

Longitudinal data collection based on local hashing for local differential privacy^[1]

Carlos Pinzón

8th Franco-Japanese Cybersecurity Workshop, WG1, formal methods

ENS de Cognitique, Bordeaux.

November 29, 2023.



Outline

1. Problem: collect longitudinal data preserving privacy
2. State of the art
3. Longitudinal local hashing

Problem: collect longitudinal data preserving privacy

	Single individual	Population
Single time	(one value)	Cross-sectional
Many times	Time series	Longitudinal

Problem

- n users, $\tau \gg 1$ time steps, k possible values: $x_t^{(u)} \in [1..k]$
- We want to collect data to estimate population frequencies:
 $\mathbf{P}(X_t = x)$ for each $x \in [1..k]$ and each $t \in [1..\tau]$
- ...but individual values are private

Local Differential Privacy (LDP)

“Collected values will be partial, approximate or mere garbage”

$$\epsilon\text{-LDP}: \quad \forall(y, x^+, x^-), \quad \mathbf{P}(Y=y|X=x^+) \leq e^\epsilon \mathbf{P}(Y=y|X=x^-)$$

Trust barrier:



Local Differential Privacy (LDP)

“Collected values will be partial, approximate or mere garbage”

$$\epsilon\text{-LDP: } \forall(y, x^+, x^-), \quad \mathbf{P}(Y=y|X=x^+) \leq e^\epsilon \mathbf{P}(Y=y|X=x^-)$$

Trust barrier:



Local Differential Privacy (LDP)

“Collected values will be partial, approximate or mere garbage”

$$\epsilon\text{-LDP}: \quad \forall(y, x^+, x^-), \quad \mathbf{P}(Y=y|X=x^+) \leq e^\epsilon \mathbf{P}(Y=y|X=x^-)$$

Trust barrier:



Mechanisms for LDP

Categorical secret $\in [1..k]$

1. Randomized Response (RR)^[1]:

reported value =
$$\begin{cases} \text{secret} & \text{with some prob.} \\ \text{any other value (uniformly)} & \text{otherwise} \end{cases}$$

2. Unary Encoding^[2]:

apply one-hot encoding, then RR to each bit

3. Local Hashing^[3]:

partition $[1..k]$ into groups, then apply RR to the group labels

Naïve solution

$$x \xrightarrow{\epsilon} x'_1$$

$$x \xrightarrow{\epsilon} x'_2$$

⋮

- The privacy guarantee is $t\epsilon$ -LDP after t reports $[x'_1, \dots, x'_t]$.
- As $t \rightarrow \infty$, LDP breaks and we guess x .

Improvement: double randomization

$$x \xrightarrow{\epsilon_0} x' \xrightarrow{\epsilon} x''_1$$

$$x' \xrightarrow{\epsilon} x''_2$$

$$x' \xrightarrow{\epsilon} x''_3$$

⋮

- As $t \rightarrow \infty$, we guess x' but not x (the user's value).
- ϵ_0 -LDP is guaranteed.

Improvement: double randomization

$$x \xrightarrow{\epsilon_0} x' \rightsquigarrow x''_1$$

$$\textcolor{green}{x} \xrightarrow{\text{memo}} x' \rightsquigarrow x''_2$$

$$\textcolor{green}{x} \xrightarrow{\text{memo}} x' \rightsquigarrow x''_3$$

⋮

- As $t \rightarrow \infty$, we guess x' but not x (the user's value).
- ϵ_0 -LDP is guaranteed.

Improvement: double randomization

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix} \xrightarrow[\epsilon_0]{\text{memo}} \begin{bmatrix} x'_1 \\ x'_2 \\ \vdots \end{bmatrix} \xrightarrow{\epsilon} \begin{bmatrix} x''_1 \\ x''_2 \\ \vdots \end{bmatrix}$$

- Is ϵ_0 -LDP guaranteed?

