

University of Idaho

**CS 404/504**

**Special Topics: Adversarial  
Machine Learning**

*Dr. Alex Vakanski*

# Lecture 12

## Defenses against Privacy Attacks

# Lecture Outline

---

- Defenses against Privacy Attacks
  - Anonymization techniques
  - Encryption techniques
  - Differential privacy
  - Distributed learning
  - ML-specific techniques
- Introduction to differential privacy
- Differentially private SGD
- Scalable private learning with PATE

# Defenses against Privacy Attacks

---

## *Defenses against Privacy Attacks*

- **Data privacy** techniques have the goal of allowing analysts to learn about *trends* in data, without revealing information specific to *individual data instances*
  - Therefore, privacy techniques involve an **intentional** release of information, and attempt to control what can be learned from the released information
- Related to data privacy is the **Fundamental Law of Information Recovery**, which states that “*overly accurate estimates of too many statistics can completely destroy privacy*”
  - I.e., extracting useful information from a dataset (e.g., for training an ML model) poses a privacy risk to the data
- There is an inevitable trade-off between privacy and accuracy (i.e., utility)
  - Preferred privacy techniques should provide an estimate of how much privacy is lost by interacting with data

# Defenses against Privacy Attacks

---

## *Defenses against Privacy Attacks*

- Defense strategies against privacy attacks in ML can be broadly classified into:
  - Anonymization techniques
  - Encryption techniques
  - Differential privacy
  - Distributed learning
  - ML-specific techniques

# Anonymization Techniques

## *Anonymization Techniques*

- **Anonymization** techniques provide privacy protection by removing identifying information in the data
- E.g., remove personal identifiable information (PII)
  - In the example below, the Name and Address columns are masked

User ID	Name	Address	Account Type	Subscription Date
001	Alice	123 A St	Pro	01/02/20
002	Bob	234 B St	Free	02/03/21
003	Charlie	456 C St	Pro	03/04/18

User ID	Name	Address	Account Type	Subscription Date
001			Pro	01/02/20
002			Free	02/03/21
003			Pro	03/04/18

# Anonymization Techniques

## Anonymization Techniques

- Anonymization is not an efficient defense method, since the remaining information in the data can be used for identifying the individual data instances
  - For example, based on health records (including diagnoses and prescriptions) with removed personal information released by an insurance group in 1997, a researcher extracted the information for the Governor of Massachusetts
    - This is referred to as *de-anonymization*
  - The same researcher later showed that 87% of all Americans can be uniquely identified using 3 bits of information: ZIP code, birth date, and gender

Dataset 1: Users medical database

User ID	Name	Address	Zip Code	Birth date	Gender	Probable disease ID
001	Alice	123 A St	83401	01/02/1997	F	120
002	Bob	234 B St	83402	02/03/1995	M	35
003	Charlie	456 C St	83403	03/04/1999	M	240

Dataset 2: Users medical database with name and address removed

User ID	Zip Code	Birth date	Gender	Probable disease ID
001	83401	01/02/1997	F	120
002	83402	02/03/1995	M	35
003	83403	03/04/1999	M	240

# Linkage Attack

## Anonymization Techniques

- De-anonymization of data by using connections to external sources of information is referred to as *linkage attack*
  - For example:
    - In 2006, Netflix published anonymized 10 million movie rankings by 500,000 customers
    - Two researchers showed later that by using movie recommendations on IMDb (Internet Movie Database) they could identify the customers in the Netflix data

Dataset 1: Anonymized dataset with removed personal information

User ID	Name	Address	Account Type	Subscription Date
001	<input type="text"/>	<input type="text"/>	Pro	01/15/20
002	<input type="text"/>	<input type="text"/>	Pro	02/03/21
003	<input type="text"/>	<input type="text"/>	Free	03/04/18

Dataset 2: External public dataset that reveals the users in Dataset 1

User ID	Product Name	Product Price	Purchase Date
001	TV	400	01/02/20
002	Iphone	1,199	02/02/21
003	Watch	130	02/22/18

# $k$ -anonymity

---

## Anonymization Techniques

- $k$ -anonymity is an approach for protecting data privacy by suppressing certain identifying data features
  - This approach removes fields of data for individuals who have unique characteristics
    - E.g., students at UI who are from Latvia and are enrolled in Architecture
- A dataset is  $k$ -anonymous if for any person's record, there are at least  $k - 1$  other records that are indistinguishable
  - Therefore, a linkage attack will result in a group of  $k$  records that can belong to a person of interest
- Limitation: this approach is mostly applicable to large datasets with low-dimensional input features
  - The more input features there are for each record, the higher the possibility of unique records

# Encryption Techniques

## Encryption Techniques

- **Encryption** is a cryptography approach, which converts the original representation of information into an alternative form
  - The sender of encrypted information shares the decoding technique only with the intended recipients of the information
- **Encrypting the training data** has been applied in ML
  - Common techniques for data encryption include:
    - Homomorphic encryption (HE)
    - Secure multi-party computation (SMPC)
- **Encrypting ML models** is less common approach
  - Homomorphic encryption has been applied to the model gradients in collaborative DL setting to protect the model privacy

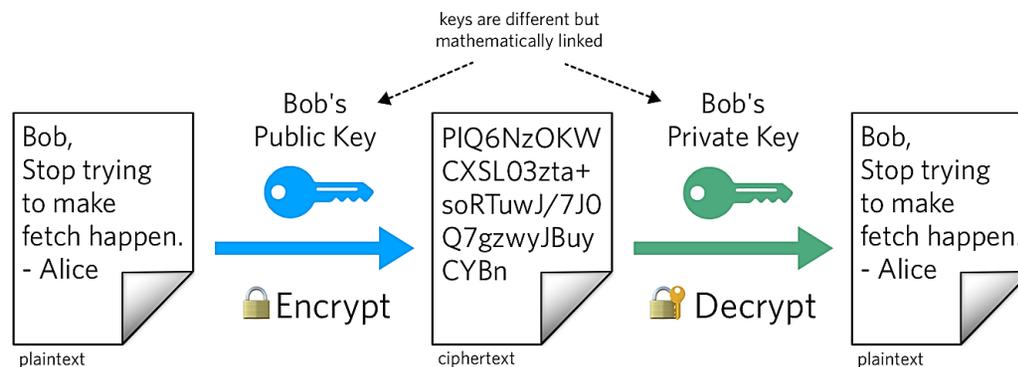


Figure form: [What is Public Key Cryptography?](#)

# Homomorphic Encryption

---

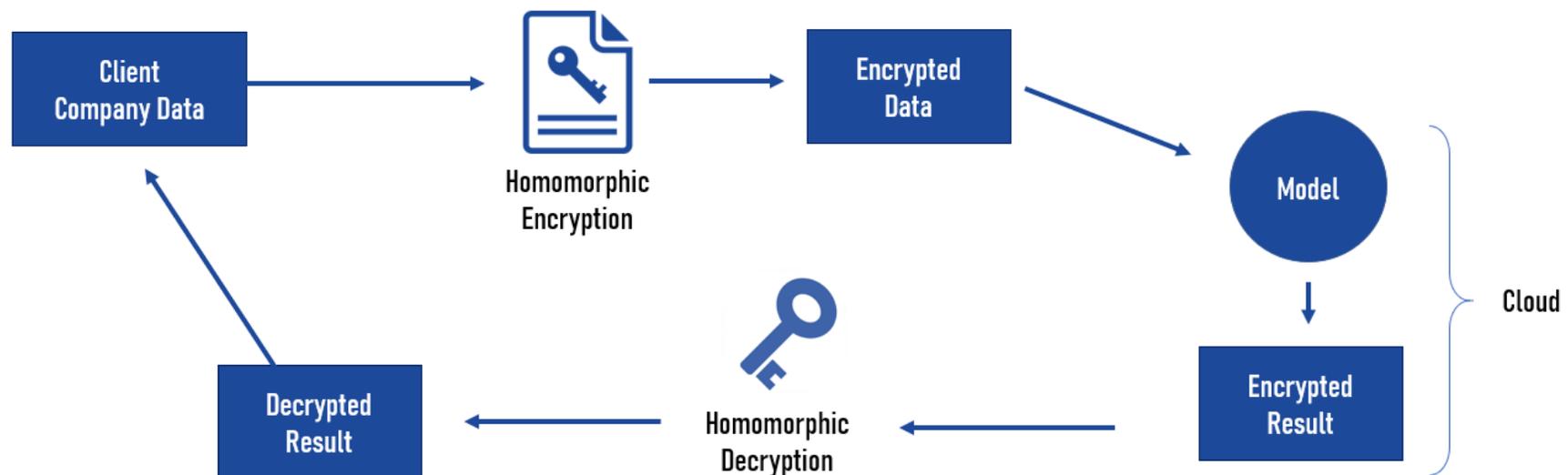
## *Encryption Techniques*

- *Homomorphic encryption (HE)* allows users to perform **computations on encrypted data** (without decrypting it)
  - Encrypted data can be analyzed and manipulated without revealing the original data
- HE uses a public key to encrypt the data, and applies an algebraic system (e.g., additions and multiplications) to allow computations while the data is still encrypted
  - Only the person who has a matching **private key** can access the decrypted results

# Homomorphic Encryption

## Encryption Techniques

- In ML, training data can be encrypted and send to a server for model training
  - Even if the server is untrusted or it is compromised, confidentiality of the data will remain preserved
  - One main limitation of HE is the slowing down of the training process
- HE has been applied to traditional ML approaches, such as Naïve Bayes, decision trees
  - Training DNNs over encrypted data is still challenging, due to the increased computational complexity



# Privacy versus Confidentiality

---

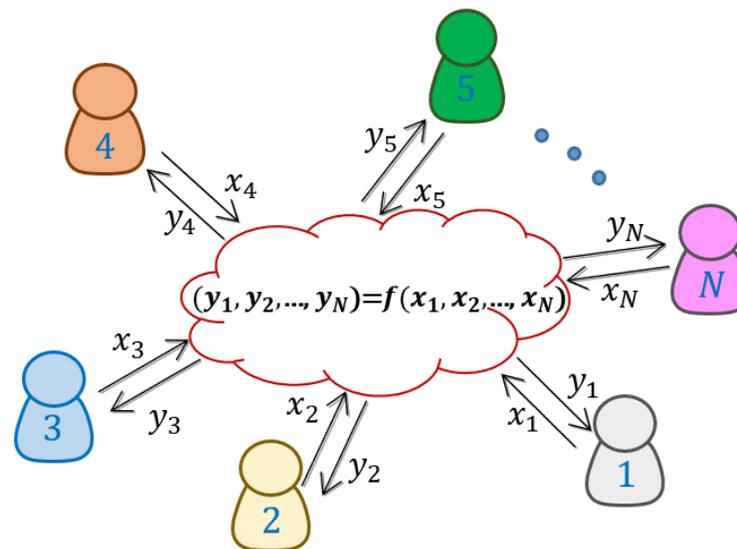
## *Encryption Techniques*

- Encryption techniques in ML are mainly applied to protect the confidentiality of the data or model
- **Confidentiality** refers to keeping the information (training data, model parameters) hidden from the clients and the public
  - It is ensuring that only authorized parties have access to the information
  - E.g., a server has an ML model trained on private data and provides the model to a client for inference
    - It is preferred to preserve the confidentiality of the model parameters from the client
- **Privacy** refers to intentional release of information in a controlled manner to prevent unintended information leakage
  - It is ensuring that released data cannot uniquely identify individual inputs
  - E.g., a server applies DP to a trained ML model to prevent memorization of information about individual inputs
- Protecting privacy is more challenging than protecting confidentiality

# Secure Multi-Party Computation

## Encryption Techniques

- **Secure Multi-Party Computation** (SMPC) is an extension of encryption in multi-party setting
  - SMPC allows two or more parties to jointly perform computation over their private data, without sharing the data
  - E.g., two banks want to know if they have both flagged the same individuals and learn about the activities by those individuals
    - The banks can share encrypted tables of flagged individuals, and they can decrypt only the matched records, but not the information for individuals that are not in both tables



