

Lecture 1: Privacy I

Franziska Boenisch and Adam Dziedzic
Course on Trustworthy Machine Learning

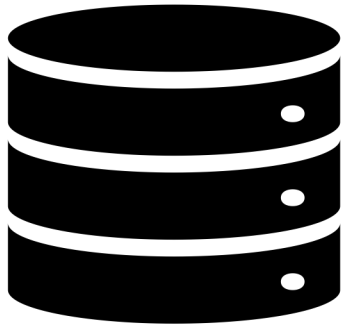


Outline

- I. Privacy Leakage in Machine Learning
 - I. Adversary
 - II. Treat-Space
- II. Attribute Inversion Attacks
- III. Model Inversion Attacks
- IV. Membership Inference Attacks
 - I. Shadow Models
 - II. Loss-based Attacks
 - III. Likelihood Ratio Attack
- V. Intro to Differential Privacy
 - I. Intuition
 - II. Formula

Motivation: Extraction of Training Data

Training Data



train



Diffusion
model

generate



(ℓ_2 distance = 0.031)



Prompt: Ann Graham Lotz

Diffusion models memorize training images and emit them at test time.

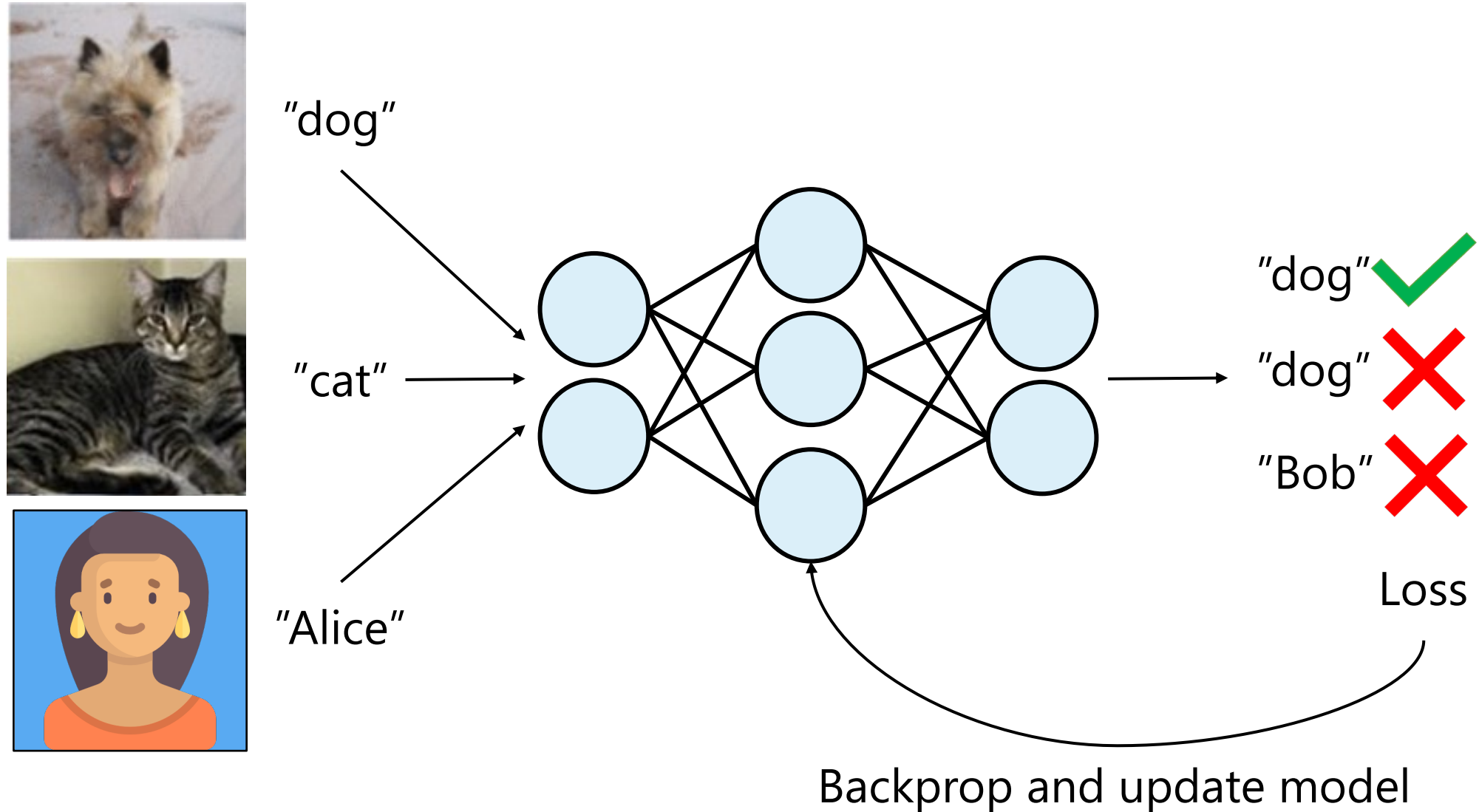
1. As the **quality of models increases so does privacy leakage**.
2. Extraction methodology:
 - Generate many examples using the diffusion model.
 - Perform **membership inference** to find training samples.



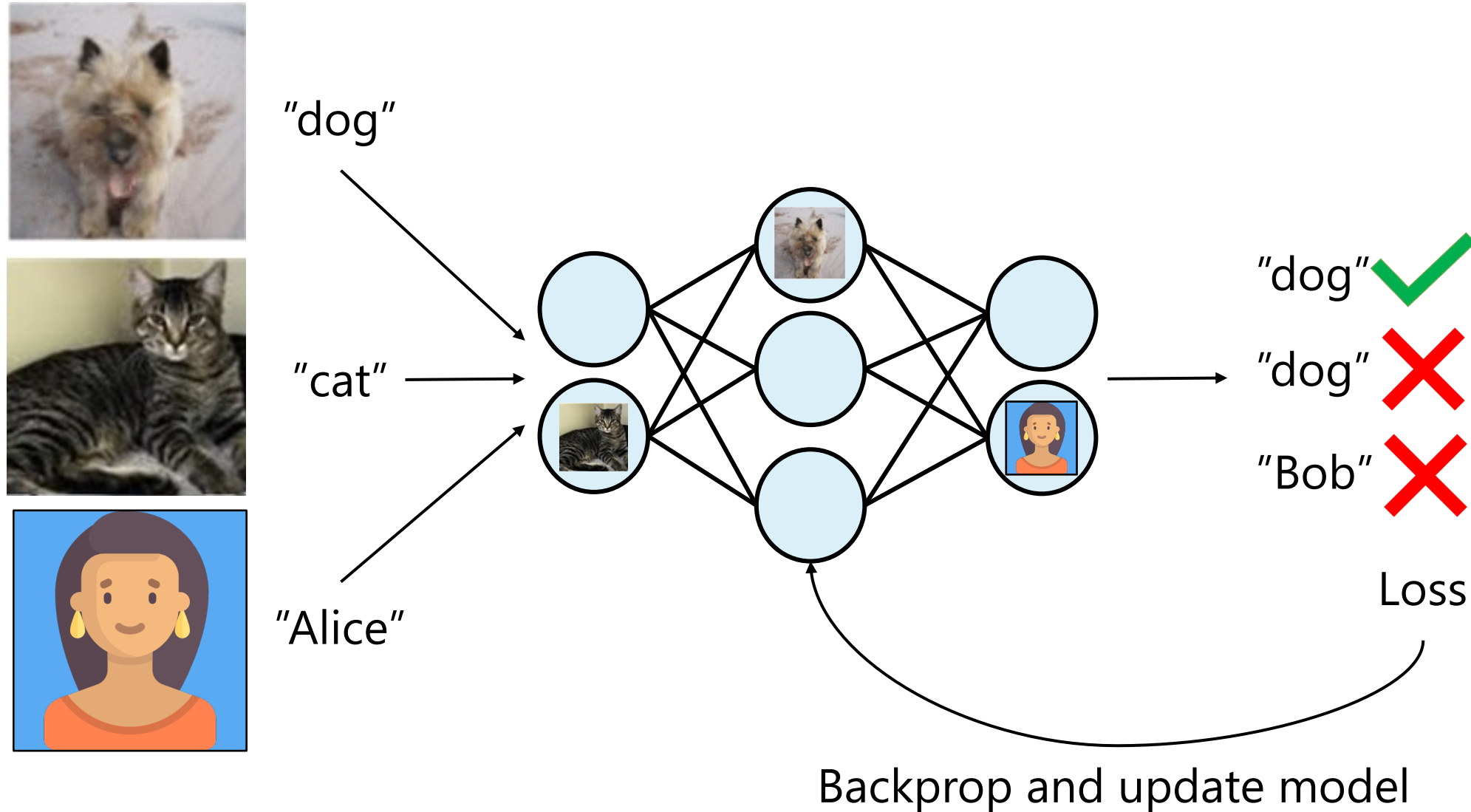
Caption: Living in the light with Ann Graham Lotz

[Wikipedia profile picture](#)

Private data in model training

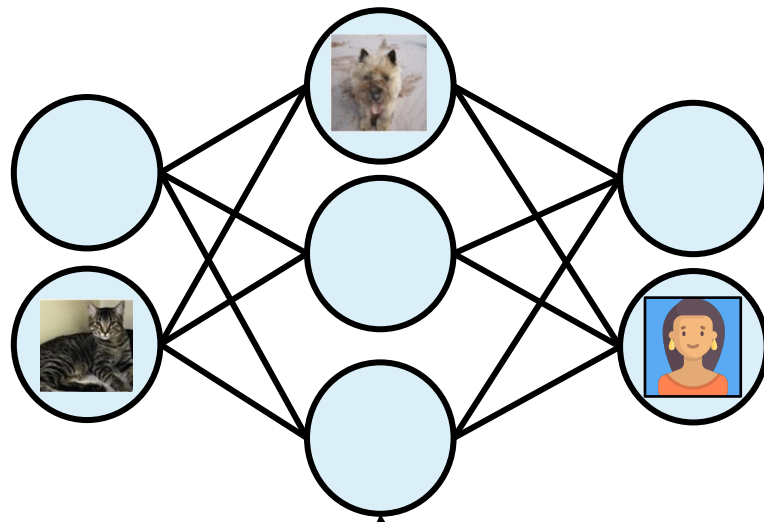


Private data in model training



Where can privacy leak?

