# The Price of Privately Releasing Contingency Tables and the Spectra of Random Matrices with Correlated Rows[*]

Shiva Prasad
Kasiviswanathan
CCS-3, Los Alamos National
Laboratory
kasivisw@gmail.com

Mark Rudelson
Department of Mathematics
University of Missouri
rudelsonm@missouri.edu

Adam Smith
Department of Computer
Science and Engineering
Pennsylvania State University
asmith@cse.psu.edu

Jonathan Ullman
SEAS, Harvard University
jullman@seas.harvard.edu

## ABSTRACT

Marginal (contingency) tables are the method of choice for government agencies releasing statistical summaries of categorical data. In this paper, we derive lower bounds on how much distortion (noise) is necessary in these tables to ensure the privacy of sensitive data. We extend a line of recent work on impossibility results for private data analysis [9, 12, 13, 15] to a natural and important class of functionalities.

Consider a database consisting of $n$ rows (one per individual), each row comprising $d$ binary attributes. For any subset of $T$ attributes of size $|T| = k$, the marginal table for $T$ has $2^k$ entries; each entry counts how many times in the database a particular setting of these attributes occurs. We provide lower bounds for releasing all $\binom{d}{k}$ $k$-attribute marginal tables under several different notions of privacy.

(1) We give efficient polynomial time attacks which allow an adversary to reconstruct sensitive information given insufficiently perturbed marginal table releases. In particular, for a constant $k$, we obtain a tight bound of $\widetilde{\Omega}(\min\{\sqrt{n}, \sqrt{d^{k-1}}\})$[1] on the average distortion per entry for any mechanism that releases all $k$-attribute marginals while providing "attribute" privacy (a weak notion implied by most privacy definitions).

(2) Our reconstruction attacks require a new lower bound on the least singular value of a random matrix with correlated rows. Let $M^{(k)}$ be a matrix with $\binom{d}{k}$ rows formed by taking all possible $k$-way entry-wise products of an underlying set of $d$ random vectors from $\{0, 1\}^n$. For constant $k$, we show that the least singular value of $M^{(k)}$ is $\widetilde{\Omega}(\sqrt{d^k})$ with high probability (the same asymptotic bound as for independent rows).

(3) We obtain stronger lower bounds for marginal tables satisfying *differential privacy*. We give a lower bound of $\widetilde{\Omega}(\min\{\sqrt{n}, \sqrt{d^k}\})$, which is tight for $n = \widetilde{\Omega}(d^k)$. We extend our analysis to obtain stronger results for mechanisms that add *instance-independent* noise and weaker results when $k$ is super-constant.

## Categories and Subject Descriptors

F.2.0 [**Analysis of Algorithms and Problem Complexity**]: General

## General Terms

Algorithms, Security, Theory

## 1. INTRODUCTION

The goal of *private data analysis* is to provide global, statistical properties of a data set of sensitive information while protecting the privacy of the individuals whose records the data set contains. There is a vast body of work on this problem in statistics and computer science. However, until recently, most schemes proposed in the literature lacked rigor: typically, the schemes had either no formal privacy guarantees or ensured security only against a specific suite of attacks.

The seminal results of Dinur and Nissim [9] initiated a rigorous study of the tradeoff between privacy and utility. The notion of *differential privacy* [12] that emerged from this line of work provides rigorous guarantees even in the presence of a malicious adversary with access to arbitrary side information. Differential privacy requires, roughly, that any single individual's data have little effect on the outcome of the analysis. Recently, many techniques have been developed for designing differentially private algorithms (see [10, 11] for two recent surveys). A typical objective is to release as accurate an approximation as possible to some function $f$ evaluated on the database $D$.

A complementary line of work seeks to establish lower bounds on how much distortion (noise) is necessary for particular functions $f$. Some of these bounds apply only to differential privacy (e.g., [12, 18, 19]); other bounds rule out *any* reasonable notion of privacy by showing how to reconstruct almost all of the data $D$ given sufficiently accurate approximations to $f(D)$ [9, 13, 15]. We refer to the latter works as lower bounds for *minimal* privacy.

---

[*]The full version [20] contains proofs omitted here.

[1]The $\widetilde{\Omega}(\cdot)$ notation hides inverse polylogarithmic factors in $d, n, k$.

In this paper, we investigate lower bounds on the distortion necessary for releasing a set of *marginal contingency tables* (marginal tables, in short), under both minimal and differential privacy. A database $D$ in our setting consists of $n$ rows, each row comprising values for $d$ binary attributes $x_1, \ldots, x_d$. For any subset of $T$ attributes of size $|T| = k$, the marginal table for $T$ has $2^k$ entries; each entry counts how many times in the database a particular setting of these attributes occurs. Alternatively, we may think of the table as counting the number of rows in the database that satisfy each of the $2^k$ possible *conjunctions* on the $k$ attributes in $T$. We call a marginal table for a set of $k$ attributes a $k$-way marginal table. The $d$-attribute marginal table is the "full" contingency table for the data set.

Marginal tables are the workhorses of categorical data analysis and, in particular, of data analysis in the medical, social and behavioral sciences (e.g., clinical trials, public health studies, and education statistics). In addition to being easy to interpret, they are sufficient statistics for popular classes of probabilistic models [4]. (As a simple example: for binary data, the mean vector and covariance matrix, which capture linear dependencies among attributes, are equivalent to the set of all 2-attribute marginal tables.) Because of this, they are the format of choice for data release by government statistical bureaus [3]. However, many of the fields in which categorical data are used generate highly sensitive data. Researchers and government agencies have ethical and legal responsibilities to protect the confidentiality of the individuals whose data they collect. Consequently, the confidentiality of contingency table releases has been an active topic of research in statistics for over thirty years (see, for example, [17, 30]). Understanding the extent to which marginal tables can be released while guaranteeing a rigorous, meaningful notion of privacy is an important problem.

## 1.1 Our Contributions

Let $\mathcal{C}_k(D)$ be the set of all $k$-way marginal tables (equivalently, the frequencies of all possible $k$-way conjunctions) for a database $D \in (\{0,1\}^d)^n$. There are $\binom{d}{k}$ such tables; however, it is convenient to think of $\mathcal{C}_k(D)$ as a single real vector of length $2^k \binom{d}{k}$. We give lower bounds for simultaneously estimating all the entries of $\mathcal{C}_k(D)$ privately. As a point of reference, for constant $k$, the best-known differentially private algorithms [5, 6, 14] add an $\widetilde{O}(\min\{n, (n^2 d)^{1/3}, 2^{\sqrt{\log d}}\sqrt{nd}, \sqrt{d^k}\})$ average distortion per entry. Here, the $\widetilde{O}(\cdot)$ notation hides polylogarithmic factors in $d, n, k$. Our lower bounds match this upper bound in different respects.

**(1) Lower Bounds for Minimal Privacy:** We show that algorithms that do not sufficiently distort the marginal tables fail to satisfy a large class of "privacy" definitions. We define two violations of privacy[2]: *attribute non-privacy* and *row non-privacy*.

Each of these rules out a large class of popular definitions of privacy. Row non-privacy rules out definitions that protect an entire row of the database even given leakage of other rows; such definitions include differential privacy as well as several definitions popular in the *randomized response* literature [33, 1, 16]. Attribute non-privacy rules out any definition that guarantees the se-

crecy of a particular "sensitive" attribute even when all other attributes are known to an attacker; such definitions include $K$-anonymity [32] and its variants [23, 22, 7, 24, 34], as well as the notions ruled out by row non-privacy.

Using a "reconstruction" attack outlined below **(2)**, we show that for any constant $k$, releasing $\mathcal{C}_k(D)$ with distortion $o(\min\{\sqrt{n}, \sqrt{d^{k-1}}\})$ per entry allows an adversary to efficiently reconstruct large fraction of the sensitive attribute entries given the nonsensitive values, thus *violating* attribute privacy (the bound holds even for releasing only all those $k$-way tables that involve the sensitive attribute and $k-1$ other attributes). Moreover, releasing $\mathcal{C}_k(D)$ with distortion $o(\min\{\sqrt{n}, \sqrt{d^k}\})$ per entry allows an adversary to efficiently reconstruct large fraction of the rows of $D$, even though this would not be possible without the release, thus *violating* row privacy. Both these bounds are (almost) tight, as there is an algorithm which is neither attribute non-private nor row non-private and which for every database $D$ adds $\widetilde{O}(\min\{\sqrt{n}, \sqrt{d^k}\})$ distortion per entry of $\mathcal{C}_k(D)$. The formal bounds for these privacy notions are stated in Table 1 and discussed in Section 2.

**(2) Reconstruction Attack & the Least Singular Value of Random Matrices with Correlated Rows:** The bounds on minimal privacy **(1)** above require significantly different techniques from previous work. Previous lower bounds [9, 13, 15] were based on variants of the following reconstruction problem: given a real-valued matrix $M$, and a corrupted "codeword" $Ms + e$, the goal is to compute an approximation $\hat{s}$ to $s$ such that the "reconstruction error" $\hat{s} - s$ is somehow bounded in terms of the noise vector $e$. Typically, assuming some norm $\|e\|_p$ is small, one can bound a related norm of $\hat{s} - s$.

Returning to data privacy: if $s \in \mathbb{R}^n$ is a database with one number assigned per person, we can think of $y = Ms + e$ as a vector of (distorted) estimates of the quantities $\langle M_i, s \rangle$, where $M_i$ is the $i$th row of $M$. Any private data release that allows a user to estimate $\langle M_i, s \rangle$, allows an attacker to obtain $y$. Therefore, an algorithm for approximating $s$ from $y$ can be used to infer sensitive data from the release.

