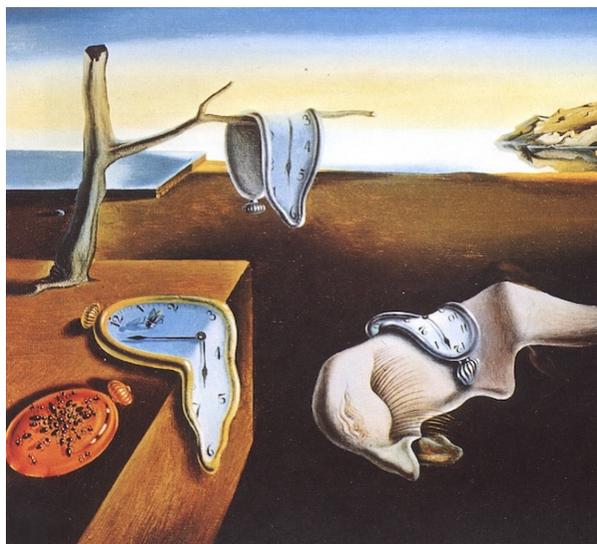




Salvador Dalí, "The Persistence of Memory," 1931



Membership Inference Attack & Differential Privacy

马兴军, 复旦大学 计算机学院



Recap: week 9

- ❑ Data Extraction Attack & Defense
- ❑ Model Stealing Attack
- ❑ Future Research

Register Final Project Team

可信机器学习组队注册

COMP737022 : Trustworthy Machine Learning

<https://trust-ml.github.io/>

🕒 11-01至11-08

*01 队伍名称

*02 选题名称

*03 成员1 (学号)

*04 成员2 (学号)

<https://docs.qq.com/form/page/DU3N1ZUFBUrJVNfJ#/result>

◆ 分组大作业 (占比60%)

- 研究主题 “每个算法/模型/论文都有它的缺陷”， **具体方向自主选择**
- **组队研究实践：3-5人，每组不多于2个博士**
 - **用实验说话**
 - **以组为单位进行期末汇报，每个组5分钟**