From Parameters



"dog" ✓

"dog" ✗

"Bob" ✗

From Predictions

Loss

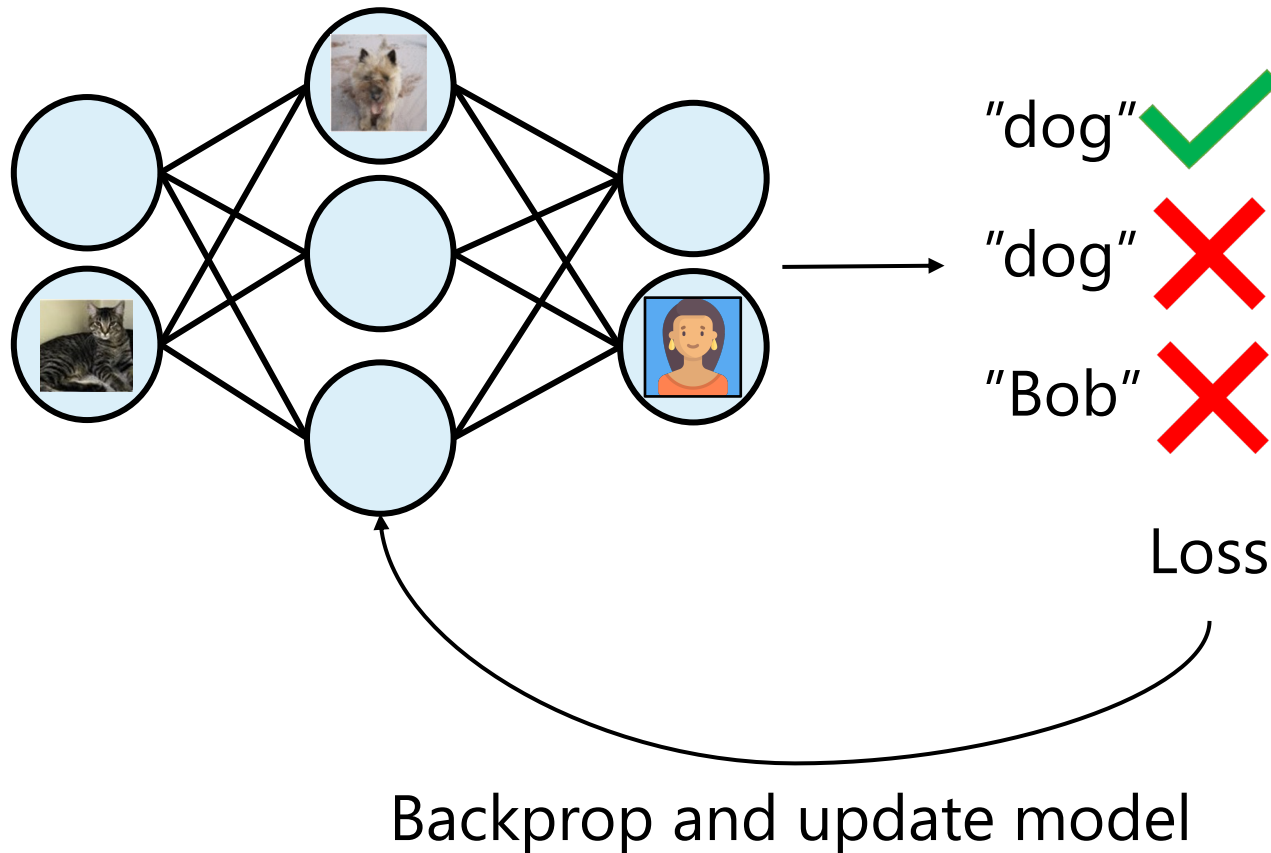
Backprop and update model

During Training

What are the adversary's abilities?

Observe/change the parameters

From Parameters

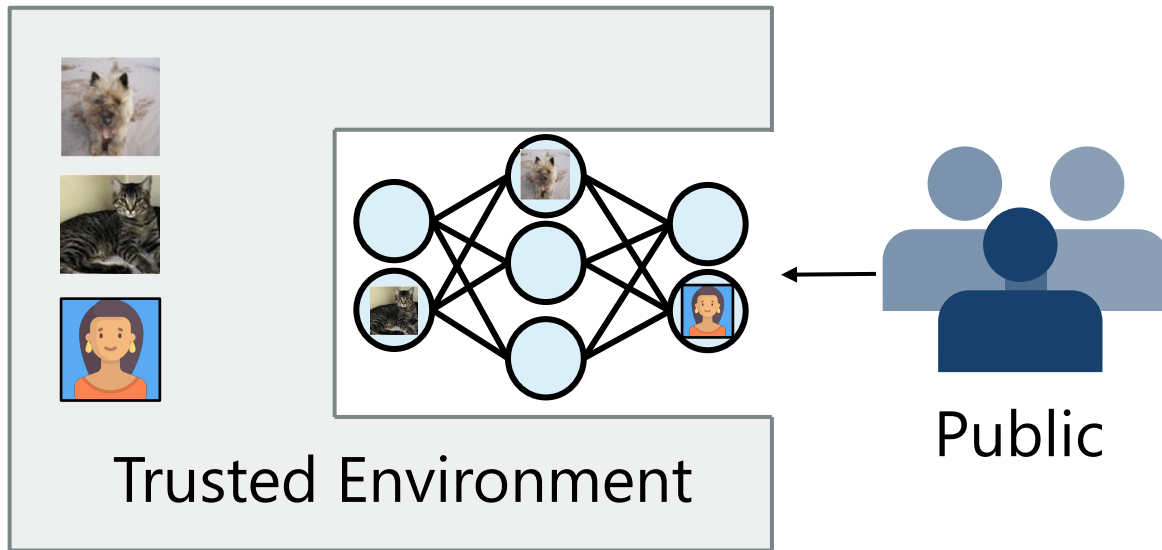


From Predictions

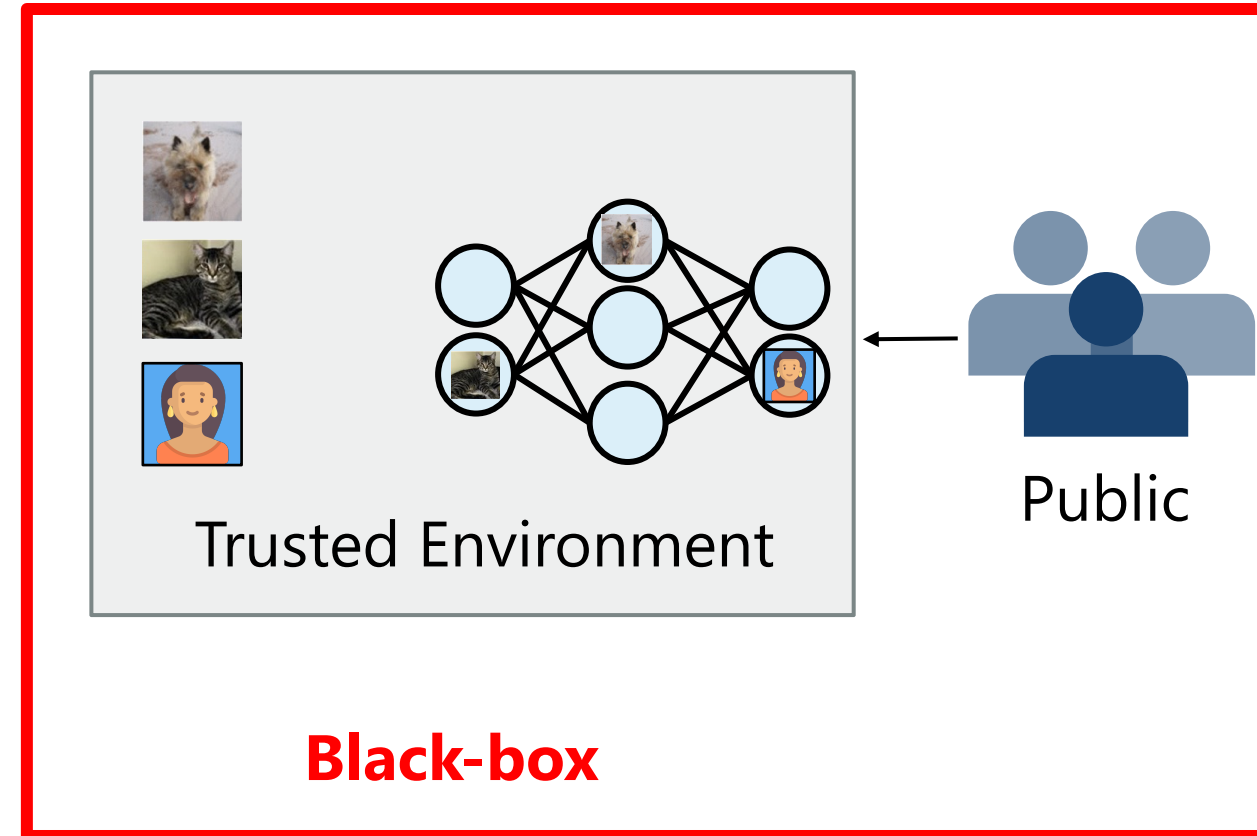
- Knowledge about:**
- **Model architecture**
 - **Data attributes**
 - **Data distribution**
 - **Hyperparameters**

Observe/manipulate training During Training

What is the threat space?



White-box



Black-box

Attribute Inversion

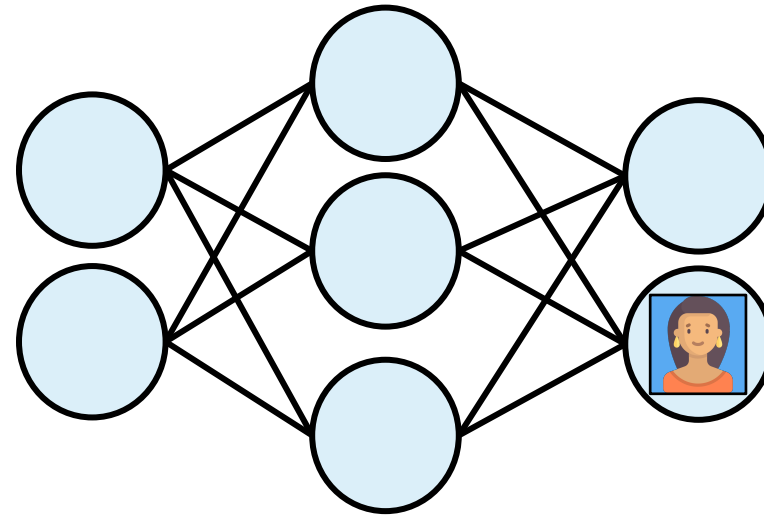
Goal: Disclose the secret attribute of a training data point.



Name: Alice
Age: 34
Height: 1,72m
Smoker: Yes

Risk: "High"

Train



Attribute Inversion

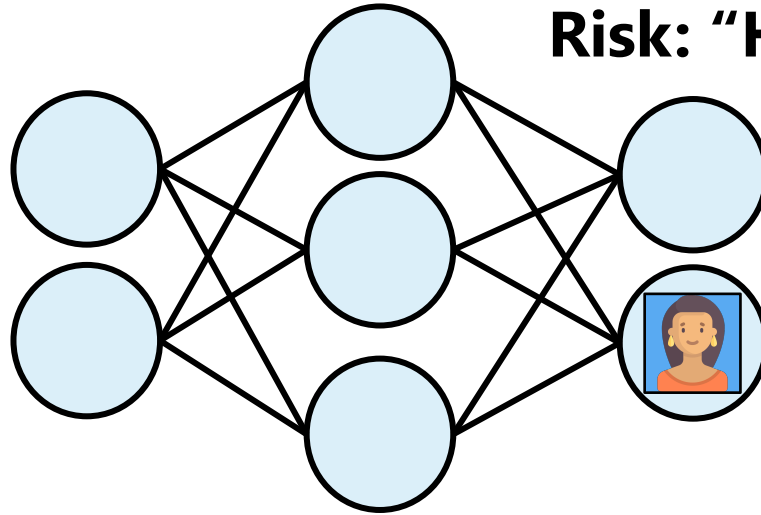


Goal: Disclose the secret attribute of a training data point.



Name: Alice
Age: 34
Height: 1,72m
Smoker: No

Predict



Confidence:
0.74

Attribute Inversion



Goal: Disclose the secret attribute of a training data point.

