
Private estimation algorithms for stochastic block models and mixture models

Hongjie Chen
ETH Zürich

Vincent Cohen-Addad
Google Research

Tommaso d’Orsi
Bocconi*

Alessandro Epasto
Google Research

Jacob Imola
UC San Diego

David Steurer
ETH Zürich

Stefan Tiegel
ETH Zürich

Abstract

We introduce general tools for designing efficient private estimation algorithms, in the high-dimensional settings, whose statistical guarantees almost match those of the best known non-private algorithms. To illustrate our techniques, we consider two problems: recovery of stochastic block models and learning mixtures of spherical Gaussians.

For the former, we present the first efficient (ε, δ) -differentially private algorithms for both weak recovery and exact recovery. Previously known algorithms achieving comparable guarantees required quasi-polynomial time. We complement these results with an information-theoretic lower bound that highlights how the guarantees of our algorithms are almost tight.

For the latter, we design an (ε, δ) -differentially private algorithm that recovers the centers of the k -mixture when the minimum separation is at least $O(k^{1/t}\sqrt{t})$. For all choices of t , this algorithm requires sample complexity $n \geq k^{O(1)}d^{O(t)}$ and time complexity $(nd)^{O(t)}$. Prior work required either an additional additive $\Omega(\sqrt{\log n})$ term in the minimum separation or an explicit upper bound on the Euclidean norm of the centers.

1 Introduction

Computing a model that best matches a dataset is a fundamental question in machine learning and statistics. Given a set of n samples from a model, how to find the most likely parameters of the model that could have generated this data? This basic question has been widely studied for several decades, and recently revisited in the context where the input data has been partially corrupted (i.e., where few samples of the data have been adversarially generated—see for instance [37, 18, 22, 20]). This has led to several recent works shedding new lights on classic model estimation problems, such as the Stochastic Block Model (SBM) [28, 46, 44, 25, 21, 40] and the Gaussian Mixture Model (GMM) [30, 36, 9, 10] (see Definitions 1.1 and 1.2).

Privacy in machine learning and statistical tasks has recently become of critical importance. New regulations, renewed consumer interest as well as privacy leaks, have led the major actors to adopt privacy-preserving solutions for the machine learning [1, 2, 3]. This new push has resulted in a flurry of activity in algorithm design for private machine learning, including very recently for SBMs and GMMs [55, 32, 16, 60]. Despite this activity, it has remain an open challenge to fully understand how privacy requirements impact model estimation problems and in particular their recovery thresholds and the computational complexity. This is the problem we tackle in this paper.

*Much of this work was done while the author was at ETH Zürich.

While other notions of privacy exist (e.g. k -anonymity), the de facto privacy standard is the differential privacy (DP) framework of Dwork, McSherry, Nissim, and Smith [24]. In this framework, the privacy quality is governed by two parameters, ϵ and δ , which in essence tell us how the probability of seeing a given output changes (both multiplicatively and additively) between two datasets that differ by any individual data element. This notion, in essence, quantifies the amount of information *leaked* by a given algorithm on individual data elements. The goal of the algorithm designer is to come up with differentially private algorithms for ϵ being a small constant and δ being of order $1/n^{\Theta(1)}$.

Differentially private analysis of graphs usually considers two notions of neighboring graphs. The weaker notion of *edge-DP* defines two graphs to be neighboring if they differ in one edge. Under the stronger notion of *node-DP*, two neighboring graphs can differ arbitrarily in the set of edges connected to a single vertex. Recently, there is a line of work on node-DP parameter estimation in random graph models, e.g. Erdős-Rényi models [61] and Graphons [12, 13]. However, for the more challenging task of graph clustering, node-DP is sometimes impossible to achieve.² Thus it is a natural first step to study edge-DP graph clustering.

Very recently, Seif, Nguyen, Vullikanti, and Tandon [55] were the first to propose differentially private algorithms for the Stochastic Block Model, with edge-DP. Concretely, they propose algorithms achieving exact recovery (exact identification of the planted clustering) while preserving privacy of individual edges of the graph. The proposed approach either takes $n^{\Theta(\log n)}$ time when ϵ is constant, or runs in polynomial time when ϵ is $\Omega(\log n)$.

Gaussian Mixture Models have also been studied in the context of differential privacy by [32, 16, 60] using the subsample-and-aggregate framework first introduced in [52] (see also recent work for robust moment estimation in the differential privacy setting [35, 8, 29]). The works of [32, 16] require an explicit bound on the euclidean norm of the centers as the sample complexity of these algorithms depends on this bound. For a mixture of k Gaussians, if there is a non-private algorithm that requires the minimum distance between the centers to be at least Δ , then [16, 60] can transform this non-private algorithm into a private one that needs the minimum distance between the centers to be at least $\Delta + \sqrt{\log n}$, where n is the number of samples.

In this paper, we tackle both clustering problems (graph clustering with the SBM and metric clustering with the GMM) through a new general privacy-preserving framework that brings us significantly closer to the state-of-the-art of non-private algorithms. As we will see, our new perspective on the problems appear to be easily extendable to many other estimation algorithms.

From robustness to privacy In recent years a large body of work (see [19, 22, 40, 18, 37, 14, 36, 30] and references therein) has advanced our understanding of parameter estimation in the presence of adversarial perturbations. In these settings, an adversary looks at the input instance and modifies it arbitrarily, under some constraints (these constraints are usually meant to ensure that it is still information theoretically possible to recover the underlying structure). As observed in the past [23, 35, 41], the two goals of designing privacy-preserving machine learning models and robust model estimation are tightly related. The common objective is to design algorithms that extract global information without over-relying on individual data samples.

Concretely, robust parameter estimation tends to morally follow a two-steps process: (i) argue that typical inputs are well-behaved, in the sense that they satisfy some property which can be used to accurately infer the desired global information, (ii) show that adversarial perturbations cannot significantly alter the quality of well-behaved inputs, so that it is still possible to obtain an accurate estimate. Conceptually, the analysis of private estimation algorithms can also be divided in two parts: *utility*, which is concerned with the accuracy of the output, and *privacy*, which ensures there is no leak of sensitive information. In particular, the canonical differential privacy definition can be interpreted as the requirement that, for any distinct inputs Y, Y' , the change in the output is *proportional* to the distance³ between Y and Y' .

It is easy to see this as a generalization of robustness: while robust algorithm needs the output to be stable for typical inputs, private algorithms requires this stability for *any possible input*. Then,

²In particular, we cannot hope to achieve *exact recovery*. We could isolate a vertex by removing all of its adjacent edges. Then it is impossible to cluster this vertex correctly. See below for formal definitions.

³The notion of distance is inherently application dependent. For example, it could be Hamming distance.

stability of the output immediately implies adding a small amount of noise to the output yields privacy. If the added noise is small enough, then utility is also preserved.

Our work further tightens this connection between robustness and privacy through a simple yet crucial insight: if two strongly convex functions over constrained sets –where both the function and the set may depend on the input– are point-wise close (say in a ℓ_2 -sense), their minimizers are also close (in the same sense). The alternative perspective is that projections of points that are close to each other, onto convex sets that are point-wise close, must also be close. This observation subsumes previously known sensitivity bounds in the empirical risk minimization literature (in particular in the output-perturbation approach to ERM, see Section 2 for a comparison).

The result is a clean, user-friendly, *framework to turn robust estimation algorithms into private algorithms*, while keeping virtually the same guarantees. We apply this paradigm to stochastic block models and Gaussian mixture models, which we introduce next.

Stochastic block model The stochastic block model is an extensively studied statistical model for community detection in graphs (see [4] for a survey).

Model 1.1 (Stochastic block model). *In its most basic form, the stochastic block model describes the distribution⁴ of an n -vertex graph $\mathbf{G} \sim \text{SBM}_n(d, \gamma, x)$, where x is a vector of n binary⁵ labels, $d \in \mathbb{N}$, $\gamma > 0$, and for every pair of distinct vertices $i, j \in [n]$ the edge $\{i, j\}$ is independently added to the graph \mathbf{G} with probability $(1 + \gamma \cdot x_i \cdot x_j) \frac{d}{n}$.*

For balanced label vector x , i.e., with roughly the same number of +1's and -1's, parameter d roughly corresponds to the average degree of the graph. Parameter γ corresponds to the *bias* introduced by the community structure. Note that for distinct vertices $i, j \in [n]$, the edge $\{i, j\}$ is present in \mathbf{G} with probability $(1 + \gamma) \frac{d}{n}$ if the vertices have the same label $x_i = x_j$ and with probability $(1 - \gamma) \frac{d}{n}$ if the vertices have different labels $x_i \neq x_j$.⁶

Given a graph \mathbf{G} sampled according to this model, the goal is to recover the (unknown) underlying vector of labels as well as possible. In particular, for a chosen algorithm returning a partition $\hat{x}(\mathbf{G}) \in \{\pm 1\}^n$, there are two main objective of interest: *weak recovery* and *exact recovery*. The former amounts to finding a partition $\hat{x}(\mathbf{G})$ correlated with the true partition. The latter instead corresponds to actually recovering the true partition with high probability. As shown in the following table, by now the statistical and computational landscape of these problems is well understood [17, 42, 48, 49, 28]:

	Objective	can be achieved (and efficiently so) iff
<i>weak recovery</i>	$\mathbb{P}_{\mathbf{G} \sim \text{SBM}_n(d, \gamma, x)} \left(\frac{1}{n} \langle x, \hat{x}(\mathbf{G}) \rangle \geq \Omega_{d, \gamma}(1) \right) \geq 1 - o(1)$	$\gamma^2 \cdot d \geq 1$
<i>exact recovery</i>	$\mathbb{P}_{\mathbf{G} \sim \text{SBM}_n(d, \gamma, x)} \left(\hat{x}(\mathbf{G}) \in \{x, -x\} \right) \geq 1 - o(1)$	$\frac{d}{\log n} \left(1 - \sqrt{1 - \gamma^2} \right) \geq 1$

Learning mixtures of spherical Gaussians The Gaussian Mixture Model we consider is the following.

Model 1.2 (Mixtures of spherical Gaussians). *Let D_1, \dots, D_k be Gaussian distributions on \mathbb{R}^d with covariance Id and means μ_1, \dots, μ_k satisfying $\|\mu_i - \mu_j\| \geq \Delta$ for any $i \neq j$. Given a set $\mathbf{Y} = \{y_1, \dots, y_n\}$ of n samples from the uniform mixture over D_1, \dots, D_k , estimate μ_1, \dots, μ_k .*

It is known that when the minimum separation is $\Delta = o(\sqrt{\log k})$, superpolynomially many samples are required to estimate the means up to small constant error [54]. Just above this threshold, at separation $k^{O(1/\gamma)}$ for any constant γ , there exist efficient algorithms based on the sum-of-squares hierarchy recovering the means up to accuracy $1/\text{poly}(k)$ [30, 36, 59]. In the regime

⁴We use **bold** characters to denote random variables.

⁵More general versions of the stochastic block model allow for more than two labels and general edge probabilities depending on the label assignment. However, many of the algorithmic phenomena of the general version can in their essence already be observed for the basic version that we consider in this work.

⁶At times we may write d_n, γ_n to emphasize that these may be functions of n . We write $o(1), \omega(1)$ for functions tending to zero (resp. infinity) as n grows.

where $\Delta = O(\sqrt{\log k})$ these algorithms yield the same guarantees but require quasipolynomial time. Recently, [39] showed how to efficiently recover the means as long as $\Delta = O(\log(k)^{1/2+c})$ for any constant $c > 0$.

1.1 Results

Stochastic block model We present here the first (ε, δ) -differentially private efficient algorithms for exact recovery. In all our results on stochastic block models, we consider the *edge privacy* model, in which two input graphs are adjacent if they differ on a single edge (cf. Definition C.1).

Theorem 1.3 (Private exact recovery of SBM). *Let $x \in \{\pm 1\}^n$ be balanced⁷. For any $\gamma, d, \varepsilon, \delta > 0$ satisfying*

$$\frac{d}{\log n} \left(1 - \sqrt{1 - \gamma^2}\right) \geq \Omega(1) \quad \text{and} \quad \frac{\gamma d}{\log n} \geq \Omega\left(\frac{1}{\varepsilon^2} \cdot \frac{\log(1/\delta)}{\log n} + \frac{1}{\varepsilon}\right),$$

there exists an (ε, δ) -differentially edge private algorithm that, on input $\mathbf{G} \sim \text{SBM}_n(d, \gamma, x)$, returns $\hat{x}(\mathbf{G}) \in \{x, -x\}$ with probability $1 - o(1)$. Moreover, the algorithm runs in polynomial time.

For any constant $\varepsilon > 0$, Theorem 1.3 states that (ε, δ) -differentially private exact recovery is possible, in polynomial time, already a constant factor close to the non-private threshold. Previous results [55] could only achieve comparable guarantees in time $O(n^{O(\log n)})$. It is also important to observe that the theorem provides a trade-off between signal-to-noise ratio of the instance (captured by the expression on the left-hand side with γ, d) and the privacy parameter ε . In particular, we highlight two regimes: for $d \geq \Omega(\log n)$ one can achieve exact recovery with high probability and privacy parameters $\delta = n^{-\Omega(1)}$, $\varepsilon = O(1/\gamma + 1/\gamma^2)$. For $d \geq \omega(\log n)$ one can achieve exact recovery with high probability and privacy parameters $\varepsilon = o(1)$, $\delta = n^{-\omega(1)}$. Theorem 1.3 follows by a result for private weak recovery and a boosting argument (cf. Theorem C.3 and Appendix C.2).

Further, we present a second, exponential-time, algorithm based on the exponential mechanism [43] which improves over the above in two regards. First, it gives *pure* differential privacy. Second, it provides utility guarantees for a larger range of graph parameters. In fact, we will also prove a lower bound which shows that its privacy guarantees are information theoretically optimal.⁸ All hidden constants are absolute and do not depend on any graph or privacy parameters unless stated otherwise. In what follows we denote by $\text{err}(\hat{x}, x)$ the minimum of the hamming distance of \hat{x} and x , and the one of $-\hat{x}$ and x , divided by n .

Theorem 1.4 (Slightly informal, see Theorem C.18 in the supplements for full version). *Let $\gamma\sqrt{d} \geq \Omega(1)$, $x \in \{\pm 1\}^n$ be balanced, and $\zeta \geq \exp(-\Omega(\gamma^2 d))$. For any $\varepsilon \geq \Omega\left(\frac{\log(1/\zeta)}{\gamma d}\right)$, there exists an algorithm which on input $\mathbf{G} \sim \text{SBM}_n(\gamma, d, x)$ outputs an estimate $\hat{x}(\mathbf{G}) \in \{\pm 1\}^n$ satisfying $\text{err}(\hat{x}(\mathbf{G}), x) \leq \zeta$ with probability at least $1 - \zeta$. In addition, the algorithm is ε -differentially edge private. Further, we can achieve error $\Theta\left(1/\sqrt{\log(1/\zeta)}\right)$ with probability $1 - e^{-n}$.*

A couple of remarks are in order. First, our algorithm works across all degree-regimes in the literature and matches known non-private thresholds and rates up to constants.⁹ In particular, for $\gamma^2 d = \Theta(1)$, we achieve weak/partial recovery with either constant or exponentially high success probability. Recall that the optimal non-private threshold is $\gamma^2 d > 1$. For the regime, where $\gamma^2 d = \omega(1)$, it is known that the optimal error rate is $\exp(-(1 - o(1))\gamma^2 d)$ [64] even non-privately which we match up to constants - here $o(1)$ denotes a function that tends to zero as $\gamma^2 d$ tends to infinity. Moreover, our algorithm achieves exact recovery as soon as $\gamma^2 d = \Omega(\log n)$ since then $\zeta < \frac{1}{n}$. This also matches known non-private thresholds up to constants [5, 47]. We remark that [55] gave an ε -DP exponential time algorithm which achieved exact recovery and has inverse polynomial success probability in the utility case as long as $\varepsilon \geq \Omega\left(\frac{\log n}{\gamma d}\right)$. We recover this result as a special case (with slightly worse constants). In fact, their algorithm is also based on the exponential mechanism, but their analysis only applies to the setting of exact recovery, while our result holds much more generally. Another

⁷A vector $x \in \{\pm 1\}^n$ is said to be balanced if $\sum_{i=1}^n x_i = 0$.

⁸It is optimal in the "small error" regime, otherwise it is almost optimal. See Theorem C.22 for more detail.

⁹For ease of exposition we did not try to optimize these constants.

crucial difference is that we show how to privatize a known boosting technique frequently used in the non-private setting, allowing us to achieve error guarantees which are optimal up to constant factors.

It is natural to ask whether, for a given set of parameters γ, d, ζ one can obtain better privacy guarantees than Theorem 1.4. Our next result implies that our algorithmic guarantees are almost tight.

Theorem 1.5 (Informal, see Theorem C.22 in the supplements for full version). *Suppose there exists an ε -differentially edge private algorithm such that for any balanced $x \in \{\pm 1\}^n$, on input $\mathbf{G} \sim \text{SBM}_n(d, \gamma, x)$, outputs $\hat{x}(\mathbf{G}) \in \{\pm 1\}^n$ satisfying*

$$\mathbb{P}(\text{err}(\hat{x}(\mathbf{G}), x) < \zeta) \geq 1 - \eta.$$

Then,

$$\varepsilon \geq \Omega\left(\frac{\log(1/\zeta)}{\gamma d} + \frac{\log(1/\eta)}{\zeta n \gamma d}\right). \quad (1.1)$$

This lower bound is tight for ε -DP exact recovery. By setting $\zeta = 1/n$ and $\eta = 1/\text{poly}(n)$, Theorem 1.5 implies no ε -DP exact recovery algorithm exists for $\varepsilon \leq O(\frac{\log n}{\gamma d})$. There exist ε -DP algorithms (Algorithm C.19 in the supplements and the algorithm in [55]) exactly recover the community for any $\varepsilon \geq \Omega(\frac{\log n}{\gamma d})$.

Notice Theorem 1.5 is a lower bound for a large range of error rates (partial to exact recovery). For failure probability $\eta = \zeta$, the lower bound simplifies to $\varepsilon \geq \Omega(\frac{\log(1/\zeta)}{\gamma d})$ and hence matches Theorem 1.4 up to constants. For exponentially small failure probability, $\eta = e^{-n}$, it becomes $\varepsilon \geq \Omega(\frac{1}{\zeta \gamma d})$. To compare, Theorem 1.4 requires $\varepsilon \geq \Omega(\frac{1}{\zeta^2 \gamma d})$ in this regime, using the substitution $\sqrt{\log(1/\zeta)} \rightarrow \zeta$.

Further, while formally incomparable, this ε -DP lower bound also suggests that the guarantees obtained by our efficient (ε, δ) -DP algorithm in Theorem 1.3 might be close to optimal. Note that setting $\zeta = \frac{1}{n}$ in Theorem 1.5 requires the algorithm to exactly recover the partitioning. In this setting, Theorem 1.3 implies that there is an efficient $(\varepsilon, n^{-\Theta(1)})$ -DP exact recovery algorithm for $\varepsilon \leq O(\sqrt{\frac{\log n}{\gamma d}})$. Theorem 1.5 states any ε -DP exact recovery algorithm requires $\varepsilon \geq \Omega(\frac{\log n}{\gamma d})$. Further, for standard privacy parameters that are required for real-world applications, such as $\varepsilon \approx 1$ and $\delta = n^{-10}$, Theorem 1.3 requires that $\gamma d \geq \Omega(\log n)$. Theorem 1.5 shows that for pure-DP algorithms with the same setting of ε this is also necessary. We leave it as fascinating open questions to bridge the gap between upper and lower bounds in the context of (ε, δ) -DP.

Learning mixtures of spherical Gaussians Our algorithm for privately learning mixtures of k spherical Gaussians provides statistical guarantees matching those of the best known non-private algorithms.

Theorem 1.6 (Privately learning mixtures of spherical Gaussians). *Consider an instance of Model 1.2. Let $t > 0$ be such that $\Delta \geq O(\sqrt{tk}^{1/t})$. For $n \geq \Omega(k^{O(1)} \cdot d^{O(t)})$, $k \geq (\log n)^{1/5}$, there exists an algorithm, running in time $(nd)^{O(t)}$, that outputs vectors $\hat{\mu}_1, \dots, \hat{\mu}_k$ satisfying*

$$\max_{\ell \in [k]} \|\hat{\mu}_\ell - \mu_{\pi(\ell)}\|_2 \leq O(k^{-12}),$$

with high probability, for some permutation $\pi : [k] \rightarrow [k]$. Moreover, for $\varepsilon \geq k^{-10}$, $\delta \geq n^{-10}$, the algorithm is (ε, δ) -differentially private¹⁰ for any input Y .

The conditions $\varepsilon \geq k^{-10}$, $\delta \geq n^{-10}$ in Theorem 1.6 are not restrictive and should be considered a formality. Moreover, setting $\varepsilon = 0.01$ and $\delta = n^{-10}$ already provides meaningful privacy guarantees in practice. The condition that $k \geq (\log n)^{1/5}$ is a technical requirement by our proofs.

Prior to this work, known differentially private algorithms could learn a mixture of k -spherical Gaussian either if: (1) they were given a ball of radius R containing all centers [32, 16];¹¹ or (2) the minimum separation between centers needs an additional additive $\Omega(\sqrt{\log n})$ term [16, 60]¹².

¹⁰Two input datasets are adjacent if they differ on a single sample. See Definition D.1 in the supplements.

¹¹In [32, 16] the sample complexity of the algorithm depends on this radius R .

¹²For $k \leq n^{o(1)}$ our algorithm provides a significant improvement as $\sqrt{\log k} = o(\sqrt{\log n})$.

To the best of our knowledge, Theorem 1.6 is the first to get the best of both worlds. That is, our algorithm requires no explicit upper bounds on the means (this also means the sample complexity does not depend on R) and only minimal separation assumption $O(\sqrt{\log k})$. Furthermore, we remark that while previous results only focused on mixtures of Gaussians, our algorithm also works for the significantly more general class of mixtures of Poincaré distributions. Concretely, in the regime $k \geq \sqrt{\log d}$, our algorithm recovers the state-of-the-art guarantees provided by non-private algorithms which are based on the sum-of-squares hierarchy [36, 30, 59]:¹³

- If $\Delta \geq k^{1/t^*}$ for some constant t^* , then by choosing $t \geq \Omega(t^*)$ the algorithm recovers the centers, up to a $1/\text{poly}(k)$ error, in time $\text{poly}(k, d)$ and using only $\text{poly}(k, d)$ samples.
- If $\Delta \geq \Omega(\sqrt{\log k})$ then choosing $t = O(\log k)$ the algorithm recovers the centers, up to a $1/\text{poly}(k)$ error, in quasi-polynomial time $\text{poly}(k^{O(t)}, d^{O(t^2)})$ and using a quasi-polynomial number of samples $\text{poly}(k, d^{O(t)})$.

For simplicity of exposition we will limit the presentation to mixtures of spherical Gaussians. We reiterate that separation $\Omega(\sqrt{\log k})$ is information-theoretically necessary for algorithms with polynomial sample complexity [54].

Subsequently and independently of our work, the work of [6] gives an algorithm that turns any non-private GMM learner into a private one based on the subsample and aggregate framework. They apply this reduction to the classical result of [45] to give the first finite-sample (ϵ, δ) -DP algorithm that learns mixtures of unbounded Gaussians, in particular, the covariance matrices of their mixture components can be arbitrary.

2 Techniques

We present here our general tools for designing efficient private estimation algorithms in the high-dimensional setting whose statistical guarantees almost match those of the best known non-private algorithms. The algorithms we design have the following structure in common: First, we solve a convex optimization problem with constraints and objective function depending on our input Y . Second, we round the optimal solution computed in the first step to a solution X for the statistical estimation problem at hand.

We organize our privacy analyses according to this structure. In order to analyze the first step, we prove a simple sensitivity bound for strongly convex optimization problems, which bounds the ℓ_2 -sensitivity of the optimal solution in terms of a uniform sensitivity bound for the objective function and the feasible region of the optimization problem.

For bounded problems –such as recovery of stochastic block models– we use this sensitivity bound, in the second step, to show that introducing small additive noise to standard rounding algorithms is enough to achieve privacy.

For unbounded problems –such as learning GMMs– we use this sensitivity bound to show that on adjacent inputs, either most entries of X only change slightly, as in the bounded case, or few entries vary significantly. We then combine different privacy techniques to hide both type of changes.

Privacy from sensitivity of strongly convex optimization problems Before illustrating our techniques with some examples, it is instructive to explicit our framework. Here we have a set of inputs \mathcal{Y} and a family of strongly convex functions $\mathcal{F}(\mathcal{Y})$ and convex sets $\mathcal{K}(\mathcal{Y})$ parametrized by these inputs. The generic *non-private* algorithm based on convex optimization we consider works as follows:

1. Compute $\hat{X} := \operatorname{argmin}_{X \in \mathcal{K}(Y)} f_Y(X)$;
2. Round \hat{X} into an integral solution.

For an estimation problem, a distributional assumption on \mathcal{Y} is made. Then one shows how, for typical inputs \mathbf{Y} sampled according to that distribution, the above scheme recovers the desired structured information.

¹³We remark that [39] give a polynomial time algorithm for separation $\Omega(\log(k)^{1/2+c})$ for constant $c > 0$ in the non-private setting but for a less general class of mixture distributions.

We can provide a privatized version of this scheme by arguing that, under reasonable assumptions on $\mathcal{F}(\mathcal{Y})$ and $\mathcal{K}(\mathcal{Y})$, the output of the function $\operatorname{argmin}_{X \in \mathcal{K}(Y)} f_Y(X)$ has low ℓ_2 -sensitivity. The consequence of this crucial observation is that one can combine the rounding step 2 with some standard privacy mechanism and achieve differential privacy. That is, the second step becomes:

2. Add random noise \mathbf{N} and round $\hat{X} + \mathbf{N}$ into an integral solution.

Our sensitivity bound is simple, yet it generalizes previously known bounds for strongly convex optimization problems (we provide a detailed comparison later in the section). For adjacent $Y, Y' \in \mathcal{Y}$, it requires the following properties:

- (i) For each $X \in \mathcal{K}(Y) \cap \mathcal{K}(Y')$ it holds $|f_Y(X) - f_{Y'}(X)| \leq \alpha$;
- (ii) For each $X \in \mathcal{K}(Y)$ its projection Z onto $\mathcal{K}(Y) \cap \mathcal{K}(Y')$ satisfies $|f_Y(X) - f_{Y'}(Z)| \leq \alpha$.

Here we think of α as some small quantity (relatively to the problem parameters). Notice, we may think of (i) as Lipschitz-continuity of the function $g(Y, X) = f_Y(X)$ with respect to Y and of (ii) as a bound on the change of the constrained set on adjacent inputs. In fact, these assumptions are enough to conclude low ℓ_2 sensitivity. Let \hat{X} and \hat{X}' be the outputs of the first step on inputs Y, Y' . Then using (i) and (ii) above and the fact that \hat{X} is an optimizer, we can show that there exists $Z \in \mathcal{K}(Y) \cap \mathcal{K}(Y')$ such that

$$|f_Y(\hat{X}) - f_Y(Z)| + |f_{Y'}(\hat{X}') - f_{Y'}(Z)| \leq O(\alpha).$$

By κ -strong convexity of $f_Y, f_{Y'}$ this implies

$$\left\| \hat{X} - Z \right\|_2^2 + \left\| \hat{X}' - Z \right\|_2^2 \leq O(\alpha/\kappa)$$

which ultimately means $\|\hat{X} - \hat{X}'\|_2^2 \leq O(\alpha/\kappa)$ (see Lemma B.1). Thus, starting from our assumptions on the point-wise distance of $f_Y, f_{Y'}$ we were able to conclude low ℓ_2 -sensitivity of our output!

A simple application: weak recovery of stochastic block models The ideas introduced above, combined with existing algorithms for weak recovery of stochastic block models, immediately imply a private algorithm for the problem. To illustrate this, consider Model 1.1 with parameters $\gamma^2 d \geq C$, for some large enough constant $C > 1$. Let $x \in \{\pm 1\}^n$ be balanced. Here Y is an n -by- n matrix corresponding to the rescaled centered adjacency matrix of the input graph:

$$Y_{ij} = \begin{cases} \frac{1}{\gamma d} \left(1 - \frac{d}{n}\right) & \text{if } ij \in E(G) \\ -\frac{1}{\gamma n} & \text{otherwise.} \end{cases}$$

The basic semidefinite program [28, 46] can be recast¹⁴ as the strongly convex constrained optimization question of finding the orthogonal projection of the matrix Y onto the set $\mathcal{K} := \{X \in \mathbb{R}^{n \times n} \mid X \succeq \mathbf{0}, X_{ii} = \frac{1}{n} \forall i\}$. That is

$$\hat{X} := \operatorname{argmin}_{X \in \mathcal{K}} \|Y - X\|_F^2.$$

Let $f_Y(X) := \|X\|_F^2 - 2\langle X, Y \rangle$ and notice that $\hat{X} = \operatorname{argmin}_{X \in \mathcal{K}} f_Y(X)$. It is a standard fact that, if our input was $\mathbf{G} \sim \operatorname{SBM}_n(d, \gamma, x)$, then with high probability $\hat{X}(\mathbf{G}) = \operatorname{argmin}_{X \in \mathcal{K}} f_{Y(\mathbf{G})}(X)$ would have leading eigenvalue-eigenvector pair satisfying

$$\lambda_1(\mathbf{G}) \geq 1 - O(1/\gamma^2 d) \quad \text{and} \quad \langle v_1(\mathbf{G}), x/\|x\| \rangle^2 \geq 1 - O(1/\gamma^2 d).$$

This problem fits perfectly the description of the previous paragraph. Note that since the constraint set does not depend on Y , Property (ii) reduces to Property (i). Thus, it stands to reason that the projections \hat{X}, \hat{X}' of Y, Y' are close whenever the input graphs generating Y and Y' are adjacent. By Hölder's Inequality with the entry-wise infinity and ℓ_1 -norm, we obtain $|f_Y(X) - f_{Y'}(X)| \leq 2\|X\|_\infty \|Y - Y'\|_1$. By standard facts about positive semidefinite matrices, we have $\|X\|_\infty \leq \frac{1}{n}$

¹⁴The objective function in [28, 46] is linear in X instead of quadratic. However, both programs have similar utility guarantees and the utility proof of our program is an adaption of that in [28] (see Lemma C.8). We use the quadratic objective function to achieve privacy via strong convexity.

for all $X \in \mathcal{K}$. Also, Y and Y' can differ on at most 2 entries and hence $\|Y - Y'\|_1 \leq O(\frac{1}{\gamma d})$. Thus, $\|\hat{X} - \hat{X}'\|_F^2 \leq O(\frac{1}{n\gamma d})$.