# Secure Multi-Party Computation

---

## *Encryption Techniques*

- SMPC versus HE
  - **SMPC** protects the **privacy** of the data in collaborative learning
    - E.g., participants in collaborative learning do not trust the other participants or the central server
  - **HE** protects the **confidentiality** of the data from external adversaries
    - E.g., a data owner wants to use a **MLaaS (Machine Learning as a Service)**, but does not trust the service provider: the owner sends encrypted data, the provider process encrypted data and sends back encrypted results, the owner decrypts the results
    - Or, a bank can store encrypted banking information in the cloud, and use HE to ensure that only the employees of the bank can access the data

# Secure Multi-Party Computation

---

## *Encryption Techniques*

- In ML, SMPC can be used to compute updates of the model parameters by multiple parties that have access to their private data
  - For examples, SMPC has been applied to federated learning, where participants encrypt their updates, and the central server can recover only the sum of the updates from all participants
  - Beside the data privacy, SMPC also offers protection against adversarial participants
    - Either all parties are honest and can jointly compute the correct output, or if a malicious party is dishonest the joint output will be incorrect
- SMPC has been applied to traditional ML models, such as decision trees, linear regression, logistic regression, Naïve Bayes,  $k$ -means clustering
  - Application of SMPC to deep NNs is also challenging, due to increased computational costs

# Differential Privacy

---

## *Differential Privacy*

- **Differential privacy** is based on employing obfuscation mechanisms for privacy protection
  - A **randomization mechanism**  $\mathcal{M}(D)$  applies noise  $\xi$  to the outputs of a function  $f(D)$  to protect the privacy of individual data instances, i.e.,  $\mathcal{M}(D) = f(D) + \xi$
  - Commonly used randomization mechanisms include Laplacian, Gaussian, and Exponential mechanism
- DP is often implemented in practical applications
- Examples include:
  - 2014: Google's RAPPOR, for statistics on unwanted software hijacking users' settings
  - 2015: Google, for sharing historical traffic statistics
  - 2016: Apple, for improving its Intelligent personal assistant technology
  - 2017: Microsoft, for telemetry in Windows
  - 2020: LinkedIn, for advertiser queries
  - 2020: U.S. Census Bureau, for demographic data

# Differential Privacy

---

## *Differential Privacy*

- In ML, DP is achieved by adding noise to:
  - *Model parameters*
    - Several works applied DP to conventional ML methods
    - **Differentially private SGD** (Abadi, 2016) clips and adds noise to the gradients of deep NNs during training
      - This reduces the memorization of individual input instances by the model
    - The approaches that apply obfuscation to the model parameters via DP are also referred to as **differentially private ML**
  - *Model outputs*
    - **PATE** (Private Aggregation of Teacher Ensembles) approach (Papernot, 2018) employs an ensemble of models trained on disjoint subsets of the training data, called teacher models
    - Noise is added to the outputs of the teacher models, and the aggregated outputs are used to train another model, called student model
  - *Training data*
    - Obfuscation of training data in ML has been also investigated in several works

# Differential Privacy

---

## *Differential Privacy*

- DP is typically applied in a *centralized learning setting*, where the data and model are at the same location
  - In this scenario, all data is gathered in one central location for model training
  - E.g., MLaaS typically requires that the users upload their data to a cloud-based server for training a model
- Recently, DP has also been applied in a *distributed learning setting*, where the data are kept at separate locations from the model
  - **DP-FedAvg** (McMahan, 2018) is applied to federated learning
  - It introduced the Federated Averaging algorithm to limit the contributions by the individual users data to the learning model

# Distributed Learning

## *Distributed Learning*

- **Distributed learning** allows multiple parties to train a global model without releasing their private data
- Some form of **aggregation** is applied to the local updates of the model parameters by the users in distributed learning to create a global model
  - E.g., averaging is one common form of aggregation
- **Federated learning** is the most popular distributed learning scheme

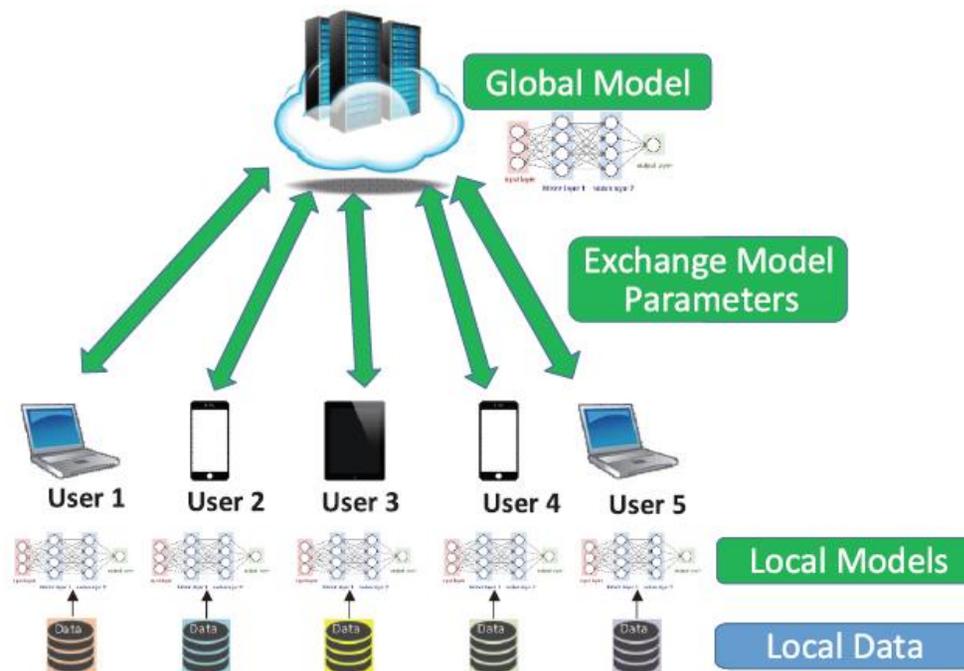


Figure from: Liu et al. (2020) When Machine Learning Meets Privacy: A Survey and Outlook

# Distributed Learning

---

## *Distributed Learning*

- ***Federated learning*** or ***collaborative learning*** – learn one global model using data stored at multiple locations (e.g., remote devices)
  - The data are processed locally, and used to update the model
    - The data does not leave the remote devices, remains private
  - The central server aggregates the updates and creates the global model
- ***Decentralized Peer-to-Peer (P2P) learning*** – the remote devices communicate and exchange the updates directly, **without a central server**
  - Removes the need to send updates to a potentially untrusted central server
- ***Split learning*** – each remote device is used to **train several layers** of the global model, and send the outputs to a central server
  - The remote devices can train the initial layers of a DNN, and the central server can train the final layers
    - The gradient is back-propagated from the central server to each user to sequentially complete the back-propagation through all layers of the model
  - The devices send the intermediate layers outputs, rather than model parameters
  - Split learning is more common for IoT devices with limited computational resources

# ML-Specific Techniques

---

## *ML-Specific Techniques*

- In the lecture on privacy attacks in ML, we mentioned that overfitting is one of the reasons for information leakage
- *Regularization techniques* in ML can therefore be used to reduce overfitting, as well as a defense strategy
  - Different regularization techniques in NNs include:
    - **Explicit regularization**: dropout, early stopping, weight decay
    - **Implicit regularization**: batch normalization
- Other ML-specific techniques include:
  - Dimensionality reduction – removing inputs with features that occur rarely in the training set
  - Weight-normalization – rescaling the weights of the model during training
  - Selective gradient sharing – in federated learning, the users share a fraction of the gradient at each update

# Introduction to Differential Privacy

Koffi Anderson Koffi



University of Idaho

## Outline

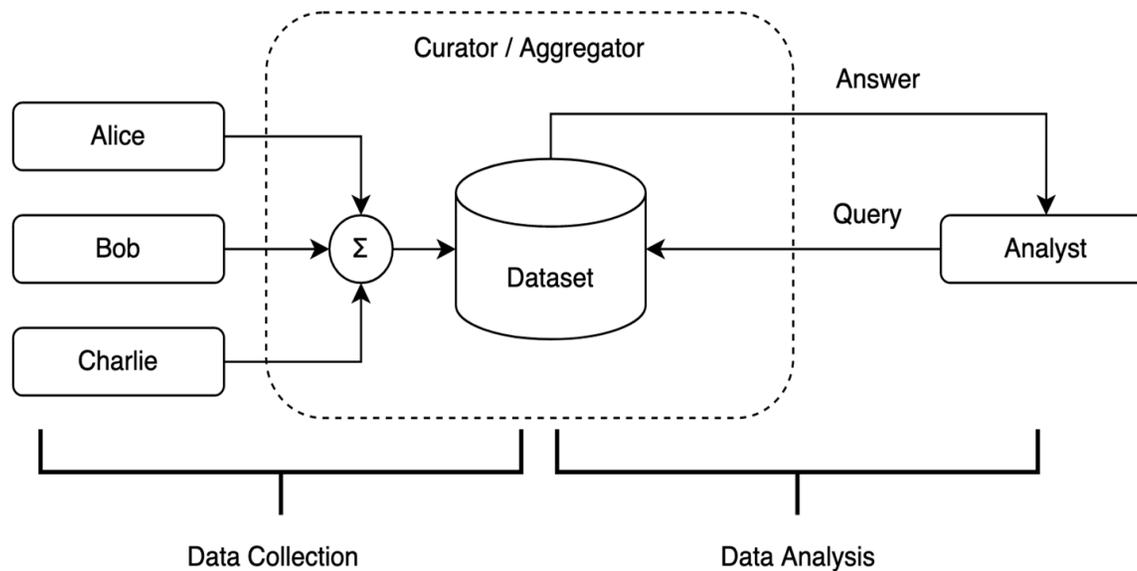
- ❖ Traditional Data Analysis
- ❖ Privacy Models
- ❖ Traditional Privacy Models
- ❖ Differential Privacy Models
- ❖ Differential Privacy
- ❖ Differential Privacy Mechanism
- ❖ Applications in Machine Learning



University of Idaho

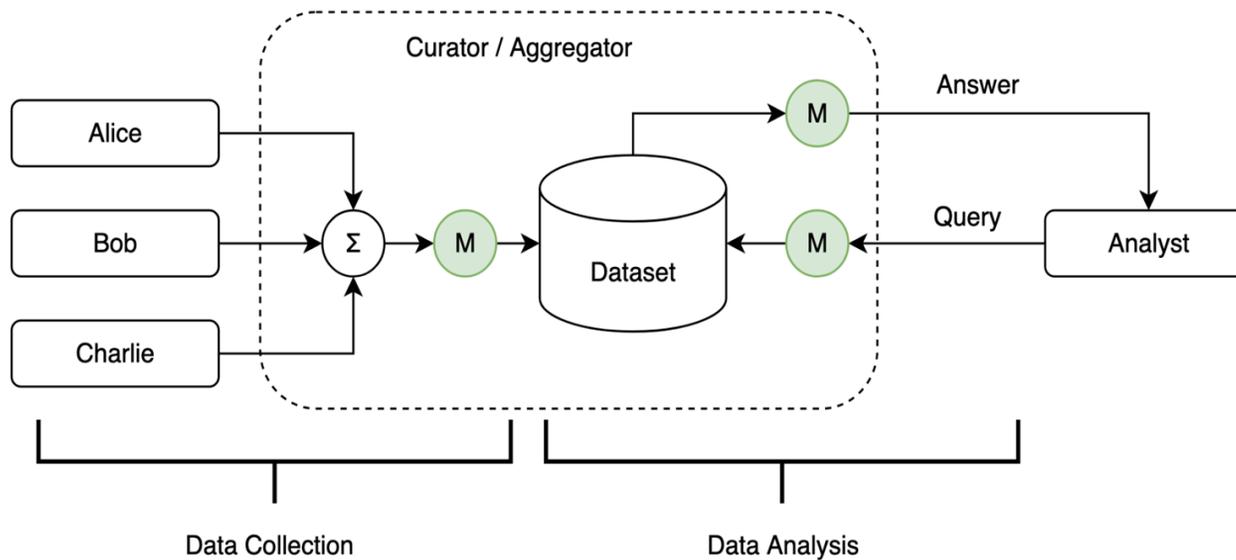
# Traditional data analysis

- Users provide data to data curator
- Data curator de-anonymizes and aggregates the data into a dataset
- Data curator makes the data available to data analysts
- Data analysts can query and retrieve information from the dataset



