > 0$  and  $\sigma_{W_{\text{mean}}} > 0$ . Define the constants

$$\Delta := \frac{5b\sqrt{B}}{2n} \exp\left(3\sigma_{\text{top}}\log\frac{2n}{b\gamma}\right), \ \beta_2 := \frac{1}{2}e^{-(k/3-1)/\sigma_{W_{\text{mean}}}} + \gamma + n^2 2^{1-b}$$

and

$$\sigma_{\mathsf{N}} = \begin{cases} \frac{\Delta}{\alpha} \sqrt{2 \log \frac{5}{4\beta_1}} & \text{if } 0 \le \alpha \le 1\\ \frac{\Delta}{\sqrt{2 \log \frac{1}{\beta_1} + 2\alpha} - \sqrt{2 \log \frac{1}{\beta_1}}} & \text{otherwise.} \end{cases}$$

With these, we have a mean-sample stability result from which Condition (C.iii) develops:

**Lemma 4.3.** Let the conditions above hold and let  $Z_j, Z_j' \stackrel{\text{iid}}{\sim} \mathsf{Lap}(\sigma_{\mathrm{top}}), W \sim \mathsf{Lap}(\sigma_{W_{\mathrm{mean}}}), Z_j^{\mathsf{N}} \stackrel{\text{iid}}{\sim} \mathsf{N}(0, \sigma_{\mathsf{N}}^2)$  in (3). If x and x' are adjacent, then

$$\widetilde{\mu}(x,A) \stackrel{d}{=}_{\alpha+1/\sigma_{W_{\text{mean}}},\beta_1+\beta_2} \widetilde{\mu}(x',A).$$

See Section 7.1 for a proof.

The last building block in the argument is to demonstrate Condition (C.iv), that the estimates  $\widetilde{\mu}(x,A)$  and  $\widetilde{\mu}(x,A')$  are close when A,A' are close. For this, we give the following lemma with general noise parameters.

**Lemma 4.4.** Let  $b, k \in \mathbb{N}$ ,  $\beta \in (0,1)$ , and  $a, \sigma_{\mathbb{N}}, \alpha_2 > 0$ . Define  $\alpha_1 = \frac{6a}{n} \log \frac{2}{\beta}$ , and define the noise scale  $\sigma_{\text{top}} = \frac{ka}{n\alpha_2}$ . Then for  $Z_j, Z'_j \stackrel{\text{iid}}{\sim} \mathsf{Lap}(\sigma_{\text{top}}), Z_j^{\mathbb{N}} \stackrel{\text{iid}}{\sim} \mathsf{N}(0, \sigma_{\mathbb{N}}^2)$ , if

$$d_{\mathrm{psd}}(A, A') \leq \frac{a}{n}$$
 then  $\widetilde{\mu}(x, A) \stackrel{d}{=}_{\alpha_1 + \alpha_2, \beta} \widetilde{\mu}(x, A')$ .

See Section 7.2 for a proof.

For the final step, we put all the pieces together to prove the theorem. We give each of the lemmas so the associated condition (of (C.i)–(C.iv)) holds with parameters ( $\varepsilon/4$ ,  $\delta$ ), after which we can then apply Proposition 1 directly. We do this in a somewhat odd order because of the dependence on the noise scale between the different lemmas, beginning with

Condition (C.ii). For  $Z_j \stackrel{\text{iid}}{\sim} \mathsf{Lap}(\sigma_Z)$  and  $W \sim \mathsf{Lap}(\sigma_{W_{\text{cov}}})$ , we use Lemma 4.2 to guarantee Condition (C.ii) that  $\widehat{\Sigma}_{-i} \stackrel{d}{=}_{\varepsilon/4,\delta} \widehat{\Sigma}'_{-i}$ . From the lemma statement, we have  $\widehat{\Sigma}_{-i} \stackrel{d}{=}_{2\alpha,(1+e^{\alpha})\beta}$   $\widehat{\Sigma}'_{-i}$ , where  $\alpha = \frac{1}{\sigma_{W_{\text{cov}}}} + \frac{2\sqrt{e}B(m_{\text{max}}+1)}{n\sigma_Z}$ ,  $\beta = \frac{\gamma}{2} + \frac{n}{2}\exp(-\frac{1}{4\sigma_Z})$ , and  $m_{\text{max}} = m + \sigma_{W_{\text{cov}}}\log\frac{1}{\gamma}$  for the m in line 1 of PRIVMEAN (though privacy does not depend on its value). We first achieve

 $2\alpha \leq \frac{\varepsilon}{4}$ . Setting  $\sigma_{W_{\text{cov}}} = \frac{16}{\varepsilon}$ , it is sufficient that  $\sigma_Z$  is large enough that  $\frac{2\sqrt{e}B(m_{\text{max}}+1)}{n\sigma_Z} \leq \frac{\varepsilon}{16}$ , i.e.,

$$\sigma_Z \ge \frac{32\sqrt{e}B(m_{\max}+1)}{n\varepsilon} = \frac{32\sqrt{e}B}{n\varepsilon}\left(m+1+\frac{16\log\frac{1}{\gamma}}{\varepsilon}\right).$$

For the  $\delta$  privacy component, we wish to have  $(1+e^{\alpha})\beta \leq \delta$ . As we have guaranteed  $\alpha \leq \frac{\varepsilon}{4}$ , taking  $\gamma = \frac{\delta}{1+e^{\varepsilon/4}}$  and making sure  $\sigma_Z$  is small enough that  $n \exp(-\frac{1}{4\sigma_Z}) \leq \gamma = \frac{\delta}{1+e^{\varepsilon/4}}$  suffices. For this, it is evidently sufficient that  $\frac{1}{\sigma_Z} \geq 4\log\frac{n(1+e^{\varepsilon/4})}{\delta}$ , i.e., (substituting for  $\sigma_Z$ )

$$n \ge \frac{128\sqrt{e}B\log\frac{n(1+e^{\varepsilon/4})}{\delta}}{\varepsilon} \left(m+1+\frac{16}{\varepsilon}\log\frac{1+e^{\varepsilon/4}}{\delta}\right)$$

guarantees  $\widehat{\Sigma}_{-i} \stackrel{d}{=}_{\varepsilon/4,\delta} \widehat{\Sigma}'_{-i}$ .

Condition (C.i). In Lemma 4.1, if the scale of the noise  $\sigma_Z$  on  $Z_j \stackrel{\text{iid}}{\sim} \mathsf{Lap}(\sigma_Z)$  satisfies  $\exp(-\frac{1}{4\sigma_Z}) \leq \gamma$ , Condition (C.i) holds. The choice of  $\sigma_Z$  to satisfy Condition (C.ii) above and the lower bound on n are evidently sufficient.

Condition (C.iii). Lemma 4.3 guarantees that if  $Z_j^{\mathsf{N}} \stackrel{\text{iid}}{\sim} \mathsf{N}(0, \sigma_{\mathsf{N}}^2)$ ,  $W \sim \mathsf{Lap}(\sigma_{W_{\mathrm{mean}}})$ , and  $\mathcal{S} \sim \mathsf{Uni}(\mathcal{P}_{n,b})$ , the call to MEANSAFE in line 10 of PRIVMEAN gives  $\widetilde{\mu}(x,A) \stackrel{d}{=}_{\alpha+1/\sigma_{W_{\mathrm{mean}}},\beta_1+\beta_2}$   $\widetilde{\mu}(x',A)$ , with  $\Delta$ ,  $\beta_2$  and  $\sigma_{\mathsf{N}}$  as defined in the lemma. To achieve  $\alpha + \frac{1}{\sigma_{W_{\mathrm{mean}}}} = \frac{\varepsilon}{4}$ , take  $\sigma_{W_{\mathrm{mean}}} = \frac{\varepsilon}{\varepsilon}$  and choose  $\alpha = \frac{\varepsilon}{8}$ . To achieve  $\beta_1 + \beta_2 \leq \delta$ , choose  $\beta_1 = \frac{\delta}{2}$  and then recognize that  $\beta_2 \leq \frac{\delta}{2}$  as long as  $\gamma \leq \frac{\delta}{6}$ ,  $n^2 2^{1-b} \leq \frac{\delta}{6}$  (or  $b \geq \log_2 \frac{6n^2}{\delta} + 1$ ) and  $\frac{1}{2} \exp(-\frac{k/3+1}{\sigma_{W_{\mathrm{mean}}}}) \leq \frac{\delta}{6}$  (or  $k \geq \frac{24}{\varepsilon} \log \frac{3}{\delta} - 3$ ). Thus, we arrive at

$$\sigma_{\mathsf{N}} = \frac{8\Delta}{\varepsilon} \sqrt{\log \frac{5}{2\delta}} = \frac{20\sqrt{B}b}{n\varepsilon} \exp\left(3\sigma_{\mathsf{top}} \log \frac{12n}{b\delta}\right)$$

for (any)  $b \ge 2\log\frac{6n}{\delta} + 1$  so long as  $\frac{\varepsilon}{8} \le 1$ . (Otherwise we may use the alternative value for  $\sigma_N$  preceding Lemma 4.3, which the  $(\varepsilon, \delta)$ -differential privacy guarantee of Lemma 2.5 justifies.)

Condition (C.iv). The last condition to verify is that  $\widetilde{\mu}(x,A) \stackrel{d}{=}_{\varepsilon/4,\delta} \widetilde{\mu}(x,A')$  for close enough A,A'. For this, we use Lemma 4.4, which guarantees that  $\widetilde{\mu}(x,A) \stackrel{d}{=}_{\alpha_1+\alpha_2,\delta} \widetilde{\mu}(x,A')$  for  $\alpha_1 = \frac{6a}{n}\log\frac{2}{\delta}$ , where we take  $a = \frac{B\sqrt{\varepsilon}}{1-B\sqrt{\varepsilon}/n}$  via Lemma 4.1, and arbitrary  $\alpha_2 > 0$ . Set  $\alpha_2 = \frac{\varepsilon}{8}$  and obtain  $\sigma_{\text{top}} = \frac{8ka}{n\varepsilon}$ . When  $n \geq \frac{48a}{\varepsilon}\log\frac{2}{\delta}$ , we have  $\alpha_1 \leq \frac{\varepsilon}{8}$ , and so the desired privacy holds.

Making appropriate substitutions gives that each of conditions (C.i)–(C.iv) holds with parameters  $(\varepsilon/4, \delta)$ . Proposition 1 gives the theorem.

# 5 Accuracy analysis

The second important component of our analysis of PRIVMEAN is its accuracy. We provide two accuracy results: the first (Theorem 3) covers the case in which the data is sub-Gaussian,

where we assume the method has some knowledge of the sub-Gaussian parameter of the sampling distribution. Of course, it is unreasonable to assume that a given distribution is sub-Gaussian or that we know its sub-Gaussian norms, and thus we extend PRIVMEAN via a procedure that adapts to the actual scale of the data in Section 5.2.

Throughout this section, we let P be a distribution on  $\mathbb{R}^d$  with mean  $\mu$  and covariance  $\Sigma$ , and we assume  $X_i \stackrel{\text{iid}}{\sim} P$ ,  $i = 1, \ldots, n$ . The classical (non-private) sample mean and covariance are  $\overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  and  $\overline{\Sigma}_n = \frac{1}{n} \sum_{i=1}^{n/2} (X_i - X_{n/2+i})(X_i - X_{n/2+i})^T$ . We assume throughout that P enjoys certain concentration properties, though we emphasize that our methods will be adaptive to the parameters we specify here.

**Assumption A1** (Sample concentration). Let  $c_1 \geq 64e$  and  $\beta \in (0,1)$ . For  $X_i \stackrel{\text{iid}}{\sim} P$  with  $\mathbb{E}[X] = \mu$  and  $\mathsf{Cov}(X) = \Sigma$ , there exists  $M < \infty$  such that the event

$$\mathcal{E}_{\text{samp}} := \left\{ \max_{i \in [n]} \|X_i - \mu\|_{\Sigma}^2 \le M^2 / c_1 \text{ and } \frac{1}{2}\Sigma \le \overline{\Sigma}_n \le \frac{3}{2}\Sigma \right\}$$

occurs with probability at least  $1 - \beta$ .

It is useful to give some context for the values of M we expect under various distributional assumptions. Because  $\mathbb{E}[\|X_i - \mu\|_{\Sigma}^2] = d$ , the constant  $M^2$  typically scales at least as d. We now give more detailed examples. In each, we let  $Z_i = \Sigma^{-1/2}(X_i - \mu)$  be the whitened data, defining the sample covariance  $\overline{\Sigma}_Z = \frac{1}{n} \sum_{i=1}^{n/2} (Z_i - Z_{n/2+i}) (Z_i - Z_{n/2+i})^T$ . Because  $\|X_i - \mu\|_{\Sigma} = \|Z_i\|_2$  and  $\overline{\Sigma}_n = \Sigma^{1/2} \overline{\Sigma}_Z \Sigma^{1/2}$ , we have the equivalence

$$\mathcal{E}_{\text{samp}} = \left\{ \max_{i \in [n]} \|Z_i\|_2^2 \le M^2/c_1 \text{ and } \|\overline{\Sigma}_Z - I\|_{\text{op}} \le \frac{1}{2} \right\}.$$

**Example 1** (Sub-Gaussian random vectors): If for all v satisfying  $||v||_2 \le 1$  the scalar random variable  $\langle Z, v \rangle$  is  $\tau^2$ -sub-Gaussian,

$$M^2 \le O(1)\tau^2 \left[d + \log \frac{n}{\beta}\right].$$

Indeed, a standard covering argument (see, e.g., [39, Ch. 5] or [38, Ch. 4]) gives that for all  $t \geq 0$ ,  $\mathbb{P}(\|Z\|_2 \geq t) \leq 4^d \exp(-ct^2/\tau^2)$ , where c > 0 is a numerical constant. Replacing  $t^2$  with  $O(1)(d\tau^2 + \tau^2 t^2)$  gives that  $\mathbb{P}(\|Z\|_2 \geq C\tau\sqrt{d+t^2}) \leq \exp(-t^2)$ , and for any  $\gamma > 0$ , setting  $t^2 = \log \frac{n}{\gamma}$  yields that with probability at least  $1 - \gamma$ ,

$$\max_{i \le n} \|Z_i\|_2^2 \le O(1)\tau^2 \left[ d + \tau^2 \log \frac{n}{\gamma} \right].$$

To control the covariance, we use Vershynin [37, Theorem 5.39], which gives that with probability at least  $1-2e^{-ct^2}$ ,  $\|\overline{\Sigma}_Z - I\|_{\text{op}} \leq O(1)\tau^2 \max\{\sqrt{d/n} + t/\sqrt{n}, d/n + t^2/n\}$ , so that (igoring the sub-Gaussian constant) for  $n \gtrsim d$ , setting  $t^2 = O(1)\log\frac{1}{\gamma}$  gives  $\|\overline{\Sigma}_Z - I\|_{\text{op}} \leq \frac{1}{2}$  with probability at least  $1-\gamma$ . Set  $\gamma = \beta/2$ .  $\Diamond$ 

**Example 2** (General moment bounds): Suppose for some  $p \ge 4$  we have  $\mathbb{E}[\|X_i - \mu\|_{\Sigma}^p] = \mathbb{E}[\|Z_i\|_2^p] \le \tau^p d^{p/2}$ , where necessarily  $\tau \ge 1$ . Then we can give two results: the first being that asymptotically  $M = o(n^{1/p})$  and the second more quantitative. For the first, we claim that  $\max_{i \le n} \|Z_i\|_2 / n^{1/p} \stackrel{a.s.}{\to} 0$ . To see this, note that for any  $\varepsilon > 0$ ,

$$\infty > \frac{1}{\varepsilon^p} \mathbb{E}[\|Z_1\|_2^p] = \int_0^\infty \mathbb{P}(\|Z_1\|_2^p \ge \varepsilon^p t) dt \ge \sum_{i=1}^\infty \mathbb{P}(\|Z_i\|_2^p \ge \varepsilon^p i).$$

By the Borel-Cantelli lemma, the event  $||Z_n||_2 \ge \varepsilon n^{1/p}$  occurs only finitely often, and so the claim follows. Meanwhile, the strong law of large numbers guarantees that  $\overline{\Sigma}_Z \stackrel{a.s.}{\to} I$ .