Previous lower bounds rely heavily on the freedom to design $M$ by selecting the rows of $M$ independently (either at random [9, 13, 15] or from an algebraic code [15]). When $k = 1$ a similar flexibility is available in our lower bounds; the matrix $M^{(1)}$ that arises in our lower bounds is a $\{0,1\}^{d \times n}$ matrix with independent random entries. However, for $k > 1$ the rows of the matrix $M^{(k)}$ that arises in our lower bounds are highly correlated: the matrix $M^{(k)}$ has $d^k$ rows which are formed by taking all possible $k$-way *entry-wise products*[3] of the rows of the random matrix $M^{(1)}$. The techniques of previous work, from the literature on both privacy and random matrices, break down. We show that reconstruction procedures using these matrices can in fact be analyzed, by showing for any constant $k$ that a random $(0,1)$-matrix with *correlated* rows has approximately the same least singular value as a random $(0,1)$-matrix with *independent* rows.

Tight bounds are known on the least singular values of various types of matrices (e.g., square, rectangular) with independent random entries (see, e.g., [28, 29, 27]

---

[2]Alternatively, we might call these "attribute leakage" and "row leakage". We use "non-privacy" for consistency with the previous works [9, 13, 15].

[3]The entry-wise product of $k$ vectors $p_1, \ldots, p_k \in \mathbb{R}^n$ is the vector in $q \in \mathbb{R}^n$ with entries $q_i = \prod_{j=1}^k p_{j_i}$.

and references therein). The least singular value of an $N \times n$ matrix with $(0,1)$ independent random entries and $N \geq n$ is $\Theta(\sqrt{N})$ with exponentially high probability (in fact, even non-asymptotic bounds are known, see [29]). To deal with the dependencies, we develop several new tools, which may be of independent interest. We show that for any constant $k$ if the random matrix $M^{(1)}$ has less than $d^k/\log^{k-2} n$ columns, then the least singular value of $M^{(k)}$ is $\widetilde{\Omega}(\sqrt{d^k})$ with exponentially high probability. Therefore, the least singular value of $M^{(k)}$ is asymptotically comparable to that of a $d^k \times n$ random matrix with independent entries, but $M^{(k)}$ (constructed out of $M^{(1)}$) uses far lower randomness.

The proof is challenging because correlations make powerful measure concentration tools hard to apply. We first reduce the problem to bounding the least singular value of a (related) random centered matrix $\tilde{\Pi}$. The smallest singular value of $\tilde{\Pi}$ is the minimum of $\|\tilde{\Pi}x\|$, over $x$ from the unit sphere. An important tool in the proof is bounding the *small ball probability*, which is the probability that $\|\tilde{\Pi}x\|$ is small for a *fixed* vector $x$. To obtain a uniform lower bound for $\|\tilde{\Pi}x\|$, we decompose the unit sphere into many pieces, and for each piece use epsilon-net arguments tailored according to the small ball probability. Then, we obtain a uniform lower estimate on the net, which is then extended to the whole unit sphere by approximation.

In the privacy context, our spectral lower bound allows for a reconstruction algorithm of the form

$$\hat{s} = \text{round}(M_{inv}^{(k)} \cdot (M^{(k)}s + e)),$$

where $s$ is a $(0,1)$-vector, $M_{inv}^{(k)}$ is an appropriate pseudoinverse of $M^{(k)}$ and round$(z)$ rounds the entries of a vector $z$ to the nearer of 0 and 1. We show that releasing $M^{(k)}s$ with distortion $o(\sqrt{n})$ per entry allows the adversary to reconstruct $n - o(n)$ bits of $s$ (that is, to find $\hat{s}$ that agrees in almost all entries of s), as long as $n = o(d^k)$. One can extend the result to get a lower bound of $\widetilde{\Omega}(\min\{\sqrt{n}, \sqrt{d^k}\})$ for all $n$.

**(3) Lower Bounds for Differential Privacy:** Using a disjoint set of techniques, we show a stronger lower bound for releasing $k$-way marginal tables under the notion of $(\epsilon, \delta)$-differential privacy. The precise bounds are stated in Table 1 and discussed in Section 4. Here, we treat $\epsilon$ and $\delta$ as constants.

For constant $k$, the best-known $(\epsilon, \delta)$-differentially private algorithms [5, 6, 14] yield an average distortion per entry of $\widetilde{O}(\min\{n, (n^2 d)^{1/3}, 2^{\sqrt{\log d}}\sqrt{nd}, \sqrt{d^k}\})$, while our lower bound is $\widetilde{\Omega}(\min\{\sqrt{n}, \sqrt{d^k}\})$. Our bounds imply that the technique of Blum *et al.* [5], which adds Gaussian noise to each entry in $\mathcal{C}_k$ is tight for large databases (when $n = \widetilde{\Omega}(d^k)$). Moreover, for a natural and popular class of algorithms based on adding *instance-independent* noise [5, 12, 3], we strengthen this bound to $\Omega(\sqrt{d^k})$, which is tight for all $n$.

Our lower bounds for differential privacy extend even to a non-constant $k$, and here we show a lower bound of $\widetilde{\Omega}(\min\{\sqrt{n}, \sqrt{\binom{d}{k}}\}/2^k)$ on the average distortion per entry. For $n = \widetilde{\Omega}(\binom{d}{k})$, this is loose by a factor of $\sqrt{2^k}$ when compared to the best-known upper bound. This

lower bound can again be strengthened for the instance-independent case (see Table 1).

Let $\mathcal{A}$ be a differentially private algorithm for $\mathcal{C}_k$. The rough idea behind these differential privacy lower bounds is to start with a particular database $D$ and then bound the projection of the *mean squared error*(MSE) matrix of $\mathcal{A}(D)$ along a large set of directions. If the algorithm adds instance-independent noise then we show that this set of directions contains an (almost) orthonormal basis, allowing us to lower bound the trace of the MSE matrix, and hence the average distortion per entry. In the general case (when the distortion is instance-dependent), we use concentration inequalities for matrix-valued random variables to show that for appropriately chosen *random* databases, the trace of the MSE matrix is large with high probability.

We expect the linear algebraic techniques developed for this bound to be useful for bounding the required distortion of a wide range of differentially private releases.

## 1.2 Significance of the Privacy Lower Bounds

Dinur and Nissim [9] showed that if a mechanism answers (or allows the user to compute) $O(n \log n)$ arbitrary inner product queries on a database (vector) $s \in \{0,1\}^n$ with noise $o(\sqrt{n})$ per response, then an adversary can reconstruct $n - o(n)$ entries of $s$. Their attack was subsequently extended to use a linear number of queries [13], allow a small fraction of answers to be arbitrarily distorted [13], and run significantly more quickly [15]. These reconstruction attacks provide lower bounds for various minimal notions of privacy; our results extend the scope of these bounds significantly.

There were also several known lower bounds specific to differential privacy, though they are not directly relevant to marginal tables [12, 26, 18]. Subsequently to our work, Hardt and Talwar [19] gave upper and lower bounds for releasing a variety of linear functions (including marginal tables) for the special case of "pure" $\epsilon$-differential privacy (with $\delta = 0$). For the case of 1-attribute marginal tables, their bound of $\Omega(d/\epsilon)$ improves on ours; we conjecture that their techniques lead to a bound of $\widetilde{\Omega}(d^k/\epsilon)$ for releasing constant $k$-way marginal tables under $\epsilon$-differential privacy. However, their techniques break down for even slightly relaxed privacy notions such as $(\epsilon, \delta)$-differential privacy.

We see our new lower bounds as interesting for several reasons.

**Natural symmetric functions.** In their simplest form, the inner product queries considered by [9, 13, 15] require the adversary to be able to "name rows", that is, specify a coefficient for each entry of the vector $s$.

Thus, the lower bound does not apply directly to any functionality that is symmetric in the rows of the data set such as marginal tables. It was pointed out in [8] that in databases with more than one entry per row, random inner product queries (on, say, attribute $x_d$) can be simulated via hashing: for example, the adversary could ask for the sum of the function $H(x_1, \ldots, x_{d-1}) \cdot x_d$ over the whole database, where $H : \{0,1\}^{d-1} \to \{0,1\}$ is an appropriate hash function. This is a symmetric query, but it might seem odd to a statistician (with, e.g., a 2-wise independent hash function). The lower bounds we give for marginal table releases are the first for symmetric functions regularly released by official statistical agencies; one can think of our reconstruction attacks as using *conjunctions* as weak hash functions to implement the idea of [8].

| Privacy Guarantee | Upper Bound on Noise | Lower Bound on Noise |
|---|---|---|
| Attribute privacy | $\widetilde{O}(\min\{\sqrt{n}, \sqrt{d^k}\})$ | $\widetilde{\Omega}\left(\min\left\{\sqrt{n}, \sqrt{d^{k-1}}\right\}\right)$ |
| Row privacy | $\widetilde{O}(\min\{\sqrt{n}, \sqrt{d^k}\})$ | $\widetilde{\Omega}\left(\min\left\{\sqrt{n}, \sqrt{d^k}\right\}\right)$ |
| $(\epsilon, \delta)$-differential privacy (instance indep. noise) | $O\left(\frac{\sqrt{\binom{d}{k}}\log(1/\delta)}{\sqrt{2^k}\epsilon}\right)$ [5, 3] | $\Omega\left(\frac{\sqrt{\binom{d}{k}}(1-\delta/\epsilon)}{2^k\epsilon}\right)$ |
| $(\epsilon, \delta)$-differential privacy (with $\delta = 1/\mathrm{poly}(n)$) | $\widetilde{O}\left(\min\left\{n, \left(\frac{n^2 dk}{\epsilon}\right)^{\frac{1}{3}}, \frac{\sqrt{nd}}{\epsilon}\cdot 2^{\sqrt{k\log d}}, \frac{\sqrt{\binom{d}{k}}}{\epsilon\sqrt{2^k}}\right\}\right)$ [3, 6, 14] | $\widetilde{\Omega}\left(\min\left\{\frac{\sqrt{n}}{2^k\sqrt{\epsilon}}, \frac{\sqrt{\binom{d}{k}}}{2^k\epsilon}\right\}\right)$ |

**Table 1:** *Upper and lower bounds on the average noise per cell entry for releasing all $k$-way marginal tables (or equivalently all $k$-way conjunction predicates) under various privacy guarantees. The results on attribute non-privacy and row non-privacy are for constant $k$. The $n$ term in the upper bound for $(\epsilon, \delta)$-differential privacy (last row) comes from an algorithm that releases a vector of $n/2$'s for all $D$'s. All the uncited results appear in this paper.*

**When is distortion acceptably low?** It is natural to ask at what point the distortion required for privacy interferes with statistical analysis. There is no simple answer, but for the "predicate queries" considered here, where each entry counts the number of occurrences of a predicate in the underlying data set, there is a large class of statistical models which inherently have "sampling error", that is standard deviation of the observed statistics, of $\Omega(\sqrt{n})$. A crude rule of thumb, then, is that the distortion interferes seriously when it is not $o(\sqrt{n})$ [12, 3, 31]. Our lower bounds of $\widetilde{\Omega}(\min(\sqrt{n}, \sqrt{d^k}))$ show that for even modest values of $d$ and $k$, the data set $n$ must be very large to get distortion $o(\sqrt{n})$.

