

CSCI 699: Privacy Preserving Machine Learning - Week 2

Differential Privacy

Sai Praneeth Karimireddy, Sep 6 2024

Quantifying Privacy Leakage



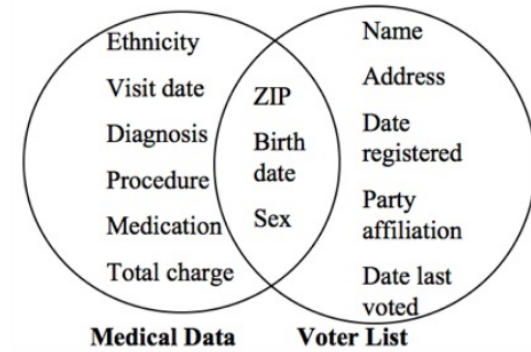
Last week recap

- We saw many definitions of privacy
 - De-identification / suppression
 - K-anonymity
 - L-diversity
- We saw none of them really protected privacy and were easily broken
- Hinted at a more widely accepted definition.

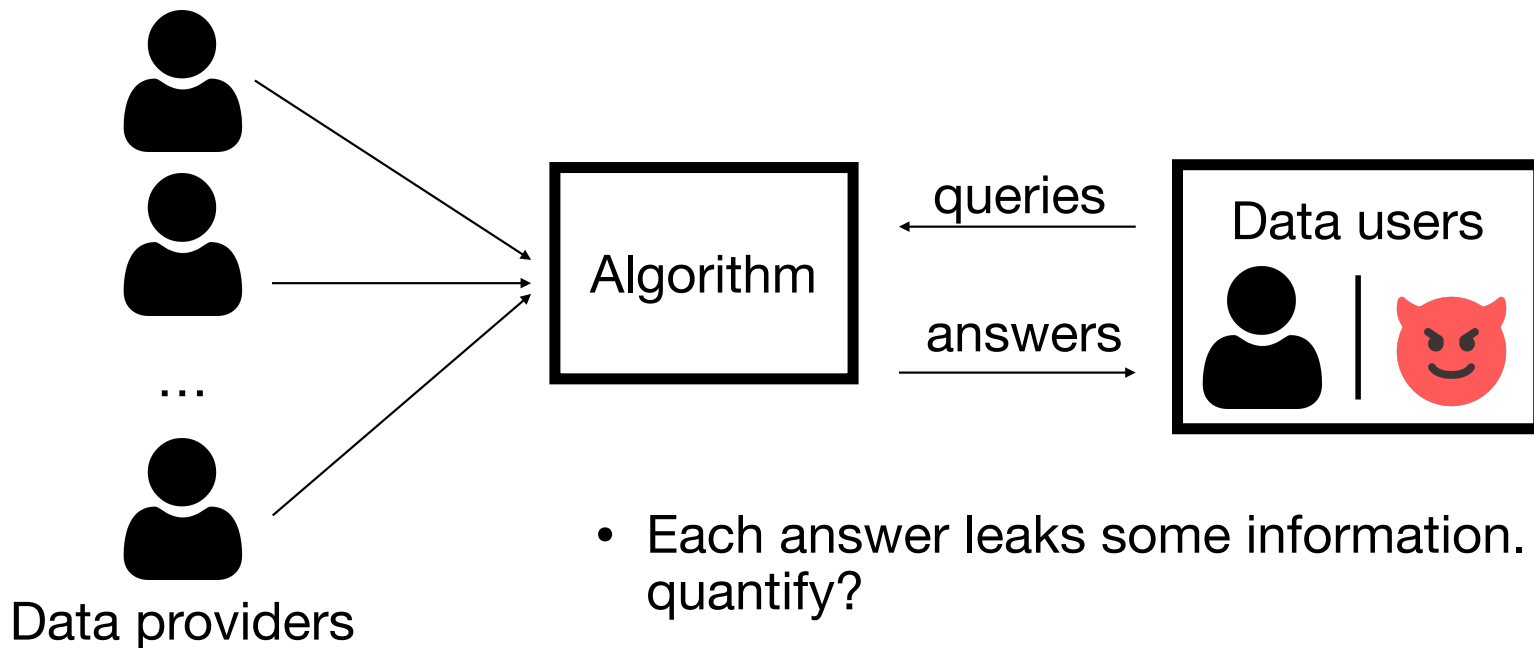
Last week takeaways

Requirements for privacy definition

- **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.



Goals of PPML



- Each answer leaks some information. How to quantify?
- How to balance usefulness of answers vs. privacy being leaked?

Quantifying Privacy Leakage

Attempt 1

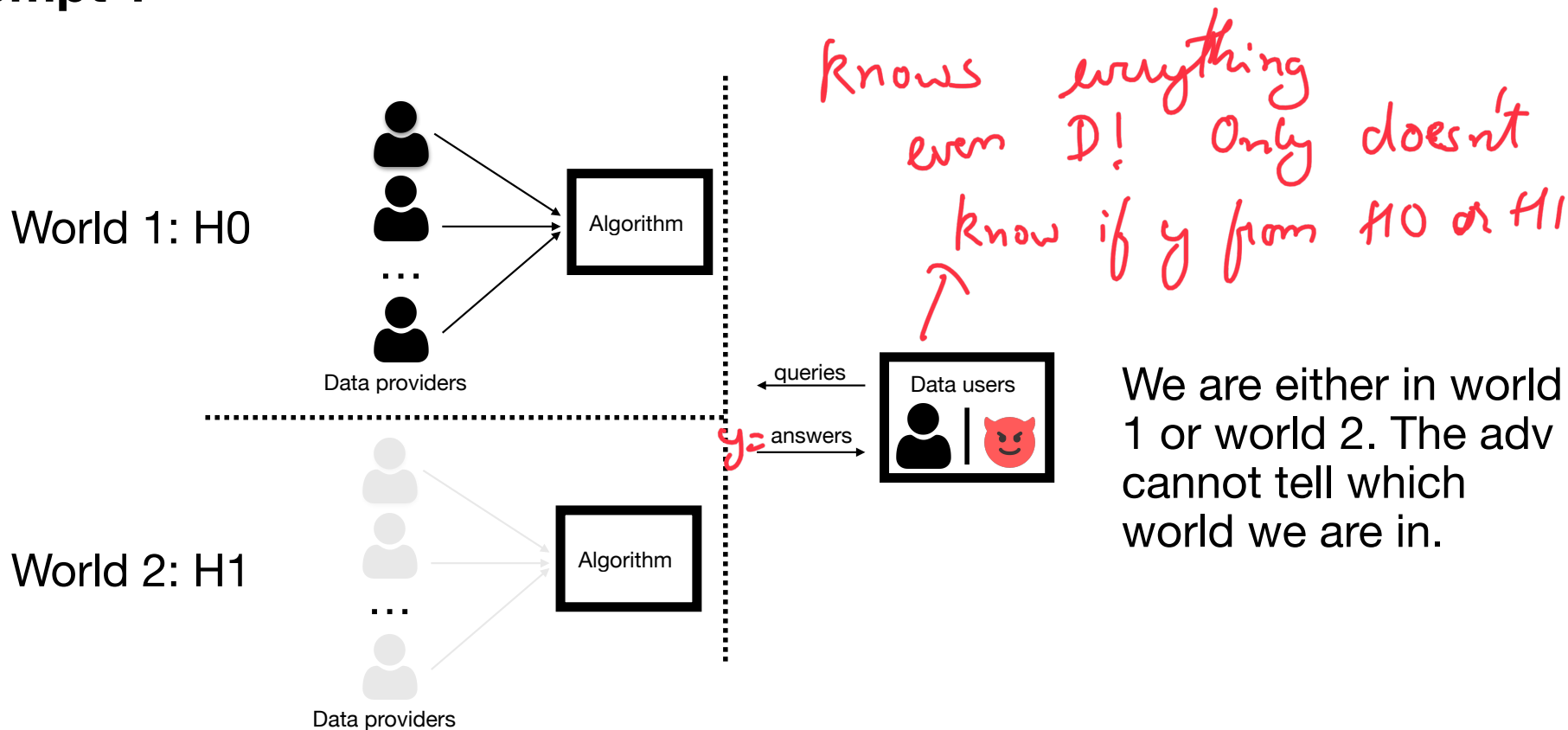
Absolute Privacy: quantify **total** information leaked

“An answer to a query is private if the response reveals no more than was already known about the individuals in the data”

- Bayesian version: the posterior and prior are identical

Quantifying Privacy Leakage

Attempt 1



Quantifying Privacy Leakage

Attempt 1

Absolute Privacy: quantify **total** information leaked

“An answer to a query is private if the response **reveals no more than was already known** about the individuals in the data”

- **Problem 1:** **Impossible to reveal anything** useful about data since any useful answer will provide some previously unknown information.

