

A Multiplicative Weights Mechanism for Privacy-Preserving Data Analysis

Moritz Hardt

Guy N. Rothblum

Abstract

We consider statistical data analysis in the interactive setting. In this setting a trusted curator maintains a database of sensitive information about individual participants, and releases privacy-preserving answers to queries as they arrive. Our primary contribution is a new differentially private multiplicative weights mechanism for answering a large number of interactive counting (or linear) queries that arrive online and may be adaptively chosen.

This is the first mechanism with worst-case accuracy guarantees that can answer large numbers of interactive queries and is *efficient* (in terms of the runtime's dependence on the data universe size). The error is asymptotically roughly *optimal* in its dependence on the number of participants, and depends only logarithmically on the number of queries being answered. The running time is nearly *linear* in the size of the data universe.

As a further contribution, when we relax the utility requirement and require accuracy only for databases drawn from a rich class of databases, we obtain exponential improvements in running time. Even in this relaxed setting we continue to guarantee privacy for *any* input database. Only the utility requirement is relaxed. Specifically, we show that when the input database is drawn from a *smooth* distribution — a distribution that does not place too much weight on any single data item — accuracy remains as above, and the running time becomes *poly-logarithmic* in the data universe size.

The main technical contributions are the application of multiplicative weights techniques to the differential privacy setting, a new privacy analysis for the interactive setting, and a technique for reducing data dimensionality for databases drawn from smooth distributions.

1 Introduction

Statistical analysis of sensitive information about individuals comes with important benefits. However, since the results of the analysis are often made public, these benefits might come at a serious cost to the privacy of individuals' sensitive data. A recent line of work, starting with the seminal works of Dinur and Nissim [DN03] and Dwork and Nissim [DN04] aims to provide a rigorous mathematical foundation for protecting the privacy of individuals in the setting of statistical analysis. This research has yielded the robust privacy guarantee of *differential privacy*, due to Dwork *et al.* [DMNS06], which guarantees that the outcome of the analysis on adjacent databases (databases that differ only in one participant's information) is "very similar" (in a strong sense). In particular, differential privacy guarantees that participation in the analysis does not incur significant additional risk for individuals (vs. non-participation). Throughout this paper and most of the prior work, the focus is on the setting where a trusted curator, holding a database of potentially sensitive information about n individuals, wishes to release statistics about the data while protecting individuals' privacy.

A central question in this line of research regards the tradeoff between utility and privacy. Namely, *what kinds of statistical queries can be answered, and with what accuracy, while protecting the privacy of individuals?* Early results seemed mixed: on one hand, moderate numbers of queries (smaller than the number of individuals in the database) could be answered with differential privacy and excellent accuracy by adding independent noise to the answers [DN03, DN04], and a growing body of subsequent work. On the other hand, it was shown that there exist specific families of simple queries such that answering "too many" of the queries (more than the database size) with "overly accurate" responses (smaller error than the sampling error¹) leads to blatant privacy violations [DN03, DMT07, DY08]. Still, these negative results left open the possibility that very rich statistical analyses, say answering huge families of queries, could be run in a privacy-preserving way, as long as their accuracy was slightly worse than the sampling error.

Non-interactive mechanisms. A beautiful work of Blum, Ligett and Roth [BLR08] showed that huge numbers of queries could in fact be answered in a privacy-preserving way. They considered *counting queries*: counting what fraction (between 0 and 1) of the participants in the analysis satisfy some property, where the "property" is taken to be a boolean predicate over the data universe U . They designed a privacy-preserving mechanism where for any set Q of counting queries specified non-interactively (i.e. in advance), the error scales only *logarithmically* with the number of queries being answered. Specifically, the error's dependence on the database size n and query set size k was roughly $(1/n^{1/3}) \cdot \log k$ (we ignore for now, and throughout the introduction, the constants and the many other parameters in the error expression). Moreover, the mechanism's output was a *synthetic database*: a (privacy-preserving) database of entries from the data universe U , where for each counting query the fraction of participants who satisfy it in the input and output database is within the error bound. This is a useful output format, as it is compatible with existing tools for analyzing databases and guarantees consistency. While this result showed that (information

¹The "sampling error" we refer to throughout the introduction is the error incurred by inferring statistical information about an underlying distribution of the population from n samples. This error is inherent to statistical data analysis. In fractional additive terms, the sampling error is asymptotically roughly $\tilde{O}(1/\sqrt{n})$ w.h.p.

theoretically at least) massive query sets could be answered in a privacy preserving way, it suffered from several disadvantages. First, this was a *non-interactive* mechanism: the query set had to be specified in advance, whereas previous mechanisms such as [DN04] were *interactive*, and allowed for answering arbitrary queries specified in an interactive and adaptive manner.² Moreover, the mechanism’s error was significantly larger than the $1/\sqrt{n}$ sampling error. This meant that, for fixed fractional accuracy, the number of participants in a data analysis needed to be significantly higher. Finally, the running time was very high: super-polynomial in the size $N = |U|$ of the data universe, and in k , the query set size.

These last two concerns were considered in a work of Dwork *et al.* [DNR⁺09]. They gave a non-interactive mechanism, whose running time was *polynomial* in N and k , where the error’s dependence on the query set size grew (roughly) as $(1/\sqrt{n}) \cdot k^{o(1)}$. This output was also a synthetic database. They showed that, under strong enough (exponential) cryptographic hardness assumptions, no general mechanism for answering counting queries that outputs synthetic data could have sublinear running-time in N or in k . In later work, Dwork, Rothblum and Vadhan [DRV10] obtained a mechanism with similar running time whose error was $(1/\sqrt{n}) \cdot \text{polylog}k$. Moreover, they showed how to obtain similar error bounds for *arbitrary* low-sensitivity queries (previous work was restricted to counting queries), for general queries the running time was no longer polynomial in N and k . Both of these mechanisms provide a slightly relaxed privacy guarantee known as (ϵ, δ) -differential privacy [DMNS06].

Interactive Mechanisms. For the *interactive* setting, where the queries are specified in an interactive and adaptive manner, it remained unclear whether large numbers of queries could be answered accurately in a privacy-preserving way. In a beautiful recent work, Roth and Roughgarden [RR10] presented a new mechanism for answering *interactive* counting queries, whose error scaled as $(1/n^{1/3}) \cdot \text{polylog}(k)$. They gave a super-polynomial time (in N and k) mechanism, and a separate polynomial-time mechanism that guaranteed similar error bound w.h.p. over a database drawn from a random distribution.

Several important questions remained unanswered, even for the case of counting queries:

1. Is there a an interactive mechanism that runs in time $\text{poly}(N)$ on each of the k queries with non-trivial error on all databases?
2. Could its error (in terms of k, n) match the sampling error $O(\sqrt{\log k/n})$?³
3. The result of [RR10] gives only (ϵ, δ) -differential privacy. Is there an interactive mechanism for handling many counting queries that achieves $(\epsilon, 0)$ -differential privacy?
4. Given the negative results of [DNR⁺09], we cannot hope for sub-linear running time in N . Do there exist mechanisms that match or nearly-match this hardness result?
5. What are open avenues for side-stepping the negative results of [DNR⁺09]? Namely, are there meaningful relaxations that permit mechanisms whose running time is sub-linear or even poly-logarithmic in N ?

²Note that in this setting, all other things being equal, an interactive mechanism is preferable to a non-interactive one: if we have an interactive mechanism, even if the queries are all specified in advance, we can still run the interactive mechanism on the queries, one by one, and obtain privacy-preserving answers to all of them.

³It can be shown that the maximum sampling error on k statistics is $O(\sqrt{\log k/n})$

1.1 This Work

Our main contribution is a new privacy-preserving interactive mechanism for answering counting queries, which we will refer to as the private multiplicative weights (PMW) mechanism. It allows us to give positive answers to the first four questions above, and to make partial progress on the last question. We proceed with a summary of our contributions, see Figure 1 for a comparison with the related work. We note that throughout this section, when we refer to a mechanism’s running time as being polynomial or linear, we are measuring the running time as a function of the data universe size N (which may be quite large for high-dimensional data).

An Interactive Mechanism. The PMW mechanism runs in linear-time and provides a worst-case accuracy guarantees for *all* input databases. The mechanism is presented Figure 2, its performance stated in the theorem below. The proof is in Section 4. See Section 3 for the formal definitions of accuracy and differential privacy for interactive mechanisms.

Theorem 1.1. *Let \mathcal{U} be a data universe of size N . For any $k, \epsilon, \delta, \beta > 0$, the Private Multiplicative Weights Mechanism of Figure 2, is an (ϵ, δ) -differentially private interactive mechanism.*

For any database of size n , the mechanism is (α, β, k) -accurate for (adaptive) counting queries over \mathcal{U} , where

$$\alpha = O\left(\frac{\sqrt{\log(k/\beta)\log(1/\delta)}\log^{1/4}N}{\sqrt{\epsilon n}}\right)$$

The running time in answering each query is $N \cdot \text{poly}(n) \cdot \text{polylog}(1/\beta, 1/\epsilon, 1/\delta)$.

The error as a function of n and k grows roughly as $\sqrt{\frac{\log k}{n}}$. Even for blatant non-privacy in the non-interactive setting this dependence on k and n is necessary [DN03]. Thus in terms of k and N our upper bound matches a lower bound that holds for a much weaker notion of privacy. In fact, it can be argued that $\sqrt{\log k/n}$ is just the statistical sampling error observed when computing the maximum error of k insensitive statistics on a sample of size n .

Moreover, the running time is only *linear* in N (for each of the k queries), nearly tight with the cryptographic hardness results of [DNR⁺09]. Previous work even in the non-interactive setting had higher polynomial running time. Finally, we remark that this mechanism can also be used to generate a synthetic database with similar error and running time bounds (in the non-interactive setting), see below for this extensions.

Achieving $(\epsilon, 0)$ -differential privacy. Prior to our work it was conceivable that there was no $(\epsilon, 0)$ -differentially private interactive release mechanism handling, say, n^2 counting queries with non-trivial error. However, using our multiplicative weights framework, we can achieve the guarantees of 5.1 also with $(\epsilon, 0)$ -differential privacy except for a somewhat worse dependence on n and $\log N$.

Theorem 1.2. *Let \mathcal{U} be a data universe of size N . For any $k, \epsilon, \beta > 0$, there is an $(\epsilon, 0)$ -differentially private interactive mechanism such that: For any database of size n , the mechanism is (α, β, k) -accurate for (adaptive) counting queries over \mathcal{U} , where*

$$\alpha = O\left(\frac{\log(k/\beta)^{1/3}\log^{1/3}N}{(\epsilon n)^{1/3}}\right).$$

The running time in answering each query is $N \cdot \text{poly}(n) \cdot \text{polylog}(1/\beta, 1/\epsilon, 1/\delta)$.

We also show a new lower bound on the error that any $(\epsilon, 0)$ -differentially private mechanism must have even in the non-interactive setting when answering $k \gg n$ counting queries. A novelty of our lower bound is that it simultaneously depends on $n, \log k, \log N$.

Theorem 1.3. *Let n be sufficiently large and let $\epsilon > 0$ be a constant independent of n . Then, for every $k \geq n^{1.1}$ there is a set of k counting queries over a universe of size N such that every $(\epsilon, 0)$ -differentially private mechanism for databases of size n must have error*

$$\alpha \geq \Omega(1) \cdot \left(\frac{\log k \cdot \log\left(\frac{N}{n}\right)}{\epsilon n} \right)^{1/2}$$

with probability $1/2$.