**The "dimension" of marginal tables.** The reconstruction attacks [9, 13, 15] above show a lower bound of roughly $\min\{\sqrt{n}, \sqrt{m}\}$ on the distortion required to answer a set of $m$ random, independent queries about a data set of size $n$. However, the bounds heavily rely on independence of the queries. This raises the question of whether certain interesting classes of queries could be answered with much less noise. For example, if a set of queries is linearly dependent, then one can compute noisy answers to only a few queries (a spanning set), and deduce the rest using the linear relationships. Both of our bounds can be interpreted as showing that the marginal statistics of a data set are, in a sense that depends on the notion of privacy, *far* from any low dimensional subspace. In particular, we show that the $\binom{d}{k}2^k$ different entries of the $k$-way marginal tables hide a set of $\Omega(d^k)$ "nearly independent" underlying features – as far as privacy is concerned, they have dimension close to $\binom{d}{k}$. It is natural to ask: what properties of a set of queries lead to this type of behavior, in general? Our techniques suggest that the right notion is related to inapproximability by low-dimensional linear spaces, however, it is unclear how to formulate this notion precisely.

## 1.3 Known Upper Bounds for Diff. Privacy

In [5, 12] it was shown that addition of carefully calibrated noise to functions satisfying a Lipschitz condition is enough to ensure differential privacy. Applied to conjunctions, they show that random noise drawn from a normal distribution with mean 0 and standard deviation $\sqrt{2\binom{d}{k}\log(1/\delta)}/\epsilon$ to each entry in $\mathcal{C}_k(D)$ guarantees $(\epsilon, \delta)$-differential privacy [5], while adding random noise drawn from a Laplacian distribution with mean 0 and standard deviation $2\binom{d}{k}/\epsilon$ to each entry in $\mathcal{C}_k(D)$ guarantees $\epsilon$-differential privacy (with

$\delta = 0$) [12]. Barak *et al.* [3] improve the dependency on $k$ in these results, saving a factor of approximately $\sqrt{2^k}$ in the required distortion.

In a different vein, Blum *et al.* [6] adapt the exponential sampling technique of [25] to release a synthetic data set. One can use their techniques to release $\mathcal{C}_k(D)$ with distortion $\widetilde{O}((n^2 dk/\epsilon)^{2/3})$ in each entry . The dependency on $d$ and $k$ in [6] is much better than in the additive noise mechanisms, but the dependency on $n$ is significantly worse; in particular, our results show that one cannot significantly reduce the dependency on $n$ without incurring a dependency on $d^k$. Finally, Dwork *et al.* [14] provide a similar, but incomparable, synthetic data mechanism, which allows one to release all $k$-way conjunctions with distortion $\widetilde{O}(\frac{\sqrt{nd}}{\epsilon}2^{\sqrt{k\log d}})$.

## 1.4 Preliminaries

We use $[n]$ to denote the set $\{1, 2, \ldots, n\}$. We use $negl(n)$ denotes a function that is asymptotically smaller than $1/n^c$ for all $c > 0$. Vectors used in the paper are by default column vectors. For a vector $v$, $v^\top$ denotes its transpose (row vector), $\|v\|$ denotes its Euclidean norm, and $v_i$ denotes its $i$th entry. We use $u_v$ to denote the unit vector corresponding to $v$ (i.e., $u_v = v/\|v\|$). For two vectors $v_1$ and $v_2$, $\langle v_1, v_2 \rangle$ denotes the inner product of $v_1$ and $v_2$. The length of projection of $v_1$ onto $v_2$ is then $\langle v_1, v_2 \rangle/\|v_2\|$. For a matrix $M$, $tr(M)$ denotes the trace and $\|M\|_\infty$ denotes the operator norm. We use $diag(a_1, \ldots, a_n)$ to denote an $n \times n$ diagonal matrix with entries $a_1, \ldots, a_n$ along the main diagonal. Let $\mathbb{I}_d$ denote the identity matrix of dimension $d$. We use $m_k$ to denote $\binom{d}{k}$.

Let $M$ be an $N \times n$ real matrix with $N \geq n$. The singular values $\sigma_j(M)$ are the eigenvalues of $\sqrt{M^\top M}$ arranged in non-increasing order. Of particular importance in this paper is the smallest singular value $\sigma_n(M) = \inf_{z:\|z\|=1} \|Mz\|$.

**Boolean Conjunctions.** It is convenient to describe the results in terms of releasing conjunction predicates over the domain $\{0, 1\}^d$. Each $x \in \{0, 1\}^d$ is interpreted as an assignment to $d$ Boolean variables $x_1, \ldots, x_d$. A conjunction predicate $c_v : \{0, 1\}^d \rightarrow \{0, 1\}$ for $v \in \{-1, 0, 1\}^d$ is defined as $c_v(x) = 1$ iff for all $i \in [d], x_i = 1$ if $v_i = 1$ and $x_i = 0$ if $v_i = -1$. The value of $v_i$ indicates whether the variable $x_i$ appears as not negated (if $v_i = 1$), negated (if $v_i = -1$), or absent (if $v_i = 0$). The length of a conjunction predicate is the number of coordinates of $v$ that are non-zero. We will refer to a conjunction predicate of length $k$ as a $k$-way conjunction. Let $\mathcal{C}_k$ be the function class

of all $k$-way conjunction predicates on variables $x_1, \ldots, x_d$. The size of $\mathcal{C}_k$, $|\mathcal{C}_k| = 2^k \binom{d}{k}$. Let $D \in (\{0,1\}^d)^n$ be a database. Each row of $D$ represents information contributed by one individual. The $i$th column of $D$ contains the assignments to variable $x_i$. For a predicate $c_v \in \mathcal{C}_k$, define $c_v(D) = \sum_{x \in D} c_v(x)$. We use $\mathcal{C}_k(D)$ to represent the vector of all predicates in $\mathcal{C}_k$ evaluated on $D$.

## 2. LOWER BOUNDS – MINIMAL PRIVACY

In this section, we introduce a reconstruction attack based on analyzing the least singular value of a random matrix with correlated entries. We then use the reconstruction attack to establish lower bounds on the noise needed for releasing $k$-way marginal tables under the notions of attribute non-privacy and row non-privacy. We treat $k$ as a constant in this section.

The lower bounds for our minimal privacy definitions proceed by "reducing" an instance of the reconstruction problem for a matrix with correlated rows into a marginal table release problem. To define the reduction from the reconstruction problem, we need the following definition of entry-wise product of vectors and matrices.

DEFINITION 2.1. *The entry-wise product of vectors $p, q \in \mathbb{R}^n$ is the vector in $p \odot q \in \mathbb{R}^n$ with entries $(p \odot q)_i = p_i \cdot q_i$. If $A$ is an $N_1 \times n$ matrix, and $B$ is an $N_2 \times n$ matrix, denote by $A \odot B$ an $N_1 \cdot N_2 \times n$ matrix, whose rows are entry-wise products of the rows of $A$ and $B$: $(A \odot B)_{j,k} = A_j \odot B_k$, where $(A \odot B)_{j,k}, A_j, B_k$ denote rows of the corresponding matrices.*

We now formally define attribute non-privacy and row non-privacy, and explain the reductions from the reconstruction problem.

DEFINITION 2.2 (ATTRIBUTE NON-PRIVACY). *An algorithm $\mathcal{A}$ for releasing all $k$-way conjunction predicates is* attribute non-private *if there exists a polynomial time adversary such that for every $s \in \{0,1\}^n$ there exists a database $D_{at}(s) \in (\{0,1\}^d)^n$ whose last column is $s$, such that the adversary with input $\mathcal{A}(D_{at}(s))$ and the first $d-1$ columns of $D_{at}(s)$, can reconstruct at least $\widetilde{\Omega}(\min\{n, d^{k-1}\})$ entries of $s$ with probability $1 - negl(d)$.*

This definition captures a common model in the data privacy literature (e.g., [32, 23, 22, 7, 24, 34]) where one assumes that a database (with $d$ attributes) consists of $d-1$ *nonsensitive* attributes (e.g., demographic information), which can be learned from other sources, and one *sensitive* attribute (e.g., disease). The attribute non-privacy lower bound applies to any notion of privacy that purports to protect individual values of the sensitive attribute (the lower bound applies in particular, to differential privacy but also, e.g., to the notion of privacy implicit in the popular "$K$-anonymization" scheme [32] and its recent variants [23, 22, 7, 24, 34]).

**Reduction from the Reconstruction Problem to Attribute Non-Privacy:** Consider a matrix $M \in \{0,1\}^{d \times n}$. Let $M^{(k)} = M \odot M \odot \ldots \odot M$ be the dimension $d^k \times n$ matrix obtained by applying $\odot$ operator $k-1$ times. Let $s \in \{0,1\}^n$. Consider the database $D = (M^\top | s) \in (\{0,1\}^{d+1})^n$. That is, the first $d$ columns of $D$ are given by the rows of $M$, and the last column is $s$. Then $\mathcal{C}_k(D)$ contains the vector

$M^{(k-1)}s$, i.e., $M^{(k-1)}s \subset \mathcal{C}_k(D)$. Now, any algorithm that allows an adversary to reconstruct a large fraction of $s$ given $M^{(k-1)}s$ is attribute non-private. This reduction holds for every $M \in \{0,1\}^{d \times n}$, but in the analysis, we use a random matrix $M$.