Quantifying Privacy Leakage

Attempt 1: Problems

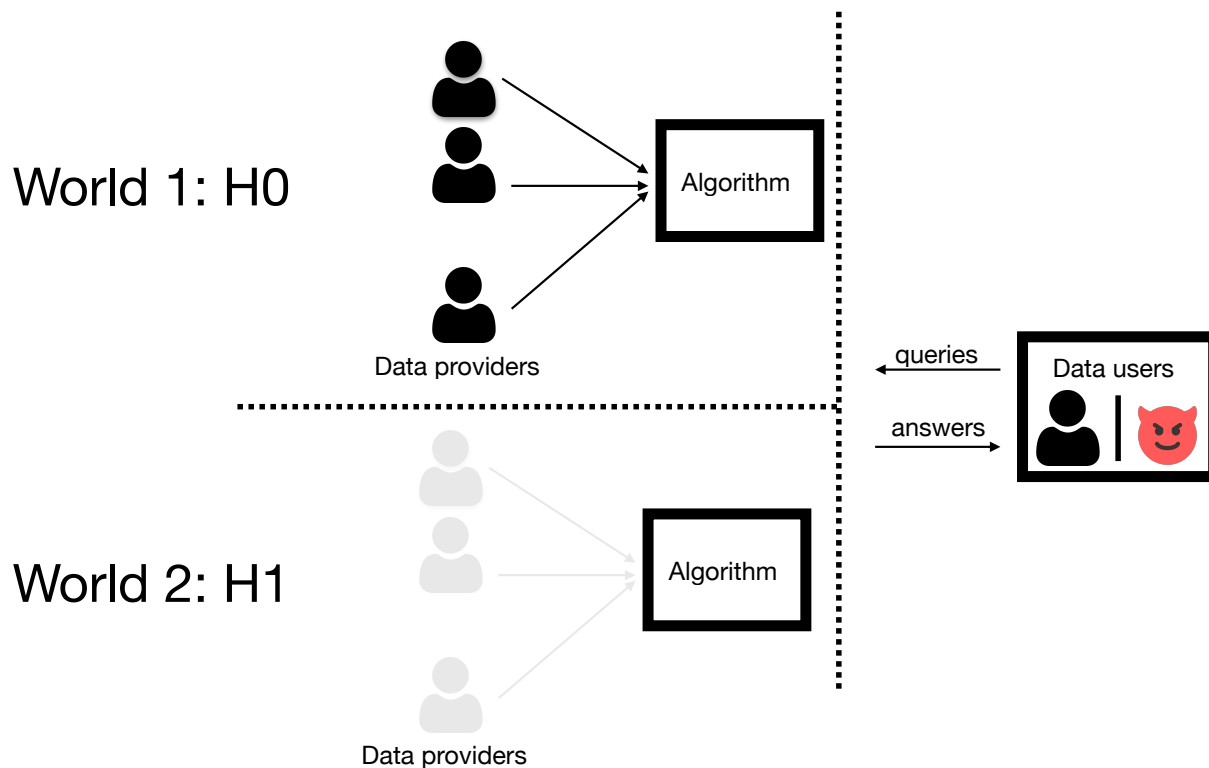
Absolute Privacy: quantify **total** information leaked

“An answer to a query is private if the response reveals no more than was already known about the individuals in the data”

- **Problem 2:** What I know before changes with auxiliary information.
- Did the model leak information about Bob?
 - Bob is a smoker, but his data was not used to train the model.
 - The model said smokers have higher risk of disease.
 - Bob’s insurance premiums were raised.

Quantifying Privacy Leakage

Attempt 1: Problems



Any information about the distribution reveals which world we are in.

Quantifying Privacy Leakage

Attempt 1: Problems

Absolute Privacy: quantify **total** information leaked

“An answer to a query is private if the response **reveals no more than was already known** about the individuals in the data”

- **Problem 2:** What I know before **changes with auxiliary information**.
- We want to safeguard individual information (**privacy**) while revealing distributional/aggregate information (**utility**)

Quantifying Privacy Leakage

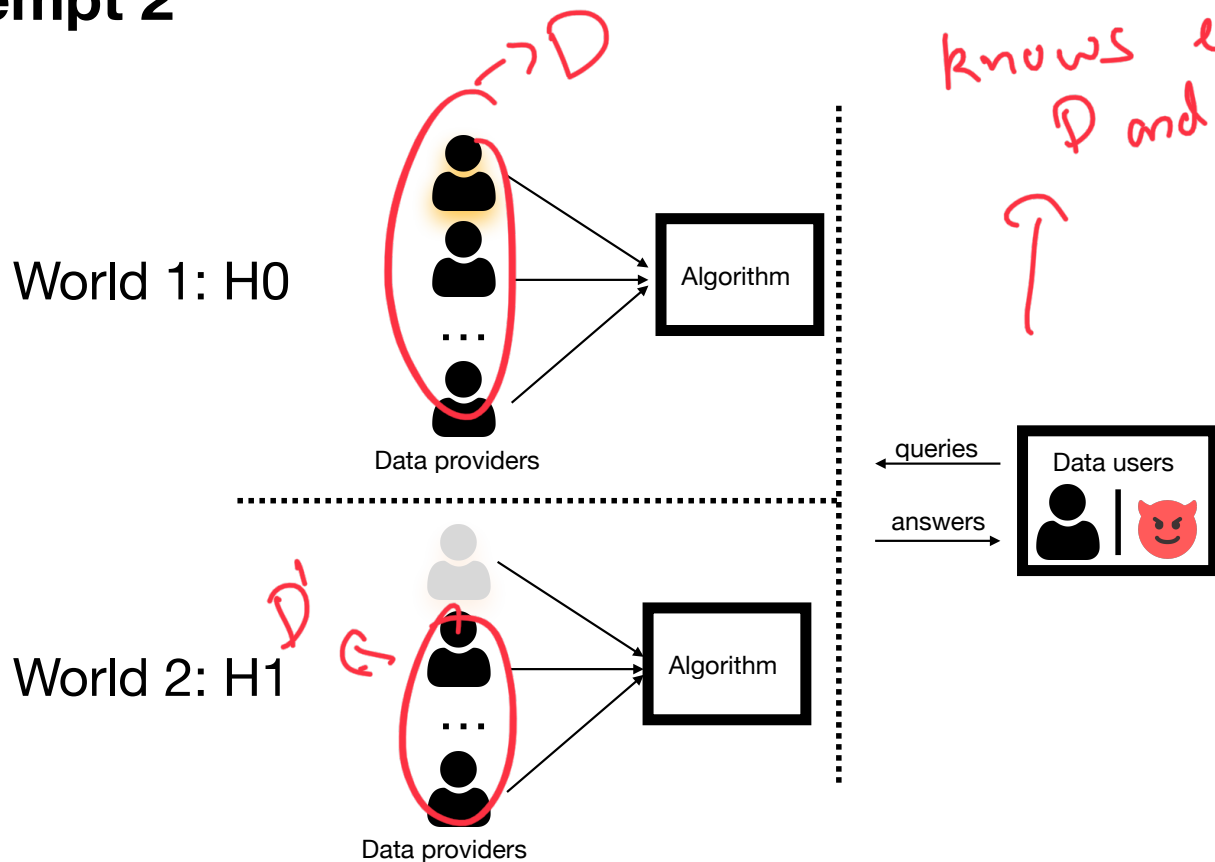
Attempt 2

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”

Quantifying Privacy Leakage

Attempt 2



knows everything including
D and D. Only doesn't
know H0 or H1.

- In world 2 only Bob is removed/replaced.
- Now from the answer, how easily can guess the correct world?

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.

Aside: how is Putin's popularity calculated?

Plausible deniability as privacy

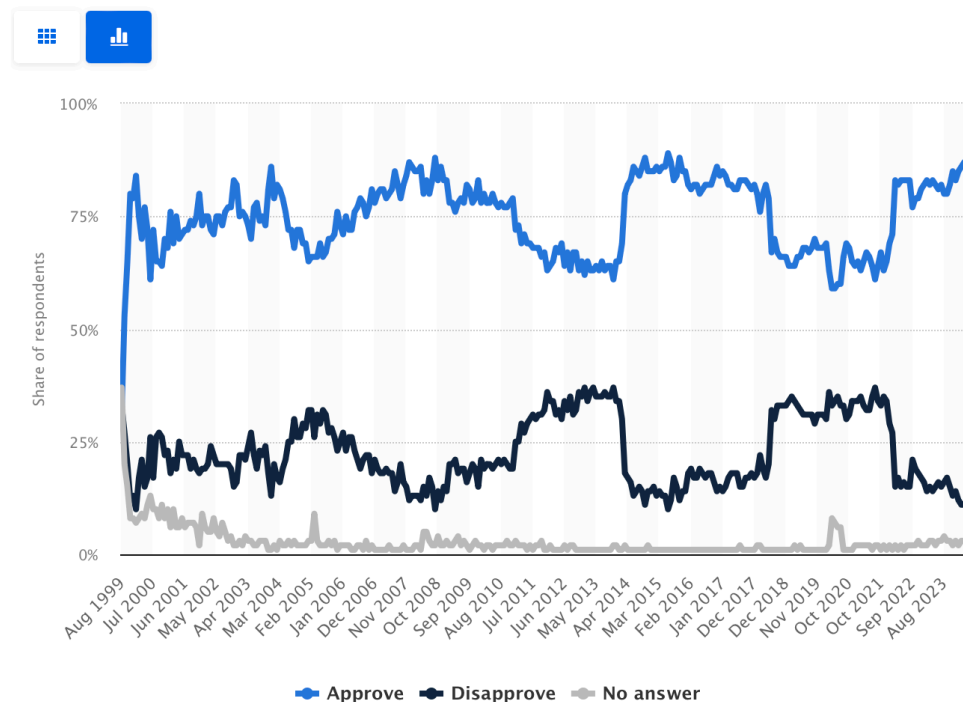
Poll: Russians Still Like Putin and Back the Ukraine War – but Are Anxious at Home

Most Russian survey respondents see the war in Ukraine as a broader conflict with the West and support it amid concerns about their own country's economy.