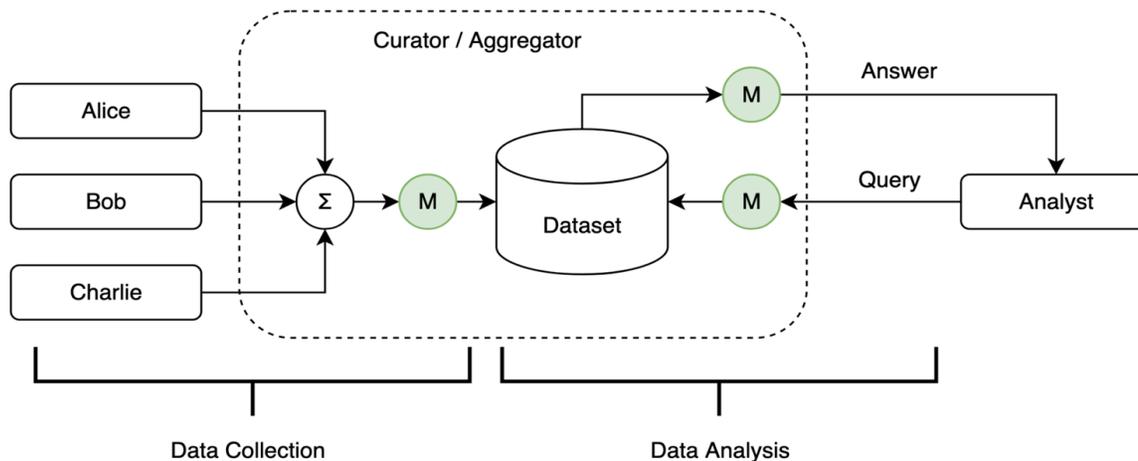
# Privacy Models

- A set of methods and techniques to preserve the privacy of data
- And their associated privacy evaluation metrics



# Privacy Models

- To protect personal data of users, a privacy model preprocesses the aggregated data before their insertion into the dataset
- To ensure the queries do not lead to loss of privacy, the queries are preprocessed before they are applied to the dataset
- To protect leakage of private data in answers to queries, the answers are preprocessed by the privacy model before submission to the analyst



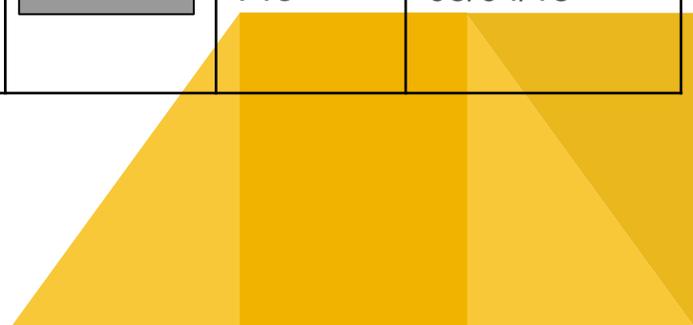
# Traditional Privacy Models

## Data Anonymization / De-identification

- Removing identifiable information from the dataset column-wise
- Methods consist of stripping, masking, swapping, perturbing, some columns in the dataset

User ID	Name	Address	Account Type	Subscription Date
001	Alice	123 A St	Pro	01/02/20
002	Bob	234 B St	Free	02/03/21
003	Charlie	456 C St	Pro	03/04/18

User ID	Name	Address	Account Type	Subscription Date
001	█	█	Pro	01/02/20
002	█	█	Free	02/03/21
003	█	█	Pro	03/04/18



# Traditional Privacy Models

## Data Anonymization / De-identification

### Privacy leakage attack through background information (Linking attack)

Dataset 1: Users database with personal information

User ID	Name	Address	Zip Code	Birth date	Gender
001	Alice	123 A St	83401	01/02/1997	F
002	Bob	234 B St	83402	02/03/1995	M
003	Charlie	456 C St	83403	03/04/1999	M

Dataset 2: Users medical database with personal information stripped (name and address)

User ID	Zip Code	Birth date	Gender	Probable disease ID
001	83401	01/02/1997	F	120
002	83402	02/03/1995	M	35
003	83403	03/04/1999	M	240

It is possible to uniquely identify a user with the triplet (Zip Code, Birth date, Gender)!



# Traditional Privacy Models

## Data Anonymization / De-identification

- Anonymized datasets have less useful information than non-anonymized datasets
- Use of anonymized dataset in business applications violates some regulations
- It is possible to infer personal information of users from de-anonymized dataset

Table 3

User ID	Name	Address	Account Type	Subscription Date
001	[REDACTED]	[REDACTED]	Pro	01/02/20
002	[REDACTED]	[REDACTED]	Free	02/03/21
003	[REDACTED]	[REDACTED]	Pro	03/04/18

Table 4

User ID	Product Name	Product Price	Purchase Date
001	TV	400	04/01/21
002	Iphone	1199	10/10/21
003	Watch	130	05/09/20

Tracking users with information from Table 1 (User ID) and Table 2 (User ID)



# Traditional Privacy Models

## k-Anonymity [Emam et al.]

- A dataset is said to be k-anonymous if, for any individual record in the dataset, there are at least  $k-1$  other records which are indistinguishable from it.
- K-anonymity ensures a linkage attack against a k-anonymous dataset can only identify a group of k records which are indistinguishable one from another
- However, k-anonymity only works for very large dataset with simple fields which makes it impracticable in real world scenarios.



# Traditional Privacy Models

## Randomized Response: A first attempt at Differential Privacy [Warner, 1965]

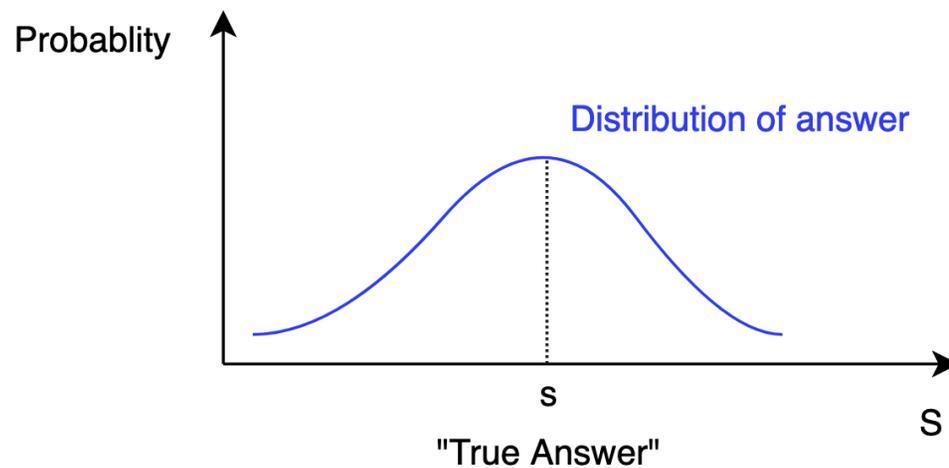
- **Problem:** Given a survey which contains a sensitive or embarrassing question (**query**), how to protect the privacy of participants' responses while performing statistics on the **answers**.
- **Approach:** Use coin flips
  - When participants are asked a sensitive question, they flip a coin before answering
  - If they get heads, they are asked to answer YES to the question regardless of their experience.
  - When they get tails, they are asked respond truthfully according to their experience.
- **Results:** Double the resulting statistics
  - At the end of the survey, half of the participant would have answer correctly, while the other half answer falsely.
  - To account for the half false answers they got, the organizers double the statistics at the end of the survey,
- **Alternative:** [Greenberg, 1969]
  - When participants get a tail, they are asked to flip a coin again to answer YES when it is heads, and NO when tails.



# Traditional Privacy Models

Randomized Response: A first attempt at Differential Privacy [Warner, 1965]

- The query answer is no longer a deterministic value, but a sampling from a distribution



# Differential Privacy Model

- A privacy model that ensures privacy of individual data with provable privacy guarantees.
- Differential privacy solves the problem of learning useful information about all the records while learning nothing about individual records
- It ensures that an adversary cannot infer an unknown data record from known data records

User ID	Name	Address	Account Type	Subscription Date
001	Alice	123 A St	Pro	01/02/20
[REDACTED]				
003	Charlie	456 C St	Pro	03/04/18

John knows two rows

John cannot infer middle row



# Differential Privacy

[Dwork et al.]

## Some Definitions

- **Dataset:** a collection of attributes and values (records)

Any subset  $D \subseteq \mathcal{X}^n$  of a set of n rec  $\mathcal{X} = \{x_1, x_2, \dots, x_n\}$

- **Neighboring datasets:**

D and D' are said to be neighboring datasets, denoted as  $D \sim D'$ , if they differ in one record.

- **Query:** A function to be applied to a dataset

A real-valued query is any function  $q: \mathcal{X}^n \rightarrow \mathbb{R}$ ,  $q \in F$ .

- **Sensitivity of a query (l1-sensitivity)  $\Delta q$ :**

$$\Delta q := \max_{x \sim x'} \|q(x) - q(x')\|_1$$

- **Answer:** The output of a query  $q(x)$  where x is a dataset



# Differential Privacy

[Dwork et al.]

## Some Definitions

- **Privacy Mechanism:** any algorithm to be applied to a dataset which ensure some privacy guarantees

A privacy mechanism  $M$  is an algorithm that has as input a dataset  $D \subseteq X^n$  and optionally a set of queries  $F$  and that outputs answers  $A \subseteq X^n$  (possible dataset rows) to queries [Dwork et al.]

$$\mathcal{M} : \mathcal{X}^n \times F \rightarrow \mathcal{X}^n, M(D) = q(x)$$

- **Randomized Privacy Mechanism:** any privacy mechanism obtained by coin flips (sampling from a distribution),  $M(D) = q(x) + (v \sim \text{Dist})$
- **Privacy Loss:** The privacy loss incurred by running a randomized privacy mechanism  $M$  on neighboring datasets  $D \sim D'$  is given by:

$$L(\mathcal{M}) = \ln\left(\frac{P[\mathcal{M}(D) \in \mathcal{S}]}{P[\mathcal{M}(D') \in \mathcal{S}]}\right)$$

In practice, we expect a low privacy loss bounded by the privacy budget.

- **Privacy Budget:** the privacy budget is the maximum privacy loss incurred by  $M$ .
- **Sensitivity  $\Delta q$  (local or global):** determines how much perturbation is required in the mechanism.



# Differential Privacy

[Dwork et al.]

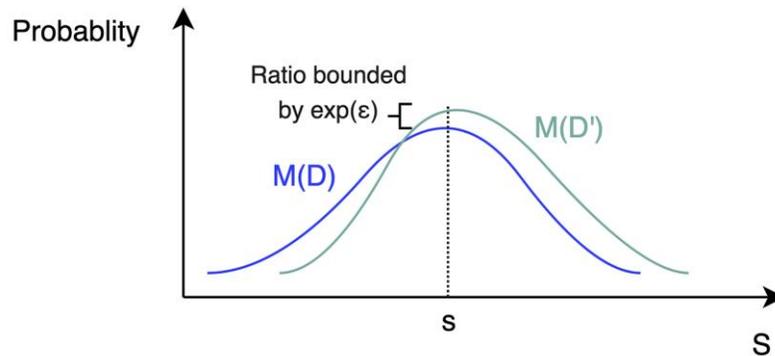
## Definition of Differential Privacy:

A randomized privacy mechanism  $M$  is a  $(\epsilon, \delta)$ -differential privacy mechanism if for every set of answers  $S = \{A_1, A_1, \dots, A_m\}$  and for any neighboring datasets  $D \sim D'$  the following inequality holds:

$$P[\mathcal{M}(D) \in S] \leq \exp(\epsilon) \cdot P[\mathcal{M}(D') \in S] + \delta$$

where the randomness results from coin flips of  $M$ .

