Privacy in Web Advertising: Analytics & Modeling

Badih Ghazi & Pasin Manurangsi Google Research

Joint work with colleagues at Google

Ghazi & Manurangsi

Privacy in Web Advertising: Analytics and Modeling

We will discuss tools for private analytics and learning for web advertising applications.

Outline

Part I: Basic & Analytics

- Privacy Attacks
- Differential Privacy: Basics
- DP Properties
- Noise Addition Mechanisms
- Threat Models
 - Central vs Local DP
 - Shuffle DP, SMPC and TEEs
- DP Ad Analytics
 - Conversion Measurement
 - Reach and Frequency

Part II: Learning & Other Topics

- DP Ad Modeling
 - Full DP
 - Label DP
 - DP with partially known features
- Future Directions

Online Advertising: Terminology

Entities

Advertiser

An entity / website / app that pays to have their ads shown to the users

Publisher

The website that shows the ads to the users

Will not focus on ad-tech in this tutorial

Ad-Tech

Entity that helps advertiser / publisher buying / selling / optimizing ads **Events**

Impression

Event indicating that an ad is viewed or clicked by a user

Conversion

Action by the users that is valuable to advertisers

Online Advertising: Example



Privacy in Web Advertising: Analytics and Modeling

Privacy Attacks

Privacy in Web Advertising: Analytics and Modeling

Ghazi & Manurangsi

Basic Setting



Reconstruction Attacks



Attacks Against Query-Answering System

Privacy in Web Advertising: Analytics and Modeling

Ghazi & Manurangsi

Query-Answering System



Reconstruction Attacks

name	zipcode	age	#convs
Dave	10520	40	150
Bob	10520	35	50
Carol	10500	30	30
Alice	10500	41	20

- Simplifying assumptions: only allow **SUM** (or **COUNT**) queries
- Only one sensitive column

Reconstruction Attacks



- Simplifying assumptions: only allow **SUM** (or **COUNT**) queries
- Only one sensitive column





The adversary can construct the entire dataset!

Mitigation I: "k-Anonymity"



k-Anonymity

Algorithm rejects if query involves < k individuals

Mitigation I: "k-Anonymity"



k-Anonymity

Algorithm rejects if query involves < k individuals

Differencing Attacks



Ghazi & Manurangsi

Mitigation II: Add Noise



Noise Addition

Add random noise to the answer

Mitigation II: Add Noise



Add random noise to the answer

An approach taken by differential privacy







 $x_1 + x_4 = 170$



 $x_1 + x_4 = 170$ $x_2 + x_3 = 80$







Attacks Against Anonymized Dataset

Privacy in Web Advertising: Analytics and Modeling

Ghazi & Manurangsi

Raw dataset



Raw dataset



Common techniques

• **Suppression:** erasing the entries

Raw dataset



Common techniques

• **Suppression:** erasing the entries

Raw dataset



Common techniques

- **Suppression:** erasing the entries
- Generalization: replacing entries with more general value

Raw dataset

Anonymized dataset



Common techniques

- **Suppression:** erasing the entries
- Generalization: replacing entries with more general value

Linkage Attacks

Adversary's Auxiliary data

Anonymized dataset



Linkage Attack (aka De-Anonymization Attack)

Adversary links record of anonymized dataset to auxiliary dataset

Linkage Attacks



Linkage Attack (aka De-Anonymization Attack)

Adversary links record of anonymized dataset to auxiliary dataset

Ghazi & Manurangsi



Linkage Attack (aka De-Anonymization Attack)

Adversary links record of anonymized dataset to auxiliary dataset

Linkage Attacks: Practical Examples

Example I: Massachusetts Health Data

Medical data of state employees in the state of Massachusetts

• Personal identifiers (eg, name, SSN) are removed

[Sweeney] managed to identify the governor's medical records

Example II: Netflix Prize

"Anonymized" user ratings

- Each row is (user, movie, rating, time)
- User id is re-randomized

99% re-identified by [Narayanan, Shmatikov'08]

Example III: EdX Dataset

5-anonymized student data

- Each row has demographics, course activities/outcomes
- Anonymization using subset of attributes

Using LinkedIn side information, a few students re-identified [Cohen'22]
k-Anonymity

2-Anonymity



Limitations

- Requires adversary to know *nothing* about unknown columns
- k-Anonymity does not hold if consider two releases

k-Anonymity

Every combination of known values appear at least k times



Limitations

- Requires adversary to know *nothing* about unknown columns
- k-Anonymity does not hold if consider two releases

k-Anonymity

Every combination of known values appear at least k times



Limitations

- Requires adversary to know *nothing* about unknown columns
- k-Anonymity does not hold if consider two releases
- Computing "optimal" anonymized dataset is NP-hard

k-Anonymity

Every combination of known values appear at least k times

What is a privacy violation?



Inference ≠ Privacy Violation

Public			Sensitive
name	Heavy smoker	age	Has Cancer
Dave	NO	40	-
Bob	YES	60	-
Carol	NO	55	-
Alice	NO	41	-

Study in 1950s

"Heavy Smoking Cause Cancer"

Inference ≠ Privacy Violation



Conclusion of the study affects Bob not his presence or absence in the dataset

Privacy in Web Advertising: Analytics and Modeling

Inference ≠ Privacy Violation



- If a user chooses not to be in the dataset, then need to be more careful
- Can model as user is in the dataset, but with a column saying "OPT OUT"

If a user's data is not in the input, ⇒ releasing the output does *not* violate the user's privacy



Differential Privacy

Privacy in Web Advertising: Analytics and Modeling

Why Differential Privacy (DP)?

- DP is only privacy notion that is robust (unlike k-anonymity) and has intuitive, practically appealing properties
 - DP give rigourous guarantees that an individual user will not be adversely affected by allowing their data to be used
 - Learns little about an individual but useful things about a population
- Other notions of security/privacy (eg, FL or MPC) have to be combined with DP to provide anonymity guarantees

- DP is the only *anonymity* notion broadly adopted by multiple organizations
 - Google, <u>Apple</u>, <u>Microsoft</u>, <u>Meta</u>, <u>LinkedIn</u>, <u>Uber</u>, <u>Amazon</u>
 - o <u>US Census Bureau</u>
 - Academia
- Privacy research is dominated by DP
 - Eg, NeurIPS'22 had 50 papers with
 "priv" in title and 49 were on DP

Differential Privacy: Basics

Privacy in Web Advertising: Analytics and Modeling



Differential Privacy: Definition

Intuition:

Adding or removing a single user should not change the output distribution too much



Differential Privacy: Definition

Intuition:

Adding or removing a single user should not change the output distribution too much

Smaller ε, δ \$
More privacy (ε , δ)-Differential Privacy [Dwork et al.'06] For every neighboring datasets X, X' and every set S of outputs, Pr[M(X) \in S] $\leq e^{\varepsilon} \cdot Pr[M(X') \in S] + \delta$ Pure-DP ε-DP **Ξ** (ε, 0)-DP



- "neighboring" notion can be broader than intuition
- e ≍: abbrv for neighboring

 ϵ = small constant

 δ = negligible in #users

User-Level DP

X is X' with an individual's data added / removed

Dataset X

	name	zipcode	age	Purchase value
x ₁	Dave	10520	40	150
x ₂	Bob	10520	35	50
x ₃	Carol	10500	30	30
	Carol	10500	30	70

Dataset X'

	name	zipcode	age	#convs
x ₁	Dave	10520	40	150
x ₂	Bob	10520	35	50

Event-Level DP



Dataset X

	name	zipcode	age	Purchase value
x ₁	Dave	10520	40	150
x ₂	Bob	10520	35	50
x ₃	Carol	10500	30	30
X ₄	Carol	10500	30	70

Dataset X

	name	zipcode	age	Conv Value
x ₁	Dave	10520	40	150
x ₂	Bob	10520	35	50
x ₃	Carol	10500	30	30

Privacy in Web Advertising: Analytics and Modeling

Add/Remove-DP



Dataset X

	name	zipcode	age	#convs
x ₁	Dave	10520	40	150
x ₂	Bob	10520	35	50
x ₃	Carol	10500	30	30

Dataset X'

	name	zipcode	age	#convs
x ₁	Dave	10520	40	150
x ₂	Bob	10520	35	50

Substitution-DP



Dataset X

	name	zipcode	age	#convs
x ₁	Dave	10520	40	150
x ₂	Bob	10520	35	50
х ₃	Carol	10500	30	30

Dataset X'

	name	zipcode	age	#convs
х ₁	Dave	10520	40	150
×2	Bob	10520	35	50
x ₃	Carol	10520	32	150

Add/Remove-DP

 $X \asymp^r X' \longrightarrow X'$ is X' with an individual's data added / removed

Substitution-DP



These notions are mostly interchangeable up to parameters

Unless otherwise stated, assume *substitution-DP*

Pure vs Approximate DP

(ε , δ)-Differential Privacy [Dwork et al.'06] For every neighboring datasets X, X' and every set S of outputs, Pr[M(X) \in S] $\leq e^{\varepsilon} \cdot Pr[M(X') \in$ S] + δ

Pure-DP
$$\delta = 0$$

- Desirable
- Worst privacy-utility trade-off for many problems
- Weaker composability



- Definition allows catastrophic leak
- δ ideally o(1/n)
- Most practical mechanisms behave gracefully

DP Properties

Privacy in Web Advertising: Analytics and Modeling

Post-Processing

"Result of DP algorithm can be used in arbitrary manner and it remains DP."