By Elliott Davis Jr. | Jan. 9, 2024



Do you approve of the activities of Vladimir Putin as the president (prime minister) of Russia?



Aside: how is Putin's popularity calculated?

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 confer plausible deniability?

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

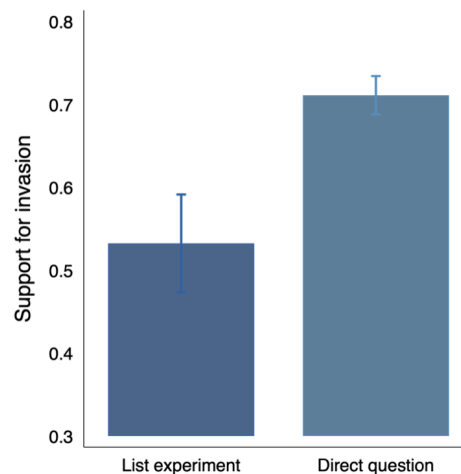
2

3 of these things

Aside: how is Putin's popularity calculated?

List Experiment

Figure 2: Support for the Russian invasion of Ukraine



Note: Bars show averages, vertical lines show 95% confidence intervals.

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”

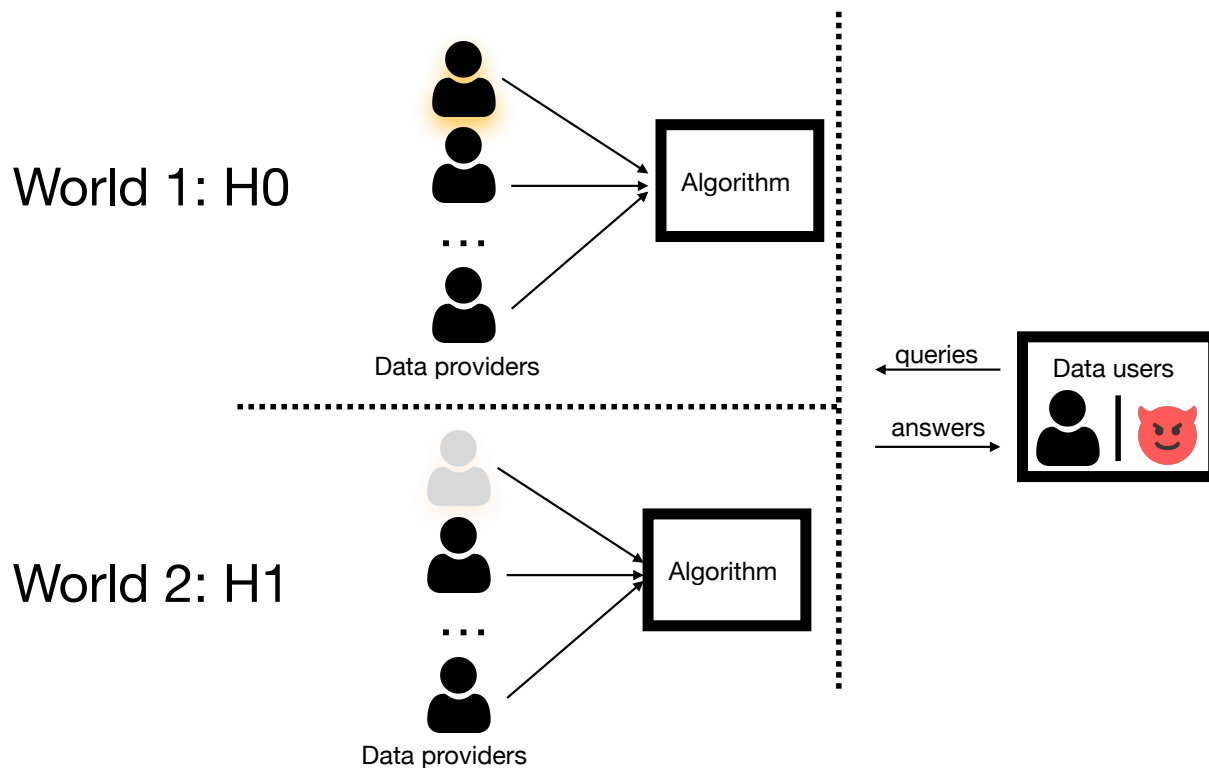
- **Question:** Can a deterministic algorithm be private?
- What if Bob is the only data point? Then can easily reverse-engineer Bob's data.

$$\min_x \ell(f(x), y)$$

- Only randomized algorithms can be private.

Membership Inference

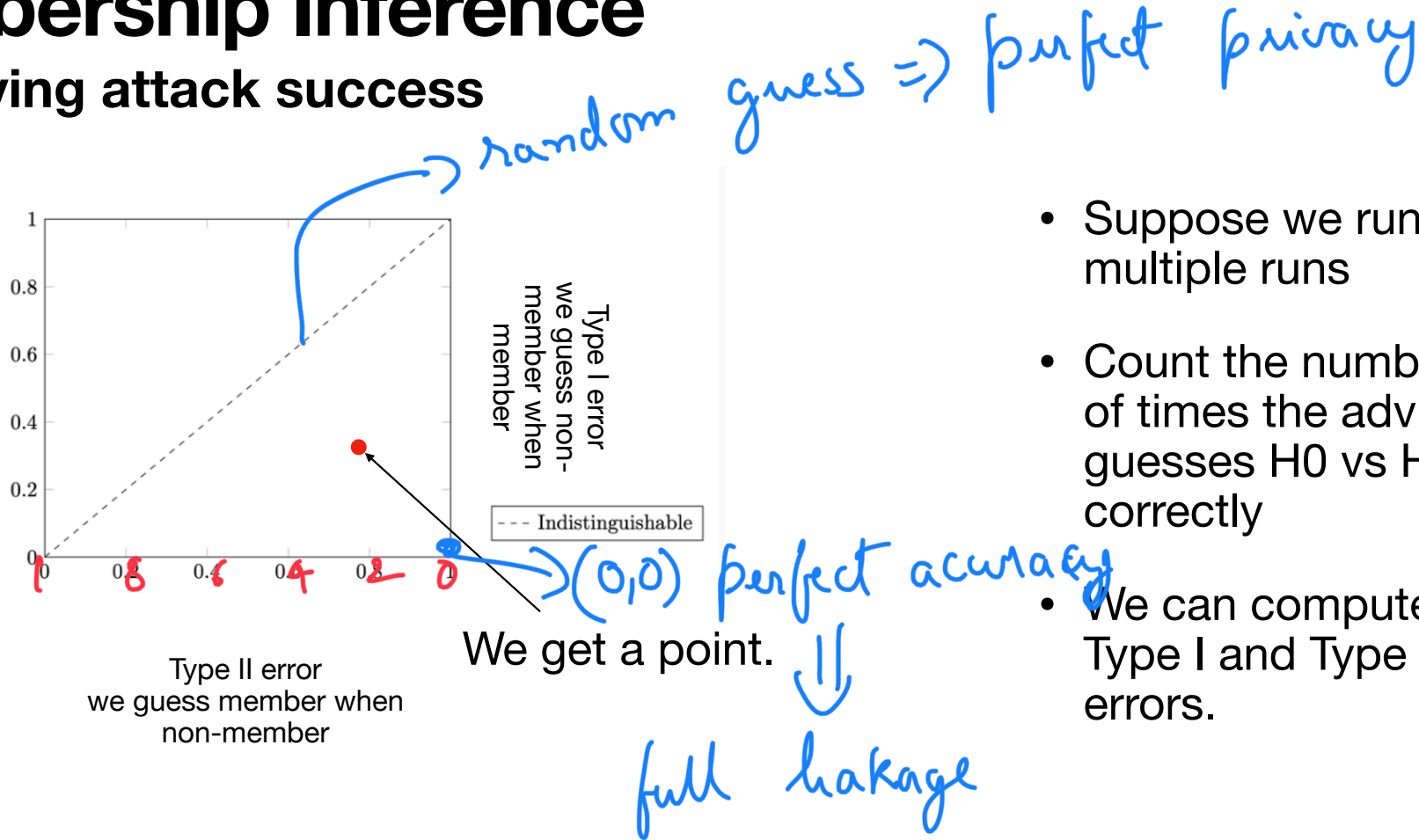
As a definition of privacy



- In world 2 only Bob is removed/replaced.
- Now from the answer, how easily can guess the correct world?
- Can have false positives, false negatives