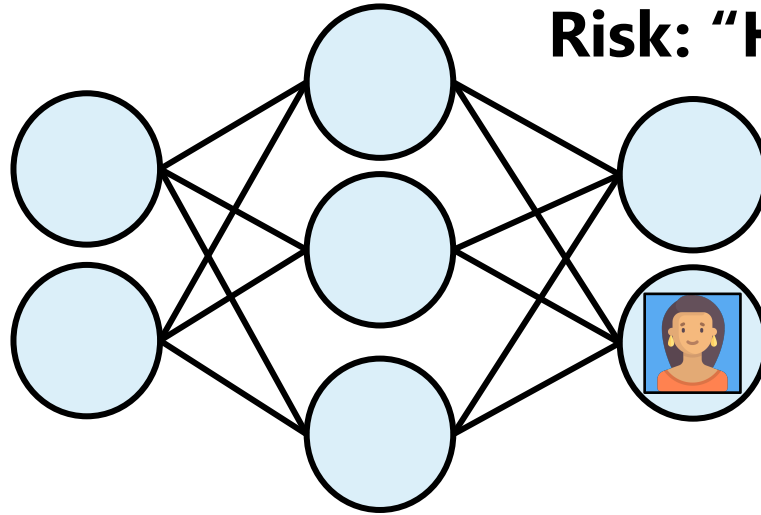


Name:
Age:
Height:
Smoker:

Alice
34
1,72m
Yes



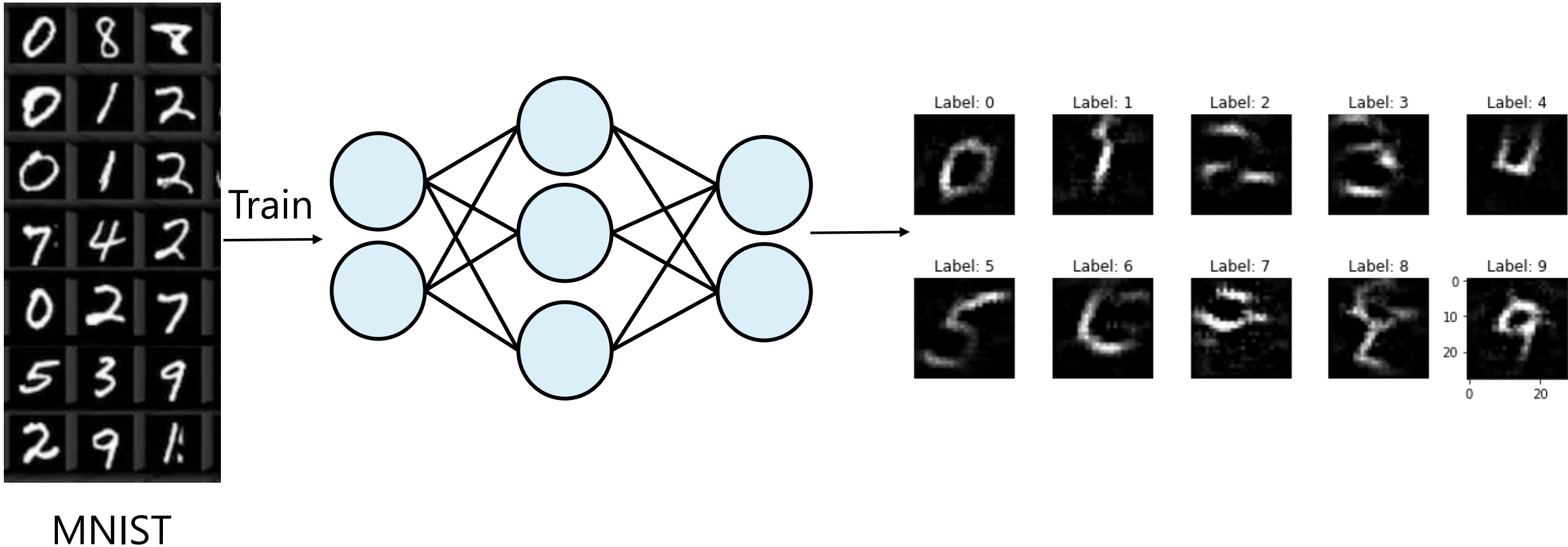
Predict



Confidence:
0.99

Model Inversion

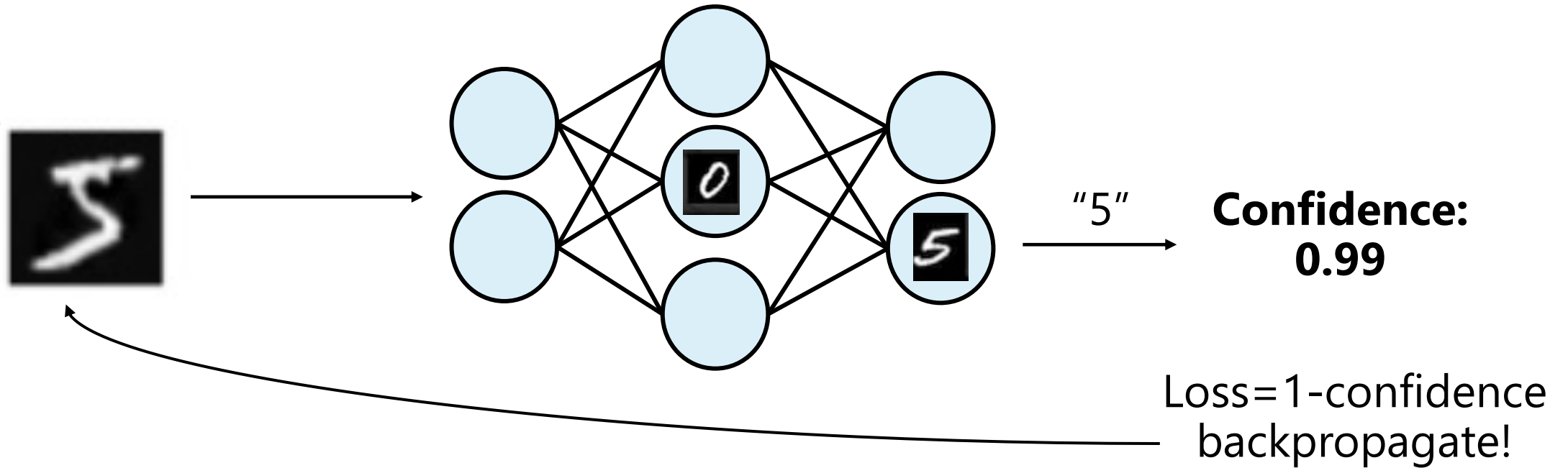
Goal: Disclose a “prototype” of each training class.



Model Inversion

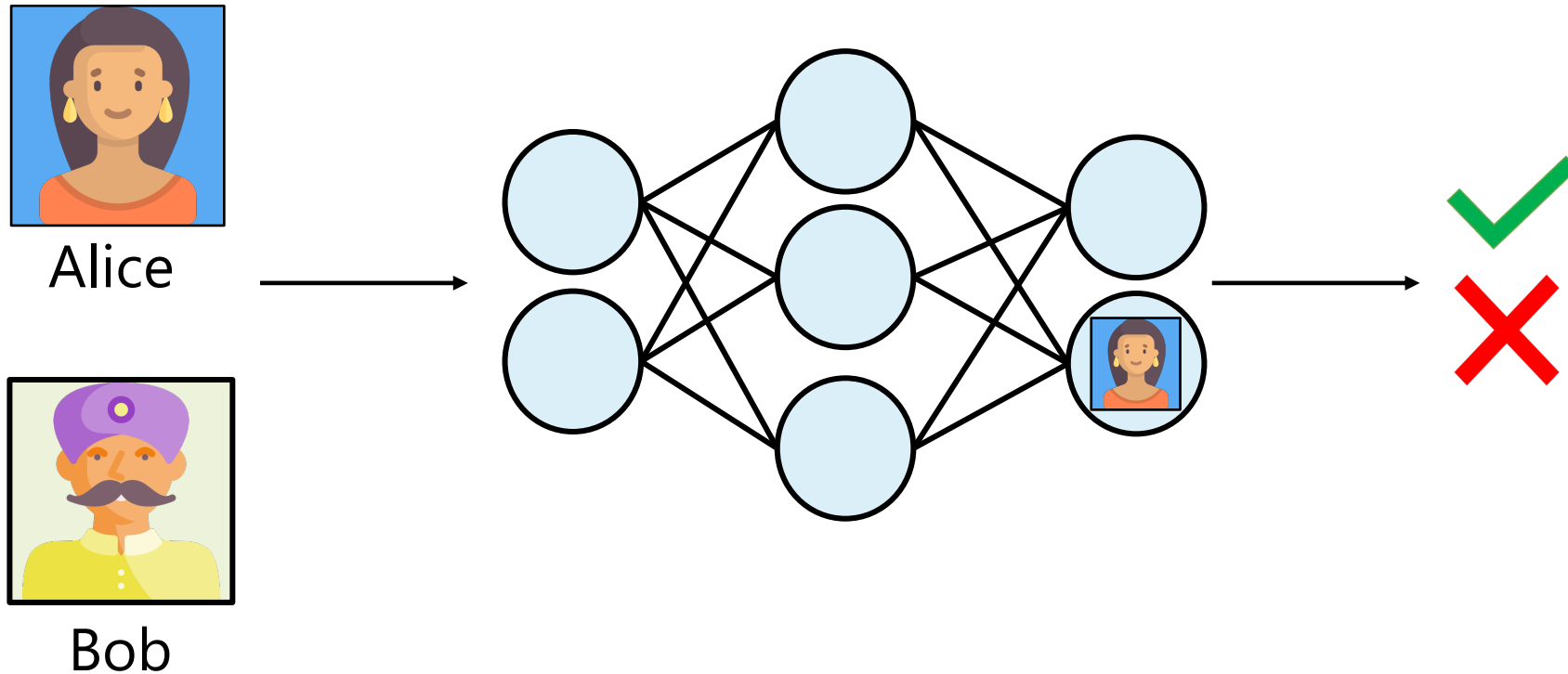


Goal: Disclose a "prototype" of each training class.



Membership Inference Attacks (MIA)

Goal: Disclose whether a given data point was used to train the model.



The Membership Inference "Game"



Challenger C

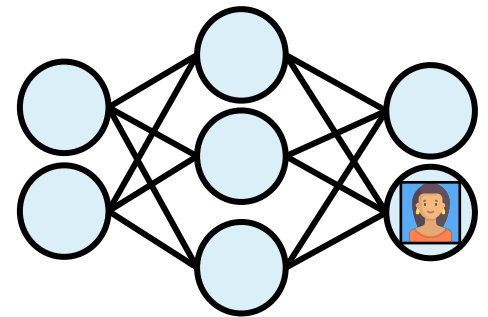
Adversary A



C samples training dataset $D \leftarrow \mathcal{D}$ and trains model $f \leftarrow \mathcal{A}(D)$ with algorithm \mathcal{A} .

1. C flips a coin b , and samples a point $(x, y) \in \mathcal{D} \setminus D$ if $b = 0$. Otherwise, if $b = 1$, sample $(x, y) \in D$.
2. C sends (x, y) to A .
3. A gets access to \mathcal{D} and model f and outputs \hat{b} .
4. A wins if $\hat{b} = b$.

Shadow-Model Based Membership Inference

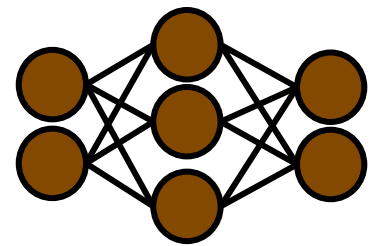


Was Alice a member of the training dataset?