■ 得分：结合**选题新颖度**、**实验创新性**、**发现独特性**和**报告质量**四个方面综合评分



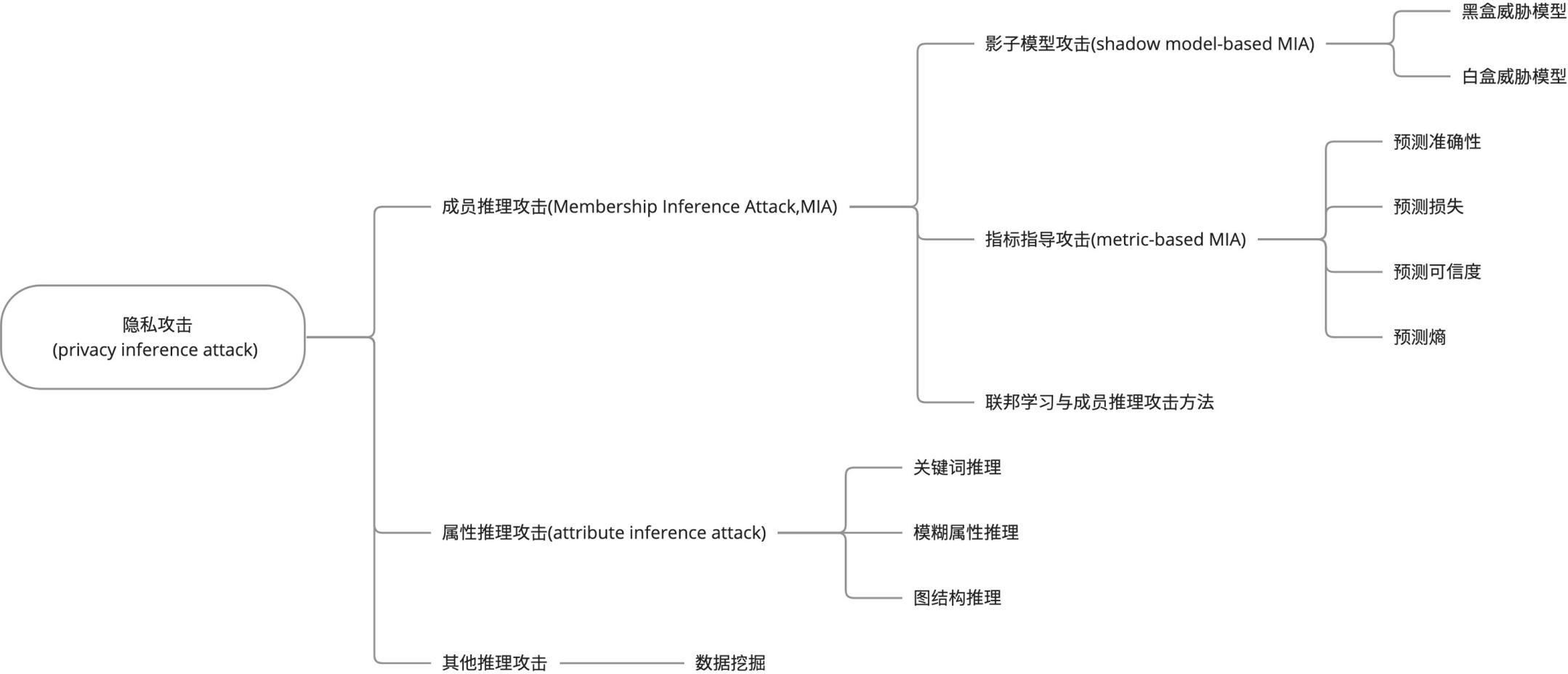
This Week

- ❑ Membership Inference Attack
- ❑ Differential Privacy

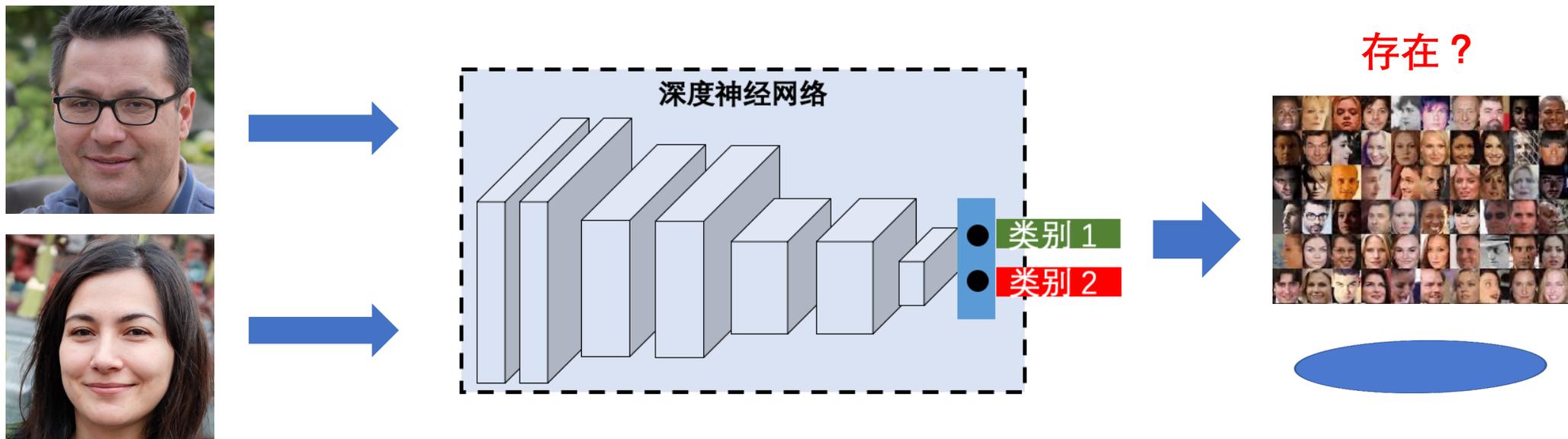
Membership Inference Attack

Differential Privacy

Membership Inference Attack



Membership Inference Attack



推理一个输入样本是否存在于训练数据集中

Privacy and Ethical Problems

- MIA could cause the following harms:
 - Leak private info: someone has been to some place or having an unspeakable illness
 - Expose info about the training data
 - MIA sensitivity also indicates data leakage risk

An Early Work

OPEN ACCESS Freely available online

PLOS GENETICS

Resolving Individuals Contributing Trace Amounts of DNA to Highly Complex Mixtures Using High-Density SNP Genotyping Microarrays

Nils Homer^{1,2}, Szabolcs Szelinger¹, Margot Redman¹, David Duggan¹, Waibhav Tembe¹, Jill Muehling¹, John V. Pearson¹, Dietrich A. Stephan¹, Stanley F. Nelson², David W. Craig^{1*}

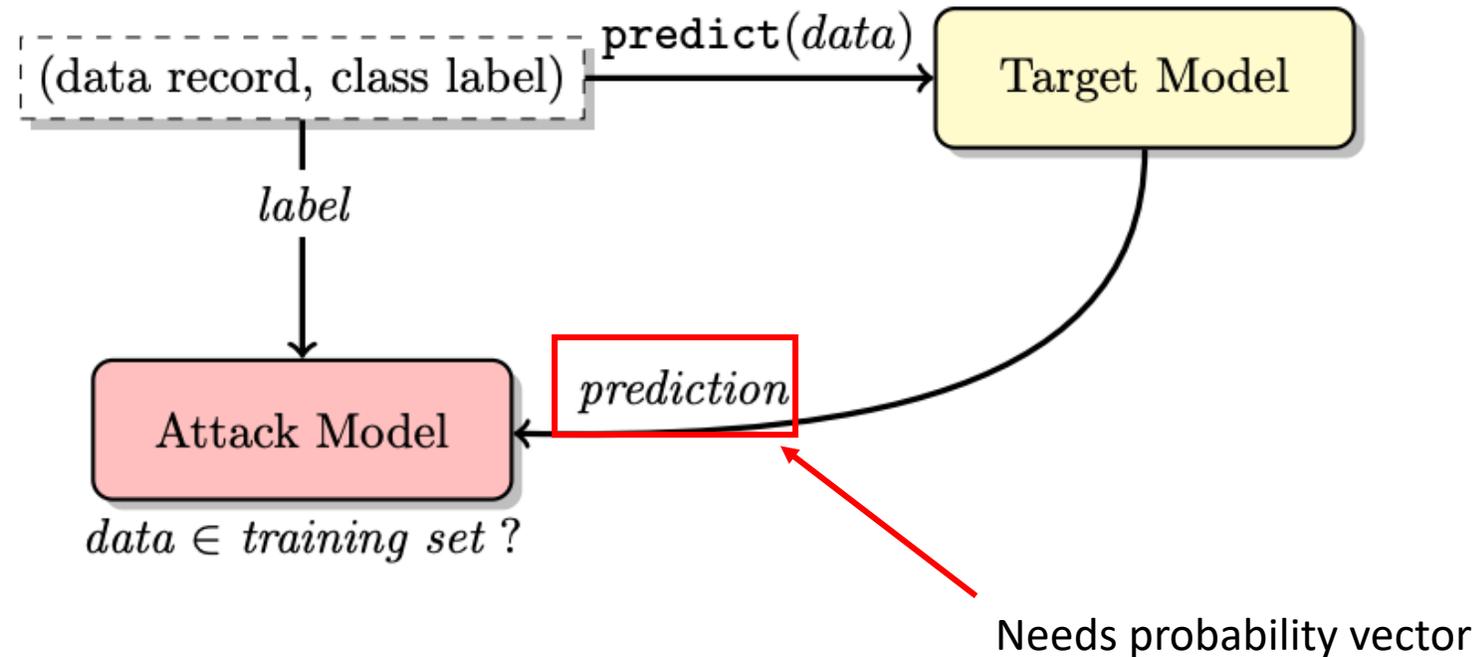
¹Translational Genomics Research Institute (TGen), Phoenix, Arizona, United States of America, ²University of California Los Angeles, Los Angeles, California, United States of America