For any pair of answers, the ratio of their probabilities distributions is bounded by  $\exp(\epsilon)$ , i.e., the **Privacy Budget** (lower represents stronger privacy).



# Differential Privacy

[Dwork et al.]

Two flavors of Differential Privacy:

1. For  $\delta=0$ ,  $\epsilon$ -differential privacy or **pure differential privacy**
2. For  $\delta > 0$ ,  $(\epsilon, \delta)$ -differential privacy or **approximate differential privacy**

Utility:  $(\lambda, \delta)$ -useful

- Given a randomized query  $q$  of a query  $f$  is  $(\lambda, \delta)$ -useful, if for any database  $X \in D$ , with probability at least  $1 - \delta$ ,  $|f(X) - q(X)| \leq \lambda$ , with  $\lambda > 0$  and  $0 < \delta < 1$ .
- The utility parameter  $\lambda$ , ensures that the answer to a randomized query is at  $\lambda$  far from the true answer from the deterministic version of the query.

Privacy/Utility trade off:

- $\epsilon=0$  (and  $\delta=0$ ): the mechanism achieves **absolute privacy, no utility (very low accuracy)**.
- $\epsilon > 0$  and  $\epsilon$  very large: the mechanism provides **no privacy**, but achieves **perfect utility (high accuracy)**.



# Differential Privacy

## Properties of Differential Privacy

### Properties of Differential Privacy

#### 1. Sequential Composition:

- If  $M_1(x)$  satisfies  $\epsilon_1$ -differential privacy
- And  $M_2(x)$  satisfies  $\epsilon_2$ -differential privacy
- Then the mechanism  $M(x)=(M_1(x),M_2(x))$  satisfies  $(\epsilon_1+\epsilon_2)$ -differential privacy

#### 2. Parallel Composition:

- If  $F(x)$  satisfies  $\epsilon$ -differential privacy
- And we split a dataset  $X$  into  $k$  disjoint sets such that  $x_1 \cup \dots \cup x_k = X$
- Then the mechanism which releases all of the results  $F(x_1), \dots, F(x_k)$  satisfies  $\epsilon$ -differential privacy

#### 3. Group Privacy:

- If we split a dataset  $X$  into  $k$  disjoint sets such that  $x_1 \cup \dots \cup x_k = X$
- Then, any  $(\epsilon, 0)$ -differentially private mechanism  $M$  is  $(k\epsilon, 0)$ -differentially private for groups of size  $k$  disjoint subsets of  $k$ .

#### 4. Post-Processing (robustness to auxiliary information):

- If  $F(X)$  satisfies  $\epsilon$ -differential privacy
- Then for any (deterministic or randomized) function  $g$ ,  $g(F(X))$  satisfies  $\epsilon$ -differential privacy



# Differential Privacy

[Dwork et al.]

**Differential Privacy via noise adding** : For a specific distribution **Dist**

Input: dataset **D**, query **q**

Steps:

1. Compute **q(x)**
2. Sample noise **v**  $\sim$  **Dist**( $\Delta q$ ,  $\epsilon$ )

Output: **q(x) + v**

**Popular Privacy Sampling Mechanisms:**

1. Laplace mechanism
2. Gaussian mechanism
3. Exponential mechanism



# Differential Privacy Mechanisms

Randomized Response: A first attempt at Differential Privacy [Warner, 1965]

- Randomized response is  $(\ln(3), 0)$ -Differential Private [Greenberg, 1969]:
  - $P[\text{Response} = \text{Yes} | \text{Truth} = \text{Yes}] / P[\text{Response} = \text{Yes} | \text{Truth} = \text{No}] = (\frac{3}{4}) / (\frac{1}{4}) = 3$
- The Privacy loss is fixed  $\epsilon = \ln(3)$
- A more general approach to achieve differential privacy from a deterministic mechanism is to add noise, sampled from a known distribution, to the output of the deterministic mechanism.

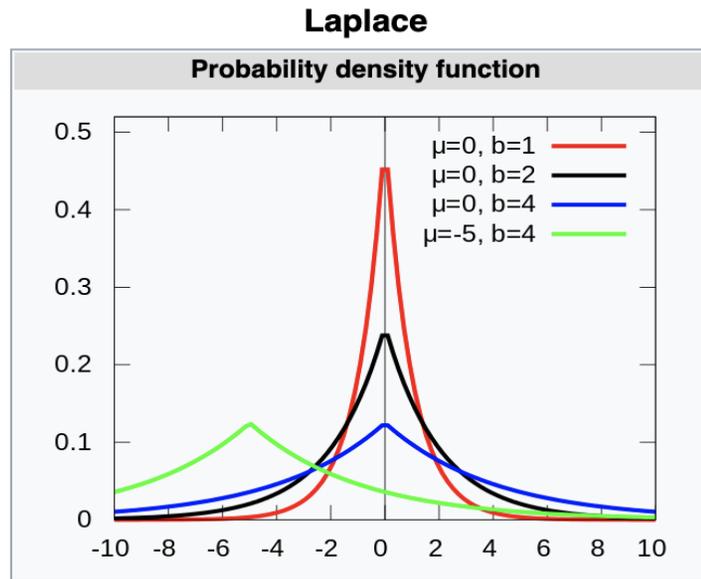


# Differential Privacy Mechanisms

## Laplace Mechanism: Achieving differential privacy by adding Laplace noise

- Privacy mechanism which outputs  $\mathbf{q}(\mathbf{x}) + (\mathbf{v} \sim \text{Lap}(\Delta\mathbf{q}, \epsilon))$
- Laplace distribution is centered at 0 with a std of  $b\sqrt{2}$ .
- Its probability density function is given by:

$$f(x | \mu, b) = \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right)$$



Source: Wikipedia

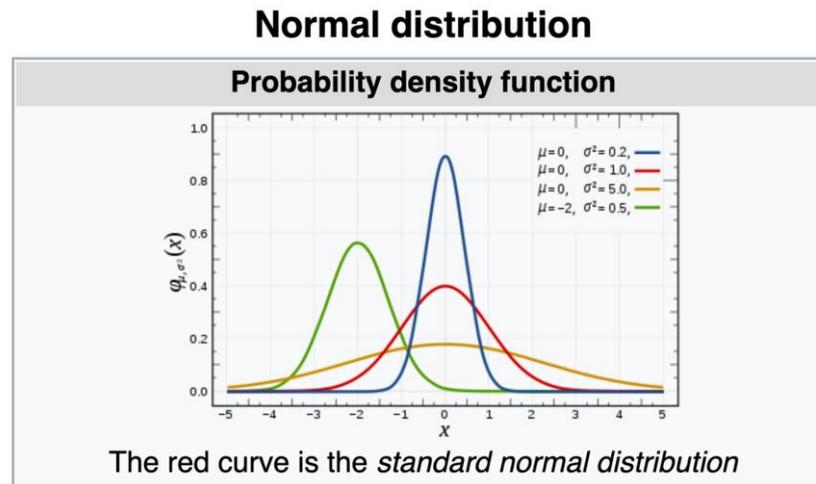


# Differential Privacy Mechanisms

Gaussian Mechanism: Achieving differential privacy by adding Gaussian noise

- Privacy mechanism which outputs  $\mathbf{q}(\mathbf{x}) + (\mathbf{v} \sim \mathbf{N}(\Delta\mathbf{q}, \epsilon))$
- The Gaussian or Normal Distribution has a probability density function given by:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$



Source: Wikipedia



# Differential Privacy Mechanisms

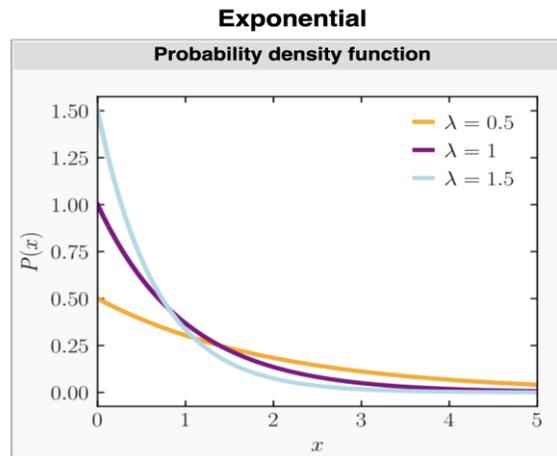
## Exponential Mechanism: Achieving differential privacy by selecting lowest sensitivity score

- Given a sensitivity score function  $\mathbf{H}$ , selects answer  $\mathbf{a}$  with the lowest sensitivity score such that  $P(\mathbf{a} \in A \text{ is selected}) \propto e^{\epsilon H(D, \mathbf{a}) / 2s(H, \|\cdot\|)}$

$$s(H, \|\cdot\|) = \max_{d(D, D')=1, a \in A} \|H(D, a) - H(D', a)\|$$

- The Exponential distribution is parameterized by with a probability density function given by:

$$f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0, \\ 0 & x < 0. \end{cases}$$



Source: Wikipedia



# Differential Privacy Mechanisms

## Differential Privacy Mechanisms Summary

### Randomized Response:

- Pros:
  - Easy to implement
- Cons:
  - Fixed privacy budget
  - Not flexible

### Laplace mechanism: ( $\epsilon$ )-differential privacy (pure differential privacy)

- Pros:
  - More accurate query answers
  - The most general mechanism
- Cons:
  - Noise can take extreme values leading
  - Requires large privacy budget to work in practice which produce less accurate answers.
  - Can only be used for for numeric queries only



# Differential Privacy Mechanisms

## Differential Privacy Mechanisms Summary

**Gaussian mechanism:**  $(\epsilon, \delta)$ -differential privacy (approximate differential privacy)

- Pros:
  - Add less noise than Laplace mechanism
  - better suited for multivariate problems.
- Cons:
  - requires the use of the the relaxed  $(\epsilon, \delta)$ -differential privacy definition
  - less accurate than the Laplace mechanism

**Exponential mechanism:**

- Pros:
  - The Laplace mechanism can be derived from the exponential mechanism
  - The exponential mechanism is extremely general
  - Can provide answers to non numeric queries
- Cons:
  - It is very difficult to implement



# Differential Privacy Mechanisms

## Differential Privacy and KL-Divergence [Ilya, 2017]

Kullback-Leibler Divergence

$$D_{KL}(P||Q) = \mathbb{E}_{x \sim P}[\log \frac{P(x)}{Q(x)}]$$

Relaxed (pure) differential privacy

$$\frac{P[\mathcal{M}(D) \in \mathcal{S}]}{P[\mathcal{M}(D') \in \mathcal{S}]} \leq e^\epsilon.$$

**Renyi Divergence:** For two probability distributions  $P$  and  $Q$  defined over  $R$ , the Renyi divergence of order  $\alpha > 1$  is

$$D_\alpha(P||Q) \triangleq \frac{1}{\alpha - 1} \log \mathbb{E}_{x \sim Q} \left( \frac{P(x)}{Q(x)} \right)^\alpha$$

- Renyi divergence is KL-divergence for  $\alpha$  approaching 1:  $D_1(P||Q) = \lim_{\alpha \rightarrow 1} D_\alpha(P||Q)$

$$D_1(P||Q) = \mathbb{E}_{x \sim P} \log \frac{P(x)}{Q(x)}.$$

**$(\alpha, \epsilon)$ -Renyi Differential Privacy:**

- A randomized mechanism  $f : D \rightarrow R$  is said to have  $(\epsilon)$ -Renyi differential privacy of order  $\alpha$ , if for any neighboring dataset  $D$  and  $D'$ :

$$D_\alpha (f(D)||f(D')) \leq \epsilon.$$



# Applications in Machine Learning

## Issues with Deep Learning without privacy mechanisms

- Neural network models memorize training examples [Carlini et al., 2019]
- Given a public model, it is possible to transform (blurred) image through Generative Model-Inversion (GMI) attacks [Zhang et al.]