DEFINITION 2.3 (ROW NON-PRIVACY). *An algorithm $\mathcal{A}$ for releasing all $k$-way conjunction predicates is* row non-private *if there exists a distribution of databases $\mathbb{D}$ over the domain $(\{0,1\}^d)^n$ under which the rows of the databases are statistically independent and there exists a set $S \subseteq [n]$ whose size is at least $\widetilde{\Omega}(\min\{n, d^k\})$ satisfying the following:*

1. *For any (not-necessarily polynomial time) adversary if $D \sim \mathbb{D}$, the adversary can output any row of $D$ indexed by the elements of $S$ with probability at most a constant (say $2/3$);*

2. *There exists a polynomial time adversary such that if $D \sim \mathbb{D}$, the adversary on input $\mathcal{A}(D)$ can output $1 - o(1)$ fraction of the rows of $D$ indexed by the elements of $S$ with probability $1 - negl(d)$.*

The row non-privacy lower bound applies, roughly, to any notion of privacy that seeks to protect any complete row of the database (as opposed to only individual entries). This includes differential privacy as well as its relaxations to metrics on probability distributions such as total variation distance or KL divergence [12, 33, 1, 16].

**Reduction from the Reconstruction Problem to Row Non-Privacy:** Consider a matrix $M \in \{0,1\}^{d \times n}$. Consider the database $D = diag(s) \cdot M^\top \in (\{0,1\}^d)^n$, where $diag(s)$ is an $n \times n$ diagonal matrix with diagonal $s$. Because $s$ is a $(0,1)$-vector, this corresponds to a world where person $i$'s data is either $M_i^\top$ or $0^d$, according to the $i$th bit of $s$ (where $M_i$ is the $i$th row of $M$). Then $\mathcal{C}_k(D)$ contains the vector $M^{(k)}s$. Again, this reduction holds for every $M \in \{0,1\}^{d \times n}$, but in the analysis, we use a random matrix $M$.

### 2.1 Lower Bounds – Reconstruction Problem

Let $s = (s_1, \ldots, s_n) \in \{0,1\}^n$ be some (secret) vector. Let $c_k$ be a constant (we will define it later in Theorem 2.5). Let the integer $a = \min\{n, c_k d^k / \log^{2k-2} n\}$ (to simplify the exposition, we shall ignore rounding issues). Let $s|_a = (s_1, \ldots, s_a)$ be the first $a$ entries of $s$. Let $\Phi \in \{0,1\}^a$ be a vector with independent entries taking values 0 and 1 with probability $1/2$. Let $\Phi_1, \ldots, \Phi_d \in \{0,1\}^a$ be $d$ independent copies of $\Phi$. Let $M$ be a $d \times a$ matrix whose rows are $\Phi_1, \ldots, \Phi_d$. Again, we consider the matrix $M^{(k)}$ constructed from $M$.

The attack works as follows: for every row $R$ in $M^{(k)}$, the adversary asks inner product of $R$ with $s|_a$, and receives noisy responses. Consider a privacy mechanism $\mathcal{A}$. Let $p = \mathcal{A}(M^{(k)}s|_a)$ be the vector of noisy answers generated by $\mathcal{A}$. Define the error (noise) vector as $e = p - M^{(k)}s|_a$. Let $M^{(k)} = P\Gamma Q$ be the singular value decomposition of $M^{(k)}$. Here, $P$ is a $d^k \times d^k$ orthogonal matrix, $\Gamma$ is a $d^k \times a$ diagonal matrix, and $Q$ is an $a \times a$ orthogonal matrix. Let $\mathbf{0}$ be a $(d^k - a) \times a$ matrix with all entries zero. Define $\Gamma^{-1} = (diag(\sigma_1(M^{(k)})^{-1}, \ldots, \sigma_a(M^{(k)})^{-1}) | \mathbf{0}^\top)$. The dimension of $\Gamma^{-1}$ is $a \times d^k$. Define $M_{inv}^{(k)} = Q^\top \Gamma^{-1} P^\top$.

Now, given $p$, the adversary uses $M_{inv}^{(k)}$ to construct $\hat{s} = (\hat{s}_1, \ldots, \hat{s}_a)$ as follows: $\hat{s}_i = 1$ if the $i$th entry in $M_{inv}^{(k)}p \geq$

1/2, and 0 otherwise. Now, the claim is that $\hat{s}$ is a good reconstruction of $s|_a$. The idea behind the analysis is that $M_{inv}^{(k)}p = s|_a + M_{inv}^{(k)}e$, and therefore (as $P$ and $Q$ are orthogonal matrices),

$$\|M_{inv}^{(k)}e\| = \|Q^\top \Gamma^{-1} P^\top e\| = \|\Gamma^{-1} P^\top e\|$$
$$\leq \|\Gamma^{-1}\|_\infty \|P^\top e\| = \|e\|/\sigma_a(M^{(k)}).$$

Corollary 2.6 shows that with high probability $\sigma_a(M^{(k)}) = \widetilde{\Omega}(\sqrt{d^k})$. If an algorithm only adds $o(\sqrt{n})$ noise to each query (i.e., all the entries in $e$ are $o(\sqrt{n})$) then $\|e\| = o(\sqrt{d^k n})$, and therefore, $\|M_{inv}^{(k)}e\| \approx o(\sqrt{n})$ with high probability. In particular, if $a = n$, then this implies that with high probability $M_{inv}^{(k)}e$ cannot have $\Omega(n)$ entries with absolute value above $1/2$, and therefore, the Hamming distance between $\hat{s}$ and $s|_a$ is $o(a) = o(n)$ (as the adversary only fails to recover those entries of $s|_a$ whose corresponding $M_{inv}^{(k)}e$ entries are greater than $1/2$). The following proposition formalizes this observation. The proof uses some ideas from a recent reconstruction attack proposed by Dwork and Yekhanin [15].

PROPOSITION 2.4. *Let $k$ be a constant. If an algorithm adds*

$$o(\min\{\sqrt{n}/\log^{(k^2+k+1)} n, \sqrt{d^k}/\log^{(k^2+3k-1)} n\})$$

*noise to each entry in $M^{(k)}s|_a$, then there exists an adversary that can reconstruct $1 - o(1)$ fraction of $s|_a$ with probability at least $1 - negl(d)$.*

The proof of the above proposition relies heavily on the following theorem that lower bounds the least singular value of a random matrix with correlated rows. A proof outline of the theorem is given in the Section 3.

THEOREM 2.5. *Let $k, m, d$ be natural numbers such that $m \leq c_k d^k / \log^{2k-2} m$ where $c_k$ depends only on $k$, and let $A$ be a $d \times m$ matrix with independent entries taking values $0$ and $1$ with probability $1/2$. Then there exists numbers $C_k, c_k'$ which depend only on $k$ such that the $k$-times entry-wise product $\tilde{A} = A \odot A \odot \ldots \odot A$ is a $d^k \times m$ matrix satisfying*

$$\Pr\left[\sigma_m(\tilde{A}) \leq \frac{c_k' \sqrt{d^k}}{\log^{(k^2+k+1)} m}\right] \leq 2 \exp\left(-C_k d/\log^{2k-2} m\right).$$

Now, as $M$ is a $d \times a$ matrix with independent entries taking values 0 and 1 with probability $1/2$, and $a \leq c_k d^k / \log^{2k-2} n$, we can apply the above theorem to conclude the following.

COROLLARY 2.6. $\Pr\left[\sigma_a(M^{(k)}) \leq \frac{\sqrt{d^k}}{\log^{(k^2+k+1)} n}\right] \leq negl(d)$.

## 2.2 Attribute Non-Privacy

We use the reduction from the reconstruction attack described earlier. Define, a vector $\hat{\Phi}_j \in \{0,1\}^n$ as follows: the $i$th entry in $\hat{\Phi}_j$ is the $i$th entry of $\Phi_i$ if $i \leq a$, and is 0 otherwise. For a vector $s \in \{0,1\}^n$, define a database $D_{at}(s) \in (\{0,1\}^d)^n$ as follows: the first $d$ columns are $\hat{\Phi}_1, \ldots, \hat{\Phi}_{d-1}$, the last column is $s$. The following theorem uses Proposition 2.4 to show that there exists an adversary that can reconstruct $1 - o(1)$ fraction of the first $\widetilde{\Omega}(\min\{n, d^{k-1}\})$ entries of $s$ if given too accurate vector $\mathcal{C}_k(D_{at}(s))$.

THEOREM 2.7. *Let $k$ be a constant. Any algorithm $\mathcal{A}$ for releasing all $k$-way marginal tables (or equivalently all $k$-way conjunction predicates) that for every database $D \in (\{0,1\}^d)^n$ adds*

$$o\left(\min\left\{\sqrt{n}/\log^{(k^2-k+1)} n, \sqrt{d^{k-1}}/\log^{(k^2+k-3)} n\right\}\right)$$

*noise to each entry in $\mathcal{A}(D)$ is attribute non-private.*

## 2.3 Row Non-Privacy

Again, we use the reduction from the reconstruction attack described earlier. For a vector $s \in \{0,1\}^n$, define a database $D_{st}(s) \in (\{0,1\}^d)^n$ as follows: $(i,j)$th entry of $D_{st}(s)$ is $s_i$ if the $i$th entry in $\hat{\Phi}_j = 1$, and 0 otherwise. We assume that the adversary has access to $\Phi_1, \ldots, \Phi_d$.