Relaxed Notions of Utility. To answer Question 5 that was raised in the introduction, we begin with a discussion of the negative results of [DNR⁺09] and possible avenues for side-stepping them. The negative results for producing synthetic data can be side-stepped by a mechanism whose output has a different format. This is a promising avenue, but synthetic data is a useful output format. It is natural to try to side-step hardness while continuing to output synthetic data. One possibility is working for restricted query classes, but recent work of Ullman and Vadhan [UV11] shows hardness even for very simple and natural query classes such as conjunctions. In the known hardness results, however, the *databases* (or rather database distributions) that are hard to sanitize are (arguably) “unnatural”, containing cryptographic data in [DNR⁺09] and PCP proofs for the validity of digital signatures in [UV11]. Thus, a natural approach to side-stepping hardness is relaxing the utility requirement, and not requiring accuracy for *every* input database.

A mechanism that works only for some input databases is only as interesting as the class of databases for which accuracy is guaranteed. For example, getting accuracy w.h.p. for *most* databases is simple, since (speaking loosely and informally) *most* databases behave like a uniformly random database. Thus, we can get privacy and accuracy by ignoring the input database (which gives perfect privacy) and answering according a new database drawn uniformly at random (which, for most input databases, will give fairly accurate answers).

Smooth databases and sublinear time. We consider accuracy guarantees for the class of (*pseudo*)-*smooth* databases. Intuitively, we think of these as databases sampled i.i.d. from *smooth* underlying distributions over the data universe U . I.e., underlying distributions that do not put too much weight on any particular data item (alternatively, they have high min-entropy). We say that a histogram or distribution y over U is ξ -*smooth*, if for every $u \in U$, the probability of u by y is at most ξ . We say that a histogram or database $x \in U^n$ is (ξ, ϕ) -*pseudo-smooth w.r.t a set \mathcal{Q} of queries* if there exists some ξ -smooth y that approximates it well w.r.t every query in \mathcal{Q} . I.e., for every $f \in \mathcal{Q}$, $|f(y) - f(x)| \leq \phi$ (where by $f(y)$ we mean the expectation of f over data items drawn from y). See Section 6 for formal definitions.

The PMW mechanism yields a mechanism with improved running time—sub-linear, or even polylogarithmic in N —for pseudo-smooth databases. The new mechanism (with smoothness parameter ξ) runs in time that depends linearly on ξN rather than N . It

guarantees differential privacy for *any* input database. Its error is similar to that of the mechanism of Theorem 5.1 (up to an additional ϕ error), but this accuracy guarantee is only: (i) for a set \mathcal{Q} of interactive counting queries that are fixed in advance (i.e. non-adaptively). We note that the mechanism is interactive in the sense that it need not know the queries in advance, but accuracy is not guaranteed for adversarially chosen queries (see the discussion in Section 2 for motivation for this relaxation), and (ii) for input databases that are (ξ, ϕ) -smooth with respect to the query class \mathcal{Q} . The performance guarantees are in Theorem 1.4 below. The proof is in Section 6

Theorem 1.4 (Smooth PMW). *Let U be a data universe of size N . For any $\varepsilon, \delta, \beta, \xi, \phi > 0$, the Private Multiplicative Weights Mechanism of Figure 2 is an (ε, δ) -differentially private interactive mechanism. For any sequence \mathcal{Q} of k interactive counting queries over U that are fixed in advance (non-adaptively), for any database of size n that is (ξ, ϕ) -pseudo-smooth w.r.t \mathcal{Q} , the mechanism is (α, β, k) -non adaptively accurate w.r.t. \mathcal{Q} , where*

$$\alpha = \tilde{O}\left(\phi + \frac{\log(1/\delta) \log^{1/4}(\xi N) \cdot (\log k + \log(1/\beta))}{\sqrt{n} \cdot \varepsilon}\right)$$

The running time in answering each query is $(\xi N) \cdot \text{poly}(n) \cdot \text{polylog}(1/\beta, 1/\varepsilon, 1/\delta, 1/\xi, 1/\phi)$.

In particular, for very good smoothness $\xi = \text{polylog} N/N$, the running time will depend only poly-logarithmically on N . The main observation for achieving this improved running time is that for (pseudo)-smooth databases we can effectively reduce the data universe size by sub-sampling, and then apply our algorithm to the smaller data universe. The mechanism does not require knowledge of the histogram which certifies that the given input database is pseudosmooth.

The privacy guarantee is the standard notion of differential privacy. I.e., privacy holds always and for every database. The accuracy guarantee is only for pseudosmooth databases, and we interpret it as follows. The dataset is drawn i.i.d from an unknown underlying distribution D (the standard view in statistics). The mechanism guarantees accuracy and sub-linear efficiency as long as the underlying data distribution is smooth. If the underlying distribution is ξ -smooth, then w.h.p. the database x (which we think of as being drawn i.i.d from D and of large enough size) is “close” to D on every query $f \in \mathcal{Q}$, and so w.h.p. x is (ξ, ϕ) -smooth and the mechanism is accurate. An important aspect of this guarantee is that *there is no need to know what the underlying distribution is*, only that it is smooth. A promising approach in practice may be to run this mechanism as a very efficient heuristic. The heuristic *guarantees* privacy, and also has a rigorous accuracy guarantee under assumptions about the underlying distribution. We note that Dwork and Lei [DL09] also proposed mechanisms that always guarantee privacy, but guarantee accuracy only for a subset of databases (or underlying distributions).

We also note that [RR10] considered databases drawn from a distribution that was itself picked randomly from the set of all distributions. Such “random distributions” are indeed very smooth (w.h.p. $\xi \leq O(\log N/N)$) and therefore a special case of our model.

An interesting direction for future work is finding differentially private mechanisms for other and more useful or well motivated families of databases, or finding natural applications where pseudo-smooth databases are of particular interest. We note that (as one would expect

Mechanism	interactive?	error in terms of n, N, k	runtime	privacy	remark
[DMNS06]	✓	\sqrt{k}	—	(ϵ, δ)	
[BLR08]	—	$n^{2/3} \log^{1/3} N \cdot \log^{1/3} k$	$N^{O(n)}$	$(\epsilon, 0)$	
[DNR ⁺ 09]	—	$\sqrt{n \log N} \cdot k^{o(1)}$	$\text{poly}(N)$	(ϵ, δ)	
[DRV10]	—	$\sqrt{n \log N} \cdot \log^2 k$	$\text{poly}(N)$	(ϵ, δ)	
[RR10]	✓	$n^{2/3} \log^{1/3} N \cdot \log k$	$N^{O(n)}$	(ϵ, δ)	
[RR10]	✓	$n^{2/3} \log^{1/3} N \cdot \log k$	$\text{poly}(N)$	(ϵ, δ)	random DBs
This work	✓	$\sqrt{n \log k} \log^{1/4} N$	$\tilde{O}(N)$	(ϵ, δ)	
This work	✓	$n^{2/3} \log^{1/3} N \log^{1/3} k$	$\tilde{O}(N)$	$(\epsilon, 0)$	
This work	✓	$\sqrt{n \log k} \log \log N \cdot \log k$	$\text{polylog} N$	(ϵ, δ)	smooth DBs

Figure 1: Comparison to previous work for k linear queries each of sensitivity 1. For simplicity the dependence on δ is omitted from the comparison. Runtime stated in terms of N omitting other factors. Error bounds are a factor n larger than throughout the paper and accurate up to polylog factors. Note that random databases are a special case of smooth databases (see Section 6).

given these positive results) the negative results for producing synthetic data are for databases that are neither smooth nor pseudo-smooth.

Figure 1 summarizes and compares relevant past work on answering counting queries. We proceed with an overview of techniques.

2 Overview of proof and techniques

Multiplicative Weights. We use a (privacy-preserving) multiplicative weights mechanism (see [LW94, AHK05]). The mechanism views databases as histograms or distributions (also known as “fractional” databases) over the data universe U (as was done in [DNR⁺09]). At a high level, the mechanism works as follows. The real database being analyzed is x (we view x as distribution or histogram over U , with positive weight on the data items in x). The mechanism also maintains an updated fractional database, denoted as x_t at the end of round t . In each round t , after the t -th counting query f_t has been specified, x_{t-1} is updated to obtain x_t . The initial database x_0 is simply the uniform distribution over the data universe. I.e., each coordinate $u \in U$ has weight $1/N$.

In the t -th round, after the t -th query f_t has been specified, we compute a noisy answer \widehat{a}_t by adding (properly scaled) Laplace noise to $f_t(x)$ —the “true” answer on the real database. We then compare this noisy answer with the answer given by the previous round’s database $f_t(x_{t-1})$. If the answers are “close”, then this is a “lazy” round, and we simply output $f_t(x_{t-1})$ and set $x_t \leftarrow x_{t-1}$. If the answers are “far”, then this is an “update” round and we need to update or “improve” x_t using a multiplicative weights re-weighting. The intuition is that the re-weighting brings x_t “closer” to an accurate answer on f_t . In a nutshell, this is all the algorithm does. The only additional step required is bounding the number of “update” rounds: if the total number of update rounds grows to be larger than (roughly) n , then the mechanism fails and terminates. This will be a low probability event. See Figure 2 for the details. Given this overview of the algorithm, it remains to specify how to: (i) compute $f_t(x_{t-1})$, and (ii) re-weight or improve the database on update rounds. We proceed with an

overview of the arguments for accuracy and privacy.

For this exposition, we think of the mechanism as explicitly maintaining the x_t databases, resulting in complexity that is roughly linear in $N = |U|$. Using standard techniques we can make the memory used by the mechanism logarithmic in N (computing each coordinate of x_t as it is needed). Either way, it is possible to compute $f_t(x_{t-1})$ in linear time.

The re-weighting (done only in update rounds), proceeds as follows. If in the comparison we made, the answer according to x_{t-1} was “too small”, then we increase by a small multiplicative factor the weight of items $u \in U$ that satisfy the query f_t ’s predicate, and decrease the weight of those that do not satisfy it by the same factor. If the answer was “too large” then do the reverse in terms of increasing and decreasing the weights. We then normalize the resulting weights to obtain a new database whose entries sum to 1. The intuition, again, is that we are bringing x_t “closer” to an accurate answer on f_t . The computational work scales linearly with N .

To argue accuracy, observe that as long as the number of update rounds stays below the (roughly n) threshold, our algorithm ensures bounded error (assuming the Laplace noise we add is not too large). The question is whether the number of update rounds remains small enough. This is in fact the case, and the proof is via a multiplicative weights potential argument. Viewing databases as distributions over U , we take the *potential* of database y to be the relative entropy $\text{RE}(x||y)$ between y and the real database x . We show that if the error of x_{t-1} on query f_t is large (roughly larger than $1/\sqrt{n}$), then the potential of the re-weighted x_t is smaller by at least (roughly) $1/n$ than the potential of x_{t-1} . Thus, in every “update” round, the potential drops, and the drop is significant. By bounding the potential of x_0 , we get that the number of update rounds is at most (roughly) n .