Membership Inference

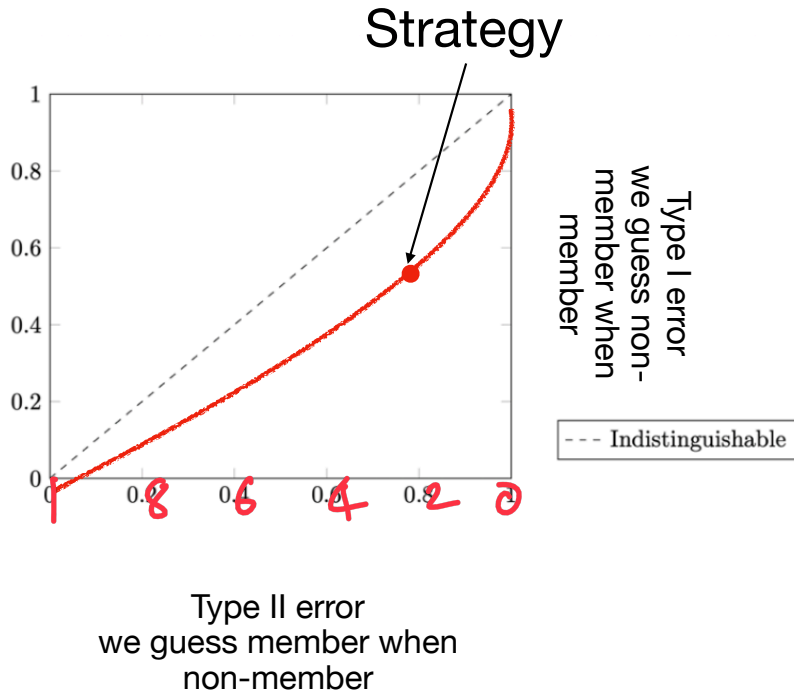
Quantifying attack success



- Suppose we run multiple runs
- Count the number of times the adv guesses H_0 vs H_1 correctly
- We can compute Type I and Type II errors.

Membership Inference

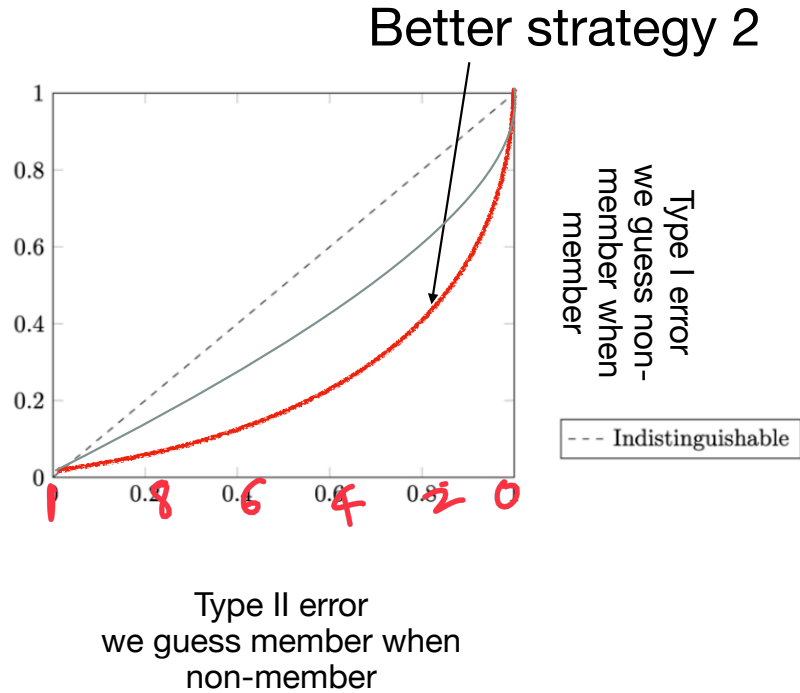
Tradeoff curve



- But sometimes we care asymmetrically
- E.g. its important not to miss anyone e.g. sending cat ads to pet owners
- Not ok if we are accusing them of a crime

Membership Inference

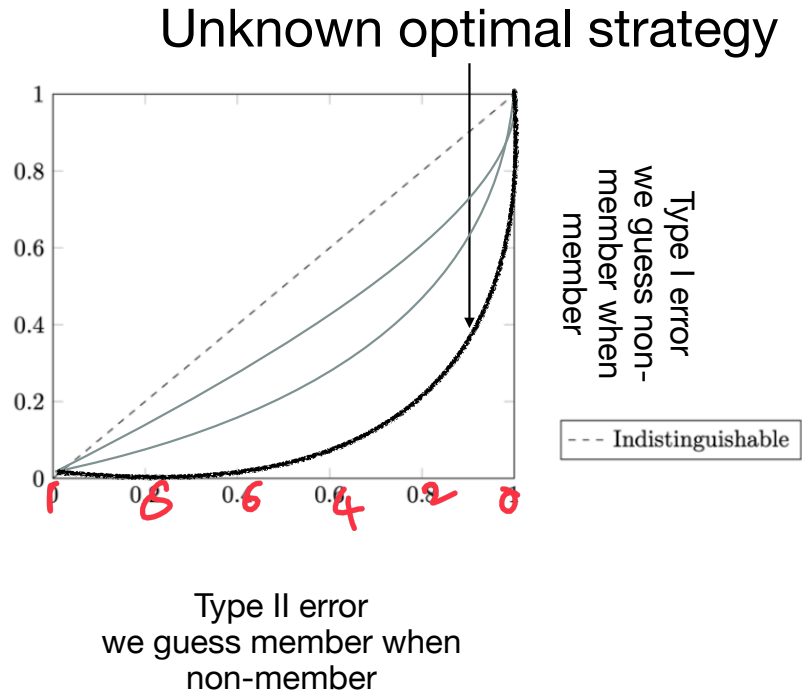
Comparing tradeoff curves



- Strategy 1 is better than Strategy 2 if the curve is uniformly above.
- Lower curve means we've found more privacy leakage

Membership Inference

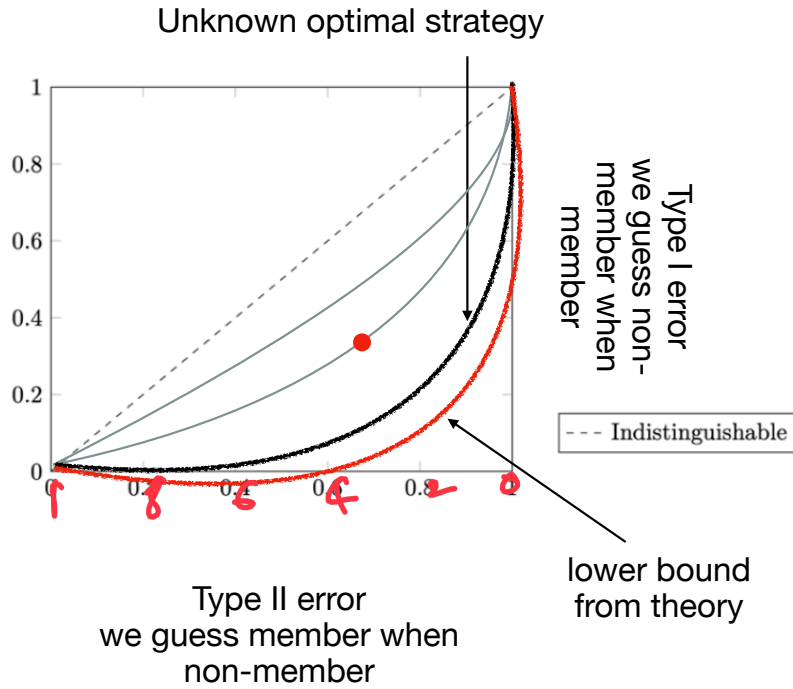
Optimal tradeoff curve



- 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.

