

ESTIMATING THE TRUE DISTRIBUTION OF DATA COLLECTED WITH RANDOMIZED RESPONSE



Code and paper

Inria
carlos.pinzon@inria.fr

Carlos Pinzón¹, Ehab ElSalamouny^{1,5}, Lucas Massot³,
Alexis Miller⁴, Héber H. Arcolezi², Catuscia Palamidessi¹

¹INRIA Saclay, France; ²INRIA Grenoble, France; ³École Polytechnique, France;
⁴Ecole Normale Supérieure de Lyon, France; ⁵Suez Canal University, Egypt

Introduction

Local Differential Privacy (LDP) provides privacy guarantees for data collection without requiring trust in the data collector.

Randomized Response (RR) is a fundamental LDP mechanism:

- Each user reports their true value with probability p or a random value otherwise
- Satisfies ϵ -LDP with $\epsilon = \log(p/q)$ where $q = (1-p)/(K-1)$
- Used by major tech companies for telemetry data collection
- Optimal for small domains: $K < 3e^\epsilon + 2$

Problem: The standard debiasing rule can produce *invalid* histograms with negative values. What is the best fix?

Problem Formulation

- N users, each with secret value $x_u \in \{1, \dots, K\}$
- True distribution (unknown): $\theta_i = \frac{| \{u: x_u=i\} |}{N}$. $\theta \in \Delta := \{\theta \in \mathbb{R}^K : \sum_i \theta_i = 1, \theta_i \geq 0\}$
- Each user applies RR mechanism: $y_u = \mathcal{M}(x_u)$
- Observed histogram: $\phi_i = \frac{| \{u: y_u=i\} |}{N}$. $\phi \in \Delta$
- **Goal:** Estimate θ from ϕ

RR Mechanism:

$$\Pr(\mathcal{M}(x) = y) = \begin{cases} p & \text{if } y = x \\ q & \text{otherwise} \end{cases} \quad p = \frac{e^\epsilon}{e^\epsilon + K - 1}, \quad q = \frac{1}{e^\epsilon + K - 1}$$

Existing Estimators

Method Complexity Properties

Method	Complexity	Properties
Inv	$O(K)$	Unbiased but invalid
InvN	$O(K)$	Valid, simple workaround
InvP	$O(K \log K)$	Valid, closest to Inv
IBU	$O(K N_{\text{iters}})$	Valid, MLE as $N_{\text{iters}} \rightarrow \infty$
MLE*	$O(K \log K)$	Valid, exact MLE

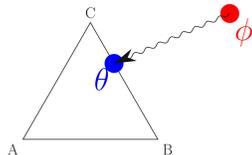
Linear Inversion (Inv) — debiasing rule:

$$\hat{\theta}_i^{\text{Inv}} = \frac{\phi_i - q}{p - q}$$

- Unbiased: $E[\hat{\theta}^{\text{Inv}}] = \theta$
- Can have $\hat{\theta}_i < 0$ (invalid!)
- Probability of invalidity is non-negligible for any N

Workarounds:

- **InvN:** Set negatives to zero and normalize
- **InvP:** Project to simplex $\hat{\theta}^{\text{InvP}} = \arg \min_{\theta \in \Delta} \|\theta - \hat{\theta}^{\text{Inv}}\|$



Iterative Bayesian Update (IBU):

- Converges to MLE as iterations $\rightarrow \infty$
- Requires many iterations in practice
- No theory on stopping conditions

Our Contribution: Closed Formula for MLE

Key Idea: Apply threshold transformation to ϕ , then use Inv

For threshold τ , define:

$$\phi_i^\tau = \begin{cases} q & \text{if } \phi_i < \tau \\ c_\tau \phi_i & \text{otherwise} \end{cases} \quad \text{where } c_\tau = \frac{1 - \sum_{\phi_i < \tau} q}{\sum_{\phi_i \geq \tau} \phi_i}$$

Our Estimator:

$$\text{MLE}^*(\phi) = \text{Inv}(\phi^{\tau^*}) \\ \tau^* = \min \{ \tau \mid \forall i, \phi_i < \tau \vee c_\tau \phi_i \geq q \} \dots \text{ can be found in } O(K \log K)$$

Theorem 1. The MLE for the RR mechanism is **unique** and given by **MLE***.

Proof Sketch:

1. Using Lagrange multipliers: $\hat{\theta}_i = 0$ or $\hat{\theta}_i = \frac{\phi_i - q}{p - q}$
2. MLE is monotonic: $\hat{\theta}_i \leq \hat{\theta}_j \Leftrightarrow \phi_i \leq \phi_j$
3. Zeros occur in smallest components of ϕ
4. Find optimal number of zeros via threshold τ^*

Theoretical Comparison

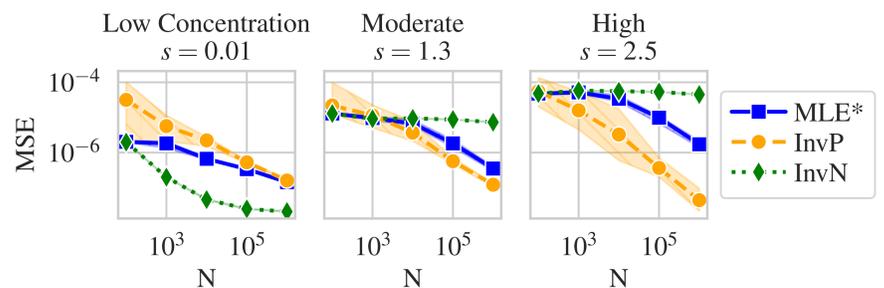
Property	Inv	InvN	InvP	IBU	MLE*
Valid	✗	✓	✓	✓	✓
Unbiased	✓	✗	✗	✗	✗
Is the MLE	✗	✗	✗	Only if $N_{\text{iters}} \rightarrow \infty$	✓
Complexity	$O(K)$	$O(K)$	$O(K \log K)$	$O(K N_{\text{iters}})$	$O(K \log K)$
Consistent	✓	✓	✓	✓	✓

Key Insights:

- All estimators are consistent: $\text{MSE} = O(K/N) \rightarrow 0$ as $N \rightarrow \infty$
- Valid estimators are necessarily biased (cannot be both valid and unbiased)
- **MLE*** achieves exact MLE with complexity $O(K \log K)$ vs. $O(K N_{\text{iters}})$ for IBU

Experimental Results

Setup: Varied $\epsilon \in \{1, \dots, 10\}$, $N \in \{10^2, \dots, 10^6\}$, $K \in \{50, \dots, 5000\}$, Zipf skewness $s \in \{0.01, 1.3, 2.5\}$. For $\epsilon = 4.0$, $K = 10000$:



Key Finding: MLE* is a **robust default estimator**

- InvP and InvN alternate as best depending on distribution
- MLE* is never the worst estimator
- Always achieves lowest negative log-likelihood (by construction)
- Safe choice when true distribution is unknown

Real-World Data: Tested on Kosarak (clickstream) and ACSIncome (census) datasets with K up to 1.4M. MLE* remains consistently robust (non-worst) across all distributions.

Conclusion

- **Exact and mathematically proven formula** for the MLE of Randomized Response
- **Efficient algorithm:** $O(K \log K)$ complexity, significantly faster than iterative IBU
- **Robust performance:** Consistently close to best estimator across diverse data distributions
- **Practical impact:** Reliable default choice for practitioners deploying LDP systems

Future Work: Extension to other mechanisms: (1) arbitrary channel matrix and (2) longitudinal protocols (RAPPOR, Local Hashing)