“Pay as you go” privacy analysis. At first glance, privacy might seem problematic: we access the database and compute a noisy answer *in every round*. Since the number of queries we want to answer (number of rounds) might be huge, unless we add a huge amount of noise this collection of noisy answers *is not privacy preserving*. The point, however, is that in most rounds we don’t release the noisy answer. All we do is check whether or not our current database x_{t-1} is accurate, and if so we use it to generate the mechanism’s output. In all but the few update rounds, the perturbed true answer is not released, and we want to argue that privacy in all those lazy rounds comes (essentially) “for free”. The argument builds on ideas from privacy analyses in previous works [DNR⁺09, DNPR10, RR10]).

A central concern is arguing that the “locations” of the update rounds be privacy-preserving (there is an additional, more standard, concern that the noisy answers in the few update rounds also preserve privacy). Speaking intuitively (and somewhat inaccurately), for any two adjacent databases, there are w.h.p. only roughly n “borderline” rounds, where the noise is such that on one database this round is update and on another this round is lazy. This is because, conditioning on a round being “borderline”, with constant probability it is actually an “update” round. Since the number of update rounds is at most roughly n , with overwhelming probability the number of borderline rounds also is roughly n . For non-borderline rounds, those rounds’ being an update or a lazy round is determined similarly for the two databases, and so privacy for these rounds come “for free”. The borderline rounds are few, and so the total privacy hit incurred for them is small.

Given this intuition, we want to argue that the “privacy loss”, or “confidence gain” of

an adversary, is small. At a high level, if we bound the worst-case confidence gain in each update round by roughly $O(\varepsilon/\sqrt{n})$, then by an “evolution of confidence” argument due to [DN03, DN04, DRV10], the total confidence gain of an adversary over the roughly n update rounds will be only ε w.h.p. To bound the confidence gain, we define “borderline” rounds as an event over the noise values on a database x , and show that: (1) Conditioned on a round being borderline on x , it will be an update round on x w.h.p. This means borderline rounds are few. (2) Conditioned on a round being borderline on x , the worst-case confidence gain of an adversary viewing the mechanism’s behavior in this round on x vs. an adjacent x' is bounded by roughly ε/\sqrt{n} . This means the privacy hit in borderline rounds isn’t too large, and we can “afford” roughly n of them. (3) Conditioned on a round *not* being borderline, there is no privacy loss in this round on x vs. any adjacent x' . I.e., non-borderline rounds come for free (in terms of privacy).

This analysis allows us to add less noise than previous works, while still maintaining (ε, δ) differential privacy. It may find other applications in interactive or adaptive privacy settings. Details are in Section 4.2.

Sublinear Time Mechanism for Smooth Databases. We observe that we can modify the PMW mechanism to work over a smaller data universe $V \subseteq U$, as long as there *exists* a database x^* whose support is only over V , and gives close answers to those of x on every query we will be asked. We modify the algorithm to maintain multiplicative weights only over the smaller set V , and increase slightly the inaccuracy threshold for declaring a round as “update”. For the analysis, we modify the potential function: it measures relative entropy to x^* rather than x . In update rounds, the distance between x_{t-1} and this new x^* on the current query is large (since x^* is close to x , and x_{t-1} is far from x). This means that re-weighting will reduce $\text{RE}(x^*||x_{t-1})$, and even though we maintain multiplicative weights only over a smaller set V , the number of update rounds will be small. Maintaining multiplicative weights over V rather than U reduces the complexity from linear in $|U|$ to linear in $|V|$.

To use the above observation, we argue that for any large set of counting queries \mathcal{Q} and any (ξ, ϕ) -pseudo-smooth database x , if we choose a uniformly random small (but not too small) sub-universe $V \subseteq U$, then w.h.p there *exists* x^* whose support is in V that is close to x on all queries in \mathcal{Q} . In fact, sampling a sub=universe of size roughly $\xi N \cdot n \cdot \log|\mathcal{Q}|$ suffices. This means that indeed PMW can be run on the reduced data universe V with reduced computational complexity. See Section 6.1 for this argument.

Utility here is for a fixed non-adaptive set \mathcal{Q} of queries (that need not be known in advance). We find this utility guarantee to still be well motivated—note that, privacy aside, the input database itself, which is sampled i.i.d from an underlying distribution, isn’t guaranteed to yield good answers for adaptively chosen queries). Finally, we remark that this technique for reducing the data universe size (the data dimensionality) may be more general than the application to PMW. In particular, previous mechanisms such as [DNR⁺09, DRV10] can also be modified to take advantage of this sampling and obtain improved running time for smooth databases (the running time will be polynomial, rather than linear as it is for the PMW mechanism).

Synthetic databases. We conclude by noting that the PMW mechanism can be used to generate synthetic data (in the non-interactive setting). To do this, iterate the mechanism

over a set of queries \mathcal{Q} , repeatedly processing all the queries in \mathcal{Q} and halting when either (i) we made roughly $n + 1$ iterations, i.e. have processed every query in \mathcal{Q} n times, or (ii) we have made a complete pass over all the queries in \mathcal{Q} without any update rounds (whichever of these two conditions occurs first). If we make a complete pass over \mathcal{Q} without any update rounds, then we know that the x_t we have is accurate for all the queries in \mathcal{Q} and we can release it (or a subsample from it) as a privacy-preserving synthetic database. By the potential argument, there can be at most roughly n update rounds. Thus, after $n + 1$ iterations we are guaranteed to have a pass without any update rounds. Previous mechanisms for generating synthetic databases involved linear programming and were more expensive computationally.

3 Preliminaries

Probability tools. Let $x, y \in \mathbb{R}^N$. We define the *relative entropy* or *Kullback-Leibler divergence* between x and y as:

$$\text{RE}(x||y) = \sum_{i \in [N]} x_i \log\left(\frac{x_i}{y_i}\right) + y_i - x_i. \quad (1)$$

This reduces to the more familiar expression $\sum_i x_i \log\left(\frac{x_i}{y_i}\right)$ when $\sum_i x_i = \sum_i y_i = 1$ (in particular this happens when x, y correspond to distributions over $[N]$).

The following facts about relative entropy are well-known and easy to verify.

Fact 3.1. For every $x, y \in \mathbb{R}^N$, we have $\text{RE}(x||y) \geq 0$. Equality holds if and only if $x = y$.

We let $\text{Lap}(\sigma)$ denote the one-dimensional Laplacian distribution centered at 0 with scaling σ and corresponding density $f(x) = \frac{1}{2\sigma} \exp\left(-\frac{|x|}{\sigma}\right)$. We utilize the following lemma about the convexity of the KL-divergence.

Lemma 3.2. Let P, Q be arbitrary distributions over a common probability space. Suppose there are distributions P_1, P_2, Q_1, Q_2 and $\lambda \in [0, 1]$, so that $P = \lambda P_1 + (1 - \lambda)P_2$ and $Q = \lambda Q_1 + (1 - \lambda)Q_2$. Then,

$$\mathbb{E}_{\mathbf{v} \sim P} \left[\log\left(\frac{P(\mathbf{v})}{Q(\mathbf{v})}\right) \right] \leq \lambda \mathbb{E}_{\mathbf{v} \sim P_1} \left[\log\left(\frac{P_1(\mathbf{v})}{Q_1(\mathbf{v})}\right) \right] + (1 - \lambda) \mathbb{E}_{\mathbf{v} \sim P_2} \left[\log\left(\frac{P_2(\mathbf{v})}{Q_2(\mathbf{v})}\right) \right]. \quad (2)$$

In other words, $\text{RE}(P||Q) \leq \lambda \text{RE}(P_1, Q_1) + (1 - \lambda) \text{RE}(P_2, Q_2)$.

The next lemma shows that ε -differential privacy translates into relative entropy $2\varepsilon^2$. This lemma was shown in [DRV10].

Lemma 3.3 ([DRV10]). Let P, Q be any two distributions on a common support \mathcal{S} with density functions dP and dQ , respectively. Suppose that

$$\sup_{v \in \mathcal{S}} \log\left(\frac{dP(v)}{dQ(v)}\right) \leq \varepsilon_0.$$

Then,

$$\mathbb{E}_{\mathbf{v} \sim P} \left[\log\left(\frac{P(\mathbf{v})}{Q(\mathbf{v})}\right) \right] \leq 2\varepsilon_0^2.$$

We also use the following general large deviation bound.

Lemma 3.4 (Method of Bounded Differences). *Let X_1, \dots, X_m be an arbitrary set of random variables and let f be a function satisfying the property that for every $j \in [m]$ there is a number $c_j \geq 0$ such that*

$$|\mathbb{E}[f \mid X_1, X_2, \dots, X_j] - \mathbb{E}[f \mid X_1, X_2, \dots, X_{j-1}]| \leq c_j.$$

Then,

$$\Pr\{f > \mathbb{E}f + \lambda\} \leq \exp\left(-\frac{\lambda^2}{2 \sum_{j \in [m]} c_j^2}\right). \quad (3)$$

Notation. We denote by $\langle x, y \rangle = \sum_{i \in [N]} x_i y_i$ the real valued inner product between two vectors $x, y \in \mathbb{R}^N$. When $x \in \mathbb{R}^S$ is a vector supported on a subset of the coordinates $S \subseteq [N]$ and $y \in \mathbb{R}^N$, we still write $\langle x, y \rangle = \sum_{i \in S} x_i y_i$.

Histograms and linear queries. A *histogram* $x \in \mathbb{R}^N$ represents a database or data distribution over a universe U of size $|U| = N$. We will assume that x is normalized so that $\sum_{i \in U} x_i = 1$. We use histograms in the natural way to denote standard databases of size n (n -item multisets in U), and also to denote distributions over the data universe. The only difference is that databases have support size n , whereas distributions do not necessarily have small support.

In this work we focus on *linear queries* $f: \mathbb{R}^N \rightarrow [0, 1]$. As usual we may view a linear query as a vector $f \in [0, 1]^N$. We then use the equality $f(x) = \langle f, x \rangle$, where the histogram x can be either an n -item database or a data distribution. By our normalization of x , the sensitivity of a linear query is $1/n$. While we assume that $f_t \in [0, 1]^N$, our algorithm applies to any linear query $f_t \in [-c, c]^N$ by considering the query defined as $1/2 + f_t[i]/2c$ in coordinate i . In this case, the error of the algorithm scales linearly in c .

A special case of linear queries are *counting queries*. A counting query associated with a predicate from U to $\{0, 1\}$, outputs what fraction of the items in its input database satisfy the predicate. We view a counting query f as a vector over $\{0, 1\}^N$ specifying which data items satisfy the query's predicate.

Accuracy and privacy in the interactive setting. Formally, an *interactive mechanism* $M(x)$ is a stateful randomized algorithm which holds a histogram $x \in \mathbb{R}^N$. It receives successive counting queries $f_1, f_2, \dots \in \mathcal{F}$ one by one, and in each round t , on query f_t , it outputs a (randomized) answer a_t (a function of the input histogram, the internal state, and the mechanism's coins). For privacy guarantees, we *always* assume that the queries are given to the mechanism in an adversarial and adaptive fashion by a randomized algorithm A called the *adversary*. For accuracy guarantees, while we usually consider adaptive adversarial, we will also consider non-adaptive adversarial queries chosen in advance—we still consider such a mechanism to be interactive, because it does not know in advance what these queries will be. The main query class we consider throughout this work is the class \mathcal{F} of all counting queries, as well as sub-classes of it.