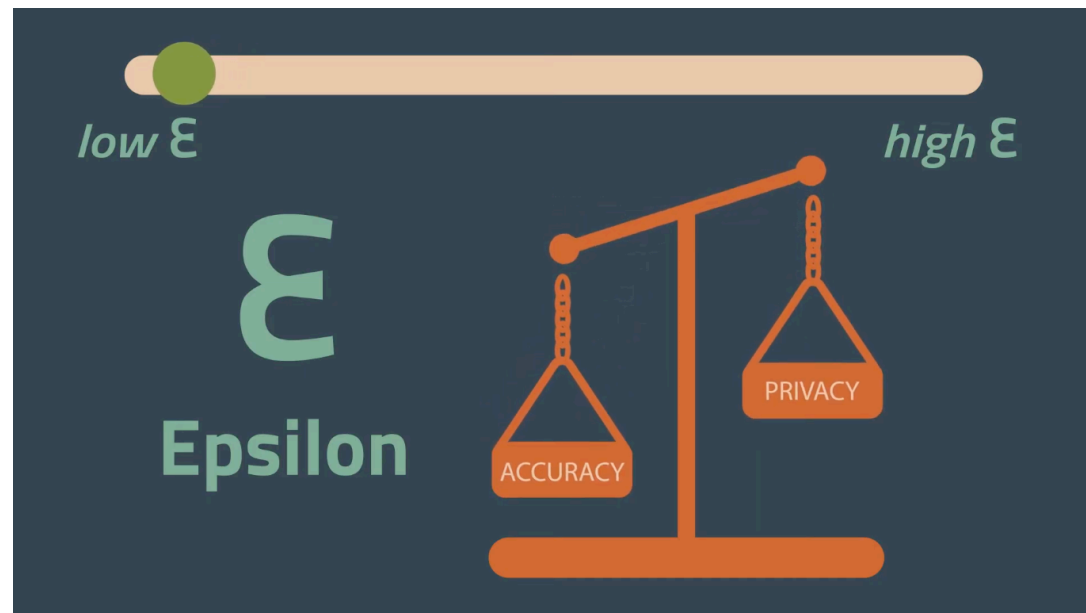
Membership Inference

Privacy from tradeoff curve



- Use optimal strategy to quantify privacy
- But empirical tests only give an upper-bound
- Need theory to give lower-bound

Differential Privacy



Differential Privacy

Calibrating Noise to Sensitivity in Private Data Analysis

2006

Cynthia Dwork¹, Frank McSherry¹, Kobbi Nissim², 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 χ be a the domain of training data
- A dataset $D \in \chi^n$ is a multiset of n records/rows of χ
- D (sensitive data) \longrightarrow algorithm $\longrightarrow Y$ (answers)
- Attacker wants to infer some information about $D \in \chi^n$
 - observes Y
 - knows algorithm, domain χ , 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 χ , and prior information.
 - can compute likelihood of dataset:

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

algorithm \swarrow prior knowledge \swarrow

Differential Privacy

Performing membership inference

- Attacker wants to infer presence of $x \in X$?
 - observes Y , knows algorithm, domain \mathcal{X} , and even $D \setminus x \in \mathcal{X}^{n-1}$
 - can compute likelihood of x in dataset

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

algorithm \swarrow prior knowledge \swarrow

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

- Attacker wants to infer some information about $D \in \mathcal{X}^n$
 - can compute likelihood of seeing some dataset

algorithm prior knowledge


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

- We design a private algorithm by controlling $Pr[Y | D]$

Differential Privacy

Strict definition

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

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

- 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

- Some indistinguishability: For all **similar inputs**, the **output probabilities are bounded**.

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

- It means by observing any Y , adversary is NOT able to distinguish between inputs x and x' beyond a bounded certainty.
- What does **similar inputs** mean?
 - Depends on use case
 - location positions that are within some range
 - datasets that differ in one individual row

Differential Privacy

Formal definition

ϵ -Differential Privacy:

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

$$\frac{\Pr[Y = y | D]}{\Pr[Y = y | D']} \leq \exp(\epsilon)$$

- Recall that D (sensitive data) \longrightarrow algorithm $\longrightarrow Y$ (answers)
- So we have, $\Pr[Y | D] = \Pr[A(D) = Y]$

Differential Privacy

Formal definition

ϵ -Differential Privacy:

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

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

- $\epsilon = 0$ means perfect privacy
- $\epsilon \gg 0$ means not private

Differential Privacy

Source of randomness

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

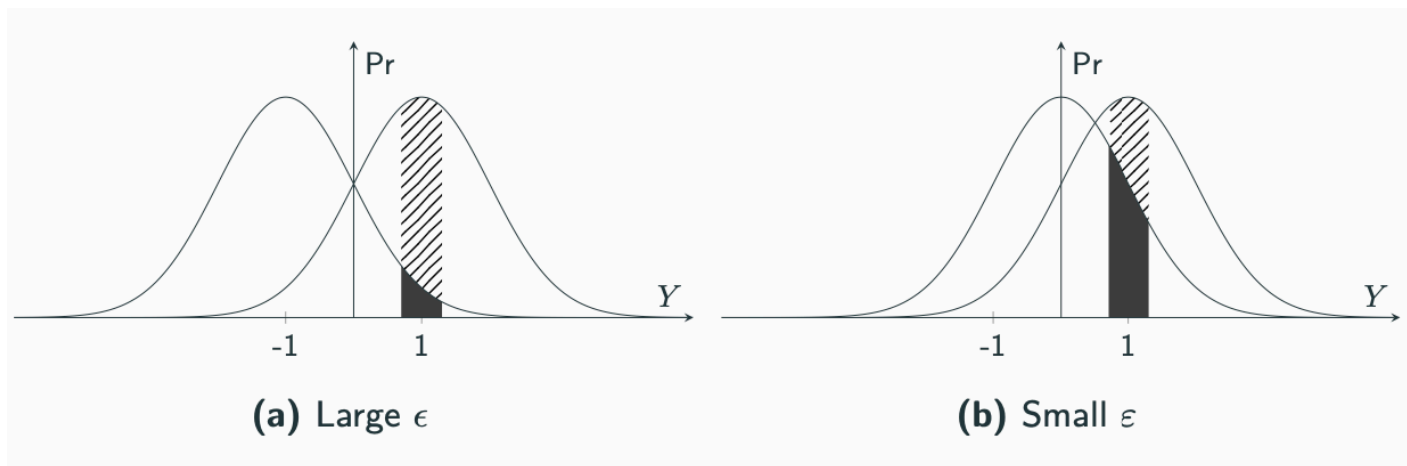
- In $\Pr[A(D) = y]$, over what randomness is the probability defined?
 - The randomness of the algorithm?
 - Yes
 - Randomness of the data $D \in \mathcal{X}^n$?
 - No.
 - We look at all possible values of D, D' i.e. worst case

Differential Privacy

Visual representation

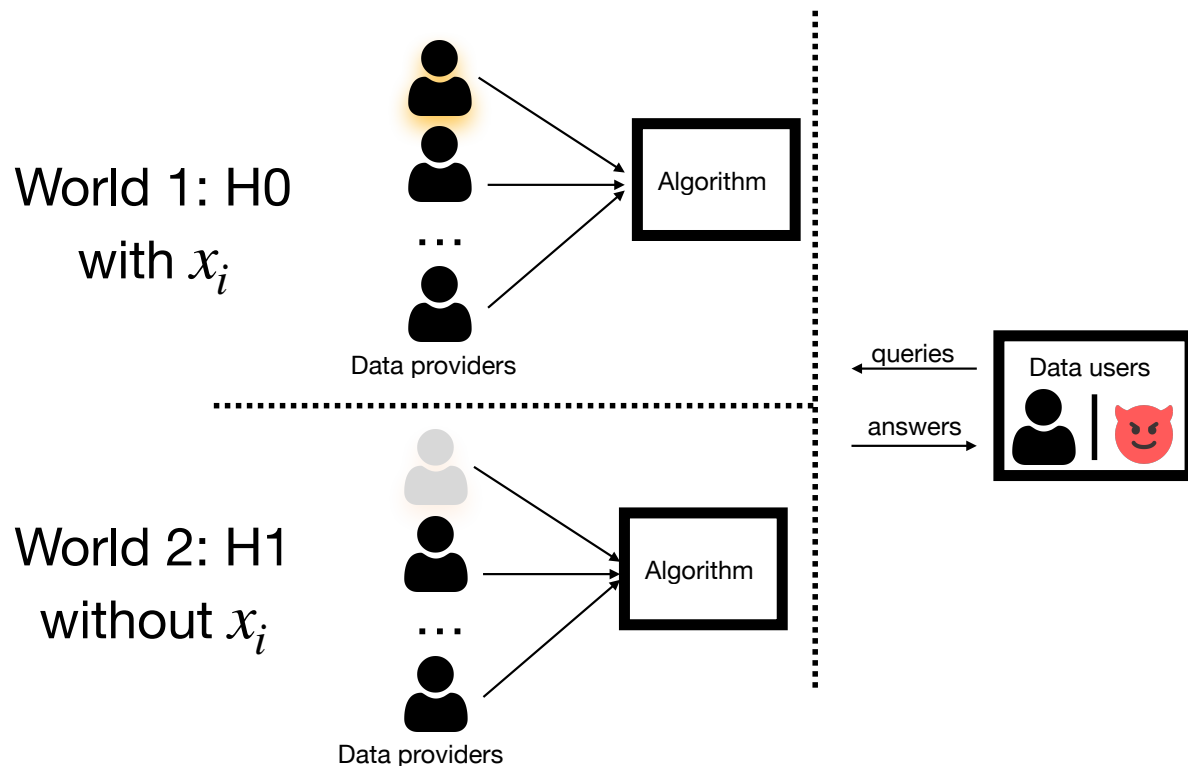
- Consider $D = \langle x_1, \dots, x_i, \dots, x_n \rangle$, and a similar dataset $D' = \langle x_1, \dots, x'_i, \dots, x_n \rangle$