Example of double randomization

$$x_1 = A \xrightarrow{\epsilon_0} A' \xrightarrow{\epsilon} x_1''$$

$$x_2 = C \xrightarrow{\epsilon_0} C' \xrightarrow{\epsilon} x_2''$$

$$x_3 = C \xrightarrow{\text{memo}} C' \xrightarrow{\epsilon} x_3''$$

$$x_5 = A \xrightarrow{\text{memo}} A' \xrightarrow{\epsilon} x_5''$$

$$x_9 = A \xrightarrow{\text{memo}} A' \xrightarrow{\epsilon} x_9''$$

$$x_4 = G \xrightarrow{\epsilon_0} G' \xrightarrow{\epsilon} x_4''$$

$$x_6 = T \xrightarrow{\epsilon_0} T' \xrightarrow{\epsilon} x_6''$$

$$x_7 = T \xrightarrow{\text{memo}} T' \xrightarrow{\epsilon} x_7''$$

$\xrightarrow{\text{memo}}$:

$\xrightarrow{\text{memo}}$:

$$\forall x \in [1..k] = \{A, G, T, C\}$$

- the fixed value x' is exposed to multiple queries.
- x is exposed only once.

But there are leakages about data changes and time patterns

Outline

1. Problem: collect longitudinal data preserving privacy
2. **State of the art**
3. Longitudinal local hashing

Real world deployments

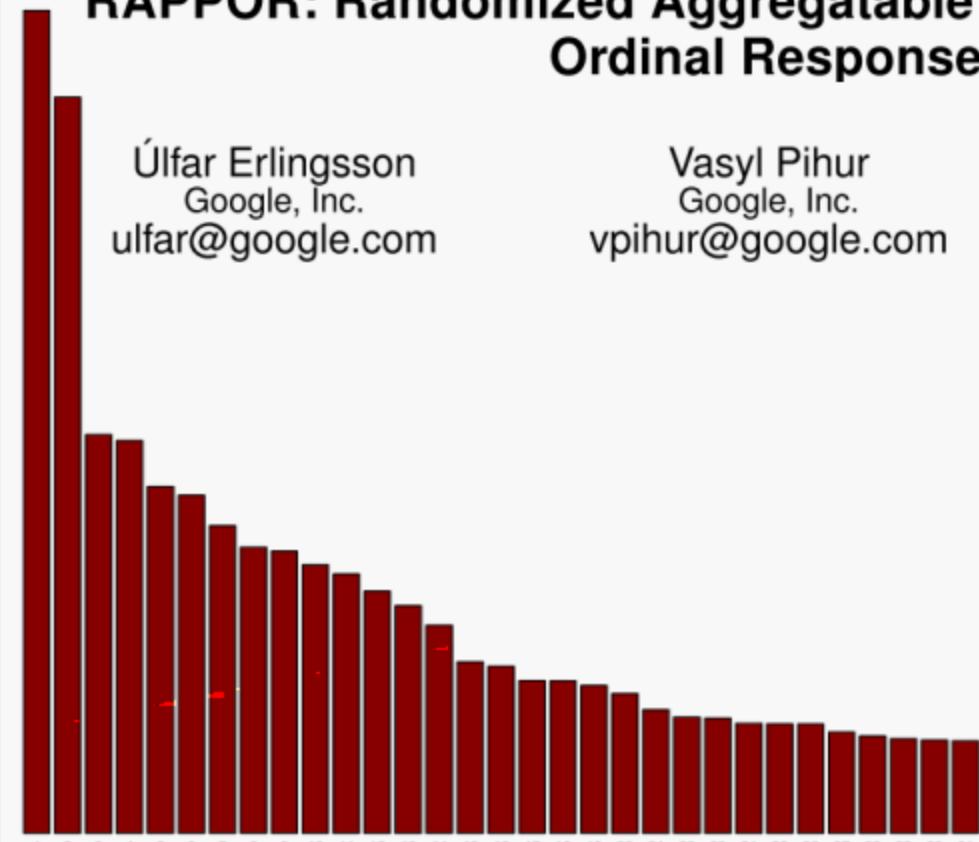


RAPPOR: Randomized Aggregatable Privacy-Preserving Ordinal Response

Úlfar Erlingsson
Google, Inc.
ulfar@google.com

Vasyl Pihur
Google, Inc.
vpihur@google.com

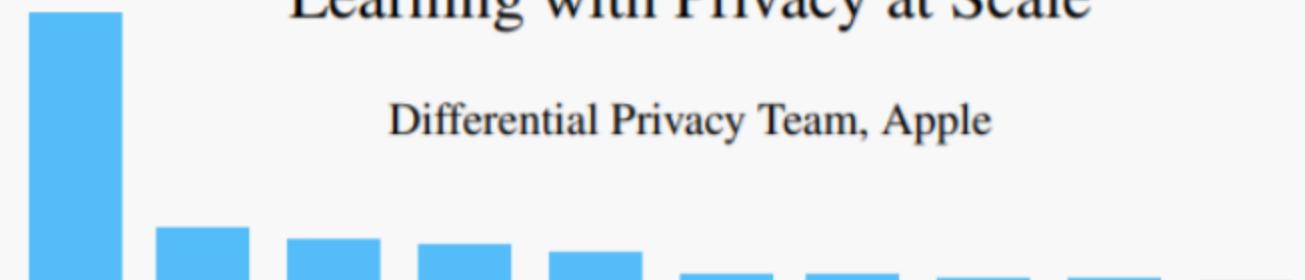
Aleksandra Korolova
University of Southern California
korolova@usc.edu



Rank	Domain	Relative Frequency (approx.)
1	google.com	100
2	www.google.com	~80
3	www.youtube.com	~60
4	www.facebook.com	~50
5	www.gmail.com	~45
6	www.g.doubleclick.net	~40
7	www.abcnews.go.com	~35
8	www.cnn.com	~30
9	www.huffingtonpost.com	~28
10	www.msnbc.msn.com	~25
11	www.espn.com	~22
12	www.nytimes.com	~20
13	www.hbo.com	~18
14	www.pewinternet.org	~15
15	www.sfgate.com	~12
16	www.usatoday.com	~10
17	www.foxnews.com	~8
18	www.huffingtonpost.ca	~6
19	www.abc.com	~5
20	www.espn.go.com	~4
21	www.cbsnews.com	~3
22	www.foxnews.ca	~2
23	www.fox59.com	~1.5