Abstract

We use high-density single nucleotide polymorphism (SNP) genotyping microarrays to demonstrate the ability to accurately and robustly determine whether individuals are in a complex genomic DNA mixture. We first develop a theoretical framework for detecting an individual's presence within a mixture, then show, through simulations, the limits associated with our method, and finally demonstrate experimentally the identification of the presence of genomic DNA of specific individuals within a series of highly complex genomic mixtures, including mixtures where an individual contributes less than 0.1% of the total genomic DNA. These findings shift the perceived utility of SNPs for identifying individual trace contributors within a forensics mixture, and suggest future research efforts into assessing the viability of previously sub-optimal DNA sources due to sample contamination. These findings also suggest that composite statistics across cohorts, such as allele frequency or genotype counts, do not mask identity within genome-wide association studies. The implications of these findings are discussed.

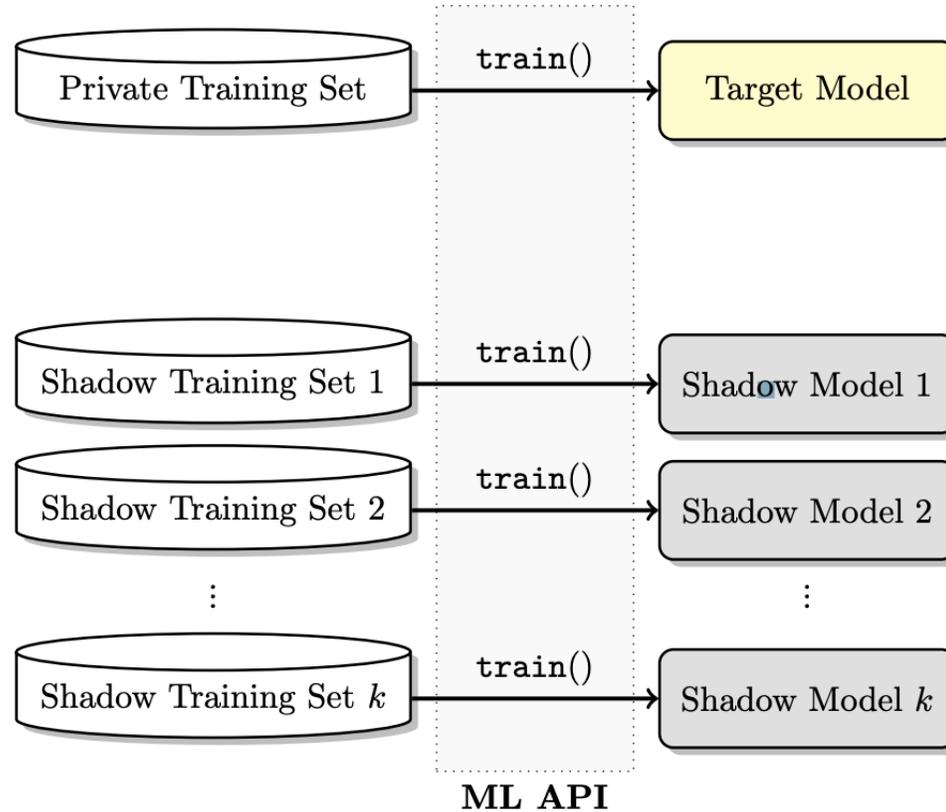
- 判断个人基因是否出现在一个复杂的混合基因里
- 可用于调查取证

MIA : The Most Well-known Work



Black-box attack pipeline

MIA : The Most Well-known Work



- ① Sample a number of subsets from D
- ② Train a model on each of the subset
- ③ Take one model as the target
- ④ Take the rest models as shadow models

Train k shadow models on disjoint datasets

MIA : The Most Well-known Work

Algorithm 1 Data synthesis using the target model

```
1: procedure SYNTHESIZE(class :  $c$ )
2:    $\mathbf{x} \leftarrow \text{RANDRECORD}()$   $\triangleright$  initialize a record randomly
3:    $y_c^* \leftarrow 0$ 
4:    $j \leftarrow 0$ 
5:    $k \leftarrow k_{max}$ 
6:   for iteration = 1  $\dots$  itermax do
7:      $\mathbf{y} \leftarrow f_{\text{target}}(\mathbf{x})$   $\triangleright$  query the target model
8:     if  $y_c \geq y_c^*$  then  $\triangleright$  accept the record
9:       if  $y_c > \text{conf}_{min}$  and  $c = \arg \max(\mathbf{y})$  then
10:        if  $\text{rand}() < y_c$  then  $\triangleright$  sample
11:          return  $\mathbf{x}$   $\triangleright$  synthetic data
12:        end if
13:      end if
14:       $\mathbf{x}^* \leftarrow \mathbf{x}$ 
15:       $y_c^* \leftarrow y_c$ 
16:       $j \leftarrow 0$ 
17:    else
18:       $j \leftarrow j + 1$ 
19:      if  $j > \text{rej}_{max}$  then  $\triangleright$  many consecutive rejects
20:         $k \leftarrow \max(k_{min}, \lceil k/2 \rceil)$ 
21:         $j \leftarrow 0$ 
22:      end if
23:    end if
24:     $\mathbf{x} \leftarrow \text{RANDRECORD}(\mathbf{x}^*, k)$   $\triangleright$  randomize  $k$  features
25:  end for
26:  return  $\perp$   $\triangleright$  failed to synthesize
27: end procedure
```

□ Different ways to get the training data :
Random Synthesis

□ Data synthesis

- **Phase 1:** searching for high confidence data points in the data space
- **Phase 2:** sample synthetic data from these points
- Repeat the above for each class c