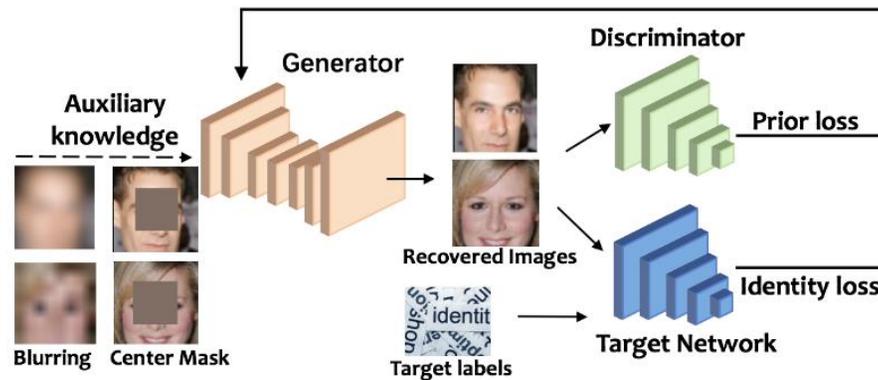


Figure 1: Overview of the proposed GMI attack method.



# Applications in Machine Learning

## Differentially Private Deep Learning

### Differential Privacy in Gradient Descent

- Update rule [Xie et al.]

$$w_{t+1} = w_t - \eta_t (\nabla f(w_t; D_{m(t)}) + N_t(0, \lambda))$$

- Due to sequential composition, the privacy loss is unbounded when we perform many iterations of the SGD algorithm
- More iterations lead to a larger privacy cost. However, in practical Deep Learning, more iterations generally result in a better model.
- But, in differentially private SGD version, more iterations can make the model worse, i.e., the noise increases with each iteration.
- Tradeoff: Find the right balance between the number of iterations and the scale of the noise added.
- Some Techniques: **Gradient clipping** before adding the noise [Abadi et al.]



# Applications in Machine Learning

## Differentially Private Deep Learning with Opacus [FAIR, Opacus]

- Implements a  $(\epsilon)$ -Renyi differential privacy of order  $\alpha$  in Pytorch

```
model = Net()
optimizer = SGD(model.parameters(), lr=0.05)
privacy_engine = PrivacyEngine(
    model,
    sample_rate=0.01,
    alphas=[1, 10, 100],
    noise_multiplier=1.3,
    max_grad_norm=1.0,
)
privacy_engine.attach(optimizer)
# Now it's business as usual
```

Source: <https://opacus.ai/>



# Conclusion

- Traditional data analysis is done without any privacy mechanism which leads to privacy violation
- Traditional privacy models are vulnerable to linking attacks
- Randomized Response was the first attempt to offer differential privacy.
- The fixed privacy loss in Randomized Response makes it unsuitable for some application
- Methods based on sampling from a distribution (Laplace, Gaussian, Exponential, etc) offer better privacy guarantees.
- Differential Privacy can be applied to ML to achieve Privacy Preserving ML.



# References

1. [Dwork et al.] Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Found. Trends Theor. Comput. Sci.*, 9(3-4):211–407, 2014.
2. [Warner, 1965] Warner, S. L. (March 1965). "Randomised response: a survey technique for eliminating evasive answer bias". *Journal of the American Statistical Association*. Taylor & Francis. 60 (309): 63–69. doi:10.1080/01621459.1965.10480775. JSTOR 2283137. PMID 12261830.
3. [Zhu et al.] Tianqing Zhu, Gang Li, Wanlei Zhou, and S Yu Philip. *Differential privacy and applications*. Springer, 2017.
4. [Near et al.] Joseph P. Near and Chik'e Abuah. *Programming Differential Privacy*, volume 1. 2021.
5. [Xie et al.] Yun Xie, Peng Li, Chao Wu, and Qiuling Wu. Differential privacy stochastic gradient descent with adaptive privacy budget allocation. In *2021 IEEE International Conference on Consumer Electronics and Computer Engineering(ICCECE)*, pages 227–231, 2021
6. [Abdai et al.] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pages 308–318, 2016
7. [Ilya, 2017] Mironov, Ilya. "Rényi differential privacy." *2017 IEEE 30th Computer Security Foundations Symposium (CSF)*. IEEE, 2017.
8. [FAIR, Opacus] <https://opacus.ai/>
9. [Emam et al.] El Emam, Khaled, and Fida Kamal Dankar. "Protecting privacy using k-anonymity." *Journal of the American Medical Informatics Association* 15.5 (2008): 627-637.
10. [Carlini et al., 2019] Carlini, Nicholas, et al. "The secret sharer: Evaluating and testing unintended memorization in neural networks." *28th {USENIX} Security Symposium ({USENIX} Security 19)*. 2019.
11. [Zhang et al.] Zhang, Yuheng, et al. "The secret revealer: Generative model-inversion attacks against deep neural networks." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2020.



Thank you  
Questions?



University of Idaho

# Differentially Private SGD

---

## *Differentially Private SGD*

- [Abadi \(2016\) Deep Learning with Differential Privacy](#)
- This work introduced differential privacy (DP) for training ML models for protecting the privacy of the training data
  - *Differential privacy (DP)* is applied to *Stochastic Gradient Descent (SGD)* during model training
  - DP-SGD clips the gradients and adds Gaussian noise to the gradients with respect to the model parameters
  - This approach controls the amount of information from the training data that is memorized by the model during training
  - The goal is to produce ML models which provide approximately the same privacy when an individual input instance is removed from the training dataset
- The paper also introduces a method for calculating the privacy loss, called *moments accountant*

# DP Example

## Differentially Private SGD

- Consider two databases  $D_1$  and  $D_2$  that show if a person has diabetes or not
  - The only difference between the two databases is that  $D_2$  does not include the last record in  $D_1$  (for Bob)
- Let's assume that the databases are publicly available for making queries
  - To protect patient identities, it is not allowed to query the patient names
- However, an adversary can query the sum of the persons with diabetes in the first database (e.g.,  $f(D_1) = 64$ ), and the sum in the second database (e.g.,  $f(D_2) = 63$ )
  - Based on the difference  $f(D_1) - f(D_2) = 64 - 63 = 1$ , the adversary can infer that Bob has diabetes
  - Alternatively, if  $f(D_1) = 63$  and  $f(D_2) = 63$ , the adversary can infer that Bob does not have diabetes

$D_1$  (includes Bob)

Name	Has Diabetes
Don	1
Monica	0
...	
...	
Chris	1
Bob	1

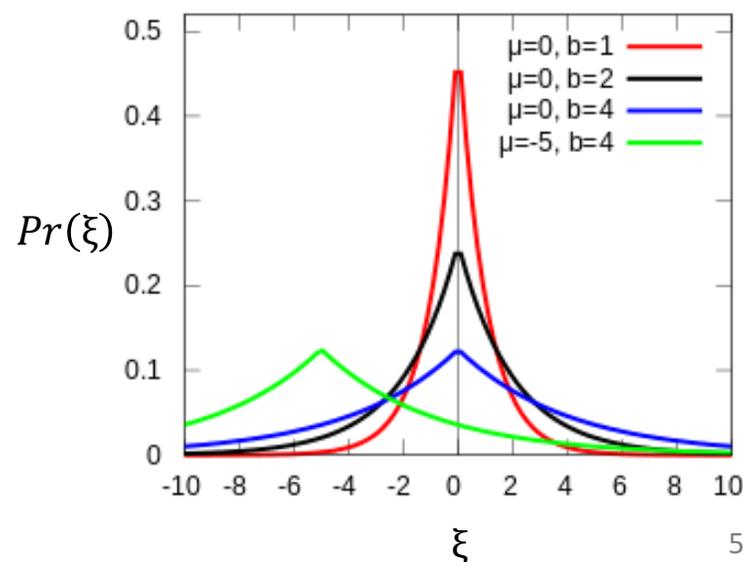
$D_2$  (without Bob)

Name	Has Diabetes
Don	1
Monica	0
...	
...	
Chris	1

# DP Example (cont'd)

## Differentially Private SGD

- An algorithm that is *differentially private* adds noise to the answers for  $f(D_1)$  and  $f(D_2)$  to make it difficult to infer the information about Bob
  - I.e., a **randomization mechanism**  $\mathcal{M}(D)$  is selected to add noise  $\xi$  to the output answers to queries  $f(D)$ , that is,  $\mathcal{M}(D) = f(D) + \xi$
- Additive noise  $\xi$  from a **Laplacian distribution** (shown) is commonly applied
  - E.g., let's assume a **privacy budget**  $\epsilon = 0.5$  and let's sample noise from a Laplacian distribution with  $\mu = 0$  and scale  $b = 1/\epsilon = 1/0.5 = 2$
  - 6 random noise samples are:  $\xi \in \{-0.13, 2.06, -1.67, -2.49, -0.52, 0.37\}$
  - Consider 3 queries by the adversary having the outputs  $f(D_1) = 64$  and  $f(D_2) = 63$  with added Laplacian noise  $\xi$ :
    - $\mathcal{M}(D_1) - \mathcal{M}(D_2) = 63.87 - 65.06 = -1.19$
    - $\mathcal{M}(D_1) - \mathcal{M}(D_2) = 62.33 - 60.51 = 1.82$
    - $\mathcal{M}(D_1) - \mathcal{M}(D_2) = 63.48 - 63.37 = 0.11$
  - Based on the differences between the randomized outputs from the queries for  $D_1$  and  $D_2$ , now it is impossible for the adversary to tell if Bob has diabetes



# DP Mechanism

---

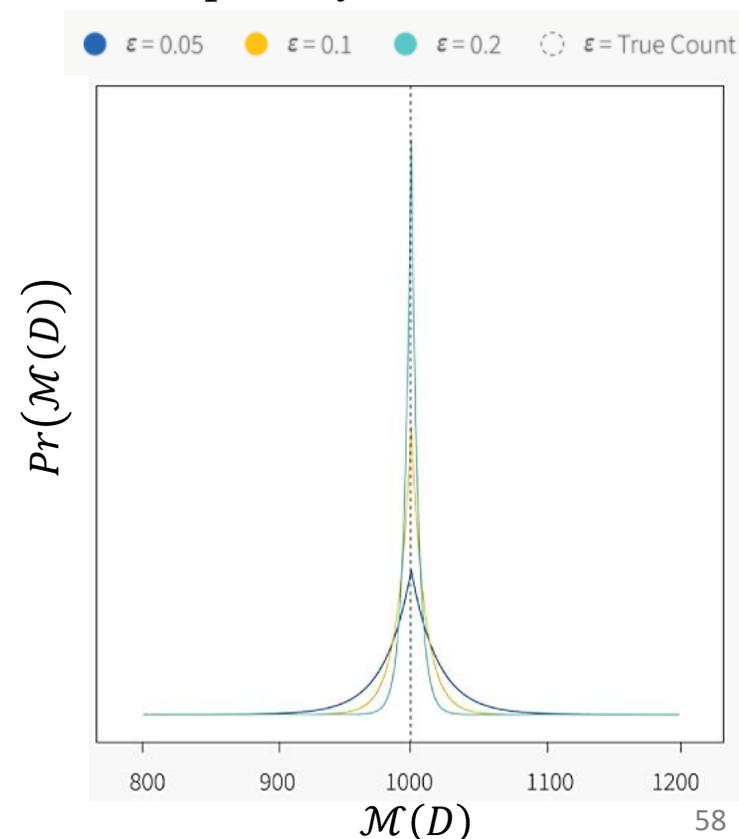
## Differentially Private SGD

- The important question in DP is: **how much noise to add?**
  - The amount of noise  $\xi$  depends on the data, and it needs to be adjusted
    - E.g., a function  $f_1(D)$  that provides the yearly income of people in thousands of dollars would require different level of noise than a function  $f_2(D)$  that provides the height in feet
- The **sensitivity** of the function  $f$  determines how much the output  $f(D)$  changes by adding a single data instance
  - Sensitivity is defined as  $\Delta f = \max \|f(D_1) - f(D_2)\|_1$  for all possible datasets  $D_1$  and  $D_2$  differing in one data instance, where  $\|\cdot\|_1$  denotes  $\ell_1$ -norm
    - E.g., for the example with medical diabetes records, the sensitivity is  $\Delta f = 1$ , since the sum of the people with diabetes can change only by 1 when a single input is added
- A Laplacian mechanism that is  **$\epsilon$ -differentially private** adds a Laplacian noise with scale  $b = \Delta f / \epsilon$
- Note that **if the privacy budget  $\epsilon$  has smaller values**, this will result in larger amount of Laplacian noise  $\xi$  added to  $f(D)$ 
  - Thus, the noisy outputs  $\mathcal{M}(D)$  will reveal less private information about the inputs (i.e., provide better privacy protection), but also the noisy answers to the queries  $\mathcal{M}(D)$  will be less accurate