The rounding step is now straightforward. Using the Gaussian mechanism we return the leading eigenvector of $\hat{X} + \mathbf{N}$ where $\mathbf{N} \sim N(0, \frac{1}{n\gamma d} \cdot \frac{\log(1/\delta)}{\epsilon^2})^{n \times n}$. This matrix has Frobenius norm significantly larger than \hat{X} but its spectral norm is only

$$\|\mathbf{N}\| \leq \frac{\sqrt{n \log(1/\delta)}}{\epsilon} \cdot \sqrt{\frac{1}{n\gamma d}} \leq \frac{1}{\epsilon} \cdot \sqrt{\frac{\log(1/\delta)}{\gamma d}}.$$

Thus by standard linear algebra, for typical instances $\mathbf{G} \sim \text{SBM}_n(d, \gamma, x)$, the leading eigenvector of $\hat{X}(\mathbf{G}) + \mathbf{N}$ will be highly correlated with the true community vector x whenever the average degree d is large enough. In conclusion, a simple randomized rounding step is enough!

Remark 2.1 (From weak recovery to exact recovery). *In the non-private setting, given a weak recovery algorithm for the stochastic block model, one can use this as an initial estimate for a boosting procedure based on majority voting to achieve exact recovery. We show that this can be done privately. See Appendix C.2 in the supplements.*

An advanced application: learning mixtures of Gaussians In the context of SBMs our argument greatly benefited from two key properties: first, on adjacent inputs $Y - Y'$ was bounded in an appropriate norm; and second, the convex set \mathcal{K} was fixed. In the context of learning mixtures of spherical Gaussians as in Model 1.2, *both* this properties are *not* satisfied (notice how one of this second properties would be satisfied assuming bounded centers!). So additional ideas are required.

The first observation, useful to overcome the first obstacle, is that before finding the centers, one can first find the n -by- n membership matrix $W(Y)$ where $W(Y)_{ij} = 1$ if i, j were sampled from the same mixture component and 0 otherwise. The advantage here is that, on adjacent inputs, $\|W(Y) - W(Y')\|_F^2 \leq 2n/k$ and thus one recovers the first property.¹⁵ Here early sum-of-squares algorithms for the problem [30, 36] turns out to be convenient as they rely on minimizing the function $\|W\|_F^2$ subject to the following system of polynomial inequalities in variables $z_{11}, \dots, z_{1k}, \dots, z_{nk}$, with $W_{ij} = \sum_{\ell} z_{i\ell} z_{j\ell}$ for all $i, j \in [n]$ and a parameter $t > 0$.

$$\left. \begin{cases} z_{i\ell}^2 = z_{i\ell} & \forall i \in [n], \ell \in [k] \quad (\text{indicators}) \\ \sum_{\ell \in [k]} z_{i\ell} \leq 1 & \forall i \in [n] \quad (\text{cluster membership}) \\ z_{i\ell} \cdot z_{i\ell'} = 0 & \forall i \in [n], \ell \in [k] \quad (\text{unique membership}) \\ \sum_i z_{i\ell} = n/k & \forall \ell \in [k] \quad (\text{size of clusters})^{16} \\ \mu'_\ell = \frac{k}{n} \sum_i z_{i\ell} \cdot y_i & \forall \ell \in [k] \quad (\text{means of clusters}) \\ \frac{k}{n} \sum_i z_{i\ell} \langle y_i - \mu'_\ell, u \rangle^{2t} \leq (2t)^t \cdot \|u\|_2^t & \forall u \in \mathbb{R}^d, \ell \in [k] \quad (\text{subgaussianity of } t\text{-moment}) \end{cases} \right\} (\mathcal{P}(Y))$$

For the scope of this discussion,¹⁷ we may disregard computational issues and assume we have access to an algorithm returning a point from the convex hull $\mathcal{K}(Y)$ of all solutions to our system of inequalities.¹⁸ Each indicator variable $z_{i\ell} \in \{0, 1\}$ is meant to indicate whether sample y_i is believed

¹⁵Notice for typical inputs \mathbf{Y} from Model 1.2 one expect $\|W(\mathbf{Y})\|_F^2 \approx n^2/k$.

¹⁶Formally, we would replace the constraint on the size of the clusters by one which requires them to be of size $(1 \pm \alpha) \frac{n}{k}$, for some small α .

¹⁷While this is far from being true, it turns out that having access to a pseudo-distribution satisfying $\mathcal{P}(Y)$ is enough for our subsequent argument to work, albeit with some additional technical work required.

¹⁸We remark that a priori it is also not clear how to encode the subgaussian constraint in a way that we could recover a degree- t pseudo-distribution satisfying $\mathcal{P}(Y)$ in polynomial time. By now this is well understood, we discuss this in Appendix A in the supplements.

to be in cluster C_ℓ . In the non-private setting, the idea behind the program is that –for typical \mathbf{Y} sampled according to Model 1.2 with minimum separation $\Delta \geq k^{1/t} \sqrt{t}$ – any solution $W(\mathbf{Y}) \in \mathcal{K}(\mathbf{Y})$ is close to the ground truth matrix $W^*(\mathbf{Y})$ in Frobenius norm: $\|W(\mathbf{Y}) - W^*(\mathbf{Y})\|_F^2 \leq 1/\text{poly}(k)$. Each row $W(\mathbf{Y})_i$ may be seen as inducing a uniform distribution over a subset of \mathbf{Y} .¹⁹ Combining the above bound with the fact that subgaussian distributions at small total variation distance have means that are close, we conclude the algorithm recovers the centers of the mixture.

While this program suggests a path to recover the first property, it also possesses a fatal flaw: the projection W' of $W \in \mathcal{K}(Y)$ onto $\mathcal{K}(Y) \cap \mathcal{K}(Y')$ may be *far* in the sense that $|\|W\|_F^2 - \|W'\|_F^2| \geq \Omega(\|W\|_F^2 + \|W'\|_F^2) \geq \Omega(n^2/k)$. The reason behind this phenomenon can be found in the constraint $\sum_i z_{i\ell} = n/k$. The set indicated by the vector $(z_{1\ell} \dots, z_{n\ell})$ may be subgaussian in the sense of $\mathcal{P}(Y)$ for input Y but, upon changing a single sample, this may no longer be true. We work around this obstacle in two steps:

1. We replace the above constraint with $\sum_i z_{i\ell} \leq n/k$.
2. We compute $\hat{W} := \text{argmin}_{W \text{ solving } \mathcal{P}(Y)} \|J - W\|_F^2$, where J is the all-one matrix.²⁰

The catch now is that the program is satisfiable for *any* input Y since we can set $z_{i\ell} = 0$ whenever necessary. Moreover, we can guarantee property (ii) (required by our sensitivity argument) for $\alpha \leq O(n/k)$, since we can obtain $W' \in \mathcal{K}(Y) \cap \mathcal{K}(Y')$ simply zeroing out the row/column in W corresponding to the sample differing in Y and Y' . Then for typical inputs \mathbf{Y} , the correlation with the true solution is guaranteed by the new strongly convex objective function.

We offer some more intuition on the choice of our objective function: Recall that W_{ij} indicates our guess whether the i -th and j -th datapoints are sampled from the same Gaussian component. A necessary condition for W to be close to its ground-truth counterpart W^* , is that they roughly have the same number of entries that are (close to) 1. One way to achieve this would be to add the lower bound constraint $\sum_\ell z_{i\ell} \gtrsim \frac{n}{k}$. However, such a constraint could cause privacy issues: There would be two neighboring datasets, such that the constraint set induced by one dataset is satisfiable, but the constraint set induced by the other dataset is not satisfiable. We avoid this issue by noticing that the appropriate number of entries close to 1 can also be induced by minimizing the distance of W to the all-one matrix. This step is also a key difference from [35], explained in more detail below.

From low sensitivity of the indicators to low sensitivity of the estimates For adjacent inputs Y, Y' let \hat{W}, \hat{W}' be respectively the matrices computed by the above strongly convex programs. Our discussion implies that, applying our sensitivity bound, we can show $\|\hat{W} - \hat{W}'\|_F^2 \leq O(n/k)$. The problem is that simply applying a randomized rounding approach here cannot work. The reason is that even though the vector \hat{W}_i induces a subgaussian distribution, the vector $\hat{W}_i + v$ for $v \in \mathbb{R}^n$, *might not*. Without the subgaussian constraint we cannot provide any meaningful utility bound. In other words, the root of our problem is that there exists heavy-tailed distributions that are arbitrarily close in total variation distance to any given subgaussian distribution.

On the other hand, our sensitivity bound implies $\|\hat{W} - \hat{W}'\|_1^2 \leq o(\|\hat{W}\|_1)$ and thus, all but a vanishing fraction of rows $i \in [n]$ must satisfy $\|\hat{W}_i - \hat{W}'_i\|_1 \leq o(\|\hat{W}_i\|_1)$. For each row i , let μ_i, μ'_i be the means of the distributions induced respectively by \hat{W}_i, \hat{W}'_i . We are thus in the following setting:

1. For a set of $(1 - o(1)) \cdot n$ good rows $\|\mu_i - \mu'_i\|_2 \leq o(1)$,
2. For the set \mathcal{B} of remaining bad rows, the distance $\|\mu_i - \mu'_i\|_2$ may be unbounded.

We hide differences of the first type as follows: pick a random subsample \mathcal{S} of $[n]$ of size n^c , for some small $c > 0$, and for each picked row use the Gaussian mechanism. The subsampling step is useful as it allows us to decrease the standard deviation of the entry-wise random noise by a factor n^{1-c} . We hide differences of the second type as follows: Note that most of the rows are clustered together in space. Hence, we aim to privately identify the regions which contain many of the rows.

¹⁹More generally, we may think of a vector $v \in \mathbb{R}^n$ as the vector inducing the distribution given by $v/\|v\|_1$ onto the set Y of n elements.

²⁰We remark that for technical reasons our function in Appendix D.1 in the supplements will be slightly different. We do not discuss it here to avoid obfuscating our main message.

Formally, we use a classic high dimensional (ϵ, δ) -differentially private histogram learner on \mathcal{S} and for the k largest bins of highest count privately return their average (cf. Lemma A.13). The crux of the argument here is that the cardinality of $\mathcal{B} \cap \mathcal{S}$ is sufficiently small that the privacy guarantees of the histogram learner can be extended even for inputs that differ in $|\mathcal{B} \cap \mathcal{S}|$ many samples. Finally, standard composition arguments will guarantee privacy of the whole algorithm.

Comparison with the framework of Kothari-Manurangsi-Velingker Both Kothari-Manurangsi-Velingker [35] and our work obtained private algorithms for high-dimensional statistical estimation problems by privatizing strongly convex programs, more specifically, sum-of-squares (SoS) programs. The main difference between KMV and our work lies in how we choose the SoS program. For the problem of robust moment estimation, KMV considered the canonical SoS program from [30, 36] which contains a minimum cardinality constraint (e.g., $\sum_l z_{il} \gtrsim \frac{n}{k}$ in the case of GMMs). Such a constraint is used to ensure good utility. However, as alluded to earlier, this is problematic for privacy: there will always exist two adjacent input datasets such that the constraints are satisfiable for one but not for the other. KMV and us resolve this privacy issue in different ways.

KMV uses an exponential mechanism to pick the lower bound of the minimum cardinality constraint. This step also ensures that solutions to the resulting SoS program will have low sensitivity. In contrast, we simply drop the minimum cardinality constraint. Then the resulting SoS program is always feasible for any input dataset! To still ensure good utility, we additionally pick an appropriate objective function. For example, in Gaussian mixture models, we chose the objective $\|W - J\|_F^2$. Our approach has the following advantages: First, the exponential mechanism in KMV requires computing $O(n)$ scores. Computing each score requires solving a large semidefinite program, which can significantly increase the running time. Second, proving that the exponential mechanism in KMV works requires several steps: 1) defining a (clever) score function, 2) bounding the sensitivity of this score function and, 3) showing existence of a large range of parameters with high score. Our approach bypasses both of these issues.

Further, as we show, our general recipe can be easily extended to other high dimensional problems of interest: construct a strongly convex optimization program and add noise to its solution. This can provide significant computational improvements. For example, in the context of SBMs, the framework of [35] would require one to sample from an exponential distribution over matrices. Constructing and sampling from such distributions is an expensive operation. However, it is well-understood that an optimal fractional solution to the basic SDP relaxation we consider can be found in *near quadratic time* using the standard matrix multiplicative weight method [7, 58], making the whole algorithm run in near-quadratic time. Whether our algorithm can be sped up to near-linear time, as in [7, 58], remains a fascinating open question.

Comparison with previous works on empirical risk minimization Results along the lines of the sensitivity bound described at the beginning of the section (see Lemma B.1 for a formal statement) have been extensively used in the context of empirical risk minimization [15, 34, 57, 11, 63, 50]. Most results focus on the special case of unconstrained optimization of strongly convex functions. In contrast, our sensitivity bound applies to the significantly more general settings where both the objective functions and the constrained set may depend on the input.²¹ Most notably for our settings of interest, [15] studied unconstrained optimization of (smooth) strongly convex functions depending on the input, with bounded gradient. We recover such a result for $X' = X$ in (ii). In [50], the authors considered constraint optimization of objective functions where the domain (but *not* the function) may depend on the input data. They showed how one can achieve differential privacy while optimize the desired objective function by randomly perturbing the constraints. It is important to remark that, in [50], the notion of utility is based on the optimization problem (and their guarantees are tight only up to logarithmic factors). In the settings we consider, even in the special case where f does not depend on the input, this notion of utility may not correspond to the notion of utility required by the estimation problem, and thus, the corresponding guarantees can turn out to be too loose to ensure the desired error bounds.

²¹The attentive reader may argue that one could cast convex optimization over a constrained domain as unconstrained optimization of a new convex function with the appropriate penalty terms. In practice however, this turns out to be hard to do for constraints such as Definition A.19.

Acknowledgments and Disclosure of Funding

This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement No 815464).

References

- [1] Tackling urban mobility with technology. <https://europe.googleblog.com/2015/11/tackling-urban-mobility-with-technology.html>, 2015. Accessed: 2022-11-06.
- [2] Learning with privacy at scale. <https://docs-assets.developer.apple.com/ml-research/papers/learning-with-privacy-at-scale.pdf>, 2017. Accessed: 2022-11-06.
- [3] Disclosure avoidance for the 2020 census: An introduction. <https://www2.census.gov/library/publications/decennial/2020/2020-census-disclosure-avoidance-handbook.pdf>, 2021. Accessed: 2022-11-06.
- [4] Emmanuel Abbe. Community detection and stochastic block models: recent developments. *The Journal of Machine Learning Research*, 18(1):6446–6531, 2017.
- [5] Emmanuel Abbe, Afonso S Bandeira, and Georgina Hall. Exact recovery in the stochastic block model. *IEEE Transactions on information theory*, 62(1):471–487, 2015.
- [6] Jamil Arbas, Hassan Ashtiani, and Christopher Liaw. Polynomial time and private learning of unbounded gaussian mixture models. *arXiv preprint arXiv:2303.04288*, 2023.
- [7] Sanjeev Arora and Satyen Kale. A combinatorial, primal-dual approach to semidefinite programs. In *Proceedings of the thirty-ninth annual ACM symposium on Theory of computing*, pages 227–236, 2007.
- [8] Hassan Ashtiani and Christopher Liaw. Private and polynomial time algorithms for learning gaussians and beyond. In Po-Ling Loh and Maxim Raginsky, editors, *Proceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pages 1075–1076. PMLR, 02–05 Jul 2022.
- [9] Ainesh Bakshi, Ilias Diakonikolas, Samuel B. Hopkins, Daniel Kane, Sushrut Karmalkar, and Pravesh K. Kothari. Outlier-robust clustering of gaussians and other non-spherical mixtures. In Sandy Irani, editor, *61st IEEE Annual Symposium on Foundations of Computer Science, FOCS 2020, Durham, NC, USA, November 16-19, 2020*, pages 149–159. IEEE, 2020.
- [10] Ainesh Bakshi, Ilias Diakonikolas, He Jia, Daniel M. Kane, Pravesh K. Kothari, and Santosh S. Vempala. Robustly learning mixtures of k arbitrary gaussians. In Stefano Leonardi and Anupam Gupta, editors, *STOC '22: 54th Annual ACM SIGACT Symposium on Theory of Computing, Rome, Italy, June 20 - 24, 2022*, pages 1234–1247. ACM, 2022.
- [11] Raef Bassily, Adam Smith, and Abhradeep Thakurta. Private empirical risk minimization: Efficient algorithms and tight error bounds. In *2014 IEEE 55th annual symposium on foundations of computer science*, pages 464–473. IEEE, 2014.
- [12] Christian Borgs, Jennifer Chayes, and Adam Smith. Private graphon estimation for sparse graphs. *Advances in Neural Information Processing Systems*, 28, 2015.
- [13] Christian Borgs, Jennifer Chayes, Adam Smith, and Ilias Zadik. Revealing network structure, confidentially: Improved rates for node-private graphon estimation. In *2018 IEEE 59th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 533–543. IEEE, 2018.
- [14] Moses Charikar, Jacob Steinhardt, and Gregory Valiant. Learning from untrusted data. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, pages 47–60, 2017.

- [15] Kamalika Chaudhuri, Claire Monteleoni, and Anand D Sarwate. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12(3), 2011.
- [16] Edith Cohen, Haim Kaplan, Yishay Mansour, Uri Stemmer, and Eliad Tsfadia. Differentially-private clustering of easy instances. In *International Conference on Machine Learning*, pages 2049–2059. PMLR, 2021.
- [17] Aurelien Decelle, Florent Krzakala, Cristopher Moore, and Lenka Zdeborová. Asymptotic analysis of the stochastic block model for modular networks and its algorithmic applications. *Physical Review E*, 84(6):066106, 2011.
- [18] Ilias Diakonikolas, Gautam Kamath, Daniel Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. Robust estimators in high-dimensions without the computational intractability. *SIAM Journal on Computing*, 2019.
- [19] Ilias Diakonikolas and Daniel M Kane. Recent advances in algorithmic high-dimensional robust statistics. *arXiv preprint arXiv:1911.05911*, 2019.
- [20] Ilias Diakonikolas, Daniel M. Kane, Sushrut Karmalkar, Ankit Pensia, and Thanasis Pittas. Robust sparse mean estimation via sum of squares. In Po-Ling Loh and Maxim Raginsky, editors, *Conference on Learning Theory, 2-5 July 2022, London, UK*, volume 178 of *Proceedings of Machine Learning Research*, pages 4703–4763. PMLR, 2022.
- [21] Jingqiu Ding, Tommaso d’Orsi, Rajai Nasser, and David Steurer. Robust recovery for stochastic block models. In *2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)*, pages 387–394. IEEE, 2022.
- [22] Tommaso d’Orsi, Pravesh K. Kothari, Gleb Novikov, and David Steurer. Sparse PCA: algorithms, adversarial perturbations and certificates. In Sandy Irani, editor, *61st IEEE Annual Symposium on Foundations of Computer Science, FOCS 2020, Durham, NC, USA, November 16-19, 2020*, pages 553–564. IEEE, 2020.
- [23] Cynthia Dwork and Jing Lei. Differential privacy and robust statistics. In *STOC’09—Proceedings of the 2009 ACM International Symposium on Theory of Computing*, pages 371–380. ACM, New York, 2009.
- [24] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam D. Smith. Calibrating noise to sensitivity in private data analysis. In Shai Halevi and Tal Rabin, editors, *Theory of Cryptography, Third Theory of Cryptography Conference, TCC 2006, New York, NY, USA, March 4-7, 2006, Proceedings*, volume 3876 of *Lecture Notes in Computer Science*, pages 265–284. Springer, 2006.
- [25] Yingjie Fei and Yudong Chen. Achieving the Bayes error rate in synchronization and block models by SDP, robustly. *IEEE Trans. Inform. Theory*, 66(6):3929–3953, 2020.
- [26] Noah Fleming, Pravesh Kothari, Toniann Pitassi, et al. Semialgebraic proofs and efficient algorithm design. *Foundations and Trends® in Theoretical Computer Science*, 14(1-2):1–221, 2019.
- [27] M. Grötschel, L. Lovász, and A. Schrijver. The ellipsoid method and its consequences in combinatorial optimization. *Combinatorica*, 1(2):169–197, 1981.
- [28] Olivier Guédon and Roman Vershynin. Community detection in sparse networks via Grothendieck’s inequality. *Probab. Theory Related Fields*, 165(3-4):1025–1049, 2016.
- [29] Samuel B Hopkins, Gautam Kamath, and Mahbod Majid. Efficient mean estimation with pure differential privacy via a sum-of-squares exponential mechanism. In *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1406–1417, 2022.
- [30] Samuel B. Hopkins and Jerry Li. Mixture models, robustness, and sum of squares proofs. In Ilias Diakonikolas, David Kempe, and Monika Henzinger, editors, *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2018, Los Angeles, CA, USA, June 25-29, 2018*, pages 1021–1034. ACM, 2018.

- [31] William B Johnson. Extensions of lipschitz mappings into a hilbert space. *Contemp. Math.*, 26:189–206, 1984.
- [32] Gautam Kamath, Or Sheffet, Vikrant Singhal, and Jonathan Ullman. Differentially private algorithms for learning mixtures of separated gaussians. *Advances in Neural Information Processing Systems*, 32, 2019.
- [33] Vishesh Karwa and Salil P. Vadhan. Finite sample differentially private confidence intervals. In Anna R. Karlin, editor, *9th Innovations in Theoretical Computer Science Conference, ITCS 2018, January 11-14, 2018, Cambridge, MA, USA*, volume 94 of *LIPICs*, pages 44:1–44:9. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2018.
- [34] Daniel Kifer, Adam Smith, and Abhradeep Thakurta. Private convex empirical risk minimization and high-dimensional regression. In *Conference on Learning Theory*, pages 25–1. JMLR Workshop and Conference Proceedings, 2012.
- [35] Pravesh Kothari, Pasin Manurangsi, and Ameya Velingker. Private robust estimation by stabilizing convex relaxations. In Po-Ling Loh and Maxim Raginsky, editors, *Conference on Learning Theory, 2-5 July 2022, London, UK*, volume 178 of *Proceedings of Machine Learning Research*, pages 723–777. PMLR, 2022.
- [36] Pravesh K. Kothari, Jacob Steinhardt, and David Steurer. Robust moment estimation and improved clustering via sum of squares. In Ilias Diakonikolas, David Kempe, and Monika Henzinger, editors, *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2018, Los Angeles, CA, USA, June 25-29, 2018*, pages 1035–1046. ACM, 2018.
- [37] Kevin A. Lai, Anup B. Rao, and Santosh S. Vempala. Agnostic estimation of mean and covariance. In Irit Dinur, editor, *IEEE 57th Annual Symposium on Foundations of Computer Science, FOCS 2016, 9-11 October 2016, Hyatt Regency, New Brunswick, New Jersey, USA*, pages 665–674. IEEE Computer Society, 2016.
- [38] Jean B. Lasserre. New positive semidefinite relaxations for nonconvex quadratic programs. In *Advances in convex analysis and global optimization (Pythagorion, 2000)*, volume 54 of *Nonconvex Optim. Appl.*, pages 319–331. Kluwer Acad. Publ., Dordrecht, 2001.
- [39] Allen Liu and Jerry Li. Clustering mixtures with almost optimal separation in polynomial time. In *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1248–1261, 2022.
- [40] Allen Liu and Ankur Moitra. Minimax rates for robust community detection. *CoRR*, abs/2207.11903, 2022.
- [41] Xiyang Liu, Weihao Kong, and Sewoong Oh. Differential privacy and robust statistics in high dimensions. In *Conference on Learning Theory*, pages 1167–1246. PMLR, 2022.
- [42] Laurent Massoulié. Community detection thresholds and the weak ramanujan property. In *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, pages 694–703, 2014.
- [43] Frank McSherry and Kunal Talwar. Mechanism design via differential privacy. In *48th Annual IEEE Symposium on Foundations of Computer Science (FOCS'07)*, pages 94–103. IEEE, 2007.
- [44] Ankur Moitra, William Perry, and Alexander S Wein. How robust are reconstruction thresholds for community detection? In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pages 828–841, 2016.
- [45] Ankur Moitra and Gregory Valiant. Settling the polynomial learnability of mixtures of gaussians. In *2010 IEEE 51st Annual Symposium on Foundations of Computer Science*, pages 93–102. IEEE, 2010.
- [46] Andrea Montanari and Subhabrata Sen. Semidefinite programs on sparse random graphs and their application to community detection. In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pages 814–827, 2016.

- [47] Elchanan Mossel, Joe Neeman, and Allan Sly. Consistency thresholds for the planted bisection model. In *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, pages 69–75, 2015.
- [48] Elchanan Mossel, Joe Neeman, and Allan Sly. Reconstruction and estimation in the planted partition model. *Probability Theory and Related Fields*, 162(3):431–461, 2015.
- [49] Elchanan Mossel, Joe Neeman, and Allan Sly. A proof of the block model threshold conjecture. *Combinatorica*, 38(3):665–708, 2018.
- [50] Andres Munoz, Umar Syed, Sergei Vassilvtiskii, and Ellen Vitercik. Private optimization without constraint violations. In *International Conference on Artificial Intelligence and Statistics*, pages 2557–2565. PMLR, 2021.
- [51] Yurii Nesterov. Squared functional systems and optimization problems. In *High performance optimization*, volume 33 of *Appl. Optim.*, pages 405–440. Kluwer Acad. Publ., Dordrecht, 2000.
- [52] Kobbi Nissim, Sofya Raskhodnikova, and Adam Smith. Smooth sensitivity and sampling in private data analysis. In *STOC’07—Proceedings of the 39th Annual ACM Symposium on Theory of Computing*, pages 75–84. ACM, New York, 2007.
- [53] Pablo A Parrilo. *Structured semidefinite programs and semialgebraic geometry methods in robustness and optimization*. PhD thesis, California Institute of Technology, 2000.
- [54] Oded Regev and Aravindan Vijayaraghavan. On learning mixtures of well-separated gaussians. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 85–96. IEEE, 2017.
- [55] Mohamed M. Seif, Dung Nguyen, Anil Vullikanti, and Ravi Tandon. Differentially private community detection for stochastic block models. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato, editors, *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 15858–15894. PMLR, 2022.
- [56] N. Z. Shor. Quadratic optimization problems. *Izv. Akad. Nauk SSSR Tekhn. Kibernet.*, (1):128–139, 222, 1987.
- [57] Shuang Song, Kamalika Chaudhuri, and Anand D Sarwate. Stochastic gradient descent with differentially private updates. In *2013 IEEE global conference on signal and information processing*, pages 245–248. IEEE, 2013.
- [58] David Steurer. Fast sdp algorithms for constraint satisfaction problems. In *Proceedings of the twenty-first annual ACM-SIAM symposium on Discrete Algorithms*, pages 684–697. SIAM, 2010.
- [59] David Steurer and Stefan Tiegel. Sos degree reduction with applications to clustering and robust moment estimation. In *Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 374–393. SIAM, 2021.
- [60] Eliad Tsfadia, Edith Cohen, Haim Kaplan, Yishay Mansour, and Uri Stemmer. Friendlycore: Practical differentially private aggregation. In *International Conference on Machine Learning*, pages 21828–21863. PMLR, 2022.
- [61] Jonathan Ullman and Adam Sealfon. Efficiently estimating erdos-renyi graphs with node differential privacy. *Advances in Neural Information Processing Systems*, 32, 2019.
- [62] Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2019.
- [63] Di Wang, Minwei Ye, and Jinhui Xu. Differentially private empirical risk minimization revisited: Faster and more general. *Advances in Neural Information Processing Systems*, 30, 2017.
- [64] Anderson Y Zhang and Harrison H Zhou. Minimax rates of community detection in stochastic block models. *The Annals of Statistics*, 44(5):2252–2280, 2016.

A Preliminaries

We use **boldface** characters for random variables. We hide multiplicative factors *logarithmic* in n using the notation $\tilde{O}(\cdot), \tilde{\Omega}(\cdot)$. Similarly, we hide absolute constant multiplicative factors using the standard notation $O(\cdot), \Omega(\cdot), \Theta(\cdot)$. Often times we use the letter C to denote universal constants independent of the parameters at play. We write $o(1), \omega(1)$ for functions tending to zero (resp. infinity) as n grows. We say that an event happens with high probability if this probability is at least $1 - o(1)$. Throughout the paper, when we say "an algorithm runs in time $O(q)$ " we mean that the number of basic arithmetic operations involved is $O(q)$. That is, we ignore bit complexity issues.