For more quantitative parameters, we first get by Markov's inequality that

$$\mathbb{P}(\max_{i \in [n]} \|Z_i\|_2 > t) \le \frac{n\mathbb{E}[\|Z_1\|_2^p]}{t^p} \le \frac{n\tau^p d^{p/2}}{t^p},$$

so setting  $M \simeq \tau \sqrt{d} n^{1/p}/\beta^{1/p}$ , we have  $\max_i ||Z_i||_2 \leq M/c_1$  with probability at least  $1 - \beta$ . To show concentration of the covariance matrix, we apply Chen et al. [8, Theorem A.1 Part 2], treating p as a constant, obtaining

$$n\mathbb{E} \left[ \| \overline{\Sigma}_{Z} - I \|_{\text{op}}^{p/2} \right]^{2/p} \lesssim \sqrt{n \log d} \sqrt{\mathbb{E}[\|Z\|_{2}^{4}]} + (n^{2/p} \log d) \mathbb{E}[\|Z_{1}\|_{2}^{p}]^{2/p}$$

$$\lesssim \max\{\sqrt{n \log d}, n^{2/p} \log d\} \mathbb{E}[\|Z_{1}\|_{2}^{p}]^{2/p},$$

and so by Markov's inequality

$$\mathbb{P}(\|\overline{\Sigma}_Z - I\|_{\text{op}} > \frac{1}{2}) \lesssim \max\{n^{-p/4} \log^{p/4} d, n^{1-p/2} \log^{p/2} d\} \mathbb{E}[\|Z_1\|_2^p] \lesssim \frac{\tau^p (d \log d)^{p/2}}{n^{p/2-1}},$$

which has bound  $\beta$  when  $n \gtrsim (\tau^2 d \log d)^{p/(p-2)} \beta^{-2/(p-2)}$ .  $\Diamond$ 

## 5.1 Accuracy of PRIVMEAN

We give our promised accuracy guarantee whenever Assumption A1 holds. Though not strictly necessary, we state the theorem assuming that  $\delta$  is not too small to allow for a cleaner result. Throughout, c denotes a numerical constant whose value can change from line to line.

**Theorem 3.** Let  $\varepsilon > 0$  and  $e^{-d} \le \delta \le \frac{1}{n}$  be privacy parameters and let Assumption A1 hold. Let  $B \ge M^2$  and suppose  $n \ge \frac{c}{\varepsilon^2} B \log^2 \frac{1}{\delta}$ . Let  $\widetilde{\mu} = \mathtt{PRIVMEAN}_{B,\varepsilon,\delta}(X_{1:n})$ . Then with probability at least  $1 - (\beta + 5\delta)$ ,  $\widetilde{\mu} \ne \bot$  and

$$\|\widetilde{\mu} - \overline{X}_n\|_{\Sigma} \le \frac{c\sqrt{Bd}\log(\frac{1}{\delta})}{n\varepsilon}.$$

Proof. We first show under the event  $\mathcal{E}_{samp}$  that with probability at least  $1-4\delta$  over the randomness in PRIVMEAN, both COVSAFE and MEANSAFE prune no observations, meaning the sets of removed indices  $R=\emptyset$  in both procedures (so that Line 7 in COVSAFE and Line 6 in MEANSAFE never fail), and thus  $\widetilde{\mu}=\overline{X}_n+\overline{\Sigma}_n^{1/2}Z^N$ . As  $\overline{\Sigma}_n\preceq\frac{3}{2}\Sigma$  on  $\mathcal{E}_{samp}$ , we have that  $\|\overline{\Sigma}_n^{1/2}Z^N\|_{\Sigma}^2\leq\frac{3}{2}\|Z^N\|_2^2$ . The result then follows follows once we show that  $\|Z^N\|_2\leq\frac{c\sqrt{Bd}\log(\frac{1}{\delta})}{n\varepsilon}$  with probability at least  $1-\delta$  and take a union bound over these events and  $\mathcal{E}_{samp}$ .

Rearranging the condition in Line 7 of COVSAFE, the element  $X_i - X_{n/2+i}$  is pruned in the first iteration only if

$$(Z_{\text{cov}})_i + (Z_{\text{cov}})_{n/2+1} > \log(B) - \log(\|X_i - X_{n/2+i}\|_{\overline{\Sigma}_n}^2)$$
  
 $\stackrel{(\star)}{\geq} \log(c_1 B/8M^2) \geq \log(c_1/8),$ 

where  $(\star)$  holds for all  $i \in [n/2]$  on event  $\mathcal{E}_{\text{samp}}$  because

$$||X_i - X_{n/2+i}||_{\Sigma_n}^2 \le 2||X_i - X_{n/2+i}||_{\Sigma}^2 \le 4||X_i - \mu||_{\Sigma}^2 + 4||X_{n/2+i} - \mu||_{\Sigma}^2 \le 8M^2/c_1.$$
 (5)

As  $c_1 \geq 64e$  by Assumption A1, if  $\|Z_{cov}\|_{\infty} \leq 1/2$  then COVSAFE in line 7 prunes no entries, instead simply passing  $\overline{\Sigma}_n$  to MEANSAFE so long as  $W_{\text{cov}} + m > 0$  (see line 14). Recall that  $(Z_{\text{cov}})_j \stackrel{\text{iid}}{\sim} \mathsf{Lap}(\sigma_Z)$  for  $j = 1, \ldots, n/2 + 1$  and  $\sigma_Z = \frac{32\sqrt{e}B(m_{\text{max}} + 1)}{n\varepsilon} \leq \frac{cB\log\frac{1}{\delta}}{n\varepsilon^2}$ , so by taking a union bound over the entries, we have with probability at least  $1 - \delta$  that

$$||Z_{\text{cov}}||_{\infty} \le \frac{cB \log \frac{1}{\delta}}{n\varepsilon^2} \log \left(\frac{n/2+1}{\delta}\right) \le \frac{1}{2},$$

where the last inequality is by the assumptions that  $n \geq \frac{cB\log^2(\frac{1}{\delta})}{\varepsilon^2}$  and  $\delta \leq \frac{1}{n}$ . Also recall that  $W_{\rm cov} \sim {\sf Lap}(\frac{16}{\varepsilon})$  and  $m = \frac{16}{\varepsilon}\log\frac{1}{\delta}$ , so  $W_{\rm cov} + m > 0$  with probability at least  $1 - \frac{\delta}{2}$ . Continuing to the next phase of PRIVMEAN, MEANSAFE with input  $A = \overline{\Sigma}_n$  prunes the

indices  $S_j$  only if

$$\widetilde{D}_j = D_j + (Z'_{\text{top}})_j > \log(\sqrt{B}/4).$$

By the same argument we used to obtain inequality (5), on  $\mathcal{E}_{samp}$  we have for all  $j \in [n/b]$ that

$$D_j = \log(\operatorname{diam}_{\overline{\Sigma}_n}(X_{S_j})) \le \log(\sqrt{8M^2/c_1}),$$

and so if  $\|Z'_{\text{top}}\|_{\infty} \leq 1/2$  then

$$\widetilde{D}_j \le \log(\sqrt{8M^2/c_1}) + \frac{1}{2} \le \log(\sqrt{B}/4),$$

where the last inequality follows from  $c_1 \geq 64e$  and  $B \geq M^2$ . Thus, MEANSAFE prunes no entries, and  $\widetilde{\mu} = \overline{X}_n + \overline{\Sigma}_n^{1/2} Z^{\mathsf{N}}$  so long as  $W_{\text{mean}} + \frac{2k}{3} > 0$  (see Line 10). Recall that  $(Z'_{\text{top}})_j \stackrel{\text{iid}}{\sim} \mathsf{Lap}(\sigma_{\text{top}}) \text{ for } j = 1, \dots, n/b \text{ and } \sigma_{\text{top}} = \frac{8k}{n\varepsilon} \frac{B\sqrt{e}}{1 - B\sqrt{e}/n} \leq \frac{cB\log(\frac{1}{\delta})}{n\varepsilon^2}.$  Another union bound gives that with probability at least  $1 - \delta$ ,

$$||Z_{\text{cov}}||_{\infty} \le \frac{cB\log(\frac{1}{\delta})}{n\varepsilon^2}\log\frac{n}{b\delta} \le \frac{1}{2},$$

where the last inequality follows from the assumption  $n \geq \frac{cB \log^2 \frac{1}{\delta}}{\varepsilon^2}$  and  $\delta \leq \frac{1}{n}$ . Also,  $W_{\text{mean}} \sim \text{Lap}(\frac{8}{\varepsilon})$  and  $\frac{2k}{3} = \frac{16}{\varepsilon} \log \frac{3}{\delta} - 2$ , so  $W_{\text{mean}} - k > \frac{8}{\varepsilon} \log \frac{3}{\delta} - 2 \geq 0$  with probability at least  $1 - \delta$ .

Therefore, PRIVMEAN returns  $\overline{X}_n + \overline{\Sigma}_n^{1/2} Z^{\mathsf{N}}$  with probability at least  $1 - 4\delta$  on the event  $\mathcal{E}_{\mathrm{samp}}$ . Recall that  $Z^{\mathsf{N}} \sim \mathsf{N}(0, \sigma_{\mathsf{N}}^2 I)$  with  $\sigma_{\mathsf{N}} = \frac{20b\sqrt{B}}{n\varepsilon} \exp(3\sigma_{\mathrm{top}}\log\frac{12n}{b\delta})$ , and because  $\sigma_{\mathrm{top}} \leq \frac{cB\log(\frac{1}{\delta})}{n\varepsilon^2}$ ,  $\delta \leq \frac{1}{n}$ , and  $n \geq \frac{cB\log^2(\frac{1}{\delta})}{\varepsilon^2}$ , we have that  $\sigma_{\mathsf{N}} \leq \frac{c\sqrt{B}\log(\frac{1}{\delta})}{n\varepsilon}$ . Classical tail bounds on the  $\chi^2$ -distribution [29, Lemma 1] give with probability at least  $1 - \delta$  that

$$||Z^{\mathsf{N}}||_2^2 \le \sigma_{\mathsf{N}}^2 \left[ d + 2\sqrt{d\log\frac{1}{\delta}} + 2\log\frac{1}{\delta} \right] \le \frac{cBd\log^2\frac{1}{\delta}}{n^2\varepsilon^2},$$

where the last inequality follows from the bound on  $\sigma_N$  and assumption that  $e^{-d} \leq \delta$ . 

#### 5.2Adapting to heavy-tailed data

In practice, we may not have a priori knowledge of the concentration properties of the data. Given the necessarily slowed rates of convergence for private estimators of means of random variables with only p moments [5], it is essential to be adaptive to the actual scale (and number of moments) of the problem. We therefore develop ADAMEAN, which automatically tunes the threshold parameter B by repeatedly calling PRIVMEAN and doubling B until  $\widetilde{\mu} \neq \bot$ . The key is that upon termination of ADAMEAN, the effective B is at most twice the realized scale  $O(1)\max_{i\leq n}\|Z_i\|^2$  of the random variables. To ensure privacy irrespective of the number of calls to PRIVMEAN, with each successive call ADAMEAN progressively decreases the privacy budget allocated to PRIVMEAN; in particular, as  $\sum_{t=1}^{\infty} 1/t^2 = \pi^2/6$ , via basic composition we can bound the total privacy loss of ADAMEAN by a factor  $\frac{\pi^2}{6}$  over the privacy loss of PRIVMEAN. Aside from an extra factor polylogarithmic in B/d, ADAMEAN matches the accuracy of PRIVMEAN, as we show presently.

## Algorithm 5: Fully Adaptive Private Mean Estimation (ADAMEAN)

```
Input : data x_{1:n}

Params: privacy budget (\varepsilon, \delta)

Output: mean estimate \widetilde{\mu}

1 for t = 1, 2, ... do

2 \widetilde{\mu}_t \leftarrow \texttt{PRIVMEAN}_{d2^{t-1}, (\varepsilon/t^2, \delta/t^2)}(x)

3 if \widetilde{\mu}_t \neq \bot then

4 return \widetilde{\mu}_t
```

**Theorem 4** (Accuracy of ADAMEAN). Let  $\varepsilon > 0$  and  $e^{-d} \le \delta \le \frac{1}{n}$  be privacy parameters and let event  $\mathcal{E}_{\mathrm{samp}}$  hold. Let  $s = \max\{1, \log_2 \frac{4M^2}{d}\}$  and assume  $n \ge \frac{c \max\{d, M^2\}s^2}{\varepsilon^2} \log^2 \frac{s^2}{\delta}$ . Let  $\widetilde{\mu} = \mathrm{ADAMEAN}_{\varepsilon, \delta}(X_{1:n})$ . Then with probability at least  $1 - (5 + \pi^2/3)\delta$ ,

$$\|\widetilde{\mu} - \overline{X}_n\|_{\Sigma} \le \frac{cs^2 \log(\frac{s^2}{\delta}) \max\{M\sqrt{d}, M\log(\frac{s^2}{\delta}), d\}}{n\varepsilon}.$$
 (6)

*Proof.* Let  $t^\star$  be the smallest positive integer such that  $M^2 \leq d2^{t^\star-1}$  and let  $t_{\rm stop}$  be the iteration when ADAMEAN terminates (which may be infinite). Note that  $t^\star \leq s$ . The proof comes in two parts: on the event  $\mathcal{E}_{\rm samp}$ , we first show that either  $t_{\rm stop} > t^\star$  or  $\widetilde{\mu}$  satisfies the claim (6) with probability at least  $1 - \pi^2 \delta/6$ ; secondly, we show ADAMEAN terminates with  $t_{\rm stop} \leq t^\star$  with probability at least  $1 - 5\delta$ . The result then follows via a union bound.

We carry out the first part with the help of the following lemma.

**Lemma 5.1.** Let  $\varepsilon > 0$  and  $e^{-d} \le \delta \le \frac{1}{n}$  be privacy parameters and let  $B \ge 0$ . Suppose the event  $\mathcal{E}_{samp}$  holds and let  $\widetilde{\mu} = \mathtt{PRIVMEAN}_{B,(\varepsilon,\delta)}(X_{1:n})$ . Then with probability at least  $1-2\delta$  over the randomness of PRIVMEAN,  $\widetilde{\mu} = \bot$  or

$$\|\widetilde{\mu} - \overline{X}_n\|_{\Sigma} \le \frac{c \log \frac{1}{\delta} \max\{M \log \frac{1}{\delta}, \sqrt{Bd}\}}{n\varepsilon}.$$

Proof. Suppose  $\widetilde{\mu} \neq \bot$  as otherwise the claim is trivial. Let  $\widehat{\Sigma}$  denote the covariance estimate of COVSAFE (that is,  $\Sigma_T$  at the final iteration of COVSAFE), and let  $\widehat{\mu}$  denote the empirical mean of the observations not pruned by MEANSAFE so that  $\widetilde{\mu} = \widehat{\mu} + \widehat{\Sigma}^{1/2} Z^{\mathbb{N}}$ . Then by the condition for returning  $\bot$  in Line 11 of MEANSAFE, MEANSAFE prunes at most  $b(\frac{2k}{3} + W_{\text{mean}})$  points and so

$$\begin{aligned} \left\| \widehat{\mu} - \overline{X}_n \right\|_{\Sigma} &\leq \frac{b(\frac{2k}{3} + W_{\text{mean}}) \max_i \|X_i - \widehat{\mu}\|_{\Sigma}}{n} \\ &\leq \frac{b(\frac{2k}{3} + W_{\text{mean}}) (\max_i \|X_i - \mu\|_{\Sigma} + \|\mu - \widehat{\mu}\|_{\Sigma})}{n} \end{aligned}$$

$$\leq \frac{2b(\frac{2k}{3} + W_{\text{mean}}) \max_i \|X_i - \mu\|_{\Sigma}}{n} \stackrel{(\star)}{\leq} \frac{2b(\frac{2k}{3} + W_{\text{mean}})M}{\sqrt{c_1}n},$$

with  $(\star)$  following directly from  $\mathcal{E}_{\text{samp}}$ . Recalling that  $k = \frac{24}{\varepsilon} \log \frac{3}{\delta} - 3$  and  $W \sim \text{Lap}(\frac{8}{\varepsilon})$ , it follows that  $\frac{2k}{3} + W < \frac{24}{\varepsilon} \log \frac{3}{\delta}$  with probability at least  $1 - \frac{\delta}{6}$ . Recalling also that  $b = 1 + \log_2 \frac{6n^2}{\delta}$ , it follows that on  $\mathcal{E}_{\text{samp}}$ ,

$$\|\widehat{\mu} - \overline{X}_n\|_{\Sigma} \le \frac{2(1 + \log_2 \frac{6n^2}{\delta})(\frac{24}{\varepsilon} \log \frac{3}{\delta})M}{\sqrt{c_1}n} \le \frac{cM \log^2 \frac{1}{\delta}}{n\varepsilon}$$

with probability at least  $1 - \frac{\delta}{6}$ .