Theorem If M is (ε, δ) -DP and $h(\cdot)$ is any function, then h(M(X)) remains (ε, δ) -DP.



Can we use the output safely in downstream applications?

Composition



Composition



Composition



"Running DP algorithms multiple times remain DP, but with worse parameters"

Basic Composition



"Adaptive" setting: previous outputs are used in subsequent algorithms

Basic Composition Theorem [Dwork et al.] All the outputs combined remain $(\varepsilon_1 + \cdots + \varepsilon_k, \delta_1 + \cdots + \delta_k)$ -DP

Works even for "adaptive" setting

Advanced Composition



Privacy in Web Advertising: Analytics and Modeling

Parallel Composition



Parallel Composition Theorem [McSherry]

All the outputs combined remain (ε, δ) -DP

Example:

Email	zipcode	Purchase value	
dave123@gmail.com	10520	500	
dave1993@gmail.com	10520	25	ε-DP Alg Outpu
alice@gmail.com	10500	15	
bob@gmail.com	10500	40	

k-Neighboring

X, X' are k-neighboring, denoted by \asymp_k iff there

exists a sequence $X = X_0 \rightleftharpoons X_1 \rightleftharpoons \cdots \rightleftharpoons X_t = X', t \le k$.

Theorem If M is (ε, δ) -DP under neighboring notion \asymp , then M(X) is (ε', δ') -DP for neighboring notion $\asymp_{\mathbf{k}}$, where $\varepsilon' = k\varepsilon$ and $\delta' = \delta \cdot (e^{k\varepsilon} - 1) / (e^{\varepsilon} - 1)$.



k-Neighboring

X, X' are k-neighboring, denoted by \asymp_k iff there

exists a sequence $X = X_0 \rightleftharpoons X_1 \rightleftharpoons \cdots \rightleftharpoons X_t = X'$, $t \le k$.

Theorem If M is (ε, δ) -DP under neighboring notion \asymp , then M(X) is (ε', δ') -DP for neighboring notion \asymp_k , where $\varepsilon' = k\varepsilon$ and $\delta' = \delta \cdot (e^{k\varepsilon} - 1) / (e^{\varepsilon} - 1)$.

Example:

Add/remove $X \Rightarrow^r X'$ X is X' with an individual's
data added / removedSubstitution $X \Rightarrow^s X'$ X is X' with an individual's
data changed

Observation If $X \equiv^{s} X'$, then $X \equiv^{r}_{2} X'$

"Every substitution neighbor is a 2-add/remove neighbor"

k-Neighboring

X, X' are k-neighboring, denoted by \asymp_k iff there

exists a sequence $X = X_0 \rightleftharpoons X_1 \rightleftharpoons \cdots \rightleftharpoons X_t = X'$, $t \le k$.

Theorem If M is (ε, δ) -DP under neighboring notion \asymp , then M(X) is (ε', δ') -DP for neighboring notion \asymp_k , where $\varepsilon' = k\varepsilon$ and $\delta' = \delta \cdot (e^{k\varepsilon} - 1) / (e^{\varepsilon} - 1)$.

Example:

Lemma If M is ε -add/remove-DP, then it is 2ε -substitution-DP



Observation If $X \cong^{s} X'$, then $X \cong^{r} X'$

"Every substitution neighbor

is a 2-add/remove neighbor"

Amplification by Subsampling

"Subsampling makes the algorithm more private."



Amplification by Subsampling

"Subsampling makes the algorithm more private."



Amplification by Subsampling

"Subsampling makes the algorithm more private."

 $\begin{array}{l} \mbox{Amplification-by-subsampling Thm} \\ \mbox{The output combined is } (\epsilon', \delta') \mbox{-} DP \mbox{ where } \\ \epsilon' = \ln(1 + p(e^{\epsilon} - 1)), \quad \delta' = p\delta \\ \mbox{with } p = B \slash n \end{array}$



Basic Mechanism: Noise Addition

Privacy in Web Advertising: Analytics and Modeling
Noise Addition Mechanism $X = (x_1, ..., x_n) \rightarrow Analyzer \rightarrow Estimate of g(X)$ Noise Addition Mechanism $X = (x_1, ..., x_n) \rightarrow Analyzer \rightarrow g(X) + Random Noise$ Intuition: Noise should be large enough to hide a user's contribution

What noise distribution should we use?

- Depends on Range(g)
- Depends on how "sensitive" g is

Assumption: Range(g) $\subseteq \mathbb{Z}$

(ie, g is integer-valued)

Sensitivity $\Delta(g) = \max_{X \neq X'} |g(X) - g(X')|$ Larger sensitivity \Rightarrow More noise required Examples: COUNT Query

g = COUNT(*) WHERE #convs > 40

	name	zipcode	age	#convs
x ₁	Dave	10520	40	150
×2	Bob	10520	35	50
×3	Carol	10500	30	30

g(X

Assumption: Range(g) $\subseteq \mathbb{Z}$

(ie, g is integer-valued)

Sensitivity $\Delta(g) = \max_{X = X'} |g(X) - g(X')|$ Larger sensitivity \Rightarrow More noise required Examples: COUNT Query

g = COUNT(*) WHERE #convs > 40

	name	zipcode	age	#convs
x ₁	Dave	10520	40	150
x ₂	Bob	10520	35	50
x ₃	Carol	10500	30	60



 $\Delta(g) = 1$

Assumption: Range(g) $\subseteq \mathbb{Z}$

(ie, g is integer-valued)

Sensitivity $\Delta(g) = \max_{X \succeq X'} |g(X) - g(X)| = 0$ Larger sensitivity \Rightarrow More noise required

Examples: SUM Query

g = SUM(#convs)

	name	zipcode	age	#convs
x ₁	Dave	10520	40	150
x ₂	Bob	10520	35	50
x ₃	Carol	10500	30	60

g(X) = 260

Assumption: Range(g) $\subseteq \mathbb{Z}$

(ie, g is integer-valued)

Sensitivity $\Delta(g) = \max_{X \succeq X'} |g(X) - g(X)| = 0$ Larger sensitivity \Rightarrow More noise required

Assume #convs \leq c

g = SUM(#convs)

	name	zipcode	age	#convs
x ₁	Dave	10520	40	150
x ₂	Bob	10520	35	50
x ₃	Carol	10500	30	190

$$g(X) = 390$$

Privacy in Web Advertising: Analytics and Modeling

 $\Delta(g) \leq c$

Assumption: Range(g) $\subseteq \mathbb{Z}$

(ie, g is integer-valued)

Sensitivity $\Delta(g) = \max_{X = X'} |g(X) - g(X')|$ Larger sensitivity \Rightarrow More noise required **Examples: AVERAGE Query**

g = AVG(#convs)

	name	zipcode	age	#convs
x ₁	Dave	10520	40	150
x ₂	Bob	10520	35	50
x ₃	Carol	10500	30	190

$$g(X) = 130$$

Assume #convs \leq c

$$\Delta(g) \leq c/n$$

Assumption: Range(g) $\subseteq \mathbb{Z}$

(ie, g is integer-valued)

Sensitivity $\Delta(g) = \max_{X = X'} |g(X) - g(X')|$ Larger sensitivity \Rightarrow More noise required



Assumption: Range(g) $\subseteq \mathbb{Z}$

(ie, g is integer-valued)

Sensitivity $\Delta(g) = \max_{X = X'} |g(X) - g(X')|$ Larger sensitivity \Rightarrow More noise required



Example:

 $\Delta(g) = 1$, $\epsilon = \ln(4/3)$, g(X) = 3

Output: 3 + DLap(1/ln(4/3))



Example: $\Delta(g) = 1$, $\epsilon = \ln(4/3)$, g(X) = 33 + DLap(1/ln(4/3))Output: w.p. 3/28 = 0.1071... 3 w.p. 1/7 = 0.1428... output = w.p. 3/28 = 0.1071...

Discrete Laplace Distribution For every integer i, $Pr[i = DLap(b)] \propto e^{-|i|/b}$



Ghazi & Manurangsi

Theorem Assuming Range(g) $\subseteq \mathbb{Z}$, Discrete Laplace Mechanism is ε -DP.

