# CSCI 699: Privacy Preserving Machine Learning - Week 2

**Differential Privacy**

Sai Praneeth Karimireddy, Sep 6 2024

# Recap
*— why privacy*

- We saw many definitions of privacy

  - De-identification / suppression     *side information "linkage attacks"*

  - K-anonymity      *— "negation"*

  - L-diversity

- We saw none of them really protected privacy and were easily broken

- Hinted at a more widely accepted definition.

# Takeaways

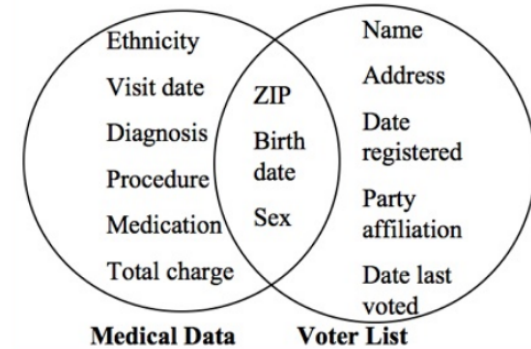## Requirements for privacy definition

*Differential Privacy*

*Post-processing*

- Unaffected by auxiliary information: we should not be able to combine extra data to undo privacy.

- Composition: We should understand what happens when data is continuously released.

- Today we will come with such a privacy definition.



| Medical Data | | Voter List |
|---|---|---|
| Ethnicity | | Name |
| Visit date | ZIP | Address |
| Diagnosis | Birth date | Date registered |
| Procedure | | Party affiliation |
| Medication | Sex | Date last voted |
| Total charge | | |

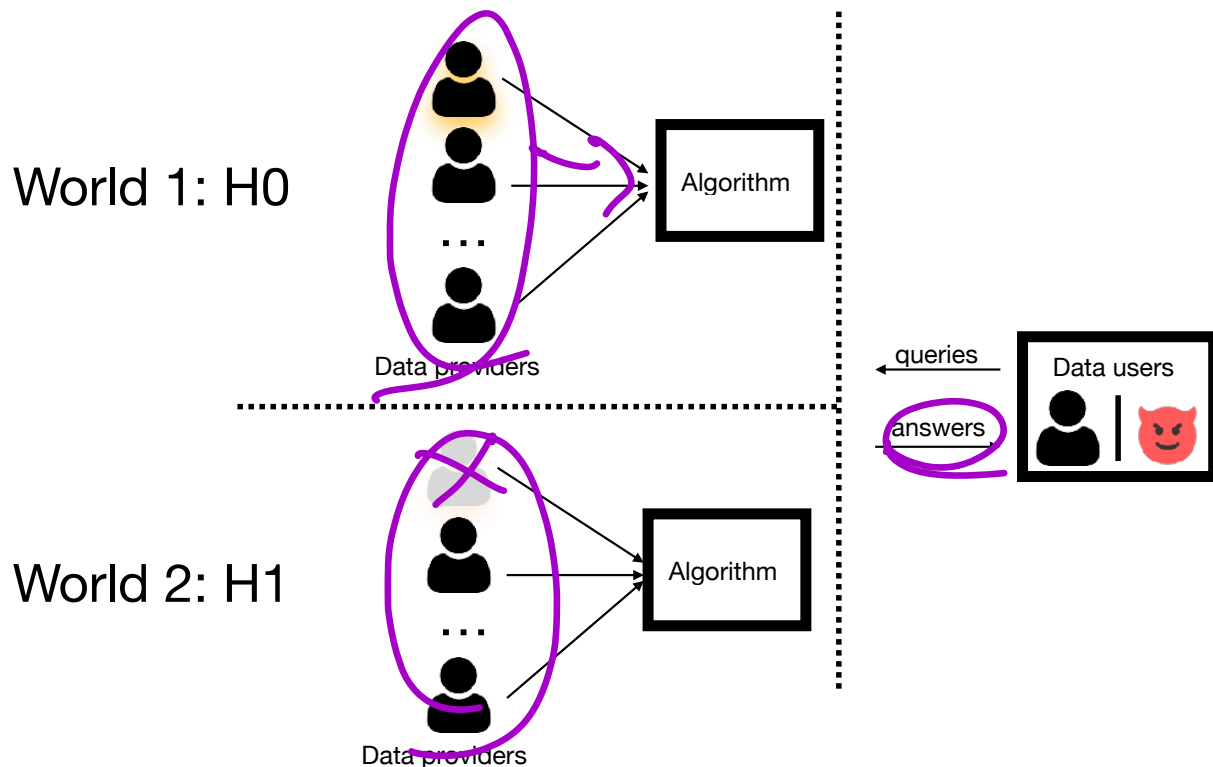# Quantifying Privacy Leakage
## Attempt 2

Relative Privacy: quantify **new** information leaked

"An analysis of a dataset is private if what can be learned about an individual in the dataset is not much more than what would be learned if the same analysis was conducted without them in the dataset"

- **Intuition**: Whether Bob is present in the data or not, the answer should not change much.

- Then, from looking at the answer, we will not learn whether Bob was present in the data or not.

- Gives Bob plausible deniability.

# Quantifying Privacy Leakage
## Attempt 2



World 1: H0

World 2: H1

Algorithm

Data providers

Data users

queries

answers

- In world 2 only Bob is removed/replaced.

- Now from the answer, how easily can guess the correct world?
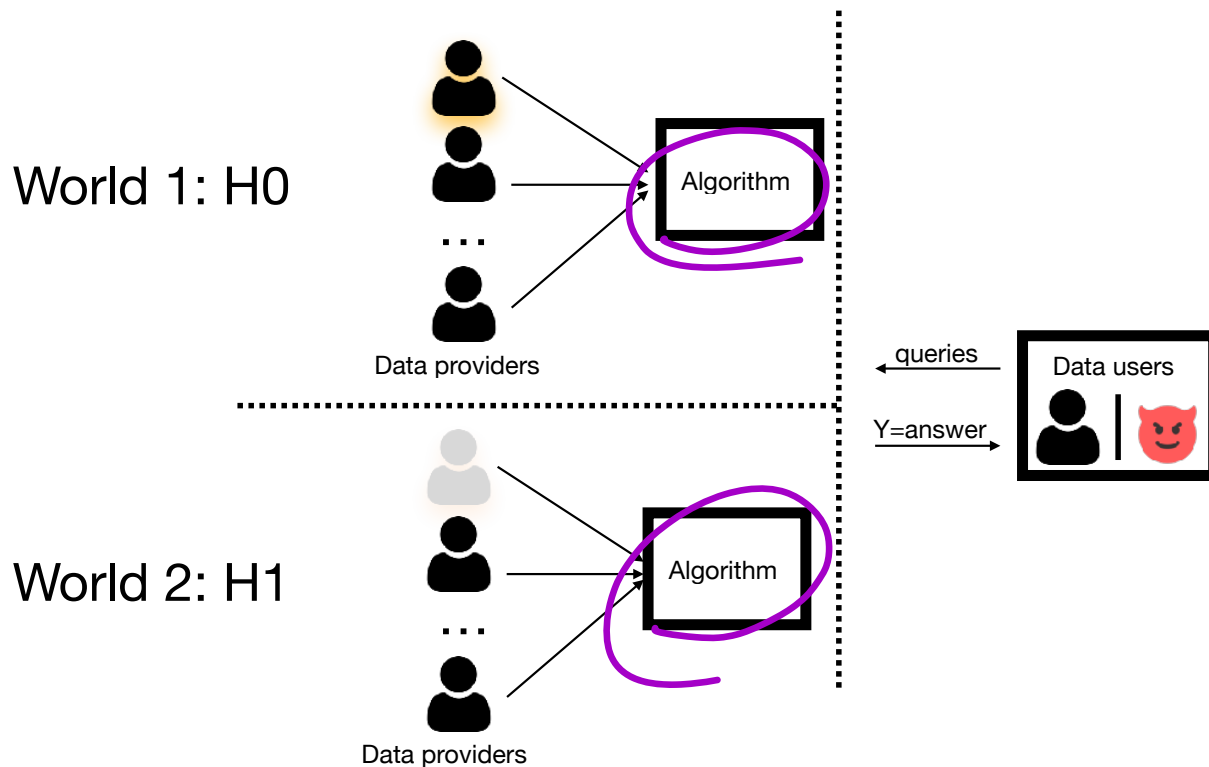
# Quantifying Privacy Leakage

# Membership Inference

**As a definition of privacy**

attack

$D = \{x_1, \ldots, x_n\}$

$x_i$

World 1: H0
with $x_i$



Data providers

World 2: H1
without $x_i$

Data providers

queries

Data users

Y=answer

- We know everything about the algorithm and even $D, x_i$,

- Only 1 bit unknown - $H_0$ or $H_1$?

