# Optimal Private Halfspace Counting via Discrepancy

S. Muthukrishnan    *Aleksandar Nikolov*

Rutgers U.

# Private Range Counting

- **Public Input**: A *ground set* $P \subseteq \mathbb{R}^d$; a *range space*, i.e. collection of sets $\mathcal{R} \subseteq 2^P$ induced by some natural geometric sets
- **Private input**: Integer weight $x_p$ for each $p \in P$.
- **Goal**: For all ranges $R \in \mathcal{R}$, approximate privately

$$R(\mathbf{x}) = \sum_{p \in R} x_p$$

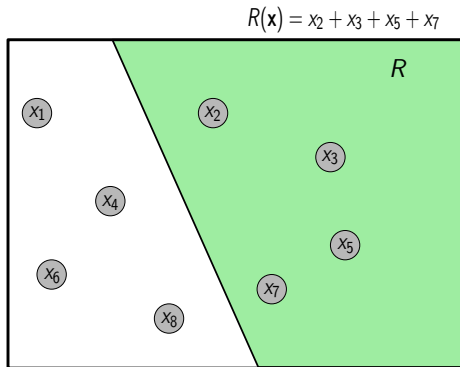- **Accuracy**: *Mean squared error of an algorithm* $\mathcal{M}$ is

$$\frac{1}{|\mathcal{R}|} \sum_{R \in \mathcal{R}} \left( R(\mathbf{x}) - \mathcal{M}(R, \mathbf{x}) \right)^2$$

## Halfspace Counting

Each $R \in \mathcal{R}$ is the points of $P$ contained in some halfspace in $\mathbb{R}^d$.

Query: what is the total weight of all points of $P$ in halfspace $R$?

Fundamental in Computational Geometry. Other range queries can be expressed as halfspace queries by "lifting" them to a higher dimension.



$R(\mathbf{x}) = x_2 + x_3 + x_5 + x_7$

## Private Linear Queries

More general algebraic problem:

- **Public Input**: A query matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$
  - In range counting: each row of $\mathbf{A}$ is the indicator of a range

- **Private Input**: A vector $\mathbf{x} \in \mathbb{Z}^n$
  - In range counting: the private point weights

- **Goal**: An algorithm $\mathcal{M}$ that approximates $\mathbf{Ax}$ and satisfies a privacy guarantee ($(\varepsilon, \delta)$-differential privacy).

- **Accuracy**: *Mean squared error* is $\frac{1}{m}\|\mathbf{Ax} - \mathcal{M}(\mathbf{A}, \mathbf{x})\|_2^2$

# Differential Privacy

### Definition

An algorithm $\mathcal{M}$ with input domain $\mathbb{Z}^n$ and output range $Y$ is $(\varepsilon, \delta)$-*differentially private* if for every $n$, every $\mathbf{x}, \mathbf{x}'$ with $\|\mathbf{x} - \mathbf{x}'\|_1 \leq 1$, and every measurable $S \subseteq Y$, $\mathcal{M}$ satisfies

$$\Pr[\mathcal{M}(\mathbf{x}) \in S] \leq e^{\varepsilon} \Pr[\mathcal{M}(\mathbf{x}') \in S] + \delta.$$

# What is Known about Halfspace Counting?

- **Lower Bounds**:
    - $\Omega(n)$ squared error necessary for arbitrary 0-1 **A** when $m > n$ [DN03]
    - *Does not apply to halfspace counting!* No superconstant lower bound known.
- **Upper Bounds**:
    - Randomized response gives $O(n \log m)$.
    - For halfspaces $m = O(n^d)$, therefore $O(nd \log n)$ error is sufficient.

# Our Results

- **Lower bounds**
  - Private halfspace counting in $\mathbb{R}^d$ requires $\Omega(n^{1-1/d})$ mean squared error.
  - *More generally*: linear queries **A** require noise lower bounded by the *(hereditary) combinatorial discrepancy* of **A** (up to a log factor).

- **Upper bounds**
  - Halfspace counting can be approximated privately with $O(n^{1-1/d})$ mean squared error.
  - *More generally*: range counting for ranges with shatter functions exponent $d$ can be approximated with the same error.
  - Bounds also extend to *worst case error* (up to polylog factors).