Illustration for $\Delta(g) = 1$:



Theorem Discrete Laplace Mechanism has MSE = $2p / (1 - p)^2$ where $p = exp(-\epsilon/\Delta(g))$

Utility Measures

- Mean Square Error (MSE): E[(output true)²]
 - o Root Mean Square Error (RMSE): √MSE
- Mean Absolute Error (MAE): E[|output true|]
- We will mostly compute MSE since it is easier to deal with
- MAE is always ≤ RMSE

Theorem Discrete Laplace Mechanism has MSE = $2p / (1 - p)^2$ where $p = exp(-\epsilon/\Delta(g))$

Theorem Discrete Laplace Mechanism has MSE = $\frac{2p}{(1-p)^2}$ where $p = \exp(-\epsilon/\Delta(g))$



 $O((\Delta(g)/\varepsilon)^2)$ for $\varepsilon \le 1$

Privacy in Web Advertising: Analytics and Modeling

Assumption: Range(g) $\subseteq \mathbb{Z}^d$

(ie, g is vector-valued with integer entries)

 $\boldsymbol{\ell_{p}}_{p}\text{-Sensitivity} \\ \boldsymbol{\Delta}_{p}^{}(g) = \max_{X = X'} \left\| g(X) - g(X') \right\|_{p}$

 $||v||_{p} = (|v_{1}|^{p} + \dots + |v_{d}|^{p})^{1/p}$

Examples: Histogram

	name	zipcode	age	#convs	
x ₁	Dave	10520	40	150	
x ₂	Bob	10520	35	50	
x _n	Carol	10500	30	30	

Assumption: Range(g) $\subseteq \mathbb{Z}^d$

(ie, g is vector-valued with integer entries)

 $\boldsymbol{\ell}_{p}\text{-Sensitivity} \\ \boldsymbol{\Delta}_{p}^{r}(g) = \max_{X = X'} \left\| g(X) - g(X') \right\|_{p}$

 $\boldsymbol{\ell}_{p}$ -Norm $||v||_{p} = (|v_{1}|^{p} + \dots + |v_{d}|^{p})^{1/p}$

Examples: Histogram

Histogram of #convs

#individuals





Discrete Laplace Mech: Multi-Dimension $X = (x_1, ..., x_n) + Analyzer + g(X) + DLap(\Delta(g)/\varepsilon)$ Assumption: Range(g) $\subseteq \mathbb{Z}^d$ Examples: Histogram Histogram of #convs

(ie, g is vector-valued with integer entries)

 $\boldsymbol{\ell_{p}}^{-} \textbf{Sensitivity} \\ \boldsymbol{\Delta_{p}}^{}(g) = \max_{\boldsymbol{X} = \boldsymbol{X}'} \left\| \boldsymbol{g}(\boldsymbol{X}) - \boldsymbol{g}(\boldsymbol{X}') \right\|_{p}$

 $\boldsymbol{\ell}_{p}$ -Norm $\|v\|_{p} = (|v_{1}|^{p} + \dots + |v_{d}|^{p})^{1/p}$





Privacy in Web Advertising: Analytics and Modeling

Theorem Assuming Range(g) $\subseteq \mathbb{Z}$, Discrete Laplace Mechanism is ε -DP.

Assumption: Range(g) $\subseteq \mathbb{Z}$

(ie, g is integer-valued)

Sensitivity $\Delta(g) = \max_{X=X'} |g(X) - g(X')|$

> Larger sensitivity ⇒ More Noise Required

Discrete Laplace Mechanism $X = (x_1, ..., x_n) \rightarrow Analyzer \rightarrow g(X) + DLap(\Delta(g)/\epsilon)$ Theorem Assuming Range(g) $\subseteq \mathbb{Z}$, Discrete Laplace Mechanism is ϵ -DP.

Assumption: Range(g) $\subseteq \mathbb{Z}$

(ie, g is integer-valued)

Sensitivity

 $\Delta(g) = \max_{X = X'} |g(X) - g(X')|$

Larger sensitivity ⇒ More Noise Required Discrete Laplace Distribution For every integer i, $Pr[i = DLap(b)] \propto e^{-|i|/b}$

• Necessary!

• If
$$g(X) = 0.1$$
, $g(X') = 0.2$

- Noise added is integer
- Does not change the fractional part
- Adversary can tell exactly which dataset it comes from

Assumption: Range(g) $\subseteq \mathbb{R}$

(ie, g is real-valued)

Laplace Distribution For every real number z, $f_{Lap(b)}(z) \propto e^{-|z|/b}$



Laplace Mechanism

$$X = (x_1, ..., x_n) \rightarrow Analyzer \rightarrow g(X) + Lap(\Delta(g)/\epsilon)$$

Theorem Assuming $Range(g) \subseteq \mathbb{R}$, Laplace Mechanism is ε -DP.

Assumption: Range(g) $\subseteq \mathbb{R}$

(ie, g is real-valued)

Laplace Distribution For every real number z, $f_{Lap(b)}(z) \propto e^{-|z|/b}$



Laplace Mechanism

$$X = (x_1, ..., x_n) \rightarrow Analyzer \rightarrow g(X) + Lap(\Delta(g)/\epsilon)$$

Theorem Assuming Range(g) $\subseteq \mathbb{R}$, Laplace Mechanism is ε -DP.

Theorem MSE = $2(\Delta(g)/\epsilon)^2$.

Assumption: Range(g) $\subseteq \mathbb{R}$

(ie, g is real-valued)

Laplace Distribution For every real number z, $f_{Lap(b)}(z) \propto e^{-|z|/b}$



Laplace Mechanism: Multi-Dimension

$$X = (x_1, ..., x_n) \rightarrow \text{Analyzer} \Rightarrow g(X) + \text{Lap}(\Delta_1(g)/\epsilon)^{\otimes d}$$

Use L-sensitivity

Each coordinate is independent

Theorem Assuming Range(g) $\subseteq \mathbb{R}^d$, Laplace Mechanism is ε -DP.

Assumption: Range(g) $\subseteq \mathbb{R}^d$

(ie, g is vector-valued with real entries)

Laplace Distribution For every real number z, $f_{Lap(b)}(z) \propto e^{-|z|/b}$

Privacy in Web Advertising: Analytics and Modeling

Gaussian Mechanism $X = (x_1, ..., x_n) \rightarrow Analyzer \rightarrow g(X) + \mathcal{N}(0, \sigma^2)$

Assumption: Range(g) $\subseteq \mathbb{R}$

(ie, g is real-valued)

Gaussian DistributionFor every real number z, $PDF_{\mathcal{N}(0, \hat{\sigma})}(z) \propto e^{-(z/\hat{\sigma})}$



$$Gaussian Mechanism \sigma = \frac{2\sqrt{2\ln(2/\delta)}}{\epsilon} \cdot \Delta(g)$$
$$X = (x_1, ..., x_n) + Analyzer + g(X) + \mathcal{N}(0, \sigma^2)$$

Theorem Assuming $\text{Range}(f) \subseteq \mathbb{R}$ and $\varepsilon \leq 1$, Gaussian Mechanism is (ε, δ) -DP.

Assumption: Range(g) $\subseteq \mathbb{R}$

(ie, g is real-valued)

Gaussian DistributionFor every real number z, $PDF_{\mathcal{N}(0, \hat{\sigma})}(z) \propto e^{-(z/\hat{\sigma})}$

Probability Density



$$Gaussian Mechanism \sigma = \frac{2\sqrt{2\ln(2/\delta)}}{\epsilon} \cdot \Delta(g)$$
$$X = (x_1, ..., x_n) + Analyzer + g(X) + \mathcal{N}(0, \sigma^2)$$

Theorem Assuming Range(f) $\subseteq \mathbb{R}$ and $\varepsilon \leq 1$, Gaussian Mechanism is (ε, δ) -DP.

Assumption: Range(g) $\subseteq \mathbb{R}$ (ie, g is real-valued)Gaussian Distribution
For every real number z,
PDF_{N(0, d)}(z) $\propto e^{-(z/d)}$

Theorem Assuming Range(f) $\subseteq \mathbb{R}^d$ and $\varepsilon \leq 1$, Gaussian Mechanism is (ε, δ) -DP.





$$\begin{array}{l} Gaussian Mechanism \\ X = (x_1, ..., x_n) \end{array} \bullet \begin{array}{l} Analyzer \end{array} \bullet g(X) + \mathcal{N}(0, \sigma^2)^{\otimes d} \end{array} \cdot \begin{array}{l} \Delta_2(g) \\ \bullet \\ Bach coordinate is independent \end{array}$$

Theorem Assuming Range(f) $\subseteq \mathbb{R}^d$ and $\varepsilon \leq 1$, Gaussian Mechanism is (ε, δ) -DP.

Assumption: Range(g) $\subseteq \mathbb{R}^d$

(ie, g is vector-valued with real entries)

Gaussian Distribution For every real number z, $PDF_{\mathcal{N}(0, \sigma^2)}(z) \propto e^{-(z/\sigma^2)}$ **Examples: Vector summation**

$$\left.\begin{array}{c} \left(0.2, -1, 1\right) \\ \vdots \\ \left(0, 2, -0.1\right) \\ \text{User n} \end{array}\right\} g(X) = (50.1, 2.3, 14.7)$$

$$\left.\begin{array}{c} \left(0, 2, -0.1\right) \\ \text{User n} \end{array}\right\}$$

$$\left.\begin{array}{c} \left(0, 2, -0.1\right) \\ \text{User n} \end{array}\right]$$

$$\left.\begin{array}{c} \left(0, 2, -0.1\right) \\ \text{User n} \end{array}\right]$$

Threat Models

Privacy in Web Advertising: Analytics and Modeling

Ghazi & Manurangsi

Central vs Local DP

Privacy in Web Advertising: Analytics and Modeling

Ghazi & Manurangsi

Model studied so far

$$X = (x_1, ..., x_n) \rightarrow Analyzer \rightarrow Output$$

"Central DP"

- Analyzer gets to see raw data
- Undesirable if:
 - There is no trusted central authority
 - The analysis is done in a distributed manner

Distributed Analytics



Privacy in Web Advertising: Analytics and Modeling

Central DP Model



Local DP Model [Kasiviswanathan et al.]