Train Binary Classifier

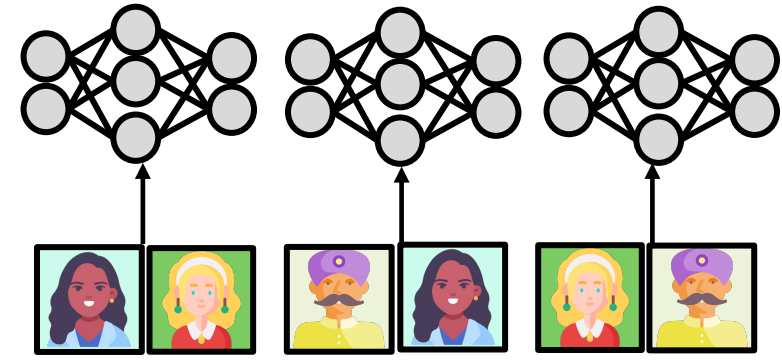
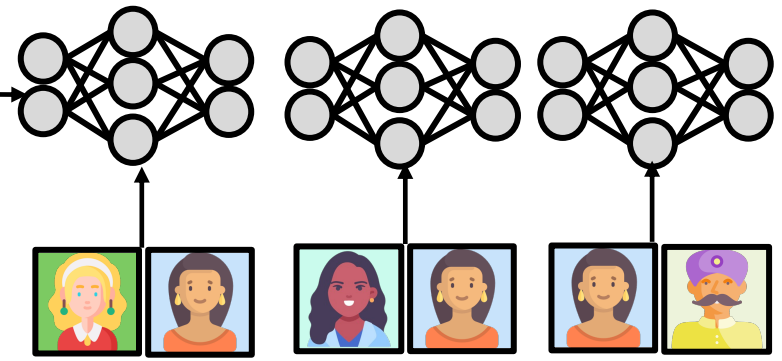
(0.14)



→ ("In"/"Out")

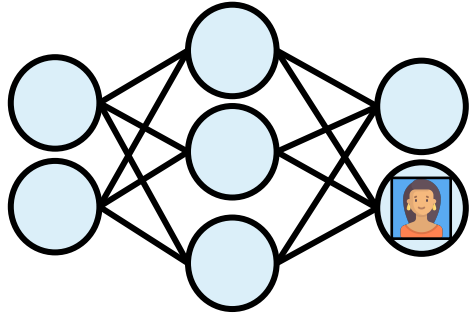
(0.79,"In") (0.83,"In") (0.92,"In") (0.44,"Out") (0.24,"Out") (0.32,"Out")

Query

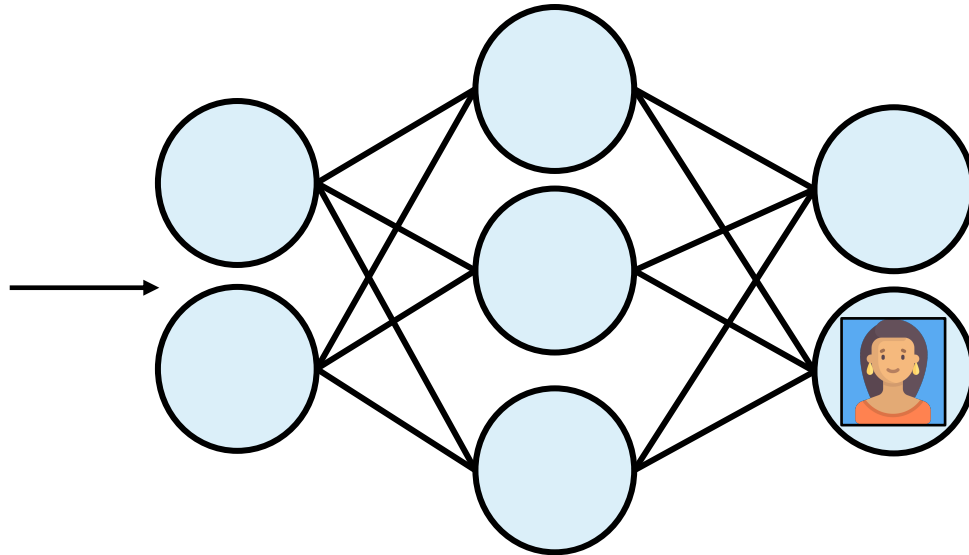


Shadow Models

Shadow-Model Based Membership Inference

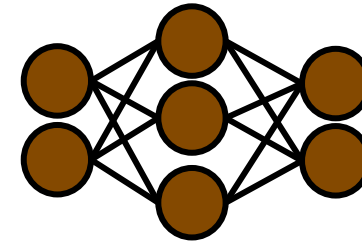


Was Alice a member of the training dataset?



Target Model

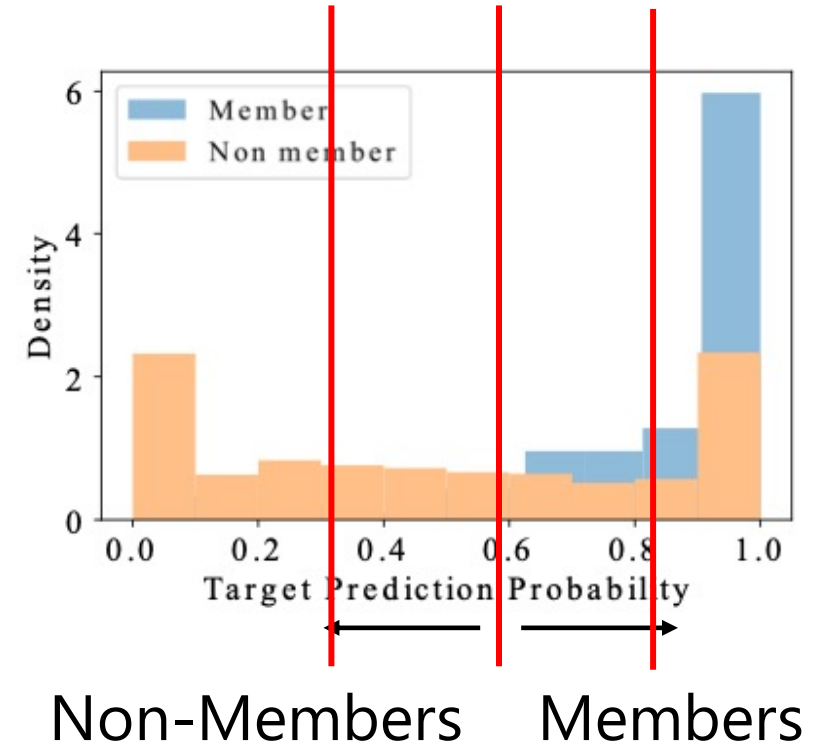
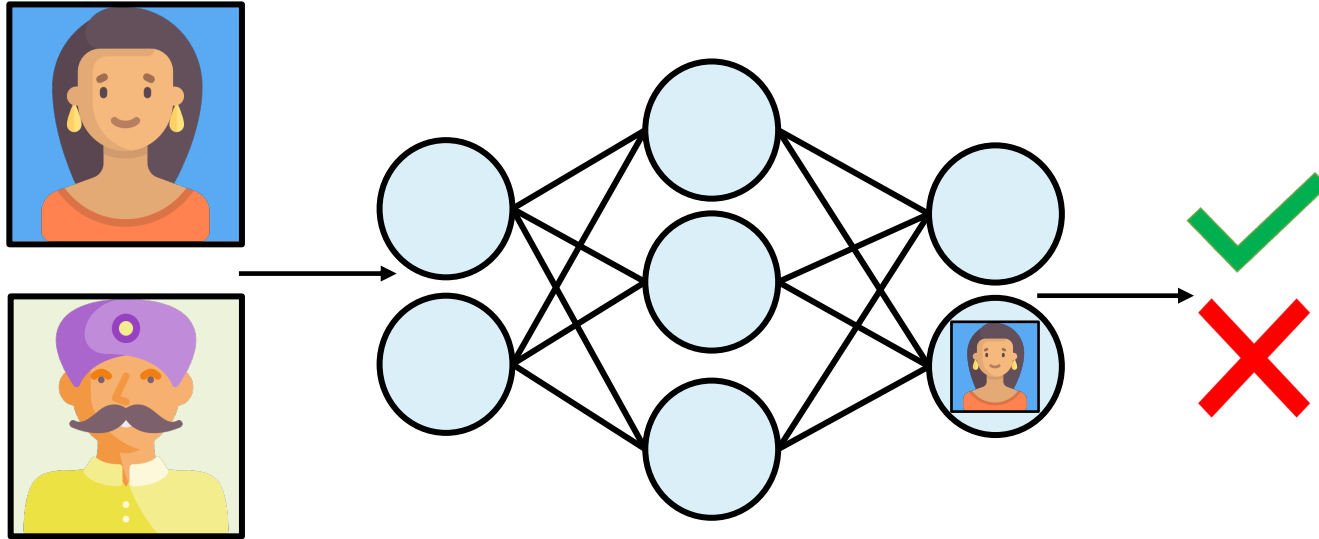
0.76



Binary Classifier

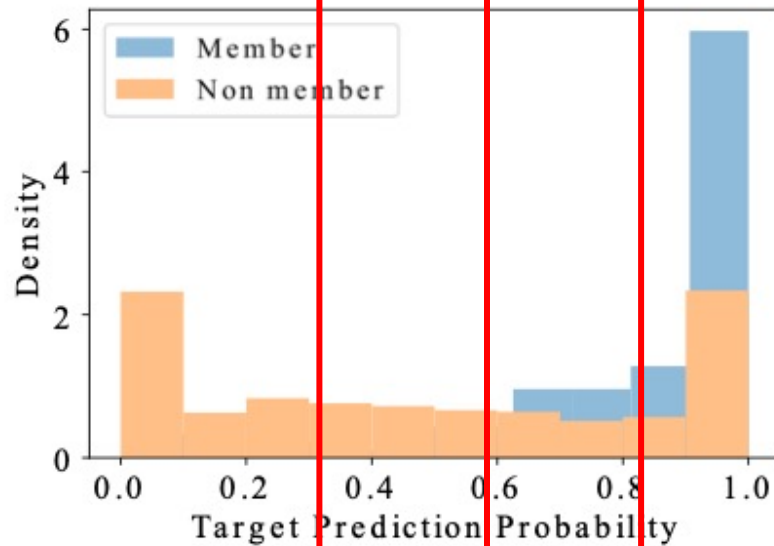
"In"

Threshold-based Membership Inference



Attacker's Prediction Success

FPR=, TRP=

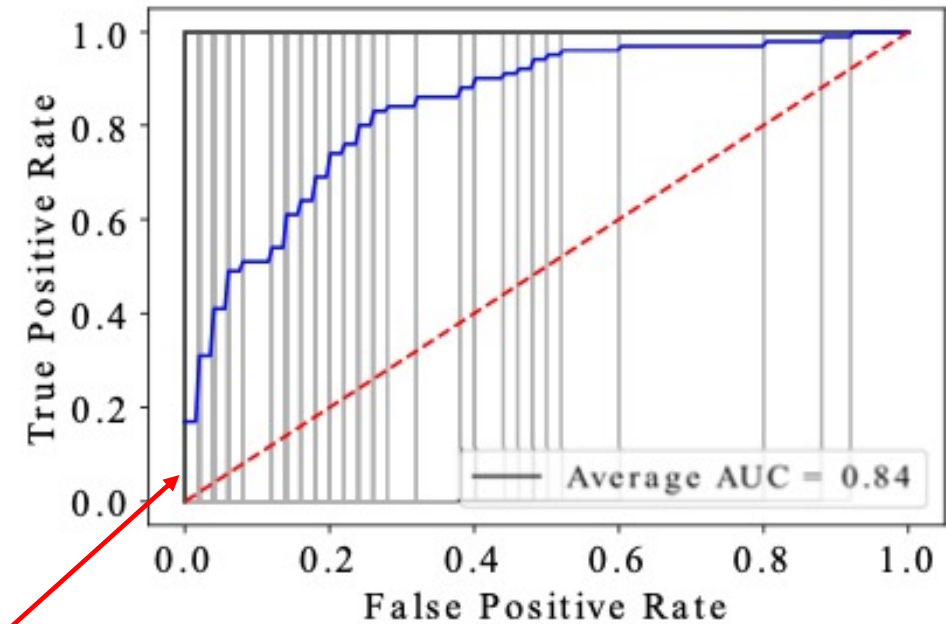


FPR=, TRP=

FPR=, TRP=

Important!
TPR at low FPR

Receiver operating characteristic curve
(ROC Curve)



AUC: Area Under the Curve

Average 😞

FPR: False Positive Rate: Non-member is classified as member.

TPR: True Positive Rate: Member is classified as member.