Vectors, matrices, tensors We use Id_n to denote the n -by- n dimensional identity matrix, $J_n \in \mathbb{R}^{n \times n}$ the all-ones matrix and $\mathbf{0}_n, \mathbf{1}_n \in \mathbb{R}^n$ to denote respectively the zero and the all-ones vectors. When the context is clear we drop the subscript. For matrices $A, B \in \mathbb{R}^{n \times n}$ we write $A \succeq B$ if $A - B$ is positive semidefinite. For a matrix M , we denote its eigenvalues by $\lambda_1(M), \dots, \lambda_n(M)$, we simply write λ_i when the context is clear. We denote by $\|M\|$ the spectral norm of M . We denote by $\mathbb{R}^{d^{\otimes t}}$ the set of real-valued order- t tensors. for a $d \times d$ matrix M , we denote by $M^{\otimes t}$ the t -fold Kronecker product $\underbrace{M \otimes M \otimes \dots \otimes M}_{t \text{ times}}$. We define the *flattening*, or *vectorization*, of M to

be the d^t -dimensional vector, whose entries are the entries of M appearing in lexicographic order. With a slight abuse of notation we refer to this flattening with M , ambiguities will be clarified from context. We denote by $N(0, \sigma^2)^{d^{\otimes t}}$ the distribution over Gaussian tensors with d^t entries with standard deviation σ . Given $u, v \in \{\pm 1\}^n$, we use $\text{Ham}(u, v) := \sum_{i=1}^n \mathbb{1}_{[u_i \neq v_i]}$ to denote their Hamming distance. Given a vector $u \in \mathbb{R}^n$, we let $\text{sign}(u) \in \{\pm 1\}^n$ denote its sign vector. A vector $u \in \{\pm 1\}^n$ is said to be *balanced* if $\sum_{i=1}^n u_i = 0$.

Graphs We consider graphs on n vertices and let \mathcal{G}_n be the set of all graphs on n vertices. For a graph G on n vertices we denote by $A(G) \in \mathbb{R}^{n \times n}$ its adjacency matrix. When the context is clear we simply write A . Let $V(G)$ (resp. $E(G)$) denote the vertex (resp. edge) set of graph G . Given two graphs G, H on the same vertex set V , let $G \setminus H := (V, E(G) \setminus H(G))$. Given a graph $H, H' \subseteq H$ means H' is a subgraph of H such that $V(H') = V(H)$ and $E(H') \subseteq E(H)$. The Hamming distance between two graphs G, H is defined to be the size of the symmetric difference between their edge sets, i.e. $\text{Ham}(G, H) := |E(G) \Delta E(H)|$.

A.1 Differential privacy

In this section we introduce standard notions of differential privacy [24].

Definition A.1 (Differential privacy). *An algorithm $\mathcal{M} : \mathcal{Y} \rightarrow \mathcal{O}$ is said to be (ε, δ) -differentially private for $\varepsilon, \delta > 0$ if and only if, for every $S \subseteq \mathcal{O}$ and every neighboring datasets $Y, Y' \in \mathcal{Y}$ we have*

$$\mathbb{P}[\mathcal{M}(Y) \in S] \leq e^\varepsilon \cdot \mathbb{P}[\mathcal{M}(Y') \in S] + \delta.$$

To avoid confusion, for each problem we will exactly state the relevant notion of neighboring datasets. Differential privacy is closed under post-processing and composition.

Lemma A.2 (Post-processing). *If $\mathcal{M} : \mathcal{Y} \rightarrow \mathcal{O}$ is an (ε, δ) -differentially private algorithm and $\mathcal{M}' : \mathcal{Y} \rightarrow \mathcal{Z}$ is any randomized function. Then the algorithm $\mathcal{M}'(\mathcal{M}(Y))$ is (ε, δ) -differentially private.*

In order to talk about composition it is convenient to also consider DP algorithms whose privacy guarantee holds only against subsets of inputs.

Definition A.3 (Differential Privacy Under Condition). *An algorithm $\mathcal{M} : \mathcal{Y} \rightarrow \mathcal{O}$ is said to be (ε, δ) -differentially private under condition Ψ (or (ε, δ) -DP under condition Ψ) for $\varepsilon, \delta > 0$ if and only if, for every $S \subseteq \mathcal{O}$ and every neighboring datasets $Y, Y' \in \mathcal{Y}$ both satisfying Ψ we have*

$$\mathbb{P}[\mathcal{M}(Y) \in S] \leq e^\varepsilon \cdot \mathbb{P}[\mathcal{M}(Y') \in S] + \delta.$$

It is not hard to see that the following composition theorem holds for privacy under condition.

Lemma A.4 (Composition for Algorithm with Halting, [35]). *Let $\mathcal{M}_1 : \mathcal{Y} \rightarrow \mathcal{O}_1 \cup \{\perp\}$, $\mathcal{M}_2 : \mathcal{O}_1 \times \mathcal{Y} \rightarrow \mathcal{O}_2 \cup \{\perp\}$, \dots , $\mathcal{M}_t : \mathcal{O}_{t-1} \times \mathcal{Y} \rightarrow \mathcal{O}_t \cup \{\perp\}$ be algorithms. Furthermore, let \mathcal{M} denote the algorithm that proceeds as follows (with \mathcal{O}_0 being empty): For $i = 1 \dots, t$ compute $o_i = \mathcal{M}_i(o_{i-1}, Y)$ and, if $o_i = \perp$, halt and output \perp . Finally, if the algorithm has not halted, then output o_t . Suppose that:*

- For any $1 \leq i \leq t$, we say that Y satisfies the condition Ψ_i if running the algorithm on Y does not result in halting after applying $\mathcal{M}_1, \dots, \mathcal{M}_i$.
- \mathcal{M}_1 is $(\varepsilon_1, \delta_1)$ -DP.
- \mathcal{M}_i is $(\varepsilon_i, \delta_i)$ -DP (with respect to neighboring datasets in the second argument) under condition Ψ_{i-1} for all $i = \{2, \dots, t\}$.

Then \mathcal{M} is $(\sum_i \varepsilon_i, \sum_i \delta_i)$ -DP.

A.1.1 Basic differential privacy mechanisms

The Gaussian and the Laplace mechanism are among the most widely used mechanisms in differential privacy. They work by adding a noise drawn from the Gaussian (respectively Laplace) distribution to the output of the function one wants to privatize. The magnitude of the noise depends on the sensitivity of the function.

Definition A.5 (Sensitivity of function). *Let $f : \mathcal{Y} \rightarrow \mathbb{R}^d$ be a function, its ℓ_1 -sensitivity and ℓ_2 -sensitivity are respectively*

$$\Delta_{f,1} := \max_{\substack{Y, Y' \in \mathcal{Y} \\ Y, Y' \text{ are adjacent}}} \|f(Y) - f(Y')\|_1 \quad \Delta_{f,2} := \max_{\substack{Y, Y' \in \mathcal{Y} \\ Y, Y' \text{ are adjacent}}} \|f(Y) - f(Y')\|_2.$$

For function with bounded ℓ_1 -sensitivity the Laplace mechanism is often the tool of choice to achieve privacy.

Definition A.6 (Laplace distribution). *The Laplace distribution with mean μ and parameter $b > 0$, denoted by $\text{Lap}(\mu, b)$, has PDF $\frac{1}{2b} e^{-|x-\mu|/b}$. Let $\text{Lap}(b)$ denote $\text{Lap}(0, b)$.*

A standard tail bound concerning the Laplace distribution will be useful throughout the paper.

Fact A.7 (Laplace tail bound). *Let $\mathbf{x} \sim \text{Lap}(\mu, b)$. Then,*

$$\mathbb{P}[|\mathbf{x} - \mu| > t] \leq e^{-t/b}.$$

The Laplace distribution is useful for the following mechanism

Lemma A.8 (Laplace mechanism). *Let $f : \mathcal{Y} \rightarrow \mathbb{R}^d$ be any function with ℓ_1 -sensitivity at most $\Delta_{f,1}$. Then the algorithm that adds $\text{Lap}\left(\frac{\Delta_{f,1}}{\varepsilon}\right)^{\otimes d}$ to f is $(\varepsilon, 0)$ -DP.*

It is also useful to consider the "truncated" version of the Laplace distribution where the noise distribution is shifted and truncated to be non-positive.

Definition A.9 (Truncated Laplace distribution). *The (negatively) truncated Laplace distribution with mean μ and parameter b on \mathbb{R} , denoted by $t\text{Lap}(\mu, b)$, is defined as $\text{Lap}(\mu, b)$ conditioned on the value being non-positive.*

Lemma A.10 (Truncated Laplace mechanism). *Let $f : \mathcal{Y} \rightarrow \mathbb{R}$ be any function with ℓ_1 -sensitivity at most $\Delta_{f,1}$. Then the algorithm that adds $t\text{Lap}\left(-\Delta_{f,1}\left(1 + \frac{\log(1/\delta)}{\varepsilon}\right), \Delta_{f,1}/\varepsilon\right)$ to f is (ε, δ) -DP.*

The following tail bound is useful when reasoning about truncated Laplace random variables.

Lemma A.11 (Tail bound truncated Laplace). *Suppose $\mu < 0$ and $b > 0$. Let $\mathbf{x} \sim t\text{Lap}(\mu, b)$. Then, for $y < \mu$ we have that*

$$\mathbb{P}[\mathbf{x} < y] \leq \frac{e^{(y-\mu)/b}}{2 - e^{\mu/b}}.$$

In contrast, when the function has bounded ℓ_2 -sensitivity, the Gaussian mechanism provides privacy.

Lemma A.12 (Gaussian mechanism). *Let $f : \mathcal{Y} \rightarrow \mathbb{R}^d$ be any function with ℓ_2 -sensitivity at most $\Delta_{f,2}$. Let $0 < \varepsilon, \delta \leq 1$. Then the algorithm that adds $N\left(0, \frac{\Delta_{f,2}^2 \cdot 2 \log(2/\delta)}{\varepsilon^2} \cdot \text{Id}\right)$ to f is (ε, δ) -DP.*

A.1.2 Private histograms

Here we present a classical private mechanism to learn a high dimensional histogram.

Lemma A.13 (High-dimensional private histogram learner, see [33]). *Let $q, b, \varepsilon > 0$ and $0 < \delta < 1/n$. Let $\{I_i\}_{i=-\infty}^{\infty}$ be a partition of \mathbb{R} into intervals of length b , where $I_i := \{x \in \mathbb{R} \mid q + (i-1) \cdot b \leq x < q + i \cdot b\}$. Consider the partition of \mathbb{R}^d into sets $\{B_{i_1, \dots, i_d}\}_{i_1, \dots, i_d=1}^{\infty}$ where*

$$B_{i_1, \dots, i_d} := \{x \in \mathbb{R}^d \mid \forall j \in [d], x_j \in I_{i_j}\}$$

Let $Y = \{y_1, \dots, y_n\} \subseteq \mathbb{R}^d$ be a dataset of n points. For each B_{i_1, \dots, i_d} , let $p_{i_1, \dots, i_d} = \frac{1}{n} |\{j \in [n] \mid y_j \in B_{i_1, \dots, i_d}\}|$. For $n \geq \frac{8}{\varepsilon \alpha} \cdot \log \frac{2}{\delta \beta}$, there exists an efficient (ε, δ) -differentially private algorithm that returns $\hat{p}_{1, \dots, 1}, \dots, \hat{p}_{i_1, \dots, i_d}, \dots$ satisfying

$$\mathbb{P} \left[\max_{i_1, \dots, i_d \in \mathbb{N}} |p_{i_1, \dots, i_d} - \hat{p}_{i_1, \dots, i_d}| \geq \alpha \right] \leq \beta.$$

Proof. We consider the following algorithm, applied to each $i_1, \dots, i_d \in \mathbb{N}$ on input Y :

1. If $p_{i_1, \dots, i_d} = 0$ set $\hat{p}_{i_1, \dots, i_d} = 0$, otherwise let $\hat{p}_{i_1, \dots, i_d} = p_{i_1, \dots, i_d} + \tau$ where $\tau \sim \text{Lap}(0, \frac{2}{n\varepsilon})$.
2. If $\hat{p}_{i_1, \dots, i_d} \leq \frac{3 \log(2/\delta)}{\varepsilon n}$ set $\hat{p}_{i_1, \dots, i_d} = 0$.

First we argue utility. By construction we get $\hat{p}_{i_1, \dots, i_d} = 0$ whenever $p_{i_1, \dots, i_d} = 0$, thus we may focus on non-zero p_{i_1, \dots, i_d} . There are at most n non zero p_{i_1, \dots, i_d} . By choice of n, δ and by Fact A.7 the maximum over n independent trials $\tau \sim \text{Lap}(0, \frac{2}{n\varepsilon})$ is bounded by α in absolute value with probability at least β .

It remains to argue privacy. Let $Y = \{y_1, \dots, y_n\}$, $Y' = \{y'_1, \dots, y'_n\}$ be adjacent datasets. For $i_1, \dots, i_d \in \mathbb{N}$, let

$$p_{i_1, \dots, i_d} = |\{j \in [n] \mid y_j \in B_{i_1, \dots, i_d}\}|$$

$$p'_{i_1, \dots, i_d} = |\{j \in [n] \mid y'_j \in B_{i_1, \dots, i_d}\}|.$$

Since Y, Y' are adjacent there exists only two set of indices $\mathcal{I} := \{i_1, \dots, i_d\}$ and $\mathcal{J} := \{j_1, \dots, j_d\}$ such that $p_{\mathcal{I}} \neq p'_{\mathcal{I}}$ and $p_{\mathcal{J}} \neq p'_{\mathcal{J}}$. Assume without loss of generality $p_{\mathcal{I}} > p'_{\mathcal{I}}$. Then it must be $p_{\mathcal{I}} = p'_{\mathcal{I}} + 1/n$ and $p_{\mathcal{J}} = p'_{\mathcal{J}} - 1/n$. Thus by the standard tail bound on the Laplace distribution in Fact A.7 and by Lemma A.8, we immediately get that the algorithm is (ε, δ) -differentially private. \square

A.2 Sum-of-squares and pseudo-distributions

We introduce here the sum-of-squares notion necessary for our private algorithm learning mixtures of Gaussians. We remark that these notions are not needed for Appendix C.

Let $w = (w_1, w_2, \dots, w_n)$ be a tuple of n indeterminates and let $\mathbb{R}[w]$ be the set of polynomials with real coefficients and indeterminates w, \dots, w_n . We say that a polynomial $p \in \mathbb{R}[w]$ is a *sum-of-squares (sos)* if there are polynomials q_1, \dots, q_r such that $p = q_1^2 + \dots + q_r^2$.

A.2.1 Pseudo-distributions

Pseudo-distributions are generalizations of probability distributions. We can represent a discrete (i.e., finitely supported) probability distribution over \mathbb{R}^n by its probability mass function $D: \mathbb{R}^n \rightarrow \mathbb{R}$ such that $D \geq 0$ and $\sum_{w \in \text{supp}(D)} D(w) = 1$. Similarly, we can describe a pseudo-distribution by its mass function. Here, we relax the constraint $D \geq 0$ and only require that D passes certain low-degree non-negativity tests.

Concretely, a *level- ℓ pseudo-distribution* is a finitely-supported function $D: \mathbb{R}^n \rightarrow \mathbb{R}$ such that $\sum_w D(w) = 1$ and $\sum_w D(w) f(w)^2 \geq 0$ for every polynomial f of degree at most $\ell/2$. (Here, the summations are over the support of D .) A straightforward polynomial-interpolation argument shows

that every level- ∞ -pseudo distribution satisfies $D \geq 0$ and is thus an actual probability distribution. We define the *pseudo-expectation* of a function f on \mathbb{R}^d with respect to a pseudo-distribution D , denoted $\tilde{\mathbb{E}}_{D(w)} f(w)$, as

$$\tilde{\mathbb{E}}_{D(w)} f(w) = \sum_w D(w) f(w). \quad (\text{A.1})$$

The degree- ℓ moment tensor of a pseudo-distribution D is the tensor $\mathbb{E}_{D(w)}(1, w_1, w_2, \dots, w_n)^{\otimes \ell}$. In particular, the moment tensor has an entry corresponding to the pseudo-expectation of all monomials of degree at most ℓ in w . The set of all degree- ℓ moment tensors of probability distribution is a convex set. Similarly, the set of all degree- ℓ moment tensors of degree d pseudo-distributions is also convex. Key to the algorithmic utility of pseudo-distributions is the fact that while there can be no efficient separation oracle for the convex set of all degree- ℓ moment tensors of an actual probability distribution, there's a separation oracle running in time $n^{O(\ell)}$ for the convex set of the degree- ℓ moment tensors of all level- ℓ pseudodistributions.

Fact A.14 ([56, 53, 51, 38]). *For any $n, \ell \in \mathbb{N}$, the following set has a $n^{O(\ell)}$ -time weak separation oracle (in the sense of [27]):*

$$\left\{ \tilde{\mathbb{E}}_{D(w)}(1, w_1, w_2, \dots, w_n)^{\otimes d} \mid \text{degree-}d \text{ pseudo-distribution } D \text{ over } \mathbb{R}^n \right\}. \quad (\text{A.2})$$

This fact, together with the equivalence of weak separation and optimization [27] allows us to efficiently optimize over pseudo-distributions (approximately)—this algorithm is referred to as the sum-of-squares algorithm.

The *level- ℓ sum-of-squares algorithm* optimizes over the space of all level- ℓ pseudo-distributions that satisfy a given set of polynomial constraints—we formally define this next.

Definition A.15 (Constrained pseudo-distributions). *Let D be a level- ℓ pseudo-distribution over \mathbb{R}^n . Let $\mathcal{A} = \{f_1 \geq 0, f_2 \geq 0, \dots, f_m \geq 0\}$ be a system of m polynomial inequality constraints. We say that D satisfies the system of constraints \mathcal{A} at degree r , denoted $D \models_r \mathcal{A}$, if for every $S \subseteq [m]$ and every sum-of-squares polynomial h with $\deg h + \sum_{i \in S} \max\{\deg f_i, r\} \leq \ell$,*

$$\tilde{\mathbb{E}}_D h \cdot \prod_{i \in S} f_i \geq 0.$$

We write $D \models \mathcal{A}$ (without specifying the degree) if $D \models_0 \mathcal{A}$ holds. Furthermore, we say that $D \models_r \mathcal{A}$ holds approximately if the above inequalities are satisfied up to an error of $2^{-n^\ell} \cdot \|h\| \cdot \prod_{i \in S} \|f_i\|$, where $\|\cdot\|$ denotes the Euclidean norm²² of the coefficients of a polynomial in the monomial basis.

We remark that if D is an actual (discrete) probability distribution, then we have $D \models \mathcal{A}$ if and only if D is supported on solutions to the constraints \mathcal{A} .

We say that a system \mathcal{A} of polynomial constraints is *explicitly bounded* if it contains a constraint of the form $\{\|w\|^2 \leq M\}$. The following fact is a consequence of Fact A.14 and [27],

Fact A.16 (Efficient Optimization over Pseudo-distributions). *There exists an $(n + m)^{O(\ell)}$ -time algorithm that, given any explicitly bounded and satisfiable system²³ \mathcal{A} of m polynomial constraints in n variables, outputs a level- ℓ pseudo-distribution that satisfies \mathcal{A} approximately.*

A.2.2 Sum-of-squares proof

Let f_1, f_2, \dots, f_r and g be multivariate polynomials in w . A *sum-of-squares proof* that the constraints $\{f_1 \geq 0, \dots, f_m \geq 0\}$ imply the constraint $\{g \geq 0\}$ consists of sum-of-squares polynomials $(p_S)_{S \subseteq [m]}$ such that

$$g = \sum_{S \subseteq [m]} p_S \cdot \prod_{i \in S} f_i. \quad (\text{A.3})$$

²²The choice of norm is not important here because the factor 2^{-n^ℓ} swamps the effects of choosing another norm.

²³Here, we assume that the bit complexity of the constraints in \mathcal{A} is $(n + m)^{O(1)}$.

We say that this proof has *degree* ℓ if for every set $S \subseteq [m]$, the polynomial $p_S \prod_{i \in S} f_i$ has degree at most ℓ . If there is a degree ℓ SoS proof that $\{f_i \geq 0 \mid i \leq r\}$ implies $\{g \geq 0\}$, we write:

$$\{f_i \geq 0 \mid i \leq r\} \Big|_{\ell} \{g \geq 0\}. \quad (\text{A.4})$$

Sum-of-squares proofs satisfy the following inference rules. For all polynomials $f, g: \mathbb{R}^n \rightarrow \mathbb{R}$ and for all functions $F: \mathbb{R}^n \rightarrow \mathbb{R}^m, G: \mathbb{R}^n \rightarrow \mathbb{R}^k, H: \mathbb{R}^p \rightarrow \mathbb{R}^n$ such that each of the coordinates of the outputs are polynomials of the inputs, we have:

$$\frac{\mathcal{A} \Big|_{\ell} \{f \geq 0, g \geq 0\}}{\mathcal{A} \Big|_{\ell} \{f + g \geq 0\}}, \frac{\mathcal{A} \Big|_{\ell} \{f \geq 0\}, \mathcal{A} \Big|_{\ell'} \{g \geq 0\}}{\mathcal{A} \Big|_{\ell+\ell'} \{f \cdot g \geq 0\}} \quad (\text{addition and multiplication})$$

$$\frac{\mathcal{A} \Big|_{\ell} \mathcal{B}, \mathcal{B} \Big|_{\ell'} C}{\mathcal{A} \Big|_{\ell, \ell'} C} \quad (\text{transitivity})$$

$$\frac{\{F \geq 0\} \Big|_{\ell} \{G \geq 0\}}{\{F(H) \geq 0\} \Big|_{\ell \cdot \deg(H)} \{G(H) \geq 0\}}. \quad (\text{substitution})$$

Low-degree sum-of-squares proofs are sound and complete if we take low-level pseudo-distributions as models.

Concretely, sum-of-squares proofs allow us to deduce properties of pseudo-distributions that satisfy some constraints.

Fact A.17 (Soundness). *If $D \Big|_{r'} \mathcal{A}$ for a level- ℓ pseudo-distribution D and there exists a sum-of-squares proof $\mathcal{A} \Big|_{r'} \mathcal{B}$, then $D \Big|_{r, r'+r'} \mathcal{B}$.*

If the pseudo-distribution D satisfies \mathcal{A} only approximately, soundness continues to hold if we require an upper bound on the bit-complexity of the sum-of-squares $\mathcal{A} \Big|_{r'} \mathcal{B}$ (number of bits required to write down the proof).

In our applications, the bit complexity of all sum of squares proofs will be $n^{O(\ell)}$ (assuming that all numbers in the input have bit complexity $n^{O(1)}$). This bound suffices in order to argue about pseudo-distributions that satisfy polynomial constraints approximately.

The following fact shows that every property of low-level pseudo-distributions can be derived by low-degree sum-of-squares proofs.

Fact A.18 (Completeness). *Suppose $d \geq r' \geq r$ and \mathcal{A} is a collection of polynomial constraints with degree at most r , and $\mathcal{A} \vdash \{\sum_{i=1}^n w_i^2 \leq B\}$ for some finite B .*

Let $\{g \geq 0\}$ be a polynomial constraint. If every degree- d pseudo-distribution that satisfies $D \Big|_{r'} \mathcal{A}$ also satisfies $D \Big|_{r'} \{g \geq 0\}$, then for every $\varepsilon > 0$, there is a sum-of-squares proof $\mathcal{A} \Big|_{d} \{g \geq -\varepsilon\}$.

A.2.3 Explicitly bounded distributions

We will consider a subset of subgaussian distributions denoted as certifiably subgaussians. Many subgaussians distributions are known to be certifiably subgaussian (see [36]).

Definition A.19 (Explicitly bounded distribution). *Let $t \in \mathbb{N}$. A distribution D over \mathbb{R}^d with mean μ is called $2t$ -explicitly σ -bounded if for each even integer s such that $1 \leq s \leq t$ the following equation has a degree s sum-of-squares proof in the vector variable u*

$$\Big|_{u}^{2s} \left\{ \mathbb{E}_{\mathbf{x} \sim D} \langle \mathbf{x} - \mu, u \rangle^{2s} \leq (\sigma s)^s \cdot \|u\|_2^{2s} \right\}$$

Furthermore, we say that D is explicitly bounded if it is $2t$ -explicitly σ -bounded for every $t \in \mathbb{N}$. A finite set $X \subseteq \mathbb{R}^d$ is said to be $2t$ -explicitly σ -bounded if the uniform distribution on X is $2t$ -explicitly σ -bounded.

Sets that are $2t$ -explicitly σ -bounded with large intersection satisfy certain key properties. Before introducing them we conveniently present the following definition.

Definition A.20 (Weight vector inducing distribution). Let Y be a set of size n and let $p \in [0, 1]^n$ be a vector satisfying $\|p\|_1 = 1$. We say that p induces the distribution D with support Y if

$$\mathbb{P}_{\mathbf{y} \sim D} [\mathbf{y} = y_i] = p_i.$$

Theorem A.21 ([36, 30]). Let $Y \subseteq \mathbb{R}^d$ be a set of cardinality n . Let $p, p' \in [0, 1]^n$ be weight vectors satisfying $\|p\|_1 = \|p'\|_1 = 1$ and $\|p - p'\|_1 \leq \beta$. Suppose that p (respectively p') induces a $2t$ -explicitly σ_1 -bounded (resp. σ_2) distribution over Y with mean $\mu_{(p)}$ (resp. $\mu_{(p')}$). There exists an absolute constant β^* such that, if $\beta \leq \beta^*$, then for $\sigma = \sigma_1 + \sigma_2$:

$$\|\mu_{(p)} - \mu_{(p')}\| \leq \beta^{1-1/2t} \cdot O(\sqrt{\sigma t}).$$

In the context of learning Gaussian mixtures, we will make heavy use of the statement below.

Theorem A.22 ([36, 30]). Let Y be a $2t$ -explicitly σ -bounded set of size n . Let $p \in \mathbb{R}^n$ be the weight vector inducing the uniform distribution over Y . Let $p' \in \mathbb{R}^n$ be a unit vector satisfying $\|p - p'\|_1 \leq \beta$ for some $\beta \leq \beta^*$ where β^* is a small constant. Then p' induces a $2t$ -explicitly $(\sigma + O(\beta^{1-1/2t}))$ -bounded distribution over Y .

B Stability of strongly-convex optimization

In this section, we prove ℓ_2 sensitivity bounds for the minimizers of a general class of (strongly) convex optimization problems. In particular, we show how to translate a uniform point-wise sensitivity bound for the objective functions into a ℓ_2 sensitivity bound for the minimizers.

Lemma B.1 (Stability of strongly-convex optimization). Let \mathcal{Y} be a set of datasets. Let $\mathcal{K}(\mathcal{Y})$ be a family of closed convex subsets of \mathbb{R}^m parametrized by $Y \in \mathcal{Y}$ and let $\mathcal{F}(\mathcal{Y})$ be a family of functions $f_Y : \mathcal{K}(Y) \rightarrow \mathbb{R}$, parametrized by $Y \in \mathcal{Y}$, such that:

- (i) for adjacent datasets $Y, Y' \in \mathcal{Y}$ and $X \in \mathcal{K}(Y)$ there exist $Z \in \mathcal{K}(Y) \cap \mathcal{K}(Y')$ satisfying $|f_Y(X) - f_{Y'}(Z)| \leq \alpha$ and $|f_Y(Z) - f_{Y'}(Z)| \leq \alpha$.
- (ii) f_Y is κ -strongly convex in $X \in \mathcal{K}(Y)$.

Then for $Y, Y' \in \mathcal{Y}$, $\hat{X} := \arg \min_{X \in \mathcal{K}(Y)} f_Y(X)$ and $\hat{X}' := \arg \min_{X' \in \mathcal{K}(Y')} f_{Y'}(X')$, it holds

$$\|\hat{X} - \hat{X}'\|_2^2 \leq \frac{12\alpha}{\kappa}.$$

Proof. Let $Z \in \mathcal{K}(Y) \cap \mathcal{K}(Y')$ be a point such that $|f_Y(\hat{X}) - f_{Y'}(Z)| \leq \alpha$ and $|f_Y(Z) - f_{Y'}(Z)| \leq \alpha$. By κ -strong convexity of f_Y and $f_{Y'}$ (Proposition G.2) it holds

$$\begin{aligned} \|\hat{X} - \hat{X}'\|_2^2 &\leq 2\|\hat{X} - Z\|_2^2 + 2\|Z - \hat{X}'\|_2^2 \\ &\leq \frac{4}{\kappa} \left(f_Y(Z) - f_Y(\hat{X}) + f_{Y'}(Z) - f_{Y'}(\hat{X}') \right). \end{aligned}$$

Suppose w.l.o.g. $f_Y(\hat{X}) \leq f_{Y'}(\hat{X}')$, for a symmetric argument works in the other case. Then

$$f_Y(Z) \leq f_{Y'}(Z) + \alpha \leq f_Y(\hat{X}) + 2\alpha$$

and

$$f_Y(\hat{X}) \leq f_{Y'}(\hat{X}') \leq f_{Y'}(Z) \leq f_Y(\hat{X}) + \alpha.$$

It follows as desired

$$f_Y(Z) - f_Y(\hat{X}) + f_{Y'}(Z) - f_{Y'}(\hat{X}') \leq 3\alpha.$$

□

C Private recovery for stochastic block models

In this section, we present how to achieve exact recovery in stochastic block models privately and thus prove Theorem 1.3. To this end, we first use the stability of strongly convex optimization (Lemma B.1) to obtain a private weak recovery algorithm in Appendix C.1. Then we show how to privately boost the weak recovery algorithm to achieve exact recovery in Appendix C.2. In Appendix C.4, we complement our algorithmic results by providing an almost tight lower bound on the privacy parameters. We start by defining the relevant notion of adjacent datasets.

Definition C.1 (Adjacent graphs). *Let G, G' be graphs with vertex set $[n]$. We say that G, G' are adjacent if $|E(G) \Delta E(G')| = 1$.*

Remark C.2 (Parameters as public information). *We remark that we assume the parameters n, γ, d to be public information given in input to the algorithm.*

C.1 Private weak recovery for stochastic block models

In this section, we show how to achieve weak recovery privately via stability of strongly convex optimization (Lemma B.1). We first introduce one convenient notation. The error rate of an estimate $\hat{x} \in \{\pm 1\}^n$ of the true partition $x \in \{\pm 1\}^n$ is defined as $\text{err}(\hat{x}, x) := \frac{1}{n} \cdot \min\{\text{Ham}(\hat{x}, x), \text{Ham}(\hat{x}, -x)\}$.²⁴ Our main result is the following theorem.