24	www.fox2now.com	~1.2
25	www.fox32chicago.com	~1
26	www.fox2now.com	~0.8
27	www.fox2now.com	~0.6
28	www.fox2now.com	~0.5
29	www.fox2now.com	~0.4
30	www.fox2now.com	~0.3
31	www.fox2now.com	~0.2

Learning with Privacy at Scale

Differential Privacy Team, Apple



The Count Mean Sketch technique allows Apple to determine the most popular emoji to help design better ways to find and use our favorite emoji. The top emoji for US English speakers contained some surprising favorites.

emojis: 😂 ❤️ 😢 😍 😜 😏 😔 💀 😊 😢 😐

Collecting Telemetry Data Privately

Bolin Ding, Janardhan Kulkarni, Sergey Yekhanin
Microsoft Research
{bolind, jakul, yekhanin}@microsoft.com

Windows Insiders in Windows 10 Fall Creators Update to protect users' privacy while collecting application usage statistics.

State of the art protocols

	Dom.	Preprocessing	Perm. rand.	Inst. rand.	Output dom.
L-RR	[1..k]		[1..k]	RR	RR
RAPPOR ^[1]	[1..k]	One hot	$\{0, 1\}^k$	SUE	SUE
L-OSUE	[1..k]	One hot	$\{0, 1\}^k$	OUE	SUE
dBitFlipPM ^[2] , $d \leq k$	[1..k]	Ad-hoc	$\{0, 1\}^d$	SUE	$\{0, 1\}^d$
LOLOHA	[1..k]	Local hash	[1..g]	RR	RR

Outline

1. Problem: collect longitudinal data preserving privacy
2. State of the art
3. **Longitudinal local hashing**

LOLOHA: Longitudinal Local Hashing

The best of both worlds:

Protocol	Privacy strength	Parameters
RAPPOR ^[1]	double sanitization	e^{ϵ_∞} and $e^{\alpha\epsilon_\infty}$
d -BitFlipPM ^[2]	reduced domain	e^{ϵ_∞} and $d \ll k$
LOLOHA	(both)	e^{ϵ_∞} and $e^{\alpha\epsilon_\infty}$ and $g \ll k$

Parameter g :

1. Tuned for privacy $g = 2$
2. Tuned for utility $g = 1 + \max \left(1, \frac{1 - e^{2\epsilon_\infty} + \sqrt{e^{4\epsilon_\infty} - 14e^{2\epsilon_\infty} + 12e^{(1+\alpha)\epsilon_\infty}(1 - e^{(1+\alpha)\epsilon_\infty}) + 12e^{(3+\alpha)\epsilon_\infty} + 1}}{6e^{\epsilon_\infty} - e^{\alpha\epsilon_\infty}} \right)$