Both results use discrepancy theory.

## Lower bound: Dinur-Nissim attack

**Assume**: There exists $\mathcal{M}$ such that for any $\mathbf{x}$ w.h.p.
$\|\mathbf{A}\mathbf{x} - \mathcal{M}(\mathbf{A}, \mathbf{x})\|_2 \leq E$.

**Adversary's Goal**: Given output of $\mathcal{M}(\mathbf{A}, \mathbf{x})$, compute $\mathbf{x}'$, $\|\mathbf{x} - \mathbf{x}'\|_1 \ll n$.
So $\mathcal{M}$ is not private.

**Procedure**: Output any $\mathbf{x}'$ s.t. $\|\mathbf{A}\mathbf{x}' - \mathcal{M}(\mathbf{A}, \mathbf{x})\|_2 \ll E$ (succeeds w.p. $1 - \beta$).

We have $\|\mathbf{A}\mathbf{x} - \mathbf{A}\mathbf{x}'\|_2 \leq \|\mathbf{A}\mathbf{x}' - \mathcal{M}(\mathbf{A}, \mathbf{x})\|_2 + \|\mathbf{A}\mathbf{x} - \mathcal{M}(\mathbf{A}, \mathbf{x})\|_2 \ll E$.

**Needed**: $E$ such that $\|\mathbf{A}\mathbf{x}' - \mathbf{A}\mathbf{x}\|_2 \ll E \Rightarrow \|\mathbf{x} - \mathbf{x}'\|_1 \ll n$.

## Discrepancy connection

**Discrepancy**: The adversary can succeed when

$$E \ll \mathrm{disc}_\alpha(\mathbf{A}) = \min_{\substack{\mathbf{b} \in \{0, \pm 1\}^n \\ \|\mathbf{b}\|_1 \geq \alpha n}} \|\mathbf{A}\mathbf{b}\|_2$$

When $\alpha = 0$, this is trivially 0.
When $\alpha = 1$, this is the classical combinatorial $\ell_2$ discrepancy.

Can we connect $disc_\alpha$ to $\mathrm{disc}_1$ when $\alpha \in (0, 1)$?

## A More Robust Lower Bound

$\mathsf{herdisc}_\alpha(\mathbf{A}) = \max_{S \subseteq [n]} \mathsf{disc}_\alpha(\mathbf{A}|_S)$

*Weaker success condition* for the adversary: choose a subset $S$ of $[n]$ (based only on $\mathbf{A}$) and then guess most of $\mathbf{x}$ restricted to $S$:

- still implies a contradiction with $(\varepsilon, \delta)$-differential privacy
- adversary can succeeds when $E \ll \mathsf{herdisc}_\alpha(\mathbf{A})$
- $\mathsf{herdisc}_\alpha(\mathbf{A}) \geq \mathsf{herdisc}_1(\mathbf{A})/O(\log n)$ (for constant $\alpha$)

# Putting it together

Theorem (Main Lower Bound)

*No algorithm $\mathcal{M}$ that satisfies*

$$\forall \mathbf{x} \in \{0,1\}^n : \Pr[\|\mathbf{Ax} - \mathcal{M}(\mathbf{A}, \mathbf{x})\|_2 = o(\text{herdisc}_1(\mathbf{A})/\log n)] \geq 1 - \beta,$$

*is $(\varepsilon, \delta)$-differentially private for $\varepsilon = O(1)$, and constant $\delta < 1$ and $\beta < 1$.*

Halfspace counting:

- Mean squared error for private halfspace queries is $\Omega(n^{1-1/d}/\log n)$
- Using the hereditary structure of halfspace range spaces, we can show mean squared error is $\Omega(n^{1-1/d})$.

# Two Tools: Input and Output Perturbation

- **Input perturbation**: Compute $\tilde{\mathbf{x}} = \mathbf{x} + \mathrm{Lap}(1/\varepsilon)^n$ and output $\mathbf{A}\tilde{\mathbf{x}}$.