Likelihood Ratio Attack (LiRA)

Consider two distributions over models:

$\mathbb{Q}_{in} = \{f \leftarrow \mathcal{A}(D \cup \{(x, y)\}) \mid D \leftarrow \mathcal{D}\}$ ← Models trained with data point (x, y)

$\mathbb{Q}_{out} = \{f \leftarrow \mathcal{A}(D \setminus \{(x, y)\}) \mid D \leftarrow \mathcal{D}\}$ ← Models trained without data point (x, y)

Thresholding the Likelihood-ratio Test between the two hypotheses:

$$\Lambda(f; x, y) = \frac{p(f | \mathbb{Q}_{in}(x, y))}{p(f | \mathbb{Q}_{out}(x, y))} \quad \textbf{Intractable!}$$

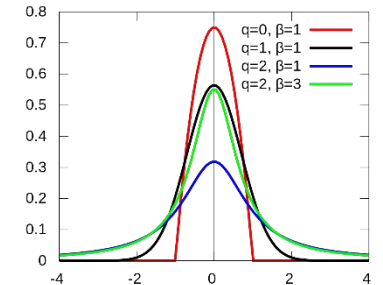
p probability density function of f under distribution of model parameters

Simplify by using loss instead: $\mathbb{Q}'_{in \setminus out}$ distribution of losses on (x, y)

$$\Lambda(f; x, y) = \frac{p(\ell(f(x), y) | \mathbb{Q}'_{in}(x, y))}{p(\ell(f(x), y) | \mathbb{Q}'_{out}(x, y))}$$

Likelihood Ratio Attack

Train shadow models to estimate \mathbb{Q}'_{in} and \mathbb{Q}'_{out} .
Simplify by assumption that they follow Gaussian distribution

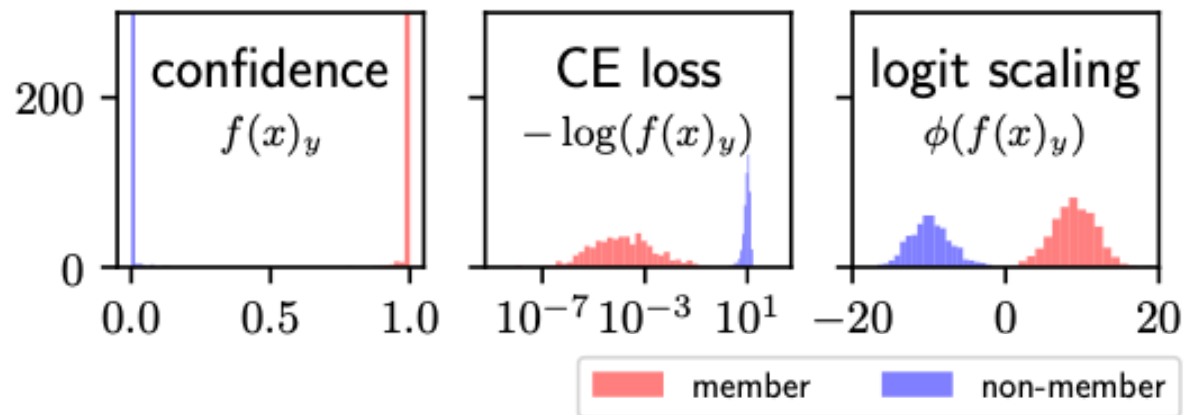


Four variables to be estimated:

means: μ_{in}, μ_{out}

standard deviations: $\sigma_{in}, \sigma_{out}$

How to ensure losses are indeed Gaussians?



Logit scaling to the model confidence:
 $\phi(p) = \log \frac{p}{1-p}$, for $p = f(x)_y$

Likelihood Ratio Attack

Require: model f , example (x, y) , data distribution \mathbb{D}

```
1:  $\text{confs}_{\text{in}} = \{\}$ 
2:  $\text{confs}_{\text{out}} = \{\}$ 
3: for  $N$  times do
4:    $D_{\text{attack}} \leftarrow^{\$} \mathbb{D}$  ▷ Sample a shadow dataset
5:    $f_{\text{in}} \leftarrow \mathcal{T}(D_{\text{attack}} \cup \{(x, y)\})$  ▷ train IN model
6:    $\text{confs}_{\text{in}} \leftarrow \text{confs}_{\text{in}} \cup \{\phi(f_{\text{in}}(x)_y)\}$ 
7:    $f_{\text{out}} \leftarrow \mathcal{T}(D_{\text{attack}} \setminus \{(x, y)\})$  ▷ train OUT model
8:    $\text{confs}_{\text{out}} \leftarrow \text{confs}_{\text{out}} \cup \{\phi(f_{\text{out}}(x)_y)\}$ 
9: end for
10:  $\mu_{\text{in}} \leftarrow \text{mean}(\text{confs}_{\text{in}})$ 
11:  $\mu_{\text{out}} \leftarrow \text{mean}(\text{confs}_{\text{out}})$ 
12:  $\sigma_{\text{in}}^2 \leftarrow \text{var}(\text{confs}_{\text{in}})$ 
13:  $\sigma_{\text{out}}^2 \leftarrow \text{var}(\text{confs}_{\text{out}})$ 
14:  $\text{conf}_{\text{obs}} = \phi(f(x)_y)$  ▷ query target model
15: return  $\Lambda = \frac{p(\text{conf}_{\text{obs}} \mid \mathcal{N}(\mu_{\text{in}}, \sigma_{\text{in}}^2))}{p(\text{conf}_{\text{obs}} \mid \mathcal{N}(\mu_{\text{out}}, \sigma_{\text{out}}^2))}$ 
```

Effective but computationally costly



How to defend against MIA?

1. Add noise to confidence vector
2. Do not output prediction probability, just output labels
3. Reduce overfitting: Regularization, different losses

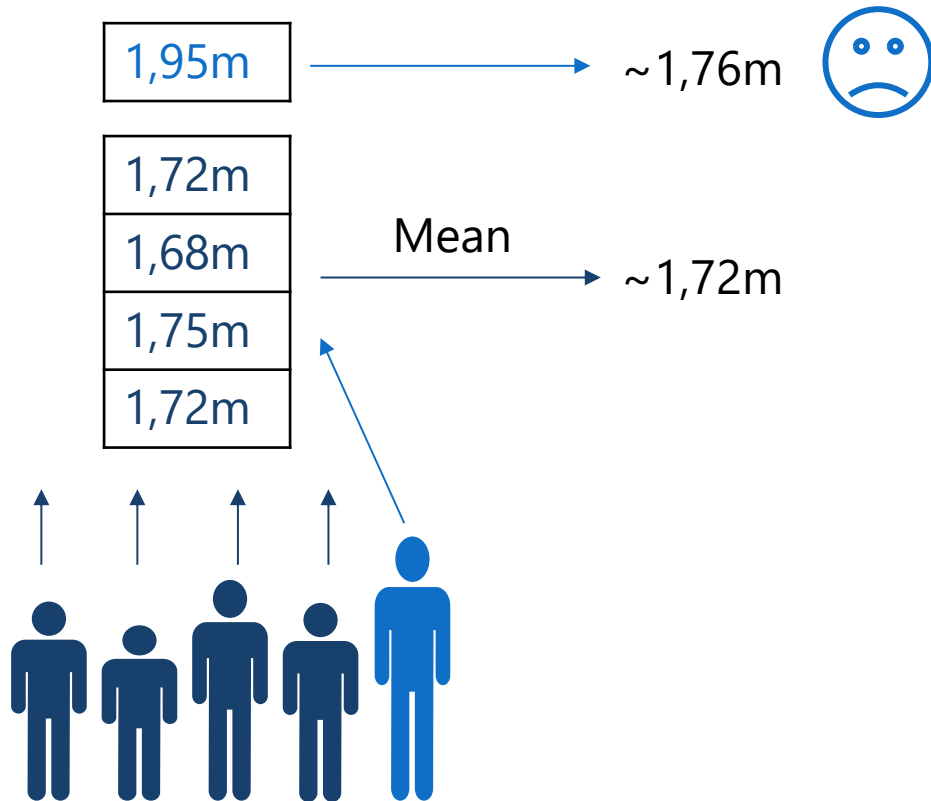


Empirical!

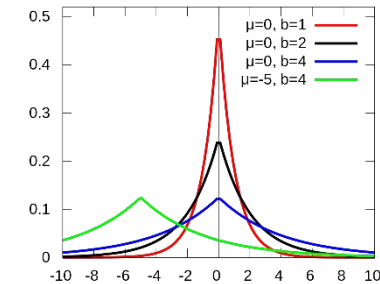
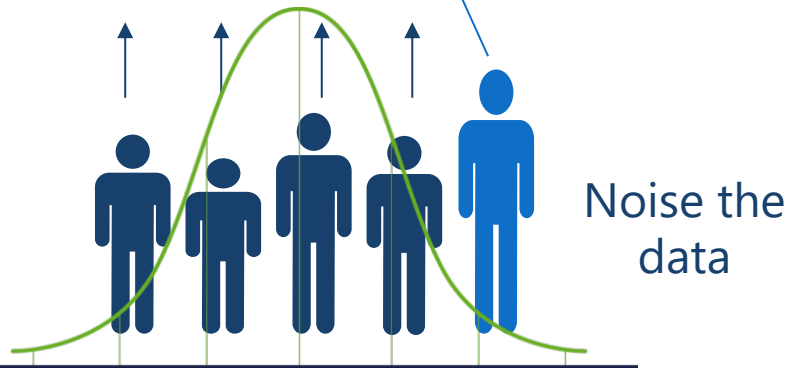
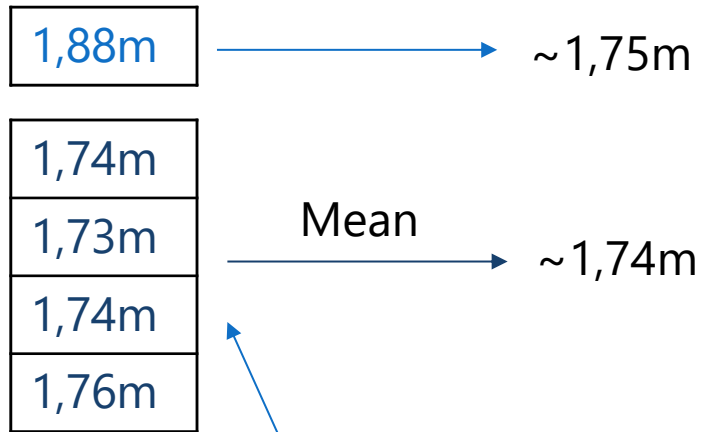


Can we get guarantees?