Theorem C.3. *Suppose $\gamma\sqrt{d} \geq 12800, \varepsilon, \delta \geq 0$. There exists an (Algorithm C.4) such that, for any $x \in \{\pm 1\}^n$, on input $\mathbf{G} \sim \text{SBM}_n(\gamma, d, x)$, outputs $\hat{x}(\mathbf{G}) \in \{\pm 1\}^n$ satisfying*

$$\text{err}(\hat{x}(\mathbf{G}), x) \leq O\left(\frac{1}{\gamma\sqrt{d}} + \frac{1}{\gamma d} \cdot \frac{\log(2/\delta)}{\varepsilon^2}\right)$$

with probability $1 - \exp(-\Omega(n))$. Moreover, the algorithm is (ε, δ) -differentially private for any input graph and runs in polynomial time.

Before presenting the algorithm we introduce some notation. Given a graph G , let $Y(G) := \frac{1}{\gamma d}(A(G) - \frac{d}{n}J)$ where $A(G)$ is the adjacency matrix of G and J denotes all-one matrices. Define $\mathcal{K} := \{X \in \mathbb{R}^{n \times n} \mid X \succeq 0, X_{ii} = \frac{1}{n} \forall i\}$. The algorithm starts with projecting matrix $Y(G)$ to set \mathcal{K} . To ensure privacy, then it adds Gaussian noise to the projection X_1 and obtains a private matrix X_2 . The last step applies a standard rounding method.

Algorithm C.4 (Private weak recovery for SBM).

Input: Graph G .

Operations:

1. *Projection:* $X_1 \leftarrow \text{argmin}_{X \in \mathcal{K}} \|Y(G) - X\|_F^2$.
2. *Noise addition:* $\mathbf{X}_2 \leftarrow X_1 + \mathbf{W}$ where $\mathbf{W} \sim \mathcal{N}\left(0, \frac{24}{n\gamma d} \frac{\log(2/\delta)}{\varepsilon^2}\right)^{n \times n}$.
3. *Rounding:* Compute the leading eigenvector \mathbf{v} of \mathbf{X}_2 and return $\text{sign}(\mathbf{v})$.

In the rest of this section, we will show Algorithm C.4 is private in Lemma C.7 and its utility guarantee in Lemma C.8. Then Theorem C.3 follows directly from Lemma C.7 and Lemma C.8.

Privacy analysis Let \mathcal{Y} be the set of all matrices $Y(G) = \frac{1}{\gamma d}(A(G) - \frac{d}{n}J)$ where G is a graph on n vertices. We further define $q: \mathcal{Y} \rightarrow \mathcal{K}$ to be the function

$$q(Y) := \text{argmin}_{X \in \mathcal{K}} \|Y - X\|_F^2. \tag{C.1}$$

We first use Lemma B.1 to prove that function q is stable.

Lemma C.5 (Stability). *The function q as defined in Eq. (C.1) has ℓ_2 -sensitivity $\Delta_{q,2} \leq \sqrt{\frac{24}{n\gamma d}}$.*

²⁴Note $|\langle \hat{x}, x \rangle| = (1 - 2 \text{err}(\hat{x}, x)) \cdot n$ for any $\hat{x}, x \in \{\pm 1\}^n$.

Proof. Let $g : \mathcal{Y} \times \mathcal{K} \rightarrow \mathbb{R}$ be the function $g(Y, X) := \|X\|_F^2 - 2\langle Y, X \rangle$. Applying Lemma B.1 with $f_Y(\cdot) = g(Y, \cdot)$, it suffices to prove that g has ℓ_1 -sensitivity $\frac{4}{n\gamma d}$ with respect to Y and that it is 2-strongly convex with respect to X . The ℓ_1 -sensitivity bound follows by observing that adjacent Y, Y' satisfy $\|Y - Y'\|_1 \leq \frac{2}{\gamma d}$ and that any $X \in \mathcal{K}$ satisfies $\|X\|_\infty \leq \frac{1}{n}$. Thus it remains to prove strong convexity with respect to $X \in \mathcal{K}$. Let $X, X' \in \mathcal{K}$ then

$$\begin{aligned} \|X'\|_F^2 &= \|X\|_F^2 + 2\langle X' - X, X \rangle + \|X - X'\|_F^2 \\ &= \|X\|_F^2 + 2\langle X' - X, X + Y - Y \rangle + \|X - X'\|_F^2 \\ &= g(Y, X) + \langle X' - X, \nabla g(X, Y) \rangle + 2\langle X', Y \rangle + \|X - X'\|_F^2. \end{aligned}$$

That is $g(Y, X)$ is 2-strongly convex with respect to X . Note any $X \in \mathcal{K}$ is symmetric. Then the result follows by Lemma B.1. \square

Remark C.6. In the special case where the constraint set \mathcal{K} does not depend on input dataset Y (e.g. stochastic block models), the proof can be cleaner as follows. Let $f_Y(X) := \|X\|_F^2 - 2\langle X, Y \rangle$. Let $\hat{X} := \operatorname{argmin}_{X \in \mathcal{K}} f_Y(X)$ and $\hat{X}' := \operatorname{argmin}_{X \in \mathcal{K}} f_{Y'}(X)$. Suppose without loss of generality $f_Y(\hat{X}) \leq f_{Y'}(\hat{X}')$, for a symmetric argument works in the other case. Then

$$f_Y(\hat{X}) \leq f_{Y'}(\hat{X}') \leq f_{Y'}(\hat{X}) \leq f_Y(\hat{X}) + \alpha.$$

Then it is easy to show the algorithm is private.

Lemma C.7 (Privacy). The weak recovery algorithm (Algorithm C.4) is (ε, δ) -DP.

Proof. Since any $X \in \mathcal{K}$ is symmetric, we only need to add a symmetric noise matrix to obtain privacy. Combining Lemma C.5 with Lemma A.12, we immediately get that the algorithm is (ε, δ) -private. \square

Utility analysis Now we show the utility guarantee of our private weak recovery algorithm.

Lemma C.8 (Utility). For any $x \in \{\pm 1\}^n$, on input $\mathbf{G} \sim \text{SBM}_n(\gamma, d, x)$, Algorithm C.4 efficiently outputs $\hat{x}(\mathbf{G}) \in \{\pm 1\}^n$ satisfying

$$\operatorname{err}(\hat{x}(\mathbf{G}), x) \leq \frac{6400}{\gamma\sqrt{d}} + \frac{7000}{\gamma d} \cdot \frac{\log(2/\delta)}{\varepsilon^2},$$

with probability $1 - \exp(-\Omega(n))$.

To prove Lemma C.8, we need the following lemma which is an adaption of a well-known result in SBM [28, Theorem 1.1]. Its proof is deferred to Appendix H.

Lemma C.9. Consider the settings of Lemma C.8. With probability $1 - \exp(-\Omega(n))$,

$$\left\| X_1(\mathbf{G}) - \frac{1}{n}xx^\top \right\|_F^2 \leq \frac{800}{\gamma\sqrt{d}}.$$

Proof of Lemma C.8. By Lemma C.9, we have

$$\left\| X_1(\mathbf{G}) - \frac{1}{n}xx^\top \right\| \leq \left\| X_1(\mathbf{G}) - \frac{1}{n}xx^\top \right\|_F \leq \sqrt{\frac{800}{\gamma\sqrt{d}}} =: r(\gamma, d)$$

with probability $1 - \exp(-\Omega(n))$. We condition our following analysis on this event happening.

Let \mathbf{u} be the leading eigenvector of $X_1(\mathbf{G})$. Let λ_1 and λ_2 be the largest and second largest eigenvalues of $X_1(\mathbf{G})$. By Weyl's inequality (Lemma F.1) and the assumption $\gamma\sqrt{d} \geq 12800$, we have

$$\lambda_1 - \lambda_2 \geq 1 - 2r(\gamma, d) \geq \frac{1}{2}.$$

Let \mathbf{v} be the leading eigenvector of $X_1(\mathbf{G}) + \mathbf{W}$. By Davis-Kahan's theorem (Lemma F.2), we have

$$\|\mathbf{u} - \mathbf{v}\| \leq \frac{2\|\mathbf{W}\|}{\lambda_1 - \lambda_2} \leq 4\|\mathbf{W}\|,$$

$$\|\mathbf{u} - x/\sqrt{n}\| \leq 2 \left\| X_1(\mathbf{G}) - \frac{1}{n}xx^\top \right\| \leq 2r(\gamma, d).$$

Putting things together and using Fact E.1, we have

$$\|\mathbf{v} - x/\sqrt{n}\| \leq \|\mathbf{u} - \mathbf{v}\| + \|\mathbf{u} - x/\sqrt{n}\| \leq \frac{24\sqrt{6}}{\sqrt{\gamma d}} \frac{\sqrt{\log(2/\delta)}}{\varepsilon} + 2r(\gamma, d)$$

with probability $1 - \exp(-\Omega(n))$.

Observe $\text{Ham}(\text{sign}(y), x) \leq \|y - x\|^2$ for any $y \in \mathbb{R}^n$ and any $x \in \{\pm 1\}^n$. Then with probability $1 - \exp(-\Omega(n))$,

$$\frac{1}{n} \cdot \text{Ham}(\text{sign}(\mathbf{v}), x) \leq \|\mathbf{v} - x/\sqrt{n}\|^2 \leq \frac{6400}{\gamma\sqrt{d}} + \frac{7000}{\gamma d} \cdot \frac{\log(2/\delta)}{\varepsilon^2}.$$

□

Proof of Theorem C.3. By Lemma C.7 and Lemma C.8. □

C.2 Private exact recovery for stochastic block models

In this section, we prove Theorem 1.3. We show how to achieve exact recovery in stochastic block models privately by combining the private weak recovery algorithm we obtained in the previous section and a private majority voting scheme.

Since exact recovery is only possible with logarithmic average degree (just to avoid isolated vertices), it is more convenient to work with the following standard parameterization of stochastic block models. Let $\alpha > \beta > 0$ be fixed constants. The intra-community edge probability is $\alpha \cdot \frac{\log n}{n}$, and the inter-community edge probability is $\beta \cdot \frac{\log n}{n}$. In the language of Model 1.1, it is $\text{SBM}_n(\frac{\alpha+\beta}{2} \cdot \log n, \frac{\alpha-\beta}{\alpha+\beta}, x)$. Our main result is the following theorem.

Theorem C.10 (Private exact recovery of SBM, restatement of Theorem 1.3). *Let $\varepsilon, \delta \geq 0$. Suppose α, β are fixed constants satisfying²⁵*

$$\sqrt{\alpha} - \sqrt{\beta} \geq 16 \quad \text{and} \quad \alpha - \beta \geq \Omega\left(\frac{1}{\varepsilon^2} \cdot \frac{\log(2/\delta)}{\log n} + \frac{1}{\varepsilon}\right), \quad (\text{C.2})$$

Then there exists an algorithm (Algorithm C.12) such that, for any balanced²⁶ $x \in \{\pm 1\}^n$, on input $\mathbf{G} \sim \text{SBM}_n(\frac{\alpha+\beta}{2} \cdot \log n, \frac{\alpha-\beta}{\alpha+\beta}, x)$, outputs $\hat{x}(\mathbf{G}) \in \{x, -x\}$ with probability $1 - o(1)$. Moreover, the algorithm is (ε, δ) -differentially private for any input graph and runs in polynomial time.

Remark C.11. *In a standard regime of privacy parameters where $\varepsilon \leq O(1)$ and $\delta = 1/\text{poly}(n)$, the private exact recovery threshold Eq. (C.2) reads*

$$\sqrt{\alpha} - \sqrt{\beta} \geq 16 \quad \text{and} \quad \alpha - \beta \geq \Omega(\varepsilon^{-2} + \varepsilon^{-1}),$$

Recall the non-private exact recovery threshold is $\sqrt{\alpha} - \sqrt{\beta} > \sqrt{2}$. Thus the non-private part in Eq. (C.2), i.e. 16, is close to optimal.

Algorithm C.12 starts with randomly splitting the input graph G into two subgraphs \mathbf{G}_1 and \mathbf{G}_2 . Setting the graph-splitting probability to $1/2$, each subgraph will contain about half of the edges of G . Then we run an (ε, δ) -DP weak recovery algorithm (Algorithm C.4) on \mathbf{G}_1 to get a rough estimate $\tilde{x}(\mathbf{G}_1)$ of accuracy around 90%. Finally, we boost the accuracy to 100% by doing majority voting (Algorithm C.13) on \mathbf{G}_2 based on the rough estimate $\tilde{x}(\mathbf{G}_1)$. That is, if a vertex has more neighbors from the opposite community (according to $\tilde{x}(\mathbf{G}_1)$) in \mathbf{G}_2 , then we assign this vertex to the opposite community. To make the majority voting step private, we add some noise to the vote.

²⁵In the language of Model 1.1, for any t we have $\sqrt{\alpha} - \sqrt{\beta} \geq t$ if and only if $\frac{d}{\log n}(1 - \sqrt{1 - \gamma^2}) \geq \frac{t^2}{2}$.

²⁶Recall a vector $x \in \{\pm 1\}^n$ is said to be balanced if $\sum_{i=1}^n x_i = 0$.

Algorithm C.12 (Private exact recovery for SBM).

Input: Graph G

Operations:

1. *Graph-splitting:* Initialize \mathbf{G}_1 to be an empty graph on vertex set $V(G)$. Independently put each edge of G in \mathbf{G}_1 with probability $1/2$. Let $\mathbf{G}_2 = G \setminus \mathbf{G}_1$.
2. *Rough estimation on \mathbf{G}_1 :* Run the (ε, δ) -DP partial recovery algorithm (Algorithm C.4) on \mathbf{G}_1 to get a rough estimate $\tilde{x}(\mathbf{G}_1)$.
3. *Majority voting on \mathbf{G}_2 :* Run the $(\varepsilon, 0)$ -DP majority voting algorithm (Algorithm C.13) with input $(\mathbf{G}_2, \tilde{x}(\mathbf{G}_1))$ and get output $\hat{\mathbf{x}}$.
4. Return $\hat{\mathbf{x}}$.

Algorithm C.13 (Private majority voting).

Input: Graph G , rough estimate $\tilde{x} \in \{\pm 1\}^n$

Operations:

1. For each vertex $v \in V(G)$, let $\mathbf{Z}_v = \mathbf{S}_v - \mathbf{D}_v$ where

$$\begin{aligned} \bullet \mathbf{D}_v &= \sum_{\{u,v\} \in E(G)} \mathbb{1}_{[\tilde{x}_u \neq \tilde{x}_v]}, \\ \bullet \mathbf{S}_v &= \sum_{\{u,v\} \in E(G)} \mathbb{1}_{[\tilde{x}_u = \tilde{x}_v]}. \end{aligned}$$

Set $\hat{\mathbf{x}}_v = \text{sign}(\mathbf{Z}_v + \mathbf{W}_v) \cdot \tilde{x}(\mathbf{G}_1)_v$ where $\mathbf{W}_v \sim \text{Lap}(2/\varepsilon)$.

2. Return $\hat{\mathbf{x}}$.

In the rest of this section, we will show Algorithm C.12 is private in Lemma C.15 and it recovers the hidden communities exactly with high probability in Lemma C.17. Then Theorem C.10 follows directly from Lemma C.15 and Lemma C.17.

Privacy analysis. We first show the differential privacy of the majority voting algorithm (Algorithm C.13) with respect to input graph G (i.e. assuming fixed the input rough estimate).

Lemma C.14. *Algorithm C.13 is $(\varepsilon, 0)$ -DP with respect to input G .*

Proof. Observing the ℓ_1 -sensitivity of the degree count function Z in step 2, the $(\varepsilon, 0)$ -DP follows directly from Laplace mechanism (Lemma A.12) and post-processing (Lemma A.2). \square

Then the privacy of the private exact recovery algorithm (Algorithm C.12) is a consequence of composition.

Lemma C.15 (Privacy). *Algorithm C.12 is (ε, δ) -DP.*

Proof. Let $\mathcal{A}_1 : \mathcal{G}_n \rightarrow \{\pm 1\}^n$ denote the (ε, δ) -DP recovery algorithm in step 2. Let $\mathcal{A}_2 : \mathcal{G}_n \times \{\pm 1\}^n \rightarrow \{\pm 1\}^n$ denote the $(\varepsilon, 0)$ -DP majority voting algorithm in step 3. Let \mathcal{A} be the composition of \mathcal{A}_1 and \mathcal{A}_2 .

We first make several notations. Given a graph H and an edge e , H_e is a graph obtained by adding e to H . Given a graph H , $\mathbf{G}_1(H)$ is a random subgraph of H by keeping each edge of H with probability $1/2$ independently.

Now, fix two adjacent graphs G and G_e where edge e appears in G_e but not in G . Also, fix two arbitrary possible outputs $x_1, x_2 \in \{\pm 1\}^n$ of algorithm \mathcal{A} .²⁷ It is direct to see,

$$\mathbb{P}(\mathcal{A}(G) = (x_1, x_2)) = \sum_{H \subseteq G} \mathbb{P}(\mathcal{A}_1(H) = x_1) \mathbb{P}(\mathcal{A}_2(G \setminus H, x_1) = x_2) \mathbb{P}(\mathbf{G}_1(G) = H). \quad (\text{C.3})$$

²⁷We can imagine that algorithm \mathcal{A} first outputs (x_1, x_2) and then outputs x_2 as a post-processing step.

Since $\mathbb{P}(\mathbf{G}_1(G) = H) = \mathbb{P}(\mathbf{G}_1(G_e) = H) + \mathbb{P}(\mathbf{G}_1(G_e) = H_e)$ for any $H \subseteq G$, we have

$$\begin{aligned} \mathbb{P}(\mathcal{A}(G_e) = (x_1, x_2)) &= \sum_{H \subseteq G} \mathbb{P}(\mathcal{A}_1(H) = x_1) \mathbb{P}(\mathcal{A}_2(G_e \setminus H, x_1) = x_2) \mathbb{P}(\mathbf{G}_1(G_e) = H) \\ &\quad + \mathbb{P}(\mathcal{A}_1(H_e) = x_1) \mathbb{P}(\mathcal{A}_2(G_e \setminus H_e, x_1) = x_2) \mathbb{P}(\mathbf{G}_1(G_e) = H_e) \end{aligned} \tag{C.4}$$

Since both \mathcal{A}_1 and \mathcal{A}_2 are (ε, δ) -DP, we have for each $H \subseteq G$,

$$\mathbb{P}(\mathcal{A}_1(H_e) = x_1) \leq e^\varepsilon \mathbb{P}(\mathcal{A}_1(H) = x_1) + \delta, \tag{C.5}$$

$$\mathbb{P}(\mathcal{A}_2(G_e \setminus H, x_1) = x_2) \leq e^\varepsilon \mathbb{P}(\mathcal{A}_2(G \setminus H, x_1) = x_2) + \delta. \tag{C.6}$$

Plugging Eq. (C.5) and Eq. (C.6) into Eq. (C.4), we obtain

$$\begin{aligned} \mathbb{P}(\mathcal{A}(G_e) = (x_1, x_2)) &\leq \sum_{H \subseteq G} [e^\varepsilon \mathbb{P}(\mathcal{A}_1(H) = x_1) \mathbb{P}(\mathcal{A}_2(G \setminus H, x_1) = x_2) + \delta] \mathbb{P}(\mathbf{G}_1(G) = H) \\ &= e^\varepsilon \mathbb{P}(\mathcal{A}(G) = (x_1, x_2)) + \delta. \end{aligned}$$

Similarly, we can show

$$\mathbb{P}(\mathcal{A}(G) = (x_1, x_2)) \leq e^\varepsilon \mathbb{P}(\mathcal{A}(G_e) = (x_1, x_2)) + \delta. \tag{C.7}$$

□

Utility analysis. We first show the utility guarantee of the private majority voting algorithm.

Lemma C.16. *Suppose \mathbf{G} is generated by first sampling $\mathbf{G} \sim \text{SBM}_n(\frac{\alpha+\beta}{2} \cdot \log n, \frac{\alpha-\beta}{\alpha+\beta}, x)$ for some balanced x and then for each vertex removing at most $\Delta \leq O(\log^2 n)$ adjacent edges arbitrarily. Then on input \mathbf{G} and a balanced rough estimate \tilde{x} satisfying $\text{Ham}(\tilde{x}, x) \leq n/16$, Algorithm C.13 efficiently outputs $\hat{x}(\mathbf{G})$ such that for each vertex v ,*

$$\mathbb{P}(\hat{x}(\mathbf{G})_v \neq x_v) \leq \exp\left(-\frac{1}{64} \cdot \varepsilon(\alpha - \beta) \cdot \log n\right) + 2 \cdot \exp\left(-\frac{1}{16^2} \cdot \frac{(\alpha - \beta)^2}{\alpha + \beta} \cdot \log n\right).$$

Proof. Let us fix an arbitrary vertex v and analyze the probability $\mathbb{P}(\hat{x}(\mathbf{G})_v \neq x_v)$. Let $r := \text{Ham}(\tilde{x}, x)/n$. Then it is not hard to see

$$\mathbb{P}(\hat{x}(\mathbf{G})_v \neq x_v) \leq \mathbb{P}(\mathbf{B} + \mathbf{A}' - \mathbf{A} - \mathbf{B}' + \mathbf{W} > 0) \tag{C.8}$$

where

- $\mathbf{A} \sim \text{Binomial}((1/2 - r)n - \Delta, \alpha \frac{\log n}{n})$, corresponding to the number of neighbors that are from the same community and correctly labeled by \tilde{x} ,
- $\mathbf{B}' \sim \text{Binomial}(rn - \Delta, \beta \frac{\log n}{n})$, corresponding to the number of neighbors that are from the different community but incorrectly labeled by \tilde{x} ,
- $\mathbf{B} \sim \text{Binomial}((1/2 - r)n, \beta \frac{\log n}{n})$, corresponding to the number of neighbors that are from the different community and correctly labeled by \tilde{x} ,
- $\mathbf{A}' \sim \text{Binomial}(rn, \alpha \frac{\log n}{n})$, corresponding to the number of neighbors that are from the same community but incorrectly labeled by \tilde{x} ,
- $\mathbf{W} \sim \text{Lap}(0, 2/\varepsilon)$, independently.

The Δ term appearing in both \mathbf{A} and \mathbf{B}' corresponds to the worst case where Δ “favorable” edges are removed. If $r \geq \Omega(1)$, then $\Delta = O(\log^2 n)$ is negligible to $rn = \Theta(n)$ and we can safely ignore the effect of removing Δ edges. If $r = o(1)$, then we can safely assume \tilde{x} is correct on all vertices and ignore the effect of removing Δ edges as well. Thus, we will assume $\Delta = 0$ in the following analysis.

For any t, t' , we have

$$\begin{aligned} \mathbb{P}(\mathbf{A}' + \mathbf{B} - \mathbf{A} - \mathbf{B}' + \mathbf{W} > 0) &\leq \mathbb{P}(\mathbf{A}' + \mathbf{B} + \mathbf{W} > t) + \mathbb{P}(\mathbf{A} + \mathbf{B}' \leq t) \\ &\leq \mathbb{P}(\mathbf{A}' + \mathbf{B} \geq t - t') + \mathbb{P}(\mathbf{W} \geq t') + \mathbb{P}(\mathbf{A} + \mathbf{B}' \leq t). \end{aligned}$$

We choose t, t' by first picking two constants $a, b > 0$ satisfying $a + b < 1$ and then solving

- $\mathbb{E}[\mathbf{A}' + \mathbf{B}] - t = a \cdot (\mathbb{E}[\mathbf{A} + \mathbf{B}'] - \mathbb{E}[\mathbf{A}' + \mathbf{B}])$ and
- $t' = (1 - a - b) \cdot (\mathbb{E}[\mathbf{A} + \mathbf{B}'] - \mathbb{E}[\mathbf{A}' + \mathbf{B}])$.

By Fact A.7,

$$\mathbb{P}(\mathbf{W} > t') \leq \exp\left(-\frac{t'\varepsilon}{2}\right) \leq \exp\left(-\frac{(1/4-r)(1-a-b)}{2} \cdot \varepsilon(\alpha-\beta) \cdot \log n\right).$$

By Fact E.4 and the assumption $r \leq 1/16$, we have

$$\mathbb{P}(\mathbf{A} + \mathbf{B}' \leq t) \leq \exp\left(-\frac{(\mathbb{E}[\mathbf{A} + \mathbf{B}'] - t)^2}{2\mathbb{E}[\mathbf{A} + \mathbf{B}']}\right) \leq \exp\left(-\frac{(1/4-r)^2 a^2 \cdot (\alpha-\beta)^2}{\alpha+\beta} \cdot \log n\right).$$

Setting $b = 1/2$, by Fact E.4 and the assumption $r \leq 1/16$, we have

$$\mathbb{P}(\mathbf{A}' + \mathbf{B} \geq t - t') \leq \exp\left(-\frac{(t - t' - \mathbb{E}[\mathbf{A}' + \mathbf{B}])^2}{t - t' + \mathbb{E}[\mathbf{A}' + \mathbf{B}]}\right) \leq \exp\left(-\frac{2(1/4-r)^2}{7} \cdot \frac{(\alpha-\beta)^2}{\alpha+\beta} \cdot \log n\right).$$

Further setting $a = 1/3$, we have

$$\mathbb{P}(\hat{x}(\mathbf{G})_v \neq x_v) \leq \exp\left(-\frac{1/4-r}{12} \cdot \varepsilon(\alpha-\beta) \cdot \log n\right) + 2 \cdot \exp\left(-\frac{(1/4-r)^2}{9} \cdot \frac{(\alpha-\beta)^2}{\alpha+\beta} \cdot \log n\right).$$

Finally, plugging the assumption $r \leq 1/16$ to conclude. \square

Then it is not difficult to show the utility guarantee of our private exact recovery algorithm.

Lemma C.17 (Utility). *Suppose α, β are fixed constants satisfying*

$$\sqrt{\alpha} - \sqrt{\beta} \geq 16 \quad \text{and} \quad \alpha - \beta \geq \Omega\left(\frac{1}{\varepsilon^2} \cdot \frac{\log(2/\delta)}{\log n} + \frac{1}{\varepsilon}\right).$$

Then for any balanced $x \in \{\pm 1\}^n$, on input $\mathbf{G} \sim \text{SBM}_n(\frac{\alpha+\beta}{2} \cdot \log n, \frac{\alpha-\beta}{\alpha+\beta}, x)$, Algorithm C.12 efficiently outputs $\hat{x}(\mathbf{G})$ satisfying $\hat{x}(\mathbf{G}) \in \{x, -x\}$ with probability $1 - o(1)$.

Proof. We will show the probability of a fixed vertex being misclassified is at most $o(1/n)$. Then by union bound, exact recovery can be achieved with probability $1 - o(1)$.

As the graph-splitting probability is $1/2$, \mathbf{G}_1 follows $\text{SBM}_n(\frac{\alpha}{2} \cdot \frac{\log n}{n}, \frac{\beta}{2} \cdot \frac{\log n}{n}, x)$. By Theorem C.3, the rough estimate $\tilde{x}(\mathbf{G}_1)$ satisfies²⁸

$$\text{err}(\tilde{x}(\mathbf{G}_1), x) \leq r := o(1) + \frac{14000}{(\alpha-\beta)\varepsilon^2} \cdot \frac{\log(2/\delta)}{\log n}. \quad (\text{C.9})$$

with probability at least $1 - \exp(-\Omega(n))$. Without loss of generality, we can assume $\text{Ham}(\tilde{x}(\mathbf{G}_1), x) \leq rn$, since we consider $-x$ otherwise. By Fact E.2, the maximum degree of \mathbf{G}_1 is at most $\Delta := 2 \log^2 n$ with probability at least $1 - n \exp(-(\log n)^2/3)$. In the following, we condition our analysis on the above two events regarding $\tilde{x}(\mathbf{G}_1)$ and \mathbf{G}_1 .

Now, let us fix a vertex and analyze the probability p_e that it is misclassified after majority voting. With G_1 being fixed, \mathbf{G}_2 can be thought of as being generated by first sampling \mathbf{G} and then removing G_1 from \mathbf{G} . To make $r \leq 1/16$, it suffices to ensure $\alpha - \beta > \frac{500^2}{\varepsilon^2} \cdot \frac{\log(2/\delta)}{\log n}$ by Eq. (C.9). Then by Lemma C.16, we have

$$p_e \leq \exp\left(-\frac{1}{64} \cdot \varepsilon(\alpha-\beta) \cdot \log n\right) + 2 \cdot \exp\left(-\frac{1}{16^2} \cdot \frac{(\alpha-\beta)^2}{\alpha+\beta} \cdot \log n\right).$$

To make p_e at most $o(1/n)$, it suffices to ensure

$$\frac{1}{64} \cdot \varepsilon(\alpha-\beta) > 1 \quad \text{and} \quad \frac{1}{16^2} \cdot \frac{(\alpha-\beta)^2}{\alpha+\beta} > 1.$$

²⁸It is easy to make the output of Algorithm C.4 balanced at the cost of increasing the error rate by a factor of at most 2.