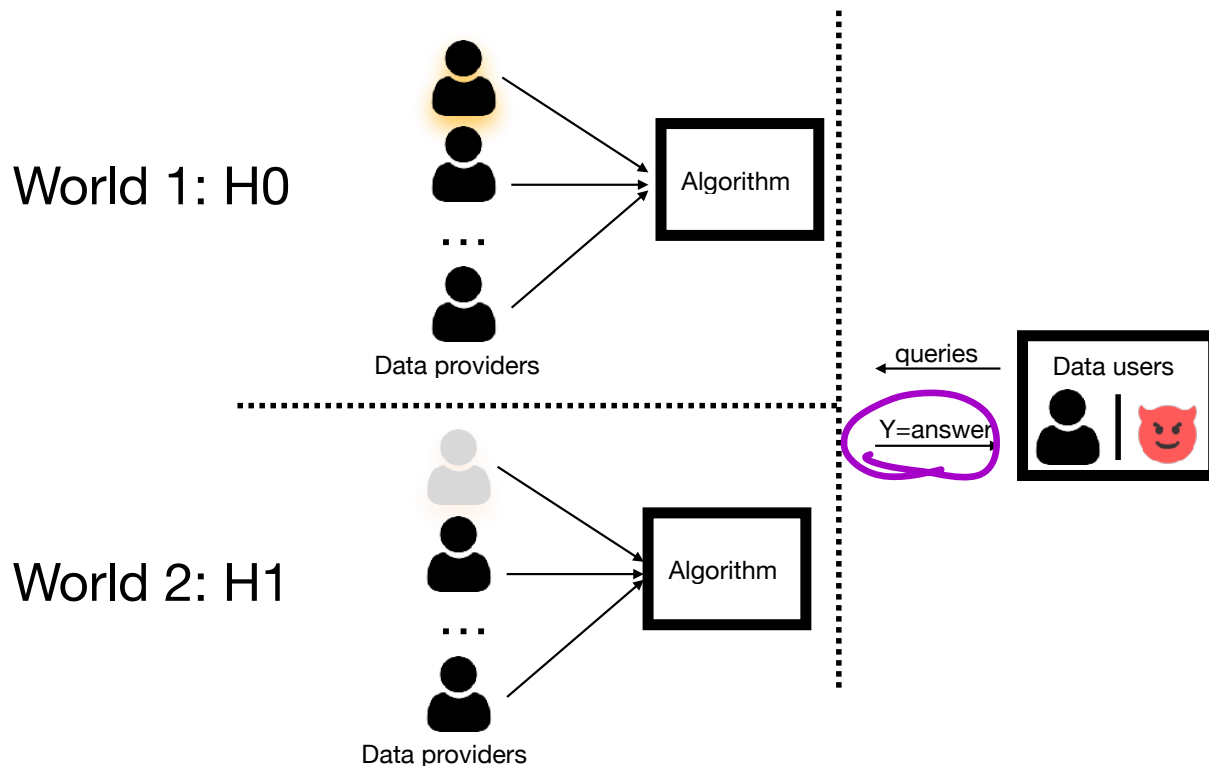- We observe an output $Y$

- Need to guess if it came from H0 or H1

# Membership Inference

World 1: H0

Algorithm

Data providers

World 2: H1

Algorithm

Data providers

queries

Y=answer

Data users

- Can a deterministic algorithm be private?

# Membership Inference

→ Y is a random variable
→ acc is also a random



World 1: H0

World 2: H1

Data providers

Data providers

Algorithm

Algorithm

queries

Y=answer

Data users

- Can a deterministic algorithm be private?

- No - adversary can simply compute $Y = f(D)$ or $f(D \backslash x_i)$?

- Need randomness - adversary will have type I and type II errors
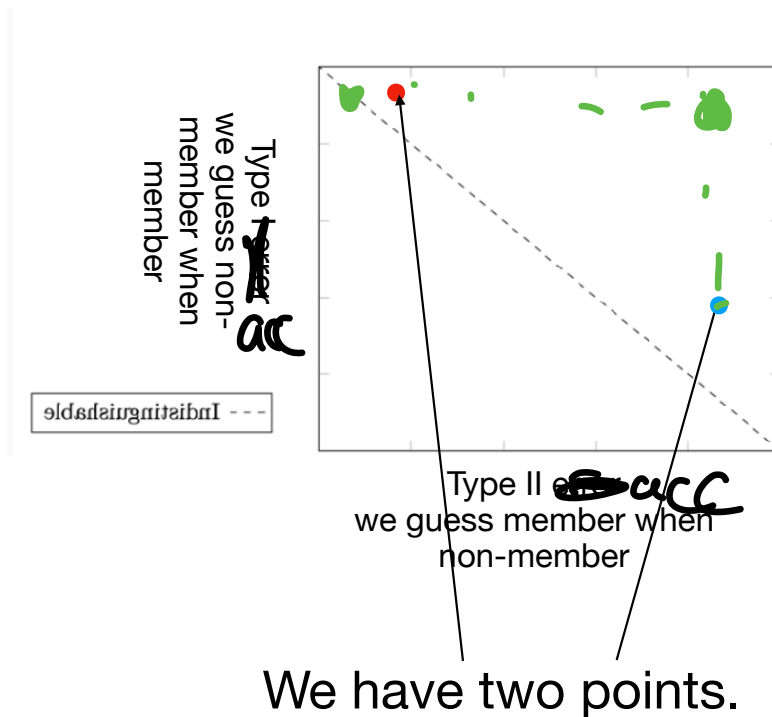
# Membership Inference
## Quantifying attack success



We get a point.

- Suppose we run multiple runs

- Count the number of times the adv guesses H0 vs H1 correctly

- We can compute Type I and Type II errors.
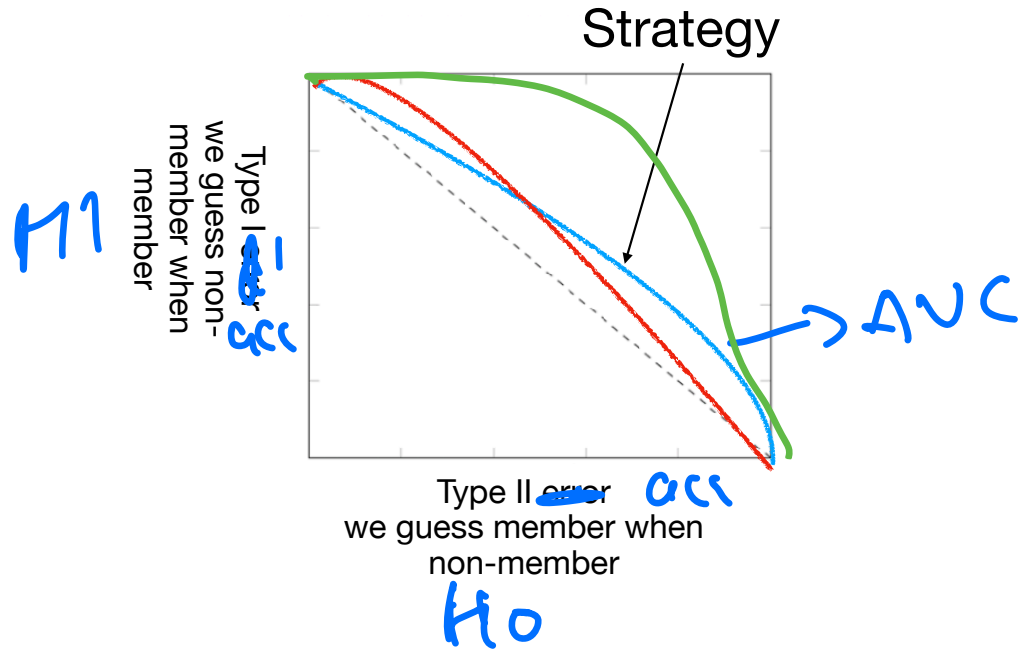
# Membership Inference
## Quantifying attack success



Type I error, we guess non-member when member

Type II error, we guess member when non-member

Indistinguishable

We have two points.

- Suppose we have two algorithms, each with different type I and type II errors.