Distributed Analytics: Counting



Privacy in Web Advertising: Analytics and Modeling

Counting in Central DP: Laplace Mechanism



Ghazi & Manurangsi

Counting in Local DP: Laplace Mechanism

 $z_1, ..., z_n \sim Lap(1/\epsilon)$



Local Laplace Mechanism: Utility

Theorem Estimator from Local Laplace Mechanism has MSE $\frac{2n}{\epsilon^2}$

 $\mathsf{RMSE} = \sqrt{(2n)/\varepsilon}$

Proof

$$= \mathbf{E}[(\sum_{i \in [n]} y_i - \sum_{i \in [n]} x_i)^2]$$

$$= \mathbf{E}[(\sum_{i \in [n]} z_i)^2]$$

$$= \sum_{i \in [n]} \mathbf{E}[z_i^2]$$

$$= \sum_{i \in [n]} \operatorname{Var}(z_i)$$

$$= n (2/\epsilon^2)$$

QED

Randomized Response (RR) [Warner]



Randomized Response (RR): Estimator



Privacy in Web Advertising: Analytics and Modeling

Randomized Response (RR): Estimator



Randomized Response (RR): Estimator



Ghazi & Manurangsi

Randomized Response: Utility

Theorem Estimator from RR has MSE $n e^{\epsilon} / (e^{2\epsilon} - 1)^2$ RMSE = $O(\sqrt{n/\epsilon})$



•
$$\mathbf{E}[\hat{\mathbf{x}}_i] = \mathbf{x}_i$$

•
$$\operatorname{Var}[\hat{x}_{i}] = e^{\varepsilon} / (e^{2\varepsilon} - 1)^{2}$$

Thus,

MSE =
$$\mathbf{E}[(\sum_{i \in [n]} \hat{x}_i - \sum_{i \in [n]} x_i)^2]$$

= $\sum_{i \in [n]} \mathbf{E}[(\hat{x}_i - x_i)^2]$
= $\sum_{i \in [n]} \operatorname{Var}(\hat{x}_i)$
= $n e^{\epsilon} / (e^{2\epsilon} - 1)^2$

QED

Recall: Histogram Problem



RAPPOR Algorithm [Erlingsson et al.'14]

<mark>(ε/2)</mark>-local DP

randomizer

R

Binary summation

$$\mathbf{x} = 2$$
 $\mathbf{B} = 3$



*B-length bit string of zeros, except one at x

Theorem RAPPOR is ε-DP.

Utility: Exactly the same as RR with $\epsilon/2$ -DP

Theorem MSE of each entry of histogram is $n e^{\epsilon/2} / (e^{\epsilon} - 1)^2$.

Shuffle DP, SMPC, TEEs

Privacy in Web Advertising: Analytics and Modeling



Recall: Central DP Model



Recall: Local DP Model



Privacy in Web Advertising: Analytics and Modeling

Shuffle DP Model [Bittau et al., Erlingsson et al.]



Shuffle DP Model: Multi-Message Setting



Privacy in Web Advertising: Analytics and Modeling

Ghazi & Manurangsi

Summation in Shuffle DP Model

Single-Message

Randomizer:

$$x \rightarrow User \rightarrow y = x + Noise$$

Analyzer:

Just sum all messages up!

Benefit over Local DP: Amplification by Shuffling Multi-Message:

Split-and-Mix Protocol [Balle et al.'19]

Randomizer:

Х

User
$$z = x + Noise$$

 y_1
 y_{m-1}
 y_m
 y_1 , ..., y_m randomly selected

so that $y_1 + \dots + y_m = z$

As good accuracy as central DP!

. V

Summation in Shuffle DP Model

Experiment (IPUMS 1940 City Dataset)



Parameters: $n \approx 60M$, B = 915, $\Box = 2 * 10^{-9}$

Multi-Message:

Split-and-Mix Protocol [Balle et al.'19]

Randomizer: $x \rightarrow User \rightarrow z = x + Noise \rightarrow y_1$ $\vdots \qquad y_{m-1} \qquad y_m$ y_1, \dots, y_m randomly selected so that $y_1 + \dots + y_m = z$

As good accuracy as central DP!

Secure Multi-Party Computation (SMPC)



Trusted Execution Environments (TEEs)



Higher performance Remote attestation of code and data Offered by Cloud providers

DP Ad Analytics

Privacy in Web Advertising: Analytics and Modeling

Ghazi & Manurangsi

Reach & Frequency

Reach

• # unique users who view the ads

k-Frequency

 Fraction of users that view the ads exactly k times among those who see ads at least once

(≥ k)-Frequency

• Fraction of users that view the ads *at least* k times among those who see ads at least once

User 1: See ads once, User 2: Doesn't see ads

User 3: See ads twice, User 4: See ads twice



Ghazi & Manurangsi

Privacy in Web Advertising: Analytics and Modeling

Slicing Queries

Example queries

- Reach for entire ads campaign
- Reach for people in MA
- Reach for people in CA
- ...

•

...

- Reach for people in MA aged 20-30
- Reach for people in MA aged 30-40

Hierarchical Queries



•••

Online Advertising: Example



Privacy in Web Advertising: Analytics and Modeling

"Allocate the credit to the different impressions leading to a conversion"



Which impressions are more important for this

conversion?

Privacy in Web Advertising: Analytics and Modeling

Ghazi & Manurangsi

"Allocate the credit to the different impressions leading to a conversion"

Last Touch Attribution (LTA): Assigns all credits to most recent impression



"Allocate the credit to the different impressions leading to a conversion"

First Touch Attribution (LTA): Assigns all credits to first impression



"Allocate the credit to the different impressions leading to a conversion"

Linear (aka Uniform): Assigns credits equally among all impressions



"Allocate the credit to the different impressions leading to a conversion"



"Allocate the credit to the different impressions leading to a conversion"

Exponential Time Decay: Assigns credits that *decreases exponentially* for older impressions



"Allocate the credit to the different impressions leading to a conversion"

Prioritization: Each impression has priority; assign all credits to highest-priority impression



"Allocate the credit to the different impressions leading to a conversion"

Data-driven attribution (DDA): Train an *ML model* and use the model to assign credits



Privacy in Web Advertising: Analytics and Modeling

For simplicity, focus on *single-touch* attribution, i.e. when one impression gets all the credits

"Allocate the credit to the different impressions leading to a conversion"

Last Touch Attribution (LTA): Assigns all credits to most recent impression



"Allocate the credit to the different impressions leading to a conversion"

Multiple publishers can be involved

Last Touch Attribution (LIA): Assigns a credits to most recent impression



Ads Attribution Queries

Aggregated Statistics

- Conversion rate
 - % views or clicks that leads to conversion
- Average conversion value
- Can be sliced / hierarchical similar to reach / frequency queries

Event-Level Statistics

- Each impression:
 - Whether it leads to conversion
 - Total conversion value

Example queries

- Conversion value for entire ad campaign
- Conversion value for people in MA
- Conversion value for people in CA
- ...
- Conversion value for people in MA aged 20-30
- Conversion value for people in MA aged 30-40
- •

...

Ads Attribution Queries

Aggregated Statistics

- Conversion rate
 - % views or clicks that leads to conversion
- Average conversion value
- Can be sliced / hierarchical similar to reach / frequency queries

Event-Level Statistics

- Each impression:
 - Whether it leads to conversion
 - Total conversion value

Used to:

- Business decision
 - Increase / decrease ads budget
- Build models
 - Predict conversion rate
 - Predict conversion value
 - → Used to determine bidding price in automatic ads auctions
Privacy & Measurements

Concerns around leakage of information about individual users through

- Released statistics
- Trained **ML models**
- Intermediate tools, e.g. 3PC, used for above purposes

Privacy & Measurements

Concerns around leakage of information about individual users through

- Released **statistics**
- Trained **ML models**
- Intermediate tools, e.g. 3PC, used for above purposes

Ad measurement goals are *aligned* with differential privacy

- Interested in patterns that are large enough to impact business decisions, not individual users
- Prefer to be robust to outliers

Anti-Tracking Efforts



+ other measures

APIs that are more "privacy-preserving" for different use cases



+ other use cases

For list of other proposals, see <u>W3C github repo</u>.

Third-Party Cookies (3PC) & Measurement

 $3PC \Rightarrow$ helps determine which impressions / conversions are from the same user

- \Rightarrow Reach / frequency measurements
- \Rightarrow Ads Attribution measurements

Privacy Concern

- Allows bad actors to track users across sites
 - I.e. know that the user visits publisher1 then \cap publisher2 then ... etc.



High-Level Overview of APIs: Aggregate

Device 1

Impr 1 Impr 2 Conv 1 Impr 3 Conv 2 LTA
Attributed Dataset: (Impr2, Conv1), (Impr3, Conv2)

Device 2

:

Summary

• Attributions happen on device / browser

High-Level Overview of APIs: Aggregate

Device 1



Summary

- Attributions happen on device / browser
- Turned into histogram contributions with capping

High-Level Overview of APIs: Aggregate

Device 1



Example Usage of Aggregate APIs



Example queries

- Conversion value for entire ad campaign
- Conversion value for people in MA
- Conversion value for people in CA
- ...
- Conversion value for people in MA aged 20-30
- Conversion value for people in MA aged 30-40
- ...