Definition 3.5 ((α, β, k) -Accuracy in the Interactive Setting). We say that a mechanism M is (α, β, k) -(adaptively) accurate for a database x , if when it is run for k rounds, for any (adaptively chosen) counting queries, with all but β probability over the mechanism's coins $\forall t \in [k], |a_t - f_t(x)| \leq \alpha$.

We say that a mechanism M is (α, β, k) -non adaptively accurate for a query sequence \mathcal{Q} of size k and a database x , if when it is run for k rounds on the queries in \mathcal{Q} , with all but β probability over the mechanism's coins $\forall t \in [k], |a_t - f_t(x)| \leq \alpha$.

For privacy, the interaction of a mechanism $M(x)$ and an adversary A specifies a probability distribution $[M(x), A]$ over *transcripts*, i.e., sequences of queries and answers $(f_1, a_1, f_2, a_2, \dots, f_k, a_k)$. Let $\text{Trans}(\mathcal{F}, k)$ denote the set of all transcripts of any length k with queries from \mathcal{F} . We will assume that the parameter k is known to the mechanism ahead of time. Our privacy requirement asks that the entire transcript satisfies differential privacy.

Definition 3.6 ((ϵ, δ) -Differential Privacy in the Interactive Setting). We say a mechanism M provides (ϵ, δ) -differential privacy for a class of queries \mathcal{F} , if for every adversary A and every two histograms $x, x' \in \mathbb{R}^N$ satisfying $\|x - x'\|_1 \leq 1/n$, the following is true: Let $P = [M(x), A]$ denote the transcript between $M(x)$ and A . Let $Q = [M(x'), A]$ denote the transcript between $M(x')$ and A . Then, for every $S \subseteq \text{Trans}(\mathcal{F}, k)$, we have

$$P(S) \leq e^\epsilon Q(S) + \delta.$$

We will find it useful to work with the following condition, which (by Lemma 3.7 below) is no weaker than (ϵ, δ) privacy:

$$\Pr_{\mathbf{v} \sim P} \left\{ \left| \log \left(\frac{P(\mathbf{v})}{Q(\mathbf{v})} \right) \right| > \epsilon \right\} \leq \delta. \quad (4)$$

(Note that here we are identifying the distribution P with its density function dP .)

Lemma 3.7. *Condition (4) implies (ϵ, δ) -differential privacy.*

Proof. Indeed, suppose (4) is satisfied and consider $B = \{v: |\log(Pv/Qv)| > \epsilon\}$. Let $S \subseteq \text{Trans}(\mathcal{F})$ and consider $S_1 = S \cap B$ and $S_2 = S \cap B^c$. We then know that $P(S) = P(S_1) + P(S_2) \leq \delta + e^\epsilon Q(S_2) \leq e^\epsilon Q(S) + \delta$. ■

4 Private multiplicative weights mechanism for linear queries

In the PMW mechanism of Figure 2, in each round t , we are given a linear query f_t over U and x_t denotes a fractional histogram (distribution over $V \subseteq U$) computed in round t . The domain of this histogram is V rather than U . Here, V could be much smaller than U and this allows for some flexibility later, in proving Theorem 1.4, where we aim for improved efficiency. For this section, unless otherwise specified, we assume that $V = U$. In particular this is the case in the statement of Theorem 5.1, the main theorem that we prove in this section.

We use a_t to denote the true answer on the database on query t , and \widehat{a}_t denotes this same answer with noise added to it. We use d_t to denote the difference between the true answer a_t and the answer given by x_{t-1} , i.e.,

$$d_t = f_t(x_{t-1}) - f_t(x).$$

We denote by \widehat{d}_t the difference between the *noisy* answer and the answer given by x_{t-1} . The boolean variable w_t denotes whether the noisy difference was large or small. If \widehat{d}_t is smaller (in absolute value) than $\approx 1/\sqrt{n}$, then this round is *lazy* and we set $w_t = 0$. If \widehat{d}_t is larger than threshold then this is an *update* round and we set $w_t = 1$.

Parameters: A subset of the coordinates $V \subseteq U$ with $|V| = M$ (by default $V = U$), intended number of rounds $k \in \mathbb{N}$, privacy parameters $\varepsilon, \delta > 0$ and failure probability $\beta > 0$. See (6) for the setting of η, σ, T .

Input: Database $D \in U^n$ corresponding to a histogram $x \in \mathbb{R}^N$

Algorithm: Set $y_0[i] = x_0[i] = 1/M$ for all $i \in V$

In each round $t \leftarrow 1, 2, \dots, k$ when receiving a linear query f_t do the following:

1. Sample $A_t \sim \text{Lap}(\sigma)$. Compute the noisy answer $\widehat{a}_t \leftarrow f_t(x) + A_t$.
2. Compute the difference $\widehat{d}_t \leftarrow f_t(x_{t-1}) - \widehat{a}_t$:
 - If $|\widehat{d}_t| \leq T$, then set $w_t \leftarrow 0$, $x_t \leftarrow x_{t-1}$, output $f_t(x_{t-1})$, and proceed to the next iteration.
 - If $|\widehat{d}_t| > T$, then set $w_t \leftarrow 1$ and:
 - for all $i \in V$, update

$$y_t[i] \leftarrow x_{t-1}[i] \cdot \exp(-\eta \cdot r_t[i]), \quad (5)$$
 where $r_t[i] = f_t[i]$ if $\widehat{d}_t > 0$ and $r_t[i] = 1 - f_t[i]$ otherwise.
 - Normalize, $x_t[i] \leftarrow \frac{y_t[i]}{\sum_{i \in V} y_t[i]}$.
 - If $\sum_{j=1}^t w_j > \eta^{-2} \log M$, then abort and output “failure”. Otherwise, output the noisy answer \widehat{a}_t and proceed to the next iteration.

Figure 2: Private Multiplicative Weights (PMW) Mechanism

Choice of parameters. We set the parameters η, σ, T as follows:

$$\eta = \sqrt{\frac{\log^{1/2} M \log(k/\beta) \log(1/\delta)}{\varepsilon n}} \quad \sigma = \frac{10\eta}{\log(k/\beta)} \quad T = 40\eta. \quad (6)$$

To understand the choice of parameters, let $m = \eta^{-2} \log M$ denote the bound on the number of update rounds ensured by our algorithm. We chose our parameters in (6) such that the following two relations hold

$$\sigma n \geq \frac{10\sqrt{m} \log(1/\delta)}{\varepsilon} \quad \text{and} \quad T \geq 4\sigma \log(k/\beta). \quad (7)$$

Intuitively speaking, the first condition ensures that the scaling σ of the Laplacian variables used in our algorithm is large enough to handle m update rounds while providing (ε, δ) -differential privacy. The second condition ensures that the Laplacian variables are small compared to the threshold T . Subject to these two constraints expressed in (7), our goal is to minimize η and σ . This is because η controls how large the noise magnitude σ has to be which in turns determines the threshold T . The error of our algorithm must scale with T .

We are now ready to prove Theorem 5.1, i.e. the utility and privacy of the PMW mechanism. This follows directly from our utility analysis provided in in Section 4.1 and our privacy argument presented in Section 4.2.

4.1 Utility analysis

To argue utility, we need to show that even for very large total number of rounds k , the number of update rounds is at most roughly n with high probability. This is done using a potential argument. Intuitively, the potential of a database x_t is the relative entropy between the true histogram x and our estimate x_t .

Since in general $V \neq U$, we will actually define the potential with respect to a target histogram $x^* \in \mathbb{R}^N$ with support only over V . This x^* need not be equal to x , nor does it have to be known by the algorithm. This added bit of generality will be useful for us later in Section 6 when we modify the mechanism to run in sublinear time. For this section, however, unless we explicitly note otherwise the reader may think of x^* as being equal to x . The potential function is then defined as

$$\Psi_t = \text{RE}(x^*||x_t) = \sum_{i \in V} x^*[i] \log \left(\frac{x^*[i]}{x_t[i]} \right). \quad (8)$$

Note that x^* and x_t are both normalized so that we can think of them both as distributions or histograms over U . We start with two simple observations:

Lemma 4.1. $\Psi_0 \leq \log M$.

Proof. Indeed, by the nonnegativity of entropy $H(x^*)$ we get that $\Psi_0 = \log M - H(x^*) \leq \log M$. ■

Second, by the nonnegativity of relative entropy (Fact 3.1), we have $\Psi_t \geq 0$ for every t .

Lemma 4.2. For every t , we have $\Psi_t \geq 0$.

Proof. By the nonnegativity of relative entropy (Fact 3.1). ■

Our goal is to show that if a round is an update round (and $w_t = 1$), then the potential drop in that round is at least η^2 . In Lemma 4.5 we show that this is indeed the case in every round, except with β/k probability over the algorithm's coins. Taking a union bound, we conclude that with all but β probability over the algorithm's coins, there are at most $\eta^{-2} \cdot \log M$ update rounds. The next lemma quantifies the potential drop in terms of the penalty vector r_t and the parameter η using a multiplicative weights argument.

Lemma 4.3. In each update round t , we have $\Psi_{t-1} - \Psi_t \geq \eta \langle r_t, x_{t-1} - x^* \rangle - \eta^2$.

Proof. We can rewrite the potential drop as follows:

$$\begin{aligned}
\Psi_{t-1} - \Psi_t &= \sum_{i \in V} x^*[i] \left(\log \left(\frac{x^*[i]}{x_{t-1}[i]} \right) - \log \left(\frac{x^*[i]}{x_t[i]} \right) \right) \\
&= \sum_{i \in V} x^*[i] \log \left(\frac{x_t[i]}{x_{t-1}[i]} \right) \\
&= \sum_{i \in V} x^*[i] \log \left(\exp(-\eta r_t[i]) \frac{x_{t-1}[i]}{\sum_{i \in V} y_t[i]} \right) \\
&= -\eta \langle r_t, x^* \rangle - \sum_{i \in V} x^*[i] \log \left(\sum_{i \in V} y_t[i] \right) \\
&= -\eta \langle r_t, x^* \rangle - \log \left(\sum_{i \in V} \exp(-\eta r_t[i]) x_{t-1}[i] \right) \quad (\text{since } \sum x^*[i] = 1)
\end{aligned}$$

Note that

$$\exp(-\eta r_t[i]) \leq 1 - \eta r_t[i] + \eta^2 r_t[i]^2 \leq 1 - \eta r_t[i] + \eta^2.$$

Using this and $\sum x_{t-1}[i] = 1$ we get

$$\log \left(\sum_{i \in V} \exp(-\eta r_t[i]) x_{t-1}[i] \right) \leq \log(1 - \eta \langle r_t, x_{t-1} \rangle + \eta^2) \leq -\eta \langle r_t, x_{t-1} \rangle + \eta^2,$$

where we used $\log(1 + y) \leq y$ for $y > -1$. We conclude that

$$\Psi_{t-1} - \Psi_t \geq -\eta \langle r_t, x^* \rangle + \eta \langle r_t, x_{t-1} \rangle - \eta^2 = \eta \langle r_t, x_{t-1} - x^* \rangle - \eta^2. \quad \blacksquare$$

In the following lemmata, we condition on the event that $|A_t| \leq T/2$. Since A_t is a centered Laplacian with standard deviation σ and $T \geq 4\sigma(\log k + \log(1/\beta))$, this event occurs with all but β/k probability in every round t .