Phase 1: 每次只改变已找到的高置信度样本的 k 个特征

MIA : The Most Well-known Work

Algorithm 1 Data synthesis using the target model

```
1: procedure SYNTHESIZE(class :  $c$ )
2:    $\mathbf{x} \leftarrow \text{RANDRECORD}()$   $\triangleright$  initialize a record randomly
3:    $y_c^* \leftarrow 0$ 
4:    $j \leftarrow 0$ 
5:    $k \leftarrow k_{max}$ 
6:   for iteration =  $1 \dots iter_{max}$  do
7:      $\mathbf{y} \leftarrow f_{target}(\mathbf{x})$   $\triangleright$  query the target model
8:     if  $y_c \geq y_c^*$  then  $\triangleright$  accept the record
9:       if  $y_c > conf_{min}$  and  $c = \arg \max(\mathbf{y})$  then
10:        if  $\text{rand}() < y_c$  then  $\triangleright$  sample
11:          return  $\mathbf{x}$   $\triangleright$  synthetic data
12:        end if
13:      end if
14:       $\mathbf{x}^* \leftarrow \mathbf{x}$ 
15:       $y_c^* \leftarrow y_c$ 
16:       $j \leftarrow 0$ 
17:    else
18:       $j \leftarrow j + 1$ 
19:      if  $j > rej_{max}$  then  $\triangleright$  many consecutive rejects
20:         $k \leftarrow \max(k_{min}, \lceil k/2 \rceil)$ 
21:         $j \leftarrow 0$ 
22:      end if
23:    end if
24:     $\mathbf{x} \leftarrow \text{RANDRECORD}(\mathbf{x}^*, k)$   $\triangleright$  randomize  $k$  features
25:  end for
26:  return  $\perp$   $\triangleright$  failed to synthesize
27: end procedure
```

□ Statistics-based synthesis

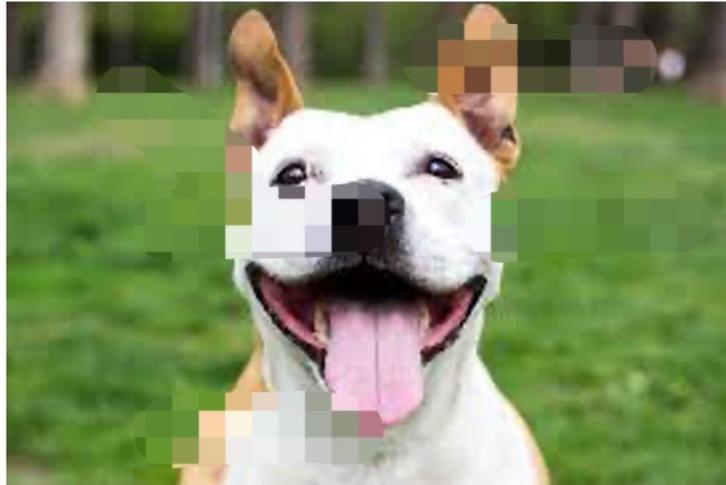
□ Prior knowledge:

- The marginal distribution w.r.t. each class

Phase 1: sample according to the statistics

MIA : The Most Well-known Work

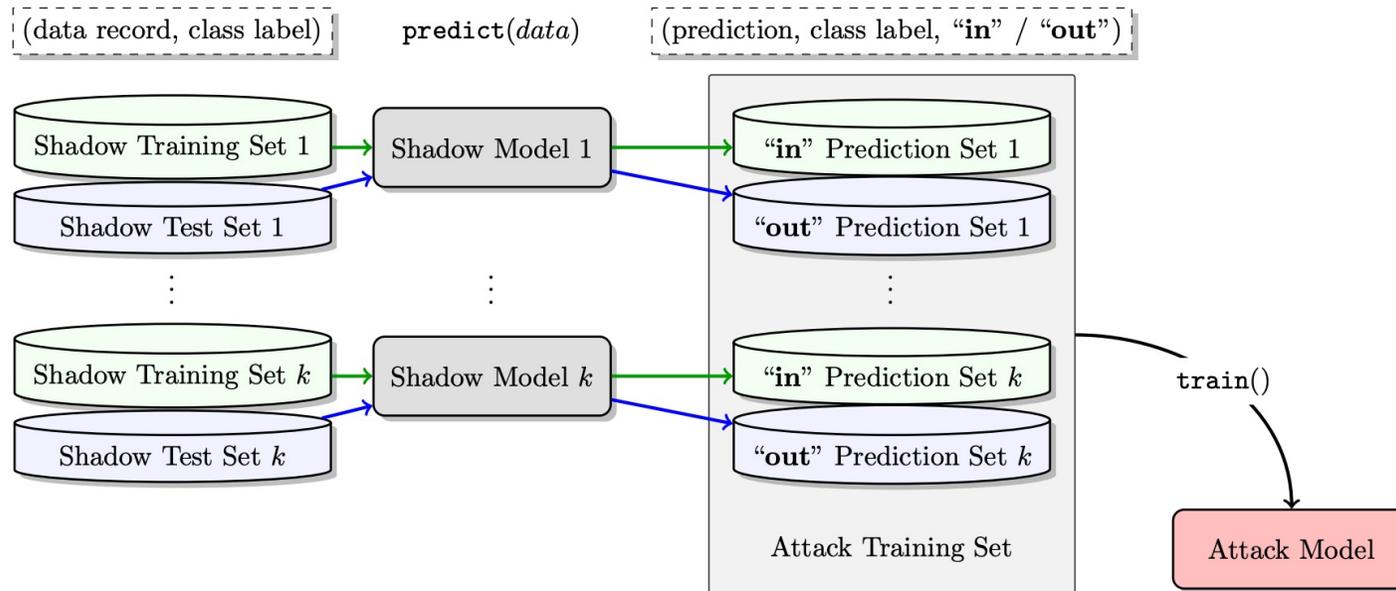
- We could also assume the attacker can access Noisy Real data: real but noisy