One bucket for each query

ARA Summary Reports

Device 1



Optimizing Utility of ARA Summary Reports

Budgeting via Rescaling

- Multiply bucket **i** contribution by a factor of **a**_i
- Rescale back after aggregation
- ⇔ assigning "privacy budget" **a**, to query **i**

Histogram Contribution:



Post-processing

- Given the noisy measurements, can post-process to reduce error further [Dawson et al.'23]
- Exploit linear constraints-parent = sum of children
 - E.g. for top query, can take linear comb.
 between itself and sum of its children



Optimizing Utility of ARA Summary Reports

Budgeting via Rescaling

- Multiply bucket **i** contribution by a factor of **a**_i
- Rescale back after aggregation
- ⇔ assigning "privacy budget" **a**, to query **i**





Post-processing

- Given the noisy measurements, can post-process to reduce error further [Dawson et al.'23]
- Exploit linear constraints-parent = sum of children
 - E.g. for top query, can take linear comb.
 between itself and sum of its children
- Extends the methods of [Hay et al., VLDB'10, Cormode et al., ICDE'12] for regular trees
- Is a special case of the matrix mechanism [Li et al., VLDB '15, Nikolov et al., STOC'13]

Privacy in Web Advertising: Analytics and Modeling

Ghazi & Manurangsi

Aggregate APIs: Challenges

- Objective
 - Aligning objective with downstream tasks
 - E.g. relative error
 - Study in DP literature has focused on additive errors
- Optimization
 - Require some historical data to perform parameter optimization

- Time Delay
 - May have to wait for a long time before getting report
- Extending APIs to support

more complex algorithms:

- Exponential mechanism
 - DP Synthetic data generation
- Sparse Vector Technique
- Propose-test-release

Cross-Media Measurement

Setting

- **p** publishers
- Publisher **i** has a multiset **S**_i of users it reaches

Goal

- **Reach**: number of unique users in the union of $S_1, ..., S_p$
- **k-Frequency**: the fraction of users that are reached k times across all publishers.



Project Halo

Background & Goals

- Led by World Federation of Advertisers (WFA)
- Allow advertisers to measure cross-platform ads campaign performances
 - Respecting user's privacy
- Scope is quite large
 - Reach / Frequency measurements
 - Deduplication of users
 - Cross-publisher data metrics, APIs
 - Planning ...
- Open-source (see <u>here</u>)
- More details <u>here</u>



Overview of Solution

Privacy Guarantee

View of all-but-one helpers (& output) is computationally user-level DP.

- Uses sketching to reduce communication complexity
- Helpers merge the sketch, add noise and estimate reach / frequency



Full protocol published and presented at PoPETs 2022

Private Ads Modeling

Privacy in Web Advertising: Analytics and Modeling

Ghazi & Manurangsi

Ad Modeling for Bidding

Bidding

- **pCVR:** predicting if an impression will result in an attributed conversion
- **pConvs:** predicting # of conversions attributed to impression
- **pValue:** predicting \$ conversion value attributed to impression
- **pCTR:** predicting if an ad will be clicked
- Prediction signals used as inputs to bidding in ad auctions
- Models are huge; e.g., billions of parameters
- Data is sparse and class-imbalanced

One Commonly Used Metric:

- AUC (Area under the Curve):
 - Plot *True Positive Rate* against *False Positive Rate*, & compute area under it.
 - AUC = probability that the classifier ranks a randomly chosen positive example higher than a randomly chosen negative one.
 - \circ Random guessing: AUC = 0.5
 - Perfect prediction utility: AUC = 1
- AUC-Loss = 1 AUC
- Relative Change in AUC-Loss (in %) = (AUC-Loss AUC-Loss_{baseline}) / AUC-Loss_{baseline} * 100
 - baseline is without DP

Full DP



Training ML Model

Ensure that the model is privacy-safe to be used in downstream tasks



DP-SGD



Gradient Descent

Training data X

Labeled Samples

 $(x_1, y_1), ..., (x_n, y_n)$

n

Loss function: $\ell(\hat{y}, y) \in \mathbb{R}$ Empirical loss

 $\mathcal{L}_{w}(X) := \Sigma_{i \in [n]} \, \boldsymbol{\ell}(f_{w}(X_{i}), y_{i}) \, / \,$

Training Objective

Find w that minimizes $\mathcal{L}_{w}(X)$

Gradient

$$\nabla_{w} \mathcal{L}(X) = [d\mathcal{L}(X) / dw_{1}, ..., d\mathcal{L}(X) / dw_{d}]$$

$$\nabla_{w} \mathcal{L}(X) = \sum_{i \in [n]} \nabla_{w} \mathcal{\ell}(f_{w}(x_{i}), y_{i}) / n$$

$$\eta_{t}: \text{ learning rate}$$

$$Gradient Descent (GD)$$

$$w_{0} \leftarrow \text{initial parameter}$$

$$For t = 1, ..., T:$$

$$w_{t} \leftarrow w_{t-1} - \eta_{t} \nabla_{w} \mathcal{L}(X)$$

$$Return w_{T}$$

Gradient Descent

Gradient Descent (GD)

 $w_{0} \leftarrow \text{initial parameter}$ For t = 1,...,T: $w_{t} \leftarrow w_{t-1} - \eta_{t} \left(\sum_{i \in [n]} \nabla_{w} \ell(f_{w}(x_{i}), y_{i}) \right) / n$ Return w_{T}

$$\begin{array}{ll} Gaussian Mechanism & \sigma = \frac{2\sqrt{2\ln(2/\delta)}}{\epsilon} \cdot \Delta_2(g) \\ X = (x_1, ..., x_n) & \rightarrow Analyzer & g(X) + \mathcal{N}(0, \sigma^2)^{\otimes d} \end{array} \xrightarrow{} \begin{array}{l} Use \ L_2 \ -sensitivity \\ Each \ coordinate \ is independent \end{array}$$

Theorem Assuming Range(f) $\subseteq \mathbb{R}^d$ and $\varepsilon \leq 1$, Gaussian Mechanism is (ε, δ) -DP.

Assumption: Range(g) $\subseteq \mathbb{R}^d$

(ie, g is vector-valued with real entries)

Gaussian Distribution For every real number z, $PDF_{\mathcal{N}(0, \sigma^2)}(z) \propto e^{-(z/\sigma^2)}$ **Examples: Vector summation**

$$\left.\begin{array}{c} & (0.2, -1, 1) \\ \vdots \\ & (0, 2, -0.1) \end{array}\right\} g(X) = (50.1, 2.3, 14.7)$$

$$\left.\begin{array}{c} & (0, 2, -0.1) \\ & (0, 2, -0.1) \end{array}\right\}$$
Assumption: Each vector has l_2 -norm $\leq C$

$$\left.\begin{array}{c} & \Delta_2(g) \leq 2C \end{array}\right.$$

Gradient Descent

Additional parameters

σ: noise standard deviation

Add Gaussian noise to average gradients!

 $\begin{array}{l} \textbf{Gradient Descent (GD)} \\ w_{0} \leftarrow \text{initial parameter} \\ \text{For } t = 1, ..., T \\ w_{t} \leftarrow w_{t-1} - \eta_{t} \left(\boldsymbol{\mathcal{N}}(0, \sigma^{2} \cdot \mathbf{I}) + \boldsymbol{\Sigma}_{i \in [n]} \boldsymbol{\nabla}_{w} \boldsymbol{\ell}(f_{w}(x_{i}), y_{i}) \right) / \\ n \\ \text{Return } w_{T} \end{array}$

Not DP: each $\nabla_{w} \ell(f_{w}(x_{i}), y_{i}))$ can be arbitrarily large!

Additional parameters **Clipping Trick** σ : noise standard deviation C: clipping norm DP-GD w₀←initial parameter For t = 1,...,T: For i = 1 n: $\mathbf{v}_{i} = \boldsymbol{\nabla}_{w} \boldsymbol{\ell}(\mathbf{f}_{w}(\mathbf{x}_{i}), \mathbf{y}_{i})$ **Enforce:** gradient $\tilde{\mathbf{v}}_{i} = \mathbf{v}_{i} \cdot \min(1, \mathbb{C} / \|\mathbf{v}_{i}\|_{2})$ norm bound at most C $W_t \leftarrow W_{t-1} - \eta_t \left(\mathcal{N}(0, \sigma^2 \cdot I) + \Sigma_{i \in [n]} \tilde{V}_i \right) / n$ Return w_T

Not DP: each $\nabla_{w} \ell(f_{w}(x_{i}), y_{i}))$ can be arbitrarily large!