- ϵ -DP means $\frac{\Pr[A(D) = y]}{\Pr[A(D') = y]} \leq \exp(\epsilon)$



Differential Privacy

Recall Membership Inference

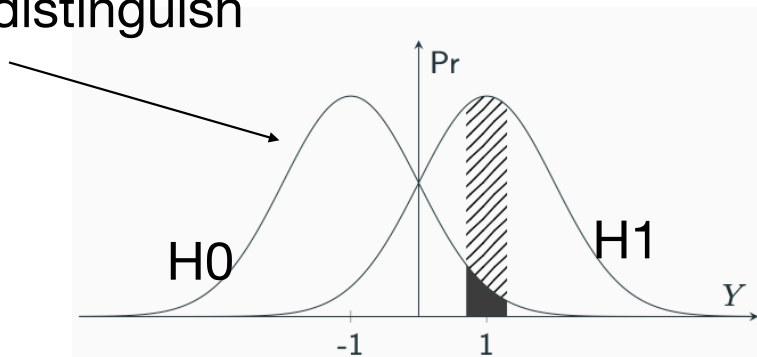


- We know everything about the algorithm and even $D \setminus x_i$
- We observe an output Y
- Need to guess if it came from H_0 or H_1

Differential Privacy

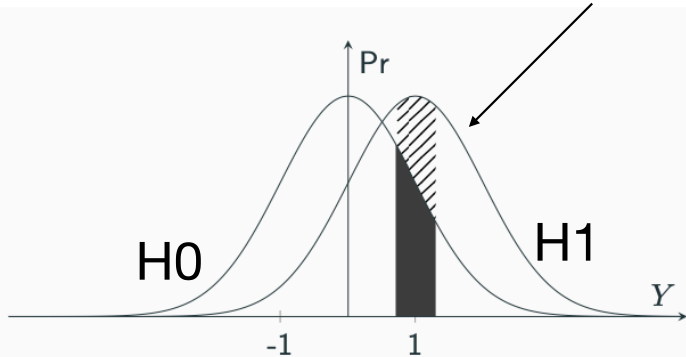
Connection to Membership Inference

Easy to distinguish



(a) Large ϵ

Hard to distinguish



(b) Small ϵ

- We observe $Y = 1$.
- Can you guess H_0 or H_1 ?

Small $\epsilon \Rightarrow$ higher error

Differential Privacy and membership inference

Quantifying connection

Can trade some Type I error for Type II, but sum of errors is lower-bounded.

Theorem

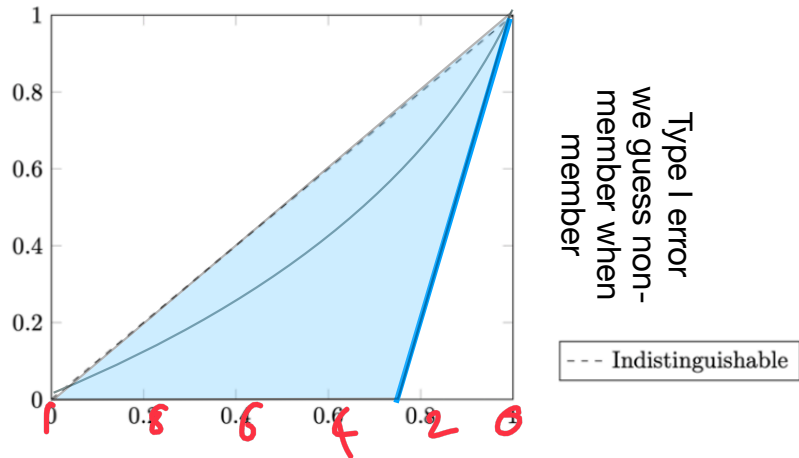
Suppose A satisfies ϵ -DP for datasets D, D' which differ by one datapoint. Then, we have

- $Pr[\text{guess } H_0 \mid H_1] + e^\epsilon Pr[\text{guess } H_1 \mid H_0] \geq 1$
- $e^\epsilon Pr[\text{guess } H_0 \mid H_1] + Pr[\text{guess } H_1 \mid H_0] \geq 1$

- Type I error = $Pr[\text{guess } H_0 \mid H_1]$
- Type II error = $Pr[\text{guess } H_1 \mid H_0]$

Differential Privacy and membership inference

Visualizing connection



Type II error
we guess member when
non-member

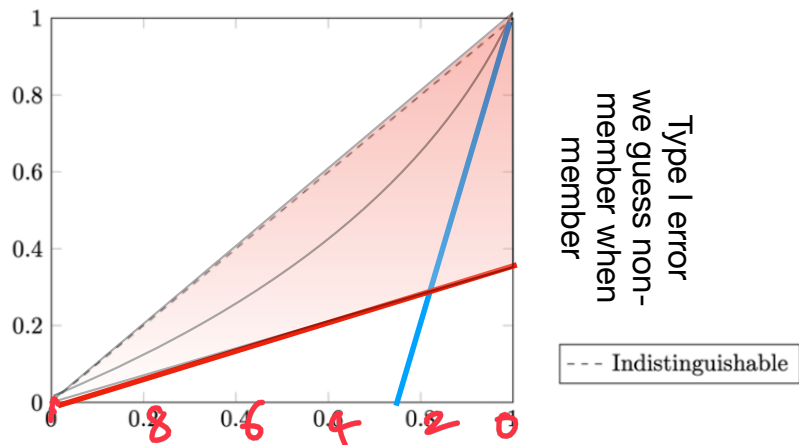
Type I error
we guess non-
member when
member

--- Indistinguishable

- $Pr[\text{guess } H0 | H1] + e^\epsilon Pr[\text{guess } H1 | H0] \geq 1$
- gives us blue line with slope e^ϵ

Differential Privacy and membership inference

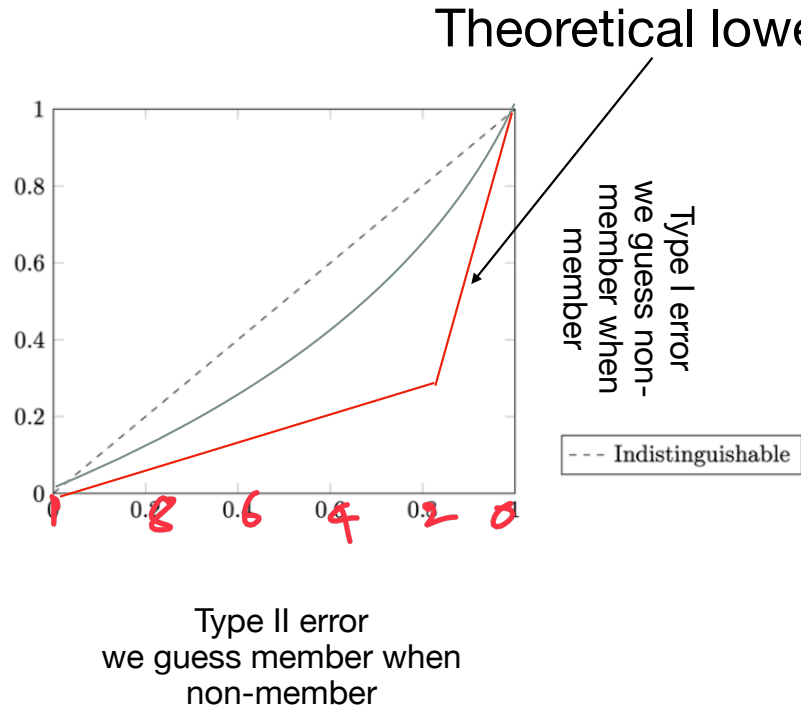
Visualizing connection



- $e^{\epsilon} Pr[\text{guess } H0 \mid H1] + Pr[\text{guess } H1 \mid H0] \geq 1$
- gives the red line with slope $e^{-\epsilon}$

