

A Huber Loss Minimization Approach to Mean Estimation under User-Level Differential Privacy

Puning Zhao¹, Lifeng Lai², Li Shen³, Qingming Li⁴, Jiafei Wu¹, Zhe Liu¹

¹ Zhejiang Lab ² University of California, Davis ³ Sun Yat-Sen University ⁴ Zhejiang University

Abstract

- Distributed system requires privacy protection of users' entire contribution of samples.
- Existing solution (two-stage) is not suitable for imbalanced users or heavy-tailed distributions.
- This work: Huber loss minimization approach.
- The new method significantly improves the performance for imbalanced users, by adjusting the connecting points of Huber loss adaptively.
- The new method significantly improves the performance for heavy-tailed distributions, by replacing the clipping operation to a moderate Huber loss penalty.
- We conduct both theoretical analysis and experiments to validate the new method.

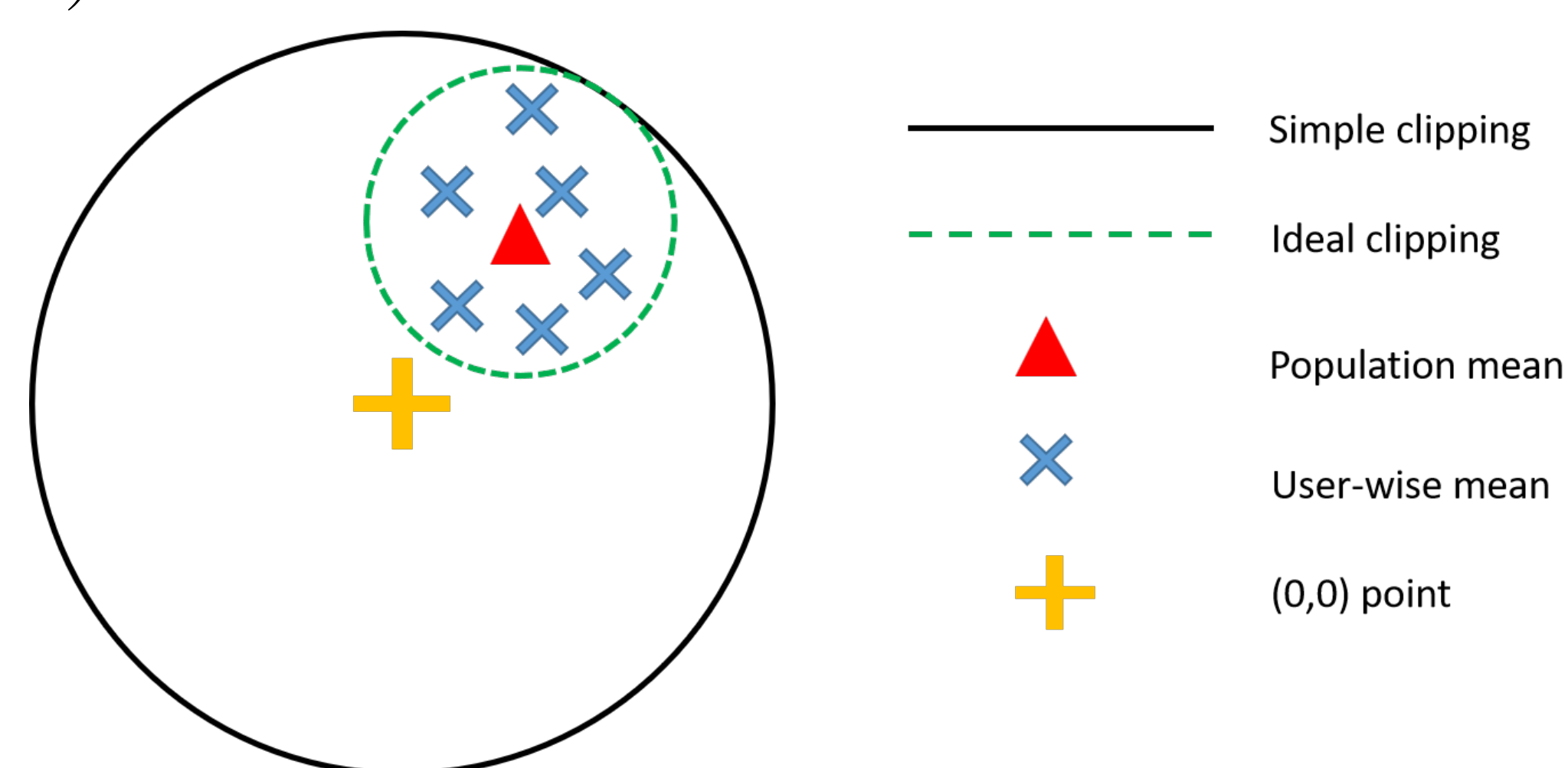
Introduction

Background

- Traditional differential privacy considers the privacy of each sample.
- Each user may contribute multiple items: In recommendation systems, an account is a user, and each visiting record can be viewed as an item. In federated learning, each client can be viewed as a user, and each sample can be viewed as an item
- We hope to protect a user's *entire contribution*.

Existing Solutions

1) *Direct calculation*:



- The local averages are already close to each other.
- Clipping radius is larger than necessary, resulting in unnecessary sacrifice of utility.

2) *Two-stage method (WME) [1]*:

- Stage I (Localization): identify a small interval that contains the truth μ with high probability
- State II (Refinement): Clip to the interval and then calculate final average with appropriate noise
- Extend to high dimensionality: Hadamard transform
- Limitation 1: Not suitable for imbalanced data
- Limitation 2: Not suitable for heavy-tailed distributions

Contributions

- Propose *Huber loss minimization* approach to address the limitations above
- Provide both theoretical analysis and numerical experiments
- Significant improvement for heavy-tailed distributions.