Note $(\alpha - \beta)^2/(\alpha + \beta) > (\sqrt{\alpha} - \sqrt{\beta})^2$ for $\alpha > \beta$. Therefore, as long as

$$\sqrt{\alpha} - \sqrt{\beta} \geq 16 \quad \text{and} \quad \alpha - \beta \geq \frac{500^2}{\varepsilon^2} \cdot \frac{\log(2/\delta)}{\log n} + \frac{64}{\varepsilon},$$

Algorithm C.12 recovers the hidden communities exactly with probability $1 - o(1)$. □

Proof of Theorem C.10. By Lemma C.15 and Lemma C.17. □

C.3 Inefficient recovery using the exponential mechanism

In this section, we will present an inefficient algorithm satisfying pure privacy which succeeds for all ranges of parameters - ranging from weak to exact recovery. The algorithm is based on the exponential mechanism [43] combined with the majority voting scheme introduced in section Appendix C.2. In particular, we will show

Theorem C.18 (Full version of Theorem 1.4). *Let $\gamma\sqrt{d} \geq 12800$ and $x \in \{\pm 1\}^n$ be balanced. Let $\zeta \geq 2 \exp\left(-\frac{\gamma^2 d}{512}\right)$. For any $\varepsilon \geq \frac{64 \log(2/\zeta)}{\gamma d}$, there exists an algorithm, Algorithm C.19, which on input $\mathbf{G} \sim \text{SBM}_n(\gamma, d, x^*)$ outputs an estimate $\hat{x}(\mathbf{G}) \in \{\pm 1\}^n$ satisfying*

$$\text{err}(\hat{x}(\mathbf{G}), x^*) \leq \zeta$$

with probability at least $1 - \zeta$. In addition, the algorithm is ε -private. Further, by slightly modifying the algorithm, we can achieve error $20/\sqrt{\log(1/\zeta)}$ with probability $1 - e^{-n}$.²⁹

A couple of remarks are in order. First, our algorithm works across all degree-regimes in the literature and matches known non-private thresholds and rates up to constants. We remark that for ease of exposition we did not try to optimize these constants. In particular, for $\gamma^2 d$ a constant we achieve weak recovery. We reiterate, that $\gamma^2 d > 1$ is the optimal non-private threshold. For the regime, where $\gamma^2 d = \omega(1)$, it is known that the optimal error rate is $\exp(-(1 - o(1))\gamma^2 d)$ even non-privately [64], where $o(1)$ goes to zero as $\gamma^2 d$ tends to infinity. We match this up to constants. Moreover, our algorithm achieves exact recovery as soon as $\gamma^2 d \geq 512 \log n$ since then $\zeta < \frac{1}{n}$. This also matches known non-private thresholds up to constants [5, 47]. Also, our dependence on the privacy parameter ε is also optimal as shown by the information-theoretic lower bounds in Appendix C.4.

We also emphasize, that if we only aim to achieve error on the order of

$$\frac{1}{\gamma\sqrt{d}} = \Theta\left(\frac{1}{\sqrt{\log(1/\zeta)}}\right),$$

we can achieve exponentially small failure probability in n , while keeping the privacy parameter ε the same. This can be achieved, by omitting the boosting step in our algorithm and will be clear from the proof of Theorem C.18. We remark that in this case, we can also handle non-balanced communities.

Again, for an input graph G , consider the matrix $Y(G) = \frac{1}{\gamma d} (A(G) - \frac{d}{n} J)$. For $x \in \{\pm 1\}^n$ we define the score function

$$s_G(x) = \langle x, Y(G)x \rangle.$$

Since the entries of $A(G)$ are in $[0, 1]$ and adjacent graphs differ in at most one edge, it follows immediately, that this score function has sensitivity at most

$$\Delta = \max_{\substack{G \sim G', \\ x \in \{\pm 1\}^n}} |s_G(x) - s_{G'}(x)| = \frac{2}{\gamma d} \cdot \max_{\substack{G \sim G', \\ x \in \{\pm 1\}^n}} |\langle x, (A(G) - A(G'))x \rangle| \leq \frac{2}{\gamma d}.$$

²⁹The first, smaller, error guarantee additionally needs the requirement that $\zeta \leq \exp(-640)$. The second one does not.

Algorithm C.19 (Inefficient algorithm for SBM).**Input:** Graph G , privacy parameter $\varepsilon > 0$ **Operations:**

1. *Graph-splitting:* Initialize \mathbf{G}_1 to be an empty graph on vertex set $V(G)$. Independently assign each edge of G to \mathbf{G}_1 with probability $1/2$. Let $\mathbf{G}_2 = G \setminus \mathbf{G}_1$.

2. *Rough estimation on \mathbf{G}_1 :* Sample \tilde{x} from the distribution with density

$$p(x) \propto \exp\left(\frac{\varepsilon}{2\Delta} \langle x, Y(\mathbf{G}_1)x \rangle\right),$$

where $\Delta = \frac{2}{\gamma d}$.

3. *Majority voting on \mathbf{G}_2 :* Run the ε -DP majority voting algorithm (Algorithm C.13) with input $(\mathbf{G}_2, \tilde{x}(\mathbf{G}_1))$. Denote its output by \hat{x} .

4. *Return \hat{x} .*

We first analyze the privacy guarantees of the above algorithm.

Lemma C.20. *Algorithm C.19 is ε -DP.*

Proof. For simplicity and clarity of notation, we will show that the algorithm satisfies 2ε -DP. Clearly, the graph splitting step is 0-DP. Step 2 corresponds to the exponential mechanism. Since the sensitivity of the score function is at most $\Delta = \frac{2}{\gamma d}$ it follows by the standard analysis of the mechanism that this step is ε -DP [43]. By Lemma C.14, the majority voting step is also ε -DP. Hence, the result follows by composition (cf. Lemma A.4). \square

Next, we will analyze its utility.

Lemma C.21. *Let $\gamma\sqrt{d} \geq 12800$ and $x \in \{\pm 1\}^n$ be balanced. Let $\exp(-640) \geq \zeta \geq 2 \exp\left(-\frac{\gamma^2 d}{512}\right)$, $\varepsilon \geq \frac{64 \log(2/\zeta)}{\gamma d}$, and $\mathbf{G} \sim \text{SBM}_n(\gamma, d, x^*)$, the output $\hat{x}(\mathbf{G}) \in \{\pm 1\}^n$ of Algorithm C.19 satisfies*

$$\text{err}(\hat{x}(\mathbf{G}), x^*) \leq \zeta$$

with probability at least $1 - \zeta$.

Proof. We will first show that the rough estimate \tilde{x} obtained in step 2 achieves

$$\text{err}(\tilde{x}, x^*) \leq \frac{20}{\sqrt{\log(1/\zeta)}}$$

with probability e^{-n} . This will prove the second part of the theorem - for this we don't need that $\zeta \leq \exp(-640)$. In fact, arbitrary ζ works. The final error guarantee will then follow by Lemma C.16. First, notice that similar to the proof of [28, Lemma 4.1], using Bernstein's inequality and a union bound, we can show that (cf. Fact H.2 for a full proof)

$$\max_{x \in \{\pm 1\}^n} \left| \langle x, \left[Y(\mathbf{G}) - \frac{1}{n} x^* (x^*)^\top \right] x \rangle \right| \leq \frac{100n}{\gamma\sqrt{d}} \leq \frac{5}{\sqrt{\log(1/\zeta)}}$$

with probability at least $1 - \exp^{-10n}$. Recall that $s_{\mathbf{G}}(x) = \langle x, Y(\mathbf{G})x \rangle$. Let $\alpha = \frac{5}{\sqrt{\log(1/\zeta)}}$. We call $x \in \{\pm 1\}^n$ *good* if $s_{\mathbf{G}}(x) \geq (1 - 3\alpha)n$. It follows that for good x it holds that

$$\frac{1}{n} \cdot \langle x, x^* \rangle^2 \geq \langle x, Y(\mathbf{G})x \rangle - \left| \left\langle x, \left[Y(\mathbf{G}) - \frac{1}{n} x^* (x^*)^\top \right] x \right\rangle \right| \geq (1 - 4\alpha)n.$$

Which implies that

$$2 \text{err}(x, x^*) \leq 1 - \sqrt{1 - 4\alpha} = 1 - \frac{1 - 4\alpha}{\sqrt{1 - 4\alpha}} \leq 1 - \frac{1 - 4\alpha}{1 - 2\alpha} = \frac{2\alpha}{1 - 2\alpha} \leq 4\alpha,$$

where we used that $\alpha \leq 1/4$ and that $\sqrt{1-4x} \leq 1-2x$ for $x \geq 0$. Hence, we have for good x that

$$\text{err}(x, x^*) \leq \frac{20}{\sqrt{\log(1/\zeta)}}.$$

Since $s_{\mathbf{G}}(x^*) \geq (1-\alpha)n$, there is at least one good candidate. Hence, we can bound the probability that we do not output a good x as

$$\frac{\exp\left(\frac{\varepsilon}{2\Delta}(1-3\alpha)n\right) \cdot e^{-n}}{\exp\left(\frac{\varepsilon}{2\Delta}(1-\alpha)n\right) \cdot 1} = \exp\left(\left(1 - \frac{2\varepsilon\alpha}{\Delta}\right)n\right) \leq e^{-n},$$

where we used that

$$\frac{2\varepsilon\alpha}{\Delta} \geq \frac{64 \log(2/\zeta)}{\gamma d} \cdot \frac{5\gamma d}{\sqrt{\log(1/\zeta)}} \geq 320\sqrt{\log(1/\zeta)} \geq 2.$$

We will use Lemma C.16 to prove the final conclusion of the theorem. In what follows, assume without loss of generality that $\text{Ham}(x, x^*) < \text{Ham}(x, -x^*)$. The above discussion implies that

$$\text{Ham}(x, x^*) \leq 8\alpha n \leq \frac{40n}{\sqrt{\log(1/\zeta)}} \leq \frac{n}{16},$$

where the last inequality uses $\zeta \leq e^{-640}$. Further, by Fact E.2 it also follows that the maximum degree of \mathbf{G}_2 is at most $O(\log^2 n)$ (by some margin). Recall that $\mathbf{G}_2 \sim \text{SBM}(d, \gamma, x^*)$. In the parametrization of Lemma C.16 this means that

$$\begin{aligned} \alpha &= \frac{(1+\gamma)d}{\log n}, & \beta &= \frac{(1-\gamma)d}{\log n}, \\ \alpha - \beta &= \frac{2\gamma d}{\log n}, & \alpha + \beta &= \frac{2d}{\log n}. \end{aligned}$$

Thus, it follows that the output \hat{x} of the majority voting step satisfies for every vertex v

$$\begin{aligned} \mathbb{P}(\hat{x}(\mathbf{G})_v \neq x_v) &\leq \exp\left(-\frac{1}{64} \cdot \varepsilon(\alpha - \beta) \cdot \log n\right) + 2 \cdot \exp\left(-\frac{1}{16^2} \cdot \frac{(\alpha - \beta)^2}{\alpha + \beta} \cdot \log n\right) \\ &\leq \exp\left(-\frac{1}{32} \cdot \varepsilon\gamma d\right) + \exp\left(-\frac{1}{16^2} \cdot \gamma^2 d\right) \\ &\leq \zeta^2/4 + \zeta^2/4 \leq \zeta^2. \end{aligned}$$

By Markov's Inequality it now follows that

$$\mathbb{P}(\text{err}(\hat{x}(\mathbf{G}), x^*) \geq \zeta) \leq \zeta.$$

□

C.4 Lower bound on the parameters for private recovery

In this section, we prove a tight lower bound for private recovery for stochastic block models. Recall the definition of error rate, $\text{err}(u, v) := \frac{1}{n} \cdot \min\{\text{Ham}(u, v), \text{Ham}(u, -v)\}$ for $u, v \in \{\pm 1\}^n$. Our main result is the following theorem.

Theorem C.22 (Full version of Theorem 1.5). *Suppose there exists an ε -differentially private algorithm such that for any balanced $x \in \{\pm 1\}^n$, on input $\mathbf{G} \sim \text{SBM}_n(d, \gamma, x)$, outputs $\hat{x}(\mathbf{G}) \in \{\pm 1\}^n$ satisfying*

$$\mathbb{P}(\text{err}(\hat{x}(\mathbf{G}), x) < \zeta) \geq 1 - \eta,$$

where³⁰ $1/n \leq \zeta \leq 0.04$ and the randomness is over both the algorithm and stochastic block models. Then,

$$e^{2\varepsilon} - 1 \geq \Omega\left(\frac{\log(1/\zeta)}{\gamma d} + \frac{\log(1/\eta)}{\zeta n \gamma d}\right). \quad (\text{C.10})$$

³⁰Error rate less than $1/n$ already means exact recovery. Thus it does not make sense to set ζ to any value strictly smaller than $1/n$. The upper bound $\zeta \leq 0.04$ is just a technical condition our proof needs for Eq. (C.12).

Remark C.23. Both terms in lower bound Eq. (C.10) are tight up to constants by the following argument. Considering typical privacy parameters $\varepsilon \leq 1$, then $e^{2\varepsilon} - 1 \approx 2\varepsilon$. For exponentially small failure probability, i.e. $\eta = 2^{-\Omega(n)}$, the lower bound reads $\varepsilon \geq \Omega(\frac{1}{\gamma d} \cdot \frac{1}{\zeta})$, which is achieved by Algorithm C.19 without the boosting step - see the discussion after Theorem C.18. For polynomially small failure probability, i.e. $\eta = 1/\text{poly}(n)$, the lower bound Eq. (C.10) reads $\varepsilon \geq \Omega(\frac{1}{\gamma d} \cdot \log \frac{1}{\zeta})$, which is achieved by Theorem C.18.

By setting $\zeta = 1/n$ in Theorem C.22, we directly obtain a tight lower bound for private exact recovery as a corollary.

Corollary C.24. Suppose there exists an ε -differentially private algorithm such that for any balanced $x \in \{\pm 1\}^n$, on input $\mathbf{G} \sim \text{SBM}_n(d, \gamma, x)$, outputs $\hat{x}(\mathbf{G}) \in \{\pm 1\}^n$ satisfying

$$\mathbb{P}(\hat{x}(\mathbf{G}) \in \{x, -x\}) \geq 1 - \eta,$$

where the randomness is over both the algorithm and stochastic block models. Then,

$$e^{2\varepsilon} - 1 \geq \Omega\left(\frac{\log(n) + \log \frac{1}{\eta}}{\gamma d}\right). \quad (\text{C.11})$$

Remark C.25. The lower bound Eq. (C.11) for private exact recovery is tight up to constants, since there exists an (inefficient) ε -differentially private exact recovery algorithm with $\varepsilon \leq O(\frac{\log n}{\gamma d})$ and $\eta = 1/\text{poly}(n)$ by Theorem C.18 and [55, Theorem 3.7].

In rest of this section, we will prove Theorem C.22. The proof applies the packing lower bound argument similar to [29, Theorem 7.1]. To this end, we first show $\text{err}(\cdot, \cdot)$ is a semimetric over $\{\pm 1\}^n$.

Lemma C.26. $\text{err}(\cdot, \cdot)$ is a semimetric over $\{\pm 1\}^n$.

Proof. Symmetry and non-negativity are obvious from the definition. We will show $\text{err}(\cdot, \cdot)$ satisfies triangle inequality via case analysis. Let $u, v, w \in \{\pm 1\}^n$ be three arbitrary sign vectors. By symmetry, we only need to consider the following four cases.

Case 1: $\text{Ham}(u, v), \text{Ham}(u, w), \text{Ham}(v, w) \leq n/2$. This case is reduced to showing Hamming distance satisfies triangle inequality, which is obvious.

Case 2: $\text{Ham}(u, v), \text{Ham}(u, w) \leq n/2$ and $\text{Ham}(v, w) \geq n/2$. We need to check two subcases. First,

$$\begin{aligned} \text{err}(u, v) \leq \text{err}(u, w) + \text{err}(v, w) &\Leftrightarrow \text{Ham}(u, v) + \text{Ham}(v, w) \leq \text{Ham}(u, w) + n \\ &\Leftrightarrow \text{Ham}(u, v) + H(u, v) + H(u, w) \leq \text{Ham}(u, w) + n \\ &\Leftrightarrow \text{Ham}(u, v) \leq n/2. \end{aligned}$$

Second,

$$\begin{aligned} \text{err}(v, w) \leq \text{err}(u, v) + \text{err}(u, w) &\Leftrightarrow n \leq \text{Ham}(v, w) + \text{Ham}(u, v) + \text{Ham}(u, w) \\ &\Leftrightarrow n \leq 2 \text{Ham}(v, w). \end{aligned}$$

Case 3: $\text{Ham}(u, v) \leq n/2$ and $\text{Ham}(u, w), \text{Ham}(v, w) \geq n/2$. This case can be reduced to case 1 by considering $u, v, -w$.

Case 4: $\text{Ham}(u, v), \text{Ham}(u, w), \text{Ham}(v, w) \geq n/2$. This case can be reduced to case 2 by considering $-u, v, w$. \square

Proof of Theorem C.22. Suppose there exists an ε -differentially private algorithm satisfying the theorem's assumption.

We first make the following notation. Given a semimetric ρ over $\{\pm 1\}^n$, a center $v \in \{\pm 1\}^n$, and a radius $r \geq 0$, define $B_\rho(v, r) := \{w \in \{\pm 1\}^n : \mathbf{1}^\top w = 0, \rho(w, v) \leq r\}$.

Pick an arbitrary balanced $x \in \{\pm 1\}^n$. Let $M = \{x^1, x^2, \dots, x^m\}$ be a maximal 2ζ -packing of $B_{\text{err}}(x, 4\zeta)$ in semimetric $\text{err}(\cdot, \cdot)$. By maximality of M , we have $B_{\text{err}}(x, 4\zeta) \subseteq \cup_{i=1}^m B_{\text{err}}(x^i, 2\zeta)$, which implies

$$|B_{\text{err}}(x, 4\zeta)| \leq \sum_{i=1}^m |B_{\text{err}}(x^i, 2\zeta)|$$

$$\begin{aligned} &\implies |B_{\text{Ham}}(x, 4\zeta)| \leq \sum_{i=1}^m 2 \cdot |B_{\text{Ham}}(x^i, 2\zeta)| = 2m \cdot |B_{\text{Ham}}(x, 2\zeta)| \\ &\implies 2m \geq \frac{|B_{\text{Ham}}(x, 4\zeta n)|}{|B_{\text{Ham}}(x, 2\zeta n)|} = \frac{\binom{n/2}{2\zeta n}^2}{\binom{n/2}{\zeta n}^2} \geq \frac{\left(\frac{1}{4\zeta}\right)^{4\zeta n}}{\left(\frac{e}{2\zeta}\right)^{2\zeta n}} = \left(\frac{1}{8e\zeta}\right)^{2\zeta n} \end{aligned} \quad (\text{C.12})$$

For each $i \in [m]$, define $Y_i := \{w \in \{\pm 1\}^n : \text{err}(w, x^i) \leq \zeta\}$. Then Y_i 's are pairwise disjoint. For each $i \in [m]$, let P_i be the distribution over n -vertex graphs generated by $\text{SBM}_n(d, \gamma, x^i)$. By our assumption on the algorithm, we have for any $i \in [m]$ that

$$\mathbb{P}_{\mathbf{G} \sim P_i} (\hat{x}(\mathbf{G}) \in Y_i) \geq 1 - \eta.$$

Combining the fact that Y_i 's are pairwise disjoint, we have

$$\sum_{i=1}^m \mathbb{P}_{\mathbf{G} \sim P_1} (\hat{x}(\mathbf{G}) \in Y_i) = \mathbb{P}_{\mathbf{G} \sim P_1} (\hat{x}(\mathbf{G}) \in \cup_{i=1}^m Y_i) \leq 1 \implies \sum_{i=2}^m \mathbb{P}_{\mathbf{G} \sim P_1} (\hat{x}(\mathbf{G}) \in Y_i) \leq \eta. \quad (\text{C.13})$$

In the following, we will lower bound $\mathbb{P}_{\mathbf{G} \sim P_1} (\hat{x}(\mathbf{G}) \in Y_i)$ for each $i \in [m] \setminus \{1\}$ using group privacy.

Note each P_i is a product of $\binom{n}{2}$ independent Bernoulli distributions. Thus for any $i, j \in [m]$, there exists a coupling ω_{ij} of P_i and P_j such that, if $(\mathbf{G}, \mathbf{H}) \sim \omega$, then

$$\text{Ham}(\mathbf{G}, \mathbf{H}) \sim \text{Binomial}(N_{ij}, p),$$

where $p = 2\gamma d/n$ and $N_{ij} = \text{Ham}(x^i, x^j) \cdot (n - \text{Ham}(x^i, x^j))$. Applying group privacy, we have for any two graphs G, H and for any $S \subseteq \{\pm 1\}^n$ that³¹

$$\mathbb{P}(\hat{x}(G) \in S) \leq \exp(\varepsilon \cdot \text{Ham}(G, H)) \cdot \mathbb{P}(\hat{x}(H) \in S). \quad (\text{C.14})$$

For each $i \in [m]$, taking expectations on both sides of Eq. (C.14) with respect to coupling ω_{i1} and setting $S = Y_i$, we have

$$\mathbb{E}_{(\mathbf{G}, \mathbf{H}) \sim \omega_{i1}} \mathbb{P}(\hat{x}(\mathbf{G}) \in Y_i) \leq \mathbb{E}_{(\mathbf{G}, \mathbf{H}) \sim \omega_{i1}} \exp(\varepsilon \cdot \text{Ham}(\mathbf{G}, \mathbf{H})) \cdot \mathbb{P}(\hat{x}(\mathbf{H}) \in Y_i). \quad (\text{C.15})$$

The left side of Eq. (C.15) is equal to

$$\mathbb{E}_{(\mathbf{G}, \mathbf{H}) \sim \omega_{i1}} \mathbb{P}(\hat{x}(\mathbf{G}) \in Y_i) = \mathbb{P}_{\mathbf{G} \sim P_i} (\hat{x}(\mathbf{G}) \in Y_i) \geq 1 - \eta.$$

Upper bounding the right side of Eq. (C.15) by Cauchy-Schwartz inequality, we have

$$\begin{aligned} &\mathbb{E}_{(\mathbf{G}, \mathbf{H}) \sim \omega_{i1}} \exp(\varepsilon \cdot \text{Ham}(\mathbf{G}, \mathbf{H})) \cdot \mathbb{P}(\hat{x}(\mathbf{H}) \in Y_i) \\ &\leq \left(\mathbb{E}_{(\mathbf{G}, \mathbf{H}) \sim \omega_{i1}} \exp(2\varepsilon \cdot \text{Ham}(\mathbf{G}, \mathbf{H})) \right)^{1/2} \cdot \left(\mathbb{E}_{(\mathbf{G}, \mathbf{H}) \sim \omega_{i1}} \mathbb{P}(\hat{x}(\mathbf{H}) \in Y_i)^2 \right)^{1/2} \\ &= \left(\mathbb{E}_{\mathbf{X} \sim \text{Binomial}(N_{i1}, p)} \exp(2\varepsilon \cdot \mathbf{X}) \right)^{1/2} \cdot \left(\mathbb{E}_{\mathbf{H} \sim P_1} \mathbb{P}(\hat{x}(\mathbf{H}) \in Y_i)^2 \right)^{1/2}. \end{aligned}$$

Using the formula for the moment generating function of binomial distributions, we have

$$\mathbb{E}_{\mathbf{X} \sim \text{Binomial}(N_{i1}, p)} \exp(2\varepsilon \cdot \mathbf{X}) = (1 - p + p \cdot e^{2\varepsilon})^{N_{i1}},$$

and it is easy to see

$$\mathbb{E}_{\mathbf{H} \sim P_1} \mathbb{P}(\hat{x}(\mathbf{H}) \in Y_i)^2 = \mathbb{E}_{\mathbf{H} \sim P_1} (\mathbb{E} \mathbb{1}_{[\hat{x}(\mathbf{H}) \in Y_i]})^2 \leq \mathbb{P}_{\mathbf{H} \sim P_1} (\hat{x}(\mathbf{H}) \in Y_i).$$

Putting things together, Eq. (C.15) implies for each $i \in [m]$ that

$$\mathbb{P}_{\mathbf{H} \sim P_1} (\hat{x}(\mathbf{H}) \in Y_i) \geq \frac{(1 - \eta)^2}{(1 - p + p \cdot e^{2\varepsilon})^{N_{i1}}}. \quad (\text{C.16})$$

³¹In Eq. (C.14), the randomness only comes from the algorithm.

Since $x^i \in B_{\text{err}}(x, 4\zeta)$ for $i \in [m]$, by assuming $\zeta \leq 1/16$, we have

$$N_{i1} = \text{Ham}(x^i, x^1) \cdot (n - \text{Ham}(x^1, x^i)) \leq 8\zeta n(n - 8\zeta n). \quad (\text{C.17})$$

Recalling $p = 2\gamma d/n$ and combining Eq. (C.12), Eq. (C.13), Eq. (C.16) and Eq. (C.17), we have

$$(m-1) \cdot \frac{(1-\eta)^2}{(1-p+p \cdot e^{2\varepsilon})^{8\zeta n(n-8\zeta n)}} \leq \eta.$$

By taking logarithm on both sides, using $t \geq \log(1+t)$ for any $t > -1$, and assuming $\zeta \leq 1/(8e)$, we have

$$e^{2\varepsilon} - 1 \gtrsim \frac{\log \frac{1}{8e\zeta}}{\gamma d} + \frac{\log \frac{1}{\eta}}{\zeta n \gamma d}.$$

□

D Private algorithms for learning mixtures of spherical Gaussians

In this section we present a private algorithm for recovering the centers of a mixtures of k Gaussians (cf. Model 1.2). Let $\mathcal{Y} \subseteq (\mathbb{R}^d)^{\otimes n}$ be the collection of sets of n points in \mathbb{R}^d . We consider the following notion of adjacency.

Definition D.1 (Adjacent datasets). *Two datasets $Y, Y' \in \mathcal{Y}$ are said to be adjacent if $|Y \cap Y'| \geq n-1$.*

Remark D.2 (Problem parameters as public information). *We consider the parameters n, k, Δ to be public information given as input to the algorithm.*

Next we present the main theorem of the section.

Theorem D.3 (Privately learning spherical mixtures of Gaussians). *Consider an instance of Model 1.2. Let $t \in \mathbb{N}$ be such that $\Delta \geq O(\sqrt{tk^{1/t}})$. For $n \geq \Omega(k^{O(1)} \cdot d^{O(t)})$, $k \geq (\log n)^{1/5}$, there exists an algorithm, running in time $(nd)^{O(t)}$, that outputs vectors $\hat{\mu}_1, \dots, \hat{\mu}_\ell$ satisfying*

$$\max_{\ell \in [k]} \|\hat{\mu}_\ell - \mu_{\pi(\ell)}\|_2 \leq O(k^{-12}),$$

with high probability, for some permutation $\pi: [k] \rightarrow [k]$.³² Moreover, for $\varepsilon \geq k^{-10}$, $\delta \geq n^{-10}$, the algorithm is (ε, δ) -differentially private for any input Y .

We remark that our algorithm not only works for mixtures of Gaussians but for all mixtures of $2t$ -explicitly bounded distributions (cf. Definition A.19).

Our algorithm is based on the sum-of-squares hierarchy and at the heart lies the following sum-of-squares program. The indeterminates $z_{11}, \dots, z_{1k}, \dots, z_{nk}$ and vector-valued indeterminates μ'_1, \dots, μ'_k , will be central to the proof of Theorem D.3. Let n, k, t be fixed parameters.

$$\left\{ \begin{array}{ll} z_{i\ell}^2 = z_{i\ell} & \forall i \in [n], \ell \in [k] \quad (\text{indicators}) \\ \sum_{\ell \in [k]} z_{i\ell} \leq 1 & \forall i \in [n] \quad (\text{cluster mem.}) \\ z_{i\ell} \cdot z_{i\ell'} = 0 & \forall i \in [n], \ell \in [k] \quad (\text{uniq. mem.}) \\ \sum_i z_{i\ell} \leq n/k & \forall \ell \in [k] \quad (\text{size of clusters}) \\ \mu'_\ell = \frac{k}{n} \sum_i z_{i\ell} \cdot y_i & \forall \ell \in [k] \quad (\text{means of clusters}) \\ \forall v \in \mathbb{R}^d: \frac{k}{n} \sum_{i=1}^n z_{i\ell} \langle y_i - \mu'_\ell, v \rangle^{2s} + \|Qv^{\otimes s}\|^2 = (2s)^s \cdot \|v\|_2^{2s} & \forall s \leq t, \ell \in [k] \quad (t \text{ moment}) \end{array} \right. \quad (\mathcal{P}_{n,k,t}(Y))$$

³²We remark that we chose constants to optimize readability and not the smallest possible ones.

We remark that the moment constraint encodes the $2t$ -explicit 2-boundedness constraint introduced in Definition A.19. Note that in the form stated above there are infinitely many constraints, one for each vector v . This is just for notational convenience. This constraint postulates equality of two polynomials in v . Formally, this can also be encoded by requiring their coefficients to agree and hence eliminating the variable v . It is not hard to see that this can be done adding only polynomially many constraints. Further, the matrix variable Q represents the SOS proof of the $2t$ -explicit 2-boundedness constraint and we can hence deduce that for all $0 \leq s \leq t$

$$\mathcal{P} \Big|_{2s} \left\{ \frac{k}{n} \sum_{i=1}^n z_{i\ell} \langle y_i - \mu'_\ell, v \rangle^{2s} \leq (2s)^s \|s\|_2^{2s} \right\}.$$

Before presenting the algorithm we will introduce some additional notation which will be convenient. We assume t, n, k to be *fixed* throughout the section and drop the corresponding subscripts. For $Y \in \mathcal{Y}$, let $\mathcal{Z}(Y)$ be the set of degree- $10t$ pseudo-distributions satisfying $\mathcal{P}(Y)$. For each $\zeta \in \mathcal{Z}(Y)$ define $W(\zeta)$ as the n -by- n matrix satisfying

$$W(\zeta)_{ij} = \tilde{\mathbb{E}}_\zeta \left[\sum_{\ell \in [k]} z_{i\ell} \cdot z_{j\ell} \right].$$

We let $\mathcal{W}(Y) := \{W(\zeta) \mid \zeta \in \mathcal{Z}(Y)\}$.