Meanwhile, observe  $\mathcal{E}_{\text{samp}}$  implies  $\widehat{\Sigma} \preceq \overline{\Sigma}_n \preceq 2\Sigma$  as pruning entries (line 7) in COVSAFE only shrinks its covariance estimate. Thus,  $\|\widehat{\Sigma}^{1/2}Z^{\mathsf{N}}\|_{\Sigma} \leq \sqrt{2}\|\widehat{\Sigma}^{1/2}Z^{\mathsf{N}}\|_{\widehat{\Sigma}} = \sqrt{2}\|Z^{\mathsf{N}}\|_{2}$ . From the same argument as in Theorem 3,

$$\left\| Z^{\mathsf{N}} \right\|_2 \leq \frac{c\sqrt{Bd}\log(\frac{1}{\delta})}{n\varepsilon}$$

with probability at least  $1 - \delta$ .

The preceding two displays together imply Lemma 5.1 after taking a union bound.  $\Box$ 

Applying Lemma 5.1 with the mapping  $B = d2^{t-1}$ ,  $\varepsilon \mapsto \varepsilon/t^2$  and  $\delta \mapsto \delta/t^2$ , we have for any  $1 \le t \le t^*$  that under the event  $\mathcal{E}_{\text{samp}}$ , with probability at least  $1 - 2\delta/t^2$ , either  $\widetilde{\mu}_t = \bot$  or

$$\|\widetilde{\mu}_t - \overline{X}_n\|_{\Sigma} \le \frac{ct^2 \log(t^2/\delta) \max\{M \log(t^2/\delta), d2^{(t-1)/2}\}}{n\varepsilon},$$

where the latter case  $\widetilde{\mu}_t$  satisfies Eq. (6) as  $t \leq t^* \leq s$ . Then via a union bound this same event holds simultaneously for all  $1 \leq t \leq t^*$  with probability at least  $1 - \pi^2 \delta/3$ , and thus either  $t_{\text{stop}} > t^*$  or ADAMEAN terminates and  $\widetilde{\mu}$  satisfies the claim (6).

Proceeding to the second part of the proof, recall that  $d2^{t^*-1} \ge M^2$  and so by applying Theorem 3 with  $B = d2^{t^*-1}$ , it follows under  $\mathcal{E}_{\text{samp}}$  that  $\widetilde{\mu}_{t^*} \ne \bot$ , and thus ADAMEAN terminates after  $t_{\text{stop}} \le t^*$  iterations, with probability at least  $1 - 5\delta/(t^*)^2 \ge 1 - 5\delta$ . The claim (6) follows.

**Example 3** (Example 1 continued): In this case,  $\Sigma^{-1/2}X_i$  is  $\tau^2$ -sub-Gaussian, so  $M^2 \lesssim \tau^2(d+\log\frac{n}{\beta})$  in Assumption A1. Thus, the sample mean concentrates as  $\|\overline{X}_n - \mu\|_{\Sigma} \lesssim \tau\sqrt{(d+\log(1/\delta))/n}$  with probability at least  $1-\delta$ , and assuming  $\delta \geq e^{-d}$ , Theorem 4 then implies with probability at least  $1-O(\delta)$  over  $\widetilde{\mu} = \text{ADAMEAN}_{\varepsilon,\delta}(X_{1:n})$  that (ignoring polylogarithmic factors in n)

$$\|\widetilde{\mu} - \mu\|_{\Sigma} = \widetilde{O}\left(\tau\sqrt{\frac{d}{n}} + \frac{\tau d\log\frac{1}{\delta}}{n\varepsilon}\right).$$

This rate is, up to a factor of  $\log \frac{1}{\delta}$  and polylogarithmic factors in n, minimax-optimal for the sub-Gaussian setting (see [35] or [26, Lemma 6.7] for a lower bound on Gaussian mean estimation with known covariance matrix).  $\Diamond$ 

**Example 4** (Example 2, continued): Recall here that  $\mathbb{E}[\|X_i - \mu\|_{\Sigma}^p] \leq \tau^p d^{p/2}$  for  $p \geq 4$  and  $\tau \geq 1$ . By Theorem 4, with probability at least  $1 - 5\delta$  over  $\widetilde{\mu} = \mathtt{ADAMEAN}_{\varepsilon,\delta}(X_{1:n})$ , we have

$$\|\widetilde{\mu} - \mu\|_{\Sigma} \le \|\overline{X}_n - \mu\|_{\Sigma} + \widetilde{O}\left(\frac{\max_i \|X_i - \mu\|_{\Sigma} \sqrt{d} \log \frac{1}{\delta}}{n\varepsilon}\right)$$

so long as the empirical covariance satisfies  $\frac{1}{2}\Sigma \leq \overline{\Sigma}_n \leq \frac{3}{2}\Sigma$ . As this occurs with constant probability and  $\|\overline{X}_n - \mu\|_{\Sigma} \lesssim \sqrt{d/n}$  with constant probability, we substitute the bounds on  $\max_i \|X_i - \mu\|_{\Sigma}$  from Example 2 to obtain that with (any) constant probability,

$$\|\widetilde{\mu} - \mu\|_{\Sigma} = \widetilde{O}\left(\sqrt{\frac{d}{n}} + \frac{\tau d \log \frac{1}{\delta}}{n^{1-1/p}\varepsilon}\right).$$

By combining minimax lower bounds of Barber and Duchi [5, Proposition 4] and Steinke and Ullman [35], the best known minimax lower bound is that with constant probability,

$$\|\widetilde{\mu} - \mu\|_{\Sigma} \gtrsim \sqrt{\frac{d}{n}} + \tau \frac{d^{\frac{2p-1}{2p}} \log^{\frac{p-1}{2p}} \frac{1}{\delta}}{(n\varepsilon)^{1-1/p}}.$$

The adaptive method thus achieves optimal scaling in n, but it may be loose in  $\varepsilon$  and off by a factor of  $d^{1/2p}$  in dimension dependence.  $\Diamond$ 

# 6 Proofs for stable covariance estimation

In this section, we provide the proofs of Lemmas 4.1 and 4.2, though we begin with a collection of preliminary results that allow us to actually prove the main two lemmas. In the proofs, we refer to each execution of the while loop beginning in Line 4 of COVSAFE as an *iteration* of COVSAFE and use the transcript  $\Gamma$  as a convenient means for tracking the full execution of COVSAFE through all iterations.

## 6.1 Properties of COVSAFE

We first formalize deterministic properties about the execution of COVSAFE, giving conditions under which outputs of COVSAFE are quite stable. In the sequel, we use these to give sets to which the noise variables Z and W belong with high probability, guaranteeing stability. Recall the notation (4) that  $\widehat{\Sigma}(x,z,w)$  is the output of COVSAFE on input sample x and noise  $z \in \mathbb{R}^{n/2+1}, w \in \mathbb{R}$ , with transcript  $\Gamma = ([R_t]_{t \leq T}, [\Sigma_t]_{t \leq T}, T)$  depending implicitly on (x,z,w), and  $\widehat{\Sigma}_{-i}(x,z,w)$  is the corresponding leave-one-out covariance. We shorthand  $\widehat{\Sigma} = \widehat{\Sigma}(x,z,w)$  and  $\widehat{\Sigma}_{-i} = \widehat{\Sigma}_{-i}(x,z,w)$  and take  $\widetilde{x} = x_{1:n/2} - x_{n/2+1:n}$  as in Line 1 of COVSAFE.

Lemma 6.1 gives necessary and sufficient conditions for pruning  $\tilde{x}_i$  in iteration t+1 of COVSAFE, i.e.,  $i \in R_{t+1}$ , and Lemma 6.2 gives similar conditions for ever pruning  $\tilde{x}_i$  (that is, whether  $i \in R_T$ ).

**Lemma 6.1.** Index 
$$i \in R_{t+1}$$
 if and only if  $\log (\|\tilde{x}_i\|_{\Sigma_t}^2) + z_i + z_{n/2+1} > \log(B)$ .

*Proof.* The "if" direction is immediate from the condition for adding an element to  $R_{t+1}$  (see line 7 of COVSAFE). For the other direction, if  $i \in R_{t+1}$  then (again from the same condition) we must have for some  $s \le t$  that

$$\log\left(\|\tilde{x}_i\|_{\Sigma_s}^2\right) + z_i + z_{n/2+1} > \log(B).$$

Because  $s \leq t$ , we have  $R_s \subset R_t$  and therefore  $\Sigma_s \succeq \Sigma_t$ , this in turn implies  $\log \left( \|\tilde{x}_i\|_{\Sigma_t}^2 \right) + z_i + z_{n/2+1} > \log(B)$ .

**Lemma 6.2.** Index  $i \notin R_T$  if and only if  $\log(\|\tilde{x}_i\|_{\Sigma_T}^2) + z_i + z_{n/2+1} \leq \log(B)$ .

*Proof.* Observe that  $\Sigma_{T-1} = \Sigma_T$  because the inner while loop of COVSAFE terminates only if the algorithm prunes no observations in the previous iteration (see line 10 of COVSAFE). Then the claim follows by applying Lemma 6.1 with t = T - 1.

Finally, we may completely characterize  $\widehat{\Sigma}_{-i}$  via the removed indices  $R_{T,-i}$  and the threshold m+w, as prescribed by the lemma below.

**Lemma 6.3.**  $\widehat{\Sigma}_{-i} = \frac{1}{n} \sum_{j \notin R_{T,-i} \cup \{i\}} \widetilde{x}_j \widetilde{x}_j^T$  if and only if  $|R_T| \leq m + w$ .

*Proof.* The claim follows immediately from the return condition in Line 14 of COVSAFE, as  $\widehat{\Sigma} \neq \bot$  implies  $\widehat{\Sigma}_{-i} = \widehat{\Sigma} - \frac{1}{n} \mathbf{1} \{ i \notin R_T \} \, \widetilde{x}_i \widetilde{x}_i^T$ , where  $\widehat{\Sigma} = \Sigma_T = \frac{1}{n} \sum_{j \notin R_T} \widetilde{x}_j \widetilde{x}_j^T$  by definition.  $\square$ 

### 6.2 Proof of Lemma 4.1

We shorthand  $\widehat{\Sigma} = \widehat{\Sigma}(x, Z, w)$  and  $\widehat{\Sigma}_{-i} = \widehat{\Sigma}_{-i}(x, Z, w)$  throughout the proof. Assume that  $\widehat{\Sigma} \neq \bot$ , as otherwise the result is trivial, and recall  $Z_j \stackrel{\text{iid}}{\sim} \mathsf{Lap}(\sigma)$ . Observe if  $i \in R_T$  we have  $\widehat{\Sigma} = \widehat{\Sigma}_{-i}$  by definition; thus, we need only consider  $i \notin R_T$ . Proceeding, Lemma 6.2 gives

$$\log\left(\|\tilde{x}_i\|_{\widehat{\Sigma}}^2\right) + Z_i + Z_{n/2+1} \le \log(B)$$

for all  $i \notin R_T$ , from which it follows that  $\|\tilde{x}_i\|_{\widehat{\Sigma}}^2 \leq B\sqrt{e}$  whenever  $Z_i + Z_{n/2+1} \geq -1/2$ . We now use that  $\|\tilde{x}_i\|_{\widehat{\Sigma}}^2 \leq B\sqrt{e}$  implies that  $d_{\text{psd}}(\widehat{\Sigma}, \widehat{\Sigma}_{-i})$  is small, which follows from the following linear algebraic lemma.

**Lemma 6.4.** Let  $A \in \mathbb{R}^{d \times d}$  be positive semi-definite and  $a \in \mathbb{R}^d$  satisfy  $||a||_A^2 < 1$ . Then

$$d_{\text{psd}}(A, A - aa^T) \le \frac{1}{1 - \|a\|_A^2} \|a\|_A^2.$$

*Proof.* Define  $C = A - aa^T$  for shorthand. We first establish that Col(C) = Col(A). Because  $||a||_A^2$  is finite, it follows that  $a \in Col(A)$  and so  $Col(C) \subset Col(A)$ . On the other hand, by expanding C we have

$$C = A^{1/2} (I - A^{\dagger/2} a a^T A^{\dagger/2}) A^{1/2} \succeq (1 - ||a||_A^2) A, \tag{7}$$

thus implying that Col(C) = Col(A).

We also have from (7) that  $C^{\dagger} \stackrel{\checkmark}{\leq} \frac{1}{1-\|a\|_A^2} A^{\dagger}$ , and so

$$\left\| C^{\dagger/2} (A - C) C^{\dagger/2} \right\|_* = \left\| C^{\dagger/2} a a^T C^{\dagger/2} \right\|_* = \|a\|_C^2 \le \frac{\|a\|_A^2}{1 - \|a\|_A^2}.$$

A parallel calculation yields  $||A^{\dagger/2}(C-A)A^{\dagger/2}||_* = ||a||_A^2$ , proving the claim.

Lemma 6.4 immediately shows that  $\widehat{\Sigma}_{-i} = \widehat{\Sigma} - \frac{1}{n} \tilde{x}_i \tilde{x}_i^T \mathbf{1} \{ i \in R_T \}$  satisfies

$$d_{\mathrm{psd}}(\widehat{\Sigma}, \widehat{\Sigma}_{-i}) \leq \frac{1}{1 - B\sqrt{e}/n} \frac{B\sqrt{e}}{n}$$

whenever  $Z_i + Z_{n/2+1} \ge -\frac{1}{2}$ . To show that this occurs with high probability, we use the following result, which follows from the observation that if  $c \ge 0$ , then for any independent variables X, Y we have  $\mathbb{P}(X + Y > c) \le \mathbb{P}(X > c/2) + \mathbb{P}(Y > c/2)$  by a union bound:

**Observation 6.1.** Let  $X, Y \stackrel{\text{iid}}{\sim} \mathsf{Lap}(\sigma)$  and  $c \geq 0$ . Then  $\mathbb{P}(X + Y > c) \leq \exp(-\frac{c}{2\sigma})$ .

We see that  $\mathbb{P}(Z_i + Z_{n/2+1} < -\frac{1}{2}) \le \exp(-\frac{1}{4\sigma})$  as claimed.

## 6.3 Proof of Lemma 4.2

The proof of Lemma 4.2 comes in four steps. The crux of the proof is a coupling argument where, via the running assumption  $\tilde{x}_{-i} = \tilde{x}'_{-i}$ , we equate the execution of COVSAFE on x to that on x' by perturbing Z in a careful way that changes the distribution of Z little. Step one in this approach, which we provide in Lemma 6.5, is a deterministic lemma relating the collections R and R' of indices COVSAFE removes on adjacent inputs x and x' via the noise input values z. In the second and third steps, which consist of Lemmas 6.6 and 6.7 respectively, we construct a map  $\pi: \mathbb{R}^{n/2+1} \to \mathbb{R}^{n/2+1}$  such that Z and  $\pi(Z)$  have similar distributions and for which  $\hat{\Sigma}_{-i}(x,z,W)$  and  $\hat{\Sigma}_{-i}(x',\pi(z),W)$  (recall the definition (4b)) likewise have similar distributions for all z, where we use the randomness in W for the latter distributional approximation. Lemma 6.6 relates the distributions of the removed indices  $R_{T,-i}$ , while Lemma 6.7 relates the probabilities that COVSAFE aborts and returns  $\bot$ . In Sec. 6.3.1, we finally synthesize the intermediate lemmas to give the proof of Lemma 4.2.

Our first step is the deterministic lemma relating the collections of removed indices.

**Lemma 6.5.** Let  $z, z' \in \mathbb{R}^{n/2+1}$  and  $w \in \mathbb{R}$ , and let

$$\begin{split} \widehat{\Sigma}, \left( [R_t]_{t=0}^T, [\Sigma_t]_{t=0}^T, T \right) &:= \mathtt{COVSAFE}_{B,m}(x; z, w), \\ \widehat{\Sigma}', \left( [R_t']_{t=0}^{T'}, [\Sigma_t']_{t=0}^{T'}, T' \right) &:= \mathtt{COVSAFE}_{B,m}(x'; z', w). \end{split}$$

Assume  $\tilde{x}'_i = 0$ . The following hold.

(a) If 
$$z'_j + z'_{n/2+1} \ge z_j + z_{n/2+1}$$
 for all  $j \in R_{T,-i}$ , then  $R_{T,-i} \subset R'_{T',-i}$ .

(b) If 
$$z'_j + z'_{n/2+1} \ge z_j + z_{n/2+1}$$
 for all  $j \notin R'_{T',-i}$ , then  $R_{T,-i} \subset R'_{T',-i}$ .

Additionally assume that  $n \geq 2B\sqrt{e}$  and  $z_i + z_{n/2+1} \geq -1/2$ . Then the following hold.

(c) If 
$$z'_j + z'_{n/2+1} \le z_j + z_{n/2+1} - 2B\sqrt{e}/n$$
 for all  $j \in R'_{T',-i}$ , then  $R'_{T',-i} \subset R_{T,-i}$ .