Define a distribution $\mathbb{D}$ over the set of databases as follows: draw a vector $s_r$ uniformly at random from $\{0,1\}^n$ and output $D_{st}(s_r)$. Let consider some $i$th row where $i \in [a]$. Let $E$ be the event that there exists a $\Phi_j$ such that $i$th entry in $\Phi_j$ is 1. Conditioned on event $E$, an adversary can only predict the $i$th row of $D_{st}(s_r)$ by guessing the $i$th entry in $s_r$. Since $s_r$ is picked uniformly at random, this implies that conditioned on $E$ no adversary can guess the $i$th row of $D_{st}(s_r)$ with probability more than $1/2$. Finally, since $\Pr[\overline{E}] = 1/2^d$, therefore, no adversary (even with access to $\Phi_1, \ldots, \Phi_d$) can guess the $i$th row of $D_{st}(s_r)$ with probability more than $1/2 + 1/2^d \leq 2/3$. Thus, $\mathbb{D}$ satisfies the first condition of Definition 2.3 for every set $S$. The following theorem uses this distribution $\mathbb{D}$ to obtain a lower bound on the noise needed for row privacy.

THEOREM 2.8. *Let $k$ be a constant. Any algorithm for releasing all $k$-way marginal tables (or equivalently all $k$-way conjunction predicates) that for every database $D \in (\{0,1\}^d)^n$ adds*

$$o\left(\min\left\{\sqrt{n}/\log^{(k^2+k+1)} n, \sqrt{d^k}/\log^{(k^2+3k-1)} n\right\}\right)$$

*noise to each entry in $\mathcal{A}(D)$ is row non-private.*

## 3. PROOF OUTLINE OF THEOREM 2.5

Estimating the smallest singular value of the matrix $\tilde{A}$ presents two challenges. The entries of this matrix are interdependent, which makes powerful measure concentration tools hard to apply. Also, the entries are non-centered, and hence its operator norm is of order $\sqrt{d^k n}$ with high probability. The norm of the matrix enters many probabilistic bounds involved in the proof, and such a large norm would render most of these bounds meaningless. To remove these obstacles, we apply a simple decoupling and symmetrization argument to reduce the problem to bounding the smallest singular value of a matrix $\tilde{\Pi}$, which is an entry-wise product of $\Pi_1, \ldots, \Pi_k$, where each $\Pi_i$ is a $d \times n$ matrix with independent random entries taking values $-1$ with probability $1/4$, 1 with probability $1/4$, and 0 with probability $1/2$. Note that the entries in $\Pi_i$ are centered. Analysis of the behavior of the least singular value of $\tilde{\Pi}$ is the core of the argument.

The first step in this analysis is obtaining a probabilistic bound for the operator norm of $\tilde{\Pi}$. This bound is proved by induction on $k$, with Talagrand's measure concentration inequality for convex functions (see [21], Corollary 4.10) applied at each step. We show that with high probability,

$$\|\tilde{\Pi}|_J\|_\infty \leq c_k''(d^{k/2} + \sqrt{|J|} \cdot \log^{k/2} n) \cdot \log^{(k-1)/2} n \quad (1)$$

for all non-empty subsets $J \subset \{1, \ldots, n\}$. Here $\tilde{\Pi}|_J$ denotes the submatrix of $\tilde{\Pi}$ with columns belonging to $J$, and $c''_k$ is a constant depending on $k$ only.

The smallest singular value of $\tilde{\Pi}$ is the minimum of $\|\tilde{\Pi}x\|$, over $x$ from the unit sphere. Before we analyze this quantity in full generality, we consider a simpler question of estimating *the small ball probability*. This is the probability that $\|\tilde{\Pi}x\|$ is small for a fixed vector $x$. Measure concentration plays a prominent role in this estimate as well. To apply measure concentration tools, we have to represent $\|\tilde{\Pi}x\|$ as a function of independent random variables. This is done by conditioning on the matrices $\Pi_1, \ldots, \Pi_{l-1}$ and $\Pi_{l+1}, \ldots, \Pi_k$, for some $l \in \{1, \ldots, k\}$. In this case the random variable $\|\tilde{\Pi}x\|$ depends only on the matrix $\Pi_l$ with independent entries.

The number $l \in \{1, \ldots, k\}$ is chosen to match the level of *compressibility* of the vector $x$. A vector is compressible if its norm is concentrated on a small number of coordinates. Note that the coordinates of the vector $\tilde{\Pi}x$ consist of $d^{k-l}$ blocks $(\Pi_1 \odot \ldots \odot \Pi_l)y_i^\top$, where $y_i$ denotes the $i$th row of the matrix $\Pi_{l+1} \odot \ldots \Pi_k \odot x^\top$. We show that with high probability most of the vectors $y_i$ have norm and degree of compressibility comparable to those of $x$. This reduces the small ball probability bound for the vector $\tilde{\Pi}x$ to a similar bound for each vector $(\Pi_1 \odot \ldots \odot \Pi_l)y_i^\top$. To obtain such a bound, we decompose the coordinates of this vector in a similar fashion in $d$ blocks $(\Pi_1 \odot \ldots \odot \Pi_{l-1} \odot \Pi_l(j))y_i^\top$, where $\Pi_l(j)$ is the $j$-th row of the matrix $\Pi_l$.

For any matrix $A$ with $n$ columns, and any vectors $u, v \in \mathbb{R}^n$, $(A \odot u^\top)v = (A \odot v^\top)u$. Using this elementary identity, we rewrite $(\Pi_1 \odot \ldots \odot \Pi_{l-1} \odot \Pi_l(j))y_i^\top$ as $(\Pi_1 \odot \ldots \odot \Pi_{l-1} \odot y_i)\Pi_l(j)^\top$. After conditioning on $B := \Pi_1 \odot \ldots \odot \Pi_{l-1}$ this becomes the product of a fixed matrix $B \odot y_i$ and a random vector $\Pi_l(j)^\top$ with independent coordinates. The small ball probability bound for such vector can be obtained by applying Talagrand's concentration theorem for convex functions to the function $F(w) = \|(B \odot y_i)w\|$. The level of concentration is determined by the Lipschitz constant of the function $F$, i.e., by the norm of the matrix $B \odot y_i$. This norm, in turn, is controlled using the inequality (1) applied to the matrix $B$.

We can obtain a uniform bound for the entire sphere via a net argument. Unfortunately, since the small ball probability for $\|\tilde{\Pi}x\|$ depends significantly on the vector $x$, it is impossible to construct one net that would work for the whole sphere. Therefore, we decompose the sphere in numerous regions, and estimate the probability that $\|\tilde{\Pi}x\|$ is small for each part separately. The regions are defined by the small ball probability, which depends on the compressibility of the vectors. For each part we apply a net argument especially tailored for a certain degree of compressibility. Namely, the region is discretized, by using a $\gamma$-net for appropriate $\gamma$. Then we obtain a uniform lower estimate on the net, using the small ball probability and the union bound. This estimate is extended to the whole region by approximation. This method requires a careful balance between the small ball probability and the size of the net. The better the small ball probability is, the bigger a net we can consider, and so the bigger region we can cover. This balance dictates the aforementioned decomposition of the sphere.

We start with obtaining a uniform estimate of $\|\tilde{\Pi}x\|$ over a set of all vectors $x$ having a given level of sparsity. We then extend the bound from the set of sparse vectors to the set of

compressible vectors with a certain level of compressibility. Finally, we show that the whole sphere can be assembled from these sets. This allows to finish the proof by using the union bound. See the full version [20] for a complete proof.

## 4. LOWER BOUNDS – DIFF. PRIVACY

In this section, we establish lower bounds on the noise needed for releasing all $k$-way marginal tables under the popular notion of differential privacy. A database $D'$ is said to be a neighbor of a database $D$ if it differs from $D$ in exactly one row. A randomized algorithm is *differentially private* if neighbor databases induce nearby distributions on the outputs.

DEFINITION 4.1 $((\epsilon, \delta)$-DIFFERENTIAL PRIVACY [12]$)$. *A randomized algorithm $\mathcal{A}$ is $(\epsilon, \delta)$-differentially private if for all neighboring databases $D, D'$, and for all sets $\mathcal{S}$ of possible outputs*

$$\Pr[\mathcal{A}(D) \in \mathcal{S}] \leq \exp(\epsilon) \cdot \Pr[\mathcal{A}(D') \in \mathcal{S}] + \delta.$$

*The probability is taken over the random coins of the algorithm $\mathcal{A}$. If $\mathcal{A}$ is $(\epsilon, 0)$-differentially private (i.e., $\delta = 0$), then we say it is $\epsilon$-differentially private.*

Let $X$ and $Y$ be random variables taking values in a set $\mathcal{O}$. We use $X \approx_{\epsilon, \delta} Y$ to indicate that random variables $X$ and $Y$ are $(\epsilon, \delta)$-indistinguishable, i.e.,

$$\forall \mathcal{S} \subseteq \mathcal{O}, \Pr[X \in \mathcal{S}] \leq \exp(\epsilon) \cdot \Pr[Y \in \mathcal{S}] + \delta \quad \text{and}$$
$$\Pr[Y \in \mathcal{S}] \leq \exp(\epsilon) \cdot \Pr[X \in \mathcal{S}] + \delta.$$

We also use $X \approx_\epsilon Y$ to indicate that random variables $X$ and $Y$ are $(\epsilon, 0)$-indistinguishable.

Our bounds are tight under a natural and popular class of differentially private algorithms based on adding instance-independent noise. This class contains algorithms that for all inputs add noise from a fixed distribution (i.e., the noise distribution is independent of the input). Formally, if an algorithm $\mathcal{A}$ for a function class $\mathcal{F}$ adds instance-independent noise from a distribution $Z$ then for all $D$, $\mathcal{A}(D) = \mathcal{F}(D) + Z$. Therefore, for $D'$ a neighbor of $D$, $\mathcal{A}(D') = \mathcal{A}(D) + \mathcal{F}(D') - \mathcal{F}(D)$. The SuLQ algorithm of Blum *et al.* [5] is an example of an algorithm that adds instance-independent noise.