- Very similar to the real dataset
- But with a few features (10% or 20%) are randomly reset

MIA : The Most Well-known Work

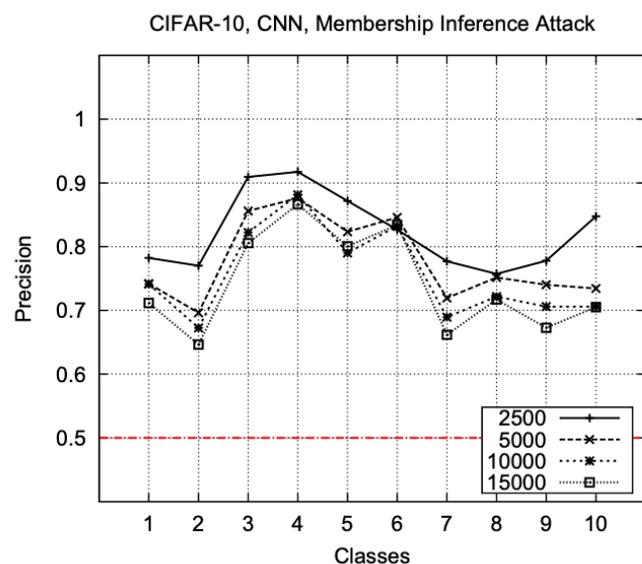
□ Finally: training the inference model



- "in": in the training set
- "out": : in the test set
- Train the inference model with dataset: **(prob1, "in"), (prob2, "in"), (prob3, "out") (prob4, "out")**

MIA : The Most Well-known Work

□ How well can MIA perform?



<i>Dataset</i>	<i>Training Accuracy</i>	<i>Testing Accuracy</i>	<i>Attack Precision</i>
Adult	0.848	0.842	0.503
MNIST	0.984	0.928	0.517
Location	1.000	0.673	0.678
Purchase (2)	0.999	0.984	0.505
Purchase (10)	0.999	0.866	0.550
Purchase (20)	1.000	0.781	0.590
Purchase (50)	1.000	0.693	0.860
Purchase (100)	0.999	0.659	0.935
TX hospital stays	0.668	0.517	0.657

数据集：CIFAR-10、CIFAR-100、Purchases、Locations、Texas hospital stays、MNIST、UCI Adult (Census Income).

White-box MIA

□ White-box vs Black-box

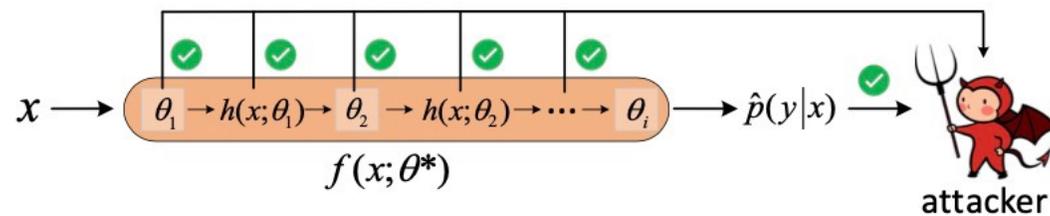


Fig. 2. Overview of white-box membership inference attacks.

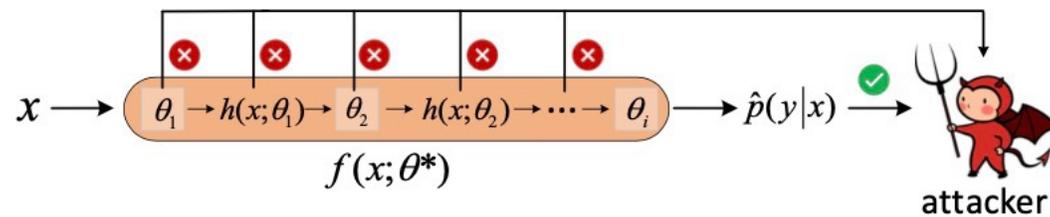
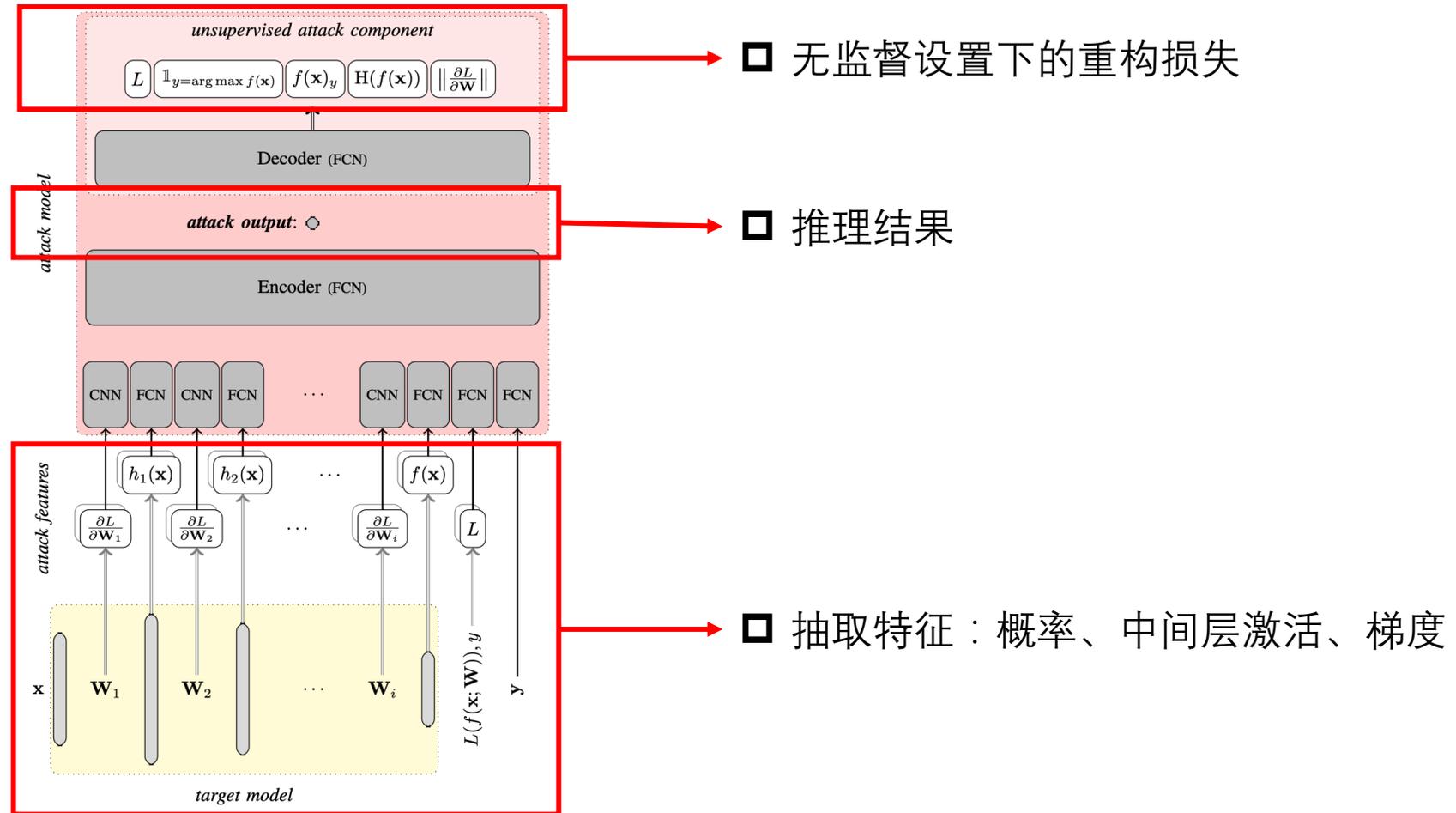