# DP with Laplacian Randomization

## Differentially Private SGD

- The figure shows the **probability distributions** of the outputs  $\mathcal{M}(D)$  for three different levels of Laplacian noise with  $\epsilon \in \{0.05, 0.1, 0.2\}$ 
  - The true output value is  $f(D) = 1,000$
  - Larger values of  $\epsilon$  have distributions that are tighter around the true value of  $f(D) = 1,000$  in the figure, and hence are more accurate, but leak more privacy
- A mechanism  $\mathcal{M}(D)$  is  **$\epsilon$ -differentially private** if for all databases  $D_1$  and  $D_2$  that differ by at most one instance, and for any subset of outputs  $S$ :
 
$$\Pr(\mathcal{M}(D_1) \in S) \leq e^\epsilon \Pr(\mathcal{M}(D_2) \in S)$$
  - In other words,  $\epsilon$ -differential privacy ensures that the probabilities of any two outputs  $\mathcal{M}(D_1)$  and  $\mathcal{M}(D_2)$  differ by at most  $e^\epsilon$
  - E.g., for  $\epsilon = 0.05$ ,  $\Pr(\mathcal{M}(D_1))/\Pr(\mathcal{M}(D_2))$  is at most  $e^{0.05} = 1.05$
  - Smaller  $\epsilon$  ensures more similar outputs  $\mathcal{M}(D_1)$  and  $\mathcal{M}(D_2)$ , and provides higher levels of privacy



# DP with Gaussian Randomization

---

## Differentially Private SGD

- There are other DP mechanisms besides the Laplacian mechanism, that are more suitable for some applications
- The *Gaussian mechanism* adds Gaussian noise instead of Laplacian noise, and the level of noise is based on the  $\ell_2$ -norm sensitivity, instead of  $\ell_1$ -norm
- A Gaussian mechanism is  $(\epsilon, \delta)$ -differentially private if for all databases  $D_1$  and  $D_2$  that differ by at most one instance, and for any subset of outputs  $S$ :

$$\Pr(\mathcal{M}(D_1) \in S) \leq e^\epsilon \Pr(\mathcal{M}(D_2) \in S) + \delta$$

- The  $(\epsilon, \delta)$ -differential privacy that is provided by the Gaussian mechanism introduces the **probability parameter  $\delta$** 
  - Informally,  $(\epsilon, \delta)$ -differential privacy is guaranteed with probability  $1 - \delta$
  - E.g., for  $\delta = 0.05$ , the method is  $\epsilon$ -differentially private with 95% probability
- The Gaussian mechanism is therefore weaker than the Laplacian mechanism, since it allows scenarios when the privacy cannot be guaranteed
  - On the other hand, additive Gaussian noise is less likely to take on extreme values than Laplacian noise

# Privacy in Machine Learning

---

## *Differentially Private SGD*

- Training ML models can be considered an extension of the previous example on querying databases
  - I.e., ML models use data to learn a function, which is afterward used for prediction
- The datasets for training ML models often contain sensitive information (e.g., medical records, personal information), so it is important to provide privacy guarantees
  - On the other hand, we know that ML models can memorize the training data, which can be exploited by adversaries to recover information about the data from a trained model
- The challenge is: how to extract enough information from data to train accurate ML models without revealing the data
- DP-SGD is an approach that adds noise to the model parameters during training, to reduce the memorization of input samples

# Differentially Private SGD

## Differentially Private SGD

- Differentially Private Stochastic Gradient Descent (DP-SGD)
  - DP-SGD adds two additional steps to SGD: Clip gradient, and Add noise

---

### Algorithm 1 Differentially private SGD (Outline)

---

**Input:** Examples  $\{x_1, \dots, x_N\}$ , loss function  $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$ . Parameters: learning rate  $\eta_t$ , noise scale  $\sigma$ , group size  $L$ , gradient norm bound  $C$ .

**Initialize**  $\theta_0$  randomly

**for**  $t \in [T]$  **do**

  Take a random sample  $L_t$  with sampling probability  $L/N$

**Compute gradient**

  For each  $i \in L_t$ , compute  $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

**Clip gradient**

$\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C})$

**Add noise**

$\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} (\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$

**Descent**

$\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

**Output**  $\theta_T$  and compute the overall privacy cost  $(\epsilon, \delta)$  using a privacy accounting method.

---

Introduced DP steps

# Gradient Clipping

## Differentially Private SGD

- ML models tend to memorize more information about some input samples than others
  - Input samples that produce large gradients can be memorized by the model, and violate privacy
- This approach proposes to clip the  $\ell_2$  norm of the gradient to a **threshold  $C$** , in order to limit the influence by the individual input samples
  - Also, since the values of the gradients cannot be estimated ahead of time, the clipping operation controls the sensitivity of the DP randomization mechanism
- If the gradient at step  $t$  by an input sample  $x_i$  is  $\mathbf{g}_t(x_i) = \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$ , the clipped gradient  $\bar{\mathbf{g}}_t(x_i)$  is:

$$\bar{\mathbf{g}}_t(x_i) = \begin{cases} \mathbf{g}_t(x_i) & \text{if } \|\mathbf{g}_t(x_i)\|_2 \leq C \\ \frac{\mathbf{g}_t(x_i)}{\|\mathbf{g}_t(x_i)\|_2 / C} & \text{if } \|\mathbf{g}_t(x_i)\|_2 > C \end{cases}$$

- That is, if the norm  $\|\mathbf{g}_t(x_i)\|_2$  is greater than  $C$ , the gradient is scaled down to have a norm equal to  $C$

# Adding Noise

---

## Differentially Private SGD

- GP-SGD approach employs a Gaussian randomization mechanism
- Gaussian noise is added to the gradients at each training step  $t$ , according to:

$$\tilde{\mathbf{g}}_t = \frac{1}{L} \left( \sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}) \right)$$

- At each step  $t$ , the average of the clipped gradient for a batch of inputs (with a batch size  $L$ ) is first calculated as  $\frac{1}{L} \sum_i \bar{\mathbf{g}}_t(x_i)$
- Gaussian noise  $\mathcal{N}(0, \sigma^2 C^2 \mathbf{I})$  with mean 0 and diagonal co-variance  $\sigma^2 C^2$  is afterward added to the batch-averaged gradient
  - Note that the co-variance is a function of the clipping threshold  $C$
  - E.g., larger value of  $C$  does less clipping, but requires more noise to achieve the same level of privacy

# Moments Accountant

---

## *Differentially Private SGD*

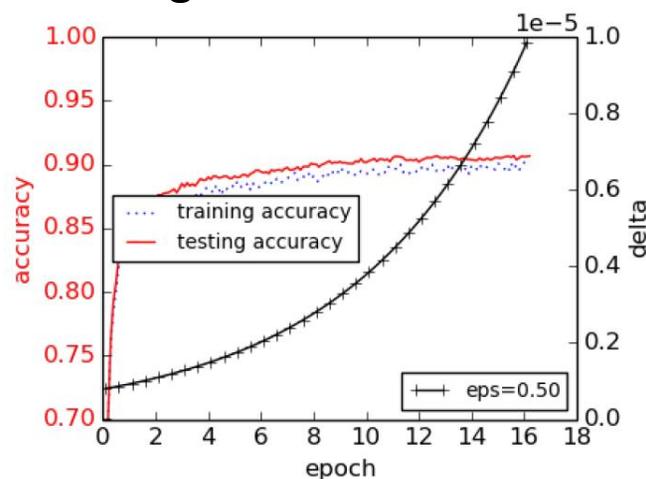
- The **composition property in DP** states that if the privacy budget for one interaction with the data is  $\epsilon_1$  and for another interaction with the data is  $\epsilon_2$ , the combined privacy budget is  $\epsilon_1 + \epsilon_2$ 
  - Therefore, by accumulating the privacy loss for each mini-batch when training an ML model, it is possible to calculate the overall privacy loss during training
- **Moments accountant** is an introduced approach in the paper that evaluates the **privacy budget** of a model training with DP-SGD
  - The privacy loss is estimated at each training step, and it is used to calculate the cumulative privacy loss over all training epochs
  - Note that increasing the number of training epochs increases the privacy loss
    - E.g., training a model for 100 epochs that achieved a privacy loss of  $\epsilon = 1.26$ , when training for 400 epochs the privacy loss increased to  $\epsilon = 2.55$
- Moments accountant employs the moments of mixtures of Gaussian distributions to calculate the upper bound of the cumulative privacy loss
  - The approach is described in more detail in the paper

# Experimental Evaluation

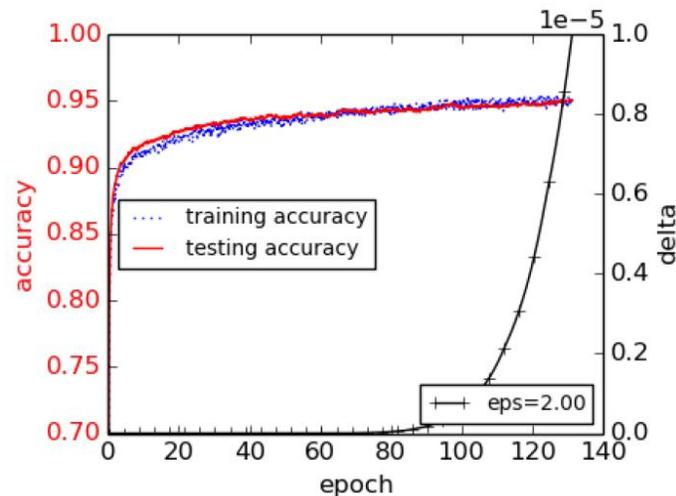
## Differentially Private SGD

- MNIST dataset
  - Training and testing accuracies for different levels of noise  $\epsilon \in \{0.5, 2, 8\}$ 
    - The corresponding Gaussian noise variances are  $\sigma \in \{8, 4, 2\}$ , the clipping threshold is  $C = 4$
  - The upper bound for the privacy probability parameter is set to  $\delta = 10^{-5}$ 
    - Thus, the corresponding  $(\epsilon, \delta)$ -differential privacies are  $(0.5, 10^{-5})$ ,  $(2, 10^{-5})$ ,  $(8, 10^{-5})$
  - The obtained test set accuracies are 90%, 95%, and 97%, respectively
  - Larger noise achieves lower test accuracy, but provides increased privacy protection