The next lemma connects the inner product $\langle r_t, x^* - x_{t-1} \rangle$ with the “error” of x_{t-1} on the query f_t . Here, error is measured with respect to the true histogram x . To relate x with x^* , we further denote

$$\text{err}(x^*, f_t) = |f_t(x^*) - f_t(x)|. \quad (9)$$

When $x^* = x$ we get that $\text{err}(x^*, f_t) = 0$ always, and in general we will be interested in x^* databases where $\text{err}(x^*, f_t)$ is small for all $t \in [k]$.

Lemma 4.4. *In each round t where $|\widehat{d}_t| \geq T$ and $|A_t| \leq T/2$ we have*

$$\langle r_t, x_{t-1} - x^* \rangle \geq |f_t(x) - f_t(x_{t-1})| - \text{err}(x^*, f_t).$$

Proof. By assumption $|\widehat{d}_t| \geq T$ and $|d_t - \widehat{d}_t| \leq |A_t| \leq T/2$. Hence, $\text{sign}(d_t) = \text{sign}(\widehat{d}_t)$. We distinguish the two cases where $\text{sign}(d_t) < 0$ and $\text{sign}(d_t) \geq 0$. First, suppose

$$0 > \text{sign}(d_t) = \text{sign}(f_t(x_{t-1}) - f_t(x)).$$

It follows that $r_t[i] = 1 - f_t[i]$. Hence,

$$\begin{aligned}
\sum_{i \in V} r_t[i](x_{t-1}[i] - x^*[i]) &= -(f_t(x_{t-1}) - f_t(x^*)) + \sum_{i \in V} x^*[i] - \sum_{i \in V} x_{t-1}[i] \\
&= -(f_t(x_{t-1}) - f_t(x^*)) \quad (\text{using } \sum_i x_{t-1}[i] = \sum_i x[i] = 1) \\
&\geq -(f_t(x_{t-1}) - f_t(x)) - \text{err}(x^*, f_t) \\
&= |f_t(x_{t-1}) - f_t(x)| - \text{err}(x^*, f_t).
\end{aligned}$$

The case where $\text{sign}(d_t) = \text{sign}(\widehat{d}_t) \geq 0$ is analogous. The claim follows. \blacksquare

Lemma 4.5. *In each round t where $|\widehat{d}_t| \geq T$ and $A_t \leq T/2$ we have*

$$\Psi_{t-1} - \Psi_t \geq \eta \left(\frac{T}{2} - \text{err}(x^*, f_t) \right) - \eta^2. \quad (10)$$

Proof. By assumption and Lemma 4.4, we have $\langle r_t, x^* - x_{t-1} \rangle \geq |f_t(x_{t-1}) - f_t(x)| - \text{err}(x^*, f_t)$. We then get from Lemma 4.3,

$$\begin{aligned}
\Psi_{t-1} - \Psi_t &\geq \eta \langle r_t, x^* - x_{t-1} \rangle - \eta^2 \\
&\geq \eta |f_t(x_{t-1}) - f_t(x)| - \eta^2 - \eta \cdot \text{err}(x^*, f_t) \\
&= \eta |d_t| - \eta^2 - \eta \cdot \text{err}(x^*, f_t).
\end{aligned}$$

On the other hand, since $|\widehat{d}_t| \geq T$ and $A_t \leq T/2$, we have that $|d_t| \geq T/2$. This proves the claim. \blacksquare

We are now ready to prove our main lemma about utility.

Lemma 4.6 (Utility for $V = U$). *When the PMW mechanism is run with $V = U$, it is an (α, β, k) -accurate interactive mechanism, where*

$$\alpha = O \left(\frac{\sqrt{\log(k/\beta) \log(1/\delta)} \log^{1/4} N}{\sqrt{\epsilon n}} \right)$$

Proof. For $V = U$, we may choose $x^* = x$ so that $\text{err}(f_t) = 0$ for all $t \in [k]$. Furthermore, with all but β probability over the algorithm's coins, the event $A_t \leq T/2$ occurs for every round $t \in [k]$. Hence, by Lemma 4.5 and $T \geq 4\eta$, the potential drop in every update round is at least

$$\Psi_{t-1} - \Psi_t \geq \eta \frac{T}{2} - \eta^2 \geq \eta^2.$$

Since $\Psi_0 \leq \log N$, the number of update rounds is bounded by $\eta^{-2} \log N$. Hence, by our termination criterion, the algorithm terminates after having answered all k queries. Furthermore, the error of the algorithm is never larger than

$$T + |A_t| \leq 2T = O \left(\frac{\sqrt{\log(k/\beta) \log(1/\delta)} \log^{1/4} N}{\sqrt{\epsilon n}} \right). \quad \blacksquare$$

We now give a utility analysis in the general case where we are working with a smaller universe $V \subseteq U$. This will be used (in Section 6) to prove the utility guarantee of Theorem 1.4. The proof is analogous that the previous one except for minor modifications.

Lemma 4.7 (Utility when $V \subsetneq U$). *Let f_1, f_2, \dots, f_k denote a sequence of k linear queries. Take*

$$\gamma = \inf_{x^*} \sup_{t \in [k]} \text{err}(x^*, f_t)$$

where x^* ranges over all histograms supported on V . When the PMW mechanism is run with V on the query sequence above, and with threshold parameter $T' = T + \gamma$, it is an (α, β, k) -non adaptively accurate interactive mechanism, where

$$\alpha = O\left(\gamma + \frac{\sqrt{\log(k/\beta) \log(1/\delta) \log^{1/4} N}}{\sqrt{\epsilon n}}\right).$$

Proof. To prove the lemma, we choose x^* as a minimizer in the definition of γ . With this choice, Lemma 4.5 implies that

$$\Psi_{t-1} - \Psi_t \geq \eta \left(\frac{T'}{2} - \gamma \right) - \eta^2 \geq \eta^2,$$

since we chose $T' \geq 4\eta + \gamma$. The argument is now the same as before. In particular, the error is bounded by $O(T') = O(\gamma + T)$ which is what we claimed. ■

4.2 Privacy analysis

Our goal in this section is to demonstrate that the interactive mechanism satisfies (ϵ, δ) -differential privacy (see Definition 3.6). We assume that all parameters such as V , σ , η , and T are publicly known. They pose no privacy threat as they do not depend on the input database. For ease of notation we will assume that $V = U$ throughout this section. The proof is the same for $V \subseteq U$ (the sub-universe V is always public information).

Simplifying the transcript. Without loss of generality, we can simplify the output of our mechanism (and hence the transcript between adversary and mechanism). We claim that the output transcript of the mechanism is determined by the following (random) vector \mathbf{v} . In particular, it is sufficient to argue that \mathbf{v} is differentially private. For every round t , the t -th entry in \mathbf{v} is defined as

$$\mathbf{v}_t = \begin{cases} \perp & \text{if } w_t = 0 \\ \widehat{a}_t & \text{if } w_t = 1 \end{cases}.$$

In other words, \mathbf{v}_t is equal to \perp if that round was a lazy round, or the noisy answer $\widehat{a}_t = f_t(x) + A_t$ if round t was an update round. This is sufficient information for reconstructing the algorithm's output: given the prefix $\mathbf{v}_{<t} = (\mathbf{v}_1, \dots, \mathbf{v}_{t-1})$, we can compute the current histogram x_{t-1} for the beginning of round t . For the lazy rounds, this is sufficient information for generating the algorithm's output. For the update rounds, $\mathbf{v}_t = \widehat{a}_t$, which is the output for round t . It is also sufficient information for re-weighting and computing the new x_t .

Note that to argue differential privacy, we need to prove that the entire transcript, including the queries of the adversary, is differentially private. Without loss of generality, we may assume that the adversary is deterministic.⁴ In this case f_t is determined by $\mathbf{v}_{<t}$. Hence, there is no need to include f_t explicitly in our transcript. It suffices to show that the vector \mathbf{v} is (ε, δ) -differentially private.

Lemma 4.8 (Privacy). *The PMW mechanism satisfies (ε, δ) -differential privacy.*

Proof. Fix an adversary and histograms $x, x' \in \mathbb{R}^N$ so that $\|x - x'\|_1 \leq 1/n$. Let $\varepsilon_0 = 1/\sigma n$ (where σ is the scaling of the Laplacian variables used in our algorithm).

Let P denote the output distribution of our mechanism when run on the input database x and similarly let Q denote the output of our mechanism when run on x' . Both distributions are supported on $\mathcal{S} = (\{\perp\} \cup \mathbb{R})^k$. For $v \in \mathcal{S}$, we define the loss function $L: \mathcal{S} \rightarrow \mathbb{R}$ as

$$L(v) := \log \left(\frac{P(v)}{Q(v)} \right). \quad (11)$$

Here and in the following we identify P with its probability density function dP (which exists by the Radon-Nikodym theorem). Henceforth $P(v)$ denotes the density of P at v .

We will then show that

$$\mathbb{P}_{\mathbf{v} \sim P} \{L(\mathbf{v}) \leq \varepsilon\} \geq 1 - \delta. \quad (12)$$

By Lemma 3.7, inequality (12) implies (ε, δ) -differential privacy and hence our claim.

Fix a transcript $v \in \mathcal{S}$ we will now proceed to analyze $L(v)$. Using the chain rule for conditional probabilities, let us rewrite $L(v)$ as

$$L(v) = \log \left(\frac{P(v)}{Q(v)} \right) = \sum_{t \in [k]} \log \left(\frac{P_t(v_t | v_{<t})}{Q_t(v_t | v_{<t})} \right), \quad (13)$$

where $P_t(v_t | v_{<t})$ denotes the conditional probability (or rather conditional density) of outputting v_t in step t on input histogram x , conditioned on $v_{<t} = (v_1, \dots, v_{t-1})$. The definition of $Q_t(v_t | v_{<t})$ is analogous with x' replacing x . Note that conditioning on $v_{<t}$ is necessary, since the coordinates of the transcript are not independent. Further, it is important to note that conditioned on $v_{<t}$, the estimate x_{t-1} in the algorithm at step t is the same regardless of whether we started from x or x' .

Borderline event. We define an event $S_t = S(v_{<t}) \subseteq \mathbb{R}$ on the noise values as follows. Let $d_t = f_t(x_{t-1}) - f_t(x)$. Note that x_{t-1} depends on $v_{<t}$ and therefore S_t will depend on it as well. We want S_t to satisfy the following properties (formally stated in Claims 4.9–4.11):

1. $\mathbb{P} \left\{ |\widehat{d}_t| > T \mid A_t \in S_t, v_{<t} \right\} \geq 1/6$. In other words, conditioned on S_t , with probability at least $1/6$ round t is a update round.
2. Conditioned on S_t not occurring, the distribution of v_t under x is *identical* to the distribution of v_t under x' , and the privacy loss is 0.

⁴We can think of a randomized adversary as a collection of deterministic adversaries one for each fixing of the adversary's randomness (which is independent of our algorithm's coin tosses).