"If $\nabla_{w} \ell(f_{w}(x_{i}), y_{i})$) is too large, rescale it to be smaller" (bias-variance tradeoff)



DP-GD: Privacy Analysis

DP-GD w₀←initial parameter For t = 1....T: For i = 1, ..., n: $\mathbf{v}_{i} = \boldsymbol{\nabla}_{w} \boldsymbol{\ell}(\mathbf{f}_{w}(\mathbf{x}_{i}), \mathbf{y}_{i})$ $\tilde{\mathbf{v}}_{i} = \mathbf{v}_{i} \cdot \min(1, \mathbb{C} / \|\mathbf{v}_{i}\|_{2})$ $\mathbf{w}_{t} \leftarrow \mathbf{w}_{t-1} - \eta_{t} \left(\boldsymbol{\mathcal{N}}(0, \sigma^{2} \cdot \mathbf{I}) + \boldsymbol{\Sigma}_{i \in [n]} \tilde{\mathbf{v}}_{i} \right) / \mathbf{n}$ Return W_{T}

Each iteration

Gaussian mechanism:

- ℓ_2 -sensitivity $\leq 2C$
- Noise std σ
- $\Rightarrow (\varepsilon_{g'}, \delta_{g})$ -DP



Entire Algorithm: (ϵ, δ) -DP

Gradient Descent & Friends



$$\begin{split} & \underset{w_{0} \leftarrow \text{initial parameter}}{\text{Mini-batch SGD}} \\ & \underset{w_{0} \leftarrow \text{initial parameter}}{\text{For } t = 1, \dots, T}: \\ & \text{For } t = 1, \dots, T: \\ & \text{S} \leftarrow \text{random index set of size } B \\ & \underset{w_{t} \leftarrow w_{t-1} - \eta_{t} \Sigma_{i \in S} \nabla_{w} \ell(f_{w}(x_{i}), y_{i}) / \\ & B \\ & \text{Return } w_{T} \end{split}$$

Modifications to achieve DP

Add Gaussian noise to average gradients!

$$\begin{array}{l} \underset{w_{0} \leftarrow \text{initial parameter}}{\text{For t} = 1, \dots, T} \\ \text{S} \leftarrow \text{random index set of size B} \\ w_{t} \leftarrow w_{t-1} - \eta_{t} \Sigma_{i \in S} \nabla_{w} \boldsymbol{\ell}(f_{w}(x_{i}), y_{i}) \ / \\ B \\ \text{Return } w_{T} \end{array}$$

Modifications to achieve DP

Add Gaussian noise to average gradients!

Additional parameters

σ: noise standard deviation

$$\begin{array}{l} \underset{w_{0} \leftarrow \text{initial parameter}}{\text{For } t = 1, \dots, T} & \textbf{Issue: gradient can be very large!} \\ & \text{S} \leftarrow \text{random index set of size B} \\ & w_{t} \leftarrow w_{t-1} - \eta_{t} \left(\boldsymbol{\mathcal{N}}(\boldsymbol{0}, \sigma^{2} \cdot \mathbf{I}) + \boldsymbol{\Sigma}_{i \in S} \boldsymbol{\nabla}_{w} \boldsymbol{\ell}(f_{w}(x_{i}), y_{i}) \right) / \\ & \text{B} \\ & \text{Return } w_{T} \end{array}$$

Modifications to achieve DP

- Add Gaussian noise to average gradients!
- Clip each gradient to bound its norm

Additional parameters

- σ: noise standard deviation
- C: clipping norm bound

Mini-batch SGD w₀←initial parameter For t = 1,...,T $S \leftarrow$ random index set of size B For $i \in S$: **Enforce:** gradient norm bound at most C $\mathbf{v}_{i} = \boldsymbol{\nabla}_{w} \boldsymbol{\ell}(\mathbf{f}_{w}(\mathbf{x}_{i}), \mathbf{y}_{i})$ $\tilde{\mathbf{v}}_{i} = \mathbf{v}_{i} \cdot \min(1, \mathbb{C} / \|\mathbf{v}_{i}\|_{2})$ $\mathbf{w}_{t} \leftarrow \mathbf{w}_{t-1} - \eta_{t} \left(\boldsymbol{\mathcal{M}}(0, \sigma^{2} \cdot \mathbf{I}) + \boldsymbol{\Sigma}_{i \in S} \tilde{\mathbf{v}}_{i} \right) / \mathbf{B}$ Return w_{T}

Mini-batch DP-SGD w_0 – initial parameter For t = 1,...,T: S — random index set of size B For $i \in S$: $\mathbf{v}_{i} = \boldsymbol{\nabla}_{w} \boldsymbol{\ell}(\mathbf{f}_{w}(\mathbf{x}_{i}), \mathbf{y}_{i})$ $\tilde{\mathbf{v}}_{i} = \mathbf{v}_{i} \cdot \min(1, C / \|\mathbf{v}_{i}\|_{2})$ $\mathbf{w}_{t} \leftarrow \mathbf{w}_{t-1} - \eta_{t} \left(\boldsymbol{\mathcal{M}}(0, \sigma^{2} \cdot \mathbf{I}) + \boldsymbol{\Sigma}_{i \in S} \tilde{\mathbf{v}}_{i} \right) / \mathbf{B}$ Return w_T

DP-SGD: Privacy Analysis

Mini-batch DP-SGD w₀←initial parameter For t = 1,...,T: $S \leftarrow$ random index set of size B For $i \in S$: $\mathbf{v}_{i} = \boldsymbol{\nabla}_{u} \boldsymbol{\ell}(\mathbf{f}_{u}(\mathbf{x}_{i}), \mathbf{y}_{i})$ $\tilde{\mathbf{v}}_{i} = \mathbf{v}_{i} \cdot \min(1, C / \|\mathbf{v}_{i}\|_{2})$ $\mathbf{w}_{t} \leftarrow \mathbf{w}_{t-1} - \eta_{t} \left(\boldsymbol{\mathcal{M}}(0, \sigma^{2} \cdot \mathbf{I}) + \boldsymbol{\Sigma}_{i \in S} \tilde{\mathbf{v}}_{i} \right) / \mathbf{B}$ Return w_T

Each iteration

Gaussian mechanism:

- ℓ_2 -sensitivity $\leq 2C$
- Noise std σ

•
$$\Rightarrow (\varepsilon_{\sigma}, \delta_{\sigma})$$
-DP



Entire Algorithm: (ϵ, δ) -DP

DP-SGD: Privacy Analysis


DP-SGD: Privacy Analysis

Mini-batch SGD w₀←initial parameter For t = 1,...,T $S \leftarrow$ random index set of size B For $i \in S$: $\mathbf{v}_{i} = \boldsymbol{\nabla}_{w} \boldsymbol{\ell}(\mathbf{f}_{w}(\mathbf{x}_{i}), \mathbf{y}_{i})$ $\tilde{\mathbf{v}}_{i} = \mathbf{v}_{i} \cdot \min(1, C / \|\mathbf{v}_{i}\|_{2})$ $\mathbf{w}_{t} \leftarrow \mathbf{w}_{t-1} - \eta_{t} \left(\boldsymbol{\mathcal{M}}(0, \sigma^{2} \cdot \mathbf{I}) + \boldsymbol{\Sigma}_{i \in S} \tilde{\mathbf{v}}_{i} \right) / \mathbf{B}$ Return w_T

Each iteration

Gaussian mechanism:

- ℓ_2 -sensitivity $\leq 2C$
- Noise std σ
- $\Rightarrow (\varepsilon_{g'}, \delta_{g})$ -DP



Entire Algorithm: (ϵ, δ) -DP

DP-SGD: Privacy Analysis

Mini-batch SGD w₀←initial parameter For t = 1,...,T $S \leftarrow$ random index set of size B For $i \in S$: $g_i = \nabla_w \ell(f_w(x_i), y_i)$ $g_i = g_i \cdot \min(1, C / ||g_i||_2)$ $W_t \leftarrow W_{t-1} - \eta_t \left(\mathcal{N}(0, \sigma^2 \cdot I) + \Sigma_{i \in S} g_i \right) / B$ Return w_T

Each iteration

Gaussian mechanism:

- ℓ_2 -sensitivity $\leq 2C$
- Noise std σ

•
$$\Rightarrow (\varepsilon_{\sigma}, \delta_{\sigma})$$
-DP





Entire Algorithm: (ϵ, δ) -DP

Merits of DP-SGD

- Generic recipe
- Broadly applicable

Possible improvements

- Gradient sparsity [Ghazi et al.'23]
- Privacy accounting

- Might incur (large) utility drop
- Increased training time

TensorFlow Privacy

Ads Modeling [Denison et al,'23]

- Noise only added once per batch
 - Bigger batches ⇒ Less noise per example
- Large batches often take more epochs to converge

Ads Modeling [Denison et al,'23]

Clipping is a bias-variance tradeoff

- Noise is scaled with clip norm
- Clipping gradients loses signal
- Tune clip norm using fixed batch size



Ads Modeling [Denison et al,'23]

Competitive Loss with DP

- +12.1% Loss @ ε = 10
- +13.5% Loss @ ε = 1
- +15.8% Loss @ ε = 0.5
- Compute needs increased by 20%

Privacy-Utility for Probability of Ad Click (pCTR)



Privacy Cost (Epsilon)

Label DP

Privacy in Web Advertising: Analytics and Modeling

Ghazi & Manurangsi

Learning with DP



Learning with Label DP [Chaudhuri-Hsu '11]

unknown data



Label DP Algorithms

- RR [Warner'65]
- RR-with-Prior [Ghazi et al. '21]
- RR-on-Bins [Ghazi et al. '23]
- Unbiased Mechanism [Badanidiyuru et al.'23]
- Logistic Regression [Meta'21]
- PATE-FM [Malek et al. '21]
- ALIBI [Malek et al. '21]
- Clustering [Esfandiari et al. '22]

Our focus

Feature-Oblivious Label DP



- Generate DP (noisy) labels $\tilde{y}_1, ..., \tilde{y}_m$ using M
- Run standard training to minimize loss using noisy labels:

 $\boldsymbol{\Sigma}_{i \in [n]} \, \boldsymbol{\ell}(\boldsymbol{f}_w(\boldsymbol{x}_i), \, \boldsymbol{\tilde{y}}_i))$