Recall that J denotes the all-ones matrix. We define the function $g : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$ as

$$g(W) = \|W\|_F^2 - (10)^{10} k^{300} \langle J, W \rangle$$

and let

$$W(\hat{\zeta}(Y)) := \operatorname{argmin}_{W \in \mathcal{W}(Y)} g(W).$$

We also consider the following function

Definition D.4 (Soft thresholding function). *We denote by $\phi : [0, 1] \rightarrow [0, 1]$ the function*

$$\phi(x) = \begin{cases} 0 & \text{if } x \leq 0.8, \\ 1 & \text{if } x \geq 0.9, \\ \frac{x-0.8}{0.9-0.8} & \text{otherwise.} \end{cases}$$

Notice that $\phi(\cdot)$ is $\frac{1}{0.9-0.8} = 10$ Lipschitz. Next we introduce our algorithm. Notice the algorithm relies on certain private subroutines. We describe them later in the section to improve the presentation.

Algorithm D.5 (Private algorithm for learning mixtures of Gaussians).

Input: Set of n points $Y \subseteq \mathbb{R}^d$, $\varepsilon, \delta > 0$, $k, t \in \mathbb{N}$, $d^* = 100 \log n$, $b = k^{-15}$.

1. Compute $W = W(\hat{\zeta}(Y))$.
2. Pick $\tau \sim tLap\left(-n^{1.6} \left(1 + \frac{\log(1/\delta)}{\varepsilon}\right), \frac{n^{1.6}}{\varepsilon}\right)$.
3. If $|\tau| \geq n^{1.7}$ or $\|\phi(W)\|_1 \leq \frac{n^2}{k} \cdot \left(1 - \frac{1}{n^{0.1}} - \frac{1}{k^{100}}\right) + \tau$ reject.
4. For all $i \in [n]$, compute the n -dimensional vector

$$\nu^{(i)} = \begin{cases} \mathbf{0} & \text{if } \|\phi(W_i)\|_1 = 0 \\ \|\phi(W_i)\|_1^{-1} \sum_j \phi(W_{ij}) \cdot y_j & \text{otherwise.} \end{cases}$$

5. Pick a set \mathcal{S} of $n^{0.01}$ indices $i \in [n]$ uniformly at random.
6. For each $i \in \mathcal{S}$ let $\bar{\nu}^{(i)} = \nu^{(i)} + \mathbf{w}$ where $\mathbf{w} \sim N\left(0, n^{-0.18} \cdot \frac{\log(2/\delta)}{\varepsilon^2} \cdot \text{Id}\right)$.
7. Pick $\Phi \sim N\left(0, \frac{1}{d^*}\right)^{d^* \times d}$, $\mathbf{q} \stackrel{u.a.r.}{\sim} [0, b]$ and run the histogram learner of Lemma A.13 with input $\Phi \bar{\nu}^{(1)}, \dots, \Phi \bar{\nu}^{(n^{0.01})}$ and parameters

$$\mathbf{q}, b, \alpha = k^{-10}, \beta = n^{-10}, \delta^* = \frac{\delta}{n}, \varepsilon^* = \varepsilon \cdot \frac{10k^{50}}{n^{0.01}}.$$

Let $\mathbf{B}_1, \dots, \mathbf{B}_k$ be the resulting d^* -dimensional bins with highest counts. Break ties randomly.

8. Reject if $\min_{i \in [k]} |\{j \mid \Phi \bar{\nu}^{(j)} \in \mathbf{B}_i\}| < \frac{n^{0.01}}{2k}$.
9. For each $l \in [k]$ output

$$\hat{\mu}_l := \frac{1}{|\{j \mid \Phi \bar{\nu}^{(j)} \in \mathbf{B}_l\}|} \cdot \left(\sum_{\Phi \bar{\nu}^{(j)} \in \mathbf{B}_l} \bar{\nu}^{(j)} \right) + \mathbf{w}',$$

where $\mathbf{w}' \sim N\left(0, N\left(0, 32 \cdot k^{-120} \cdot \frac{\log(2kn/\delta)}{\varepsilon^2} \cdot \text{Id}\right)\right)$.

For convenience, we introduce some preliminary facts.

Definition D.6 (Good Y). Let Y be sampled according to Model 1.2. We say that Y is good if:

1. for each $\ell \in [k]$, there are at least $\frac{n}{k} - n^{0.6}$ and most $\frac{n}{k} + n^{0.6}$ points sampled from D_ℓ in Y . Let $Y_\ell \subseteq Y$ be such set of points.
2. Each Y_ℓ is $2t$ -explicitly 2-bounded.

It turns out that typical instances Y are indeed good.

Lemma D.7 ([30, 36]). Consider the settings of Theorem D.3. Then Y is good with high probability. Further, in this case the sets $\mathcal{Z}(Y)$ and $\mathcal{W}(Y)$ are non-empty.

D.1 Privacy analysis

In this section we show that our clustering algorithm is private.

Lemma D.8 (Differential privacy of the algorithm). Consider the settings of Theorem D.3. Then Algorithm D.5 is (ε, δ) -differentially private.

We split our analysis in multiple steps and combine them at the end. On a high level, we will argue that on adjacent inputs Y, Y' many of the vectors $\nu^{(i)}$ by the algorithm are close to each other and a small part can be very far. We can then show that we can mask this small difference using the

Gaussian mechanism and afterwards treat this subset of the vectors as privatized (cf. Lemma I.4). Then we can combine this with known histogram learners to deal with the small set of $\nu^{(i)}$'s that is far from each other on adjacent inputs.

D.1.1 Sensitivity of the matrix W

Here we use Lemma B.1 to reason about the sensitivity of $\phi(W(\hat{\zeta}(Y)))$. For adjacent datasets $Y, Y' \in \mathcal{Y}$ we let $\hat{\zeta}, \hat{\zeta}'$ be the pseudo-distribution corresponding to $W(\hat{\zeta}(Y))$ and $W(\hat{\zeta}(Y'))$ computed in step 1 of the algorithm, respectively. We prove the following result.

Lemma D.9 (ℓ_1 -sensitivity of $\phi(W)$). *Consider the settings of Theorem D.3. Let W, W' be respectively be the matrices computed in step 1 by Algorithm D.5 on adjacent inputs $Y, Y' \in \mathcal{Y}$. Then*

$$\|\phi(W) - \phi(W')\|_1 \leq n^{1.6}.$$

For all but $n^{0.8}$ rows i of $\phi(W), \phi(W')$, it holds

$$\|\phi(W)_i - \phi(W')_i\|_1 \leq n^{0.8}.$$

Proof. The second inequality is an immediate consequence of the first via Markov's inequality. Thus it suffices to prove the first. Since $\phi(\cdot)$ is 10-Lipschitz, we immediately obtain the result if

$$\left\| W(\hat{\zeta}(Y)) - W(\hat{\zeta}(Y')) \right\|_1 \leq n^{1.55}.$$

Thus we focus on this inequality. To prove it, we verify the two conditions of Lemma B.1. First notice that g is 2-strongly convex with respect to its input W . Indeed for $W, W' \in \mathcal{W}(Y)$, since $\forall i, j \in [n], W_{ij} \geq 0$ it holds that

$$\begin{aligned} \|W'\|_F^2 &= \|W\|_F^2 + \|W - W'\|_F^2 + 2\langle W' - W, W \rangle \\ &= \|W\|_F^2 + \|W - W'\|_F^2 + 2\langle W' - W, W \rangle + \langle W' - W, (10)^{10} k^{300} (J - J) \rangle \\ &= g(W) + \|W - W'\|_F^2 + \langle W' - W, \nabla g(W) \rangle + \langle W', (10)^{10} k^{300} J \rangle, \end{aligned}$$

where we used that $\nabla g(W) = 2W - (10)^{10} k^{300} J$. Thus it remain to prove (i) of Lemma B.1.

Let $\hat{\zeta} \in \mathcal{Z}(Y), \hat{\zeta}' \in \mathcal{Z}(Y')$ be the pseudo-distributions such that $W_Y(\hat{\zeta}) = W$ and $W_{Y'}(\hat{\zeta}') = W'$. We claim that there always exists $\zeta_{\text{adj}} \in \mathcal{Z}(Y) \cap \mathcal{Z}(Y')$ such that

1. $|g(W(\zeta)) - g(W(\zeta_{\text{adj}}))| \leq \frac{2n}{k} \cdot ((10)^{10} k^{300} + 1) \leq 3 \cdot (10)^{10} k^{300} n,$
2. $|g_{Y'}(W(\zeta_{\text{adj}})) - g(W(\zeta_{\text{adj}}))| = 0.$

Note that in this case the second point is always true since g doesn't depend on Y . Together with Lemma B.1 these two inequalities will imply that

$$\left\| W(\hat{\zeta}(Y)) - W(\hat{\zeta}(Y')) \right\|_F^2 \leq 18 \cdot (10)^{10} k^{300} n.$$

By assumption on n , an application of Cauchy-Schwarz will give us the desired result.

So, let i be the index at which Y, Y' differ. We construct ζ_{adj} as follows: for all polynomials p of degree at most $10t$ we let

$$\tilde{\mathbb{E}}_{\zeta_{\text{adj}}}[p] = \begin{cases} \tilde{\mathbb{E}}_{\zeta}[p] & \text{if } p \text{ does not contain variables } z_{i\ell} \text{ for any } \ell \in [k] \\ 0 & \text{otherwise.} \end{cases}$$

By construction $\zeta_{\text{adj}} \in \mathcal{Z}(Y) \cap \mathcal{Z}(Y')$. Moreover, $W(\zeta), W(\zeta_{\text{adj}})$ differ in at most $2n/k$ entries. Since all entries of the two matrices are in $[0, 1]$, the first inequality follows by definition of the objective function. \square

D.1.2 Sensitivity of the resulting vectors

In this section we argue that if the algorithm does not reject in step 3 then the vectors $\nu^{(i)}$ are stable on adjacent inputs. Concretely our statement goes as follows:

Lemma D.10 (Stability of the $\nu^{(i)}$'s). *Consider the settings of Theorem D.3. Suppose Algorithm D.5 does not reject in step 3, on adjacent inputs $Y, Y' \in \mathcal{Y}$. Then for all but $\frac{6n}{k^{50}}$ indices $i \in [n]$, it holds:*

$$\left\| \nu_Y^{(i)} - \nu_{Y'}^{(i)} \right\|_2 \leq O(n^{-0.1}).$$

The proof of Lemma D.10 crucially relies on the next statement.

Lemma D.11 (Covariance bound). *Consider the settings of Theorem D.3. Let W be the matrix computed by Algorithm D.5 on input $Y \in \mathcal{Y}$. For $i \in [n]$, if $\|\phi(W_i)\|_1 \geq \frac{n}{k} \cdot (1 - \frac{10}{k^{50}})$ then $\nu^{(i)}$ induces a 2-explicitly 40-bounded distribution over Y .*

Proof. First, by assumption notice that there must be at least $\frac{n}{k} \cdot (1 - \frac{10}{k^{50}})$ entries of $\phi(W_i)$ larger than 0.8. We denote the set of $j \in [n]$ such that $W_{ij} \geq 0.8$ by \mathcal{G} . Let $\zeta \in \mathcal{Z}(Y)$ be the degree 10t pseudo-distribution so that $W = W(\zeta(Y))$. Since ζ satisfies $\mathcal{P}(Y)$, for $\ell \in [k]$ it follows from the moment bound constraint for $s = 1$ that for all unit vectors u it holds that

$$\mathcal{P} \Big|_{\frac{1}{4}} \left\{ 0 \leq \frac{k}{n} \sum_{j=1}^n z_{j\ell} \langle \mathbf{y}_j - \mu'_\ell, u \rangle^2 \leq 2 \right\},$$

Using the SOS triangle inequality (cf. Fact I.2) $\left| \frac{a \cdot b}{2} (a + b)^2 \leq 2(a^2 + b^2) \right|$ it now follows that

$$\mathbf{0} \preceq \tilde{\mathbb{E}}_\zeta \left[\frac{k^2}{n^2} \sum_{j, j' \in [n]} z_{j\ell} z_{j'\ell} \cdot (y_j - y_{j'})^{\otimes 2} \right] \preceq 8\text{Id}$$

and thus

$$\mathbf{0} \preceq \tilde{\mathbb{E}}_\zeta \left[\frac{k^2}{n^2} \sum_{\ell \in [k]} \sum_{j, j' \in [n]} z_{i\ell} z_{j\ell} z_{j'\ell} \cdot (y_j - y_{j'})^{\otimes 2} \right] \preceq 8\text{Id}.$$

Furthermore using $\mathcal{P}(Y) \Big|_{\frac{1}{2}} \{z_{i\ell} z_{i\ell'} = 0\}$ for $\ell \neq \ell'$ we have

$$\tilde{\mathbb{E}}_\zeta \left[\sum_{\ell \in [k]} \sum_{j, j' \in [n]} z_{i\ell} z_{j\ell} z_{j'\ell} \right] = \tilde{\mathbb{E}}_\zeta \left[\left(\sum_{\ell \in [k], j \in [n]} z_{i\ell} z_{j\ell} \right) \cdot \left(\sum_{\ell' \in [k], j' \in [n]} z_{i\ell'} z_{j'\ell'} \right) \right].$$

Now, for fixed $j, j' \in [n]$, using

$$\{a^2 = a, b^2 = b\} \Big|_{O(1)} \{1 + ab - a - b = 1 - ab - (a - b)^2 \geq 0\}$$

with $a = \sum_{\ell \in [k]} z_{i\ell} z_{j\ell}$ and $b = \sum_{\ell' \in [k]} z_{i\ell'} z_{j'\ell'}$ we get

$$\begin{aligned} \tilde{\mathbb{E}}_\zeta \left[\left(\sum_{\ell \in [k]} z_{i\ell} z_{j\ell} \right) \left(\sum_{\ell' \in [k]} z_{i\ell'} z_{j'\ell'} \right) \right] &\geq \tilde{\mathbb{E}}_\zeta \left[\sum_{\ell \in [k]} z_{i\ell} z_{j\ell} + \sum_{\ell' \in [k]} z_{i\ell'} z_{j'\ell'} \right] - 1 \\ &= W_{ij} + W_{ij'} - 1. \end{aligned}$$

Now if $j, j' \in \mathcal{G}$ we must have

$$\sum_{\ell \in [k]} \tilde{\mathbb{E}}_\zeta [z_{i\ell} z_{j\ell} z_{j'\ell}] = \tilde{\mathbb{E}}_\zeta \left[\left(\sum_{\ell \in [k]} z_{i\ell} z_{j\ell} \right) \left(\sum_{\ell' \in [k]} z_{i\ell'} z_{j'\ell'} \right) \right] \geq 0.6.$$

Since $\phi(W_{ij}) \leq 1$ by definition and $\|\phi(W_i)\|_1 \geq \frac{n}{k} \cdot (1 - \frac{10}{k^{50}})$, we conclude

$$\|\phi(W_i)\|_1^{-2} \left[\sum_{j, j' \in [n]} \phi(W_{ij}) \phi(W_{ij'}) (y_j - y_{j'})^{\otimes 2} \right]$$

$$\begin{aligned} &\preceq 5 \cdot \frac{k^2}{n^2} \sum_{j, j' \in [n], \ell \in [k]} \tilde{\mathbb{E}}_{\zeta} [z_{i\ell} z_{j\ell} z_{j'\ell}] \cdot (y_j - y_{j'})^{\otimes 2} \\ &\preceq 40 \text{Id}. \end{aligned}$$

as desired. \square

We can now prove Lemma D.10.

Proof of Lemma D.10. Let W, W' be the matrices computed by Algorithm D.5 in step 1 on input Y, Y' , respectively. Let $\mathcal{G} \subseteq [n]$ be the set of indices i such that

$$\|\phi(W)_i - \phi(W')_i\|_1 \leq n^{0.8}.$$

Notice that $|\mathcal{G}| \geq n - n^{0.8}$ by Lemma D.9. Since on input Y the algorithm did not reject in step 3 we must have

$$\|\phi(W)\|_1 \geq \frac{n^2}{k} \cdot \left(1 - \frac{1}{n^{0.1}} - \frac{1}{k^{100}}\right) - n^{1.7} \geq \frac{n^2}{k} \cdot \left(1 - \frac{2}{k^{100}}\right).$$

Let g_W be the number of indices $i \in \mathcal{G}$ such that $\|\phi(W)_i\|_1 \geq \frac{n}{k} \cdot \left(1 - \frac{1}{k^{50}}\right)$. It holds that

$$\begin{aligned} \frac{n^2}{k} \cdot \left(1 - \frac{2}{k^{100}}\right) &\leq g_W \cdot \frac{n}{k} + (n - |\mathcal{G}|) \cdot \frac{n}{k} + (|\mathcal{G}| - g_W) \frac{n}{k} \cdot \left(1 - \frac{1}{k^{50}}\right) \\ &\leq g_W \cdot \frac{n}{k} \cdot \frac{1}{k^{50}} + \frac{n^{1.8}}{k} + \frac{n^2}{k} \cdot \left(1 - \frac{1}{k^{50}}\right) \\ &\leq g_W \cdot \frac{n}{k} \cdot \frac{1}{k^{50}} + \frac{n^2}{k} \cdot \left(1 + \frac{1}{k^{100}} - \frac{1}{k^{50}}\right). \end{aligned}$$

Rearranging now yields

$$g_W \geq n \cdot \left(1 - \frac{3}{k^{50}}\right).$$

Similarly, let $g_{W'}$ be the number of indices $i \in \mathcal{G}$ such that $\|\phi(W')_i\|_1 \geq \frac{n}{k} \cdot \left(1 - \frac{1}{k^{50}}\right)$. By an analogous argument it follows that $g_{W'} \geq n \cdot \left(1 - \frac{3}{k^{50}}\right)$. Thus, by the pigeonhole principle there are at least $g_W \geq n \cdot \left(1 - \frac{6}{k^{50}}\right)$ indices i such that

1. $\|\phi(W)_i\|_1 \geq \frac{n}{k} \left(1 - \frac{1}{k^{50}}\right)$,
2. $\|\phi(W')_i\|_1 \geq \frac{n}{k} \left(1 - \frac{1}{k^{50}}\right)$,
3. $\|\phi(W)_i - \phi(W')_i\|_1 \leq n^{0.8}$.

Combining these with Lemma D.11 we may also add

4. the distribution induced by $\|\phi(W_i)\|_1^{-1} \phi(W_i)$ is 2-explicitly 40-bounded,
5. the distribution induced by $\|\phi(W'_i)\|_1^{-1} \phi(W'_i)$ is 2-explicitly 40-bounded.

Using that for non-zero vectors x, y it holds that $\left\| \frac{x}{\|x\|} - \frac{y}{\|y\|} \right\| \leq \frac{2}{\|x\|} \|x - y\|$ points 1 to 3 above imply that

$$\left\| \|\phi(W_i)\|_1^{-1} \phi(W_i) - \|\phi(W'_i)\|_1^{-1} \phi(W'_i) \right\| \leq \frac{2n^{0.8}}{\frac{n}{k} \cdot \left(1 - \frac{1}{k^{50}}\right)} = O(n^{-0.2}).$$

Hence, applying Theorem A.21 with $t = 1$ it follows that

$$\left\| \nu_Y^{(i)} - \nu_{Y'}^{(i)} \right\|_2 \leq O(n^{-0.1}).$$

\square

D.1.3 From low sensitivity to privacy

In this section we argue privacy of the whole algorithm, proving Lemma D.8. Before doing that we observe that low-sensitivity is preserved with high probability under subsampling.

Fact D.12 (Stability of \mathcal{S}). *Consider the settings of Theorem D.3. Suppose Algorithm D.5 does not reject in step 3, on adjacent inputs $Y, Y' \in \mathcal{Y}$. With probability at least $1 - e^{-n^{\Omega(1)}}$ over the random choices of \mathcal{S} , for all but $\frac{10n^{0.01}}{k^{50}}$ indices $i \in \mathcal{S}$, it holds:*

$$\left\| \nu_Y^{(i)} - \nu_{Y'}^{(i)} \right\|_2 \leq O(n^{-0.1}).$$

Proof. There are at most $\frac{6n}{k^{50}}$ such indices in $[n]$ by Lemma D.10. By Chernoff's bound, cf. Fact E.4, the claim follows. \square

Finally, we prove our main privacy lemma.

Proof of Lemma D.8. For simplicity, we will prove that the algorithm is $(5\varepsilon, 5\delta)$ -private. Let $Y, Y' \in \mathcal{Y}$ be adjacent inputs. By Lemma A.10 and Lemma D.9 the test in step 3 of Algorithm D.5 is (ε, δ) -private.

Thus suppose now the algorithm did not reject in step 3 on inputs Y, Y' . By composition (cf. Lemma A.4) it is enough to show that the rest of the algorithm is (ε, δ) -private with respect to Y, Y' under this condition. Next, let $\nu_Y^{(1)}, \dots, \nu_Y^{(n)}$ and $\nu_{Y'}^{(1)}, \dots, \nu_{Y'}^{(n)}$ be the vectors computed in step 4 of the algorithm and \mathcal{S} be the random set of indices computed in step 5.³³ By Lemma D.10 and Fact D.12 with probability $1 - e^{-n^{\Omega(1)}}$ over the random choices of \mathcal{S} we get that for all but $\frac{10n^{0.01}}{k^{50}}$ indices $i \in \mathcal{S}$, it holds that

$$\left\| \nu_Y^{(i)} - \nu_{Y'}^{(i)} \right\|_2 \leq O(n^{-0.1}).$$

Denote this set of indices by \mathcal{G} . Note, that we may incorporate the failure probability $e^{-n^{\Omega(1)}} \leq \min\{\varepsilon/2, \delta/2\}$ into the final privacy parameters using Fact I.3.

Denote by \mathbf{V}, \mathbf{V}' the $|\mathcal{S}|$ -by- d matrices respectively with rows $\nu_Y^{(i_1)}, \dots, \nu_Y^{(i_{|\mathcal{S}|})}$ and $\nu_{Y'}^{(i_1)}, \dots, \nu_{Y'}^{(i_{|\mathcal{S}|})}$, where $i_1, \dots, i_{|\mathcal{S}|}$ are the indices in \mathcal{S} . Recall, that $|\mathcal{G}|$ rows of \mathbf{V} and \mathbf{V}' differ by at most $O(n^{-0.1})$ in ℓ_2 -norm. Thus, by the Gaussian mechanism used in step 6 (cf. Lemma A.12) and Lemma I.4 it is enough to show that step 7 to step 9 of the algorithm are private with respect to pairs of inputs V and V' differing in at most 1 row.³⁴ In particular, suppose these steps are $(\varepsilon_1, \delta_1)$ -private. Then, for $m = n^{0.01} - |\mathcal{G}| \leq \frac{10n^{0.01}}{k^{50}}$, by Lemma I.4 it follows that step 6 to step 9 are (ε', δ') -differentially private with

$$\begin{aligned} \varepsilon' &:= \varepsilon + m\varepsilon_1, \\ \delta' &:= e^\varepsilon m e^{(m-1)\varepsilon_1} \delta_1 + \delta. \end{aligned}$$

Consider steps 7 and 8. Recall, that in step 7 we invoke the histogram learner with parameters

$$b = k^{-15}, \mathbf{q} \stackrel{u.a.r.}{\sim} [0, b], \alpha = k^{-10}, \beta = n^{-10}, \delta^* = \frac{\delta}{n}, \varepsilon^* = \varepsilon \cdot \frac{10k^{50}}{n^{0.01}}.$$

Hence, by Lemma A.13 this step is $(\varepsilon^*, \delta^*)$ -private since

$$\frac{8}{\varepsilon^* \alpha} \cdot \log \left(\frac{2}{\delta^* \beta} \right) \leq \frac{200 \cdot k^{10} \cdot n^{0.01}}{10 \cdot k^{50} \cdot \varepsilon} \cdot \log n = \frac{20 \cdot n^{0.01}}{k^{40} \cdot \varepsilon} \cdot \log n \leq n,$$

for $\varepsilon \geq k^{-10}$. Step 8 is private by post-processing.

³³Note that since this does not depend on Y or Y' , respectively, we can assume this to be the same in both cases. Formally, this can be shown, e.g., via a direct calculation or using Lemma A.4.

³⁴Note that for the remainder of the analysis, these do *not* correspond to \mathbf{V} and \mathbf{V}' , since those differ in m rows. Lemma I.4 handles this difference.

Next, we argue that step 9 is private by showing that the average over the bins has small ℓ_2 -sensitivity. By Lemma A.4 we can consider the bins $\mathbf{B}_1, \dots, \mathbf{B}_k$ computed in the previous step as fixed. Further, we can assume that the algorithm did not reject in step 8, i.e., that each bin contains at least $\frac{n^{0.01}}{2k}$ points of V and V' respectively. As a consequence, every bin contains at least two (projections of) points of the input V or V' respectively. In particular, it contains at least one (projection of a) point which is present in both V and V' . Fix a bin \mathbf{B}_l and let \bar{v}^* be such that it is both in V and V' and $\Phi \bar{v}^* \in \mathbf{B}_l$. Also, define

$$S_l := \left| \left\{ j \mid \Phi \bar{v}_Y^{(j)} \in \mathbf{B}_l \right\} \right|,$$

$$S'_l := \left| \left\{ j \mid \Phi \bar{v}_{Y'}^{(j)} \in \mathbf{B}_l \right\} \right|.$$

Assume V and V' differ on index j . We consider two cases. First, assume that $\Phi \bar{v}_Y^{(j)}$ and $\Phi \bar{v}_{Y'}^{(j)}$ both lie in \mathbf{B}_l . In this case, $S_l = S'_l$ and using Lemma E.5 it follows that with probability $n^{-100} \leq \min\{\varepsilon/2, \delta/2\}$ it holds that

$$\begin{aligned} \left\| \bar{v}_Y^{(j)} - \bar{v}_{Y'}^{(j)} \right\|_2 &\leq \left\| \bar{v}_Y^{(j)} - \bar{v}^* \right\|_2 + \left\| \bar{v}^* - \bar{v}_{Y'}^{(j)} \right\|_2 \leq 10 \cdot \left(\left\| \Phi \bar{v}_Y^{(j)} - \Phi \bar{v}^* \right\|_2 + \left\| \Phi \bar{v}_{Y'}^{(j)} - \Phi \bar{v}^* \right\|_2 \right) \\ &\leq 20 \cdot \sqrt{d^*} \cdot b \leq 200 \cdot k^{-12}. \end{aligned}$$

And hence we can bound

$$\left\| \frac{1}{S_l} \cdot \left(\sum_{\Phi \bar{v}_Y^{(j)} \in \mathbf{B}_l} \bar{v}_Y^{(j)} \right) - \frac{1}{S'_l} \cdot \left(\sum_{\Phi \bar{v}_{Y'}^{(j)} \in \mathbf{B}_l} \bar{v}_{Y'}^{(j)} \right) \right\|_2 \leq \frac{\left\| \bar{v}_Y^{(j)} - \bar{v}_{Y'}^{(j)} \right\|_2}{S_l} \leq \frac{400 \cdot k^{-11}}{n^{0.01}}.$$

Next, assume that $\Phi \bar{v}_Y^{(j)} \notin \mathbf{B}_l$ and $\Phi \bar{v}_{Y'}^{(j)} \in \mathbf{B}_l$ (the other case works symmetrically). It follows that $S_l = S'_l - 1$ and we can bound

$$\begin{aligned} \left\| \frac{1}{S_l} \cdot \left(\sum_{\Phi \bar{v}_Y^{(j)} \in \mathbf{B}_l} \bar{v}_Y^{(j)} \right) - \frac{1}{S'_l} \cdot \left(\sum_{\Phi \bar{v}_{Y'}^{(j)} \in \mathbf{B}_l} \bar{v}_{Y'}^{(j)} \right) \right\|_2 &= \frac{1}{S_l \cdot S'_l} \cdot \left\| S'_l \left(\sum_{\Phi \bar{v}_Y^{(j)} \in \mathbf{B}_l} \bar{v}_Y^{(j)} \right) - (S'_l - 1) \left(\sum_{\Phi \bar{v}_{Y'}^{(j)} \in \mathbf{B}_l} \bar{v}_{Y'}^{(j)} \right) \right\|_2 \\ &= \frac{1}{S_l \cdot S'_l} \cdot \left\| S'_l \cdot \bar{v}_{Y'}^{(j)} + \left(\sum_{\Phi \bar{v}_{Y'}^{(j)} \in \mathbf{B}_l} \bar{v}_{Y'}^{(j)} \right) \right\|_2 \\ &= \frac{1}{S_l} \cdot \left\| \bar{v}_{Y'}^{(j)} - \frac{1}{S'_l} \left(\sum_{\Phi \bar{v}_{Y'}^{(j)} \in \mathbf{B}_l} \bar{v}_{Y'}^{(j)} \right) \right\|_2 \\ &\leq \frac{\sqrt{d^*} \cdot b}{S_l} \leq \frac{20 \cdot k^{-11}}{n^{0.01}}. \end{aligned}$$

Hence, the ℓ_2 -sensitivity is at most $\Delta := \frac{400 \cdot k^{-11}}{n^{0.01}}$. Since

$$2\Delta^2 \cdot \frac{\log(2/(\delta^*/k))}{(\varepsilon^*/k)^2} = 32 \cdot k^{-120} \cdot \frac{\log(2kn/\delta)}{\varepsilon^2}$$

and $\mathbf{w}' \sim N\left(0, 32 \cdot k^{-120} \cdot \frac{\log(2kn/\delta)}{\varepsilon^2} \cdot \text{Id}\right)$ it follows that outputting $\hat{\mu}_l$ is $(\varepsilon^*/k, \delta^*/k)$ -DP by the Gaussian Mechanism that. By Lemma A.4 it follows step 9 is $(\varepsilon^*, \delta^*)$ -private.