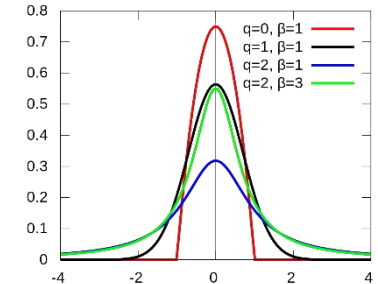
Differential Privacy – Example 1



Differential Privacy – Example 1

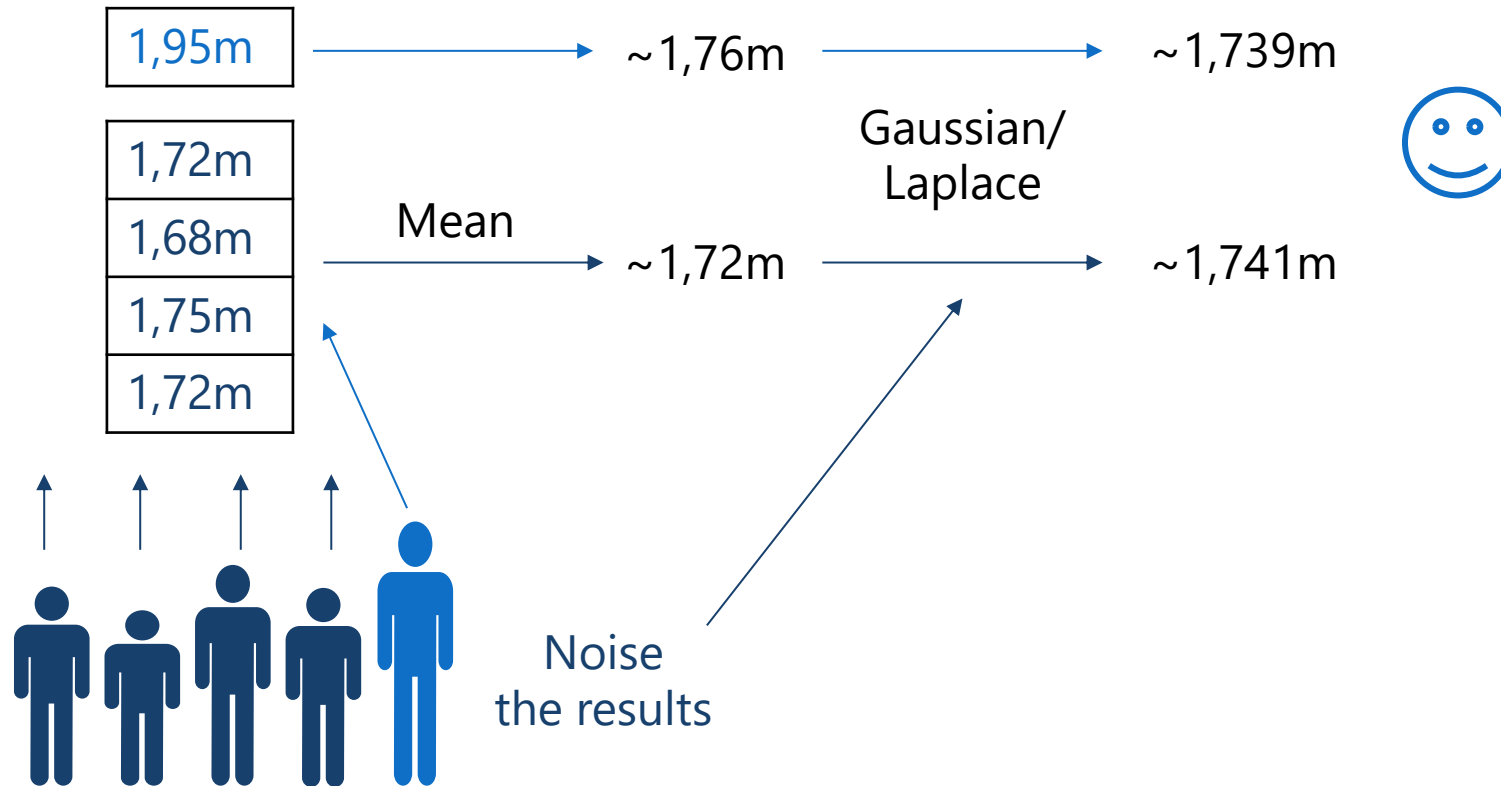


Laplace

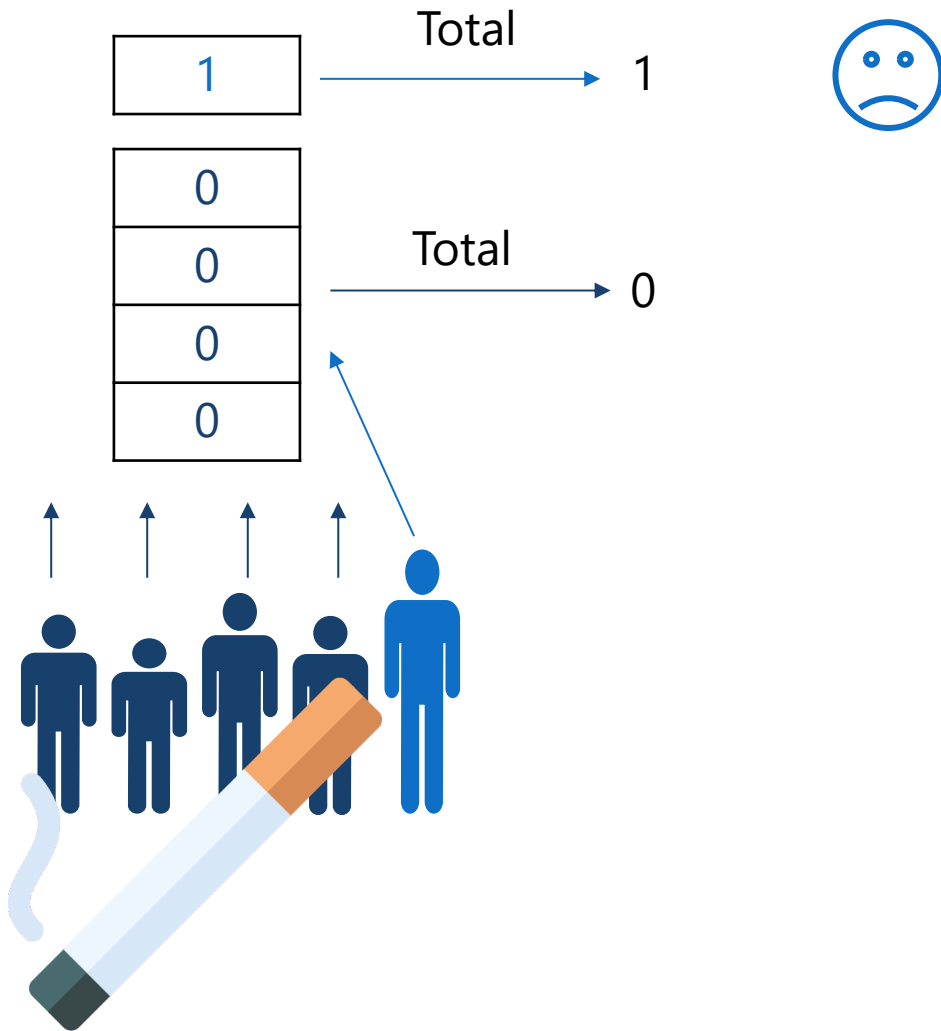


Gauss

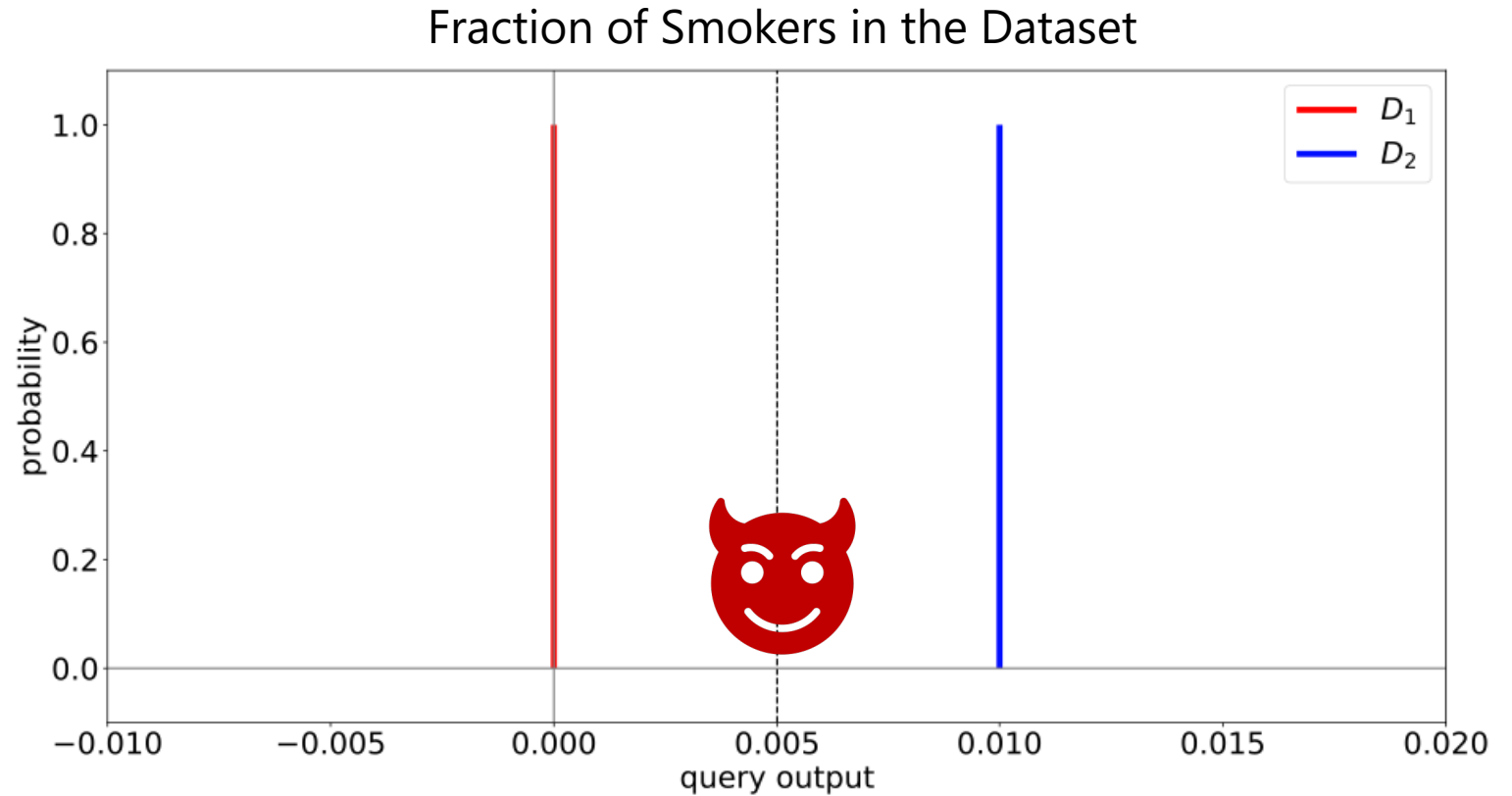
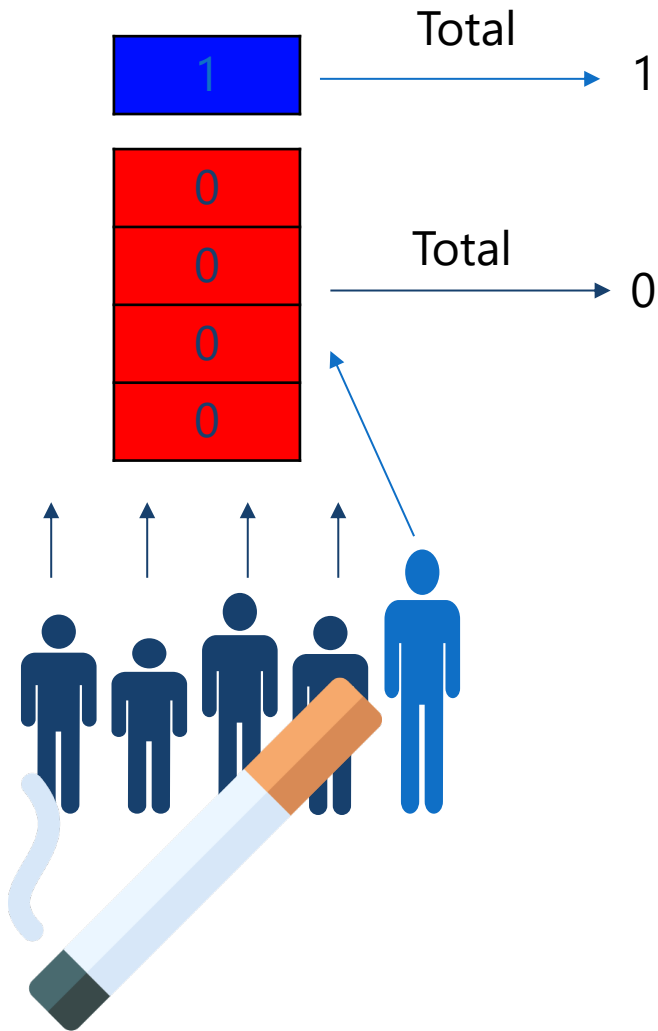
Differential Privacy – Example 1