- Laplace Mechanism on each y_i
- Randomized Response on each y_i

Choose ε -DP mechanism M to minimize $\mathbf{E}_{y, \tilde{y} \sim M(y)} L(\tilde{y}, y)$

- Privately learn a prior over labels y, to construct M
- For classification, $L(\tilde{y}, y) = \mathbf{1}[\tilde{y} \neq y]$
- For regression, eg, $L(\tilde{y}, y) = (\tilde{y} y)^2$

RR-with-Prior (Classification) [Ghazi et al'21]

Randomized Response (RR) [Warner'65]

Raw Data (x, y) Privatized Data (x, y')

$$\Pr[y' = a] = \begin{cases} e^{\varepsilon} / (e^{\varepsilon} + K - 1) & \text{if } a = y \\\\ 1 / (e^{\varepsilon} + K - 1) & \text{if } a \neq y \end{cases}$$

Labels $\in \{1, ..., K\}$



- Run RR only on the top **k** labels (eg, $k \le 3$)
- Can help reduce the noise, for same ε
- This is optimal for the objective of maximizing Pr[y' = y], assuming p is the "true prior" and k is chosen appropriately

Optimality of RR-with-Prior

Choose
$$\varepsilon$$
-DP mechanism M to minimize $\mathbf{E}_{y\sim p, \ \tilde{y}\sim M(y)}$ $\mathbf{1}[\tilde{y} \neq y]$
M can be written as a K x K matrix
 $M(\tilde{y}, y) = Pr[M \text{ outputs } \tilde{y} \text{ on input } y]$
maximize $\Sigma_{y \in [K]} p_y \cdot M(y, y)$ subject to
(i) $\Sigma_{\tilde{y} \in [K]} M(\tilde{y}, y) = 1, \ \forall y \in [K]$
(ii) $M(\tilde{y}, y) \ge 0$
(iii) $M(\tilde{y}, y) \le e^{\varepsilon} \cdot M(\tilde{y}, y'), \ \forall \ \tilde{y}, y, y' \in DP$ condition

Ghazi & Manurangsi

Privacy in Web Advertising: Analytics and Modeling

RR-on-Bins (Regression) [Ghazi et al'23]

Partition the range of possible y values into intervals ("bins"), and assign a representative value to each bin



- $M(y) := \varepsilon RR$ over $\Phi(y)$, where $\Phi(y) =$ representative value for bin containing y
- M computed efficiently using dynamic programming in time O(KD²), where
 - D = #possible y values (discretize if values are continuous)
 - K = upper bound on number of bins (\leq D)
- M(·) on distribution p is optimal for minimizing $\mathbf{E}_{v \sim p, \tilde{v} \sim M(v)} L(\tilde{y}, y)$
 - Characterize the optimal solutions to an underlying LP
- Led to <u>flexible reports</u> on Privacy Sandbox Chrome/Android APIs

Criteo Conversion Log Dataset

• Criteo Sponsored Search Conversion Log Dataset:

90 days of Criteo live traffic data, with ~15M examples.

ailab.criteo.com/criteo-sponsored-search-conversion-log-dataset/

- Goal: Predict conversion value (in €) (clipped to €400 for simplicity)
- Regression task / dataset

Results on Criteo Conversion Log Dataset



Privacy in Web Advertising: Analytics and Modeling

Ghazi & Manurangsi

Bias-Variance Trade-offs

Challenge: Sometimes RR-on-Bins do not give best results after training

- E.g. sometimes less bias \Rightarrow better accuracy [Badanidiyuru et al. '23]
- Problem can still be formulated as a linear program
- Isn't known to have an explicit solution

Unbiased Mechanism

Regression [Ghazi et al. '23]





Classification [Ghazi et al. '21]



Regression [Badanidiyuru et al.'23]



Prior P



Optimal Unbiased Mechanism M for $\epsilon = 0.5$



Privacy in Web Advertising: Analytics and Modeling

Ghazi & Manurangsi

Motivation for Unbiased Noisy Labels



Regression [Badanidiyuru et al.'23]



Minimizing variance, while having zero bias.

Zero bias preserves the Bayes Optimal Predictor



• Zero bias provides unbiased stochastic gradients

 $\mathop{\mathbb{E}}\limits_{y'\sim M(y)}
abla_ heta\ell(f_ heta(x),y')=
abla_ heta\ell(f_ heta(x),y)$

Since gradient is affine in the label: $\nabla_{\theta}\ell(f_{\theta}(x),y) = f_{\theta}(x) \cdot \nabla_{\theta}f_{\theta}(x) - y \cdot \nabla_{\theta}f_{\theta}(x)$

Final mechanism: Using privately estimated prior



Apply ε_{label} -DP randomizer M to every label in training set.

Evaluation on Criteo Conversion Log Dataset



Prediction loss on test data:
$$\frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} (f_w(x_i) - y_i)^2$$



Privacy in Web Advertising: Analytics and Modeling

Multi-Stage Training



Privacy in Web Advertising: Analytics and Modeling

Multi-Stage Training



Privacy in Web Advertising: Analytics and Modeling

Ghazi & Manurangsi

Results on CIFAR-10 (Classification)



Constructing priors

- Self-supervised training (SSL): BYOL, DINO
- Clustering and DP histograms
- Multi-stage training with data splitting
- Multi-stage training with privacy budget splitting

Significantly narrowed the performance gap between private and non-private models

Logistic Regression

- For logistic loss, gradient depends on labels Y only through additive term proportional to X^TY.
- For given features X, quantity X^TY is a linear function of labels Y that can be estimated using Discrete Laplace Mechanism (e.g., via summary reports in Privacy Sandbox ARA)
 - Would satisfy label DP in central model
- Similar idea used in Meta's winning solution to the Criteo AdKDD 2021 competition
 <u>https://medium.com/criteo-engineering/results-from-the-criteo-adkdd-2021-challenge-50</u>

 <u>abc9fa3a6</u>

Label DP Logistic Regression

Model type	Privacy ε	Optimizer, learning rate, batch size	AUC
Private logistic			
regression with			
noise	3	yogi, 0.01, 1024	0.74277
Logistic			
regression			
without noise	NA	yogi, 0.01, 1024	0.77057
Logistic			
baseline	NA	yogi, 0.01, 1024	0.7646
RR + logistic			
baseline	3	yogi, 0.01, 1024	0.75745
RR + MLP	3	yogi, 0.01, 1024	0.80212
MLP	NA	yogi, 0.01, 1024	0.805

• Results on Criteo pCTR dataset

•

For same value of ɛ, randomized response substantially outperforms label DP logistic regression

Privacy in Web Advertising: Analytics and Modeling

Training with (User) Label Differential Privacy

unknown data



and User x Time privacy units



Privacy in Web Advertising: Analytics and Modeling

Gnazi & ivianurangs

Handling Multiple Impressions Per Privacy Unit

Example: Consider User x Time privacy unit.

Cap # of impressions per user and time period to K (keeping K random impressions, or K first impressions). Then, we have multiple options including:

1. For each user, set the privacy budget per impression to ϵ/K .

2. For each user i with $K_i \le K$ impressions, set the privacy budget per impression to ϵ/K_i . Both options satisfy ϵ -Label-DP for User x Time privacy unit.

Similar options hold for User x Advertiser x Time & User x Publisher x Time privacy units.

For Impression x Time privacy unit, no capping is needed. RR is applied privacy budget ϵ .

Criteo Attribution Modeling for Bidding Dataset

https://ailab.criteo.com/criteo-attribution-modeling-bidding-dataset/

- Sample of 30 days of Criteo live traffic data.
- Each example corresponds to a click and contains:
 - **Features:** campaign ID, 9 contextual features, and the cost paid for the display.
 - Label: a 0/1 field indicating whether there was a conversion in the 30 days after the click and that is last-touch attributed to this click.
 - User ID: can be used to evaluate User x Time privacy unit.
- Number of rows is 5,947,563. Conversion rate (under last-touch attribution) is 6.74%.
- Feed-forward neural network
 - Embedding dimension is 8
 - \circ Hidden dimensions are 128 x 64.
- Tune hyperparameters with optimizer in {rmsprop, adam, sgd}, learning rate in {0.0005, 0.0008, 0.001, 0.002, 0.005}, and training epochs in {100, 200}.

Criteo Attribution Modeling for Bidding – Statistics



Criteo Attribution Modeling for Bidding – Evaluation Results



Notes

- For Impression x Time privacy unit and ε = 4, relative AUC loss is 0.79%.
- For User x Time privacy unit with $\varepsilon = 4$, smallest relative AUC loss is 8.51%.
- For User x Time privacy unit, smaller loss is achieved by increasing caps as we increase ε.

Proprietary Ads Dataset

- App install ads
- Contains data from multiple advertisers and publishers.
 - Can be used to evaluate User x Time, User x Publisher x Time, and User x Advertiser x Time privacy units
- Features: use categorical features, and pass the concatenation of their embeddings through multiple layers of a fully connected feedforward neural network.
- Labels: 0/1 corresponding to installs (= conversions)

Proprietary Ads Dataset – Evaluation Results



Notes

- For Impression x Time privacy unit and ε = 3, relative AUC loss is 0.83%.
- For User x Time privacy unit and $\varepsilon = 3$, smallest relative AUC loss is 4.50%.
- For User x Publisher x Time privacy unit and ε = 3, smallest Relative AUC loss is 2.67%.
- For User x Advertiser x Time unit with $\varepsilon = 3$, best Relative AUC loss is 1.56%.
- In this experiment, for same ε, loss for Impression x Time
 - < loss for User x Advertiser x Time
 - < loss for User x Publisher x Time
 - < loss for User x Time.