(d) If 
$$z'_j + z'_{n/2+1} \le z_j + z_{n/2+1} - 2B\sqrt{e}/n$$
 for all  $j \notin R_{T,-i}$ , then  $R'_{T',-i} \subset R_{T,-i}$ .

*Proof.* We prove each claim by induction over  $t \in \mathbb{N}$ .

Proceeding with the first claim (a), observe trivially that  $R_{0,-i} = \emptyset \subset R'_{T',-i}$ . Now suppose for the sake of induction that  $R_{t,-i} \subset R'_{T',-i}$ . If t = T then there is nothing to show, so we take t < T. Our assumption that  $\tilde{x}_{-i} = \tilde{x}'_{-i}$  and  $\tilde{x}'_{i} = 0$  implies  $\Sigma_{t} \succeq \Sigma'_{T'}$ . Thus, for all  $j \in R_{t,-i} \subset R_{T,-i}$  we have

$$\log(B) \stackrel{(i)}{<} \log\left(\|\tilde{x}_{j}\|_{\Sigma_{t}}^{2}\right) + z_{j} + z_{n/2+1} \stackrel{(ii)}{\leq} \log\left(\|\tilde{x}_{j}\|_{\Sigma'_{T'}}^{2}\right) + z_{j} + z_{n/2+1}$$

$$\stackrel{(iii)}{\leq} \log\left(\|\tilde{x}_{j}\|_{\Sigma'_{T'}}^{2}\right) + z'_{j} + z'_{n/2+1},$$

where inequality (i) follows from Lemma 6.1 (applied with noise z), inequality (ii) because  $\Sigma_t \succeq \Sigma'_{T'}$ , and inequality (iii) is by assumption in case (a). Lemma 6.2 (with data x' and noise z') then gives that  $j \in R'_{T',-i}$  if and only if  $\log B < \log(\|\tilde{x}_j\|^2_{\Sigma'_{T'}}) + z'_j + z'_{n/2+1}$ , and as  $j \in R_{t,-i}$  was arbitrary we have  $R_{t+1,-i} \subset R'_{T',-i}$ . This completes the inductive step and so  $R_{t,-i} \subset R'_{T',-i}$  for all  $t \leq T$  and the first claim holds.

The proof of claim (b) relies on a similar inductive argument as that for the first claim (a). Equivalent to the inclusion  $R_{T,-i} \subset R'_{T',-i}$  is that if  $j \notin R'_{T',-i}$ , then  $j \notin R_{T,-i}$ . Consider

 $j \notin R'_{T',-i}$ , and begin with the inductive assumption that  $R_{t,-i} \subset R'_{T',-i}$ ; it suffices to show that  $j \notin R_{t+1,-i}$ . Because  $\Sigma_t \succeq \Sigma'_{T'}$  by construction of  $\Sigma_t$ , we obtain

$$\log\left(\|\tilde{x}_{j}\|_{\Sigma_{t}}^{2}\right) + z_{j} + z_{n/2+1} \leq \log\left(\|\tilde{x}_{j}\|_{\Sigma_{T'}'}^{2}\right) + z_{j} + z_{n/2+1}$$

$$\stackrel{(i)}{\leq} \log\left(\|\tilde{x}_{j}\|_{\Sigma_{T'}'}^{2}\right) + z_{j}' + z_{n/2+1}' \leq \log(B),$$

where step (i) is by the assumption that  $z_j + z_{n/2+1} \le z'_j + z'_{n/2+1}$  in part (b) and the final inequality is Lemma 6.2. Applying Lemma 6.1 with the inequality  $\log \left(\|\tilde{x}_j\|_{\Sigma_t}^2\right) + z_j + z_{n/2+1} \le \log B$  then guarantees that  $j \notin R_{t+1,-i}$  as desired, completing the proof of claim (b).

For the proof of claim (c), we induct on  $R'_{t,-i}$  for  $t \leq T'$  and must account for the possibility that  $\Sigma_T \npreceq \Sigma'_t$  even if  $R'_{t,-i} \subset R_{T,-i}$ , because  $\Sigma_T$  may include the term  $\tilde{x}_i \tilde{x}_i^T$  (i.e.,  $i \notin R_T$ ). The base case for t=0 is trivial, so assume that  $R'_{t,-i} \subset R_{T,-i}$ . If  $i \notin R_T$ , then Lemma 6.2 and the standing assumption that  $z_i + z_{n/2+1} \geq -\frac{1}{2}$  guarantee that

$$\log\left(\|\tilde{x}_j\|_{\Sigma_T}^2\right) \le \log B - z_i - z_{n/2+1} \le \log B + \frac{1}{2} = \log(B\sqrt{e}),$$

i.e.,  $\|\tilde{x}_j\|_{\Sigma_T}^2 \leq B\sqrt{e}$ . We require the following technical observation about positive definite matrices, whose proof we temporarily defer.

**Observation 6.2.** Let  $\Sigma \in \mathbb{R}^{d \times d}$  be positive semi-definite,  $\alpha \geq 0$ , and  $u \in \mathbb{R}^d$ . Define  $\Sigma' := \Sigma - \alpha u u^T$ . If  $||u||_{\Sigma}^2 \leq \frac{1}{2\alpha}$ , then  $\Sigma' \succeq \frac{1}{2}\Sigma$  and for any  $v \in \mathsf{Col}(\Sigma)$ ,

$$\left|\log\left(\left\|v\right\|_{\Sigma}^{2}\right) - \log\left(\left\|v\right\|_{\Sigma'}^{2}\right)\right| \leq 2\alpha \left\|u\right\|_{\Sigma}^{2}.$$

As  $\|\tilde{x}_i\|_{\Sigma_T}^2 \leq B\sqrt{e}$ , Observation 6.2 applies with  $u = \tilde{x}_i$  and  $\alpha = \frac{1}{n}$  when  $n \geq 2B\sqrt{e}$ , and thus

$$\log\left(\left\|v\right\|_{\Sigma_{T}-\frac{1}{n}1\left\{i\notin R_{T}\right\}\tilde{x}_{i}\tilde{x}_{i}^{T}\right) \leq \log\left(\left\|v\right\|_{\Sigma_{T}}^{2}\right) + \frac{2B\sqrt{e}}{n}$$

$$\tag{8}$$

for all v, and in particular, for  $v = \tilde{x}_j$  for each  $j \in [n/2]$ . On the other hand, regardless of whether  $i \in R_T$ , the inductive assumption that  $R'_{t,-i} \subset R_{T,-i}$  guarantees that

$$\Sigma_t' \succeq \Sigma_T - \frac{1}{n} 1\{i \notin R_T\} \, \tilde{x}_i \tilde{x}_i^T. \tag{9}$$

Considering  $j \in R'_{t+1,-i}$ , then, Lemma 6.1 implies

$$\log B < \log(\|\tilde{x}_{j}\|_{\Sigma'_{t}}^{2}) + z'_{j} + z'_{n/2+1}$$

$$\stackrel{(i)}{\leq} \log\left(\|\tilde{x}_{j}\|_{\Sigma_{T} - \frac{1}{n} \mathbb{I}\{i \notin R_{T}\}\tilde{x}_{i}\tilde{x}_{i}^{T}}\right) + z'_{j} + z'_{n/2+1}$$

$$\stackrel{(ii)}{\leq} \log\left(\|\tilde{x}_{j}\|_{\Sigma_{T}}^{2}\right) + \frac{2B\sqrt{e}}{n} + z'_{j} + z'_{n/2+1}$$

$$\stackrel{(iii)}{\leq} \log\left(\|\tilde{x}_{j}\|_{\Sigma_{T}}^{2}\right) + z_{j} + z_{n/2+1}.$$

Here inequality (i) follows from the ordering relation (9); inequality (ii) holds because if  $i \in R_T$ , then  $\Sigma_T = \Sigma_T - \frac{1}{n} \mathbb{1}\{i \notin R_T\} \tilde{x}_i \tilde{x}_i^T$  and if  $i \notin R_T$  then inequality (8) holds; the final inequality (iii) follows by assumption under claim (c). This gives the induction that  $R'_{t+1,-i} \subset R_{T,-i}$ , as Lemma 6.1 shows that  $j \in R_{T,-i}$ .

Claim (d) follows from an essentially identical induction argument,  $mutatis\ mutandis$ , as that for claim (c).

**Proof of Observation 6.2.** We finally return to prove the claimed observation. That  $\Sigma' \succeq \frac{1}{2}\Sigma$  follows by observing that  $u \in \mathsf{Col}(\Sigma)$  and hence

$$\Sigma - \alpha u u^T = \Sigma^{1/2} \underbrace{(I - \alpha \Sigma^{\dagger/2} u u^T \Sigma^{\dagger/2})}_{\succeq (1/2)I} \Sigma^{1/2} \succeq \frac{1}{2} \Sigma.$$

This also implies that  $Col(\Sigma') = Col(\Sigma)$ .

To prove the remainder of the lemma, it suffices to show for  $v \in \mathsf{Col}(\Sigma)$  that  $\log(\|v\|_{\Sigma'}^2) \le \log(\|v\|_{\Sigma}^2) + 2\alpha \|u\|_{\Sigma}^2$ , since the other direction is immediate from  $\Sigma \succeq \Sigma'$ . Observe that

$$(\Sigma - \alpha u u^T)^{\dagger} = \Sigma^{\dagger/2} (I - \alpha \Sigma^{\dagger/2} u u^T \Sigma^{\dagger/2})^{-1} \Sigma^{\dagger/2}.$$

By the inequality  $I - \alpha \Sigma^{\dagger/2} u u^T \Sigma^{\dagger/2} \succeq (1 - \alpha \|u\|_{\Sigma}^2) I$  we have

$$(I - \alpha \Sigma^{\dagger/2} u u^T \Sigma^{\dagger/2})^{-1} \leq (1 - \alpha \|u\|_{\Sigma}^2)^{-1} I \leq (1 + 2\alpha \|u\|_{\Sigma}^2) I,$$

where the final inequality follows from the assumption that  $||u||_{\Sigma}^2 \leq \frac{1}{2\alpha}$ . Combining this with the preceding display implies that

$$\Sigma'^{\dagger} \leq (1 + 2\alpha \|u\|_{\Sigma}^{2}) \Sigma^{\dagger}$$

and so

$$\log \left( \|v\|_{\Sigma'}^2 \right) \le \log \left( (1 + 2\alpha \|u\|_{\Sigma}^2) \|v\|_{\Sigma}^2 \right) \le \log \left( \|v\|_{\Sigma}^2 \right) + 2\alpha \|u\|_{\Sigma}^2$$

as desired.  $\Box$ 

We move to the second step we outline at the beginning of this section, which relates the distributions of removed indices  $R_T$  in the execution of COVSAFE on adjacent inputs x and x'. The key idea is to construct a deterministic map  $\pi$  so that the execution of COVSAFE on input x with noise z and that on x' with noise  $\pi(z)$  is similar—leveraging Lemma 6.5—and to show that the distributions of  $\pi(Z)$  and Z are similar. Lemma 6.3 shows that the set of outlier indices  $R_{T,-i}$  completely determines  $\widehat{\Sigma}_{-i}$  except in the case that  $\widehat{\Sigma}_{-i} = \bot$ , which occurs with high probability if  $|R_{T,-i}|$  is large, so the next lemma controls the distribution of the sets of removed indices. To state the lemma, we require a few events whose probabilities we can control. Recalling that  $Z_i \stackrel{\text{iid}}{\sim} \mathsf{Lap}(\sigma_Z)$ , define

$$\mathcal{E}_{\text{prune}} := \{ z \in \mathbb{R}^{n/2+1} \mid z_j + z_{n/2+1} \ge -1/2 \text{ for all } j \in [n/2] \}.$$
 (10)

To set notation for the remainder of the proof, we shorthand the definition (4a) as

$$\widehat{\Sigma}, ([R_t]_{t=0}^T, [\Sigma_t]_{t=0}^T, T) := \text{COVSAFE}_{B,m}(x; Z, W)$$

$$\widehat{\Sigma}', ([R_t']_{t=0}^{T'}, [\Sigma_t']_{t=0}^{T'}, T') := \text{COVSAFE}_{B,m}(x'; Z, W),$$
(11)

and the definition (4b) as

$$\widehat{\Sigma}_{-i} = \widehat{\Sigma} - \frac{1}{n} \mathbb{1} \{ i \notin R_T \} \, \widetilde{x}_i \widetilde{x}_i^T \text{ and } \widehat{\Sigma}'_{-i} = \widehat{\Sigma}' - \frac{1}{n} \mathbb{1} \{ i \notin R_T \} \, \widetilde{x}_i' \widetilde{x}_i'^T,$$

where  $\perp + v = \perp$  for any vector v.

We have the following distributional guarantee on the removed indices regardless of W.

**Lemma 6.6.** Let  $S \subset [n/2] \setminus \{i\}$  and define  $\alpha = \frac{2\sqrt{e}B(|S|+1)}{n\sigma_Z}$ . If  $\tilde{x}'_i = 0$ , then

(a) 
$$\mathbb{P}(R_{T,-i} = S, Z \in \mathcal{E}_{\text{prune}}) \leq \exp(\alpha) \mathbb{P}(R'_{T',-i} = S).$$

(b) 
$$\mathbb{P}(R'_{T'-i} = S, Z \in \mathcal{E}_{prune}) \le \exp(\alpha) \mathbb{P}(R_{T,-i} = S).$$

*Proof.* The input noise Z completely determines  $R_T$  and  $R'_{T'}$  in COVSAFE (see the **while** loop constructing  $R_t$  in lines 4–12). Consequently, we may define sets of input noise Z yielding a given set of removed indices, letting

$$\mathcal{Z}(S) := \{ z \in \mathbb{R}^{n/2+1} \mid R_{T,-i} = S \text{ for } Z = z \}$$
  
$$\mathcal{Z}'(S) := \{ z \in \mathbb{R}^{n/2+1} \mid R'_{T',-i} = S \text{ for } Z = z \},$$

so  $Z \in \mathcal{Z}(S)$  is equivalent to  $R_{T,-i} = S$ . It suffices to show that

$$\mathbb{P}\left(Z \in \mathcal{Z}(S) \cap \mathcal{E}_{\text{prune}}\right) \leq e^{\alpha} \mathbb{P}\left(Z \in \mathcal{Z}'(S)\right) \text{ and} 
\mathbb{P}\left(Z \in \mathcal{Z}'(S) \cap \mathcal{E}_{\text{prune}}\right) \leq e^{\alpha} \mathbb{P}\left(Z \in \mathcal{Z}(S)\right), \tag{12}$$

as claim (a) follows via the first bound and claim (b) the second.

Proceeding with the first bound, define  $\eta \in \mathbb{R}^{n/2+1}$  and  $\pi : \mathbb{R}^{n/2+1} \to \mathbb{R}^{n/2+1}$  by

$$\eta_j := \begin{cases} 2B\sqrt{e}/n & j \in S \\ -2B\sqrt{e}/n & j = n/2 + 1 & \pi(z) := z + \eta. \\ 0 & \text{otherwise,} \end{cases}$$

The deterministic removals Lemma 6.5 shows that on  $Z \in \mathcal{E}_{prune}$ , if we let z = Z and  $z' = \pi(z)$  so that  $z_j + z_{n/2+1} \ge z_j' + z_{n/2+1}'$  for  $j \in S$ , then parts (a) and (d) of Lemma 6.5 give

$$\pi(\mathcal{Z}(S) \cap \mathcal{E}_{\text{prune}}) \subset \mathcal{Z}'(S).$$

The first bound in (12) then follows by the standard Laplacian ratio bounds in Lemma 2.7. Indeed, we have  $Z_j \stackrel{\text{iid}}{\sim} \mathsf{Lap}(\sigma_Z) = \mathsf{Lap}(\|\eta\|_1 \frac{\sigma_Z}{\|\eta\|_1})$  and  $\|\eta\|_1 = \frac{2\sqrt{e}B(|S|+1)}{n}$ . Then setting  $\beta = \|\eta\|_1$  yields  $\beta/\sigma_Z \leq \alpha$ , so we can apply Lemma 2.7 to obtain the claimed bound (12) via

$$\mathbb{P}(Z \in \mathcal{Z}(S) \cap \mathcal{E}_{\text{prune}}) \leq e^{\alpha} \mathbb{P}(Z \in \pi(\mathcal{Z}(S) \cap \mathcal{E}_{\text{prune}})) \leq e^{\alpha} \mathbb{P}(Z \in \mathcal{Z}'(S)).$$

The proof of the second bound (12) is essentially the same, only this time we let

Then the event  $Z \in \mathcal{E}_{\text{prune}}$  implies  $\pi(Z) \in \mathcal{E}_{\text{prune}}$ , as  $\pi(Z)_j + \pi(Z)_{n/2+1} \ge Z_j + Z_{n/2+1}$  for all  $j \in [n/2]$ . We may thus appeal to cases (b) and (c) of Lemma 6.5 with the settings  $z = \pi(Z)$ , z' = Z and  $R'_{T',-n} = S$  and proceed with the same argument as above.