Hence, by Lemma A.4 it follows that step 7 to step 9 are $(2\varepsilon^*, 2\delta^*)$ -differentially private. Using $m \leq \frac{10n^{0.01}}{k^{10}}$ it now follows by Lemma I.4 that step 6 to step 9 are (ε', δ') -private for

$$\begin{aligned} \varepsilon' &= \varepsilon + 2m\varepsilon^* \leq 3\varepsilon, \\ \delta' &= 2e^\varepsilon m e^{(m-1)2\varepsilon^*} \delta^* + \delta \leq 2me^{3\varepsilon} \cdot \frac{\delta}{n} + \delta \leq 3\delta. \end{aligned}$$

Thus, combined with the private check and Fact I.3 in step 3 the whole algorithm is $(5\varepsilon, 5\delta)$ -private. \square

D.2 Utility analysis

In this section we reason about the utility of Algorithm D.5 and prove Theorem D.3. We first introduce some notation.

Definition D.13 (True solution). *Let \mathbf{Y} be an input sampled from Model 1.2. Denote by $W^*(\mathbf{Y}) \in \mathcal{W}(\mathbf{Y})$ the matrix induced by the true solution (or ground truth). I.e., let*

$$W^*(\mathbf{Y})_{ij} = \begin{cases} 1 & \text{if } i, j \text{ were both sampled from the same component of the mixture,} \\ 0 & \text{otherwise.} \end{cases}$$

Whenever the context is clear, we simply write \mathbf{W}^* to ease the notation.

First, we show that in the utility case step 3 of Algorithm D.5 rejects only with low probability.

Lemma D.14 (Algorithm does not reject on good inputs). *Consider the settings of Theorem D.3. Suppose \mathbf{Y} is a good set as per Definition D.6. Then $\|W(\hat{\zeta}(\mathbf{Y}))\|_1 \geq \frac{n^2}{k} \cdot \left(1 - n^{-0.4} - \frac{1}{(10)^{10}k^{300}}\right)$ and Algorithm D.5 rejects with probability at most $\exp(-\Omega(n^{1.7}))$.*

Proof. Since \mathbf{Y} is good, there exists $\mathbf{W}^* \in \mathcal{W}(\mathbf{Y})$, corresponding to the indicator matrix of the true solution, such that

$$\begin{aligned} g(\mathbf{W}^*) &= \|\mathbf{W}^*\|_{\mathbb{F}}^2 - 10^{10}k^{300}\langle J, \mathbf{W}^* \rangle \leq \frac{n^2}{k} + n^{1.6} - (10)^{10}k^{300} \left(\frac{n^2}{k} - n^{1.6} \right) \\ &= \frac{n^2}{k} \left(1 + \frac{k}{n^{0.4}} - (10)^{10}k^{300} \left(1 - \frac{k}{n^{0.4}} \right) \right). \end{aligned}$$

Since $g(W(\hat{\zeta}(\mathbf{Y}))) \leq g(\mathbf{W}^*)$ it follows that

$$(10)^{10}k^{300}\langle J, W(\hat{\zeta}(\mathbf{Y})) \rangle \geq |g(W(\hat{\zeta}(\mathbf{Y})))| \geq \frac{n^2}{k} \left((10)^{10}k^{300} \left(1 - \frac{k}{n^{0.4}} \right) - 1 - \frac{k}{n^{0.4}} \right).$$

Since, $\|W(\hat{\zeta}(\mathbf{Y}))\|_1 \geq \langle J, W(\hat{\zeta}(\mathbf{Y})) \rangle$ the first claim follows rearranging the terms. This means that the algorithm rejects only if $|\tau| \geq n^{1.7}$. Recall that $\tau \sim \text{tLap}\left(-n^{1.6} \left(1 + \frac{\log(1/\delta)}{\varepsilon}\right), \frac{n^{1.6}}{\varepsilon}\right)$. Hence, by Lemma A.11 it follows that

$$\mathbb{P}(|\tau| \geq n^{1.7}) \leq \frac{\exp(-n^{1.7} + \varepsilon + \log(1/\delta))}{2 - \exp(-\varepsilon - \log(1/\delta))} = \exp(-\Omega(n^{1.7})).$$

□

The next step shows that on a good input \mathbf{Y} the matrix $\phi(W(\hat{\zeta}(\mathbf{Y})))$ is close to the true solution.

Lemma D.15 (Closeness to true solution on good inputs). *Consider the settings of Theorem D.3. Suppose \mathbf{Y} is a good set as per Definition D.6. Let $W(\mathbf{Y}) \in \mathcal{W}(\mathbf{Y})$ be the matrix computed by Algorithm D.5. Suppose the algorithm does not reject. Then*

$$\|\phi(W(\mathbf{Y})) - \mathbf{W}^*\|_1 \leq \frac{n^2}{k} \cdot \frac{3}{k^{98}}.$$

The proof is similar to the classical utility analysis of the sum-of-squares program found, e.g., in [30, 26]. We defer it to Appendix I.

Together, the above results imply that the vectors $\nu^{(i)}$ computed by the algorithm are close to the true centers of the mixture.

Lemma D.16 (Closeness to true centers). *Consider the settings of Theorem D.3. Suppose \mathbf{Y} is a good set as per Definition D.6. Let $\mathbf{W} \in \mathcal{W}(\mathbf{Y})$ be the matrix computed by Algorithm D.5. Suppose the algorithm does not reject in step 3. Then for each $\ell \in [k]$, there exists $\frac{n}{k} \cdot \left(1 - \frac{2}{k^{47}}\right)$ indices $i \in [n]$, such that*

$$\|\nu^{(i)}(\mathbf{W}) - \mu_\ell\|_2 \leq O(k^{-25}).$$

Proof. We aim to show that for most indices $i \in [n]$ the vectors $\|\phi(\mathbf{W}_i)\|_1^{-1} \phi(\mathbf{W}_i)$ and $\|\mathbf{W}_i^*\|_1^{-1} \mathbf{W}_i^*$ induce a 2-explicitly 40-bounded distribution over \mathbf{Y} . If additionally the two vectors are close in ℓ_1 -norm, the result will follow by Theorem A.21.

Note that $\|\mathbf{W}_i^*\|_1^{-1} \mathbf{W}_i^*$ induces a 2-explicitly 40-bounded distribution by Lemma D.7. By Markov's inequality and Lemma D.15 there can be at most n/k^{48} indices $j \in [n]$ such that

$$\|\phi(\mathbf{W})_j - \mathbf{W}_j^*\|_1 \geq \frac{n}{k} \cdot \frac{3}{k^{50}}.$$

Consider all remaining indices i . It follows that

$$\|\phi(\mathbf{W}_i)\|_1 \geq \|\mathbf{W}_i^*\|_1 - \|\phi(\mathbf{W})_i - \mathbf{W}_i^*\|_1 \geq \frac{n}{k} \cdot \left(1 - \frac{k}{n^{0.4}} - \frac{3}{k^{50}}\right) \geq \frac{n}{k} \cdot \left(1 - \frac{10}{k^{50}}\right).$$

Hence, by Lemma D.11 the distribution induced by $\|\phi(\mathbf{W}_i)\|_1^{-1} \phi(\mathbf{W}_i)$ is 2-explicitly 40-bounded distribution. Further, using $\|\mathbf{W}_i^*\|_1 \geq \frac{n}{k} \left(1 - \frac{k}{n^{0.4}}\right)$ we can bound

$$\begin{aligned} & \left\| \|\phi(\mathbf{W}_i)\|_1^{-1} \phi(\mathbf{W}_i) - \|\mathbf{W}_i^*\|_1^{-1} \mathbf{W}_i^* \right\|_1 = \|\phi(\mathbf{W}_i)\|_1^{-1} \|\mathbf{W}_i^*\|_1^{-1} \cdot \left| \|\mathbf{W}_i^*\|_1 \phi(\mathbf{W}_i) - \|\phi(\mathbf{W}_i)\|_1 \mathbf{W}_i^* \right|_1 \\ & \leq \|\phi(\mathbf{W}_i)\|_1^{-1} \|\mathbf{W}_i^*\|_1^{-1} \cdot (\|\phi(\mathbf{W}_i)\|_1 - \|\mathbf{W}_i^*\|_1) \cdot \|\phi(\mathbf{W}_i)\|_1 + \|\phi(\mathbf{W}_i)\|_1 \cdot \|\phi(\mathbf{W}_i) - \mathbf{W}_i^*\|_1 \\ & \leq \|\mathbf{W}_i^*\|_1^{-1} \cdot 2 \|\phi(\mathbf{W}_i) - \mathbf{W}_i^*\|_1 \leq \frac{6}{k^{50} \cdot \left(1 - \frac{k}{n^{0.4}}\right)} \leq \frac{7}{k^{50}}. \end{aligned}$$

Hence, by Theorem A.21 for each $l \in [k]$ there are at least $\frac{n}{k} - n^{0.6} - \frac{n}{k^{48}} \geq \frac{n}{k} \cdot \left(1 - \frac{2}{k^{47}}\right)$ indices i such that

$$\left\| \nu^{(i)}(\mathbf{W}) - \|\mathbf{W}_i^*\|_1^{-1} \sum_{j=1}^n \mathbf{W}_{i,j}^* \mathbf{Y}_j \right\|_2 \leq O(k^{-25}).$$

The result now follows by standard concentration bounds applied to the distribution induced by $\|\mathbf{W}_i^*\|_1^{-1} \mathbf{W}_i^*$. \square

An immediate consequence of Lemma D.16 is that the vectors $\bar{\nu}^{(i)}$ inherits the good properties of the vectors $\nu^{(i)}$ with high probability.

Corollary D.17 (Closeness to true centers after sub-sampling). *Consider the settings of Theorem D.3. Suppose \mathbf{Y} is a good set as per Definition D.6. Let $\mathbf{W} \in \mathcal{W}(\mathbf{Y})$ be the matrix computed by Algorithm D.5. Suppose the algorithm does not reject. Then with high probability for each $\ell \in [k]$, there exists $\frac{n^{0.01}}{k} \cdot \left(1 - \frac{150}{k^{47}}\right)$ indices $i \in \mathcal{S}$, such that*

$$\left\| \bar{\nu}^{(i)} - \mu_\ell \right\|_2 \leq O(k^{-25}).$$

Proof. For each $\ell \in [k]$, denote by \mathcal{T}_ℓ the set of indices in $[n]$ satisfying

$$\left\| \nu^{(i)}(\mathbf{W}) - \mu_\ell \right\|_2 \leq O(k^{-25}).$$

By Lemma D.16 we know that \mathcal{T}_ℓ has size at least $\frac{n}{k} \cdot \left(1 - \frac{2}{k^{47}}\right)$. Further, let \mathcal{S} be the set of indices selected by the algorithm. By Chernoff's bound Fact E.4 with probability $1 - e^{-n^{\Omega(1)}}$, we have $|\mathcal{S} \cap \mathcal{T}_\ell| \geq \frac{n^{0.01}}{k} \cdot \left(1 - \frac{150}{k^{47}}\right)$. Taking a union bound over all $\ell \in [k]$ we get that with probability $1 - e^{-n^{\Omega(1)}}$, for each $\ell \in [k]$, there exists $\frac{n^{0.01}}{k} \cdot \left(1 - \frac{150}{k^{47}}\right)$ indices $i \in \mathcal{S}$ such that

$$\left\| \nu^{(i)}(\mathbf{W}) - \mu_\ell \right\|_2 \leq O(k^{-25}).$$

Now, we obtain the corollary observing (cf. Fact E.1 with $m = 1$) that with probability at least $1 - e^{-n^{\Omega(1)}}$, for all $i \in \mathcal{S}$

$$\left\| \bar{\nu}^{(i)} - \nu^{(i)}(\mathbf{W}) \right\|_2 = \|\mathbf{w}\|_2 \leq n^{-0.05} \cdot \frac{\sqrt{\log(2/\delta)}}{\varepsilon} \cdot \sqrt{d} \leq n^{-0.04} \leq O(k^{-25}).$$

\square

For each ℓ , denote by $\mathcal{G}_\ell \subseteq \mathcal{S}$ the set of indices $i \in \mathcal{S}$ satisfying

$$\left\| \bar{\nu}^{(i)} - \mu_\ell \right\|_2 \leq O(k^{-25}).$$

Let $\mathcal{G} := \bigcup_{\ell \in [k]} \mathcal{G}_\ell$. We now have all the tools to prove utility of Algorithm D.5. We achieve this by

showing that with high probability, each bin returned by the algorithm at step 7 satisfies $\mathcal{G}_{\ell'} \subseteq \mathbf{B}_\ell$ for some $\ell, \ell' \in [k]$. Choosing the bins small enough will yield the desired result.

Lemma D.18 (Closeness of estimates). *Consider the settings of Theorem D.3. Suppose \mathbf{Y} is a good set as per Definition D.6. Let $\mathbf{W} \in \mathcal{W}(\mathbf{Y})$ be the matrix computed by Algorithm D.5. Suppose the algorithm does not reject. Then with high probability, there exists a permutation $\pi : [k] \rightarrow [k]$ such that*

$$\max_{\ell \in [k]} \left\| \mu_\ell - \hat{\mu}_{\pi(\ell)} \right\|_2 \leq O(k^{-20})$$

Proof. Consider distinct $\ell, \ell' \in [k]$. By Corollary D.17 for each $\bar{\nu}^{(i)}, \bar{\nu}^{(j)} \in \mathcal{G}_\ell$ it holds that

$$\left\| \bar{\nu}^{(i)} - \bar{\nu}^{(j)} \right\|_2 \leq C \cdot k^{-25},$$

for some universal constant $C > 0$. Moreover, by assumption on $\mu_\ell, \mu_{\ell'}$ for each $\bar{\nu}^{(i)} \in \mathcal{G}_\ell$ and $\bar{\nu}^{(j)} \in \mathcal{G}_{\ell'}$

$$\left\| \bar{\nu}^{(i)} - \bar{\nu}^{(j)} \right\|_2 \geq \Delta - O(k^{-25}).$$

Thus, by Lemma E.5 with probability at least $1 - e^{-\Omega(d^*)} \geq 1 - n^{-100}$ it holds that for each $\bar{\nu}^{(i)}, \bar{\nu}^{(j)} \in \mathcal{G}_\ell$ and $\bar{\nu}^{(r)} \in \mathcal{G}_{\ell'}$ with $\ell' \neq \ell$,

$$\left\| \Phi \bar{\nu}^{(i)} - \Phi \bar{\nu}^{(j)} \right\|_2 \leq C^* \cdot k^{-25} \quad \text{and} \quad \left\| \Phi \bar{\nu}^{(i)} - \Phi \bar{\nu}^{(r)} \right\|_2 \geq \Delta - C^* \cdot k^{-25}$$

for some other universal constant $C^* > C$. Let $Q_\Phi(\mathcal{G}_\ell) \subseteq \mathbb{R}^{d^*}$ be a ball of radius $C^* \cdot (k^{-25})$ such that $\forall i \in \mathcal{G}_\ell$ it holds $\Phi \bar{\nu}^{(i)} \in Q_\Phi(\mathcal{G}_\ell)$. That is, $Q_\Phi(\mathcal{G}_\ell)$ contains the projection of all points in \mathcal{G}_ℓ .

Recall that $d^* = 100 \log(n) \leq 100k^5$ and $b = k^{-15}$. Let $\mathcal{B} = \{\mathbf{B}_i\}_{i=1}^\infty$ be the sequence of bins computed by the histogram learner of Lemma A.13 for \mathbb{R}^{d^*} at step 7 of the algorithm. By choice of b , and since \mathbf{q} is chosen uniformly at random in $[0, b]$, the probability that there exists a bin $\mathbf{B} \in \mathcal{B}$ containing $Q_\Phi(\mathcal{G}_\ell)$ is at least

$$1 - d^* \cdot \frac{C^*}{b} \cdot (k^{-25}) \geq 1 - \frac{100C^*}{b} \cdot k^{-20} \geq 1 - O(k^{-5}),$$

where we used that $d^* = 100 \log n \leq 100k^5$. A simple union bound over $\ell \in [k]$ yields that with high probability for all $\ell \in [k]$, there exists $\mathbf{B} \in \mathcal{B}$ such that $Q_\Phi(\mathcal{G}_\ell) \subseteq \mathbf{B}$. For simplicity, denote such bin by \mathbf{B}_ℓ .

We continue our analysis conditioning on the above events, happening with high probability. First, notice that for all $l \in [k]$

$$\max_{u, u' \in \mathbf{B}_\ell} \|u - u'\|_2^2 \leq d^* \cdot b^2 \leq 100k^{-25} \leq \frac{\Delta - C^*k^{-25}}{k^{10}},$$

and thus there cannot be $\ell, \ell' \in [k]$ such that $Q_\Phi(\mathcal{G}_\ell) \subseteq \mathbf{B}_\ell$ and $Q_\Phi(\mathcal{G}_{\ell'}) \subseteq \mathbf{B}_\ell$. Moreover, by Corollary D.17 and

$$\min_{\ell \in [k]} |\mathcal{G}_\ell| \geq \frac{n^{0.01}}{k} \cdot \left(1 - \frac{150}{k^{47}}\right),$$

and hence

$$|\mathcal{S} \setminus \mathcal{G}| \leq n^{0.01} \cdot \frac{150}{k^{47}} = \frac{n^{0.01}}{k} \cdot \frac{150}{k^{46}}$$

it must be that step 7 returned bins $\mathbf{B}_1, \dots, \mathbf{B}_k$. This also implies that the algorithm does not reject. Further, by Lemma E.5 for all $\bar{\nu}^{(i)}, \bar{\nu}^{(j)}$ such that $\Phi \bar{\nu}^{(i)}, \Phi \bar{\nu}^{(j)} \in \mathbf{B}_l$ it holds that

$$\left\| \bar{\nu}^{(i)} - \bar{\nu}^{(j)} \right\|_2 \leq C^* \cdot \left\| \Phi \bar{\nu}^{(i)} - \Phi \bar{\nu}^{(j)} \right\|_2 \leq C^* \cdot \sqrt{d^*} \cdot b \leq O(k^{-12}).$$

And hence, by triangle inequality, we get

$$\|\bar{\nu}^{(i)} - \mu_l\|_2 \leq O(k^{-12}).$$

Finally, recall that for each $\ell \in [k]$,

$$\hat{\mu}_l := \frac{1}{|\{j \mid \Phi \bar{\nu}^{(j)} \in \mathbf{B}_l\}|} \cdot \left(\sum_{\Phi \bar{\nu}^{(j)} \in \mathbf{B}_l} \bar{\nu}^{(j)} \right) + \mathbf{w}',$$

where $\mathbf{w}' \sim N\left(0, N\left(0, 32 \cdot k^{-120} \cdot \frac{\log(2kn/\delta)}{\varepsilon^2} \cdot \text{Id}\right)\right)$. Since by choice of n, k, ε it holds that

$$32 \cdot k^{-120} \cdot \frac{\log(2kn/\delta)}{\varepsilon^2} \leq O(k^{-90}),$$

we get with probability at least $1 - e^{-k^{\Omega(1)}}$ for each $\ell \in [k]$, by Fact E.1, with $m = 1$, and a union bound that

$$\|\mathbf{w}'\| \leq O(k^{-20}).$$

Since all $\bar{\nu}^{(i)}$ such that $\Phi \bar{\nu}^{(i)} \in \mathbf{B}_l$ are at most $O(k^{-12})$ -far from μ_l , also their average is. We conclude that

$$\|\hat{\mu}_l - \mu_l\|_2 \leq O(k^{-12}) + \|\mathbf{w}'\|_2 \leq O(k^{-12}).$$

This completes the proof. \square

Now Theorem D.3 is a trivial consequence.

Proof of Theorem D.3. The error guarantees and privacy guarantees immediately follows combining Lemma D.8, Lemma D.15, Lemma D.14 and Lemma D.18. The running time follows by Fact A.16. \square

E Concentration inequalities

We introduce here several useful and standard concentration inequalities.

Fact E.1 (Concentration of spectral norm of Gaussian matrices). *Let $\mathbf{W} \sim \mathcal{N}(0, 1)^{m \times n}$. Then for any t , we have*

$$\mathbb{P}(\sqrt{m} - \sqrt{n} - t \leq \sigma_{\min}(\mathbf{W}) \leq \sigma_{\max}(\mathbf{W}) \leq \sqrt{m} + \sqrt{n} + t) \geq 1 - 2 \exp\left(-\frac{t^2}{2}\right),$$

where $\sigma_{\min}(\cdot)$ and $\sigma_{\max}(\cdot)$ denote the minimum and the maximum singular values of a matrix, respectively.

Let \mathbf{W}' be an n -by- n symmetric matrix with independent entries sampled from $N(0, \sigma^2)$. Then $\|\mathbf{W}'\| \leq 3\sigma\sqrt{n}$ with probability at least $1 - \exp(-\Omega(n))$.

Fact E.2 (Maximum degree of Erdős-Rényi graphs). *Let G be an Erdős-Rényi graph on n vertices with edge probability p . Then with probability at least $1 - n \exp(-np/3)$, any vertex in G has degree at most $2np$.*

Fact E.3 (Gaussian concentration bounds). *Let $\mathbf{X} \sim \mathcal{N}(0, \sigma^2)$. Then for any $t \geq 0$,*

$$\max\{\mathbb{P}(\mathbf{X} \geq t), \mathbb{P}(\mathbf{X} \leq -t)\} \leq \exp\left(-\frac{t^2}{2\sigma^2}\right).$$

Fact E.4 (Chernoff bound). *Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be independent random variables taking values in $\{0, 1\}$. Let $\mathbf{X} := \sum_{i=1}^n \mathbf{X}_i$ and let $\mu := \mathbb{E} \mathbf{X}$. Then for any $\delta > 0$,*

$$\mathbb{P}(\mathbf{X} \leq (1 - \delta)\mu) \leq \exp\left(-\frac{\delta^2 \mu}{2}\right),$$

$$\mathbb{P}(\mathbf{X} \geq (1 + \delta)\mu) \leq \exp\left(-\frac{\delta^2 \mu}{2 + \delta}\right).$$

Lemma E.5 ([31]). *Let Φ be a d -by- n Gaussian matrix, with each entry independently chosen from $N(0, 1/d)$. Then, for every vector $u \in \mathbb{R}^n$ and every $\alpha \in (0, 1)$*

$$\mathbb{P}(\|\Phi u\| = (1 \pm \alpha)\|u\|) \geq 1 - e^{-\Omega(\alpha^2 d)}.$$

F Linear algebra

Lemma F.1 (Weyl's inequality). *Let A and B be symmetric matrices. Let $R = A - B$. Let $\alpha_1 \geq \dots \geq \alpha_n$ be the eigenvalues of A . Let $\beta_1 \geq \dots \geq \beta_n$ be the eigenvalues of B . Then for each $i \in [n]$,*

$$|\alpha_i - \beta_i| \leq \|R\|.$$

Lemma F.2 (Davis-Kahan's theorem). *Let A and B be symmetric matrices. Let $R = A - B$. Let $\alpha_1 \geq \dots \geq \alpha_n$ be the eigenvalues of A with corresponding eigenvectors v_1, \dots, v_n . Let $\beta_1 \geq \dots \geq \beta_n$ be the eigenvalues of B with corresponding eigenvectors u_1, \dots, u_n . Let θ_i be the angle between $\pm v_i$ and $\pm u_i$. Then for each $i \in [n]$,*

$$\sin(2\theta_i) \leq \frac{2\|R\|}{\min_{j \neq i} |\alpha_i - \alpha_j|}.$$

G Convex optimization

Proposition G.1. *Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$ be a convex function. Let $\mathcal{K} \subseteq \mathbb{R}^m$ be a convex set. Then $y^* \in \mathcal{K}$ is a minimizer of f over \mathcal{K} if and only if there exists a subgradient $g \in \partial f(y^*)$ such that*

$$\langle y - y^*, g \rangle \geq 0 \quad \forall y \in \mathcal{K}.$$

Proof. Define indicator function

$$I_{\mathcal{K}}(y) = \begin{cases} 0, & y \in \mathcal{K}, \\ \infty, & y \notin \mathcal{K}. \end{cases}$$

Then for $y \in \mathcal{K}$, one has

$$\partial I_{\mathcal{K}}(y) = \{g \in \mathbb{R}^m : \langle g, y - y' \rangle \geq 0 \forall y' \in \mathcal{K}\}.$$

Note y^* is a minimizer of f over \mathcal{K} , if and only if y^* is a minimizer of $f + I_{\mathcal{K}}$ over \mathbb{R}^m , if and only if $\mathbf{0}_m \in \partial(f + I_{\mathcal{K}})(y^*) = \partial f(y^*) + \partial I_{\mathcal{K}}(y^*)$, if and only if there exists $g \in \partial f(y^*)$ such that $\langle g, y - y^* \rangle \geq 0$ for any $y \in \mathcal{K}$. \square

Proposition G.2 (Pythagorean theorem from strong convexity). *Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$ be a convex function. Let $\mathcal{K} \subseteq \mathbb{R}^m$ be a convex set. Suppose f is κ -strongly convex over \mathcal{K} . Let $x^* \in \mathcal{K}$ be a minimizer of f over \mathcal{K} . Then for any $x \in \mathcal{K}$, one has*

$$\|x - x^*\|^2 \leq \frac{2}{\kappa}(f(x) - f(x^*)).$$

Proof. By strong convexity, for any subgradient $g \in \partial f(x^*)$ one has

$$f(x) \geq f(x^*) + \langle x - x^*, g \rangle + \frac{\kappa}{2} \|x - x^*\|^2.$$

By Proposition G.1, $\langle x - x^*, g \rangle \geq 0$ for some $g \in \partial f(x^*)$. Then the result follows. \square

H Deferred proofs SBM

We prove Lemma C.9 restated below.

Lemma H.1 (Restatement of Lemma C.9). *Consider the settings of Lemma C.8. With probability $1 - \exp(-\Omega(n))$ over $\mathbf{G} \sim \text{SBM}_n(\gamma, d, x)$,*

$$\left\| \hat{X}(Y(\mathbf{G})) - \frac{1}{n}xx^\top \right\|_F^2 \leq \frac{800}{\gamma\sqrt{d}}.$$

Proof. Recall $\mathcal{K} = \{X \in \mathbb{R}^{n \times n} : X \succeq 0, X_{ii} = 1/n \forall i\}$. Let $X^* := \frac{1}{n}xx^\top$. Since $\hat{X} = \hat{X}(Y(\mathbf{G}))$ is a minimizer of $\min_{X \in \mathcal{K}} \|Y(\mathbf{G}) - X\|_F^2$ and $X^* \in \mathcal{K}$, we have

$$\left\| \hat{X} - Y(\mathbf{G}) \right\|_F^2 \leq \|X^* - Y(\mathbf{G})\|_F^2 \iff \left\| \hat{X} - X^* \right\|_F^2 \leq 2 \left\langle \hat{X} - X^*, Y(\mathbf{G}) - X^* \right\rangle.$$

The infinity-to-one norm of a matrix $M \in \mathbb{R}^{m \times n}$ is defined as

$$\|M\|_{\infty \rightarrow 1} := \max \{ \langle u, Mv \rangle : u \in \{\pm 1\}^m, v \in \{\pm 1\}^n \}.$$

By [28, Fact 3.2], every $Z \in \mathcal{K}$ satisfies

$$|\langle Z, Y(\mathbf{G}) - X^* \rangle| \leq \frac{K_G}{n} \cdot \|Y(\mathbf{G}) - X^*\|_{\infty \rightarrow 1},$$

where $K_G \leq 1.783$ is Grothendieck's constant. Similar to the proof of [28, Lemma 4.1], using Bernstein's inequality and union bound, we can show (cf. Fact H.2)

$$\|Y(\mathbf{G}) - X^*\|_{\infty \rightarrow 1} \leq \frac{100n}{\gamma\sqrt{d}}$$

with probability $1 - \exp(-\Omega(n))$. Putting things together, we have

$$\left\| \hat{X}(Y(\mathbf{G})) - \frac{1}{n}xx^\top \right\|_F^2 \leq \frac{400 \cdot K_G}{\gamma\sqrt{d}},$$

with probability $1 - \exp(-\Omega(n))$. □

Fact H.2. Let $\gamma > 0, d \in \mathbb{N}, x^* \in \{\pm 1\}^n$, and $\mathbf{G} \sim \text{SBM}(\gamma, d, x^*)$. Let $Y(\mathbf{G}) = \frac{1}{\gamma d} (A(\mathbf{G}) - \frac{d}{n}J)$, where $A(\mathbf{G})$ is the adjacency matrix of (G) with entries d/n on the diagonal. Then

$$\max_{x \in \{\pm 1\}^n} |x^\top (Y(\mathbf{G}) - \frac{1}{n}x^*(x^*)^\top) x| \leq \frac{100n}{\gamma\sqrt{d}}$$

with probability at least $1 - e^{-10n}$.