In Section 4.1, we consider $(\epsilon, \delta)$-differentially private algorithms for $\mathcal{C}_k$ that add instance-independent noise. For an instance-independent differentially private algorithm $\mathcal{A}$, we can measure the perturbation introduced by $\mathcal{A}$ either by using the mean squared error matrix

$$\Sigma_{\mathcal{A}}(D) = \mathbb{E}[(\mathcal{A}(D) - \mathcal{C}_k(D))(\mathcal{A}(D) - \mathcal{C}_k(D))^\top]$$

or the covariance matrix of $\mathcal{A}(D)$, and the results are the same with either choice. Define the average mean squared error of $\mathcal{A}(D)$ as the trace of $\Sigma_{\mathcal{A}}(D)$ divided by the size of $\mathcal{A}(D)$. Let $m_k = \binom{d}{k}$. We show that if for every database $D$, $\mathcal{A}(D)$ has an average mean squared error (or variance) of $o(m_k(1 - \delta/\epsilon)^2/(2^{2k}\epsilon^2))$, then $\mathcal{A}$ is *not* $(\epsilon, \delta)$-differentially private. To do so, we analyze projections onto various directions. The idea is to show that for any neighboring databases $D$ and $D'$ with $\mathcal{C}_k(D') - \mathcal{C}_k(D) = \Delta$, the indistinguishability requirement of differential privacy forces both the expected squared length of the projection of $\mathcal{A}(D) - \mathcal{C}_k(D)$ on $\Delta$ and the expected squared length of the projection of $\mathcal{A}(D') - \mathcal{C}_k(D')$ on $\Delta$ to be at least square of the length of $\Delta$.

Of particular interest to us are the direction vectors $\Delta$'s with large lengths (close to the largest possible length of $\sqrt{m_k}$). Then, using a careful argument involving geometries of these $\Delta$ vectors we show that there exists a database $D^*$ such that the trace of $\Sigma_{\mathcal{A}}(D^*)$ is at least $m_k^2(1 - \delta/\epsilon)^2/(2^k\epsilon^2)$. The result follows by dividing the trace by the size of $\mathcal{A}(D^*)$.

In Section 4.2, we consider general $(\epsilon, \delta)$-differentially private algorithms for $\mathcal{C}_k$. For a "general" differentially private algorithm $\mathcal{A}$, we *need*[4] to use the mean squared error matrix to measure the perturbation. We show that if for every database $D$, $\mathcal{A}(D)$ has an average mean squared error of $o(\min\{m_k(1 - \delta/\epsilon)^2/(2^{2k}\epsilon^2), n(1 - \delta/\epsilon)^2/(2^{2k}\epsilon\log m_k)\})$, then $\mathcal{A}$ is *not* $(\epsilon, \delta)$-differentially private. Again for neighboring databases $D$ and $D'$ with $\mathcal{C}_k(D') - \mathcal{C}_k(D) = \Delta$, we investigate the expected squared length of the projections of $\mathcal{A}(D) - \mathcal{C}_k(D)$ and $\mathcal{A}(D') - \mathcal{C}_k(D')$ on $\Delta$. The analysis of the general case is harder, because now the indistinguishability requirement forces only one among these two projection lengths to be greater than the squared length of $\Delta$. Our proof looks at random databases and shows that for a random database $D_r$ with high probability the trace of $\Sigma_{\mathcal{A}}(D_r)$ is at least $\min\{m_k^2(1 - \delta/\epsilon)^2/(2^k\epsilon^2), nm_k(1 - \delta/\epsilon)^2/(2^k\epsilon\log m_k)\}$.

**Inner Products.** In our analysis, (for simplicity) instead of conjunctions, we consider inner products over the domain $\{-1, 1\}^d$. An inner product predicate $i_v : \{-1, 1\}^d \to \{-1, 1\}$ is defined as $i_v(x) = \prod_i x_i \cdot v_i$, where the value of $v_i$ indicates whether $x_i$ is present (if $v_i = 1$) or not (if $v_i = 0$). Similar to $\mathcal{C}_k$, let $\mathcal{I}_k$ be the class of all $k$-way inner product predicates. Let $D$ be a database from $(\{-1, 1\}^d)^n$. For a predicate $i_v \in \mathcal{I}_k$, define $i_v(D) = \sum_{x \in D} i_v(x)$. Let $\mathcal{I}_k(D)$ be the vector of all predicates in $\mathcal{I}_k$ evaluated on $D$. In the full version [20], we provide the relationship between releasing conjunctions and inner products. Informally, what we show is that if there exists a database $D_o \in (\{-1, 1\}^d)^n$ such that no $(\epsilon, \delta)$-differentially private algorithm $\mathcal{B}$ for $\mathcal{I}_k$ has $tr(\Sigma_{\mathcal{B}}(D_o)) \leq T$, then there exists a database $D_z \in (\{0, 1\}^d)^n$ such that no $(\epsilon, \delta)$-differentially private algorithm $\mathcal{A}$ for $\mathcal{C}_k$ has $tr(\Sigma_{\mathcal{A}}(D_z)) \leq T/2^k$.

### 4.1 Instance-independent Additive Case

For simplicity, we set $\delta = 0$ in the following discussion (we will introduce $\delta$ in Theorem 4.6). We start by proving a very useful property about differential privacy. We state the lemma in terms of a general function class $\mathcal{F}$ and later use if for our specific function class $\mathcal{I}_k$. Let $\mathcal{A}$ be an $\epsilon$-differentially private algorithm for $\mathcal{F}$ that adds instance-independent noise. The lemma shows that both $\mathbb{E}[\langle \mathcal{A}(D) - \mathcal{F}(D), \Delta \rangle^2]$ and $\mathbb{E}[\langle \mathcal{A}(D') - \mathcal{F}(D'), \Delta \rangle^2]$ are $\Omega(\langle \Delta, \Delta \rangle^2/\epsilon^2)$ where $\Delta = \mathcal{F}(D') - \mathcal{F}(D)$. The proof uses the fact that projections onto direction $\Delta$ need to be $\epsilon$-indistinguishable for $\mathcal{A}(D)$ and $\mathcal{A}(D')$.

LEMMA 4.2. *Let $\mathcal{F}$ be a function class of Boolean predicates, and let $\mathcal{A}$ be an $\epsilon$-differentially private algorithm for $\mathcal{F}$ that adds instance-independent noise. Let $\mathcal{A}(D) = \mathcal{F}(D) +$*

$Z$. *Let $\Delta = \mathcal{F}(D') - \mathcal{F}(D)$, and let $\mathcal{A}(D) \approx_{\epsilon} \mathcal{A}(D') = \mathcal{A}(D) + \Delta$. Then, $\mathbb{E}[\langle \mathcal{A}(D) - \mathcal{F}(D), \Delta \rangle^2] = \Omega(\langle \Delta, \Delta \rangle^2/\epsilon^2)$.*

**One-way Inner products.** If $k = 1$ (i.e., 1-way inner products), then the remaining analysis is quite simple. Let $D_e$ be any database which has at least a row of both $(-1)^d$ and $(1)^d$. Consider a vector $\Delta \in \{-2, 2\}^d$, construct $D_{\Delta}$ from $D_e$ by replacing the row $(-1)^d$ is replaced by $\Delta/2$ and the row $(1)^d$ by $\Delta/2$. The Hamming distance between $D_e$ and $D_{\Delta}$ is 2 and $\mathcal{I}_1(D_{\Delta}) - \mathcal{I}_1(D_e) = \Delta$. Also, $\langle \Delta, \Delta \rangle^2 = 16d^2$. We use the above construction to create for every $\Delta = \{-2, 2\}^d$ a corresponding database $D_{\Delta}$. The idea now is to use the fact that this set of $\Delta$'s (which contains every vector from $\{-2, 2\}^d$) contains an orthogonal (Hadamard) basis $\Delta_1, \ldots, \Delta_d$, and therefore, by invoking Lemma 4.2 for $\mathcal{I}_1$, and noting that $\sum_{i=1}^{d} \Delta_i \Delta_i^{\top} = 4d \cdot \mathbb{I}_d$, we have

$$\Omega\left(\frac{d^3}{\epsilon^2}\right) = \sum_{\Delta_i} tr(\mathbb{E}[\langle \mathcal{A}(D_e) - \mathcal{I}_1(D_e), \Delta_i \rangle^2])$$

$$= \sum_{\Delta_i} tr(\Delta_i^{\top} \Sigma_{\mathcal{A}}(D_e) \Delta_i) = \sum_{\Delta_i} tr(\Sigma_{\mathcal{A}}(D_e) \Delta_i \Delta_i^{\top})$$

$$= tr(\Sigma_{\mathcal{A}}(D_e) \cdot 4d \cdot \mathbb{I}_d) = 4d \cdot tr(\Sigma_{\mathcal{A}}(D_e)).$$

PROPOSITION 4.3. *Let $\mathcal{A}$ be an $\epsilon$-differentially private algorithm for $\mathcal{I}_1$ that adds instance-independent noise. Let $D_e$ be any database which has at least a row of both $(-1)^d$ and $(1)^d$. Then, $tr(\Sigma_{\mathcal{A}}(D_e)) = \Omega(d^2/\epsilon^2)$.*

**Extension to $k$-way inner products.** The analysis for $k > 1$ ($k$ not necessarily a constant) is trickier, as we don't get a set of orthogonal $\Delta$ vectors. Here, we start with a special database $D_c = ((1)^d)^n$ (a database of all 1's), and look at the neighbors of $D_c$ obtained by replacing a row of $D_c$ by a vector from $\{-1, 1\}^d$. Let $D_c'$ be a neighbor of $D_c$, and let $\tilde{z} = \mathcal{I}_k(D_c') - \mathcal{I}_k(D_c)$. Assume that $D_c'$ is obtained from $D_c$ by replacing the $j$th row of $D_c$ by a vector $d_c' \in \{-1, 1\}^d$. Therefore, $\tilde{z} = (n)^{m_k} - ((n-1)^{m_k} + \mathcal{I}_k(d_c')) = (1)^{m_k} - \mathcal{I}_k(d_c')$ has lots of 0 entries making $\|\tilde{z}\|$ "small". This is true for many different choices of $D_c'$. To overcome this problem we analyze projections of $\mathcal{A}(D) - \mathcal{I}_k(D)$ onto direction $\pi\tilde{z}$ where $\pi = \mathbb{I}_{m_k} - oo^{\top}/\langle o, o \rangle$ ($\pi\tilde{z}$ is the orthogonal projection of $\tilde{z}$ onto the orthogonal complement of $o = (1)^{m_k}$). These projections have the advantage that $\|\pi\tilde{z}\| = \|\pi \cdot \mathcal{I}_k(d_c')\|$ is "big" with probability at least $1/2$ over random choices of $\tilde{z}$. The idea now is to use Lemma 4.2 to show that $\mathbb{E}[\langle \mathcal{A}(D_c) - \mathcal{I}_k(D_c), \pi\tilde{z} \rangle^2] = \Omega(\langle \pi\tilde{z}, \pi\tilde{z} \rangle^2/\epsilon^2)$.