The third step we outline at the beginning of the proof of Lemma 4.2 is to relate the probabilities that COVSAFE aborts on neighboring inputs x and x'. Recall  $\widehat{\Sigma}$  and  $\widehat{\Sigma}'$  are the covariances COVSAFE outputs on inputs x and x', respectively, as in definition (11).

**Lemma 6.7.** Let  $\sigma_Z, \sigma_W > 0$ ,  $Z_j \stackrel{\text{iid}}{\sim} \mathsf{Lap}(\sigma_Z)$ , and  $W \sim \mathsf{Lap}(\sigma_W)$  in definition (11). If  $\tilde{x}_{-i} = \tilde{x}'_{-i}$  and  $\tilde{x}'_i = 0$ , then

(a) 
$$\mathbb{P}(\widehat{\Sigma} = \bot) \le \exp(\frac{1}{\sigma_W}) \mathbb{P}(\widehat{\Sigma}' = \bot)$$

(b) 
$$\mathbb{P}(\widehat{\Sigma}' = \bot, Z \in \mathcal{E}_{\text{prune}}) \leq \exp(\frac{2\sqrt{eB}}{n\sigma_Z})\mathbb{P}(\widehat{\Sigma} = \bot).$$

*Proof.* Let f denote the density of W, so that  $f(w) = \frac{1}{2\sigma_{W_{\text{cov}}}} \exp(-|w|/\sigma_{W_{\text{cov}}})$ , and thus  $|\log \frac{f(w)}{f(w-1)}| \leq \frac{1}{\sigma_{W_{\text{cov}}}}$  for all w, and recall the threshold  $m \in \mathbb{N}$  in line 14 of COVSAFE. Proceeding with the first claim of the lemma, we have the following sequence of inequalities:

$$\mathbb{P}\left(\widehat{\Sigma} = \bot\right) \stackrel{(i)}{\leq} \int \mathbb{P}\left(|R_T| > m + w\right) f(w) dw \leq \int \mathbb{P}\left(|R_{T,-i}| > m + w - 1\right) f(w) dw 
\stackrel{(ii)}{\leq} \int \mathbb{P}\left(|R'_{T',-i}| > m + w - 1\right) f(w) dw 
\leq \exp(1/\sigma_{W_{\text{cov}}}) \int \mathbb{P}\left(|R'_{T',-i}| > m + w - 1\right) f(w - 1) dw 
\stackrel{(iii)}{=} \exp(1/\sigma_{W_{\text{cov}}}) \mathbb{P}(\widehat{\Sigma}' = \bot).$$

Here, step (i) follows from Lemma 6.3 that  $\widehat{\Sigma}' = \bot$  if and only if  $|R'_T| > m+w$ . Step (ii) follows from the coupling argument in Lemma 6.5 part (a): because the noise Z is identical in both executions of COVSAFE(x; Z, W) and COVSAFE(x'; Z, W), we have  $R_{T,-i} \subset R'_{T',-i}$ . Step (iii) applies because of Lemma 6.3 again, as the assumption  $\widetilde{x}'_i = 0$  guarantees  $R'_{T',-i} = R'_{T'}$  (recall the rejection threshold in Line 7). Claim (a) follows.

For the claim (b), again applying Lemma 6.3 we have

$$\mathbb{P}\left(\widehat{\Sigma}' = \bot, Z \in \mathcal{E}_{\text{prune}}\right) = \int \mathbb{P}\left(\left|R'_{T',-i}\right| > m + w, Z \in \mathcal{E}_{\text{prune}}\right) f(w) dw$$

and

$$\int \mathbb{P}\left(|R_{T,-i}| > m + w\right) f(w) dw \leq \mathbb{P}\left(\widehat{\Sigma} = \bot\right).$$

Combining these displays, it thus suffices to show for all w that

$$\mathbb{P}\left(\left|R'_{T',-i}\right| > m + w, Z \in \mathcal{E}_{\text{prune}}\right) \le \exp\left(\frac{2\sqrt{e}B}{n\sigma_Z}\right) \mathbb{P}\left(\left|R_{T,-i}\right| > m + w\right). \tag{13}$$

To this end, we adopt a similar tack as in the proof of Lemma 6.6, defining

$$\mathcal{Z}(w) := \{ z \in \mathbb{R}^{n/2+1} \mid |R_{T,-i}| > m + w \text{ for } Z = z \}$$
  
$$\mathcal{Z}'(w) := \{ z \in \mathbb{R}^{n/2+1} \mid |R'_{T',-i}| > m + w \text{ for } Z = z \}$$

and the single coordinate perturbation  $\pi(z):=z+\eta$  for  $\eta\in\mathbb{R}^{n/2+1}$  the vector with all zeros except that  $\eta_{n/2+1}=\frac{2\sqrt{e}B}{n}$ . Similar to our proof of Lemma 6.6, the mapping  $\pi$  guarantees that  $z'=\pi(z)$  satisfies  $z'_j+z'_{n/2+1}\leq z_j+z_{n/2+1}-\frac{2B\sqrt{e}}{n}$  for all j, which is precisely the condition for case (c) of Lemma 6.5, and so  $R'_{T',-i}\subset R_{T,-i}$  irrespective of  $R'_{T',-i}$  and  $R_{T,-i}$ . It thus holds simultaneously for all  $w\in\mathbb{R}$  that

$$\pi\left(\mathcal{Z}'(w)\cap\mathcal{E}_{\text{prune}}\right)\subset\mathcal{Z}(w).$$

Noting that  $\|\eta\|_1 = \frac{2\sqrt{e}B}{n}$ , Lemma 2.7 on likelihood ratios for Laplace random variables then guarantees

$$\mathbb{P}(Z \in \mathcal{Z}'(w) \cap \mathcal{E}_{\text{prune}}) \leq \exp\left(\frac{2\sqrt{e}B}{n\sigma_Z}\right) \mathbb{P}(Z \in \pi(\mathcal{Z}'(w) \cap \mathcal{E}_{\text{prune}})) \leq \exp\left(\frac{2\sqrt{e}B}{n\sigma_Z}\right) \mathbb{P}(Z \in \mathcal{Z}(w)),$$
 which is equivalent to inequality (13).

### 6.3.1 Finalizing proof of Lemma 4.2

By combining Lemmas 6.6 and 6.7, we can prove the stability of COVSAFE. Recall the set  $\mathcal{E}_{\text{prune}}$  in (10) and that  $W \sim \mathsf{Lap}(\sigma_{W_{\text{cov}}})$ , and additionally define

$$\mathcal{E}_{ ext{thr}} := \left(-\infty, \sigma_{W_{ ext{cov}}} \log \frac{1}{\gamma}\right].$$

The key part of our argument is to show that when x and x' are adjacent but  $\tilde{x}'_i = 0$ , if the noise variables Z, W satisfy  $Z \in \mathcal{E}_{\text{prune}}$  and  $W \in \mathcal{E}_{\text{thr}}$ , then for  $\widehat{\Sigma}$  and  $\widehat{\Sigma}'$  as in the call (11) the leave-one-out covariances  $\widehat{\Sigma}_{-i}$  and  $\widehat{\Sigma}'_{-i}$  are similar. We then bound the probabilities of the individual events and use a group composition argument to give the lemma for arbitrary  $\tilde{x}'_i$ .

With this in mind, let  $A \in \mathbb{R}^{d \times d}$  and note that for any fixed sample x,  $\widehat{\Sigma}$  and  $\widehat{\Sigma}_{-i}$  can take only finitely many values. For

$$\alpha = \frac{1}{\sigma_{W_{\text{cov}}}} + \frac{2\sqrt{e}B(m_{\text{max}} + 1)}{n\sigma_Z},$$

we show for any  $A \in \mathbb{R}^{d \times d} \cup \{\bot\}$  that

$$\mathbb{P}\left(\widehat{\Sigma}_{-i} = A, Z \in \mathcal{E}_{\text{prune}}, W \in \mathcal{E}_{\text{thr}}\right) \le \exp(\alpha) \mathbb{P}\left(\widehat{\Sigma}'_{-i} = A\right) \text{ and }$$
(14)

$$\mathbb{P}\left(\widehat{\Sigma}'_{-i} = A, Z \in \mathcal{E}_{\text{prune}}, W \in \mathcal{E}_{\text{thr}}\right) \le \exp(\alpha) \mathbb{P}\left(\widehat{\Sigma}_{-i} = A\right). \tag{15}$$

Lemma 6.7 already implies both inequalities (14) and (15) hold for  $A = \bot$ , so all that remains is to show the same for  $A \in \mathbb{R}^{d \times d}$ .

Proceeding first with inequality (14), let  $f(w) = \frac{1}{2\sigma_{W_{\text{cov}}}} \exp(-|w|/\sigma_{W_{\text{cov}}})$  be the density of W and  $S(A) := \{S \subset [n/2] \setminus \{i\} \mid A = \frac{1}{n} \sum_{j \notin S \cup \{i\}} \tilde{x}_j \tilde{x}_j^T \}$ . Marginalizing over W gives

$$\mathbb{P}\left(\widehat{\Sigma}_{-i} = A, Z \in \mathcal{E}_{\text{prune}}, W \in \mathcal{E}_{\text{thr}}\right)$$

$$\stackrel{(i)}{=} \int_{\mathcal{E}_{\text{thr}}} \mathbb{P}\left(R_{T,-i} \in S(A), |R_T| \leq m + w, Z \in \mathcal{E}_{\text{prune}}\right) f(w) dw$$

$$\leq \int \mathbb{P}\left(R_{T,-i} \in S(A), |R_{T,-i}| \leq \min\{m + w, m_{\text{max}}\}, Z \in \mathcal{E}_{\text{prune}}\right) f(w) dw,$$

where step (i) follows from the condition that  $|R_T| \leq m + w$  if and only if  $\widehat{\Sigma}_{-i} \neq \bot$  from Lemma 6.3, and the final inequality follows because  $\mathcal{E}_{\text{thr}} = \{w \mid w \leq m_{\text{max}}\}$ . Continuing, note for each  $S \in S(A)$ , we can have  $S = R_{T,-i}$  with  $|R_{T,-i}| \leq \min\{m + w, m_{\text{max}}\}$  only if  $|S| \leq \min\{m + w, m_{\text{max}}\} \leq m_{\text{max}}$ , so that by case (a) of Lemma 6.6

$$\mathbb{P}(R_{T,-i} \in S(A), |R_{T,-i}| \le \min\{m+w, m_{\max}\}, Z \in \mathcal{E}_{\text{prune}})$$

$$\le \exp\left(\frac{2\sqrt{e}B(m_{\max}+1)}{n\sigma_Z}\right) \mathbb{P}\left(R'_{T',-i} \in S(A), \left|R'_{T',-i}\right| \le m+w\right).$$

Returning to the integral above, we obtain inequality (14) by integrating and applying Lemma 6.3:

$$\int \mathbb{P}\left(R'_{T',-i} \in S(A), \left| R'_{T',-i} \right| \le m + w\right) f(w) dw = \mathbb{P}(\widehat{\Sigma}'_{-i} = A)$$

as  $R'_{T',-i} = R'_{T'}$  because  $\tilde{x}'_i = 0$  by assumption.

The proof of inequality (15) is essentially the same, only now we must take additional care to account for the possibility that  $i \in R_T$ . As in the preceding integral inequalities, Lemma 6.3 gives

$$\begin{split} & \mathbb{P}\left(\widehat{\Sigma}'_{-i} = A, Z \in \mathcal{E}_{\text{prune}}, W \in \mathcal{E}_{\text{thr}}\right) \\ & \leq \int_{\mathcal{E}_{\text{thr}}} \mathbb{P}\left(R'_{T', -i} \in S(A), \left|R'_{T', -i}\right| \leq \min\{m + w, m_{\text{max}}\}, Z \in \mathcal{E}_{\text{prune}}\right) f(w) dw. \end{split}$$

In this case, with reasoning identical to that above, we apply case (b) of Lemma 6.6 to achieve

$$\mathbb{P}\left(R'_{T',-i} \in S(A), \left| R'_{T',-i} \right| \leq \min\{m+w, m_{\max}\}, Z \in \mathcal{E}_{\text{prune}}\right) \\
\leq \exp\left(\frac{2\sqrt{e}B(m_{\max}+1)}{n\sigma_Z}\right) \mathbb{P}\left(R_{T,-i} \in S(A), |R_{T,-i}| \leq m+w\right),$$

so

$$\mathbb{P}\left(\widehat{\Sigma}'_{-i} = A, Z \in \mathcal{E}_{\text{prune}}, W \in \mathcal{E}_{\text{thr}}\right)$$

$$\leq \exp\left(\frac{2\sqrt{e}B(m_{\text{max}} + 1)}{n\sigma_Z}\right) \int \mathbb{P}\left(R_{T, -i} \in S(A), |R_{T, -i}| \leq m + w\right) f(w)dw.$$

We upper bound the final integral by noting that

$$\mathbb{P}\left(\widehat{\Sigma}_{-i} = A\right) \stackrel{(\star)}{=} \int \mathbb{P}\left(R_{T,-i} \in S(A), |R_T| \le m + w\right) f(w) dw$$
$$\ge \int \mathbb{P}\left(R_{T,-i} \in S(A), |R_{T,-i}| \le m + w - 1\right) f(w) dw,$$

where  $(\star)$  follows from Lemma 6.3, and then using  $|\log \frac{f(w)}{f(w+1)}| \leq \frac{1}{\sigma_{W_{\text{COV}}}}$  for all w to see that

$$\int \mathbb{P}\left(R_{T,-i} \in S(A), |R_{T,-i}| \le m+w\right) f(w) dw \le \exp\left(\frac{1}{\sigma_{W_{\text{cov}}}}\right) \mathbb{P}\left(\widehat{\Sigma}_{-i} = A\right),$$

which gives inequality (15) once we substitute  $\alpha = \frac{1}{\sigma_{W_{cov}}} + \frac{2\sqrt{e}B(m_{max}+1)}{n\sigma_Z}$ .

We combine inequalities (14) and (15) to get Lemma 4.2. For any set  $C \subset \mathbb{R}^{d \times d} \cup \{\bot\}$ ,

$$\mathbb{P}(\widehat{\Sigma}_{-i} \in C) \leq \mathbb{P}(\widehat{\Sigma}_{-i} \in C, Z \in \mathcal{E}_{\text{prune}}, W \in \mathcal{E}_{\text{thr}}) + \mathbb{P}(Z \notin \mathcal{E}_{\text{prune}}) + \mathbb{P}(W \notin \mathcal{E}_{\text{thr}})$$
$$\leq e^{\alpha} \mathbb{P}(\widehat{\Sigma}'_{-i} \in C) + \mathbb{P}(Z \notin \mathcal{E}_{\text{prune}}) + \mathbb{P}(W \notin \mathcal{E}_{\text{thr}})$$

by inequality (14). We then have the immediate bounds  $\mathbb{P}(W \notin \mathcal{E}_{thr}) = \mathbb{P}(W > \sigma_{W_{cov}} \log \frac{1}{\gamma}) = \frac{1}{2} \exp(-\log \frac{1}{\gamma}) = \frac{\gamma}{2}$ . Similarly,  $\mathbb{P}(Z \notin \mathcal{E}_{prune}) \leq \frac{n}{2} \exp(-\frac{1}{4\sigma_Z})$  by Observation 6.1. The upper bound on  $\mathbb{P}(\widehat{\Sigma}'_{-i} \in C)$  is similar but uses inequality (15).

