

Differential Privacy in Applications

Jiaqi Shao

The Chinese University of Hong Kong, Shenzhen

June 27, 2022



PACK

1 Plan

- Frequent Pattern Mining with **heterogeneous data types**, i.e. string, numeric
- Federated Analytics with **relative small sample size**, i.e. $n = 10^3, 10^4$

2 Achievement

- Familiar with differential privacy
- Study different secure computation schemes
- Improve TrieHH **infeasible problem** by reasonably increasing privacy tolerance.
- Compare different **prefix tree structures** to implement TrieHH algorithm.

3 Challenge

- TrieHH achieves (ϵ, δ) -DP with **large privacy budget**

4 Knowledge

- **Sampling and threshold**: A simple sample-and-threshold approach provides an (ϵ, δ) -DP guarantee for histograms
- **Prefix Tree (Trie)**: predict words or bit strings.

PACK

1 Plan

- Frequent Pattern Mining with **heterogeneous data types**, i.e. string, numeric
- Federated Analytics with **relative small sample size**, i.e. $n = 10^3, 10^4$

2 Achievement

- Familiar with differential privacy
- Study different secure computation schemes
- Improve TrieHH **infeasible problem** by reasonably increasing privacy tolerance.
- Compare different **prefix tree structures** to implement TrieHH algorithm.

3 Challenge

- TrieHH achieves (ϵ, δ) -DP with **large privacy budget**

4 Knowledge

- **Sampling and threshold**: A simple sample-and-threshold approach provides an (ϵ, δ) -DP guarantee for histograms
- **Prefix Tree (Trie)**: predict words or bit strings.

Differential Privacy: Definition

Definition

A function \mathcal{M} is often called a mechanism satisfying (approximate) differential privacy, if for all neighboring datasets \mathcal{D} and \mathcal{D}' , and all possible outputs $\mathcal{S} \subseteq \mathcal{M}(\cdot)$:

$$\mathcal{P}[\mathcal{M}(\mathcal{D}) \in \mathcal{S}] \leq e^{\varepsilon} \times \mathcal{P}[\mathcal{M}(\mathcal{D}') \in \mathcal{S}] (+\delta)$$

where ε is the **privacy budget**, and δ represents a “**failure probability**”

- Key Idea: With given ε, δ , we can add **noise** to the output to preserve privacy, e.g Laplace and Gaussian Mechanisms.

Differential Privacy: Properties

Sequential Composition

Let \mathcal{M}_i each provides ε_i -differential privacy. The sequence of $\mathcal{M}_i(\mathcal{D})$: $\mathcal{M}(\mathcal{D}) = (\mathcal{M}_1, \mathcal{M}_2, \dots)$ provides $\sum_i \varepsilon_i - DP$

Parallel Composition

If \mathcal{M} satisfies $\varepsilon - DP$, and split dataset \mathcal{D} into k disjoint chunks $d_1 \cup d_2 \cup \dots \cup d_k = \mathcal{D}$. Then $\mathcal{M}(d_i)$ achieves $\varepsilon - DP$

Post-processing

If $\mathcal{M}(\mathcal{D})$ achieves $\varepsilon - DP$, then for any function g , $g(\mathcal{M}(\mathcal{D}))$ achieves $\varepsilon - DP$.

PACK

1 Plan

- Frequent Pattern Mining with **heterogeneous data types**, i.e. string, numeric
- Federated Analytics with **relative small sample size**, i.e. $n = 10^3, 10^4$

2 Achievement

- Familiar with differential privacy
- Study different secure computation schemes
- Improve TrieHH **infeasible problem** by reasonably increasing privacy tolerance.
- Compare different **prefix tree structures** to implement TrieHH algorithm.

3 Challenge

- TrieHH achieves (ϵ, δ) -DP with **large privacy budget**

4 Knowledge

- **Sampling and threshold**: A simple sample-and-threshold approach provides an (ϵ, δ) -DP guarantee for histograms
- **Prefix Tree (Trie)**: predict words or bit strings.

Secure Computation Schemes

1. Secret Sharing: Key Idea

Secret sharing schemes split a secret into multiple shares that are **meaningless unless** τ (threshold) of them are **collected** and the secret is reconstructed.

2. Homomorphic Encryption: Key Idea

Homomorphic encryption schemes allow users' data to be **protected anytime** it is sent to the cloud, because it can allow operations and functions to be preformed over **encrypted data**.

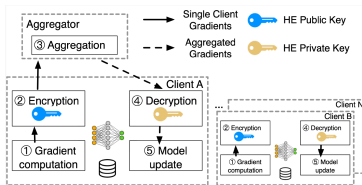


Figure: Homomorphic Encryption applied in Federated Learning

PACK

1 Plan

- Frequent Pattern Mining with **heterogeneous data types**, i.e. string, numeric
- Federated Analytics with **relative small sample size**, i.e. $n = 10^3, 10^4$

2 Achievement

- Familiar with differential privacy
- Study different secure computation schemes
- Improve TrieHH **infeasible problem** by reasonably increasing privacy tolerance.
- Compare different **prefix tree structures** to implement TrieHH algorithm.