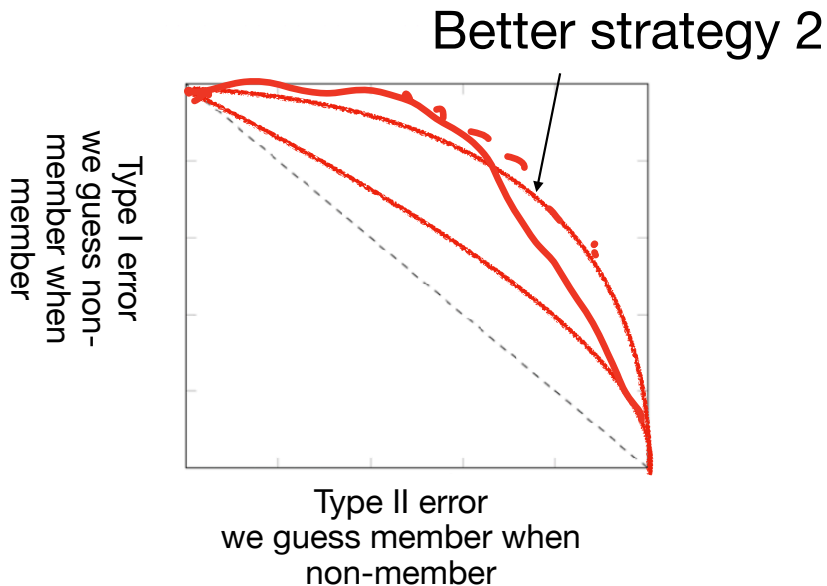- Which one has more privacy leakage?

# Membership Inference
## Tradeoff curve



Strategy

Type I error
we guess non-member when member

M1
acc

→ AUC

Type II error
we guess member when non-member

acc

H0

- Depends on what we care

- E.g. its important not to miss anyone e.g. sending cat ads to pet owners - coverage

- Not ok if we are accusing them of a crime - precision much more important

- Impossible to compare individual points - need to compare entire trade off curves.
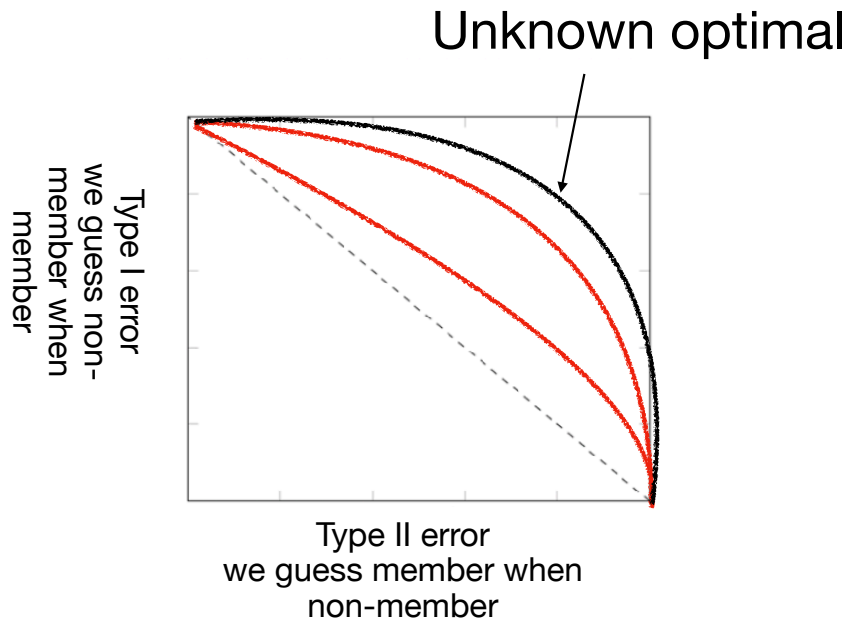
# Membership Inference
## Comparing tradeoff curves

Better strategy 2



Type I error
we guess non-
member when
member

Type II error
we guess member when
non-member

- Tradeoff curve depends on testing strategy adversary uses.

- Strategy 2is better than Strategy 1 if the curve is uniformly above.

- Higher curve means we've found more privacy leakage
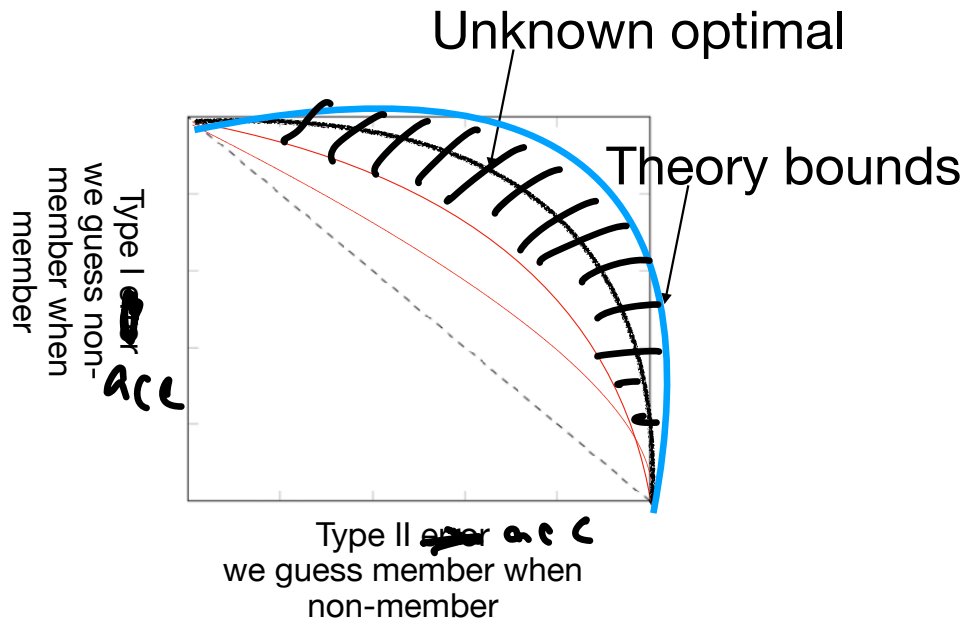
# Membership Inference
## Optimal tradeoff curve

Unknown optimal

Type I error
we guess non-
member when
member

Type II error
we guess member when
non-member

- There is an optimal strategy

- use this to quantify privacy leakage

- What if no single strategy is best?

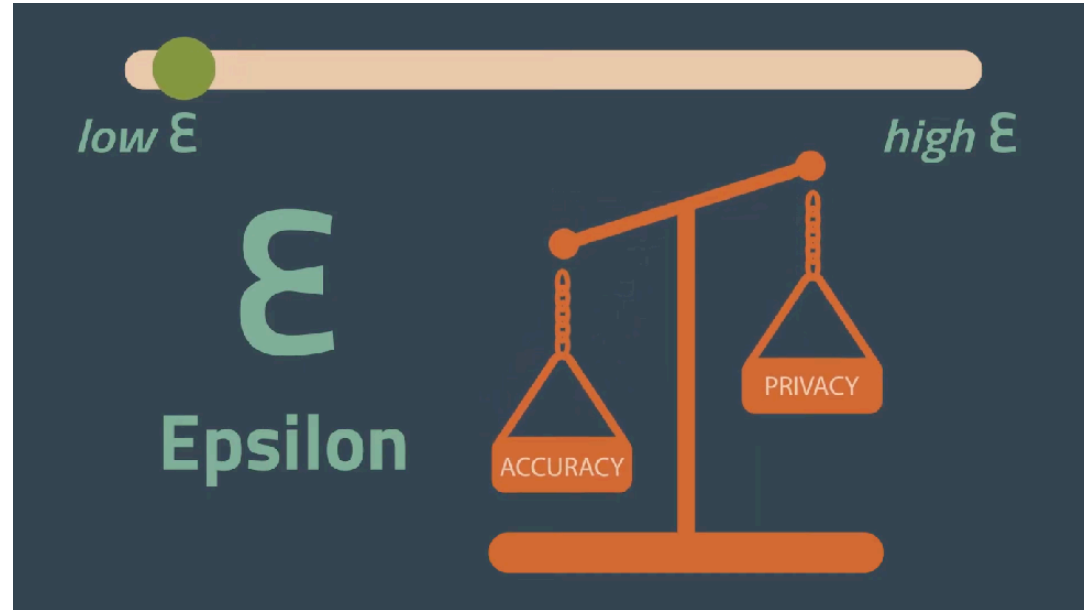- **Neyman–Pearson lemma** guarantees existence of uniformly most powerful test.

# Membership Inference
## Privacy from tradeoff curve



- Use optimal strategy to quantify privacy

- But empirical tests only give an lower-bound

- Need theory to give upper-bound

# Differential Privacy

# Differential Privacy

## Calibrating Noise to Sensitivity in Private Data Analysis

2006

Cynthia Dwork[1], Frank McSherry[1], Kobbi Nissim[2], and Adam Smith[3*]

**2017 Gödel Prize**

Differential privacy is a powerful theoretical model for dealing with the privacy of statistical data. The intellectual impact of differential privacy has been broad, influencing thinking about privacy across many disciplines. The work of Cynthia Dwork (Harvard University), Frank McSherry (independent researcher), Kobbi Nissim (Harvard University), and Adam Smith (Harvard University) launched a new line of theoretical research aimed at understanding the possibilities and limitations of differentially private algorithms. Deep connections have been exposed in other areas of theory (including learning, cryptography, discrepancy, and geometry) and have created new insights affecting multiple communities.