To this point, we have shown that if x and x'' are adjacent samples differing only in that the difference  $\tilde{x}_i = x_i - x_{n/2+i}$  may be non-zero while  $\tilde{x}_i'' = x_i'' - x_{n/2+i}'' = 0$ , then returning to the notation (4) and identifying  $\hat{\Sigma}_{-i} = \hat{\Sigma}_{-i}(x, Z, W)$  and  $\hat{\Sigma}'_{-i} = \hat{\Sigma}_{-i}(x'', Z, W)$ ,

$$\widehat{\Sigma}_{-i}(x, Z, W) \stackrel{d}{=}_{\alpha, \beta} \widehat{\Sigma}_{-i}(x'', Z, W)$$

for  $\alpha = \frac{1}{\sigma_{W_{\text{cov}}}} + \frac{2\sqrt{e}B(m_{\text{max}}+1)}{n\sigma_Z}$  and  $\beta = \frac{\gamma}{2} + \frac{n}{2}\exp(-\frac{1}{4\sigma_Z})$ . Thus we obtain that if x' is any sample satisfying  $x'_{-i} = x_{-i}$ ,

$$\widehat{\Sigma}_{-i}(x, Z, W) \stackrel{d}{=}_{\alpha, \beta} \widehat{\Sigma}_{-i}(x'', Z, W) \stackrel{d}{=}_{\alpha, \beta} \widehat{\Sigma}_{-i}(x', Z, W).$$

Using group composition (Lemma 2.2), we obtain

$$\widehat{\Sigma}_{-i}(x, Z, W) \stackrel{d}{=}_{2\alpha, \beta + e^{\alpha}\beta} \widehat{\Sigma}_{-i}(x', Z, W),$$

which is the desired Lemma 4.2.

#### Proofs for mean estimation 7

In this section, we provide the proofs of Lemmas 4.3 and 4.4. Throughout, we differentiate outputs of MEANSAFE on inputs x versus x' (or A versus A') via tick marks, so that (for example)  $\widehat{\mu}$  corresponds to the mean in Line 9 of MEANSAFE on input sample x, or  $D'_i$  corresponds to the log-diameter in Line 2 of MEANSAFE on input sample x'. We will make this precise using the function  $\Gamma(x,A)$  from (3), which is the transcript MEANSAFE outputs on input x,A.

#### 7.1Proof of Lemma 4.3

We shorthand  $\widetilde{\mu}(x,A)$  and  $\widetilde{\mu}(x',A)$  as  $\widetilde{\mu}$  and  $\widetilde{\mu}'$  respectively, and unpack the corresponding execution transcripts:

$$(D,\widetilde{D},R,t,\widehat{\mu}):=\Gamma(x,A)\quad\text{and}\quad (D',\widetilde{D}',R',t',\widehat{\mu}'):=\Gamma(x',A).$$

Throughout our arguments,  $i \in [n]$  denotes the index at which the samples x, x' differ, that is,  $x_{-i} = x'_{-i}$  while we may have  $x_i \neq x'_i$ .

The main idea in the proof of Lemma 4.3 is to first bound the sensitivity of the mean, showing that (with high probability)  $\|\widehat{\mu} - \widehat{\mu}'\|_A$  is small, unless there are too many outlying entries  $x_i$ . We do this in Lemma 7.3 by showing that for appropriate subgroup sizes b (recall the random partition S of [n] into blocks of size n/b in MEANSAFE), the MEANSAFE algorithm correctly identifies all outliers without pruning many inlying datapoints. In the second step, we finalize the proof (section 7.1.1) by combining the sensitivity bound with more or less standard distributional stability guarantees for Gaussian distributions, which we list in the preliminary section 2.

We begin by formalizing two properties that will be helpful to proving the sensitivity bound in Lemma 7.3. We recall the notation t (respectively t') for denoting the number of pruned groups in Lines 5–8 of MEANSAFE on inputs x and x', while R and R' denote the sets of all pruned indices. Of the next two lemmas, Lemma 7.1 bounds differences between R and R' and t and t', while Lemma 7.2 is a generic lemma that bounds the difference of empirical means with nested index sets. These two lemmas are combined in Lemma 7.3 to bound the difference between the estimated mean  $\widehat{\mu} = \frac{1}{n-|R|} \sum_{j \notin R} x_j$  and  $\widehat{\mu}'$ . Before stating Lemma 7.1, recall for sets S, S' that  $d_{\text{sym}}(S, S') = \max\{|S \setminus S'|, |S' \setminus S|\}$ .

**Lemma 7.1** (Stability of rejected indices). Let t, t' and R, R' be as above. Then  $|t - t'| \le 1$ and  $d_{\text{sym}}(R, R') \leq b$ .

Proof. Let the set  $J:=\{j\mid \widetilde{D}_j\neq \bot,\widetilde{D}_j\geq \log(\sqrt{B}/4)\}$  index the subgroups pruned by the execution of MEANSAFE on the sample x', and similarly define J' relative to  $\widetilde{D}'$  for the sample x'. Then t=|J| and  $R=\cup_{j\in J}S_j$ , and also t'=|J'| and  $R'=\cup_{j\in J'}S_j'$ . We show  $d_{\mathrm{sym}}(J,J')\leq 1$ , from which the claim  $|t-t'|\leq 1$  follows immediately and the claim  $d_{\mathrm{sym}}(R,R')\leq b$  follows from the fact that  $|S_j|=b$  for all  $j\in [n/b]$ .

Recalling x and x' differ only at index i, suppose that  $i \in S_{\ell}$  for  $\ell \in [n/b]$ . Then  $x_{S_j} = x'_{S_j}$  for all  $j \neq \ell$ ; in particular, diam $_A(x_{S_j}) = \operatorname{diam}_A(x'_{S_j})$  and so  $D_j = D'_j$  for  $j \neq \ell$ . Thus, the indices of the k largest elements of D+Z and D'+Z, i.e., those subgroups identified by TOPk as having the largest diameters, which we denote by  $K = \{j \mid \widetilde{D}_j \neq \bot\}$  and  $K' = \{j \mid \widetilde{D}'_j \neq \bot\}$  respectively, differ by at most one index:  $d_{\operatorname{sym}}(K,K') \leq 1$  with equality obtaining only if  $\ell$  is in exactly one of K or K'. If  $\ell$  is in neither K nor K', then J = J' and the claim  $d_{\operatorname{sym}}(J,J') \leq 1$  follows. Otherwise, supposing  $\ell \in K$ , the bound  $d_{\operatorname{sym}}(K,K') \leq 1$  implies  $K \setminus \{\ell\} \subset K'$  and thus  $\widetilde{D}_{K \setminus \{\ell\}} = \widetilde{D}'_{K \setminus \{\ell\}}$ , or vice versa if  $\ell \in K'$ ;  $d_{\operatorname{sym}}(J,J') \leq 1$  then follows from  $J \subset K$  and  $J' \subset K'$ .

**Lemma 7.2.** Let  $\{y_1, \ldots, y_n\}$  be an arbitrary collection of vectors and  $S \subset S' \subset [n]$ . Define  $\mu_S := \frac{1}{|S|} \sum_{i \in S} y_i$  and  $\mu_{S'} := \frac{1}{|S'|} \sum_{i \in S'} y_i$ . Then

$$\|\mu_S - \mu_{S'}\| \le \frac{|S' \setminus S| \operatorname{diam}_{\|\cdot\|}(y_{S'})}{|S'|}.$$

*Proof.* Observe

$$\mu_S - \mu_{S'} = \mu_S - \left(\frac{|S|}{|S'|}\mu_S + \frac{1}{|S'|}\sum_{i \in S' \setminus S} y_i\right) = \frac{1}{|S'|}\sum_{i \in S' \setminus S} (\mu_S - y_i),$$

where from the assumption that  $S \subset S'$  we have

$$\max_{i \in S' \setminus S} \|y_i - \mu_S\| \le \max_{i \in S, i \in S'} \|y_i - y_i\| \le \operatorname{diam}_{\|\cdot\|}(S').$$

The claim then follows as  $\|\mu_S - \mu_{S'}\| \leq \frac{1}{|S'|} \sum_{i \in S' \setminus S} \operatorname{diam}_{\|\cdot\|}(y_{S'}) = \frac{|S' \setminus S| \operatorname{diam}_{\|\cdot\|}(x_{S'})}{|S'|}$ .

We now turn to the first step we outline, providing an explicit bound on  $\|\widehat{\mu} - \widehat{\mu}'\|_A$  except on the event that  $\max\{t,t'\} = k$ . Recall the definition  $\Delta = \frac{5b\sqrt{B}}{2n} \exp(3\sigma_{\text{top}} \log \frac{2n}{b\gamma})$  in the statement of Lemma 4.3.

**Lemma 7.3.** With probability at least  $1 - \gamma - n^2 2^{1-b}$ ,  $\max\{t, t'\} = k$  or  $\|\widehat{\mu} - \widehat{\mu}'\|_A \leq \Delta$ .

*Proof.* We first show that with probability at least  $1 - n^2 2^{-b}$  over the random partition  $S \sim \text{Uni}(\mathcal{P}_{n,b}), S = (S_1, \ldots, S_{n/b}),$ 

$$\operatorname{diam}_{A}(x_{R^{c}}) \leq \frac{1}{2} \exp(2 \|Z\|_{\infty} + \|Z'\|_{\infty}) \sqrt{B}, \tag{16}$$

with the same bound holding for x' by symmetry. To this end, observe that for the index set

$$J := \{ j \in [n/b] \mid \widetilde{D}_j \neq \bot, \widetilde{D}_j \ge \log(\sqrt{B}/4) \},$$

MEANSAFE constructs the removed indices R in Lines 5–8 via the union  $R = \bigcup_{j \in J} S_j$ . The first step in the bound (16) is to bound the diameter of the set  $x_{R^c}$  by the diameters of the constituent sets within R, which the following generic lemma allows (see Section 7.1.2 for a proof).

Claim 7.4. Let  $\{y_1, \ldots, y_n\}$  be an arbitrary collection of vectors and  $S \sim \text{Uni}(\mathcal{P}_{n,b})$ . With probability at least  $1 - n^2 2^{-b}$ , for all index sets  $J \subset [n/b]$ , the set  $S_J := \bigcup_{j \in J} S_j$  satisfies  $\dim_{\|\cdot\|}(y_{S_J}) \leq 2 \max_{j \in J} \dim_{\|\cdot\|}(y_{S_j})$ .

In light of Claim 7.4, inequality (16) follows by showing

$$\operatorname{diam}_{A}(x_{S_{i}}) \leq \exp(2 \|Z\|_{\infty} + \|Z'\|_{\infty}) \sqrt{B}/4 \tag{17}$$

for all  $j \notin J$  on the event t < k. When t < k, there exists an index  $\ell \in [n/b]$  such that  $\widetilde{D}_{\ell} \neq \bot$  and  $\widetilde{D}_{\ell} \leq \log(\sqrt{B}/4)$ , i.e.,  $\ell$  indexes one of the k largest elements of D + Z but  $\ell \notin J$ . Thus, for  $j \notin J$  such that  $\widetilde{D}_j = \bot$ , i.e.,  $\log(\operatorname{diam}_A(x_{S_j})) + Z_j$  is not among the k largest elements of D + Z (by the construction in TOPk), we have

$$\log(\operatorname{diam}_A(x_{S_i})) \leq \log(\operatorname{diam}_A(x_{S_\ell})) + 2 \|Z\|_{\infty}$$

Meanwhile, for all  $j \notin J$  such that  $\widetilde{D}_j \neq \bot$ , including  $j = \ell$ , from the definition of J we immediately have

$$\log(\operatorname{diam}_A(x_{S_i})) + Z_i' \le \log(\sqrt{B}/4).$$

The claim (17), and hence claim (16), thus follows from the preceding two displays. Moreover, via a union bound over the two executions of MEANSAFE, Claim 7.4 gives

$$\max\{\operatorname{diam}_{A}(x_{R^{c}}), \operatorname{diam}_{A}(x'_{R^{c}})\} \leq \exp(2\|Z\|_{\infty} + \|Z'\|_{\infty}) \frac{\sqrt{B}}{2} \text{ or } \max\{t, t'\} = k$$
 (18)

with probability at least  $1 - n^2 2^{1-b}$ .

We can now bound  $\|\widehat{\mu} - \widehat{\mu}'\|_A$  for  $\widehat{\mu} = \frac{1}{n-|R|} \sum_{j \notin R} x_j$  and  $\widehat{\mu} = \frac{1}{n-|R'|} \sum_{j \notin R'} x_j'$  via the following claim (essentially, a number of applications of the triangle inequality), whose proof we also defer (see Section 7.1.3).

Claim 7.5. 
$$\|\widehat{\mu} - \widehat{\mu}'\|_A \leq \frac{4(b+1)}{n} \max\{ \operatorname{diam}_A(x_{R^c}), \operatorname{diam}_A(x'_{R'^c}) \}.$$

Using Claim 7.5, the main Lemma 7.3 follows relatively quickly. By combining the display (18) with the fact that, by elementary calculation,

$$\mathbb{P}(\max\{\|Z\|_{\infty}, \|Z'\|_{\infty}\} > \sigma_{\text{top}}\log(2n/b\gamma)) \le \gamma,$$

we obtain that with probability at least  $1 - \gamma - n^2 2^{1-b}$ ,  $\max\{t, t'\} = k$  or

$$\left\|\widehat{\mu}' - \widehat{\mu}'\right\|_{A} \le \frac{2(b+1)\sqrt{B}}{n} \exp\left(2\left\|Z\right\|_{\infty} + \left\|Z'\right\|_{\infty}\right) \le \frac{2(b+1)\sqrt{B}}{n} \exp\left(3\sigma_{\text{top}}\log\frac{2n}{b\gamma}\right).$$

Recalling the assumption that the batchsize  $b \ge 4$  (so  $2(b+1) \le \frac{5}{2}b$ ), we obtain the lemma.  $\square$ 

# 7.1.1 Finalizing proof of Lemma 4.3

We prove for any (measurable) event  $O \subset \mathbb{R}^d \cup \{\bot\}$  that

$$\mathbb{P}(\widetilde{\mu} \in O) \le e^{\alpha + 1/\sigma_{W_{\text{mean}}}} \mathbb{P}(\widetilde{\mu}' \in O) + \beta_1 + \beta_2, \tag{19}$$

where  $\alpha > 0$  and  $\beta_1 \in (0,1)$  determine the Gaussian noise scale for  $Z^{\mathbb{N}} \sim \mathbb{N}(0,\sigma_{\mathbb{N}}^2 I)$  via

$$\sigma_{\mathsf{N}} = \begin{cases} \frac{\Delta}{\alpha} \sqrt{1.25 \log \frac{1}{\beta_1}} & \text{if } \alpha \leq 1\\ \frac{\Delta}{\sqrt{2 \log \frac{1}{\beta_1} + 2\alpha} - \sqrt{2 \log \frac{1}{\beta_1}}} & \text{otherwise,} \end{cases} \text{ and } \beta_2 = \frac{1}{2} e^{-(k/3 - 1)/\sigma_{W_{\text{mean}}}} + \gamma + n^2 2^{1 - b}.$$

The other direction follows by symmetry. We treat  $O \subset \mathbb{R}^d$  and  $O = \perp$  separately, merging the two cases at the end to show the claim (19). Supposing first  $O \subset \mathbb{R}^d$ , the following observation delineates necessary and sufficient conditions for  $\widetilde{\mu} \in O$ .

**Observation 7.1.** Let  $O \subset \mathbb{R}^d$ . Then  $\widetilde{\mu} \in O$  if and only if  $t \leq 2k/3 + W$  and  $\widehat{\mu} + A^{1/2}Z^{\mathsf{N}} \in O$ .

*Proof.* From the condition for returning  $\bot$  in Line 11 of MEANSAFE, we immediately have  $\widetilde{\mu} = \bot \notin \mathbb{R}^d$  if and only if t > 2k/3 + W; thus, the condition  $t \le 2k/3 + W$  is necessary and sufficient for  $\widetilde{\mu} \in \mathbb{R}^d$ . As either  $\widetilde{\mu} = \bot$  or  $\widetilde{\mu} = \widehat{\mu} + A^{1/2}Z^{\mathbb{N}}$  by definition, it then follows trivially that  $t \le 2k/3 + W$  and  $\widehat{\mu} + A^{1/2}Z^{\mathbb{N}} \in O$  together suffice to obtain  $\widetilde{\mu} \in O$ .