3 Challenge

- TrieHH achieves (ϵ, δ) -DP with **large privacy budget**

4 Knowledge

- **Sampling and threshold**: A simple sample-and-threshold approach provides an (ϵ, δ) -DP guarantee for histograms
- **Prefix Tree (Trie)**: predict words or bit strings.

TrieHH Overview¹

Research Question: How to develop an interactive **heavy hitters discovery algorithm** that achieves **central DP** while **minimizing the data collected** from users?

Model:

- Following the Federated Learning protocol: (1) Sampling m clients per round. (2) Execute the TrieHH algorithm and the sever aggregates for multiple rounds; and (3) Server broadcasts the result.
- Achieving DP through **sampling** and **mechanism** which outputs a Trie.

Contributions

- Achieve central DP **without centralizing raw data** and **adding noise**
- Obtain **excellent utility** compared with local DP

¹Wennan Zhu et al. "Federated Heavy Hitters Discovery with Differential Privacy". In: *CoRR* abs/1902.08534 (2019). arXiv: 1902.08534. URL: <http://arxiv.org/abs/1902.08534>.

TrieHH Insights and Improvements

Insights

The author used [Corollary 1](#) to take out experiments.

Corollary 1

To achieve (ε, δ) -differential privacy, set $\gamma = \left(e^{\frac{\varepsilon}{L}} - 1\right) \sqrt{n} / \left(\theta e^{\frac{\varepsilon}{L}}\right)$ and $\theta = \max \left\{10, \left\lceil e^W (C_\delta) + 1 - \frac{1}{2} \right\rceil, \left\lceil e^{\frac{\varepsilon}{L}} - 1 \right\rceil \right\}$, where W is the Lambert W function and $C_\delta = e^{-1} \ln \left(\frac{8}{7\sqrt{2\pi}} \delta^{-1} \right)$. Further, when $n \geq 10^4$, choosing $\theta = \lceil \log_{10} n + 6 \rceil$ ensures that Algorithm 1 is $(\varepsilon, \frac{1}{300n})$ -differential private.

- The Corollary 1 is derived from the Theorem 1 (using approximation techniques)
 - \Rightarrow Corollary 1 results in **larger** failure probability.
 - \Rightarrow Using Theorem 1 to compute θ, γ

TrieHH Insights and Improvements

Theorem 1

When $4 \leq \theta \leq \sqrt{n}$ (threshold) and $1 \leq \gamma \leq \frac{\sqrt{n}}{\theta+1}$ (sampling size $m = \gamma\sqrt{n}$, where n is the total client size),
TrieHH algorithm is $(L\ln(1 + \frac{1}{\frac{\sqrt{n}}{\gamma\theta}-1}), \frac{\theta-2}{(\theta-3)\theta!})$ -differential private.

Goal: given $\varepsilon, \delta = 1/n^2, n$, calculate γ, θ . By DP's definition:

$$L\ln(1 + \frac{1}{\frac{\sqrt{n}}{\gamma\theta}-1}) \leq \varepsilon \quad (1)$$

$$\frac{\theta-2}{(\theta-3)\theta!} \leq \delta \quad (2)$$

TrieHH Insights and Improvements

- *Observation 1*: By Eq.2, $\frac{(\theta-3)\theta!}{\theta-2} > 1/\delta = n^2$, where LHS is incremental function.
- *Observation 2*: By Eq. 1, $\gamma \leq \frac{e^{\varepsilon/L}-1}{\theta * e^{\varepsilon/L}} \sqrt{n}$, and by Theorem 1, $\gamma \in [1, \frac{\sqrt{n}}{\theta+1}]$, to obtain a **feasible** γ :

$$1 \leq \frac{e^{\varepsilon/L} - 1}{\theta * e^{\varepsilon/L}} \sqrt{n}$$

$$\Rightarrow \theta \leq \frac{e^{\varepsilon/L}-1}{e^{\varepsilon/L}} \sqrt{n} \leq \sqrt{n} \text{ (THETACEIL: } \frac{e^{\varepsilon/L}-1}{e^{\varepsilon/L}} \sqrt{n})$$