Fig. 3. Overview of black-box membership inference attacks.

White-box MIA



Limitations of MIA

- Constructing shadow models
- Assuming access to some data or prior knowledge
- **Overfitting is a must**
- Limited to classification models
- Limited to small models

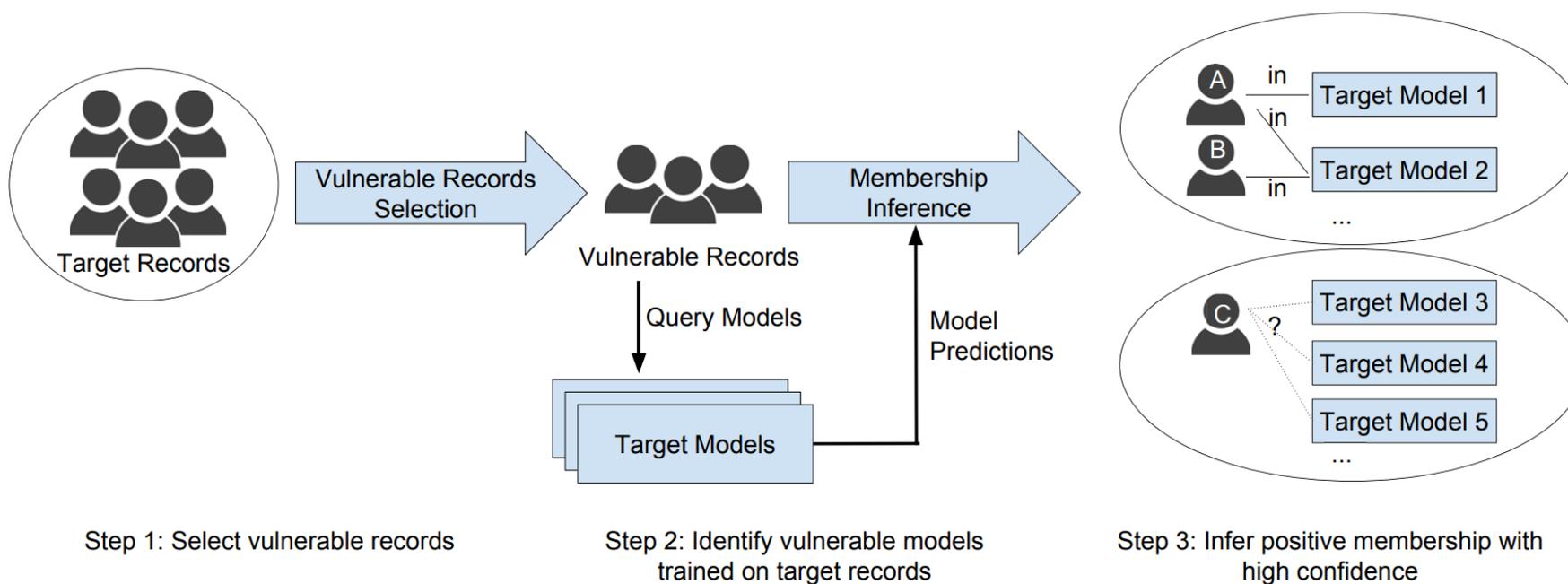
Addressing Limitations of MIA

□ Model and Data Independent MIA

Adversary type	Shadow model design		Target model's training data distribution
	No. shadow models	Target model structure	
Shokri et al. [38]	multiple	✓	✓
Our adversary 1	1	-	✓
Our adversary 2	1	-	-
Our adversary 3	-	-	-