3. Conditioned on S_t , the t -th entry v_t is ε_0 differentially private.

We will define S_t so it contains all of the noise values A_t where $|\widehat{d}_t| = |d_t + A_t|$ is “close to” (within distance σ) or larger than T . This will achieve all three of the above properties. Formally, we construct $S_t = S^+ \cup S^-$ to be made up of two intervals of noise values

$$S^- = (-\infty, -T - d_t + \sigma] \quad \text{and} \quad S^+ = [T - d_t - \sigma, \infty).$$

Note that, since $T > 2\sigma$, these two intervals never intersect. The following claims show that all three properties hold:

Claim 4.9 (Property 1). $\mathbb{P}\Pr\left\{|\widehat{d}_t| \geq T \mid A_t \in S_t, v_{<t}\right\} \geq 1/6$.

Proof. Recall that $S^+ = [T - d_t - \sigma, \infty]$. Since A_t is a Laplace random variable with magnitude σ , we get that

$$\mathbb{P}\Pr\{A_t \geq T - d_t \mid A_t \in S^+, v_{<t}\} = \mathbb{P}\Pr\{\text{Lap}(\sigma) \geq \sigma\} = 1/2e \geq 1/6.$$

And similarly, $\mathbb{P}\Pr\{A_t \leq -T - d_t \mid A_t \in S^-, v_{<t}\} \geq 1/6$. Since $|\widehat{d}_t| \geq T$ iff $A_t \geq T - d_t$ or $A_t \leq -T - d_t$, we conclude that $\mathbb{P}\Pr\left\{|\widehat{d}_t| \geq T \mid A_t \in S_t, v_{<t}\right\} \geq 1/6$. ■

Claim 4.10 (Property 2). For every $a \in \mathbb{R} \cup \{\perp\}$:

$$\log\left(\frac{P_t(a \mid A_t \notin S_t, v_{<t})}{Q_t(a \mid A_t \notin S_t, v_{<t})}\right) = 0.$$

Proof. When $A_t \notin S_t$, we know that $-T - d_t + \sigma \leq A_t \leq T - d_t - \sigma$. In particular, this means that (conditioned on S_t not occurring), v_t is always \perp , both on x and on x' . ■

Claim 4.11 (Property 3). For every $a \in \mathbb{R} \cup \{\perp\}$:

$$\log\left(\frac{P_t(a \mid A_t \in S_t, v_{<t})}{Q_t(a \mid A_t \in S_t, v_{<t})}\right) \leq 2\varepsilon_0.$$

Proof. Since A_t is a Laplace variable of scale σ , for any $a \in \mathbb{R}$ either its probability by both P_t and Q_t is 0, or otherwise its probabilities by x and x' differ by an $e^{1/\sigma n} = e^{\varepsilon_0}$ ratio. Similarly, we can bound the ratio between the probability of \perp by P and by Q . Note that

$$P_t(\perp \mid A_t, v_{<t}) = \mathbb{P}\Pr\{A_t + d_t \in (-T, -T + \sigma] \cup [T - \sigma, T)\},$$

while

$$Q_t(\perp \mid A_t, v_{<t}) = \mathbb{P}\Pr\{A_t + d'_t \in (-T, -T + \sigma] \cup [T - \sigma, T)\},$$

where $|d_t - d'_t| \leq 1/n$. Since A_t is a Laplacian variable of scale σ , it follows that the ratio of the two probabilities on the RHS is bounded by $e^{1/\sigma n} = e^{\varepsilon_0}$. ■

Bounding the Expectation. We will now bound the expected loss $\mathbb{E}[L(\mathbf{v})]$ for a random choice of \mathbf{v} sampled according to P . Applying Lemma 3.3 to Claim 4.11, we get that

$$\mathbb{E} \left[\log \left(\frac{P_t(\mathbf{v}_t | A_t \in S_t, \mathbf{v}_{<t})}{Q_t(\mathbf{v}_t | A_t \in S_t, \mathbf{v}_{<t})} \right) \right] \leq 8\varepsilon_0^2. \quad (14)$$

On the other hand, we have by Claim 4.10,

$$\mathbb{E} \left[\log \left(\frac{P_t(\mathbf{v}_t | A_t \notin S_t, \mathbf{v}_{<t})}{Q_t(\mathbf{v}_t | A_t \notin S_t, \mathbf{v}_{<t})} \right) \right] = 0. \quad (15)$$

We can express $P_t(\mathbf{v}_t | \mathbf{v}_{<t})$ as a convex combination in the form

$$P_t(\mathbf{v}_t | \mathbf{v}_{<t}) = \mathbb{P}\{A_t \in S_t | \mathbf{v}_{<t}\} P_t(\mathbf{v}_t | A_t \in S_t, \mathbf{v}_{<t}) + \mathbb{P}\{A_t \notin S_t | \mathbf{v}_{<t}\} P_t(\mathbf{v}_t | A_t \notin S_t, \mathbf{v}_{<t}),$$

and we can express $Q_t(\mathbf{v}_t | \mathbf{v}_{<t})$ similarly. We can then apply Lemma 3.2 (convexity of relative entropy) to conclude that

$$\mathbb{E} \left[\log \left(\frac{P_t(\mathbf{v}_t | \mathbf{v}_{<t})}{Q_t(\mathbf{v}_t | \mathbf{v}_{<t})} \right) \right] \leq 8\varepsilon_0^2 \mathbb{E}[\mathbb{P}\{A_t \in S_t | \mathbf{v}_{<t}\}]. \quad (16)$$

We conclude that

$$\begin{aligned} \mathbb{E} L(v) &= \sum_{t=1}^k \mathbb{E} \left[\log \left(\frac{P_t(\mathbf{v}_t | \mathbf{v}_{<t})}{Q_t(\mathbf{v}_t | \mathbf{v}_{<t})} \right) \right] \\ &\leq 8\varepsilon_0^2 \mathbb{E} \left[\sum_{t=1}^k \mathbb{P}\{A_t \in S_t | \mathbf{v}_{<t}\} \right] \\ &\leq 48\varepsilon_0^2 m \leq \varepsilon/2. \end{aligned} \quad (17)$$

Here we used that $\mathbb{E}[\sum_t \mathbb{P}\{A_t \in S_t | \mathbf{v}_{<t}\}]$ is just the expected number of borderline rounds which has to be bounded by $6m$ since every borderline round is an update round with probability at least $1/6$ and there are at most m update rounds.

Number of Borderline Rounds. With overwhelming probability, the number m' of borderline rounds (rounds t where S_t occurs) is not much larger than m (the bound on the number of update rounds). This is because every borderline round is with probability at least $1/6$ a update round (Claim 4.9). This is made formal in the claim below.

Claim 4.12. $\mathbb{P}\{m' > 32m \log^{1/2}(1/\delta)\} \leq \delta/2$

Proof. We have already argued that $\mathbb{E}[m'] \leq 6m$. Moreover, the noise in each round is independent from previous rounds. Hence, by tail bounds for Bernoulli variables, the event $m' > 32\sqrt{\log(1/\delta)m}$ has probability less than $\exp(-\log(2/\delta))$. ■

Putting it together. Condition on there being at most $m' = 32m \log^{1/2}(1/\delta)$ borderline rounds (this is the case with all but $\delta/2$ probability). We proceed by an “evolution of confidence argument” similar to [DN03, DN04].

Specifically, we will apply Azuma’s inequality to the set of m' borderline rounds. Formally, let $B \subseteq [k]$ denote the set of borderline rounds. For each $t \in B$, we view

$$X_t = \log \left(\frac{P_t(\mathbf{v}_t \mid \mathbf{v}_{<t})}{Q_t(\mathbf{v}_t \mid \mathbf{v}_{<t})} \right)$$

as a random variable. Note that $L(\mathbf{v}) = \sum_{t \in B} X_t$. Further $|X_t| \leq 2\varepsilon_0$ by Claim 4.11. Hence, by Azuma’s inequality (Lemma 3.4),

$$\Pr \{|L(\mathbf{v})| > \varepsilon\} \leq 2 \Pr \left\{ L(\mathbf{v}) > \mathbb{E}L(\mathbf{v}) + \frac{\varepsilon}{2} \right\} \leq 2 \exp \left(-\frac{\varepsilon^2}{8m' \cdot \varepsilon_0^2} \right).$$

On the other hand, by (7),

$$\frac{\varepsilon^2}{8m' \cdot \varepsilon_0^2} \geq \frac{\varepsilon^2 \sigma^2 n^2}{m'} \geq 100 \log(1/\delta) \frac{m}{m'}.$$

So, conditioning on having at most m' borderline rounds (occurs with all but $\delta/2$ probability), with all but $\delta/2$ probability the loss $L(\mathbf{v})$ deviates by at most $\varepsilon/2$ from its expectation. The expectation itself is at most $\varepsilon/2$ by (17). We conclude that with all but δ probability, the total loss $L(\mathbf{v})$ is bounded by ε in magnitude. ■

5 Achieving $(\varepsilon, 0)$ -differential privacy

Our previous mechanism satisfies (ε, δ) -differential privacy. We can achieve $(\varepsilon, 0)$ -differential privacy (or ε -differential privacy in short) by going from error $n^{-1/2}$ to error $n^{-1/3}$ (in terms of n).

Modifications to PMW. We will need to modify our algorithm in two regards. Specifically, instead of the parameter setting in (6) we use the setting

$$\eta = \frac{\log^{1/3} M \cdot \log^{1/3}(k/2\beta)}{\varepsilon^{1/3} n^{1/3}} \quad \sigma = \frac{10\eta}{\log(k/2\beta)} \quad T = 40\eta. \quad (18)$$

Furthermore, in step (2) of PMW we replace the threshold T by a randomized threshold $\widehat{T} = T + \text{Lap}(\sigma_T)$ where $\sigma_T = 10/n\varepsilon$. Our algorithm remains unchanged otherwise.

With these two modifications we can prove the next theorem.

Theorem 5.1. *Let U be a data universe of size N . For any $k, \varepsilon, \beta > 0$, there is an $(\varepsilon, 0)$ -differentially private interactive mechanism which is (α, β, k) -accurate for (adaptive) counting queries over U and data bases of size n , where*

$$\alpha = O \left(\frac{\log(k/\beta)^{1/3} \log^{1/3} N}{(\varepsilon n)^{1/3}} \right).$$

The running time in answering each query is $N \cdot \text{poly}(n) \cdot \text{polylog}(1/\beta, 1/\varepsilon, 1/\delta)$.

Proof. The algorithm stated in the theorem is given by PMW with the modifications described above. Letting $m = \eta^{-2} \log N$, we note that this setting of parameters in (18) satisfies the two properties

$$\sigma n \geq \frac{10m}{\varepsilon} \quad \text{and} \quad T \geq 4\sigma \log(k/\beta).$$

Using the second property, we can repeat the utility analysis verbatim to argue that there at most $\eta^{-2} \log M$ update rounds. Hence, with probability $1 - \beta/2$, the algorithm answers all k queries with error $\alpha = O(\widehat{T})$. Moreover it is easy to see that with probability $1 - \beta/2$, $\widehat{T} = O(T)$. Hence, with probability $1 - \beta$, we have

$$\alpha = O\left(\frac{\log(k/\beta)^{1/3} \log^{1/3} N}{(\varepsilon n)^{1/3}}\right).$$

It remains to argue that the mechanism satisfies $(\varepsilon, 0)$ -differential privacy. Fix two histograms x, x' such that $\|x - x'\|_1 \leq 1/n$. As in the proof of Lemma 4.8, we let $v \in (\mathbb{R} \cup \{\perp\})^k$ denote a transcript and we let $P(v)$ and $Q(v)$ denote the probability of this transcript when our mechanism is run on x and x' , respectively. It suffices to argue that for all transcripts v ,

$$-\varepsilon \leq \log\left(\frac{P(v)}{Q(v)}\right) \leq \varepsilon. \quad (19)$$

Let us again write

$$\log\left(\frac{P(v)}{Q(v)}\right) = \sum_{t \in [k]} \log\left(\frac{P_t(v_t | v_{<t})}{Q_t(v_t | v_{<t})}\right).$$

Further let $H = \{t \in [k]: v_t \neq \perp\}$ and $H^c = [k] \setminus H$.