Marginalizing over the number of sets of rejected indices t and  $\widehat{\mu}$  we have the following sequence of inequalities:

$$\begin{split} &\mathbb{P}(\widetilde{\mu} \in O) \\ &= \mathbb{E}\left[\mathbb{P}(\widehat{\mu} + A^{1/2}Z^{\mathsf{N}} \in O \mid \widehat{\mu})\mathbb{P}(t \leq 2k/3 + W \mid t)\right] \\ &\stackrel{(i)}{\leq} \mathbb{E}\left[\mathbb{P}(\widehat{\mu} + A^{1/2}Z^{\mathsf{N}} \in O \mid \widehat{\mu})\mathbb{P}(t \leq 2k/3 + W \mid t)1\{\|\widehat{\mu}' - \widehat{\mu}'\|_{A} \leq \Delta\}\right] \\ &\quad + \mathbb{E}\left[\mathbb{P}(t \leq 2k/3 + W \mid t)1\{\max\{t, t'\} = k\}\right] + \gamma + n^{2}2^{1-b} \\ &\stackrel{(ii)}{\leq} \mathbb{E}\left[\mathbb{P}(\widehat{\mu} + A^{1/2}Z^{\mathsf{N}} \in O \mid \widehat{\mu})\mathbb{P}(t \leq 2k/3 + W \mid t)1\{\|\widehat{\mu}' - \widehat{\mu}'\|_{A} \leq \Delta\}\right] \\ &\quad + \mathbb{P}(W \geq k/3 - 1) + \gamma + n^{2}2^{1-b} \\ &= \mathbb{E}\left[\mathbb{P}(\widehat{\mu} + A^{1/2}Z^{\mathsf{N}} \in O \mid \widehat{\mu})\mathbb{P}(t \leq 2k/3 + W \mid t)1\{\|\widehat{\mu}' - \widehat{\mu}'\|_{A} \leq \Delta\}\right] + \beta_{2} \end{split} \tag{20}$$

Here, step (i) follows because  $\|\widehat{\mu}' - \widehat{\mu}'\|_A \leq \Delta$  or  $\max\{t,t'\} = k$  occurs with probability at least  $1 - \gamma - n^2 2^{1-b}$  by Lemma 7.3; step (ii) because  $|t - t'| \leq 1$  by Lemma 7.1 and so  $\max\{t,t'\} = k$  implies  $t \geq k-1$ ; the final equality follows from the identity  $\mathbb{P}(W \geq k/3 - 1) = \frac{1}{2}e^{-(k/3-1)/\sigma_{W_{\text{mean}}}}$  and definition of  $\beta_2$ .

Continuing, we can bound the last expectation in the preceding display by

$$\mathbb{E}\left[\mathbb{P}(\widehat{\mu} + A^{1/2}Z^{\mathsf{N}} \in O \mid \widehat{\mu})\mathbb{P}(t \leq 2k/3 + W \mid t)1\{\|\widehat{\mu}' - \widehat{\mu}'\|_{A} \leq \Delta\}\right] \\
\stackrel{(i)}{\leq} \exp(\alpha)\mathbb{E}\left[\mathbb{P}(\widehat{\mu}' + A^{1/2}Z^{\mathsf{N}} \in O \mid \widehat{\mu}')\mathbb{P}(t \leq 2k/3 + W \mid t)\right] + \beta_{1} \\
\stackrel{(ii)}{\leq} \exp(\alpha + 1/\sigma_{W_{\text{mean}}})\mathbb{E}\left[\mathbb{P}(\widehat{\mu}' + A^{1/2}Z^{\mathsf{N}} \in O \mid \widehat{\mu}')\mathbb{P}(t' \leq 2k/3 + W \mid t')\right] + \beta_{1} \\
= \exp(\alpha + 1/\sigma_{W_{\text{mean}}})\mathbb{P}(\widetilde{\mu}' \in O) + \beta_{1}, \tag{21}$$

with step (i) following from the privacy of the Gaussian mechanism with noise  $\sigma_N$  and sensitivity bound  $\|\widehat{\mu} - \widehat{\mu}'\| \leq \Delta$  (Lemma 2.5); step (ii) from  $|t - t'| \leq 1$  by Lemma 7.1 and that  $W \sim \mathsf{Lap}(\sigma_{W_{\text{mean}}})$ ; and the final equality follows directly from Observation 7.1, applied

here to the execution of MEANSAFE on data x'. Combining inequalities (20) and (21) yields the claim (19) when  $O \subset \mathbb{R}^d$ .

For the case that  $O = \{\bot\}$ , we have

$$\begin{split} \mathbb{P}(\widetilde{\mu} = \bot) &= \mathbb{E}\left[\mathbb{P}(t > 2k/3 + W \mid t)\right] \\ &\leq e^{\frac{1}{\sigma_{W_{\text{mean}}}}} \mathbb{E}\left[\mathbb{P}(t' > 2k/3 + W \mid t')\right] = \exp^{\frac{1}{\sigma_{W_{\text{mean}}}}} \mathbb{P}(\widetilde{\mu} = \bot). \end{split}$$

Here, the two equalities follow from the condition for returning  $\bot$  in Line 11 of MEANSAFE, while the inequality follows because  $|t-t'| \le 1$  by Lemma 7.1 and that  $W \sim \mathsf{Lap}(\sigma_{W_{\text{mean}}})$ . The claim (19) for arbitrary O is immediate.

### 7.1.2 Proof of Claim 7.4

Consider the event  $\mathcal{E}$  that for all indices  $i_1, i_2 \in [n]$ , with  $i_1 \in S_{j_1}$  and  $i_2 \in S_{j_2}$ , we have  $||y_{i_1} - y_{i_2}|| \leq 2 \max\{\operatorname{diam}(y_{S_{j_1}}), \operatorname{diam}(y_{S_{j_2}})\}$ . The claim holds on  $\mathcal{E}$ : for any  $J \subset [n/b]$  and  $S_J := \bigcup_{j \in J} S_j$ , there exist  $j_1, j_2 \in J$  with  $i_1 \in S_{j_1}$  and  $i_2 \in S_{j_2}$  attaining  $\operatorname{diam}(y_{S_J}) = ||y_{i_1} - y_{i_2}||$ , and so

$$||y_{i_1} - y_{i_2}|| \le 2 \max\{\operatorname{diam}(y_{S_{j_1}}), \operatorname{diam}(y_{S_{j_2}})\} \le 2 \max_{j \in J} \operatorname{diam}(y_{S_j}).$$

It remains to show that  $\mathcal{E}$  occurs with probability at least  $1 - n^2 2^{-b}$ . As there are  $\binom{n}{2} \leq \frac{1}{2}n^2$  unordered pairs of distinct indices  $i_1, i_2 \in [n]$ , the result obtains from a union bound if we show that  $||y_{i_1} - y_{i_2}|| > 2 \max\{\operatorname{diam}(y_{S_{j_1}}), \operatorname{diam}(y_{S_{j_2}})\}$  occurs with probability at most  $2^{1-b}$ .

Proceeding, let  $i_1, i_2 \in [n]$  and  $i_1 \in S_{j_1}, i_2 \in S_{j_2}$  and let  $c = \frac{1}{2} ||y_{i_1} - y_{i_2}||$ . If  $i_1 = i_2$  or  $j_1 = j_2$ , there is nothing to show, so assume  $i_1 \neq i_2$  and  $j_1 \neq j_2$ . Let  $C_1 = \{i \in [n] \setminus \{i_1, i_2\} | ||y_{i_1} - y_i|| < c\}$  and  $C_2 = \{i \in [n] \setminus \{i_1, i_2\} | ||y_{i_2} - y_i|| < c\}$  be those indices i for which  $y_i$  is close to  $y_{i_1}$  or  $y_{i_2}$ , respectively. By the triangle inequality,  $C_1$  is disjoint from  $C_2$ , and so without loss of generality, we suppose that  $|C_1| \leq (n-2)/2$ .

We wish to show that  $\operatorname{diam}(y_{S_{j_1}}) \geq c$ , for which it is sufficient that there exists an index in  $S_{j_1} \setminus \{i_1\}$  not in  $C_1$ . So by showing  $S_{j_1} \setminus \{i_1\} \subset C_1$  occurs with probability at most  $2^{1-b}$ , we will be done. As  $S \sim \operatorname{Uni}(\mathcal{P}_{n,b})$ , the set  $S_{j_1} \setminus \{i_1\}$  is a uniformly distributed subset of  $[n] \setminus \{i_1, i_2\}$  of size b-1. Consequently, there are  $\binom{n-2}{b-1}$  distinct values it can take and  $\binom{|C_1|}{b-1}$  values such that  $S_{j_1} \setminus \{i_1\} \subset C_1$ . Therefore, the probability that  $S_{j_1} \setminus \{i_1\} \subset C_1$  is

$$\mathbb{P}(S_{j_1} \setminus \{i_1\} \subset C_1) = \frac{\binom{|C_1|}{b-1}}{\binom{n-2}{b-1}} = \prod_{i=0}^{b-2} \frac{(|C_1|-i)_+}{n-2-i} \le \left(\frac{|C_1|}{n-2}\right)^{b-1} \le 2^{1-b},$$

where the last inequality follows because  $|C_1| \leq (n-2)/2$ .

### 7.1.3 Proof of Claim 7.5

Recall that

$$\widehat{\mu} = \frac{1}{n - |R|} \sum_{j \notin R} x_j$$
 and  $\widehat{\mu}' = \frac{1}{n - |R'|} \sum_{j \notin R'} x'_j$ ,

and define

$$R_{\mathrm{all}} := R \cup R', \quad \widehat{\mu}_{\mathrm{all}} := \frac{1}{n - |R_{\mathrm{all}}|} \sum_{j \notin R_{\mathrm{all}}} x_j, \quad \widehat{\mu}'_{\mathrm{all}} := \frac{1}{n - |R_{\mathrm{all}}|} \sum_{j \notin R_{\mathrm{all}}} x'_j.$$

Lemma 7.1 gives  $|R^c \setminus R_{\text{all}}^c| = |R_{\text{all}} \setminus R| \le b$ , and by assumption on the batchsize b and rejection threshold k we also have  $|R_{\text{all}}| \le b + |R| \le b + kb \le \frac{n}{2}$ .

Applying Lemma 7.2 with  $S=R_{\rm all}^c$ , and  $S'=R^c$ , we get

$$\|\widehat{\mu} - \widehat{\mu}_{\text{all}}\|_A \le \frac{|R^c \setminus R_{\text{all}}^c| \operatorname{diam}_A(x_{R^c})}{|R^c|} \le \frac{2b \operatorname{diam}_A(x_{R^c})}{n}$$

as  $|R| \leq |R_{\text{all}}| \leq \frac{n}{2}$ . Applying Lemma 7.2 again, this time with dataset x',  $S = R'^c_{\text{all}}$  and  $S' = R'^c$ , we get  $\|\widehat{\mu}' - \widehat{\mu}'_{\text{all}}\|_A \leq \frac{2b}{2} \operatorname{diam}_A(x'_{R'^c})$ .

 $S' = R'^c$ , we get  $\|\widehat{\mu}' - \widehat{\mu}'_{\text{all}}\|_A \leq \frac{2b}{n} \operatorname{diam}_A(x'_{R'^c})$ . Now we bound  $\widehat{\mu}_{\text{all}} - \widehat{\mu}'_{\text{all}}$ , where recalling that index i is the sole (potentially) differing index in x, x' (that is,  $x_{-i} = x'_{-i}$ ), we can write as

$$\widehat{\mu}_{\text{all}} - \widehat{\mu}'_{\text{all}} = \frac{1}{n - |R_{\text{all}}|} \sum_{j \notin R_{\text{all}}} (x_j - x'_j) = \frac{\mathbf{1} \{i \notin R_{\text{all}}\}}{n - |R_{\text{all}}|} (x_i - x'_i).$$

If  $i \in R_{\text{all}}$ , this difference is 0. Otherwise,  $i \notin R$  and  $i \notin R'$ . As  $|R_{\text{all}}| \leq \frac{n}{2}$ , we may pick some  $j' \notin R_{\text{all}} \cup \{i\}$ . Because  $x_{j'} = x'_{j'}$ , we have both  $||x_i - x_{j'}||_A \leq \text{diam}_A(x_{R^c})$  and both  $||x'_i - x_{j'}||_A \leq \text{diam}_A(x'_{R^{c}})$ . The triangle inequality then gives  $||x_i - x'_i||_A \leq 2 \max\{\text{diam}_A(x_{R^c}), \text{diam}_A(x'_{R^{c}})\}$ , and so  $||\widehat{\mu}_{\text{all}} - \widehat{\mu}'_{\text{all}}||_A \leq \frac{4}{n} \max\{\text{diam}_A(x_{R^c}), \text{diam}_A(x'_{R^{c}})\}$ . Combining the above, the claim follows immediately from

$$\begin{split} \|\widehat{\mu} - \widehat{\mu}'\|_{A} &\leq \|\widehat{\mu} - \widehat{\mu}_{\text{all}}\|_{A} + \|\widehat{\mu}_{\text{all}} - \widehat{\mu}'_{\text{all}}\|_{A} + \|\widehat{\mu}'_{\text{all}} - \widehat{\mu}'\|_{A} \\ &\leq \frac{2b \text{diam}_{A}(x_{R^{c}})}{n} + \frac{4 \max\{\text{diam}_{A}(x_{R^{c}}), \text{diam}_{A}(x'_{R'^{c}})\}}{n} + \frac{2b \text{diam}_{A}(x'_{R'^{c}})}{n}. \end{split}$$

### 7.2 Proof of Lemma 4.4

Unpacking the execution transcripts  $\Gamma(x,A)$  and  $\Gamma(x,A')$  from (3) as

$$(D, \widetilde{D}, R, t, \widehat{\mu}) := \Gamma(x, A)$$
 and  $(D', \widetilde{D}', R', t', \widehat{\mu}') := \Gamma(x, A'),$ 

observe that given the pair  $(\widetilde{D}, A^{1/2}Z^{\mathbb{N}})$ ,  $\widetilde{\mu}(x, A)$  is independent of A (see the execution of MEANSAFE), and analogously,  $\widetilde{\mu}(x, A')$  is independent of A' given  $(\widetilde{D}', A'^{1/2}Z^{\mathbb{N}})$ . Therefore, by showing  $A^{1/2}Z^{\mathbb{N}} \stackrel{d}{=}_{\alpha_1,\beta} A'^{1/2}Z^{\mathbb{N}}$  and  $\widetilde{D} \stackrel{d}{=}_{\alpha_2,0} \widetilde{D}'$ , basic composition (Lemma 2.1) and the post-processing property (Lemma 2.3) will imply the claimed result that  $\widetilde{\mu} \stackrel{d}{=}_{\alpha_1+\alpha_2,\beta} \widetilde{\mu}'$ .

post-processing property (Lemma 2.3) will imply the claimed result that  $\widetilde{\mu} \stackrel{d}{=}_{\alpha_1 + \alpha_2, \beta} \widetilde{\mu}'$ . Recalling  $Z^{\mathsf{N}} \sim \mathsf{N}(0, \sigma_{\mathsf{N}}^2 I)$ , we have  $A^{1/2} Z^{\mathsf{N}} \sim \mathsf{N}(0, \sigma_{\mathsf{N}}^2 A)$  and  $A'^{1/2} Z^{\mathsf{N}} \sim \mathsf{N}(0, \sigma_{\mathsf{N}}^2 A')$ , and so  $A^{1/2} Z^{\mathsf{N}} \stackrel{d}{=}_{\alpha_1, \beta} A'^{1/2} Z^{\mathsf{N}}$  follows immediately from the assumption  $d_{\mathrm{psd}}(A, A') \leq \frac{a}{n}$  and the closeness of Gaussian distributions with differing covariances (Lemma 2.6).

To show  $\widetilde{D} \stackrel{d}{=}_{\alpha_2,0} \widetilde{D}'$ , we make the following observation to bound the sensitivity of the log-Mahalanobis norm for A and A'.

**Observation 7.2.** Suppose  $A, A' \in \mathbb{R}^{d \times d}$  and  $d_{\text{psd}}(A, A') \leq \gamma < \infty$ . Then for any  $v \in \text{Col}(A)$ ,  $|\log ||v||_A - \log ||v||_{A'}| \leq \gamma/2$ . For any  $v \notin \text{Col}(A)$ ,  $\log ||v||_A = \log ||v||_{A'} = \infty$ .

*Proof.* Observe  $d_{\text{psd}}(A, A') < \infty$  trivially implies A and A' are PSD and their columnspaces coincide, from which the second claim immediately follows. For  $v \in \text{Col}(A)$ , we only show  $\log \|v\|_{A'} \leq \frac{1}{2}\gamma + \log \|v\|_A$ , as the reverse inequality holds by symmetry. By assumption,

$$\left\|A'^{\dagger/2}(A-A')A'^{\dagger/2}\right\|_{\text{op}} \le d_{\text{psd}}(A,A') \le \gamma$$

and hence  $A'^{\dagger/2}(A-A')A'^{\dagger/2} \leq \gamma I$ . Conjugating by  $A'^{1/2}$  and rearranging terms, we have  $\Pi_{A'}A\Pi_{A'} \leq (1+\gamma)A'$ . Because  $\Pi_{A'} = \Pi_A$ , we have  $\Pi_{A'}A\Pi_{A'} = A$ , which yields  $A \leq (1+\gamma)A'$ , or equivalently  $A'^{\dagger} \leq (1+\gamma)A^{\dagger}$ . Therefore  $\|v\|_{A'}^2 \leq (1+\gamma)\|v\|_A^2$ . Taking square roots and logarithms on both sides proves the claim as  $\log(\sqrt{1+\gamma}) \leq \frac{\gamma}{2}$ .