# Differential Privacy

**Threat model**

- Let $\chi$ be a the domain of training data

- A dataset $D \in \chi^n$ is a multiset of n records/rows of $\chi$

- $D$ (sensitive data) $\longrightarrow$ algorithm $\longrightarrow$ $Y$ (answers)

- Attacker wants to infer some information about $D \in \chi^n$

  - observes $Y$

  - knows algorithm, domain $\chi$, and potentially more prior information

  - cannot control what attacker knows

# Differential Privacy
## Threat model

- Attacker wants to infer some information about $D \in \chi^n$

    - observes $Y$, knows algorithm, domain $\chi$, and prior information.

    - can compute likelihood of dataset:

    algorithm    prior knowledge

    $$Pr[D \,|\, Y] = \frac{Pr[Y \,|\, D] \cdot Pr[D]}{Pr[Y]}$$

# Differential Privacy
## Performing membership inference

- Attacker wants to infer presence of $x \in X$?

  - observes $Y$, knows algorithm, domain $\chi$, and even $D\backslash x \in \chi^{n-1}$

  - can compute likelihood of x in dataset

algorithm                    prior knowledge

$$Pr[x' \,|\, Y] = \frac{Pr[Y \,|\, x'] \cdot Pr[x']}{Pr[Y]}$$

# Differential Privacy
## Performing membership inference

- Attacker wants to infer presence of $x \in X$?

  - can compute likelihood of x in dataset

    algorithm                          prior knowledge

    $$Pr[x' \,|\, Y] = \frac{Pr[Y \,|\, x'] \cdot Pr[x']}{Pr[Y]}$$

  - Can even recover $x$ using max-likelihood

    $$\hat{x} = \arg \max_{x'} Pr[Y \,|\, x']Pr[x']$$

# Differential Privacy
**Goal**

$$Pr\{D \mid Y=y\}$$
$$= Pr\{D \mid \xi_x) \mid Y=y\}$$
$$\parallel$$
$$Pr\{Y=y \mid D\}$$
$$= Pr\{Y=y \mid D\{\xi_x\}$$

- Attacker wants to infer some information about $D \in \chi^n$

  - can compute likelihood of seeing some dataset

algorithm          prior knowledge

$$Pr[D \mid \theta] = \frac{Pr[\theta \mid D] \cdot Pr[D]}{Pr[\theta]}$$

- We design a private algorithm by controlling $Pr[\theta \mid D]$

# Differential Privacy
**Strict definition**

- Perfect relative indistinguishability: For all inputs, the output probability is the same.

$$\forall D, D', y : \quad \Pr[Y = y \,|\, \mathscr{D} = D] = \Pr[Y = y \,|\, \mathscr{D} = D']$$

- The mechanism does not leak any information about D

- However, achieving it is very hard, does not allow any information about D.

# Differential Privacy

## A better definition

$\mathcal{E}\text{-}\mathcal{DP}$

- Some indistinguishability: For all neighboring datasets, the output probabilities are bounded.

$$\forall y, \forall \text{ similar } D, D': \quad \frac{\Pr[Y = y \mid \mathcal{D} = D]}{\Pr[Y = y \mid \mathcal{D} = D']} \leq \text{ constant } e^{\mathcal{E}}$$

- It means by observing any $Y$, adversary is NOT able to distinguish between inputs x and x' beyond a bounded certainty.

- What does neighboring datasets mean? Depends on use case
  - location positions that are within some range
  - datasets that differ in one individual row (our focus)
  - edit distance 1

# Differential Privacy

**Formal definition**

$\varepsilon$-Differential Privacy:

An algorithm A satisfies $\varepsilon$-DP if for any neighboring datasets $D, D' \in \chi^n$ and $y \in \mathcal{Y}$

$$\text{Privacy} \quad \log \frac{Pr[A(D) = y]}{Pr[A(D') = y]} \leq \varepsilon \text{ a.s.}$$

- Recall that $D$ (sensitive data) $\longrightarrow$ algorithm $\longrightarrow Y$ (answers)

- So we have, $Pr[Y | D] = Pr[A(D) = Y]$

# Differential Privacy

**Formal definition**

$\varepsilon$-Differential Privacy:

An algorithm A satisfies $\varepsilon$-DP if for any similar datasets $D, D' \in \chi^n$ and $y \in \mathcal{Y}$

$$\log \frac{Pr[A(D) = y]}{Pr[A(D') = y]} \leq \varepsilon$$

- $\varepsilon = 0$ means perfect privacy

- $\varepsilon \gg 0$ means not private

# Differential Privacy
## Source of randomness

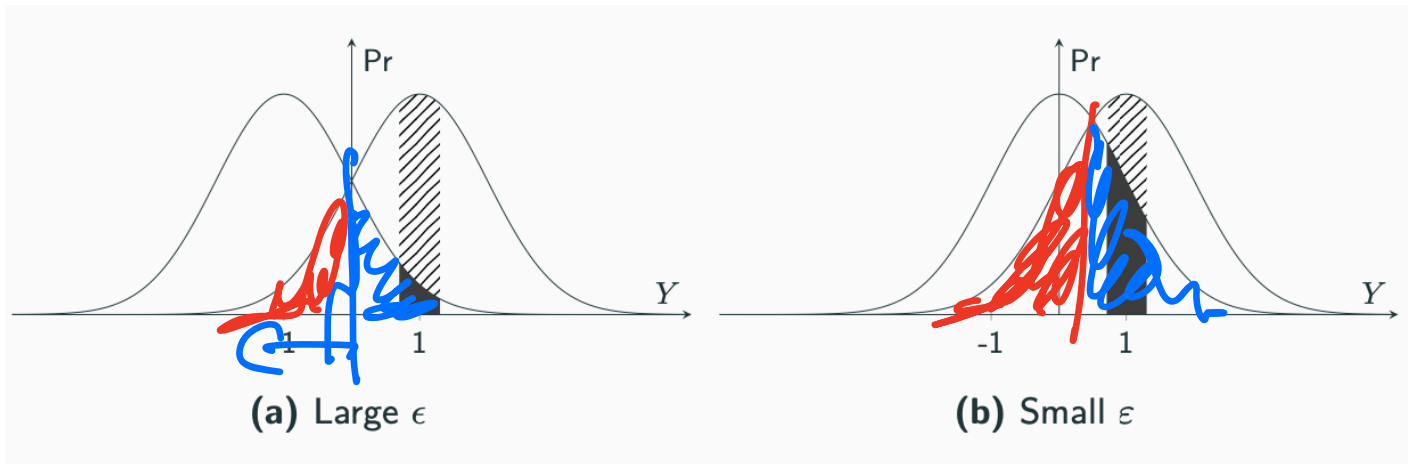$$\log \frac{Pr[A(D) = y]}{Pr[A(D') = y]} \leq \varepsilon \quad a.s.$$

- In $\Pr[A(D) = y]$, over what randomness is the probability defined?

    - The randomness of the algorithm?
        - Yes

    - Randomness of the data $D \in \chi^n$?
        - No.
        - We look at all possible values of $D, D'$ i.e. worst case