Reason: penalizing large distance yields smaller bias than simple clipping

- Significant improvement for imbalanced data.

Reason: Adaptive thresholds and weights, leading to better sensitivity-bias tradeoff

Preliminaries

Differential privacy (DP)

If for any $O \subseteq \Theta$ and any two adjacent datasets \mathcal{D} and \mathcal{D}'

$$P(\mathcal{A}(\mathcal{D}) \in O) \leq e^{\epsilon} P(\mathcal{A}(\mathcal{D}') \in O) + \delta, \quad (1)$$

then $\mathcal{A} : \Omega \rightarrow \Theta$ is (ϵ, δ) -DP

User-level DP

Two datasets $\mathcal{D}, \mathcal{D}'$ are user-level adjacent if they differ in items belonging to only one user. \mathcal{A} is user-level (ϵ, δ) -DP if (1) is satisfied for any two user-level adjacent datasets \mathcal{D} and \mathcal{D}' .

The Proposed Method

- Estimator without adding noise:

$$\hat{\mu}_0(\mathcal{D}) = \arg \min_{\mathbf{s}} \sum_{i=1}^n w_i \phi_i(\mathbf{s}, \mathbf{y}_i(\mathcal{D})), \quad (2)$$

in which w_i is the weight. ϕ_i is the Huber loss function:

$$\phi_i(\mathbf{s}, \mathbf{y}) = \begin{cases} \frac{1}{2} \|\mathbf{s} - \mathbf{y}\|^2 & \text{if } \|\mathbf{s} - \mathbf{y}\| \leq T_i \\ T_i \|\mathbf{s} - \mathbf{y}\| - \frac{1}{2} T_i^2 & \text{if } \|\mathbf{s} - \mathbf{y}\| > T_i. \end{cases} \quad (3)$$

- Final estimator:

$$\hat{\mu}(\mathcal{D}) = \text{Clip}(\hat{\mu}_0(\mathcal{D}), R_c) + \mathbf{W}, \quad (4)$$

Theoretical results

Theorem 1: Bounded support, balanced users

Under some assumptions (omitted here), for n users with m items per user,

$$\mathbb{E} [\|\hat{\mu}(D) - \mu\|^2] \lesssim \frac{R^2}{mn} + \frac{dR^2}{mn^2\epsilon^2} \ln(mnd) \ln \frac{1}{\delta}.$$

- No sacrifice of utility under this simple case

Theorem 2: Heavy-tailed distributions, balanced users

If the distribution has p -th bounded moment ($p \geq 2$), for n users with m items per user,

$$\mathbb{E} [\|\hat{\mu}(D) - \mu\|^2] \lesssim \frac{1}{mn} + \left[\frac{d \ln(nd)}{mn^2\epsilon^2} + \left(\frac{d}{m^2n^2\epsilon^2} \right)^{1-\frac{1}{p}} \ln^2(nd) \right] \ln \frac{1}{\delta}.$$