Addressing Limitations of MIA

- ❑ Attacking non-overfitting DNNs
- ❑ Focusing on minimizing false positives



目标问题：样本A/B在哪个模型的训练数据里？

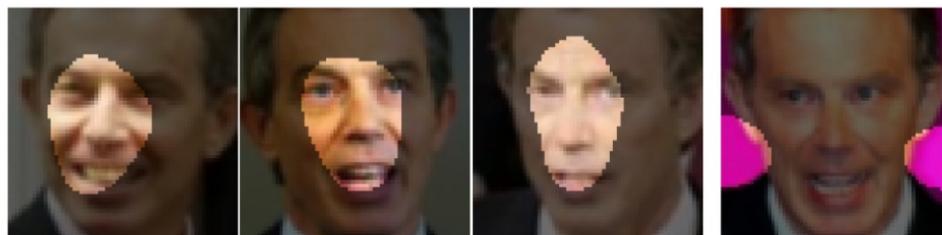
Addressing Limitations of MIA

- ❑ More practical white-box threat model
- ❑ The adversary only knows the model but not the data distribution



Training images

(a)



(b)

(c)

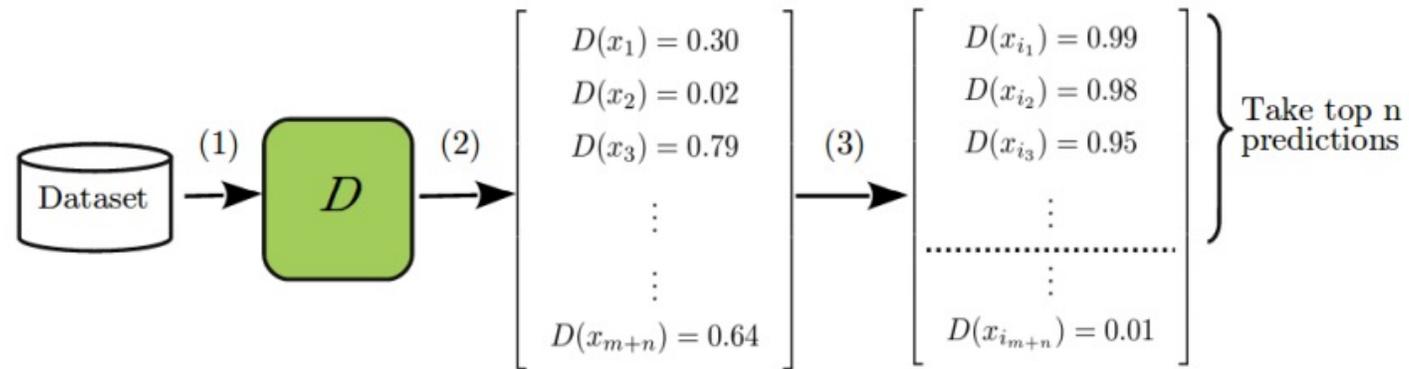
Internal explanations

Pink background explanation of Tony Blair

利用诡异的独家记忆进行成员推理

Addressing Limitations of MIA

Extension to generative models



充分利用判别器的判别能力：高置信度的大概率来自原始训练数据集

Metric-guided MIA

□ Metric based Anomaly detection



- 预测正确性 : $\mathcal{M}(\hat{\mathbf{p}}(y|\mathbf{x}), y) = \mathbb{1}[\arg \max \hat{\mathbf{p}}(y|\mathbf{x}) = y]$ 预测正确的就是成员
- 预测损失 : $\mathcal{M}(\hat{\mathbf{p}}(y|\mathbf{x}), y) = \mathbb{1}[\mathcal{L}(\hat{\mathbf{p}}(y|\mathbf{x}); y) \leq \tau]$ 高于训练样本平均损失的是成员
- 预测置信度 : $\mathcal{M}(\hat{\mathbf{p}}(y|\mathbf{x})) = \mathbb{1}[\max \hat{\mathbf{p}}(y|\mathbf{x}) \geq \tau]$ 有概率接近1的是成员
- 预测熵 : $\mathcal{M}(\hat{\mathbf{p}}(y|\mathbf{x})) = \mathbb{1}[\mathbb{H}(\hat{\mathbf{p}}(y|\mathbf{x})) \leq \tau] = \mathbb{1}[-\sum_i \mathbf{p}_i \log(\mathbf{p}_i) \leq \tau]$ 低概率熵的是成员
- 修正预测熵 : $\text{MH}(\hat{\mathbf{p}}(y|\mathbf{x}), y) = -(1 - \mathbf{p}_y) \log(\mathbf{p}_y) - \sum_{i \neq y} \mathbf{p}_i \log(1 - \mathbf{p}_i)$ 不同类别区别考虑

A Summary of Existing MIAs

□ Used Datasets

- **Image:**
 - CIFAR-10, CIFAR-100, MNIST, Fashion-MNIST, Yale Face, ChestX-ray8, SVHN, CelebA, ImageNet
- **Tabulate:**
 - Adult, Foursquare, Purchase-100, Texas100, Location, etc.
- **Audio:**
 - LibriSpeech, TIMIT, TED
- **Text:**
 - Weibo, Tweet EmoInt, SATED, Dislogs, Reddit comments, Cora, Pubmed, Citesser

A Summary of Existing MIAs

- **Target models:**

- On **image**:

- Multi-layer CNN + 1 or 2 FC (> 5 papers used 2-4 layers CNN)
 - Alexnet, ResNet18, ResNet50, VGG16, VGG19, DenseNet121, Efficient-netv2, EfficientNetB0
 - GAN: InfoGAN, PGGAN, WGANGP, DCGAN, MEDGAN, and VAEGAN

- On **tabulate data**:

- FC only models

- On **text**:

- Multi-layer CNN, multi-layer RNN/LSTM, transformers (e.g., BERT, GPT-2)

- On **audio**:

- Hybrid system: HMM-DNN model
 - End-to-end: Multi-layer LSTM/ RNN/GRU

- **MLaaS (Online):**

- Google Prediction API, Amazon ML

Membership Inference Attack

Differential Privacy

Differential Privacy

□ Finite Difference and Derivative

$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h} \quad h \text{ tends to be small (zero)}$$