# Differential Privacy

## Visual representation

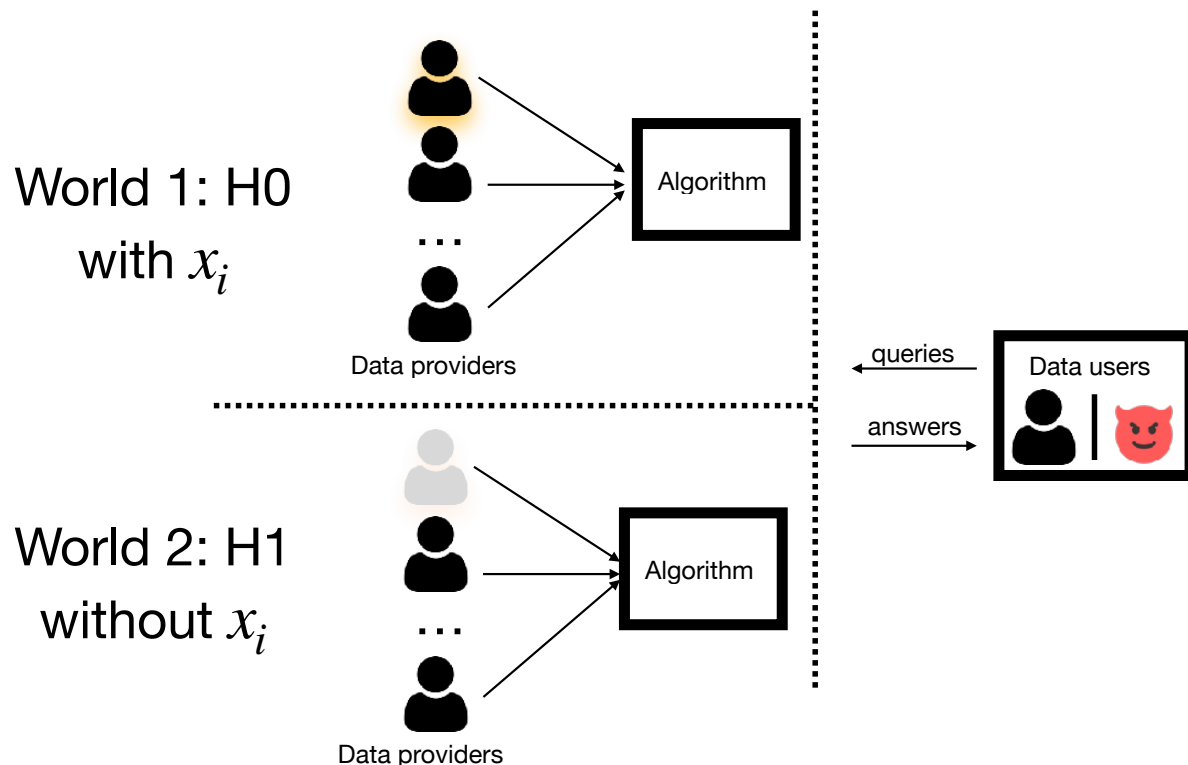- Consider $D = \langle x_1, \cdots, x_i, \cdots x_n \rangle$, and a similar dataset $D' = \langle x_1, \cdots, \cancel{x_i}, \cdots x_n \rangle$

- $\varepsilon$-DP means $\dfrac{Pr[A(D) = y]}{Pr[A(D') = y]} \leq \exp(\varepsilon)$



(a) Large $\epsilon$

(b) Small $\epsilon$

# Differential Privacy
## Recall Membership Inference



World 1: H0
with $x_i$

Algorithm

Data providers

World 2: H1
without $x_i$

Algorithm

...

Data providers

queries
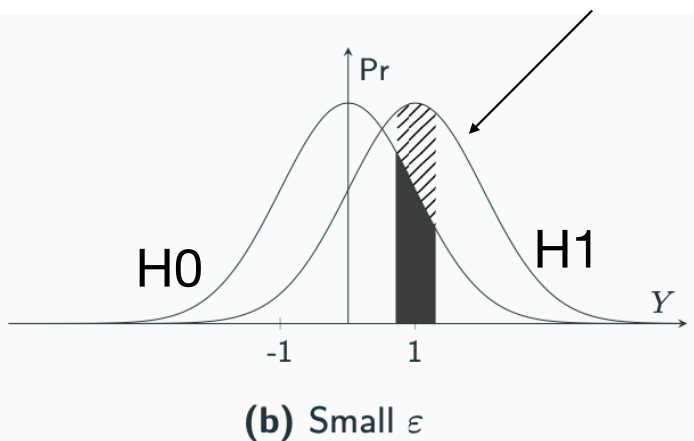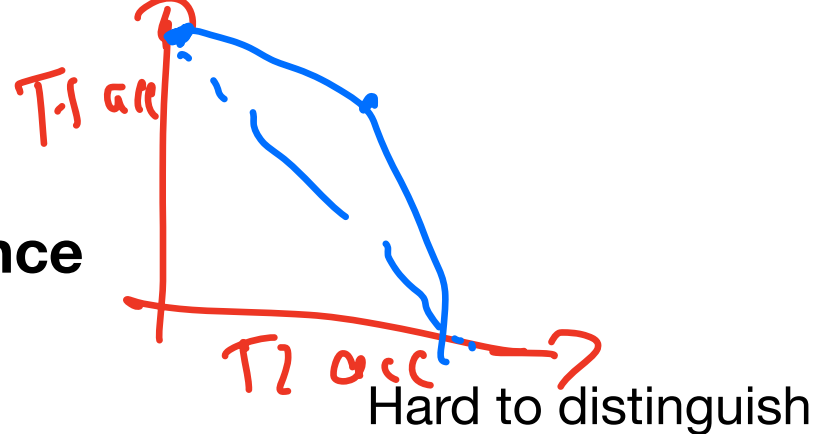
answers

Data users

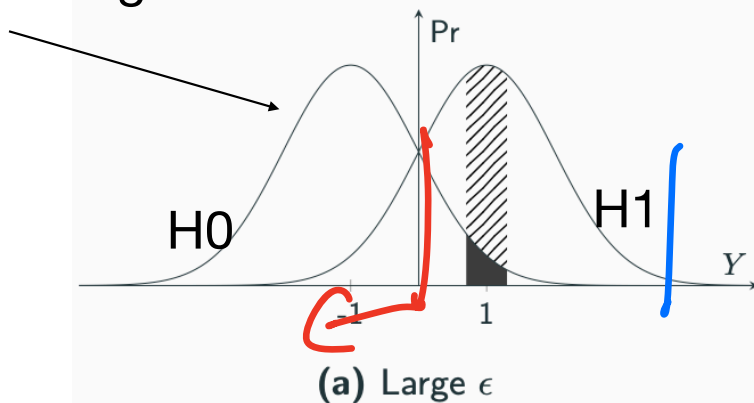- We know everything about the algorithm and even $D \backslash x_i$

- We observe an output Y

- Need to guess if it came from H0 or H1

# Differential Privacy
## Connection to Membership Inference

T·S all

T2 occ ?

Hard to distinguish

Easy to distinguish



(a) Large $\epsilon$ — H0, H1, Pr, Y, -1, 1

(b) Small $\epsilon$ — H0, H1, Pr, Y, -1, 1

- We observe Y = 1.

- Can you guess H0 or H1?

# Differential Privacy and membership inference
## Quantifying connection

Theorem

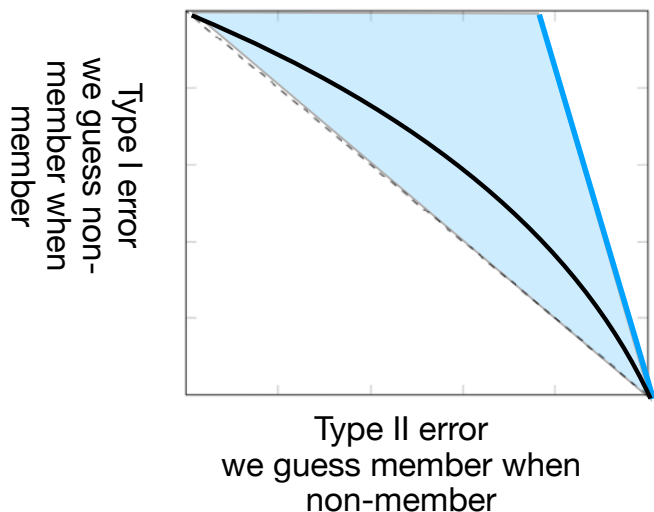Suppose A satisfies $\varepsilon$-DP for datasets $D, D'$ which differ by one datapoint. Then, we have

- $Pr[\text{guess H0} \mid H1] + e^{\varepsilon}Pr[\text{guess H1} \mid H0] \geq 1$

- $e^{\varepsilon}Pr[\text{guess H0} \mid H1] + Pr[\text{guess H1} \mid H0] \geq 1$

- Type I error $= Pr[\text{guess H0} \mid H1]$

- Type II error $= Pr[\text{guess H1} \mid H0]$

# Differential Privacy and membership inference

## Visualizing connection

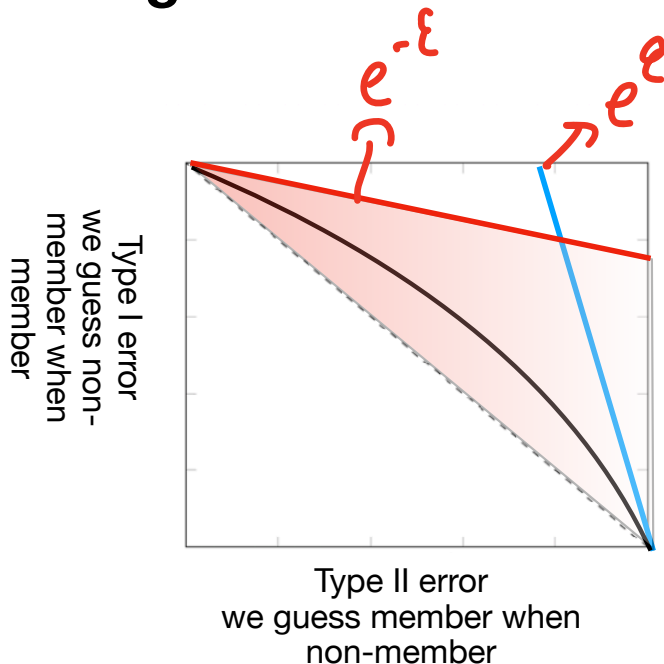$\varepsilon - DP$



Type I error
we guess non-member when member

Type II error
we guess member when non-member

- $Pr[\text{guess H0}\,|\,H1] + e^{\varepsilon}Pr[\text{guess H1}\,|\,H0] \geq 1$