- Significant improvement over existing solution [1]
- With $m = 1$, the result matches the state-of-the-art item-level DP estimators

Bounded support, imbalanced users

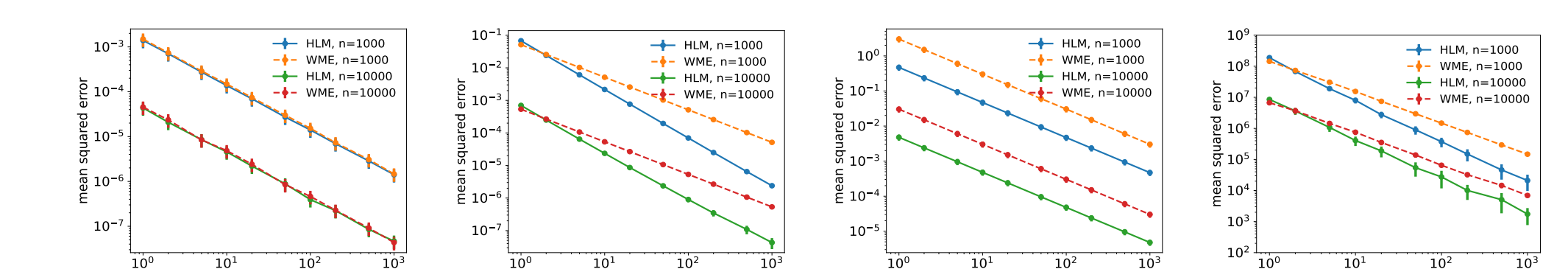
Let γ be the degree of imbalance (definition omitted here, for balanced users $\gamma = 1$; large γ indicates strong imbalance), for N total items distributed in n users,

$$\mathbb{E} [\|\hat{\mu}(\mathcal{D}) - \mu\|^2] \lesssim \frac{R^2}{N} + \frac{dR^2\gamma}{Nn\epsilon^2} \ln^2(Nnd) \ln \frac{1}{\delta}.$$

- With proper h , both ℓ_2 and ℓ_∞ bounds are nearly optimal (up to log factor)

Evaluation

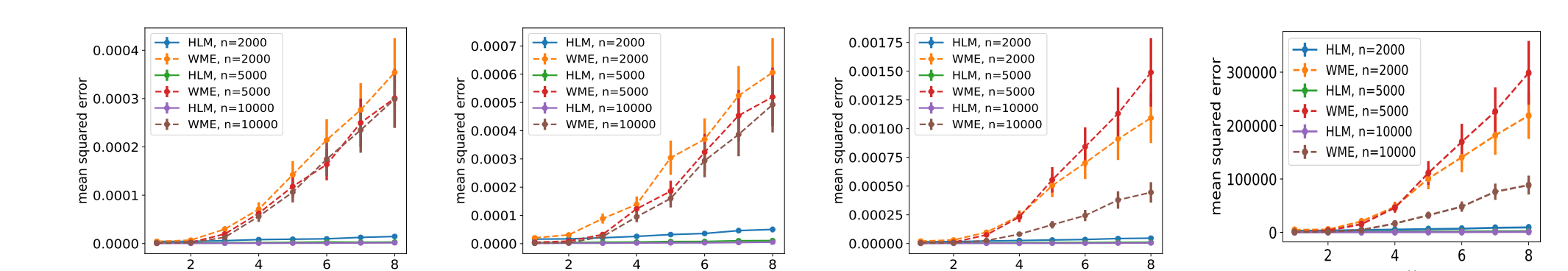
Balanced users:



(a) Uniform, $d = 1$ (b) Lomax, $d = 1$ (c) Uniform, $d = 3$ (d) IPUMS total income

Figure: Convergence of mean squared error with balanced users.

Imbalanced users:



(a) Uniform distribution. (b) Gaussian distribution. (c) Exp. distribution. (d) IPUMS total income.

Figure: Growth of mean squared error with degree of imbalance γ .

Reference

[1] Levy, Daniel, et al. "Learning with user-level privacy." NeurIPS 2021