Compute θ, γ

- 1 Initially set $\theta = 4$
 - ▶ If $\theta > \text{THETACEIL}(\varepsilon \text{ underflow})$, $\varepsilon \leftarrow \text{THETACEIL} = 4$
- 2 MINIMIZE $_{\theta}$ $h(\theta) = \frac{(\theta-3)\theta!}{\theta-2} > 1/\delta$ ($h(\theta)$ is incremental)
 - ▶ If $\theta > \text{THETACEIL}$ (θ overflow), $\theta \leftarrow \text{THETACEIL}$, $\delta \leftarrow 1/h(\theta)$
- 3 $\gamma \leftarrow \text{MIN}(\text{THETACEIL}/\theta, \frac{\sqrt{n}}{\theta+1})$

TrieHH Insights and Improvements

	$\varepsilon = 1$	$\varepsilon = 2$	$\varepsilon = 3$
$n = 1e3$	infeasible	infeasible	infeasible
$n = 5e3$	infeasible	infeasible	$\theta = 11$
$n = 1e4$	infeasible	$\theta = 12$	$\theta = 12$
$n = 2e4$	$\theta = 12$	$\theta = 12$	$\theta = 12$
$n = 3e4$	$\theta = 13$	$\theta = 13$	$\theta = 13$
$n = 1e5$	$\theta = 14$	$\theta = 14$	$\theta = 14$

Table: Choices of γ, θ to achieve $\varepsilon = 1, 2, 3$ and $\delta \leq 1/n^2$

Remark

Infeasible means θ overflow or ε underflow \Rightarrow reduce θ & increase δ or increase ε

TrieHH Insights and Improvements

```
Total number of clients: 10000
Theta overflow, infeasible r:: Reduce theta, and Increase delta
((1, 3.215020576131687e-06))-DP; Gamma used: 1.06
Theta: 9
Batch size used by TrieHH: 105
Discovered 1 heavy hitters in run #1
['the']
Total number of clients: 10000
Theta overflow, infeasible r:: Reduce theta, and Increase delta
((1, 3.215020576131687e-06))-DP; Gamma used: 1.06
Theta: 9
Batch size used by TrieHH: 105
Discovered 1 heavy hitters in run #2
['the']
```

```
Total number of clients: 10000
(1, 1e-08)-DP
Theta used by TrieHH: 12
Batch size used by TrieHH: 79
Discovered 1 heavy hitters in run #1
['the']
Discovered 0 heavy hitters in run #2
[]
Discovered 0 heavy hitters in run #3
[]
Discovered 0 heavy hitters in run #4
[]
Discovered 0 heavy hitters in run #5
[]
```

Figure: Finding heavy hitters under given $n = 10^4, \varepsilon = 1, \delta = 1/n^2$. (Left: decrease δ to solve δ overflow)

- **Remark 1:** Server goal is to **maximum utility**, i.e. finding heavy hitters. However, when client size is small, γ, θ are increased. (For example, $n = 1000 \Rightarrow \varepsilon = 31.4$, and $n = 2000 \Rightarrow \delta = 0.083$)
- **Remark 2:** $\varepsilon < 1$ is hard to achieve, since it positively relative to L (the max word length, default to 10)

PACK

1 Plan

- Frequent Pattern Mining with **heterogeneous data types**, i.e. string, numeric
- Federated Analytics with **relative small sample size**, i.e. $n = 10^3, 10^4$

2 Achievement

- Familiar with differential privacy
- Study different secure computation schemes
- Improve TrieHH **infeasible problem** by reasonably increasing privacy tolerance.
- Compare different **prefix tree structures** to implement TrieHH algorithm.

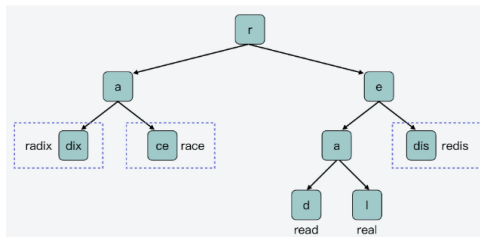
3 Challenge

- TrieHH achieves (ϵ, δ) -DP with **large privacy budget**

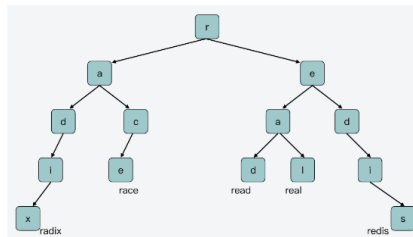
4 Knowledge

- **Sampling and threshold**: A simple sample-and-threshold approach provides an (ϵ, δ) -DP guarantee for histograms
- **Prefix Tree (Trie)**: predict words or bit strings.

Prefix Tree Structures



Radix Trie



Standard Trie

- Radix Tree is the **compact** vision of the standard trie.
- Radix Tree could reduce the tree depth to achieve **searching efficiency**.

PACK

1 Plan

- Frequent Pattern Mining with **heterogeneous data types**, i.e. string, numeric
- Federated Analytics with **relative small sample size**, i.e. $n = 10^3, 10^4$

2 Achievement

- Familiar with differential privacy
- Study different secure computation schemes
- Improve TrieHH **infeasible problem** by reasonably increasing privacy tolerance.
- Compare different **prefix tree structures** to implement TrieHH algorithm.