LEMMA 4.4. *Let $\mathcal{A}$ be an $\epsilon$-differentially private algorithm for $\mathcal{I}_k$ that adds instance-independent noise. Let $D_c = ((1)^d)^n$. Let $D_c' \in (\{-1, 1\}^d)^n$ be a neighbor of $D_c$. Let $\tilde{z} = \mathcal{I}_k(D_c') - \mathcal{I}_k(D_c)$, $z = \mathcal{I}_k(D_c') - (n - 1)^{m_k}$, and $\pi = \mathbb{I}_{m_k} - oo^{\top}/\langle o, o \rangle$ where $o = (1)^{m_k}$. Then,*

$$\mathbb{E}[\langle \mathcal{A}(D_c) - \mathcal{I}_k(D_c), \pi\tilde{z} \rangle^2] = \Omega\left(\frac{\langle \pi\tilde{z}, \pi\tilde{z} \rangle^2}{\epsilon^2}\right) = \Omega\left(\frac{\langle \pi z, \pi z \rangle^2}{\epsilon^2}\right).$$

In the proof, we use random directions $z$. The following lemma analyzes the structure of $\mathbb{E}_z[zz^{\top}]$ (where the randomness is over the choice of $z$).

LEMMA 4.5. *Let $r \in \{-1, 1\}^d$ be a random vector with independent entries taking values $-1$ and $1$ with probability $1/2$. Let $m_k = \binom{d}{k}$. Define a random vector $z_r$ of length $m_k$ as $z_r = \mathcal{I}_k(r)$. Define a matrix $B = \mathbb{E}_{z_r}[z_r z_r^{\top}]$ where*

---

[4]This is because an algorithm could always add noise such that the output released is always a 0 vector for every database. This is clearly *not* a good algorithm as the deviation from the true answer could be quite big. But this algorithm clearly satisfies all the privacy requirements and also the variance in each coordinate of $\mathcal{A}(D)$ is 0.

the randomness is over $z_r$. Then, $B = \mathbb{I}_{m_k}$ where $\mathbb{I}_{m_k}$ is an identity matrix of dimension $m_k$.

The following theorem uses Lemmata 4.4 and 4.5 to show that every instance-independent differentially private algorithm needs to add a lot of noise to release $\mathcal{I}_k(D_c)$. The idea is to analyze $\mathbb{E}_{z_r}[(\pi z_r)^\top \Sigma_{\mathcal{A}}(D_c)(\pi z_r)]$. Since, $\mathbb{E}_{z_r}[z_r z_r^\top] = \mathbb{I}_{m_k}$ (from Lemma 4.5), we have

$$\mathbb{E}_{z_r}[(\pi z_r)^\top \Sigma_{\mathcal{A}}(D_c)(\pi z_r)] = \mathbb{E}_{z_r}[tr((\pi z_r)^\top \Sigma_{\mathcal{A}}(D_c)(\pi z_r))]$$
$$= \mathbb{E}_{z_r}[tr(\pi^\top \Sigma_{\mathcal{A}}(D_c)\pi z_r z_r^\top)] = tr(\pi^\top \Sigma_{\mathcal{A}}(D_c)\pi \mathbb{I}_{m_k})$$
$$\leq tr(\Sigma_{\mathcal{A}}(D_c))\|\pi\|_\infty \|\mathbb{I}_{m_k}\|_\infty = tr(\Sigma_{\mathcal{A}}(D_c)).$$

The third equality follows as trace and expectation commute. The last equality follows as $\pi$ is a projection matrix. Now, using Lemma 4.4, we show that

$$\mathbb{E}_{z_r}[(\pi z_r)^\top \Sigma_{\mathcal{A}}(D_c)(\pi z_r)] = \Omega(m_k^2/\epsilon^2),$$

and therefore, $tr(\Sigma_{\mathcal{A}}(D_c)) = \Omega(m_k^2/\epsilon^2)$. As explained earlier, a lower bound for releasing $k$-way inner products can be converted to a lower bound for releasing $k$-way conjunctions.

THEOREM 4.6. *Let $m_k = \binom{d}{k}$. Any algorithm $\mathcal{A}$ for releasing all $k$-way marginal tables (or equivalently all $k$-way conjunction predicates) that adds instance-independent noise and that for every database $D \in (\{0,1\}^d)^n$ has a root mean squared error (or standard deviation) of*

$$o(\sqrt{m_k}(1 - \delta/\epsilon)/(2^k \epsilon))$$

*for each entry of $\mathcal{A}(D)$ is not $(\epsilon, \delta)$-differentially private.*

## 4.2   Lower Bounds – General Case

Again our analysis looks at the related problem of releasing inner products. We initially prove the lower bound by fixing $\epsilon$ to $1/2$. We start by proving an extension of Lemma 4.2 to general differentially private algorithms. Let $\mathcal{F}$ be a function class, and let $\mathcal{A}$ be an $(\epsilon, \delta)$-differentially private algorithm for $\mathcal{F}$. Let $\mathcal{A}(D) \approx_{1/2, \delta} \mathcal{A}(D')$, and let $\Delta = \mathcal{F}(D') - \mathcal{F}(D)$. Unlike in the instance-independent case (Lemma 4.2) both $\mathbb{E}[\langle \mathcal{A}(D) - \mathcal{F}(D), \Delta \rangle^2]$ and $\mathbb{E}[\langle \mathcal{A}(D') - \mathcal{F}(D'), \Delta \rangle^2]$ needn't be $\|\Delta\|^4$, but the following lemma shows that at least one of them is $\|\Delta\|^4$ (this one-sided behavior is in fact unavoidable and is explained in the full version [20]).

LEMMA 4.7. *Let $\mathcal{F}$ be a function class of Boolean predicates, and let $\mathcal{A}$ be a $(1/2, \delta)$-differentially private algorithm for $\mathcal{F}$. Let $\mathcal{A}(D) \approx_{1/2, \delta} \mathcal{A}(D')$. Let $\Delta = \mathcal{F}(D') - \mathcal{F}(D)$. Then, at least one of $\mathbb{E}[\langle \mathcal{A}(D) - \mathcal{F}(D), \Delta \rangle^2]$ or $\mathbb{E}[\langle \mathcal{A}(D') - \mathcal{F}(D'), \Delta \rangle^2]$ is $\Omega(\langle \Delta, \Delta \rangle^2 (1 - \delta)^2)$.*

For a database $D \in (\{-1, 1\}^d)^n$, consider the $n$ neighboring databases[5] $\widetilde{D}_1, \ldots, \widetilde{D}_n$ where $\widetilde{D}_i$ is obtained by replacing $i$th row of $D$ by $(1)^d$. Let $T_k(D) = \{z_1, \ldots, z_n\}$ denote the (multi) set such that $\mathcal{I}_k(\widetilde{D}_i) - \mathcal{I}_k(D) = \tilde{z}_i$ and $z_i = o - \tilde{z}_i$. Let $\pi = \mathbb{I}_{m_k} - oo^\top/\langle o, o \rangle$ be an orthogonal projection matrix. Notice that, $\pi \tilde{z}_i = -\pi z_i$. For the reasons same as in the instance-independent case, we analyze projections onto

--------
[5]If $D$ has a row of $(1)^d$ (say the $i$th), then $D = \widetilde{D}_i$. For uniformity, we will still treat $D$ and $\widetilde{D}_i$ as neighbors.

$\pi \tilde{z}_i$ (or equivalently $-\pi z_i$). Let $u_{z_i}$ be the unit vector corresponding to $z_i$. Define

$$S_k(D) = \{z \in T_k(D)| \mathbb{E}[(\pi z)^\top \Sigma_{\mathcal{A}}(D)(\pi z)] = \Omega(m_k^2(1 - \delta)^2)\},$$
$$U_k(D) = \sum_{z \in S_k(D)} u_z u_z^\top \quad \text{and} \quad V_k(D) = \sum_{z \in T_k(D)} u_z u_z^\top.$$

For a database $D$ it is possible that the expected squared length of the projection of $\mathcal{A}(D) - \mathcal{I}_k(D)$ onto $\Delta = \mathcal{I}_k(D') - \mathcal{I}_k(D)$ is $o(\|\Delta\|^2(1 - \delta)^2)$ for many neighbors $D'$ of $D$. For example, it could happen that the expected squared length of the projection $\mathcal{A}(D') - \mathcal{I}_k(D')$ onto $\Delta$ is $\Omega(\|\Delta\|^2(1 - \delta)^2)$, whereas the expected squared length of the projection of $\mathcal{A}(D) - \mathcal{I}_k(D)$ onto $\Delta$ is only $o(\|\Delta\|^2(1 - \delta)^2)$. To overcome this problem we use random databases. Let $D_r$ be a database drawn uniformly at random from $(\{-1, 1\}^d)^n$. The idea is as follows,

$$\Omega(m_k^2(1 - \delta)^2|S_k(D_r)|) = \sum_{z \in S_k(D_r)} \mathbb{E}[\langle \mathcal{A}(D_r) - \mathcal{I}_k(D_r), \pi z \rangle^2]$$
$$= \sum_{z \in S_k(D_r)} (\pi z)^\top \Sigma_{\mathcal{A}}(D_r)(\pi z) = \sum_{z \in S_k(D_r)} tr(\pi^\top \Sigma_{\mathcal{A}}(D_r)\pi z z^\top)$$
$$= tr(\pi^\top \Sigma_{\mathcal{A}}(D_r)\pi \sum_{z \in S_k(D_r)} z z^\top) = tr(\pi^\top \Sigma_{\mathcal{A}}(D_r)\pi m_k U_k(D_r)).$$

The last step follows from the definition of $U_k(D_r)$. Now, since $\pi$, $\Sigma_{\mathcal{A}}(D_r)$, and $U_k(D_r)$ are all positive semidefinite, therefore,

$$tr(\pi^\top \Sigma_{\mathcal{A}}(D_r)\pi U_k(D_r)) \leq tr(\Sigma_{\mathcal{A}}(D_r))\|\pi\|_\infty \|U_k(D_r)\|_\infty$$
$$= tr(\Sigma_{\mathcal{A}}(D_r))\|U_k(D_r)\|_\infty.$$

Putting together,

$$\Omega(m_k^2(1 - \delta)^2|S_k(D_r)|) = m_k tr(\Sigma_{\mathcal{A}}(D_r))\|U_k(D_r)\|_\infty.$$

Therefore, to lower bound $tr(\Sigma_{\mathcal{A}}(D_r))$, we need good upper bound on the largest eigenvalue of $U_k(D_r)$. Lemma 4.8 does that by using the matrix-valued Chernoff bound from Ahlswede and Winter [2].