Differential Privacy and membership inference

Visualizing tradeoff curve of DP



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

Algorithms for Differential Privacy

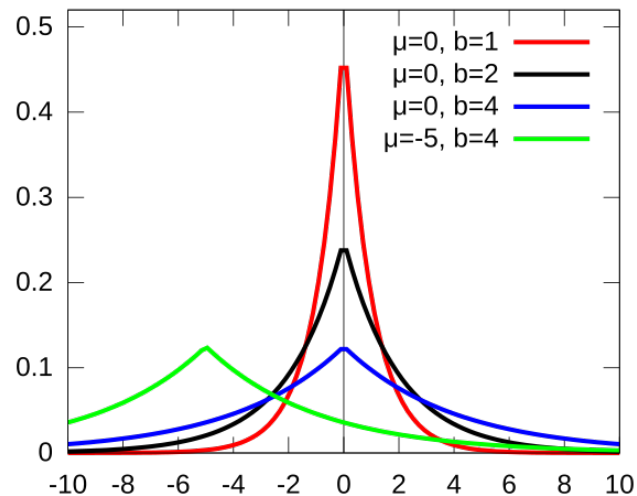


Differentially Private Algorithms

Just add Laplace noise

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

- Suppose $A(D) = 0$, $A(D') = 1$.
- Release $\hat{y} = y + \text{Laplace}(0, \epsilon^{-1})$
- $z \sim \text{Laplace}(\mu, b) \Rightarrow p(z) = \frac{1}{2b} e^{-\frac{|z-\mu|}{b}}$



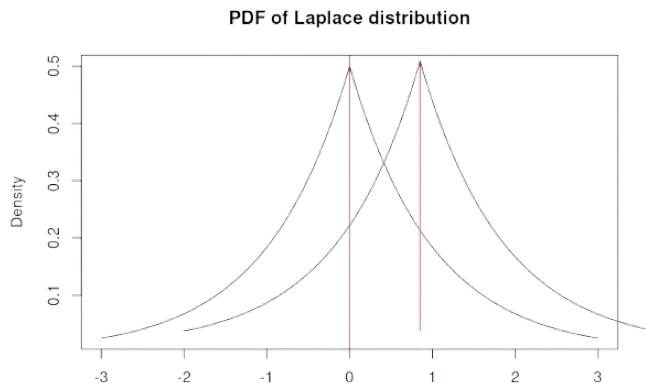
Differentially Private Algorithms

Just add Laplace noise

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

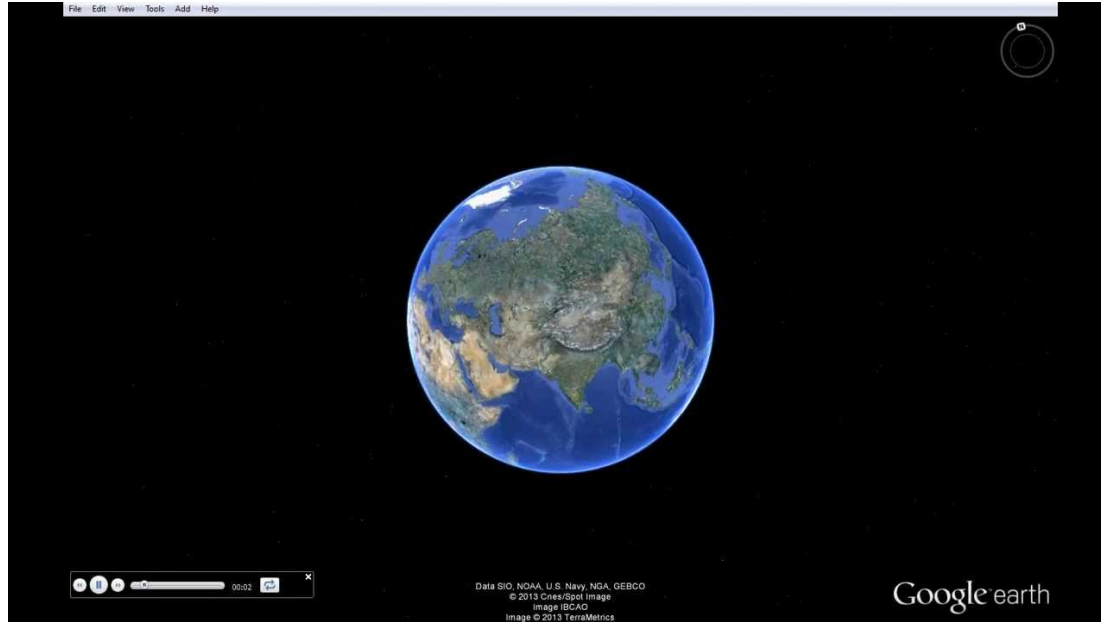
- Suppose $A(D) = 0$, $A(D') = 1$. Release $\hat{y} = y + \text{Laplace}(0, \varepsilon^{-1})$
- $\Pr[\hat{y} | y = 0] = \text{Laplace}(0, \varepsilon^{-1})$ and $\Pr[\hat{y} | y = 1] = \text{Laplace}(1, \varepsilon^{-1})$

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



Differentially Private Algorithms

Sensitivity



- I release average income at different zoom levels. Added $\text{Lap}(0,1)$.
- Do they all leak same amount of privacy?

Differentially Private Algorithms

Sensitivity and Laplace mechanism

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

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

$y \in \mathbb{R}^d$

$f(D) \in \mathbb{R}^d$

Theorem

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

$$[A(D)]_i = [f(D)]_i + \text{Laplace}(0, \Delta \epsilon^{-1}) \quad \forall i \in d$$

Proof: $P_A[A(D)] = y = P_A[f(D)] + \text{Lop}(0, \Delta \epsilon^{-1}) = y$

+ tied

$$= \frac{\epsilon}{2\Delta} \prod_{i=1}^d \exp\left(\frac{-\epsilon}{\Delta} |y_i - [f(D)]_i|\right)$$

$$\Rightarrow \frac{P_A[A(D) = y]}{P_A[A(D') = y]} = \frac{\prod_{i=1}^d \exp\left(\frac{-\epsilon}{\Delta} |y_i - [f(D)]_i|\right)}{\prod_{i=1}^d \exp\left(\frac{-\epsilon}{\Delta} |y_i - [f(D')]_i|\right)}$$

$$= \exp\left(\frac{\epsilon}{\Delta} \sum_{i=1}^d (|y_i - [f(D')]_i| - |y_i - [f(D)]_i|)\right)$$

$$(a) \leq \exp\left(\frac{\varepsilon}{\Delta} \sum_{i=1}^d |f(D')_i - [f(D)]_i|\right)$$

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

$$\leq \exp(\varepsilon)$$

(a) is because $|a| - |b| \leq |a - b|$ (triangle inequality)

(b) is because f is Δ -sensitive wrt $\|\cdot\|_1$ norm

\mathbb{R}

Differentially Private Algorithms

Sensitivity and Laplace mechanism

- **Definition: Sensitivity** of a function $f : (x_1, \dots, x_n) \mapsto (y_1, \dots, y_k)$ 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 Δ -sensitivity wrt $\|\cdot\|_\infty$
- What about Δ -sensitivity wrt $\|\cdot\|_2$ *Can convert to $\|\cdot\|_1$, then apply prev theorem.*
- Laplace mechanism is great for functions with small ℓ_1 sensitivity, not so much for small ℓ_2 sensitivity

Really bad if $d \gg 1$

$$\left\{ \begin{array}{l} \text{Use } \|v\|_1 \leq d \|v\|_\infty \\ \|v\|_2 \leq \|v\|_1 \leq \sqrt{d} \|v\|_2 \end{array} \right.$$

Differentially Private Algorithms

Gaussian mechanism

- Suppose $A(D) = 0$, $A(D') = 1$.
- Release $\hat{y} = y + \text{Gaussian}(0, \epsilon^{-1})$
- $z \sim \text{Gaussian}(\mu, \sigma^2) \Rightarrow p(z) \propto \frac{1}{\sigma} e^{-\frac{1}{2}(\frac{z-\mu}{\sigma})^2}$
- $Pr[\hat{y} | y = 0] = \text{Gaussian}(0, \epsilon^{-1})$ and $Pr[\hat{y} | y = 1] = \text{Gaussian}(1, \epsilon^{-1})$
- $\frac{Pr[A(D) = y]}{Pr[A(D') = y]} = ?$ What happens at the tails?