Claim 5.2.

$$-\frac{\varepsilon}{10} \leq \sum_{t \in H} \log\left(\frac{P_t(v_t | v_{<t})}{Q_t(v_t | v_{<t})}\right) \leq \frac{\varepsilon}{10}. \quad (20)$$

Proof. Note that $|H| \leq m$ by the termination criterion of our algorithm. It therefore follows from standard properties of the Laplacian distribution and our choice of parameters that

$$\sum_{t \in H} \log\left(\frac{P_t(v_t | v_{<t})}{Q_t(v_t | v_{<t})}\right) \leq \frac{m}{\sigma n} \leq \frac{\varepsilon}{10}. \quad (21)$$

The same argument shows a lower bound of $-\varepsilon/10$. This concludes the proof. \blacksquare

The next lemma handles the coordinates $t \in H^c$.

Claim 5.3.

$$-\frac{\varepsilon}{10} \leq \sum_{t \in H^c} \log\left(\frac{P_t(v_t | v_{<t})}{Q_t(v_t | v_{<t})}\right) \leq \frac{\varepsilon}{10}. \quad (22)$$

Proof. Let $\mathcal{A}(x)$ be the set of values for the noise variables (A_1, \dots, A_t) which lead to the event that the transcript is \perp in all rounds $t \in H^c$ when we run our algorithm on input x and

condition the transcript on being equal to v_t in all rounds $t \in H$. Define $\mathcal{A}_Z(x)$ in the same way except that we additionally condition the algorithm on the event that $\widehat{T} = Z$. Observe that

$$\mathcal{A}_{Z-1/n}(x') \subseteq \mathcal{A}_Z(x) \subseteq \mathcal{A}_{Z+1/n}(x'). \quad (23)$$

Here we used the assumption $\|x - x'\|_1 \leq 1/n$ and thus $|f(x) - f(x')| \leq 1/n$ for any possibly queries f .

Further observe that, by the product rule for conditional probabilities,

$$\begin{aligned} \prod_{t \in H^c} P_t(v_t | v_{<t}) &= \mathbb{P}\{(A_1, \dots, A_k) \in \mathcal{A}(x)\} \\ &= \int_{-\infty}^{\infty} \mathbb{P}\{\widehat{T} = Z\} \mathbb{P}\{(A_1, \dots, A_k) \in \mathcal{A}_Z(x)\} dZ. \end{aligned}$$

The first step follows from the definition of $\mathcal{A}_Z(x)$ and the second step uses independence between the random variables \widehat{T} and (A_1, \dots, A_k) . On the other hand,

$$\mathbb{P}\{\widehat{T} = Z\} \leq e^{\varepsilon/10} \mathbb{P}\{\widehat{T} = Z + 1/n\},$$

and

$$\mathbb{P}\{\widehat{T} = Z - 1/n\} \geq e^{-\varepsilon/10} \mathbb{P}\{\widehat{T} = Z - 1/n\}.$$

Therefore,

$$\begin{aligned} \prod_{t \in H^c} P_t(v_t | v_{<t}) &= \int_{-\infty}^{\infty} \mathbb{P}\{\widehat{T} = Z\} \mathbb{P}\{(A_1, \dots, A_k) \in \mathcal{A}_Z(x)\} dZ \\ &\leq e^{\varepsilon/10} \int_{-\infty}^{\infty} \mathbb{P}\{\widehat{T} = Z + 1/n\} \mathbb{P}\{(A_1, \dots, A_k) \in \mathcal{A}_{Z+1/n}(x')\} dZ \quad (\text{using (23)}) \\ &= e^{\varepsilon/10} \int_{-\infty}^{\infty} \mathbb{P}\{\widehat{T} = Z\} \mathbb{P}\{(A_1, \dots, A_k) \in \mathcal{A}_Z(x')\} dZ \\ &= e^{\varepsilon/10} \prod_{t \in H^c} Q_t(v_t | v_{<t}). \end{aligned} \quad (24)$$

Using the same reasoning we get

$$\prod_{t \in H^c} P_t(v_t | v_{<t}) \geq e^{-\varepsilon/10} \prod_{t \in H^c} Q_t(v_t | v_{<t}). \quad (25)$$

Taking logarithms on both sides of (24) and (25) shows that (22) holds which is what we wanted to show. \blacksquare

Putting together (20) and (22), it follows that the sum over all $t \in [k]$ is in the interval $[-\varepsilon/5, \varepsilon/5]$. This establishes the bound stated in (19). We conclude that the algorithm satisfies $(\varepsilon, 0)$ -differential privacy. The theorem follows. \blacksquare

5.1 Lower bound for $(\epsilon, 0)$ -differential privacy

As discussed before, there is a lower bound of roughly $\sqrt{\log(k)/n}$ that holds for blatant non-privacy [DN03]. A shortcoming is that this lower bound does not depend on the universe size. In this section we will show a lower bound on the accuracy of any mechanism that satisfies $(\epsilon, 0)$ -differential privacy even if it works in the non-interactive setting.

Theorem 5.4. *Let n be sufficiently large and let $\epsilon > 0$ be a constant independent of n . Then, for every $k \geq n^{1.1}$ there is a set of k linear queries over a universe of size N such that every $(\epsilon, 0)$ -differentially private mechanism for databases of size n must have error*

$$\alpha \geq \Omega(1) \cdot \left(\frac{\log k \cdot \log\left(\frac{N}{n}\right)}{\epsilon n} \right)^{1/2}$$

with probability $1/2$.

Proof. Let U be a universe of size N . Our lower bound uses a discrete variant of the packing argument from [HT10].

Consider the family $\mathcal{X} \subseteq (\frac{1}{n}\mathbb{Z}_+)^N$ of all histograms with exactly s nonzeros. Note that such histograms correspond to databases of size n with s distinct elements. Here, s is some parameter that we will fix shortly. Since we normalize histograms to have norm 1, this means that each nonzero coordinate is $1/s$. Further let \mathcal{F} be the uniform distribution over linear queries of the form $f \in \{0, 1\}^N$.

We say that two histograms $x, y \in \mathcal{X}$ are *half disjoint* if $\|x - y\|_1 \leq 1/2$. The next claim shows that two randomly chosen elements from \mathcal{X} are very likely half disjoint.

Claim 5.5. *The probability that $x, y \sim \mathcal{X}$ are not half disjoint is at most $\exp(-\Omega(s \cdot \log(N/n)))$.*

Proof. Note that for $x, y \sim \mathcal{X}$ not to be half disjoint it must be the case that half the nonzero coordinates of y must fall into the s nonzero coordinates of x . The probability of this event is less than

$$\binom{s}{s/2} \cdot \left(\frac{s}{N}\right)^{s/2} \leq 2^{-\Omega(s \log(N/s))}.$$

Here we used that $s \leq n$. ■

We also need to show that any two half disjoint histograms x, y are “well separated” by k random linear queries.

Claim 5.6. *Let x, y be half disjoint. Then, for k queries f_1, \dots, f_k chosen uniformly at random from \mathcal{F} , we have*

$$\mathbb{Pr} \left\{ \max_{t \in [k]} |f_t(x) - f_t(y)| \leq \frac{c}{\sqrt{s}} \right\} \leq \exp(-2^{-\Omega(c^2)} k). \quad (26)$$

Proof. If we choose a single query f at random, then $|f(x) - f(y)|$ is expected to be $\Theta(1/\sqrt{s})$. Further, $|f(x) - f(y)| > c/\sqrt{s}$ with probability at least $2^{-\Omega(c^2)}$. This follows from standard lower bounds on the tail of the binomial distribution. The probability that none out of k random queries has difference c/\sqrt{s} is therefore bounded by $(1 - 2^{-\Omega(c^2)})^k$ which implies the claim. ■

We will now put the previous two claims together. To this end, fix $s = C\epsilon n/\log(N)$ for sufficiently large $C > 0$ and put

$$\alpha_0 = \sqrt{\frac{\log k}{s}} = \sqrt{\frac{\log k \cdot \log N}{\epsilon n}}. \quad (27)$$

It then follows from Claim 5.6 that

$$\mathbb{P}\left\{\max_{t \in [k]} |f_t(x) - f_t(y)| \leq \alpha_0\right\} \leq \exp(-k^{0.99}). \quad (28)$$

On the other hand, using Claim 5.5, it follows that there exists a set $\mathcal{P} \subseteq \mathcal{X}$ such that every pair $x, y \in \mathcal{P}$ with $x \neq y$ is half disjoint and

$$|\mathcal{P}| \geq \exp(\Omega(s \log(N/n))) \geq 3\exp(\epsilon n). \quad (29)$$

In the second inequality we used the fact that we can choose C sufficiently large in the setting of s above.

Further, by our assumption $k \geq n^{1.1}$, we have $\exp(-k^{0.99}) \ll |\mathcal{P}|^{-2}$. Hence, we can take the union bound over all distinct pairs in \mathcal{P} and conclude that there must exist a set of k linear queries f_1, \dots, f_k such that for every two histograms $x, y \in \mathcal{P}$ with $x \neq y$ we have

$$\max_{t \in [k]} |f_t(x) - f_t(y)| \geq \alpha_0.$$

Now, for the sake of contradiction suppose there is an $(\epsilon, 0)$ -differentially private mechanism M for answering the k linear queries f_1, \dots, f_k which has maximum error $\alpha = \alpha_0/2$ with probability $1/2$.

For a histogram x let Fx denote the vector $(f_1(x), \dots, f_k(x)) \in \mathbb{R}^k$ and let $B(Fx, \alpha)$ denote the ℓ_∞ -ball around Fx of radius α . Note that by the accuracy guarantee of M , we have for every histogram x ,

$$\mathbb{P}\{M(x) \in B(Fx, \alpha)\} \geq \frac{1}{2}.$$

By $(\epsilon, 0)$ -differential privacy, we further have for every two histograms x, y ,

$$\mathbb{P}\{M(y) \in B(Fx, \alpha)\} \geq \frac{\exp(-\epsilon n)}{2}.$$

Fix any histogram $x \in \mathcal{P}$, then

$$\begin{aligned} 1 &\geq \mathbb{P}\left\{M(x) \in \bigcup_{y \in \mathcal{P}} B(Fy, \alpha)\right\} = \sum_{y \in \mathcal{P}} \mathbb{P}\{M(x) \in B(Fy, \alpha)\} \\ &\geq |\mathcal{P}| \frac{\exp(-\epsilon n)}{2}. \end{aligned}$$

Here we used that for $y, y' \in \mathcal{P}$ with $y \neq y'$ we have $B(Fy, \alpha) \cap B(Fy', \alpha) = \emptyset$. But $|\mathcal{P}| \geq 3\exp(\epsilon n)$. Hence we have arrived at a contradiction showing that such a mechanism M cannot exist. ■

Remark 5.7. In our proof we used databases in which individual data items occur with high multiplicity. This is not necessary as we can always move to a universe of size nN in which every data item occurs n times. Hence, we can repeat the same construction without multiplicities by losing only a factor of n in the universe size.