Differential Privacy – Example 2

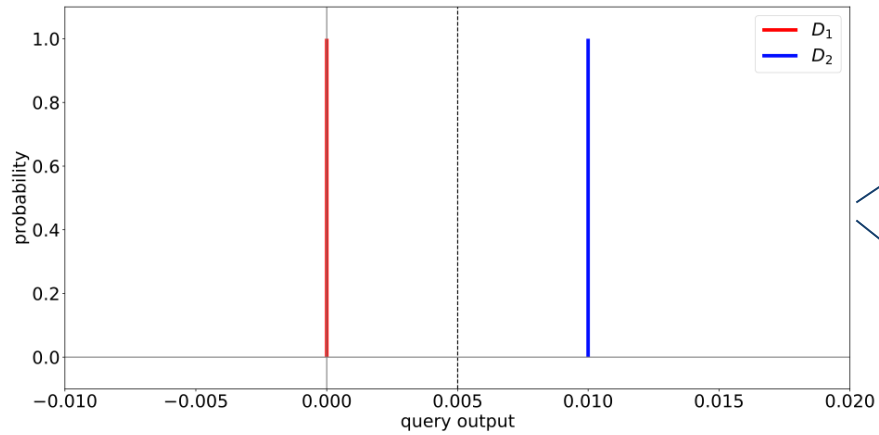


Deterministic algorithms yield no privacy

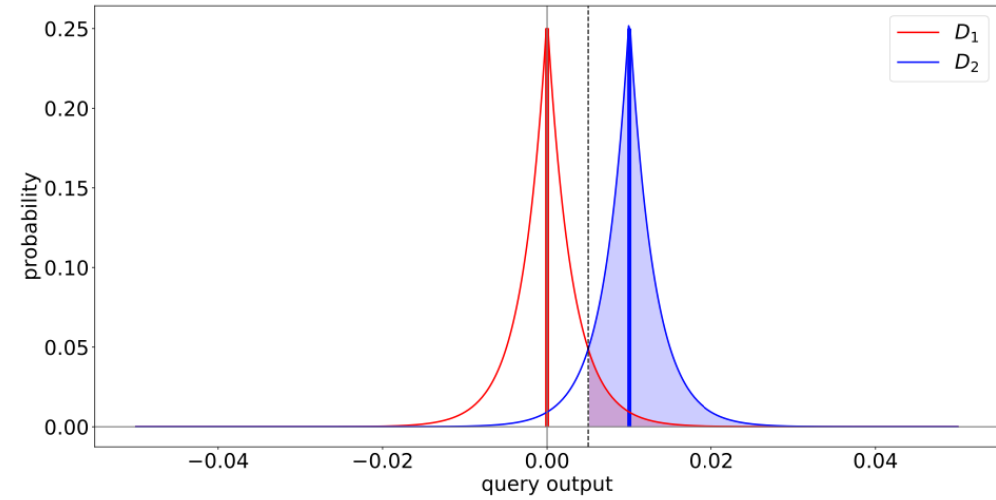


How much noise to add?

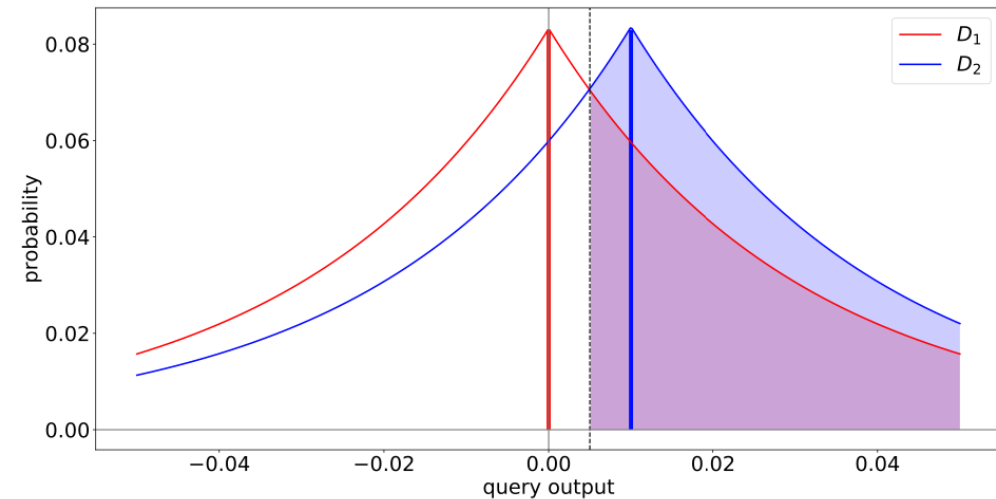
Fraction of Smokers in the Dataset



Little Noise



Much Noise



Formalizing Differential Privacy

Intuition: An algorithm \mathbf{M} provides (ϵ, δ) -Differential Privacy (DP) if it produces “roughly same” outputs on any pair of training datasets d and d' that differ only by a single data point.

$$\Pr[\mathbf{M}(d) \in S] \leq e^{\epsilon} \Pr[\mathbf{M}(d') \in S] + \delta$$

Diagram illustrating the formalization of Differential Privacy:

- Closeness**: Points to the term e^{ϵ} , representing the multiplicative factor for closeness.
- Probability of the Closeness Violation**: Points to the term δ , representing the probability of a violation.
- Randomized Algorithm**: Points to the term $\mathbf{M}(d)$, representing the algorithm applied to dataset d .
- Possible Outputs**: Points to the term S , representing the set of possible outputs.

Formalizing Differential Privacy

$$\Pr[\mathbf{M}(d) \in S] \leq e^{\epsilon} \Pr[\mathbf{M}(d') \in S] + \delta$$

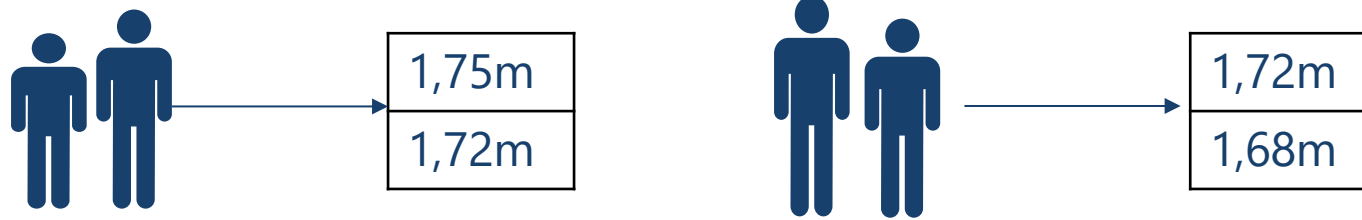
$\epsilon > 0$: Privacy budget
Smaller \rightarrow more privacy

$\delta \in [0,1]$: Probability of violating closeness
Smaller \rightarrow more privacy, usually chosen $< 1/n$ with n data points

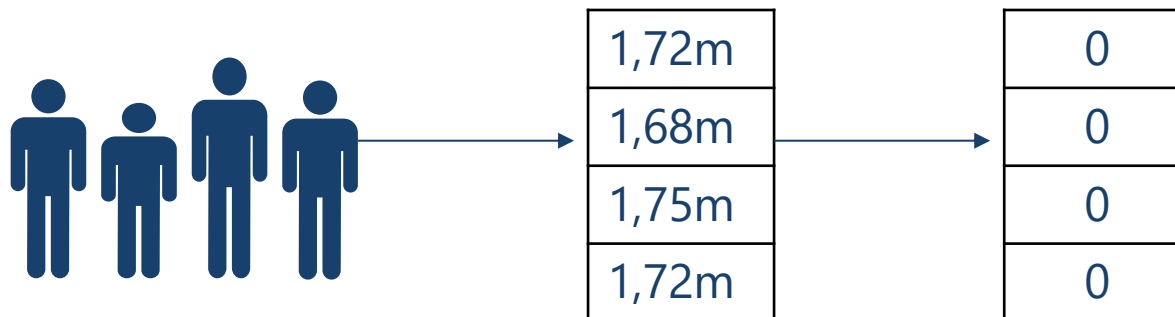
With $(\epsilon, 0)$, we fulfill pure ϵ -DP, with $\delta > 0$, we have approximate DP.

Properties Differential Privacy

Parallel Composition: If $\mathbf{M}(x)$ fulfills ϵ, δ -DP, and if we split our data \mathcal{D} into k disjoint subsets $\mathcal{D} = x_1 \cup \dots \cup x_k$, then the mechanism that releases all results $M(x_1), \dots, M(x_k)$ is ϵ, δ -DP.

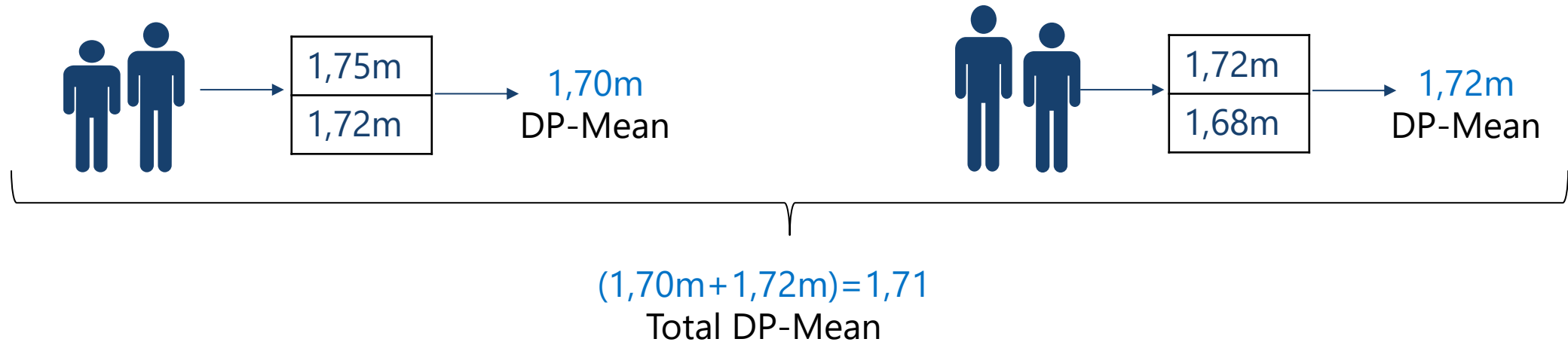


Sequential Composition: If $\mathbf{M}_1(x)$ fulfills ϵ_1, δ_1 -DP and $\mathbf{M}_2(x)$ fulfills ϵ_2, δ_2 -DP, then $\mathbf{G}(x) = (\mathbf{M}_1(x), \mathbf{M}_2(x))$ fulfills $(\epsilon_1 + \epsilon_2, \delta_1 + \delta_2)$ -DP.