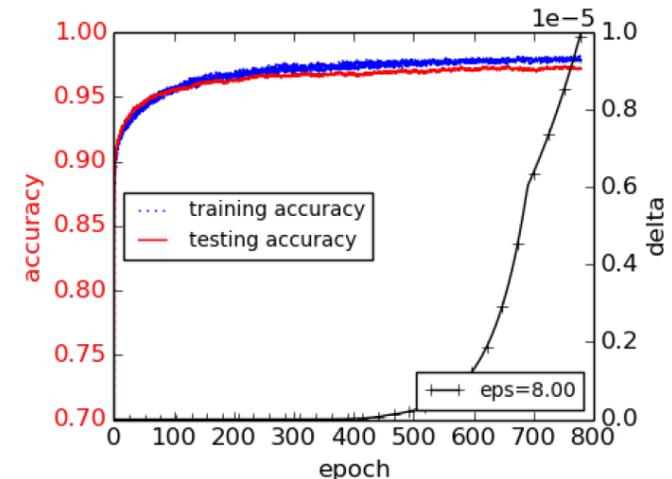
Large noise  $\epsilon = 0.5$



Medium noise  $\epsilon = 2$



Small noise  $\epsilon = 8$

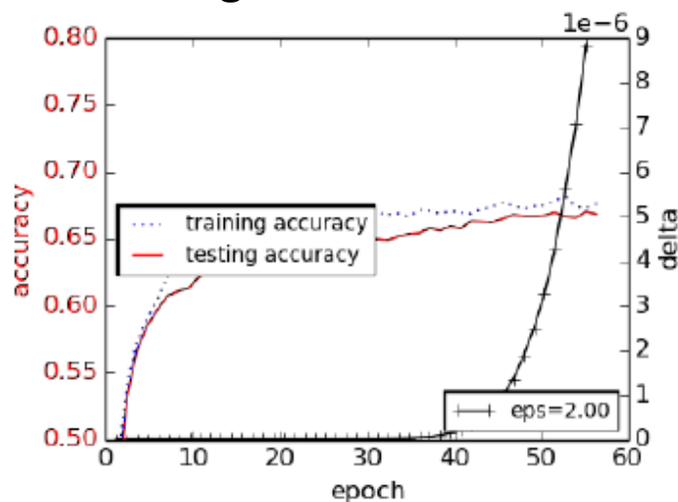


# Experimental Evaluation

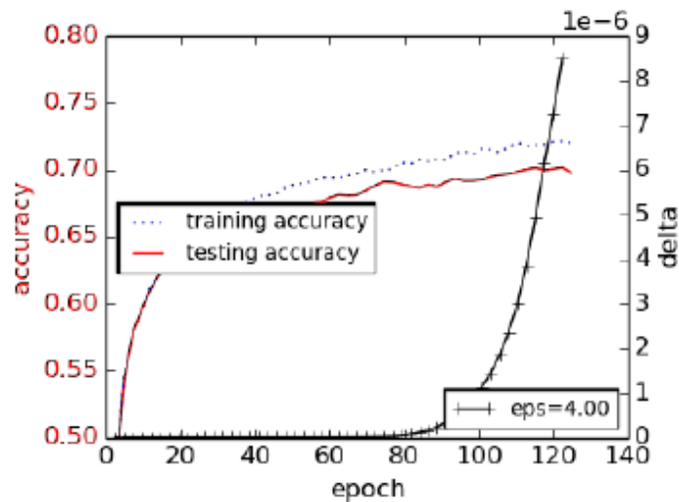
## Differentially Private SGD

- CIFAR-10 dataset
  - The results are similar to the obtained performance for MNIST
  - Training and testing accuracies for different levels of noise  $\epsilon \in \{2, 4, 8\}$ 
    - The target probability parameter is again set to  $\delta = 10^{-5}$
    - The corresponding  $(\epsilon, \delta)$ -differential privacies are  $(2, 10^{-5})$ ,  $(4, 10^{-5})$ ,  $(8, 10^{-5})$
  - The Gaussian noise variance is fixed to  $\sigma = 6$  for all experiments, the clipping threshold is  $C = 3$
  - The achieved test set accuracies are 67%, 70%, and 73%, respectively

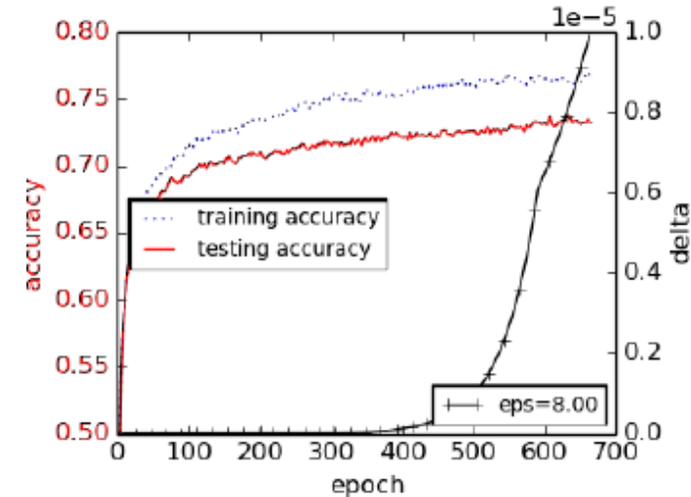
Large noise  $\epsilon = 2$



Medium noise  $\epsilon = 4$



Small noise  $\epsilon = 8$



# Privacy versus Accuracy Trade-Off

---

## *Differentially Private SGD*

- Perfect privacy in ML models is not possible
  - Adding too much noise to the model parameters would diminish the accuracy and the usefulness of the model
  - There is a trade-off between privacy protection and accuracy
- DP-SGD achieves privacy protection for deep NNs with a small decrease in the model accuracy and small increase in the training complexity
  - The approach adds Gaussian noise to the gradients in SGD to reduce the possibility for memorization of individual input instances by the model
  - This work also developed the moments accountant approach to calculate the cumulative privacy loss for the combination of the model and dataset



**University of Idaho**

College of Engineering

# **Scalable Private Learning with PATE**

**Written by: Nicolas Papernot, Shuang Song,  
Ilya Mironov, Ananth Raghunathan, Kuna  
Talwar, and Úlfar Erlingsson**

PowerPoint created by: Jacob Friedberg



# OUTLINE:

- I Background
- I Problem Statement
- I Proposed Solution
- I Design
- I Experimental Setup
- I Results
- I Conclusion
- I Questions

# BACKGROUND

## I Attacks on Privacy

- Data collection is a ~40-billion-dollar industry[1].
- Some of the data is personally identifiable information(PII)
- It takes only 3 pieces of information to identify 87% of the United States[2]
  - Zip code
  - Gender
  - Date of birth

## I What's the big deal?

- PII can be used against you(identity theft, social engineering and linkage attacks)
  - Linkage attacks identify you by cross-referencing data someone knows about you with other sources.
  - Also know as re identification attacks(from anonymized to identified)

# BACKGROUND

## I What can we do?

- Not all data collection is bad. Some collection, especially health data can be used to provide patterns for localizing sickness.
- Companies are going to collect data based on our usage of their services. To make better products
- Implement a way collect data without revealing PII.

## I Differential Privacy

- System for collecting and computing on data while maintaining privacy.
- This solves the earlier issue.
  - Allows companies to compute on data and get the same result as if the data was not anonymized.

# BACKGROUND

I How is this accomplished?

- Injecting noise into the dataset to create plausible deniability.
  - The noise shouldn't change the outcome of the computation.
- An observer should not be able to determine any PII from the output or identify whose data was trained on.

$\epsilon$  = Upper bound for the loss of privacy  
 $\delta$  = probability that privacy will not be held  
M = Model training algorithm

I Equation to know

[3]

**Definition 1.** *A randomized mechanism  $\mathcal{M}$  with domain  $\mathcal{D}$  and range  $\mathcal{R}$  satisfies  $(\epsilon, \delta)$ -differential privacy if for any two adjacent inputs  $D, D' \in \mathcal{D}$  and for any subset of outputs  $S \subseteq \mathcal{R}$  it holds that:*

$$\Pr[\mathcal{M}(D) \in S] \leq e^\epsilon \cdot \Pr[\mathcal{M}(D') \in S] + \delta. \quad (1)$$

# PROBLEM STATEMENT

- I Current approaches show potential but are untested at scale.
  - Private Aggregation of Teacher Ensembles(PATE)
    - Prior implementation of differential privacy for ML models
    - Uses Teachers to train students to answer queries.
      - Partitions data to be trained onto each teacher
      - All teachers aggregate answers to a single “aggregate teacher”
      - Students are trained on teacher responses, and queried by users
  - Issues with
    - Scalability
    - Robustness
    - Utility

# PROPOSED SOLUTION

## I Build off the PATE method

- New methods for aggregating teacher/student answers
  - Confidence Aggregator
    - Teacher consensus module where Min of T teachers guarantee a correct classification and throw out queries where teachers don't know
  - Interactive Aggregator
    - Student confidence scoring. Don't ask teachers for an answer if confidence is high that the student knows
  - Expensive Queries are high cost to privacy
- Gaussian noise instead of Laplacian
  - Less computationally expensive
  - Causes less noise overall.

# DESIGN

PATE design:

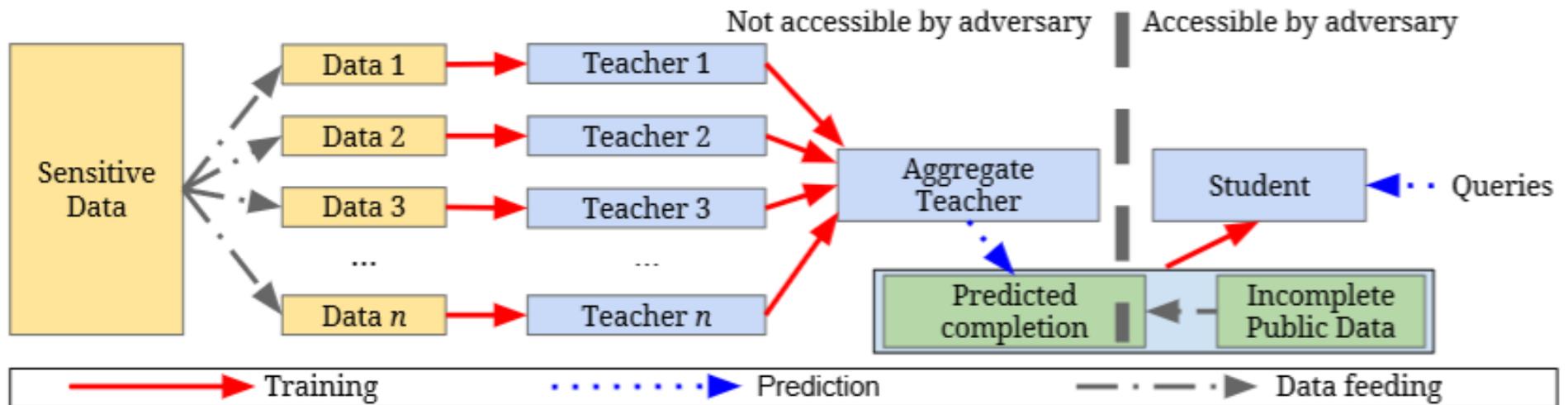


Figure 2: Overview of the approach: (1) an ensemble of teachers is trained on disjoint subsets of the sensitive data, (2) a student model is trained on public data labeled using the ensemble.

# DESIGN

Confident Aggregator:

---

**Algorithm 1 – Confident-GNMax Aggregator:** given a query, consensus among teachers is first estimated in a privacy-preserving way to then only reveal confident teacher predictions.

---

**Input:** input  $x$ , threshold  $T$ , noise parameters  $\sigma_1$  and  $\sigma_2$

- 1: **if**  $\max_i \{n_j(x)\} + \mathcal{N}(0, \sigma_1^2) \geq T$  **then** ▷ Privately check for consensus
  - 2:     **return**  $\operatorname{argmax}_j \{n_j(x) + \mathcal{N}(0, \sigma_2^2)\}$  ▷ Run the usual max-of-Gaussian
  - 3: **else**
  - 4:     **return**  $\perp$
  - 5: **end if**
-

# DESIGN

Interactive Aggregator:

---

**Algorithm 2 – Interactive-GNMax Aggregator:** the protocol first compares student predictions to the teacher votes in a privacy-preserving way to then either (a) reinforce the student prediction for the given query or (b) provide the student with a new label predicted by the teachers.