- gives us blue line with slope $e^{\varepsilon}$

# Differential Privacy and membership inference

## Visualizing connection



Type I error
we guess non-member when member

Type II error
we guess member when
non-member

- $e^{\varepsilon} Pr[\text{guess H0} \mid H1] + Pr[\text{guess H1} \mid H0] \geq 1$

  - gives the red line with slope $e^{-\varepsilon}$

# Differential Privacy and membership inference

## Visualizing tradeoff curve of DP

Theoretical upper bound $= \mathcal{E}\text{-}DP$



Type I error
we guess non-member when member

Type II error
we guess member when non-member

- $Pr[\text{guess H0} \,|\, H1] + e^{\varepsilon} Pr[\text{guess H1} \,|\, H0] \geq 1$
  - gives us blue line
- $e^{\varepsilon} Pr[\text{guess H0} \,|\, H1] + Pr[\text{guess H1} \,|\, H0] \geq 1$
  - gives the red line

what $\varepsilon$ does this satisfy?

# Aside: Is Putin's popularity calculation private?

**List Experiment**

- Split users randomly into two groups

- Design a set of options very similar to the one you actually care about

- To control only ask about the rest. To the treatment include your option.

- Does this satisfy DP?

*n = 100*

**How many of the following things do you personally support? You don't need to say which ones you support, just specify the number of them (0, 1, 2, 3, or 4).**

Actions of the Russian armed forces in Ukraine

Legalization of same-sex marriage in Russia

Increase in monthly allowances for low-income Russian families

State measures to prevent abortion

**I support:**
- ○ 0
- ○ 1
- ○ 2
- ○ 3
- ○ 4 of these things

**How many of the following things do you personally support? You don't need to say which ones you support, just specify the number of them (0, 1, 2, or 3).**

State measures to prevent abortion

Legalization of same-sex marriage in Russia

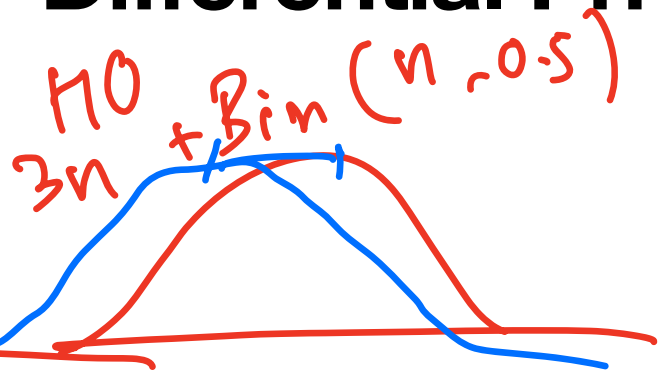Increase in monthly allowances for low-income Russian families
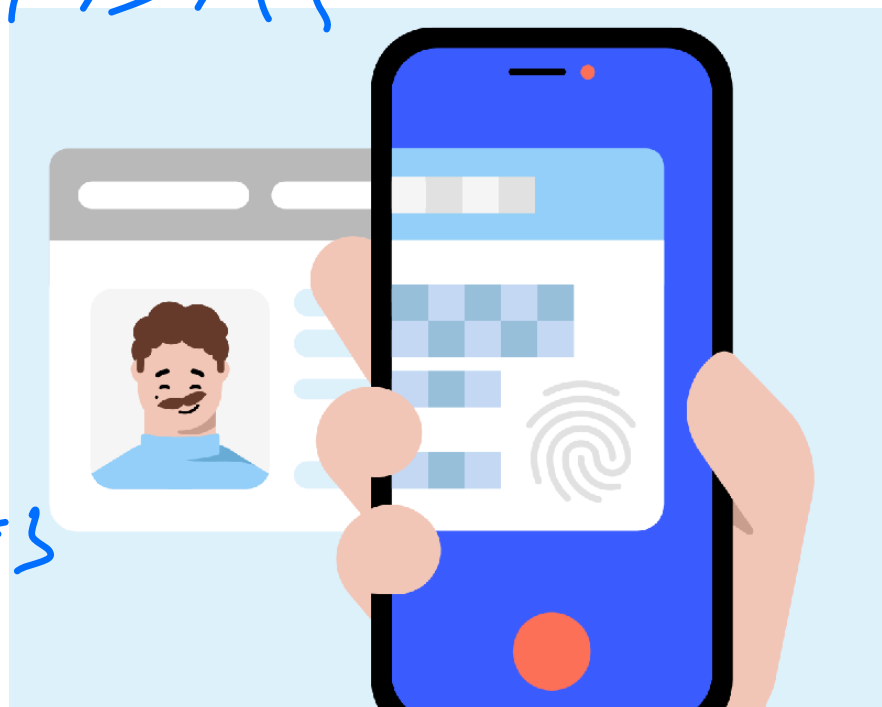
**I support:**
- ○ 0
- ○ 1
- ○ 2
- ● 3 of these things

*50 / 50*

Chapkovski and Schaub 2022. "Do Russians tell the truth when they say they support the war in Ukraine? Evidence from a list experiment" LSE Blog

$D = ?$

$?c_i = ?$

$Alg = ?$

$x = \sum_i y_i \begin{cases} 3 \\ 4 \end{cases}$

$y = \underbrace{\phantom{xxxx}}_{} \to \{0, 1, 3, 4\}_i$

# Algorithms for Differential Privacy

$HO$
$3n + Bin(n, 0.5)$

$H1 = H0 \cdot 3$

$n = 1$       $H$

$f(1$ $HO$

$0$   $1$

$Y = $   $0$   $\leftarrow D$

$Y = 1$   $1$   $\leftarrow D \setminus \{x\}$

$output = Y + Noise$

$\downarrow$

$Lap(\frac{1}{\varepsilon})$

$\propto e^{-t \cdot \varepsilon}$

$\leq e^{\varepsilon}$

Laplace

# Differentially Private Algorithms

**Just add Laplace noise**

$$\forall y, \forall \text{ similar } D, D' : \quad \frac{Pr[A(D) = y]}{Pr[A(D') = y]} \leq \exp(\varepsilon)$$

- Suppose A(D) = 0, A(D') = 1.

- Release $\hat{y} = y + \text{Laplace}(0, \varepsilon^{-1})$

- $z \sim \text{Laplace}(\mu, b) \Rightarrow p(z) = \dfrac{1}{2b} e^{\frac{-|z - \mu|}{b}}$
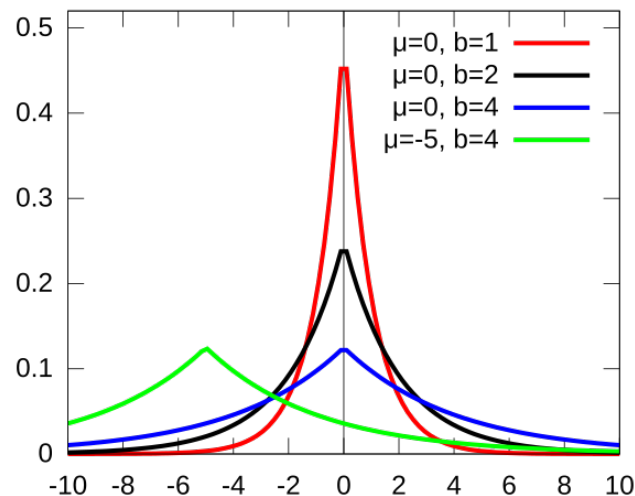
# Differentially Private Algorithms
## Just add Laplace noise

$$\forall y, \forall \text{ similar } D, D' : \quad \frac{Pr[A(D) = y]}{Pr[A(D') = y]} \leq \exp(\varepsilon)$$

- Suppose A(D) = 0, A(D') = 1. Release $\hat{y} = y + \text{Laplace}(0, \varepsilon^{-1})$

- $Pr[\hat{y} \mid y = 0] = \text{Laplace}(0, \varepsilon^{-1})$ and $Pr[\hat{y} \mid y = 1] = \text{Laplace}(1, \varepsilon^{-1})$