LOLOHA: LOngitudinal LOcal HAshing

Setup:

1. Fix integer g
2. Each user u reports a random hash function $H_u : [1..k] \rightarrow [1..g]$

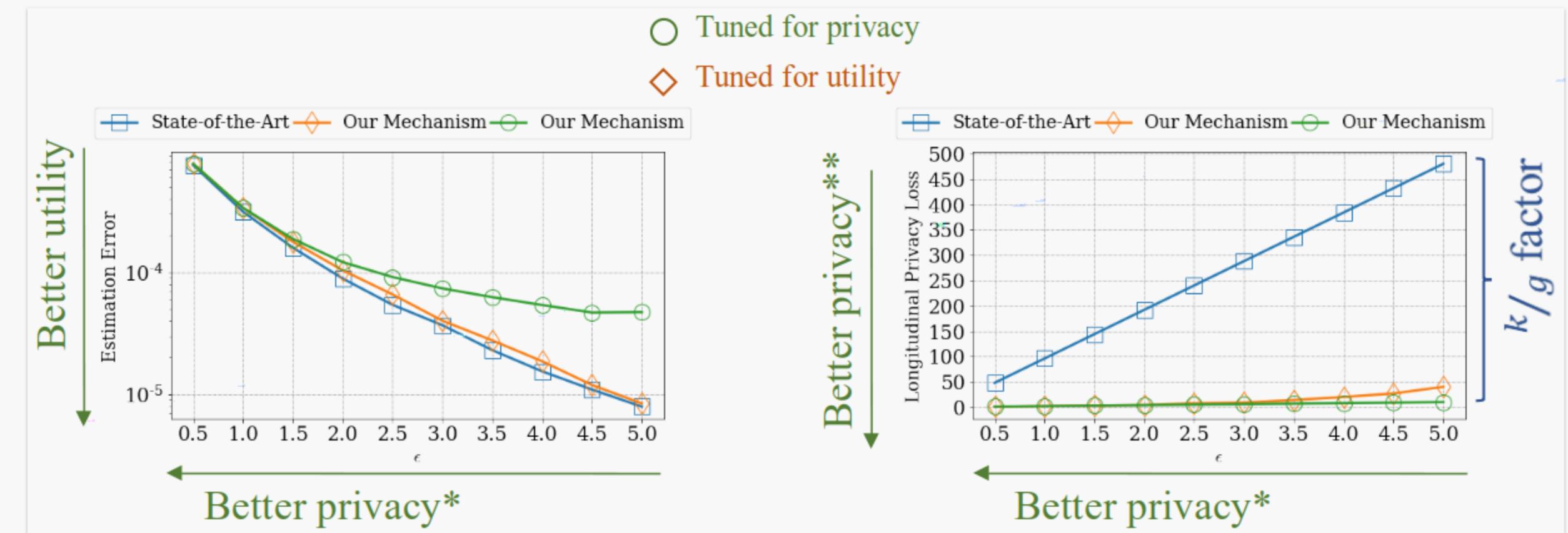
Execution (each user u at each time step t):

1. Hash the secret $x_t^{(u)}$ into $[1..g]$
2. Apply permanent g -RR and instantaneous g -RR

$$x \quad \mapsto \quad \text{RR} \left(\text{RR}_{\text{memo}} \left(H_u(x_t^{(u)}) \bmod g \right) \right)$$

Results

Adult dataset (hours_per_week attribute, $n = 45422$, $k = 96$, $\tau = 260$)



$\epsilon_\infty \in (0, \infty)$ *LDP privacy loss assuming constant input.

MSE $\in [0, 1]$ $\frac{1}{\tau} \sum_{t=1}^{\tau} \frac{1}{k} \sum_{x=1}^k (f(x) - \hat{f}(x))^2$

$\check{\epsilon}_{\text{avg}} \in (0, k\epsilon_\infty)$ **“longitudinal loss” permanent randomizations.

Conclusion

LOLOHA: a protocol for collecting evolving categorical data

Remarks

- LOLOHA combines double randomization with reduced domain
- Double randomization is decent but not LDP

Contributions

- Similar performance with less “*LDP on the users' values*”
- First performant method not based on Unary Encoding