Evaluation limitations and Future Directions

- Utility might be improvable for binary conversion models
 - Debiased loss functions
- Evaluation focused on binary conversion models
 - Label DP algorithms (and RR in particular) can be extended to non-binary predictions (such as predicting number of conversions, and conversion value etc.)
- Evaluation was done offline and assumed one reporting window
 - Online performance can be impacted by delays, which might require multiple reporting windows (and more noise for the same privacy)

High-Level Overview of APIs: Event

Device

Impr 1 Impr 2 Conv 1 Impr 3 Conv 2 LTA			
Discretized attribution information			

Summary

- Attributions happen on device / browser
- Discretize attribution information to a finite space
High-Level Overview of APIs: Event

Device

	Summary
Impr 1 Impr 2 Conv 1 Impr 3 Conv 2	 Attributions happ
LTA	Discretize attribu
Attributed Dataset: (Impr2, Conv1), (Impr3, Conv2)	• Apply RR & send
Discretized attribution information	
Randomized Response	
Noisy Discretized attribution information	Report Collector

Summary

- Attributions happen on device / browser
- Discretize attribution information to a finite space
- Apply RR & send the output to report collector

ARA Event-Level Reports



Total # of states = $\begin{pmatrix} w \cdot 2^b + C \\ C \end{pmatrix}$

- RR at impression level
- There are w (≤ 5) "reporting (time) windows"
 - E.g. reporting windows = 1d, 3d, 7d means that reports will be sent 1d, 3d, 7d after the impression.
- Up to C (\leq 20) conversions per impression:
 - For each conversion, up to b (\leq 3) bits of

"metadata" information

Report

Collector

Privacy in Web Advertising: Analytics and Modeling

Example Usage of Event-Level Reports



Privacy in Web Advertising: Analytics and Modeling

Ghazi & Manurangsi

Example Usage of Event-Level Reports



Privacy in Web Advertising: Analytics and Modeling

Ghazi & Manurangsi

Event Reports: Utility Optimization & Challenges

Discretizing the values

- Shown in [Ghazi et al.'23] that "RR-on-bins" are optimal for certain loss between randomized labels and true labels: Led to <u>Flexible Reports</u> on Chrome/Android
- Optimizing for the "bins" can be done efficiently using historical data (or private aggregate queries)
- Training methods that are robust to noisy labels:
 E.g. regularization, debiasing, self-supervised
 learning
- Challenge: Sometimes these do not give best results after training: E.g. sometimes less bias ⇒ better accuracy [Badanidiyuru et al. 23]

Extending APIs to support more

complex algorithms:

- Other Label DP algorithms
 - RRWithPrior [Ghazi et al.'21]
 - PATE-FM [Malek et al. '21]
 - ALIBI [Malek et al. '21]
 - Clustering [Esfandiari et al. '22]
- **Research question:** other better *Label DP* algorithms?

DP with Partially Known Features



Privacy Sandbox Protected Audience API on Chrome and Android

- Some features depend on cross-site information and use 3p cookies, some do not
- E.g. remarketing ads use case

DP Training with Partially Known Features

Hybrid Algorithm: Label-DP phase followed by DP-SGD phase [Chua et al. '24]

Label-DP Phase: Train truncated model with randomized response labels and unknown embeddings set to 0, with (ε_1 , 0)-DP

DP-SGD Phase: Train entire model with (ε_2 , δ)-DP

Total privacy budget (ε , δ) is split between two phases as: $\varepsilon_1 := \min \{0.6 \varepsilon, 3\}$

 $\varepsilon_2 := \varepsilon - \varepsilon_1$

Two baselines:

- 1) **RR**: labels privatized with noise, unknown features discarded ($\varepsilon_1 = \varepsilon, \varepsilon_2 = 0$)
- 2) **DP-SGD**: all features treated as unknown ($\varepsilon_1 = 0, \varepsilon_2 = \varepsilon$)

We train binary classification models with binary cross entropy loss and report the test AUC loss % relative to the non-private baseline.

Hybrid DP Algorithm Evaluation

- **Evaluation 1:** Criteo Display Ads pCTR Dataset
- 40M examples over 7 days of Criteo traffic kaggle.com/c/criteo-display-ad-challenge/overview
- Treat even-numbered features as unknown and odd-numbered features as known
- Limit # of DP-SGD epochs to ≤3 because of high cost
- Goal: Predict probability of click

Evaluation 2: Criteo Attribution Modeling for Bidding Dataset

- 16M impressions from 30 days of Criteo traffic ailab.criteo.com/criteo-attribution-modeling-bidding-dataset/
- Treat cat1, cat2 as unknown, cat[3-9] and campaign as known.
- Goal: Predict probability of conversion
- Evaluate both example-level DP and user-level DP



Ghazi & Manurangsi

Privacy in Web Advertising: Analytics and Modeling

User level DP with Partially Known Features

Translate example-level DP to user-level DP via group privacy:

If a mechanism *M* satisfies (ε, δ) -DP, then the mechanism $M^{(k)}$ obtained by retaining $\leq k$ examples per user satisfies $(k\varepsilon, \delta(e^{k\varepsilon} - 1)/(e^{\varepsilon} - 1))$ -user-level-DP

Hybrid algorithm improves over RR and DP-SGD over a wide range of privacy budgets, especially when $\varepsilon \ge 3$.



DP with a Known Set of Feature Values

Alternative notion applicable in some ad modeling settings

- Only set of values for a feature is known (but not the specific value on a per-example basis)
- E.g., the set of all ads are known to the ad-tech but not the specific feature associated with each example
- Provides stronger privacy guarantee than DP with partially known features, but by definition cannot provide higher utility
- Recent algorithm in this setting [Krichene et al. '23]

Some Practical Insights

- DP-SGD is feasible even for tasks with very sparse gradients, enabled by
 - Advances in compilation methods
 - Larger batch sizes
 - Tighter privacy accounting, eg, Privacy Loss Distributions (PLD)
 - Tailored algorithms
- Label DP often allows superior utility
- Tailoring privacy model to application provides further improvements
- Prior information on labels can be
 - Leveraged to improve utility
 - Obtained from historical data, current model predictions, or clustering and histograms
- Known features could improve utility

Modeling Research Questions

- Better algorithms for ads modeling with full DP, label DP, and DP with partially known features?
- Label DP
 - Full characterization of optimal unbiased randomizers?
 - Better algorithms for computing optimal unbiased randomizers?
- Algorithms for when some features are partially known while others belong to a known set of values?

Conclusion & Future Directions

Privacy in Web Advertising: Analytics and Modeling



Open-Source Libraries

- Some open-source DP libraries:
 - Generic DP Libraries:
 - Google DP Library
 - IBM Diffprivlib Library
 - OpenDP Library
 - DP ML Libraries:
 - Tensorflow privacy
 - Pytorch Opacus

Differential Privacy

i Note

If you are unfamiliar with differential privacy (DP), you might want to go through "A friendly, non-technical introduction to differential privacy".

This repository contains libraries to generate ε - and (ε , δ)-differentially private statistics over datasets. It contains the following tools.

- Privacy on Beam is an end-to-end differential privacy framework built on top of Apache Beam. It is intended to be easy to use, even by non-experts.
- Three "DP building block" libraries, in C++, Go, and Java. These libraries implement basic noise addition
 primitives and differentially private aggregations. Privacy on Beam is implemented using these libraries.
- A stochastic tester, used to help catch regressions that could make the differential privacy property no longer hold.
- A differential privacy accounting library, used for tracking privacy budget.
- A command line interface for running differentially private SQL queries with ZetaSQL.

To get started on generating differentially private data, we recommend you follow the Privacy on Beam codelab.

Currently, the DP building block libraries support the following algorithms:

Algorithm	C++	Go	Java
Laplace mechanism	Supported	Supported	Supported
Gaussian mechanism	Supported	Supported	Supported
Count	Supported	Supported	Supported
Sum	Supported	Supported	Supported
Mean	Supported	Supported	Supported
Variance	Supported	Supported	Supported
Standard deviation	Supported	Supported	Planned
Quantiles	Supported	Supported	Supported

github.com/google/differential-privacy

Ghazi & Manurangsi

Pointers

- Private Attribution Measurement APIs:
 - List of proposals (beyond attribution): see <u>W3C github repo</u>
 - Join WICG: <u>https://www.w3.org/community/wicg/</u>
 - Join meetings / discussions on github (e.g. for ARA see <u>here</u>)
- Project Halo
 - WFA project website
 - Open-source github repo for Cross-media measurement

Conclusion

- Formal privacy guarantees for ad analytics and modeling functionalities are possible
- Better algorithms are pushing the privacy-utility trade-off curve
- Well-defined functionalities and utility metrics are invaluable for designing privacy-preserving algorithms

Future Directions

- General-purpose DP synthetic ad data remains an important research problem
- Combining known data that the analyst might have access to with outputs of privacy-preserving algorithms and APIs is an important future direction.