- **Output perturbation**: Output $\mathbf{A}\mathbf{x} + \mathrm{Lap}(1/\varepsilon')^m$ for $\varepsilon'$ chosen to satisfy $(\varepsilon, \delta)$-differential privacy.

## When Do the Tools Work?

For range counting:

- input perturbation works well with small ranges (squared error linear in size of range)

- output perturbation works well when each point belongs to few ranges (squared error linear in maximum degree)

But for halfspaces most ranges are large and most points belong to many ranges.

*Solution from discrepancy theory*: halfspace ranges admit a nice decomposition [Mat95]. (works for range spaces with VC dimension $d$ and shatter function exponent $d$)

## Decomposition

Decompose $\mathcal{R}$ into a series of new range spaces $\{\mathcal{T}_i\}_{i=1}^{\log n}$ such that *approximating counts for each $\mathcal{T}_i$ gives the counts for $\mathcal{R}$.*

$\mathcal{R}$ is decomposed into:

- $\mathcal{T}_i$ with many small sets ($i$ large): can use input perturbation
- $\mathcal{T}_i$ with few large sets ($i$ small): can use output perturbation

Do we achieve the right balance? *No!*

- Values of $i$ s.t. noise variances is $O(n^{1-1/d})$:

Output perturbation          Input perturbation



$$i_1 = \frac{\log n}{d} - \frac{\log n}{d^2} \qquad\qquad i_0 = \frac{\log n}{d}$$

## How to Make It work

For $i \in (i_1, i_0)$:

- For any $\mathcal{T}_i$, there are points $p$ that belong to a lot of sets and incur large privacy loss
  - i.e. we need noise with variance $\Omega(n)$ to preserve their privacy

- But we control both maximum set size and number of sets in $\mathcal{T}_i$!
- *Idea*: use *average* privacy loss (privacy loss averaged over all $p$)
  - The "average" $p$ requires only $O(n^{1-1/d})$ noise to preserve its privacy

- We find a set $X$ s.t.
  - the privacy of each $p \in X$ can be preserved by noise with variance $O(n^{1-1/d})$
  - $|X| \geq |P|/2$.

A *partial coloring* style algorithm:

- For $i \geq i_0$ use input perturbation to approximate counts for $\mathcal{T}_i$ **w.r.t.** $X$

- For $i < i_0$, we add Laplace noise with variance $O(n^{1-1/d}2^{(i_0-i)(1-d)})$ to approximate counts for $\mathcal{T}_i$ **w.r.t.** $X$

- This allows us to compute halfspace counts over $X$ with squared error $O(n^{1-1/d})$.

- Recurse on $P \setminus X$ (still a halfspace range space)

*This work*:

- Optimal upper and lower bounds for private halfspace counting
- Connection between discrepancy theory and noise lower bounds for differential privacy

*Other results*: A lower bound of $\Omega((\log n)^{d-1})$ for *orthogonal range counting*. Tight up to the dependence on $d$.

*Open question*: Does discrepancy always characterize the error needed to preserve privacy of linear queries?

**Thank you!**

📄 Avrim Blum, Katrina Ligett, and Aaron Roth.
A learning theory approach to non-interactive database privacy.
In *Proceedings of the 40th annual ACM symposium on Theory of computing*, STOC '08, pages 609–618, New York, NY, USA, 2008. ACM.

📄 Irit Dinur and Kobbi Nissim.
Revealing information while preserving privacy.
In *Proceedings of the twenty-second ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, PODS '03, pages 202–210, New York, NY, USA, 2003. ACM.

📄 C. Dwork, G. N. Rothblum, and S. Vadhan.
Boosting and differential privacy.
In *Proc. 51st Annual IEEE Symp. Foundations of Computer Science (FOCS)*, pages 51–60, 2010.

📄 M. Hardt and G. N. Rothblum.
A multiplicative weights mechanism for privacy-preserving data analysis.

In *Proc. 51st Annual IEEE Symp. Foundations of Computer Science (FOCS)*, pages 61–70, 2010.

📄 J. Matoušek.
Tight upper bounds for the discrepancy of half-spaces.
*Discrete and Computational Geometry*, 13(1):593–601, 1995.