- $\frac{Pr[A(D) = y]}{Pr[A(D') = y]} = \frac{e^{-\varepsilon|y|}}{e^{-\varepsilon|y-1|}} = e^{\varepsilon}$

$$e^{\varepsilon(|y-1| - |y|)}$$



PDF of Laplace distribution

# Differentially Private Algorithms
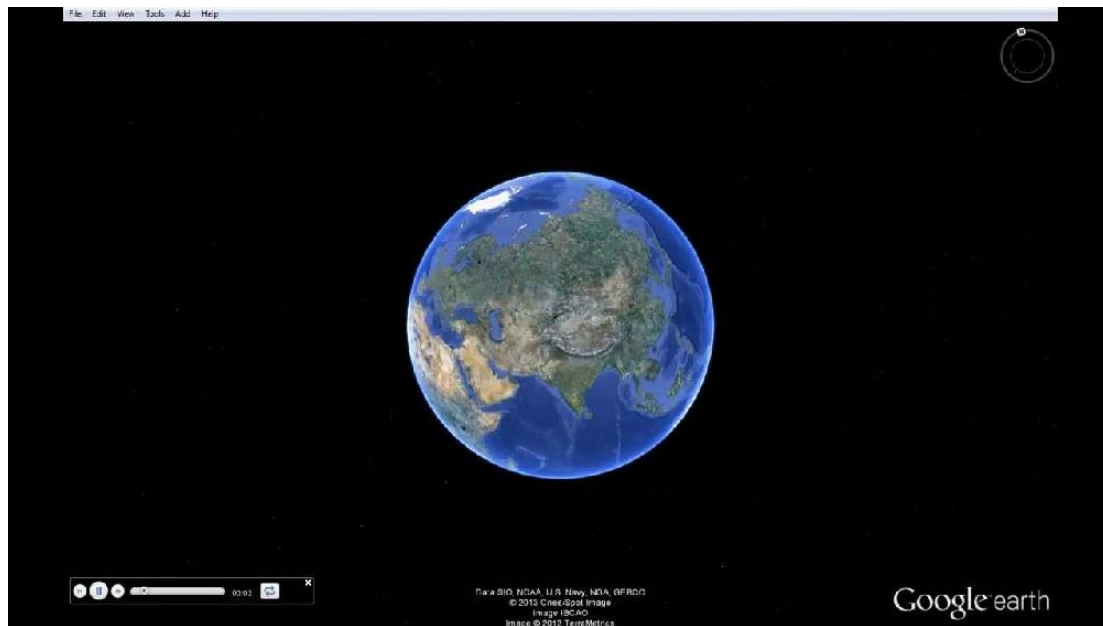**Sensitivity**



- I release average income at different zoom levels. Added Lap(0,1).

- Do they all leak same amount of privacy?

# Differentially Private Algorithms
## Sensitivity and Laplace mechanism

- **Definition: Sensitivity** of a function $f : (x_1, \cdots, x_n) \mapsto (y_1, \cdots, y_k)$ with respect to a norm $\|\cdot\|$ is

$$\Delta f = \max_{\text{similar datasets } D,D'} \|f(D) - f(D')\|$$

•

### Theorem

Suppose $f$ is $\Delta$-sensitive with respect to $\|\cdot\|_1$. Then, the following satisfies $\varepsilon$-DP:

$$[A(D)]_i = [f(D)]_i + \text{Laplace}(0, \Delta\varepsilon^{-1})$$

$$H0 : \quad f(D) + Lap\left(0, \frac{\Delta}{\varepsilon}\right)$$

$$= Lap\left(f(D), \frac{\Delta}{\varepsilon}\right)$$

$$H1 = \quad f(D') + Lap\left(0, \frac{\Delta}{\varepsilon}\right)$$

$$= Lap\left(f(D'), \frac{\Delta}{\varepsilon}\right)$$

$$\frac{P(Y=y|H0)}{P(Y=y|H1)} = \frac{\exp\left(-|y - f(D)|\, \varepsilon/\Delta\right)}{\exp\left(-|y - f(D')|\, \varepsilon/\Delta\right)}$$

$$= \exp\left(\frac{\varepsilon}{\Delta}\left(|y - f(D')| - |y - f(D)|\right)\right) \quad \curvearrowleft \, \triangle^{\prime e} \, ineq$$

$$\leq \exp\left(\frac{\varepsilon}{\Delta}\underbrace{|f(D) - f(D')|}_{\leftarrow \Delta}\right)$$

$$\leq \exp\left(\frac{\varepsilon}{\Delta} \cdot \cancel{\Delta}\right)$$

$$Y_i \mid H0 \propto Lap\left(\underline{[f(D)]_i}, \frac{\triangle}{\varepsilon}\right)$$

$$Y_i \mid H1 \propto Lap\left([f(D')]_i, \frac{\triangle}{\varepsilon}\right) d\triangle_\infty$$

$$\frac{Pr\left[Y_i = y_i \mid H0\right]}{Pr\left[Y_i = y_i \mid H1\right]} \leq \exp\left(\frac{\varepsilon}{\triangle} \left| [f(D)]_i - f(D') \right| \right)$$

$$\frac{Pr\left[\overset{\rightharpoonup}{Y} = y \mid H0\right]^{d\text{-}dim}}{Pr\left[\overset{\rightharpoonup}{Y} = y \mid H1\right]} = \prod_{i=1}^{d} \frac{Pr\left[Y_i = y_i \mid H0\right]}{Pr\left[Y_i = y_i \mid H1\right]}$$

$$\leq \exp\left(\frac{\varepsilon}{\triangle} \sum_i \left| f_i(D) - f_i(D') \right| \right)$$

$$\leq d\triangle_\infty \qquad = \exp\left(\frac{\varepsilon}{\triangle} \|f(D) - f(D')\|_1\right)$$

$$\leq \sqrt{d}\triangle_2 \qquad \leq \exp(\varepsilon) \qquad \leq \triangle$$

# Differentially Private Algorithms
**Sensitivity and Laplace mechanism**

$$\varepsilon_i \left( f_i(D) - f_i(D') \right)$$

$$\subseteq \triangle_2^2$$

- **Definition: Sensitivity** of a function $f : (x_1, \cdots, x_n) \mapsto (y_1, \cdots, y_k)$ with respect to a norm $\|\cdot\|$ is

$$\Delta f = \max_{\text{similar datasets } D, D'} \|f(D) - f(D')\|$$

- How much noise should we add if we have $\Delta$-sensitivity wrt $\| \cdot \|_{\infty}$

- What about $\Delta$-sensitivity wrt $\| \cdot \|_2$

- Laplace mechanism is great for functions with small $\ell_1$ sensitivity, not so much for small $\ell_2$ sensitivity
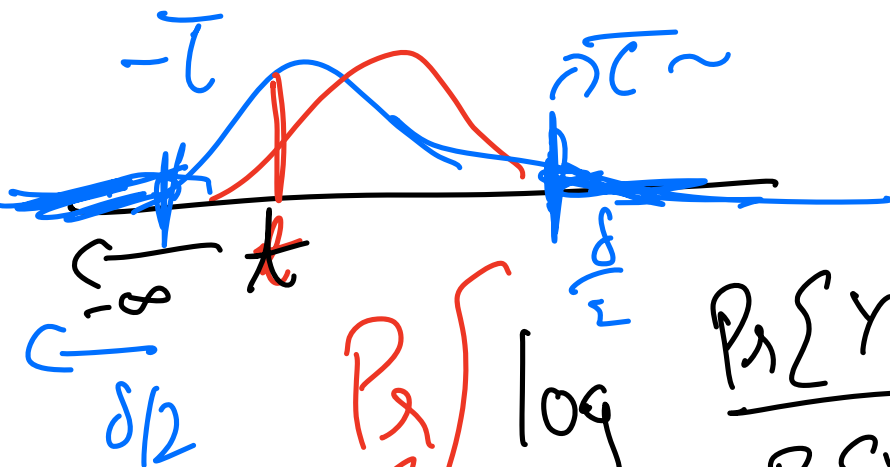
$$A(D) = 0 \qquad A(D') = 1$$

$$\hat{Y} = Y + N(0, \sigma^2)$$

$$\frac{P_A\{Y=y \mid H0\}}{P_A\{Y=y \mid H1\}} = \frac{N(0, \sigma^2)}{N(1, \sigma^2)}$$

$$= \frac{\exp\left(-\frac{1}{2} y^2/\sigma^2\right)}{\exp\left(-\frac{1}{2}(y-1)^2/\sigma^2\right)}$$

$$= \exp\left(\frac{1}{2\sigma^2}\left((y-1)^2 - y^2\right)\right)$$

$$= \exp\left(\frac{1}{2\sigma^2}\left(1 - 2y\right)\right) \qquad \varepsilon\left|f_i^2(D) - f_i(D')^2\right|$$

$$\underbrace{\phantom{=======}}_{\leq \varepsilon} \qquad \leq \varepsilon\left(f_i(D) - f_i(D')\right)$$