3 Challenge

- TrieHH achieves (ϵ, δ) -DP with **large privacy budget**

4 Knowledge

- **Sampling and threshold**: A simple sample-and-threshold approach provides an (ϵ, δ) -DP guarantee for histograms
- **Prefix Tree (Trie)**: predict words or bit strings.

PACK

1 Plan

- Frequent Pattern Mining with **heterogeneous data types**, i.e. string, numeric
- Federated Analytics with **relative small sample size**, i.e. $n = 10^3, 10^4$

2 Achievement

- Familiar with differential privacy
- Study different secure computation schemes
- Improve TrieHH **infeasible problem** by reasonably increasing privacy tolerance.
- Compare different **prefix tree structures** to implement TrieHH algorithm.

3 Challenge

- TrieHH achieves (ϵ, δ) -DP with **large privacy budget**

4 Knowledge

- **Sampling and threshold**: A simple sample-and-threshold approach provides an (ϵ, δ) -DP guarantee for histograms
- **Prefix Tree (Trie)**: predict words or bit strings.

Sample and Threshold Differential Privacy²

Research Question: How differential privacy can be obtained via a simple sample-and-threshold mechanism?

Model: Bernoulli sampling with threshold.

Contributions:

- **Bernoulli sampling** is sufficient to provide differential privacy.
- The resulting mechanism can also answer heavy hitter, quantile and range queries.
- The associated counts provide **accurate frequency estimates** for items from the input.

²Akash Bharadwaj and Graham Cormode. “Sample and Threshold Differential Privacy: Histograms and applications”. In: *CoRR* abs/2112.05693 (2021). arXiv: 2112.05693. URL: <https://arxiv.org/abs/2112.05693>.

Sample and Threshold Differential Privacy

Insights

Lemma 1

If we set the sampling rate $p_s = \alpha(1 - e^{-\varepsilon})$ for some $0 < \alpha \leq 1$ and $\varepsilon \leq 1$, then sample-and-threshold achieves (ε, δ) differential privacy for $\delta = \exp(-C_\alpha \tau)$, where $C_\alpha = \ln(1/\alpha) - 1/(1 + \alpha)$.

For example, for $\varepsilon = 1$ and $\alpha = 1/6$, the sampling rate is $p_s = 0.105 \approx 0.1$ and, choosing $\tau = 20$, $\delta < 10^{-8}$ using $C_\alpha = 0.935$

Sample and Threshold Differential Privacy

Lemma 2

The TrieHH++ protocol using L sample-and-threshold histograms with (ε, δ) -DP achieves an overall guarantee of $(L\varepsilon, L\delta)$ -DP.

Remark 1: Instead of proceeding in rounds, simply applying the basic histogram protocol to the full inputs (**without build the Trie**), and reporting the items which survive the threshold achieves (ε, δ) -DP (L can be dropped from these bounds.)

Remark 2: Hyperparameters are $\varepsilon, \theta, \alpha$
 \Rightarrow Sampling rate $p_s = \alpha(1 - e^{-\varepsilon})$
 $\Rightarrow \delta(\theta, \alpha)$

Question:

- Sample-and-threshold is not related to client size.
- Hyperparameters is not intuitive.

PACK

1 Plan

- Frequent Pattern Mining with **heterogeneous data types**, i.e. string, numeric
- Federated Analytics with **relative small sample size**, i.e. $n = 10^3, 10^4$

2 Achievement

- Familiar with differential privacy
- Study different secure computation schemes
- Improve TrieHH **infeasible problem** by reasonably increasing privacy tolerance.
- Compare different **prefix tree structures** to implement TrieHH algorithm.

3 Challenge

- TrieHH achieves (ϵ, δ) -DP with **large privacy budget**

4 Knowledge

- **Sampling and threshold**: A simple sample-and-threshold approach provides an (ϵ, δ) -DP guarantee for histograms
- **Prefix Tree (Trie)**: predict words or bit strings.

Dynamic Prefix Tree as a predict model³

Ideas

- Radix Trie can reduce deadly depth caused by long bit strings.
- Dynamic extend a Radix Trie within several rounds, then, clients **only** send a few prefix to infer heavy hitters in the rest rounds.

Remaining Questions:

- Fault Tolerance: client may lose connection.
- DP decision: ϵ, δ is decided by the sever, whose goal is maximizing the utility.

³Xiang Lisa Li and Percy Liang. "Prefix-Tuning: Optimizing Continuous Prompts for Generation". In: *CoRR* abs/2101.00190 (2021). arXiv: 2101.00190. URL: <https://arxiv.org/abs/2101.00190>.