This observation, coupled with our construction that both  $\widetilde{\mu}(x,A)$  and  $\widetilde{\mu}(x,A')$  use the same (random) partition  $\mathcal{S}=(S_1,\ldots,S_{n/b})$ , implies  $|D_j-D_j'|\leq \frac{a}{2n}$  for all  $j\in[n/b]$ ; hence  $\|D-D'\|_{\infty}\leq \frac{a}{2n}$  (the indices where the entries are infinite coincide). The closeness properties of TOPk (Lemma 2.4) and our choice  $\sigma_{\text{top}}=\frac{ka}{n\alpha_2}$  then give  $\widetilde{D}\stackrel{d}{=}_{\alpha_2,0}\widetilde{D}'$ .

# 8 Discussion

The simplicity of mean estimation in classical statistics belies the sophistication necessary to adaptively and accurately estimate a mean under differential privacy constraints. While we have developed (nearly) minimax optimal procedures for mean estimation, a number of questions remain open, and we hope that we or others will tackle them. From a practical perspective, while our procedure is implementable, the numerical constant factors we have maintained to guarantee privacy—in addition to the logarithmic factors in n and  $\log \frac{1}{\delta}$ —may make effective use of the procedure challeenging. From a theoretical perspective, it is still interesting to attempt to remove the logarithmic factors present in our bounds. Additionally, while we can adapt to weaker than sub-Gaussian moment bounds (via the method ADAMEAN), it may be possible to provide a sharper procedure or tighter analysis to achieve optimal dependence on dimension d and privacy level  $\varepsilon$ , as in the case that p moments are available, our results appear to be roughly a factor of  $(\sqrt{d}/\varepsilon)^{1/p}$  loose (recall Examples 2 and 4). It will be interesting to see the extent of possibilities for differentially private estimation in these more general cases.

# A Proofs of standard privacy results

# A.1 Proof of Lemma 2.5

The first case follows from Dwork and Roth [12, Theorem 3.22]. For the second, we use Mironov's Rényi-differential privacy [33]. The Rényi  $\alpha$ -divergence between distributions P and Q is  $D_{\alpha}\left(P\|Q\right) = \frac{1}{\alpha-1}\log\int\left(\frac{dP}{dQ}\right)^{\alpha}dQ$ , and [33, Proposition 3] if  $D_{\alpha}\left(P\|Q\right) \leq c$ , then for all measurable A and  $\delta > 0$  we have  $P(A) \leq \exp\left(c + \frac{\log(1/\delta)}{\alpha-1}\right)Q(A) + \delta$ . The Rényi divergence for Gaussians has the explicit form

$$D_{\alpha}\left(\mathsf{N}(\mu_1, \tau^2 \Sigma) \| \mathsf{N}(\mu_2, \tau^2 \Sigma)\right) = \frac{\alpha}{2\tau^2} \| \mu_1 - \mu_2 \|_{\Sigma}^2.$$

When  $\rho \geq \|\mu_1 - \mu_2\|_{\Sigma}$ , we set  $\beta = \alpha - 1$  and see that for  $\varepsilon$  satisfying

$$\varepsilon = \frac{\rho^2}{2\tau^2} + \frac{\beta\rho^2}{2\tau^2} + \frac{\log(1/\delta)}{\beta}$$

we obtain  $\mathsf{N}(\mu_1, \tau^2 \Sigma) \stackrel{d}{=}_{\varepsilon, \delta} \mathsf{N}(\mu_2, \tau^2 \Sigma)$ . Choosing  $\beta$  to minimize the preceding  $\varepsilon$  we obtain  $\varepsilon = \frac{\rho^2}{2\tau^2} + \frac{\rho}{\tau} \sqrt{2\log\frac{1}{\delta}}$ , and solving for  $\eta = \frac{1}{\tau}$  in  $\frac{\rho^2}{2}\eta^2 + \sqrt{2\log\frac{1}{\delta}}\rho\eta - \varepsilon$  yields

$$\tau = \frac{1}{\eta} = \frac{\rho}{\sqrt{2\log\frac{1}{\delta} + 2\varepsilon} - \sqrt{2\log\frac{1}{\delta}}}$$

is always sufficient to guarantee  $N(\mu_1, \tau^2 \Sigma) \stackrel{d}{=}_{\varepsilon, \delta} N(\mu_2, \tau^2 \Sigma)$ .

# A.2 Proof of Lemma 2.6

Without loss of generality, we may assume  $\mu=0$ . We first reduce to the case where  $\Sigma_1$  and  $\Sigma_2$  are full-rank. Because  $d_{\mathrm{psd}}(\Sigma_1,\Sigma_2)<\infty$ , we have immediately that there exists a vector space  $V\subset\mathbb{R}^d$  with  $V=\mathrm{Col}(\Sigma_1)=\mathrm{Col}(\Sigma_2)$ . Letting  $k=\dim(\mathrm{Col}(\Sigma_1))$ , take  $U\in\mathbb{R}^{d\times k}$  to be an orthonormal matrix such that  $V=\mathrm{Col}(U)$ . The random variables  $X\sim \mathsf{N}(0,\Sigma_1)$  and  $Y\sim\mathsf{N}(0,\Sigma_2)$  have support V and multiplication by  $U^T$  is an isomorphism between V and  $\mathbb{R}^k$ , so  $X\stackrel{d}{=}_{\varepsilon,\delta}Y$  if and only if  $U^TX\stackrel{d}{=}_{\varepsilon,\delta}U^TY$ . Of course,  $U^TX\sim\mathsf{N}\left(0,U^T\Sigma_1U\right)$  and  $U^TY\sim\mathsf{N}\left(0,U^T\Sigma_2U\right)$  and both  $U^T\Sigma_1U$  and  $U^T\Sigma_2U$  are full rank. The orthogonality of U gives U0 gives U1 gives U2. Hence, by showing the lemma for the full-rank matrices  $U^T\Sigma_1U$  and  $U^T\Sigma_2U$ , we will have shown the claim for  $\Sigma_1$  and  $\Sigma_2$ .

We proceed with the full-rank case with an argument similar to Brown et al. [7, Lemma 4.15]. Define  $D_1 = \Sigma_1^{1/2} \Sigma_2^{-1} \Sigma_1^{1/2} - I$  and  $D_2 = \Sigma_2^{1/2} \Sigma_1^{-1} \Sigma_2^{1/2} - I$ . As  $D_1$  has the same spectrum as  $\Sigma_2^{-1/2} \Sigma_1 \Sigma_2^{-1/2} - I$ , we have by assumption that  $\|D_1\|_* \leq \gamma$  and  $\|D_2\|_* \leq \gamma$ .

Let  $f_1$  be the density of  $P_1 = \mathsf{N}(0,\Sigma_1)$  and  $f_2$  that of  $P_2 = \mathsf{N}(0,\Sigma_2)$ . Then, to show  $(\varepsilon,\delta)$ -closeness, it suffices to show  $\ell(W) := \left|\log \frac{f_1(W)}{f_2(W)}\right| \le \varepsilon$  with probability at least  $1-\delta$  when W is drawn from either  $P_1$  or  $P_2$ . By symmetry, it suffices to only show this bound for the case when  $W \sim P_1$ . Expanding  $\ell$ , we have

$$\ell(w) = \left| \frac{1}{2} \log \frac{\det(\Sigma_2)}{\det(\Sigma_1)} + \frac{1}{2} w^T (\Sigma_2^{-1} - \Sigma_1^{-1}) w \right| \le \frac{1}{2} \left| w^T (\Sigma_2^{-1} - \Sigma_1^{-1}) w \right| + \frac{1}{2} \left| \log \frac{\det(\Sigma_2)}{\det(\Sigma_1)} \right|. \tag{22}$$

The final term is independent of w and has the bound

$$\left| \log \frac{\det(\Sigma_2)}{\det(\Sigma_1)} \right| = \max \left\{ \log \det(\Sigma_2^{1/2} \Sigma_1^{-1} \Sigma_2^{1/2}), \log \det(\Sigma_1^{1/2} \Sigma_2^{-1} \Sigma_1^{1/2}) \right\}$$

$$\leq \max \{ \operatorname{tr}(D_2), \operatorname{tr}(D_1) \} \leq \max \{ \|D_2\|_*, \|D_1\|_* \} \leq \gamma,$$

where the first inequality holds because  $\log \det(A) \leq \operatorname{tr}(A-I)$  for any positive definite A. Now we bound the first term on the right hand side of inequality (22) with high probability. Since  $W \sim P_1$ , the whitened random variable  $Z = \Sigma_1^{-1/2} W \sim \mathsf{N}(0,I)$ . We then have

$$|W^T(\Sigma_2^{-1} - \Sigma_1^{-1})W| = |Z^T D_1 Z|,$$

and so by the Hanson-Wright inequality [e.g. 38, Thm. 6.2.1], we have with probability at least  $1-\delta$  that

$$|Z^T D_1 Z| \le |\operatorname{tr}(D_1)| + 2 ||D_1||_{\operatorname{Fr}} \sqrt{\log \frac{2}{\delta}} + 2 ||D_1||_{\operatorname{op}} \log \frac{2}{\delta} \le 5\gamma \log \frac{2}{\delta},$$

where the inequality holds because  $||D_1||_{\text{op}} \leq ||D_1||_{\text{Fr}} \leq ||D_1||_* \leq \gamma$  and  $\log \frac{2}{\delta} \geq 1$ . Thus  $\ell(W) \leq 6\gamma \log \frac{2}{\delta} \leq \varepsilon$  with probability at least  $1 - \delta$ .

# References

[1] M. Abadi, A. Chu, I. Goodfellow, B. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep learning with differential privacy. In 23rd ACM Conference on Computer and Communications Security (ACM CCS), pages 308–318, 2016.

- [2] D. Alabi, P. K. Kothari, P. Tankala, P. Venkat, and F. Zhang. Privately estimating a gaussian: Efficient, robust and optimal. arXiv:2212.08018 [cs.DS], 2022.
- [3] Apple Differential Privacy Team. Learning with privacy at scale, 2017. Available at https://machinelearning.apple.com/2017/12/06/learning-with-privacy-at-scale.html.
- [4] H. Ashtiani and C. Liaw. Private and polynomial time algorithms for learning gaussians and beyond. In *Proceedings of the Thirty Fifth Annual Conference on Computational Learning Theory*, 2022.
- [5] R. F. Barber and J. C. Duchi. Privacy and statistical risk: Formalisms and minimax bounds. arXiv:1412.4451 [math.ST], 2014.
- [6] S. Biswas, Y. Dong, G. Kamath, and J. Ullman. CoinPress: Practical private mean and covariance estimation. In *Advances in Neural Information Processing Systems* 33, 2020.
- [7] G. Brown, M. Gaboardi, A. Smith, J. Ullman, and L. Zakynthinou. Covariance-aware private mean estimation without private covariance estimation. arXiv:2106.13329 [cs.LG], 2021.
- [8] R. Chen, A. Gittens, and J. A. Tropp. The masked sample covariance estimator: an analysis using matrix concentration inequalities. *Information and Inference*, to appear, 2012.
- [9] H. Cramér. Mathematical Methods in Statistics. Princeton University Press, 1946.
- [10] I. Diakonikolas and D. M. Kane. Algorithmic High-dimensional Robust Statistics. Cambridge University Press, 2022.
- [11] C. Dwork and J. Lei. Differential privacy and robust statistics. In *Proceedings of the Forty-First Annual ACM Symposium on the Theory of Computing*, pages 371–380, 2009.
- [12] C. Dwork and A. Roth. The algorithmic foundations of differential privacy. Foundations and Trends in Theoretical Computer Science, 9(3 & 4):211–407, 2014.
- [13] C. Dwork, K. Kenthapadi, F. McSherry, I. Mironov, and M. Naor. Our data, ourselves: Privacy via distributed noise generation. In Advances in Cryptology (EUROCRYPT 2006), 2006.
- [14] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *Proceedings of the Third Theory of Cryptography Conference*, pages 265–284, 2006.
- [15] C. Dwork, G. N. Rothblum, and S. P. Vadhan. Boosting and differential privacy. In 51st Annual Symposium on Foundations of Computer Science, pages 51–60, 2010.
- [16] C. Dwork, K. Talwar, A. Thakurta, and L. Zhang. Analyze Gauss: optimal bounds for privacy-preserving principal component analysis. In *Proceedings of the Forty-Sixth* Annual ACM Symposium on the Theory of Computing. ACM, 2014.
- [17] U. Erlingsson, V. Pihur, and A. Korolova. RAPPOR: Randomized aggregatable privacypreserving ordinal response. In *Proceedings of the 21st ACM Conference on Computer* and Communications Security (CCS), 2014.

- [18] S. Garfinkel, J. Abowd, and S. Powazek. Issues encountered deploying differential privacy. In *Proceedings of the 2018 Workshop on Privacy in the Electronic Society (WPES)*, pages 133–137, 2018.
- [19] F. Hampel, E. Ronchetti, P. Rousseeuw, and W. Stahel. Robust Statistics: The Approach Based on Influence Functions. John Wiley & Sons, 1986.
- [20] M. Hardt and G. N. Rothblum. A multiplicative weights mechanism for privacy-preserving data analysis. In 51st Annual Symposium on Foundations of Computer Science, pages 61–70, 2010.
- [21] S. B. Hopkins, G. Kamath, and M. Majid. Efficient mean estimation with pure differential privacy via a sum-of-squares exponential mechanism. In *Proceedings of the Fifty-Fourth Annual ACM Symposium on the Theory of Computing*, pages 1406–1417. ACM, 2022.
- [22] S. B. Hopkins, G. Kamath, M. Majid, and S. Narayanan. Robustness implies privacy in statistical estimation. arXiv:2212.05015 [cs.DS], 2022.
- [23] Z. Huang, Y. Liang, and K. Yi. Instance-optimal mean estimation under differential privacy. In *Advances in Neural Information Processing Systems* 34, volume 34, pages 25993–26004, 2021.
- [24] P. J. Huber. Robust estimation of a location parameter. Annals of Mathematical Statistics, 35(1):73–101, 1964.
- [25] P. J. Huber and E. M. Ronchetti. *Robust Statistics*. John Wiley and Sons, second edition, 2009.
- [26] G. Kamath, J. Li, V. Singhal, and J. R. Ullman. Privately learning high-dimensional distributions. arXiv:1805.00216 [cs.DS], 2018.
- [27] P. K. Kothari, P. Manurangsi, and A. Velingker. Private robust estimation by stabilizing convex relaxations. In *Proceedings of the Thirty Fifth Annual Conference on Computational Learning Theory*, pages 1–55, 2022.
- [28] A. N. D. A. D. Lauger, P. E. Singer, D. Kifer, J. P. Reiter, A. Machanavajjhala, S. L. Garfinkel, S. A. Dahl, M. Graham, V. Karwa, H. Kim, P. Leclerc, I. M. Schmutte, W. N. Sexton, L. Vilhuber, and J. M. Abowd. The modernization of statistical disclosure limitation at the U.S. Census Bureau. Available online at https://www2.census.gov/cac/sac/meetings/2017-09/statistical-disclosure-limitation.pdf, 2017.
- [29] B. Laurent and P. Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, 28(5):1302 1338, 2000.
- [30] L. Le Cam. Asymptotic Methods in Statistical Decision Theory. Springer-Verlag, 1986.
- [31] X. Liu, W. Kong, S. Kakade, and S. Oh. Robust and differentially private mean estimation. arXiv:2102.09159 [cs.LG], 2021.
- [32] X. Liu, W. Kong, and S. Oh. Differential privacy and robust statistics in high dimensions. In *Proceedings of the Thirty Fifth Annual Conference on Computational Learning Theory*, pages 1167–1246, 2022.

- [33] I. Mironov. Rényi differential privacy. In 30th IEEE Computer Security Foundations Symposium (CSF), pages 263–275, 2017.
- [34] G. Qiao, W. Su, and L. Zhang. Oneshot differentially private top-k selection. In Proceedings of the 38th International Conference on Machine Learning, 2021.
- [35] T. Steinke and J. Ullman. Between pure and approximate differential privacy. *Journal of Privacy and Confidentiality*, 7(2):3–22, 2017.
- [36] J. W. Tukey. A survey of sampling from contaminated distributions. Princeton University, 1960.
- [37] R. Vershynin. Introduction to the non-asymptotic analysis of random matrices. In *Compressed Sensing: Theory and Applications*, chapter 5, pages 210–268. Cambridge University Press, 2012.
- [38] R. Vershynin. High Dimensional Probability: An Introduction with Applications in Data Science. Cambridge University Press, 2019.
- [39] M. J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge University Press, 2019.