Proof. The result will follow using Bernstein's Inequality and a union bound. Define $\mathbf{E} := Y(\mathbf{G}) - \frac{1}{n}x^*(x^*)^\top$. Fix $x \in \{\pm 1\}^n$ and for $1 \leq i < j \leq n$, let $\mathbf{Z}_{i,j} := \mathbf{E}_{i,j}x_i x_j$. Then $x^\top \mathbf{E}x = 2 \sum_{1 \leq i < j \leq n} \mathbf{Z}_{i,j}$. Note that

$$\begin{aligned} \mathbb{E} \mathbf{Z}_{i,j} &= 0, \\ |\mathbf{Z}_{i,j}| &\leq \frac{1}{\gamma n} \cdot \left(\frac{n}{d} - 1\right) + \frac{1}{\gamma dn} \leq \frac{1}{\gamma d}, \\ \mathbb{E} \mathbf{Z}_{i,j}^2 &= \text{Var} [\mathbf{Y}(\mathbf{G})_{i,j}] \leq \mathbb{E} \mathbf{Y}(\mathbf{G})_{i,j}^2 \leq (1 + \gamma) \frac{d}{n} \cdot \frac{1}{\gamma^2 n^2} \left[\left(\frac{n}{d} - 1\right)^2 - \frac{1}{\gamma^2 n^2} \right] + \frac{1}{\gamma^2 n^2} \\ &\leq (1 + \gamma) \frac{1}{d\gamma^2 n} + \frac{1}{\gamma^2 n^2} \leq \frac{3}{\gamma^2 dn}. \end{aligned}$$

By Bernstein's Inequality (cf. [62, Proposition 2.14]) it follows that

$$\begin{aligned} \mathbb{P} \left(\sum_{i < j} \mathbf{Z}_{i,j} \geq \frac{50n}{\gamma\sqrt{d}} \right) &\leq \mathbb{P} \left(\sum_{i < j} \mathbf{Z}_{i,j} \geq \frac{n^2}{2} \cdot \frac{100n}{\gamma\sqrt{d}} \right) \leq 2 \exp \left(-\frac{\frac{10^4}{\gamma^2 d}}{\frac{3}{\gamma^2 dn} + \frac{100}{3\gamma^2 d^{3/2}n}} \right) \\ &= 2 \exp \left(-\frac{10^4 n}{3 + \frac{100}{\sqrt{d}}} \right) \leq \exp(-50n). \end{aligned}$$

Hence, by a union bound over all $x \in \{\pm 1\}^n$ it follows that

$$\max_{x \in \{\pm 1\}^n} |x^\top (Y(\mathbf{G}) - \frac{1}{n}x^*(x^*)^\top) x| \leq \frac{100n}{\gamma\sqrt{d}}$$

with probability at least $1 - e^{-10n}$. □

I Deferred proofs for clustering

In this section, we will prove Lemma D.15 restated below.

Lemma (Restatement of Lemma D.15). *Consider the settings of Theorem D.3. Suppose \mathbf{Y} is a good set as per Definition D.6. Let $W(\mathbf{Y}) \in \mathcal{W}(\mathbf{Y})$ be the matrix computed by Algorithm D.5. Suppose the algorithm does not reject. Then*

$$\|\phi(W(\mathbf{Y})) - \mathbf{W}^*\|_1 \leq \frac{n^2}{k} \cdot \frac{3}{k^{98}}.$$

We will need the following fact about our clustering program. Similar facts were used, e.g., in [30, 26]. One difference for us is that we don't have a constraint on the lower bound on the cluster size indicated by our SOS variables. However, since we maximize a variant of the ℓ_1 norm of the second moment matrix of the pseudo-distribution this will make up for this.

Fact I.1. *Consider the same setting as in Lemma D.15. Let $0 < \delta \leq \frac{1}{1.5 \cdot 10^{10}} \cdot \frac{1}{k^{201}}$ and denote by $\mathbf{C}_1, \dots, \mathbf{C}_k \subseteq [n]$ the indices belonging to each true cluster. Then $W(\mathbf{Y})$ satisfies the following three properties:*

1. For all $i, j \in [n]$ it holds that $0 \leq \mathbf{W}_{i,j} \leq 1$,
2. for all $i \in [n]$ it holds that $\sum_{j=1}^n \mathbf{W}_{i,j} \leq \frac{n}{k}$ and for at least $(1 - \frac{1}{1000k^{100}})n$ indices $i \in [n]$ it holds that $\sum_{j=1}^n \mathbf{W}_{i,j} \geq (1 - \frac{1}{(10)^6 k^{200}}) \cdot \frac{n}{k}$,
3. for all $r \in [k]$ it holds that $\sum_{i \in \mathbf{C}_r, j \notin \mathbf{C}_r} \mathbf{W}_{i,j} \leq \delta \cdot \frac{n^2}{k}$.

We will prove Fact I.1 at the end of this section. With this in hand, we can proof Lemma D.15.

Proof of Lemma D.15. For brevity, we write $\mathbf{W} = W(\mathbf{Y})$. Since $\phi(\mathbf{W}^*) = \mathbf{W}^*$ and ϕ is 10-Lipschitz we can also bound

$$\|\phi(\mathbf{W}) - \mathbf{W}^*\|_1 \leq 10 \cdot \|\mathbf{W} - \mathbf{W}^*\|_1.$$

Let $\delta \leq \frac{1}{1.5 \cdot 10^{10}} \cdot \frac{1}{k^{201}}$ and again let $\mathbf{C}_1, \dots, \mathbf{C}_k \subseteq [n]$ denote the indices belonging to each true cluster. Note that by assumption that \mathbf{Y} is a good sample it holds for each $r \in [k]$ that $\frac{n}{k} - n^{0.6} \leq |\mathbf{C}_r| \leq \frac{n}{k} + n^{0.6}$.

Let $r, r' \in [k]$. We can write

$$\|\mathbf{W} - \mathbf{W}^*\|_1 = \sum_{r=1}^k \sum_{i,j \in \mathbf{C}_r} |\mathbf{W}_{i,j} - 1| + \sum_{r=1}^k \sum_{i \in \mathbf{C}_r, j \notin \mathbf{C}_r} |\mathbf{W}_{i,j} - 0| \quad (\text{I.1})$$

Note that we can bound the second sum by $k \cdot \delta \frac{n^2}{k}$ using Item 3. Further, in what follows consider only indices i such that $\sum_{j=1}^n \mathbf{W}_{i,j} \geq (1 - \frac{1}{(10)^6 k^{200}}) \cdot \frac{n}{k}$. By Item 2 we can bound the contribution of the other indices by

$$\frac{1}{1000k^{100}} n \cdot \left(\frac{n}{k} + n^{0.6} \right) \leq \frac{2}{1000k^{100}} \cdot \frac{n^2}{k}.$$

Focusing only on such indices, for the first sum in Eq. (I.1), fix $r \in [k]$. We will aim to show that most entries of \mathbf{W} are large if and only if the corresponding entry of \mathbf{W}^* is 1. By Item 3 and Markov's Inequality, it follows that for at least a $(1 - \frac{1}{1000k^{100}})$ -fraction of the indices $i \in \mathbf{C}_r$ it holds that

$$\sum_{j \notin \mathbf{C}_r} \mathbf{W}_{i,j} \leq 1000k^{100} \cdot \delta \frac{n^2}{k \cdot |\mathbf{C}_r|} \leq 1000k^{100} \delta \cdot \frac{n}{1 - k \cdot n^{-0.4}} \leq 2000k^{101} \delta \cdot \frac{n}{k},$$

where we used that $|\mathbf{C}_r| \geq \frac{n}{k} - n^{0.6}$. Call such indices *good*. Notice that for good indices it follows using Item 2 that

$$\sum_{j \in \mathbf{C}_r} \mathbf{W}_{i,j} \geq \frac{n}{k} \cdot \left(1 - \frac{1}{(10)^6 k^{200}} - 2000k^{101} \delta \right).$$

Denote by G the number of $j \in \mathbf{C}_r$ such that $\mathbf{W}_{i,j} \geq 1 - \frac{1}{1000k^{100}}$. Using the previous display and that $\mathbf{W}_{i,j} \leq 1$ we obtain

$$\begin{aligned} \frac{n}{k} \cdot \left(1 - \frac{1}{(10)^6 k^{200}} - 2000k^{101}\delta\right) &\leq \sum_{j \in \mathbf{C}_r} \mathbf{W}_{i,j} \leq G \cdot 1 + (|\mathbf{C}_r| - G) \cdot \left(1 - \frac{1}{1000k^{100}}\right) \\ &\leq G \cdot \frac{1}{1000k^{100}} + \frac{n}{k} \cdot \left(1 + \frac{1}{kn^{0.4}}\right) \cdot \left(1 - \frac{1}{1000k^{100}}\right) \\ &\leq G \cdot \frac{1}{1000k^{100}} + \frac{n}{k} \cdot \left(1 + \frac{1}{kn^{0.4}}\right), \end{aligned}$$

where we also used $|\mathbf{C}_r| \leq \frac{n}{k} + n^{0.6}$. Rearranging now yields

$$G \geq \frac{n}{k} \cdot \left(1 - \frac{1}{1000k^{100}} - \frac{10^3 k^{99}}{n^{0.4}} - 2 \cdot 10^6 k^{101}\delta\right) \geq \frac{n}{k} \cdot \left(1 - \frac{2}{1000k^{100}} - 2 \cdot 10^6 k^{101}\delta\right).$$

We can now bound

$$\begin{aligned} \sum_{i,j \in \mathbf{C}_r} |\mathbf{W}_{i,j} - 1| &= \sum_{i,j \in \mathbf{C}_r, i \text{ is good}} |\mathbf{W}_{i,j} - 1| + \sum_{i,j \in \mathbf{C}_r, i \text{ is not good}} |\mathbf{W}_{i,j} - 1| \\ &\leq |\mathbf{C}_r| \cdot \left((|\mathbf{C}_r| - G) \cdot 1 + |\mathbf{C}_r| \cdot \frac{1}{1000k^{100}}\right) + \frac{1}{1000k^{100}} \cdot |\mathbf{C}_r|^2 \\ &\leq |\mathbf{C}_r|^2 \left(1 + \frac{1}{500k^{100}}\right) - G \cdot |\mathbf{C}_r| \\ &\leq \frac{n^2}{k^2} \left(1 + \frac{k}{n^{0.4}}\right)^2 \left(1 + \frac{1}{500k^{100}}\right) - \frac{n^2}{k^2} \left(1 - \frac{2}{1000k^{100}} - 2 \cdot 10^6 k^{101}\delta\right) \left(1 - \frac{k}{n^{0.4}}\right) \\ &\leq \frac{n^2}{k^2} \cdot \left(30 \cdot 10^6 k^{101}\delta + \frac{11}{500k^{100}}\right) \leq \frac{n^2}{k} \cdot \left(30 \cdot 10^6 k^{100}\delta + \frac{11}{500k^{101}}\right) \\ &\leq \frac{n^2}{k} \cdot \frac{3}{125k^{101}}. \end{aligned}$$

Putting everything together, it follows that

$$\|\phi(\mathbf{W}) - \mathbf{W}^*\|_{\mathbb{F}}^2 \leq \|\phi(\mathbf{W}) - \mathbf{W}^*\|_1 \leq 10 \cdot \frac{n^2}{k} \left(\delta k + \frac{2}{1000k^{100}} + \frac{3}{125k^{100}}\right) \leq \frac{n^2}{k} \cdot \frac{4}{k^{100}} \leq \frac{n^2}{k} \cdot \frac{3}{k^{98}}.$$

□

It remains to verify Fact I.1.

Proof of Fact I.1. Let $\mathcal{P} = \mathcal{P}_{n,k,t}(\mathbf{Y})$ be the system of Eq. $(\mathcal{P}_{n,k,t}(Y))$. Recall that $\mathbf{W}_{i,j} = \tilde{\mathbb{E}} \sum_{l \in [k]} z_{i,l} z_{j,l}$. Since

$$\mathcal{P} \Big|_4 \left\{ 0 \leq \sum_{l \in [k]} z_{i,l} z_{j,l} \leq \sum_{l \in [k]} z_{i,l} \leq 1 \right\},$$

it follows that $0 \leq \mathbf{W}_{i,j} \leq 1$. Further, for each $i \in [n]$ it holds that

$$\mathcal{P} \Big|_4 \left\{ \sum_{j \in [n], l \in [k]} z_{j,l} z_{i,l} \leq \frac{n}{k} \sum_{l \in [k]} z_{i,l} \leq \frac{n}{k} \right\}$$

implying that $\sum_{j \in [n]} \mathbf{W}_{i,j} \leq \frac{n}{k}$. Further, by Lemma D.14

$$\|\mathbf{W}\|_1 \geq \frac{n^2}{k} \cdot \left(1 - n^{-0.4} - \frac{1}{(10)^{10} k^{300}}\right) \geq \frac{n^2}{k} \cdot \left(1 - \frac{1}{(10)^9 k^{300}}\right).$$

Denote by \mathbf{W}_i the i -th row of \mathbf{W} and by L the number of rows which have ℓ_1 norm at least $\left(1 - \frac{1}{(10)^6 k^{200}}\right) \cdot \frac{n}{k}$. Since for all i it holds that $\|\mathbf{W}_i\|_1 \leq \frac{n}{k}$ it follows that

$$\begin{aligned} \frac{n^2}{k} \cdot \left(1 - \frac{1}{(10)^9 k^{300}}\right) &\leq \sum_{i \in [n]} \|\mathbf{W}_i\|_1 \leq L \cdot \frac{n}{k} + (n - L) \cdot \left(1 - \frac{1}{(10)^6 k^{200}}\right) \cdot \frac{n}{k} \\ &= L \cdot \frac{1}{(10)^6 k^{200}} \cdot \frac{n}{k} + \frac{n^2}{k} \cdot \left(1 - \frac{1}{(10)^6 k^{200}}\right) \end{aligned}$$

Rearranging then yields $L \geq (1 - \frac{1}{1000k^{100}}) \cdot n$ which proves Item 2.

It remains to verify Item 3. Fix $r, l \in [k]$ and define $z_l(\mathbf{C}_r) = \frac{k}{n} \sum_{i \in \mathbf{C}_r} z_{i,l}$. Let $t > 0$ be an integer. We aim to show that for all unit vectors v it holds that

$$\mathcal{P} \Big|_{10t} \left\{ z_l(\mathbf{C}_r) \cdot \frac{1}{\Delta^{2t}} \sum_{r' \neq r} z_l(\mathbf{C}_{r'}) \langle \mu_r - \mu_{r'}, v \rangle^{2t} \leq \frac{\delta}{k} \right\}, \quad (I.2)$$

where Δ is the minimal separation between the true means. Before proving this, let us examine how we can use this fact to prove Item 3. Note, that for all $r \neq r'$ it holds that

$$\sum_{s,u \in [k]} \left\langle \mu_r - \mu_{r'}, \frac{\mu_s - \mu_u}{\|\mu_s - \mu_u\|} \right\rangle^{2t} \geq \Delta^{2t}.$$

Hence, if the above SOS proof indeed exists, we obtain

$$\begin{aligned} \sum_{i \in \mathbf{C}_r, j \notin \mathbf{C}_r} \mathbf{W}_{i,j} &= \sum_{l=1}^k \tilde{\mathbb{E}} \sum_{i \in \mathbf{C}_r, j \notin \mathbf{C}_r} z_{i,l} z_{j,l} = \frac{n^2}{k^2} \tilde{\mathbb{E}} z_l(\mathbf{C}_r) \cdot \sum_{r' \neq r} z_l(\mathbf{C}_{r'}) \\ &\leq \frac{n^2}{\Delta^{2t} k^2} \sum_{s,u \in [k]} \tilde{\mathbb{E}} z_l(\mathbf{C}_r) \cdot \sum_{r' \neq r} z_l(\mathbf{C}_r) \left\langle \mu_r - \mu_{r'}, \frac{\mu_s - \mu_u}{\|\mu_s - \mu_u\|} \right\rangle^{2t} \\ &\leq \frac{\delta}{k} k^2 \cdot \frac{n^2}{k^2} = \delta \cdot \frac{n^2}{k}. \end{aligned}$$

In the remainder of this proof we will prove Eq. (I.2). We will use the following SOS version of the triangle inequality (cf. Fact I.2)

$$\left| \frac{x,y}{2t} (x+y)^{2t} \leq 2^{2t-1} (x^{2t} + y^{2t}). \right.$$

Recall that $\mu'_l = \frac{k}{n} \sum_{i=1}^n z_{i,l} y_i$ and denote by $\mu_{\pi(i)}$ the true mean corresponding to the i -th sample. Let v be an arbitrary unit vector, it follows that

$$\begin{aligned} \mathcal{P} \Big|_{10t} \left\{ z_l(\mathbf{C}_r) \cdot \frac{1}{\Delta^{2t}} \sum_{r' \neq r} z_l(\mathbf{C}_{r'}) \langle \mu_r - \mu_{r'}, v \rangle^{2t} \right. \\ \leq z_l(\mathbf{C}_r) \cdot \frac{2^{2t-1}}{\Delta^{2t}} \sum_{r' \neq r} z_l(\mathbf{C}_{r'}) \left(\langle \mu_r - \mu'_l, v \rangle^{2t} + \langle \mu_{r'} - \mu'_l, v \rangle^{2t} \right) \\ \leq \frac{2^{2t-1}}{\Delta^{2t}} \sum_{r=1}^k z_l(\mathbf{C}_r) \langle \mu_r - \mu'_l, v \rangle^{2t} = \frac{2^{2t-1}}{\Delta^{2t}} \cdot \frac{k}{n} \sum_{i=1}^n z_{i,l} \langle \mu_{\pi(i)} - \mu'_l, v \rangle^{2t} \left. \right\}, \end{aligned}$$

where we used that $\mathcal{P} \Big|_{\Gamma} \sum_{r=1}^k z_l(\mathbf{C}_r) \leq 1$. Using the SOS triangle inequality again and that $\mathcal{P} \Big|_{\frac{\Gamma}{2}} z_{i,l} \leq 1$ we obtain

$$\begin{aligned} \mathcal{P} \Big|_{10t} \left\{ z_l(\mathbf{C}_r) \cdot \frac{1}{\Delta^{2t}} \sum_{r' \neq r} z_l(\mathbf{C}_{r'}) \langle \mu_r - \mu_{r'}, v \rangle^{2t} \right. \\ \leq \frac{2^{4t-1}}{\Delta^{2t}} \cdot \left(k \cdot \frac{1}{n} \sum_{i=1}^n \langle \mathbf{y}_i - \mu_{\pi(i)}, v \rangle^{2t} + \frac{k}{n} \sum_{i=1}^n z_{i,l} \langle \mathbf{y}_i - \mu'_l, v \rangle^{2t} \right) \left. \right\}. \end{aligned}$$

We start by bounding the first sum. Recall that by assumption the uniform distribution over each true cluster is $2t$ -explicitly 2-bounded. It follows that

$$\left| \frac{1}{2t} \left\{ \frac{1}{n} \sum_{i=1}^n \langle \mathbf{y}_i - \mu_{\pi(i)}, v \rangle^{2t} \right\} = \frac{1}{k} \sum_{r=1}^k \frac{k}{n} \sum_{i \in \mathbf{C}_r} \langle \mathbf{y}_i - \mu_r, v \rangle^{2t} \leq \frac{1}{k} \sum_{r=1}^k \frac{k}{n} \cdot |\mathbf{C}_r| \cdot (2t)^t \cdot \|v\|_2^{2t} \right. \quad (I.3)$$

$$\left. \leq \left(1 + \frac{k}{n^{0.4}} \right) \cdot (2t)^t \leq 2(2t)^t \right\}, \quad (I.4)$$

where we used that $|\mathbf{C}_r| \leq \frac{n}{k} + n^{0.6}$. To bound the second sum, we will use the moment bound constraints. In particular, we know that

$$\mathcal{P} \Big|_{10t} \left\{ \frac{k}{n} \sum_{i=1}^n z_{i,l} \langle \mathbf{y}_i - \mu'_l, v \rangle^{2t} \leq (2t)^t \right\}. \quad (\text{I.5})$$

Combining Eq. (I.4) and Eq. (I.5) now yields

$$\mathcal{P} \Big|_{10t} \left\{ z_l(\mathbf{C}_r) \cdot \frac{1}{\Delta^{2t}} \sum_{r' \neq r} z_l(\mathbf{C}_{r'}) \langle \mu_r - \mu_{r'}, v \rangle^{2t} \leq k \frac{2^{2t+1} (2t)^t}{\Delta^{2t}} \leq k \left(\frac{8t}{\Delta^2} \right)^t \right\}.$$

Note that by assumption $\Delta \geq O(\sqrt{tk^{1/t}})$. Overloading notation, we can choose the t parameter in the SOS proof to be 202 times the t parameter in the lower bound in the separation to obtain³⁵

$$\sum_{i \in \mathbf{C}_r, j \notin \mathbf{C}_r} \mathbf{W}_{i,j} \leq \delta \cdot \frac{n^2}{k}.$$

□

I.1 Small Lemmas

Fact I.2 (Lemma A.2 in [36]). *For all integers $t > 0$ it holds that*

$$\left| \frac{x,y}{2t} (x+y)^{2t} \leq 2^{2t-1} (x^{2t} + y^{2t}) \right|.$$

Fact I.3. *Let $\varepsilon, \delta > 0$. Let $\mathcal{M}: \mathcal{Y} \rightarrow \mathcal{O}$ be a randomized algorithm that, for every pair of adjacent inputs, with probability at least $1 - \gamma \geq 1/2$ over the internal randomness of \mathcal{Y} ³⁶ satisfies (ε, δ) -privacy. Then \mathcal{M} is $(\varepsilon + 2\gamma, \delta + \gamma)$ -private.*

Proof. Let X, X' be adjacent input and let B be the event under which \mathcal{M} is (ε, δ) -private. By assumption, we know that $\mathbb{P}(B) \geq 1 - \gamma$. Let $S \in \mathcal{O}$, it follows that

$$\begin{aligned} \mathbb{P}(\mathcal{M}(X) \in S) &= \mathbb{P}(B) \cdot \mathbb{P}(\mathcal{M}(X) \in S \mid B) + \mathbb{P}(B^c) \cdot \mathbb{P}(\mathcal{M}(X) \in S \mid B^c) \\ &\leq \mathbb{P}(\mathcal{M}(X) \in S \mid B) + \gamma \\ &\leq e^\varepsilon \mathbb{P}(\mathcal{M}(X) \in S \mid B) + \delta + \gamma \\ &\leq \frac{e^\varepsilon}{\mathbb{P}(B)} \cdot \mathbb{P}(\mathcal{M}(X) \in S) + \delta + \gamma \\ &\leq e^{\varepsilon + \log\left(\frac{1}{1-\gamma}\right)} \cdot \mathbb{P}(\mathcal{M}(X) \in S) + (\delta + \gamma) \\ &\leq e^{\varepsilon + 2\gamma} \cdot \mathbb{P}(\mathcal{M}(X) \in S) + (\delta + \gamma), \end{aligned}$$

where we used that $\log(1 - \gamma) \geq -2\gamma$ for $\gamma \in [0, 1/2]$. □

I.2 Privatizing input using the Gaussian Mechanism

In this section, we will prove the following helpful lemma used in the privacy analysis of our clustering algorithm (Algorithm D.5). In summary, it says that when restricted to some set our input has small ℓ_2 sensitivity, we can first add Gaussian noise proportional to this sensitivity and afterwards treat this part of the input as "privatized". In particular, for the remainder of the privacy analysis we can treat this part as the same on adjacent inputs. Note that we phrase the lemma in terms of matrix inputs since this is what we use in our application. Of course, it also holds for more general inputs.

³⁵Note that this influences the exponent in the running time and sample complexity only by a constant factor and hence doesn't violate the assumptions of Theorem D.3.

³⁶In particular, this randomness is independent of the input

Lemma I.4. Let $V, V' \in \mathbb{R}^{n \times d}$, $m \in [n]$ and $\Delta > 0$ be such that there exists a set S of size at least $n - m$ satisfying

$$\forall i \in S. \|V_i - V'_i\|_2^2 \leq \Delta^2,$$

where V_i, V'_i denote the rows of V, V' , respectively. Let $\mathcal{A}_2: \mathbb{R}^{n \times d} \rightarrow \mathcal{O}$ be an algorithm that is $(\varepsilon_2, \delta_2)$ -differentially private in the standard sense, i.e., for all sets $S \subseteq \mathcal{O}$ and datasets $X, X' \in: \mathbb{R}^{n \times d}$ differing only in a single row it holds that

$$\mathbb{P}(\mathcal{A}_2(X) \in S) \leq e^{\varepsilon_2} \mathbb{P}(\mathcal{A}_2(X') \in S) + \delta_2.$$

Further, let $\mathcal{A}_1: \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$ be the Gaussian Mechanism with parameters $\Delta, \varepsilon_1, \delta_1$. I.e., on input M it samples $\mathbf{W} \sim N\left(0, 2\Delta^2 \cdot \frac{\log(2/\delta_1)}{\varepsilon_1^2}\right)^{n \times d}$ and outputs $M + \mathbf{W}$.

Then for

$$\begin{aligned} \varepsilon' &:= \varepsilon_1 + m\varepsilon_2, \\ \delta' &:= e^{\varepsilon_1} m e^{(m-1)\varepsilon_2} \delta_2 + \delta_1. \end{aligned}$$

$\mathcal{A}_2 \circ \mathcal{A}_1$ is (ε', δ') -differentially private with respect to V and V' , i.e., for all sets $S \subseteq \mathcal{O}$ it holds that

$$\mathbb{P}((\mathcal{A}_2 \circ \mathcal{A}_1)(V) \in S) \leq e^{\varepsilon'} \mathbb{P}((\mathcal{A}_2 \circ \mathcal{A}_1)(V') \in S) + \delta'.$$

Proof. Without loss of generality, assume that $S = \{1, \dots, m\}$. Denote by V_1, V_2 the first m and last $n - m$ rows of V respectively. Analogously for V'_1, V'_2 . We will later partition the noise \mathbf{W} of the Gaussian mechanism in the same way. Further, for a subset A of $\mathbb{R}^{n \times n}$ and $Y \in \mathbb{R}^{m \times n}$ define

$$T_{A,Y} = \left\{ X \in \mathbb{R}^{(n-m) \times n} \mid \begin{pmatrix} X \\ Y \end{pmatrix} \in A \right\} \subseteq \mathbb{R}^{(n-m) \times n}.$$

Note that $\begin{pmatrix} X \\ Y \end{pmatrix} \in A$ if and only if $X \in T_{A,Y}$.

Let $S \subseteq \mathcal{O}$. It now follows that

$$\begin{aligned} \mathbb{P}_{\mathcal{A}_2, \mathbf{W}}[(\mathcal{A}_2 \circ \mathcal{A}_1)(V) \in S] &= \mathbb{E}_{\mathcal{A}_2, \mathbf{W}}[\mathbb{1}\{V + \mathbf{W} \in \mathcal{A}_2^{-1}(S)\}] \\ &= \mathbb{E}_{\mathcal{A}_2, \mathbf{W}_2} \left[\mathbb{E}_{\mathbf{W}_1} \left[\mathbb{1}\left\{ \begin{pmatrix} V_1 + \mathbf{W}_1 \\ V_2 + \mathbf{W}_2 \end{pmatrix} \in \mathcal{A}_2^{-1}(S) \right\} \mid \mathbf{W}_2 \right] \right] \\ &= \mathbb{E}_{\mathcal{A}_2, \mathbf{W}_2} \left[\mathbb{E}_{\mathbf{W}_1} \left[\mathbb{1}\{V_1 + \mathbf{W}_1 \in T_{\mathcal{A}_2^{-1}(S), V_2 + \mathbf{W}_2}\} \mid \mathbf{W}_2 \right] \right] \\ &\leq e^{\varepsilon_1} \cdot \mathbb{E}_{\mathcal{A}_2, \mathbf{W}_2} \left[\mathbb{E}_{\mathbf{W}_1} \left[\mathbb{1}\{V'_1 + \mathbf{W}_1 \in T_{\mathcal{A}_2^{-1}(S), V_2 + \mathbf{W}_2}\} \mid \mathbf{W}_2 \right] \right] + \delta_1 \\ &= e^{\varepsilon_1} \cdot \mathbb{E}_{\mathcal{A}_2, \mathbf{W}} \left[\mathbb{1}\left\{ \begin{pmatrix} V'_1 + \mathbf{W}_1 \\ V_2 + \mathbf{W}_2 \end{pmatrix} \in \mathcal{A}_2^{-1}(S) \right\} \right] + \delta_1, \end{aligned}$$

where the inequality follows by the guarantees of the Gaussian Mechanism. Further, we can bound

$$\begin{aligned} \mathbb{E}_{\mathcal{A}_2, \mathbf{W}} \left[\mathbb{1}\left\{ \begin{pmatrix} V'_1 + \mathbf{W}_1 \\ V_2 + \mathbf{W}_2 \end{pmatrix} \in \mathcal{A}_2^{-1}(S) \right\} \right] &= \mathbb{E}_{\mathbf{W}} \left[\mathbb{E}_{\mathcal{A}_2} \left[\mathbb{1}\left\{ \mathcal{A}_2 \left(\begin{pmatrix} V'_1 + \mathbf{W}_1 \\ V_2 + \mathbf{W}_2 \end{pmatrix} \right) \in S \right\} \mid \mathbf{W} \right] \right] \\ &\leq e^{m\varepsilon_2} \cdot \mathbb{E}_{\mathbf{W}} \left[\mathbb{E}_{\mathcal{A}_2} \left[\mathbb{1}\left\{ \mathcal{A}_2 \left(\begin{pmatrix} V'_1 + \mathbf{W}_1 \\ V_2 + \mathbf{W}_2 \end{pmatrix} \right) \in S \right\} \mid \mathbf{W} \right] \right] + m e^{(m-1)\varepsilon_2} \delta_2 \\ &= e^{m\varepsilon_2} \cdot \mathbb{E}_{\mathcal{A}_2, \mathbf{W}} \left[\mathbb{1}\left\{ \begin{pmatrix} V'_1 + \mathbf{W}_1 \\ V_2 + \mathbf{W}_2 \end{pmatrix} \in \mathcal{A}_2^{-1}(S) \right\} \right] + m e^{(m-1)\varepsilon_2} \delta_2, \end{aligned}$$

where the inequality follows by the privacy guarantees of \mathcal{A}_2 combined with standard group privacy arguments.

Putting the above two displays together and plugging in the definition of ε', δ' we finally obtain

$$\mathbb{P}_{\mathcal{A}_2, \mathbf{W}}[(\mathcal{A}_2 \circ \mathcal{A}_1)(V) \in S] \leq e^{\varepsilon'} \mathbb{P}_{\mathcal{A}_2, \mathbf{W}}[(\mathcal{A}_2 \circ \mathcal{A}_1)(V') \in S] + \delta'.$$

□