---

**Input:** input  $x$ , confidence  $\gamma$ , threshold  $T$ , noise parameters  $\sigma_1$  and  $\sigma_2$ , total number of teachers  $M$

- 1: Ask the student to provide prediction scores  $\mathbf{p}(x)$
- 2: **if**  $\max_j \{n_j(x) - Mp_j(x)\} + \mathcal{N}(0, \sigma_1^2) \geq T$  **then** ▷ Student does not agree with teachers
- 3:     **return**  $\operatorname{argmax}_j \{n_j(x) + \mathcal{N}(0, \sigma_2^2)\}$  ▷ Teachers provide new label
- 4: **else if**  $\max\{p_i(x)\} > \gamma$  **then** ▷ Student agrees with teachers and is confident
- 5:     **return**  $\operatorname{argmax}_j p_j(x)$  ▷ Reinforce student's prediction
- 6: **else**
- 7:     **return**  $\perp$  ▷ No output given for this label
- 8: **end if**

---

# EXPERIMENTAL SETUP

## I Datasets

- MNIST
- Street View House Numbers (SVHN)
- US Census Income Adult (UCI Adult)
  - Table of attributes about a person. Mainly used to test privacy.
- Glyph
  - Synthetically generated computer font symbols with at most 150 different classes



Figure 3: **Some example inputs from the Glyph dataset along with the class they are labeled as.** Note the ambiguity (between the comma and apostrophe) and the mislabeled input.



# EXPERIMENTAL SETUP

## I Model Descriptions:

- Teacher layers Semi Supervised
  - Convolutional Networks
    - MNIST
    - SVHN
    - Glyph(with ResNet backbone)
- Teacher layers Fully Supervised
  - Decision tree
    - UCI Adult

# EXPERIMENTAL SETUP

I Teacher Ensembles(number of teachers, and therefore partitions of data)

- 100
- 500
- 1000
- 5000

I Queries

- 500-12000 depending on dataset

I Privacy parameters:

- $\delta = 10^{-8}$  probability that privacy will not be held

# RESULTS

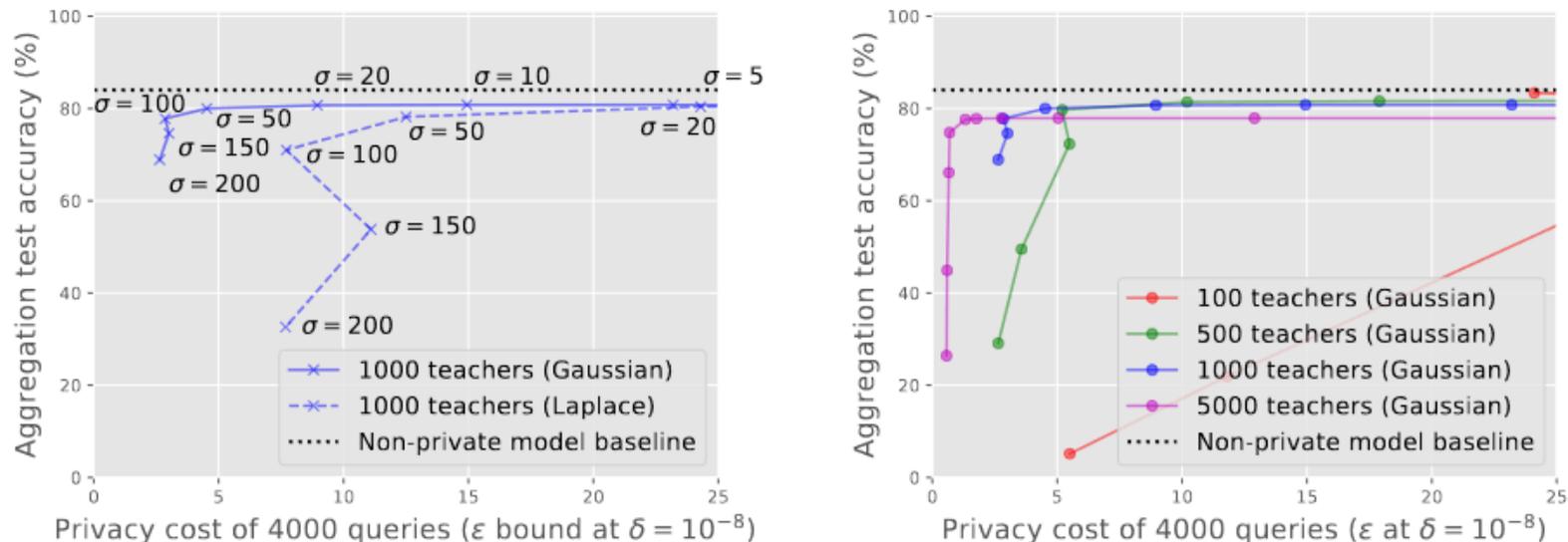


Figure 4: **Tradeoff between utility and privacy for the LNMax and GNMax aggregators on Glyph**: effect of the noise distribution (left) and size of the teacher ensemble (right). The LNMax aggregator uses a Laplace distribution and GNMax a Gaussian. Smaller values of the privacy cost  $\epsilon$  (often obtained by increasing the noise scale  $\sigma$ —see [Section 4](#)) and higher accuracy are better.

# RESULTS

Dataset	Aggregator	Queries answered	Privacy bound $\epsilon$	Accuracy	
				Student	Baseline
MNIST	LNMax (Papernot et al., 2017)	100	2.04	98.0%	99.2%
	LNMax (Papernot et al., 2017)	1,000	8.03	98.1%	
	Confident-GNMax ( $T=200, \sigma_1=150, \sigma_2=40$ )	286	<b>1.97</b>	<b>98.5%</b>	
SVHN	LNMax (Papernot et al., 2017)	500	5.04	82.7%	92.8%
	LNMax (Papernot et al., 2017)	1,000	8.19	90.7%	
	Confident-GNMax ( $T=300, \sigma_1=200, \sigma_2=40$ )	3,098	<b>4.96</b>	<b>91.6%</b>	
Adult	LNMax (Papernot et al., 2017)	500	2.66	83.0%	85.0%
	Confident-GNMax ( $T=300, \sigma_1=200, \sigma_2=40$ )	524	<b>1.90</b>	<b>83.7%</b>	
Glyph	LNMax	4,000	4.3	72.4%	82.2%
	Confident-GNMax ( $T=1000, \sigma_1=500, \sigma_2=100$ )	10,762	2.03	<b>75.5%</b>	
	Interactive-GNMax, two rounds	4,341	<b>0.837</b>	73.2%	

250 Teachers used for MNIST SVHN ADULT. 5000 for Glyph.  
 $\delta = 10^{-5}$  For MNIST ADULT,  $\delta = 10^{-6}$  SVHN,  $\delta = 10^{-8}$  Glyph

# RESULTS

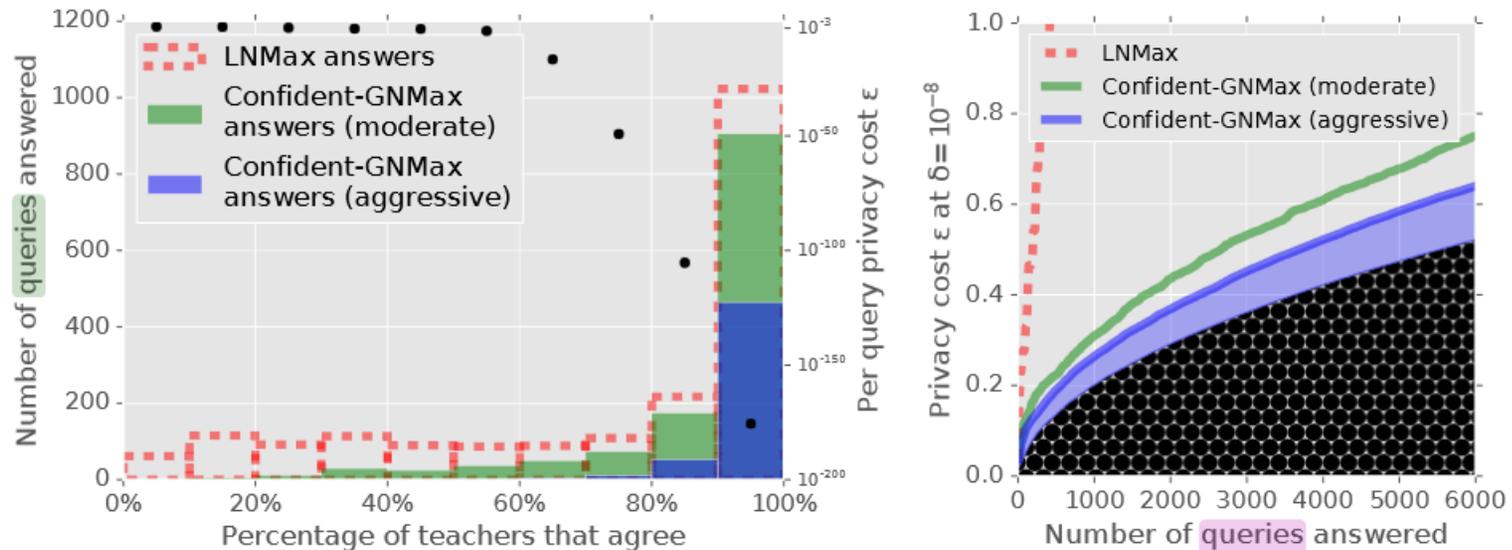


Figure 5: **Effects of the noisy threshold checking:** *Left:* The number of queries answered by LNMax, Confident-GNMax moderate ( $T=3500, \sigma_1=1500$ ), and Confident-GNMax aggressive ( $T=5000, \sigma_1=1500$ ). The black dots and the right axis (in log scale) show the expected cost of answering a single query in each bin (via GNMax,  $\sigma_2=100$ ). *Right:* Privacy cost of answering all (LNMax) vs only inexpensive queries (GNMax) for a given number of answered queries. The very dark area under the curve is the cost of selecting queries; the rest is the cost of answering them.



# CONCLUSION

## I Build off the PATE method

- New methods for aggregating teacher/student answers provide a privacy preserving technique that reduces leakage
  - Caps queries at a confidence interval
  - Stops overfitting student's queries by checking confidence of student's answers
- Improvements across the board in Privacy  $\epsilon$  loss(lower is better)
- Shows that this PATE method has the potential to be used at scale(thousands of class labels)
  - Generalization improved by changing perturbation method



**QUESTIONS?**

# REFERENCES

- [1] <https://www.statista.com/statistics/254266/global-big-data-market-forecast/>
- [2] L. Sweeney, Simple Demographics Often Identify People Uniquely. Carnegie Mellon University, Data Privacy Working Paper 3. Pittsburgh 2000.
- [3] Papernot, N., Song, S., Mironov, I., Raghunathan, A., Talwar, K., and Erlingsson, Ú., “Scalable Private Learning with PATE”, *arXiv e-prints*, 2018.

# References

---

1. Liu et al. (2020) When Machine Learning Meets Privacy: A Survey and Outlook ([link](#))
2. Rigaki and Carcia (2021) A Survey of Privacy Attacks in Machine Learning ([link](#))
3. Cristofaro (2020) An Overview of Privacy in Machine Learning ([link](#))
4. Borealis AI Tutorial #13: Differential Privacy II: Machine Learning and Data Generation ([link](#))
5. Davide Testuggine and Ilya Mironov blog: Differential Privacy Series Part 1- DP-SGD Algorithm Explained ([link](#))
6. Joseph P. Near and Chiké Abuah, Programming Differential Privacy – Chapter III: Differential Privacy ([link](#))