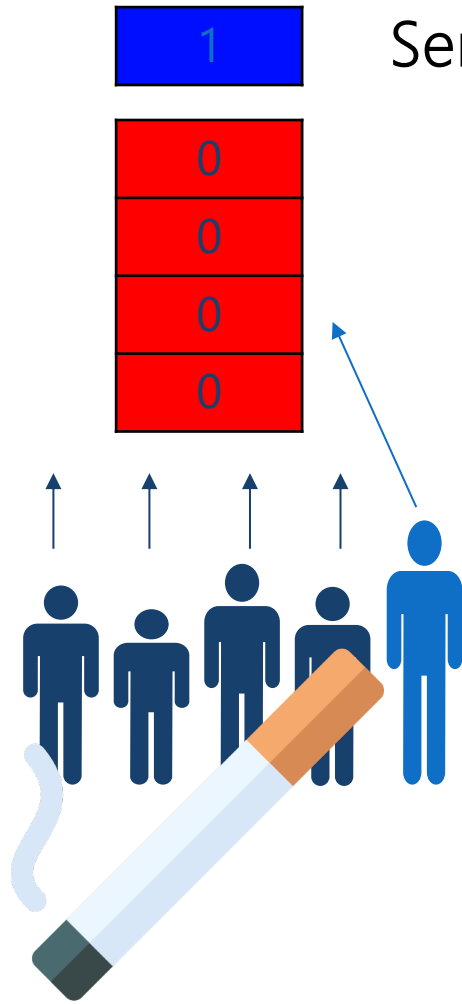


Properties Differential Privacy

Postprocessing guarantees: If an output of an (ϵ, δ) -DP mechanism is further processed or transformed, the guarantees remain.



How to find the noise level: Sensitivity



Sensitivity: By how much can a single data point change the outcome.

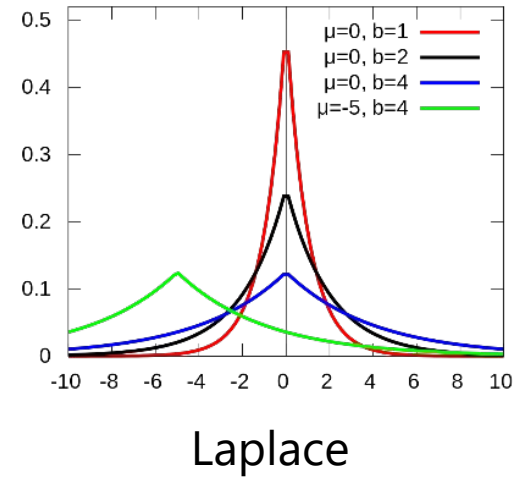
Sensitivity Δf of a function f operating on the neighboring datasets d and d' is defined as $\Delta f = \max(|f(d) - f(d')|)$.

Sensitivity of any counting function is 1.

We can use different norms to calculate the sensitivity.

Laplace Mechanism

Given a function $f: \mathcal{D} \rightarrow \mathbb{R}^d$ where \mathcal{D} is the domain of the dataset and d is the dimension of the output, the Laplace mechanism adds Laplace noise to the output of f .



$$Lap(x|b, \mu) = \frac{1}{2b} e^{-\frac{|x-\mu|}{b}} \longleftarrow \text{Laplace (noise) Distribution}$$

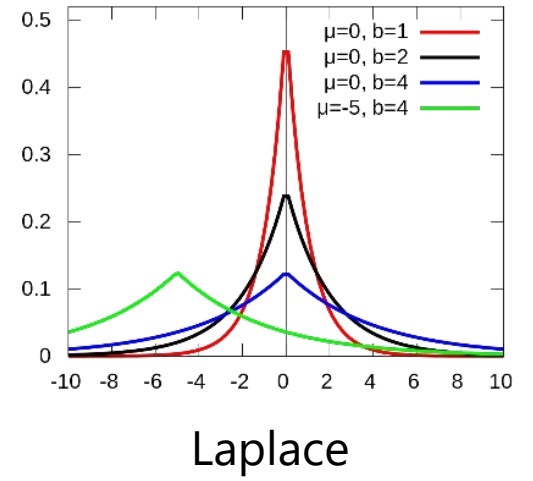
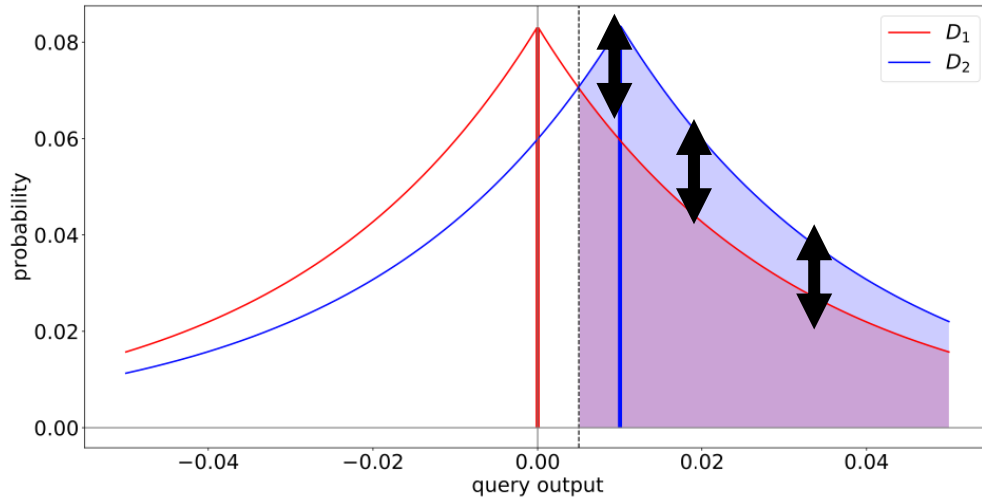
b is the scale parameter of the Laplace distribution.

The Laplace mechanism is $M(D) = f(D) + Lap(0|b)^d$.

If we choose $b = \frac{\Delta f}{\epsilon}$, this mechanism fulfills ϵ -DP.

Proof Sketch

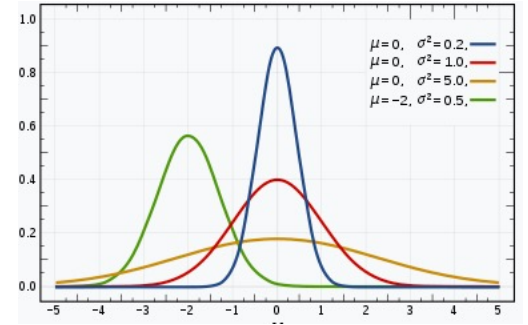
If we choose $b = \frac{\Delta f}{\epsilon}$, this mechanism fulfills ϵ -DP.



Show that $\frac{\Pr[\mathbf{M}(d) \in S]}{\Pr[\mathbf{M}(d') \in S]} \leq e^\epsilon$.

Gaussian Mechanism

Given a function $f: \mathcal{D} \rightarrow \mathbb{R}^d$ where \mathcal{D} is the domain of the dataset and d is the dimension of the output, the Gaussian mechanism adds Gaussian noise to the output of f .



Gauss

$$\mathcal{N}(x|\sigma, \mu) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \longleftarrow \begin{array}{l} \text{Gaussian (noise)} \\ \text{Distribution} \end{array}$$

σ is the standard deviation, and μ the mean.

The Gaussian mechanism is $M(D) = f(D) + \mathcal{N}(0|\sigma, \mu)^d$.

If we choose $\mu = 0$ and $\sigma^2 = \frac{2 \ln\left(\frac{1.25}{\delta}\right)(\Delta f)^2}{\epsilon^2}$, this mechanism fulfills (ϵ, δ) -DP.

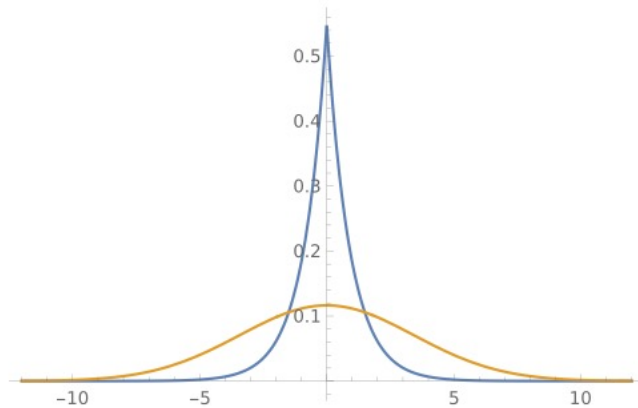
Laplace vs Gaussian Mechanism

$$\text{Lap}(x|b, \mu) = \frac{1}{2b} e^{-\frac{|x-\mu|}{b}} \quad \mathcal{N}(x|\sigma, \mu) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

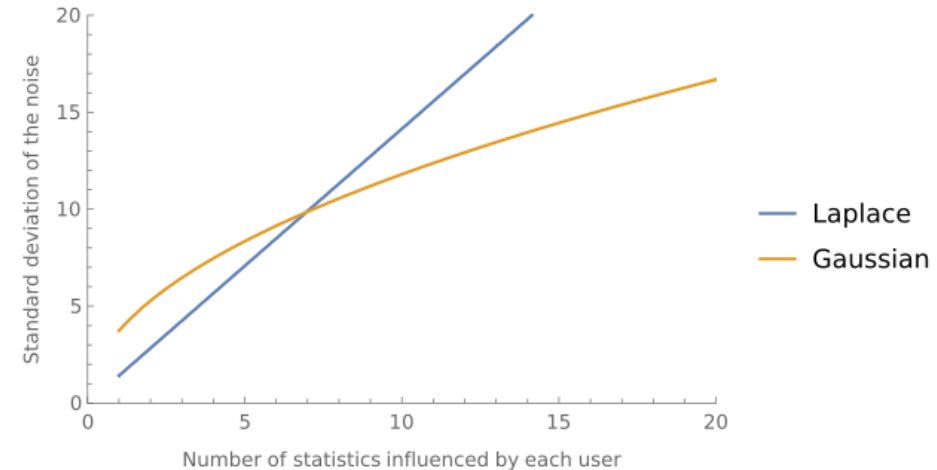
ϵ -DP vs. (ϵ, δ) -DP
Privacy Guarantees

$L1$ vs $L2$

Norms



Noise distributions
at the same $\epsilon \approx 1$



Standard deviation after
sequential execution

Further Reading

- [1] Papernot, N., McDaniel, P., Sinha, A., & Wellman, M. P. (2018, April). [Sok: Security and privacy in machine learning](#). In *2018 IEEE European symposium on security and privacy (EuroS&P)* (pp. 399-414). IEEE.
- [2] Fredrikson, Matt, Somesh Jha, and Thomas Ristenpart. "[Model inversion attacks that exploit confidence information and basic countermeasures](#)." In *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*, pp. 1322-1333. 2015.
- [3] Shokri, Reza, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. "[Membership inference attacks against machine learning models](#)." In *2017 IEEE symposium on security and privacy (SP)*, pp. 3-18. IEEE, 2017.
- [4] Yeom, Samuel, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. "[Privacy risk in machine learning: Analyzing the connection to overfitting](#)." In *2018 IEEE 31st computer security foundations symposium (CSF)*, pp. 268-282. IEEE, 2018.
- [5] Choquette-Choo, Christopher A., Florian Tramer, Nicholas Carlini, and Nicolas Papernot. "[Label-only membership inference attacks](#)." In *International conference on machine learning*, pp. 1964-1974. PMLR, 2021.
- [6] Carlini, Nicholas, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramer. "[Membership inference attacks from first principles](#)." In *2022 IEEE Symposium on Security and Privacy (SP)*, pp. 1897-1914. IEEE, 2022.
- [7] Dwork, Cynthia. "[Differential privacy](#)." In *International colloquium on automata, languages, and programming*, pp. 1-12. Berlin, Heidelberg: Springer Berlin Heidelberg, 2006.
- [8] Carlini, Nicolas, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramer, Borja Balle, Daphne Ippolito, and Eric Wallace. "[Extracting training data from diffusion models](#)." In *32nd USENIX Security Symposium (USENIX Security 23)*, pp. 5253-5270. 2023.

Thank you!

Franziska Boenisch and Adam Dziedzic
boenisch@cispa.de, adam.dziedzic@cispa.de
sprintml.com

Course on Trustworthy Machine Learning