$$t \sim A(D)$$

$$P_A\left[\log \frac{P_A\{Y=t \mid H0\}}{P_A\{Y=t \mid H1\}} \geq \varepsilon\right] \leq \delta$$

$$\exp\left( \frac{1}{2\sigma^2} \left( 1 + 2\tau_\delta \right) \right)$$

$$\leq e^{\varepsilon}$$

$$\sigma^2 = \frac{1 + 2\tau_\delta}{2\varepsilon}$$

$$\Rightarrow (\varepsilon, \delta)\text{-DP}$$

gaussian mechanism

$$\|f(D) - f(D')\|_2 \le \Delta_2 \qquad \forall D, D'$$

$$Y_i = f_i(D) + \text{Lap}\left(\frac{\sqrt{d}\,\Delta_2}{\varepsilon}\right)$$

$$Y_i = f_i(D) + N\left(0, \frac{\Delta_2}{\varepsilon}\right)$$

independent of d

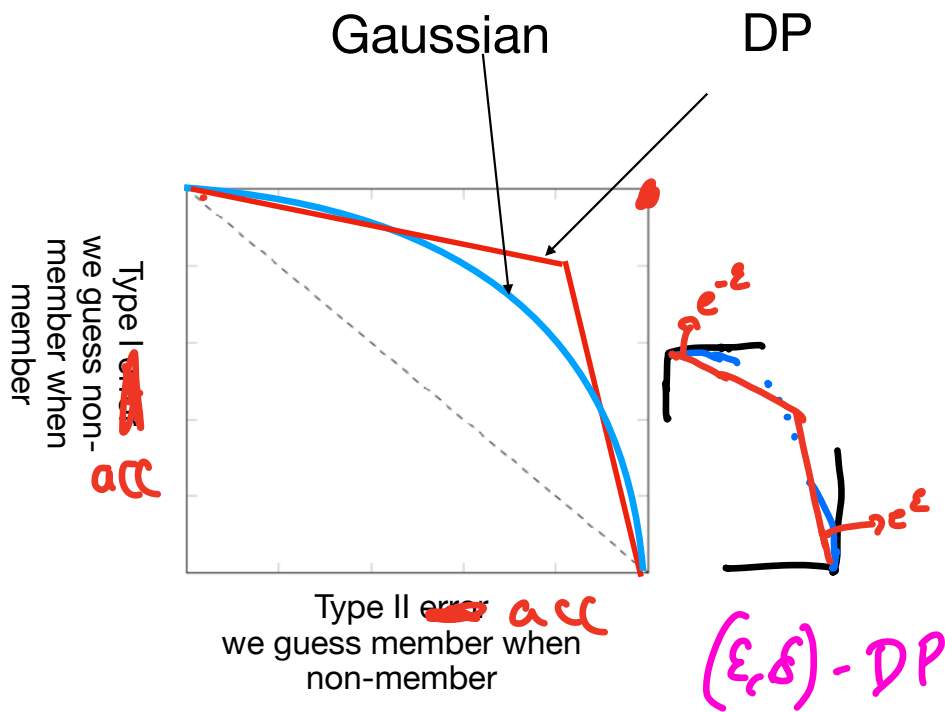# Differentially Private Algorithms
## Gaussian mechanism

- Suppose A(D) = 0, A(D') = 1.

- Release $\hat{y} = y + \text{Gaussian}(0, \varepsilon^{-1})$

- $z \sim \text{Gaussian}(\mu, \sigma^2) \Rightarrow p(z) \propto \frac{1}{\sigma} e^{-\frac{1}{2}\left(\frac{z-\mu}{\sigma}\right)^2}$

- $Pr[\hat{y} \,|\, y = 0] = \text{Gaussian}(0, \varepsilon^{-1})$ and $Pr[\hat{y} \,|\, y = 1] = \text{Gaussian}(1, \varepsilon^{-1})$

- $\dfrac{Pr[A(D) = y]}{Pr[A(D') = y]} = ?$ What happens at the tails?

$(\varepsilon, \delta)$ - Approx.

D. P

$(\varepsilon, \delta)$ -DP

# Differentially Private Algorithms
## Visualizing tradeoff curve of DP and Gaussian mechanism

Gaussian          DP



Type I error
we guess non-
member when
member

*acc*

Type II error *acc*
we guess member when
non-member

$(\varepsilon, \delta) - DP$

- At the edges, the slope of gaussian mechanism is vertical

- Impossible to get DP guarantee for any value of $\varepsilon$
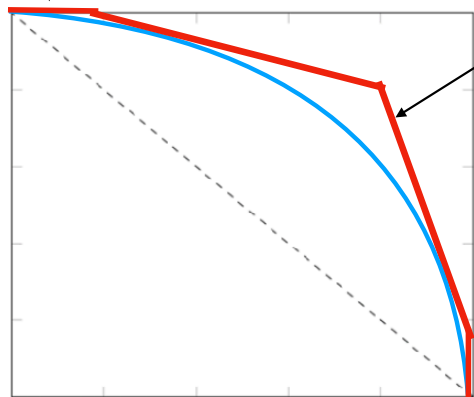
- Does this mean Gaussian mechanism is not private?

# Differentially Private Algorithms
## Approximate DP

Horizontal line of size $\delta$

Approximate $(\varepsilon, \delta)$-DP

Type I error
we guess non-member when member

Type II error
we guess member when non-member

Vertical line of size $\delta$

- Add flat lines of length $\delta$ at the edges to make some space for Gaussian mechanism

- Now chance for Gaussian mechanism to show privacy!

# Differentially Private Algorithms
## Approximate Differential Privacy
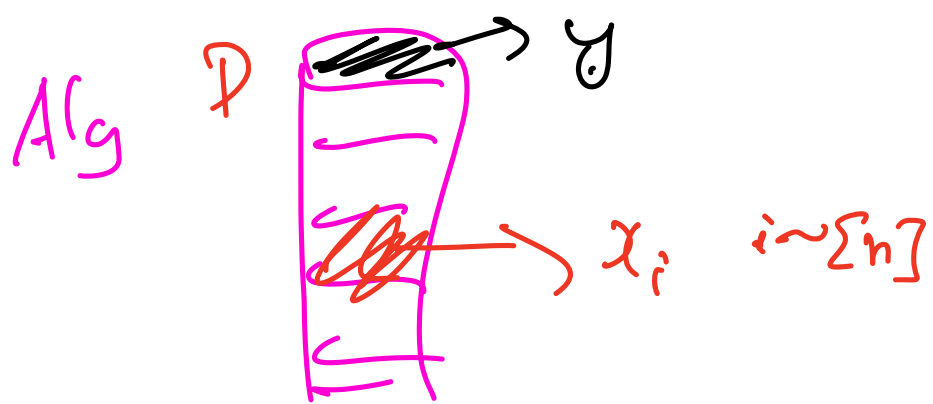
$(\varepsilon, \delta)$-Differential Privacy:

Let us draw a variable $t \sim A(D)$. Then the privacy loss random variable:

$$\mathscr{L}_{D,D'} := \ln \left( \frac{Pr[A(D) = t]}{Pr[A(D') = t]} \right)$$

A satisfies $(\varepsilon, \delta)$-DP iff for any neighboring datasets $D, D' \in \chi^n$ we have

$$Pr\left[ \mathscr{L}_{D,D'} \geq \varepsilon \right] \leq \delta$$

- With $\delta$ probability, arbitrarily bad things can happen.

- Ideally $\delta$ is chosen very small $\delta \leq n^{-1}$, or more common in fixed to $10^{-5}$.

Alg $\mathcal{D}$ $\longrightarrow$ Y

$x_i$ $i \sim [n]$

$$\Pr_A \left[ \log \frac{\Pr_A[Y=y \mid \mathcal{D}]}{\Pr_A[Y=y \mid \mathcal{D} \setminus \{x_i\}]} > 0 \right] \leq \frac{1}{n}$$

$\varepsilon = 0, \delta = \frac{1}{n}$

$(0, \frac{1}{n}) - DP$

# Differentially Private Algorithms
**Gaussian mechanism**

- Suppose A(D) = 0, A(D') = 1. Release $\hat{y} = y + \text{Gaussian}(0, \varepsilon^{-1})$

- $z \sim \text{Gaussian}(\mu, \sigma^2) \Rightarrow p(z) \propto \dfrac{1}{\sigma} e^{-\frac{1}{2}\left(\frac{z-\mu}{\sigma}\right)^2}$

- $Pr[\hat{y} \,|\, y = 0] = \text{Gaussian}(0, \varepsilon^{-1})$ and $Pr[\hat{y} \,|\, y = 1] = \text{Gaussian}(1, \varepsilon^{-1})$

- $\dfrac{Pr[A(D) = y]}{Pr[A(D') = y]} = ?$ what happens now?