$$\frac{P_A[A(D) = y]}{P_A[A(D') = y]} = \exp\left(\frac{\epsilon}{2} (y-1)^2 - y^2\right)$$

$$P_A[A(D') = y] = \exp\left(\frac{\epsilon}{2} (1-2y)\right)$$

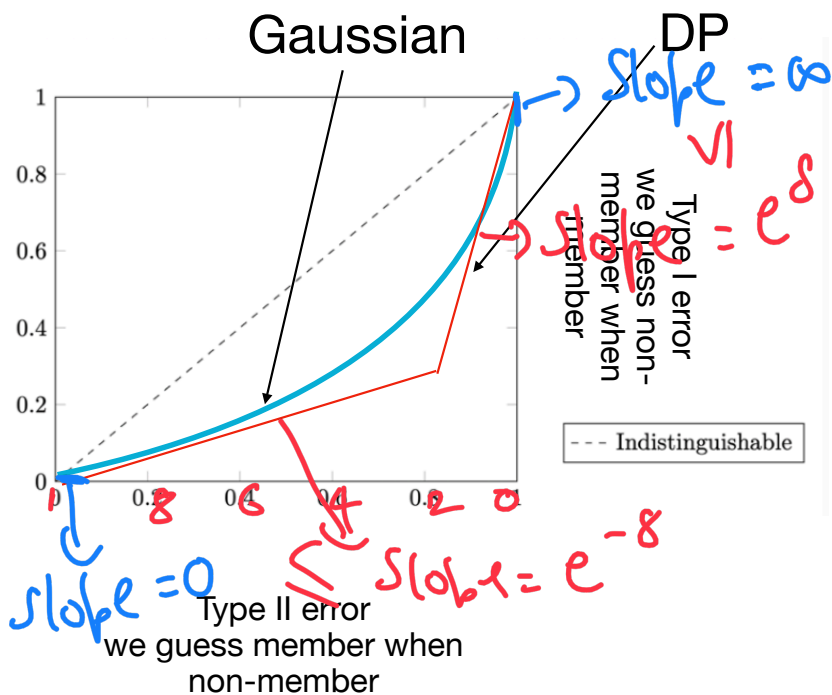
When $y \rightarrow -\infty$, \uparrow this blows up to ∞

\Rightarrow ϵ -DP impossible $\forall \epsilon!$

\Rightarrow No privacy?

Differentially Private Algorithms

Visualizing tradeoff curve of DP and Gaussian mechanism



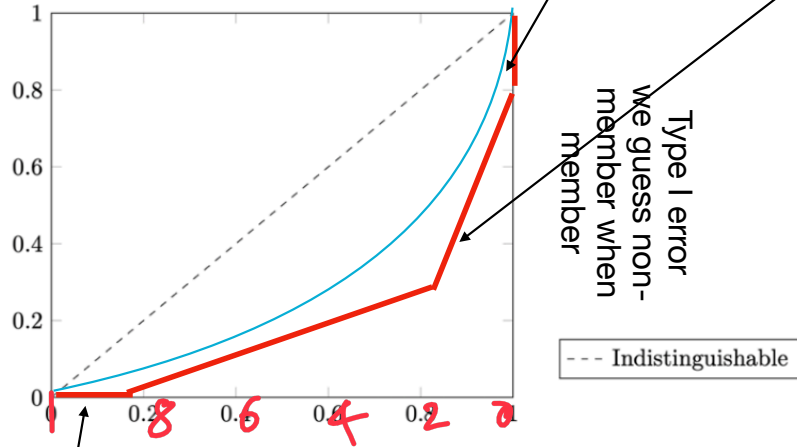
- At the edges, the slope of gaussian mechanism is vertical
- Impossible to get DP guarantee for any value of ϵ
- Does this mean Gaussian mechanism is not private?

Differentially Private Algorithms

Approximate DP

Vertical line of size δ

Approximate (ϵ, δ) -DP



Type I error
we guess non-
member when
member

Type II error
we guess member when
non-member

Horizontal line of size δ

- Add flat lines of length δ at the edges to make some space for Gaussian mechanism
- Now chance for Gaussian mechanism to show privacy!

Differentially Private Algorithms

Approximate Differential Privacy

(ϵ, δ) -Differential Privacy:

An algorithm A satisfies (ϵ, δ) -DP if for any **similar** datasets $D, D' \in \mathcal{X}^n$ and $y \in \mathcal{Y}$

$$\Pr[A(D) = y] \leq \Pr[A(D') = y] \cdot \exp(\epsilon) + \delta$$

- With δ probability anything can happen
- Typically δ is chosen very small $\delta \leq n^{-1}$

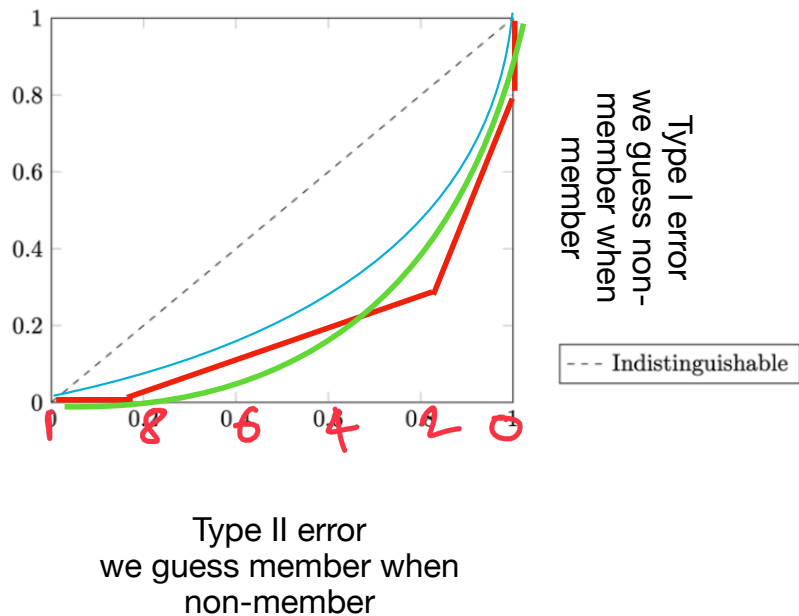
Differentially Private Algorithms

Gaussian mechanism

- Suppose $A(D) = 0$, $A(D') = 1$. Release $\hat{y} = y + \text{Gaussian}(0, \epsilon^{-1})$
- $z \sim \text{Gaussian}(\mu, \sigma^2) \Rightarrow p(z) \propto \frac{1}{\sigma} e^{-\frac{1}{2}(\frac{z-\mu}{\sigma})^2}$
- $\Pr[\hat{y} | y = 0] = \text{Gaussian}(0, \epsilon^{-1})$ and $\Pr[\hat{y} | y = 1] = \text{Gaussian}(1, \epsilon^{-1})$
- $\frac{\Pr[A(D) = y]}{\Pr[A(D') = y]} = ?$ what happens now?

Differentially Private Algorithms

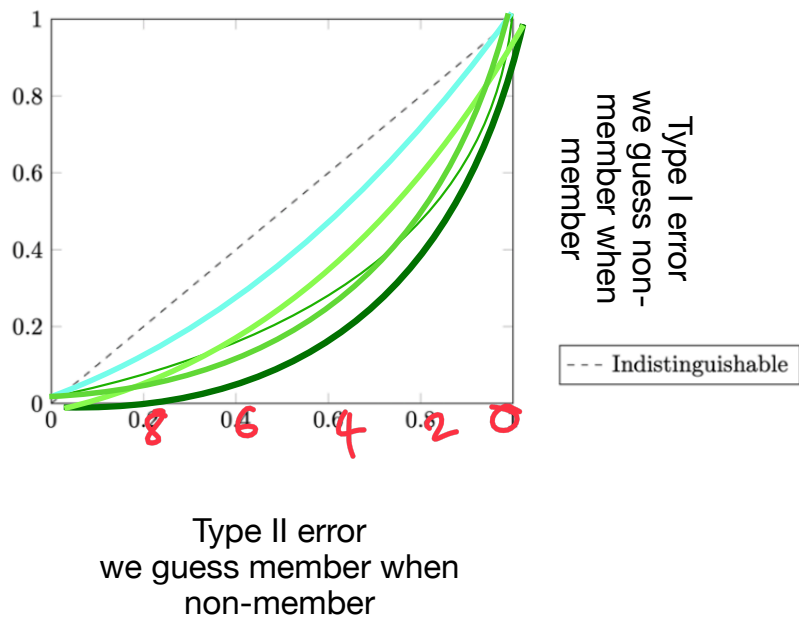
f-DP and Gaussian DP



- All this seems a bit ad-hoc. Is there a “canonical” definition of privacy?
- **Definition.** An algorithm A satisfies **f-DP** if the optimal tradeoff curve is below the function f .
- Generalizes all previous notions. What f should we pick? Both green and red curves satisfy.

Differentially Private Algorithms

f-DP and Gaussian DP



- There is a special family of curves:
Gaussian tradeoff curve
- **Definition.** An algorithm A satisfies μ -**Gaussian Differential Privacy** if it is harder to distinguish between $A(D)$ vs. $A(D')$ than $\mathcal{N}(0,1)$ vs. $\mathcal{N}(\mu,1)$

Differentially Private Algorithms

Gaussian mechanism

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

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

- What about Δ -sensitivity wrt $\|\cdot\|_2$
- Gaussian mechanism with GDP is great for ℓ_2 sensitivity!