LEMMA 4.8. *For all $D \in (\{-1, 1\}^d)^n$, we have $\|U_k(D)\|_\infty \leq \|V_k(D)\|_\infty$, and with probability at least $1 - 1/n$ over the choice of $D_r$, $\|V_k(D_r)\|_\infty = O(\max\{n/m_k, \log m_k\})$.*

Let $\widetilde{\mathcal{D}}$ be the set of all databases from $(\{-1, 1\}^d)^n$ which have at least one row of $(1)^d$. If with high probability, $|S_k(D_r)|$ is $\Omega(n)$, then the upper bound on $\|U_k(D_r)\|_\infty$ from Lemma 4.8 will give us the lower bound on $tr(\Sigma_{\mathcal{A}}(D_r))$. But it is not necessary that $|S_k(D_r)| = \Omega(n)$ (as in constructing $S_k(D_r)$ we only consider the neighbors of $D_r$ belonging to $\widetilde{\mathcal{D}}$). In that case, we show that one could pick a database $\widetilde{D}_r$ uniformly at random from $\widetilde{\mathcal{D}}$, and use a similar analysis to lower bound $tr(\Sigma_{\mathcal{A}}(\widetilde{D}_r))$. We show that if the expected size of $S_k(D_r)$ is small (less than $n/4$), then the expected size of $\widetilde{S}_k(\widetilde{D}_r)$ is greater than $n/4$ (i.e., at least one of the two expected sizes is greater than $n/4$). The following proposition uses these ideas to show that every differentially private algorithm with probability $\Omega(1 - 1/n)$ needs to add a lot of noise to either $\mathcal{I}_k(D_r)$ or $\mathcal{I}_k(\widetilde{D}_r)$.

PROPOSITION 4.9. *Let $\mathcal{A}$ be a $(1/2, \delta)$-differentially private algorithm for $\mathcal{I}_k$. Let $D_r$ be a database chosen uniformly at random from $(\{-1, 1\}^d)^n$ and $\widetilde{D}_r$ be a database*

*chosen uniformly at random from* $\widetilde{\mathcal{D}}$. *Then, with probability* $\Omega(1 - 1/n)$, *at least one of* $tr(\Sigma_{\mathcal{A}}(D_r))$ *or* $tr(\Sigma_{\mathcal{A}}(\widetilde{D}_r))$ *is* $\Omega(\min\{m_k^2(1-\delta)^2, nm_k(1-\delta)^2/(\log m_k)\})$.

This previous lower bound for $(1/2, \delta)$-differentially private algorithms can be converted into a lower bound for $(\epsilon, \delta)$-differentially private algorithms. Details are deferred to the full version [20]. We now summarize the result.

THEOREM 4.10. *Let* $m_k = \binom{d}{k}$. *Any algorithm* $\mathcal{A}$ *for releasing all k-way marginal tables (or equivalently all k-way conjunction predicates) that for every database* $D \in (\{0,1\}^d)^n$ *has a root mean squared error of*

$$o(\min\{\sqrt{m_k}(1 - \delta/\epsilon)/(2^k\epsilon), \sqrt{n}(1-\delta/\epsilon)/(2^k\sqrt{\epsilon\log m_k})\})$$

*for each entry of* $\mathcal{A}(D)$ *is not* $(\epsilon, \delta)$-*differentially private.*

# 5. REFERENCES

[1] R. Agrawal and R. Srikant. Privacy-preserving data mining. In *SIGMOD*, pages 439–450. ACM, 2000.

[2] R. Ahlswede and A. Winter. Strong converse for identification via quantum channels. *IEEE Transactions on Information Theory*, 48(3), 2002.

[3] B. Barak, K. Chaudhuri, C. Dwork, S. Kale, F. McSherry, and K. Talwar. Privacy, accuracy, and consistency too: a holistic solution to contingency table release. In *PODS*, pages 273–282. ACM, 2007.

[4] Y. Bishop, S. Fienberg, and P. Holland. *Discrete Multivariate Analysis: Theory and Practice*. MIT Press, Cambridge MA, 1975.

[5] A. Blum, C. Dwork, F. McSherry, and K. Nissim. Practical privacy: The SuLQ framework. In *PODS*, pages 128–138. ACM, 2005.

[6] A. Blum, K. Ligett, and A. Roth. A learning theory approach to non-interactive database privacy. In *STOC*, pages 609–618. ACM, 2008.

[7] B.-C. Chen, R. Ramakrishnan, and K. LeFevre. Privacy skyline: Privacy with multidimensional adversarial knowledge. In *VLDB*, pages 770–781, 2007.

[8] I. Dinur, C. Dwork, and K. Nissim. Revealing information while preserving privacy, full version of [9], in preparation, 2010.

[9] I. Dinur and K. Nissim. Revealing information while preserving privacy. In *PODS*, pages 202–210. ACM, 2003.

[10] C. Dwork. Differential privacy: A survey of results. In *TAMC*, pages 1–19. Springer, 2008.

[11] C. Dwork. The differential privacy frontier. In *TCC*, pages 496–502. Springer, 2009.

[12] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *TCC*, pages 265–284. Springer, 2006.

[13] C. Dwork, F. McSherry, and K. Talwar. The price of privacy and the limits of LP decoding. In *STOC*, pages 85–94. ACM, 2007.

[14] C. Dwork, M. Naor, O. Reingold, G. Rothblum, and S. Vadhan. On the complexity of differentially private data release: efficient algorithms and hardness results. In *STOC*, pages 381–390, 2009.

[15] C. Dwork and S. Yekhanin. New efficient attacks on statistical disclosure control mechanisms. In *CRYPTO*, pages 469–480. Springer, 2008.

[16] A. Evfimievski, R. Srikant, R. Agrawal, and J. Gehrke. Privacy preserving mining of association rules. In *KDD*, pages 217–228. ACM, 2002.

[17] S. Fienberg and L. C. R. J. Willenborg. Special issue on disclosure control. *Journal of Official Statistics*, 14(4), 1998.

[18] A. Ghosh, T. Roughgarden, and M. Sundararajan. Universally utility-maximizing privacy mechanisms. In *STOC*, pages 351–360. ACM, 2009.

[19] M. Hardt and K. Talwar. On the geometry of differential privacy. In *these proceedings*. ACM, 2010.

[20] Full version of this abstract. Available from the authors, 2010.

[21] M. Ledoux. The concentration of measure phenomenon, volume 89 of Mathematical Surveys and Monographs. *American Mathematical Society*, 208:2005–2006, 2001.

[22] N. Li, T. Li, and S. Venkatasubramanian. *t*-closeness: Privacy beyond *k*-anonymity and *l*-diversity. In *ICDE*, pages 106–115. IEEE, 2007.

[23] A. Machanavajjhala, J. Gehrke, D. Kifer, and M. Venkitasubramaniam. *l*-diversity: Privacy beyond *k*-anonymity. *ACM TKDD*, 1(1):3, 2007.

[24] D. J. Martin, D. Kifer, A. Machanavajjhala, J. Gehrke, and J. Y. Halpern. Worst-case background knowledge for privacy-preserving data publishing. In *ICDE*, pages 126–135. IEEE, 2007.

[25] F. McSherry and K. Talwar. Mechanism design via differential privacy. In *FOCS*, pages 94–103. IEEE, 2007.

[26] V. Rastogi, S. Hong, and D. Suciu. The boundary between privacy and utility in data publishing. In *VLDB*, pages 531–542, 2007.

[27] M. Rudelson and R. Vershynin. The least singular value of a random square matrix is $O(n^{-1/2})$. *Comptes rendus-Mathématique*, 2008.

[28] M. Rudelson and R. Vershynin. The Littlewood–Offord problem and invertibility of random matrices. *Advances in Mathematics*, 218(2):600–633, 2008.

[29] M. Rudelson and R. Vershynin. The smallest singular value of a random rectangular matrix. *Communications on Pure and Applied Mathematics*, pages 1707 – 1739, 2009.

[30] C. Skinner and N. Shlomo. Assessing identification risk in survey microdata using log-linear models. *Journal of the American Statistical Association*, 103(483):989–1001, 2008.

[31] A. Smith. Efficient, differentially private point estimators. *CoRR*, abs/0809.4794, 2008.

[32] L. Sweeney. *k*-anonymity: A model for protecting privacy. *International Journal on Uncertainty, Fuzziness and Knowledge-based Systems*, 10(5):557–570, 2002.

[33] S. L. Warner. Randomized response: A survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 60(309):63–69, 1965.

[34] X. Xiao and Y. Tao. M-invariance: towards privacy preserving re-publication of dynamic datasets. In *SIGMOD*, pages 689–700. ACM, 2007.