We leave it as an open problem to close the gap between our upper and lower bound. Indeed, it would be interesting to know if the optimal dependence on n is $n^{-1/3}$ or rather $n^{-1/2}$. We also point out the open problem of coming up with a lower bound for (ϵ, δ) -differential privacy, such as $\Omega(\log^{c_1} k \cdot \log^{c_2} N \cdot n^{-1/2})$, that simultaneously has a (poly-logarithmic) dependence on the universe size N and the number of queries k .

6 Average-case complexity and smooth instances

In this section, we define a notion of average case complexity for interactive (and non-interactive) mechanisms that allows us to improve the running time of the PMW mechanism as a function of the data universe size. This is done using an argument for reducing the data universe size.

We start by defining the notion of a *smooth* histogram. We think of these histograms as distributions over the data universe that do not place too much weight on any given data item. In other words, we require the histogram to have high min-entropy.

Definition 6.1 (Smooth). A histogram $x \in \mathbb{R}^U$ s.t. $\sum_{u \in U} x_u = 1$ and $\forall u \in U : x_u \geq 0$ is ξ -smooth if $\forall u \in U : x_u \leq \xi$.

In particular, a ξ -smooth histogram has *min-entropy* at least $\log(1/\xi)$. We typically think of ξ has a function of N , such as $\text{polylog}N/N$ or $1/\sqrt{N}$. Note that small databases (viewed as histograms) cannot be very smooth, since a ξ -smooth histogram has at least $1/\xi$ nonzero coordinates.

We therefore extend the notion of smoothness to the notion of pseudo-smoothness with respect to a set of queries \mathcal{Q} . A histogram is *pseudo-smooth* w.r.t a query class \mathcal{Q} roughly speaking when there exists a smooth histogram x^* that is close on every query in \mathcal{Q} . This notion allows even very sparse histograms (corresponding to small databases) to be very pseudo-smooth. The formal definition is as follows.

Definition 6.2 (Pseudo-smooth). A histogram $x \in \mathbb{R}^U$ s.t. $\sum_{u \in U} x_u = 1$ and $\forall u \in U : x_u \geq 0$ is (ξ, ϕ) -smooth w.r.t a class of linear queries \mathcal{Q} if there *exists* a ξ -smooth histogram x^* s.t.

$$\forall f \in \mathcal{Q}: \quad |f(x) - f(x^*)| \leq \phi.$$

A straightforward way of obtaining pseudo-smooth databases is by sampling from a smooth histogram.

Claim 6.3. Let U be a data universe, \mathcal{Q} a class of linear queries over U , and x^* a ξ -smooth histogram over U . For any $\alpha, \beta > 0$, sample a database x of $m = (\log(2/\beta) + \log|\mathcal{Q}|)/\alpha^2$ items i.i.d from the distribution of x^* (i.e. in each sample we independently pick each $u \in U$ with probability x_u^*). Then with all but β probability over the samples taken, $\forall f \in \mathcal{Q}: |f(x) - f(x^*)| \leq \alpha$, and so the database x is (ξ, α) -pseudosmooth w.r.t \mathcal{Q} .

Proof. The proof is by a Chernoff bound (as in [DNR⁺09]). ■

6.1 Domain reduction for pseudosmooth histograms

For a given smoothness parameter ξ , data universe U , and query class \mathcal{Q} , let $V \subseteq U$ be a sub-universe sampled uniformly and at random from U . In this section we show that (as long as V is large enough) if x was a pseudosmooth histogram over U w.r.t a query class \mathcal{Q} , then w.h.p. there will be a histogram x^* with support only over (the smaller) V that is “close” to x on \mathcal{Q} . We emphasize that sampling the sub-universe V does not require knowing x nor knowing any x^* that certifies x being pseudosmooth, we only need to know ξ . In particular, this approach is privacy-preserving. This technique for reducing the universe size can be used to improve the efficiency of the PMW mechanism for pseudosmooth input databases.

Lemma 6.4. *Let U be a data universe and \mathcal{Q} a collection of linear queries over U . Let x be (ξ, ϕ) -psuedo-smooth w.r.t \mathcal{Q} . Take $\alpha, \beta > 0$, and sample uniformly at random (with replacement) $V \subseteq U$ so that*

$$M = |V| = 4 \max\{\xi N \cdot (\log(1/\beta) + \log|\mathcal{Q}|)/\alpha^2, \log(1/\beta)\} \quad (30)$$

Then, with all but β probability over the choice of V , there exists a histogram x^ with support only over V such that*

$$\forall f \in \mathcal{Q} : |f(x) - f(x^*)| \leq \phi + \alpha. \quad (31)$$

Proof. Let y be the ξ -smooth histogram which shows that x is (ξ, ϕ) -pseudosmooth. If we sampled uniformly at random from x or from y then by Claim 6.3, we could get a database over a very small sub-universe that is (as required) close to x on all the queries in \mathcal{Q} . This is insufficient because we want the sub-universe that we find to be *independent* of the database x (and so also independent of y).

Still, let us re-examine the idea of sampling from y . One way of doing this is by *rejection sampling*. Namely, repeatedly sample $u \in U$ uniformly at random and then “keep” u with probability y_u/ξ . Otherwise reject. When we use this rejection sampling, since y is a ξ -smooth distribution, each sample that we keep is distributed by y (i.e. it is $u \in U$ w.p. y_u). Repeat this process until $m_1 = (\log(2/\beta) + \log|\mathcal{Q}|)/\alpha^2$ samples have been accepted. There is now a set of coordinates $V_1 \subseteq U$, those that were kept (of size at most m_1), and a set of coordinates $V_2 \subseteq U$, those that were rejected. By Claim 6.3 the sub-universe V_1 of samples that we keep (which are i.i.d samples from y) supports (except with probability $\beta/2$) a database x^* that is “close” to y (w.r.t \mathcal{Q}), and so it will also be “close” to x . In particular, by triangle inequality,

$$\max_{f \in \mathcal{Q}} |f(x) - f(x^*)| \leq \max_{f \in \mathcal{Q}} |f(x) - f(y)| + \max_{f \in \mathcal{Q}} |f(y) - f(x^*)| \leq \phi + \alpha.$$

But now we may take $V = V_1 \cup V_2$. Note that V is simply a uniformly random subset of the coordinates of U . And by the previous argument, V supports a histogram that satisfies (31), namely x^* . To conclude the proof it remains to argue that V has the required size. Note that the probability of accepting sample i in the rejection procedure is given by $\sum_{i=1}^N \frac{1}{N} \cdot \frac{y_i}{\xi} = 1/\xi N$. Hence, the expected number of queries in total is $\mu = 2\xi N \cdot (\log(2/\beta) + \log|\mathcal{Q}|)/\alpha^2$. Moreover, since every sample is independent, we have concentration around the expectation. A multiplicative Chernoff bound shows that the probability that V is larger than twice its expectation is bounded by $\exp(-\mu) \leq \beta/2$. ■

Finally, we use Lemma 6.4 together with Lemma 4.7 (utility of PMW for general V), to derive the accuracy guarantee of Theorem 1.4 for the performance of the PMW mechanism on pseudo-smooth databases.

Proof of Utility for Theorem 1.4. We run PMW on a uniformly chosen sub-universe V of the appropriate size M as stated in Equation (30) above, taking $\alpha = 1/\sqrt{n}$. We conclude that with all but $\beta/2$ probability over the sampling, there exists a database x^* supported on V that is $\phi + 1/\sqrt{n}$ -close to x w.r.t. the given sequence of k counting queries. Plugging this into Lemma 4.7, we obtain the accuracy bound claimed in Theorem 1.4. ■

7 Acknowledgements

We thank Boaz Barak, Cynthia Dwork, Moni Naor, Aaron Roth, Rob Schapire and Salil Vadhan for their helpful and insightful comments. We thank Aaron Roth and Jonathan Ullman for suggesting an improved parameter setting in our algorithm that lead to a $\sqrt{\log k}$ -factor improvement in our upper bound. Thanks also to the anonymous FOCS 2010 reviewers for their helpful comments.

References

- [AHK05] Sanjeev Arora, Elad Hazan, and Satyen Kale. The multiplicative weights update method: a meta algorithm and applications. Technical report, Princeton University, 2005.
- [BLR08] Avrim Blum, Katrina Ligett, and Aaron Roth. A learning theory approach to non-interactive database privacy. In *Proc. 40th Symposium on Theory of Computing (STOC)*, pages 609–618. ACM, 2008.
- [DL09] Cynthia Dwork and Jing Lei. Differential privacy and robust statistics. In *Proc. 41st Symposium on Theory of Computing (STOC)*, pages 371–380. ACM, 2009.
- [DMNS06] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Proc. 3rd TCC*, pages 265–284. Springer, 2006.
- [DMT07] Cynthia Dwork, Frank McSherry, and Kunal Talwar. The price of privacy and the limits of LP decoding. In *Proc. 39th Symposium on Theory of Computing (STOC)*, pages 85–94. ACM, 2007.
- [DN03] Irit Dinur and Kobbi Nissim. Revealing information while preserving privacy. In *Proc. 22nd PODS*, pages 202–210. ACM, 2003.
- [DN04] Cynthia Dwork and Kobbi Nissim. Privacy-preserving datamining on vertically partitioned databases. In *Proc. 24th CRYPTO*, pages 528–544. Springer, 2004.
- [DNPR10] Cynthia Dwork, Moni Naor, Toniann Pitassi, and Guy N. Rothblum. Differential privacy under continual observation. In *Proc. 42nd Symposium on Theory of Computing (STOC)*. ACM, 2010.

- [DNR⁺09] Cynthia Dwork, Moni Naor, Omer Reingold, Guy N. Rothblum, and Salil P. Vadhan. On the complexity of differentially private data release: efficient algorithms and hardness results. In *Proc. 41st Symposium on Theory of Computing (STOC)*, pages 381–390. ACM, 2009.
- [DRV10] Cynthia Dwork, Guy N. Rothblum, and Salil Vadhan. Boosting and differential privacy. In *Proc. 51st Foundations of Computer Science (FOCS)*. IEEE, 2010.
- [DY08] Cynthia Dwork and Sergey Yekhanin. New efficient attacks on statistical disclosure control mechanisms. In *Proc. 28th CRYPTO*, pages 469–480. Springer, 2008.
- [HT10] Moritz Hardt and Kunal Talwar. On the geometry of differential privacy. In *Proc. 42nd Symposium on Theory of Computing (STOC)*. ACM, 2010.
- [LW94] Nick Littlestone and Manfred K. Warmuth. The weighted majority algorithm. *Inf. Comput.*, 108(2):212–261, 1994.
- [RR10] Aaron Roth and Tim Roughgarden. Interactive privacy via the median mechanism. In *Proc. 42nd Symposium on Theory of Computing (STOC)*, pages 765–774. ACM, 2010.
- [UV11] Jonathan Ullman and Salil P. Vadhan. Pcps and the hardness of generating private synthetic data. In *TCC*, pages 400–416. Springer, 2011.