通过函数在某一点随微小扰动的变化可以估计在这一点的梯度

如果对数据集进行微小扰动呢？

Differential Privacy

□ Finite Difference -> Differential Privacy

算法/机制 \mathcal{M}

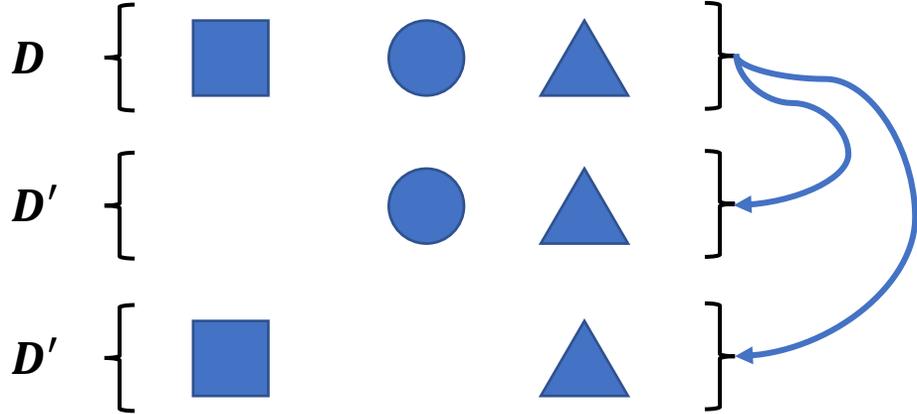
$f(x)$ 函数 \longrightarrow 算法/机制 \mathcal{M}

a 输入值 \longrightarrow 数据集 D

数据集的微小变化会导致多大的算法输出变化？

Differential Privacy

□ 邻接数据集 D 、 D'



数据集的微小变化会导致多大的算法输出变化？

Differential Privacy

定义 5.1. 差分隐私: 对于一个随机算法 M , P_m 为算法 M 所有可能输出的集合, 若算法 M 满足 $(\epsilon, \delta) - DP$, 当且仅当相邻数据集 D, D' 对 M 的所有可能输出子集 $S_m \in P_m$, 满足不等式 [Dwork et al., 2006a]:

$$P_r[M(D) \in S_m] \leq e^\epsilon P_r[M(D') \in S_m] + \delta$$

ϵ : 隐私预算 (Privacy Budget), 越小隐私越好

$e^\epsilon \approx 1 + \epsilon$ (for small ϵ), $\epsilon=1, 2, 3 < 10$ is reasonable

$\delta < 1/n$ (n is the dataset size) : 打破 $(\epsilon, \delta) - DP$ 的可能性, probability of (potential) privacy failure

$\delta=0$: pure differential privacy

$\delta>0$: approximate differential privacy

Properties of DP

性质 5.1. 顺序合成: 给定 K 个随机算法 $M_i (i = 1, \dots, K)$, 分别满足 $\epsilon_i - DP$, 如果将他们作用在同一个数据集上, 则满足 $\sum_{i=1}^K \epsilon_i - DP$ 。

性质 5.2. 平行合成: 将数据集 D 分割成 K 个不相交的子集 $\{D_1, D_2, \dots, D_K\}$, 在每个子集上分别作用满足 $\epsilon_i - DP$ 的随机算法 M_i , 则数据集 D 整体满足 $(\max\{\epsilon_1, \dots, \epsilon_K\}) - DP$ 。

性质 5.3. 交换不变性: 给定任意算法 M_1 满足 $\epsilon - DP$, 数据集 D , 对于任意算法 M_2 (M_2 不一定满足差分隐私), 则 $M_2(M_1(D))$ 满足 $\epsilon - DP$ 。

性质 5.4. 中凸性: 给定满足 $\epsilon - DP$ 的随机算法 M_1 和 M_2 , 对于任意的概率 $P \in [0, 1]$, 用 A_P 表示一种选择机制, 以 P 的概率选择算法 M_1 , 以 $1 - P$ 的概率选择算法 M_2 , 则 A_P 机制满足 $\epsilon - DP$ 。

How to Obtain a Differentially Private Model

思考：



Measuring Sensitivity

定义 5.2. 全局敏感度 (Global Sensitivity): 给定查询函数 $f : D \rightarrow R$, D 为数据集, R 为查询结果。在任意一对相邻数据集 D, D' 上, 全局敏感度定义为:

$$S(f) = \max_{D, D'} \|f(D) - f(D')\|_1$$

定义 5.3. 局部敏感度 (Local Sensitivity): 给定查询函数 $f : D \rightarrow R$, D 为数据集, R 为查询结果。在一给定的数据集 D 和它相邻的任意数据集 D' 上, 局部敏感度定义为:

$$LS(f) = \max_{D'} \|f(D) - f(D')\|_1$$



Noise Models

□ 几种噪声添加机制

- 拉普拉斯机制 (Laplacian)

$$M(D) = f(D) + \text{Lap}\left(\frac{S(f)}{\epsilon}\right) \quad \text{Lap}\left(\frac{S(f)}{\epsilon}\right) \text{ 表示位置参数为 } 0, \text{ 尺度参数为 } \frac{S(f)}{\epsilon} \text{ 的拉普拉斯分布}$$

- 高斯机制 (Gaussian)

$$M(D) = f(D) + \mathcal{N}(\delta^2) \\ \text{s.t. } \delta^2 = \frac{2S(f)^2 \log(1.25/\delta)}{\epsilon^2} \quad \mathcal{N}(\delta^2) \text{ 表示中心为 } 0, \text{ 方差为 } \delta^2 \text{ 的高斯分布}$$

- 指数机制：离散 \rightarrow 概率；确定 \rightarrow 不确定 $M(D) = \text{return}(R_i \propto \exp(\frac{\epsilon q(D, R_i)}{2S(q)})) \quad \text{Pr}(R_i) = \frac{\exp(\frac{\epsilon q(D, R_i)}{2S(q)})}{\sum_{j=1}^N \exp(\frac{\epsilon q(D, R_j)}{2S(q)})}$

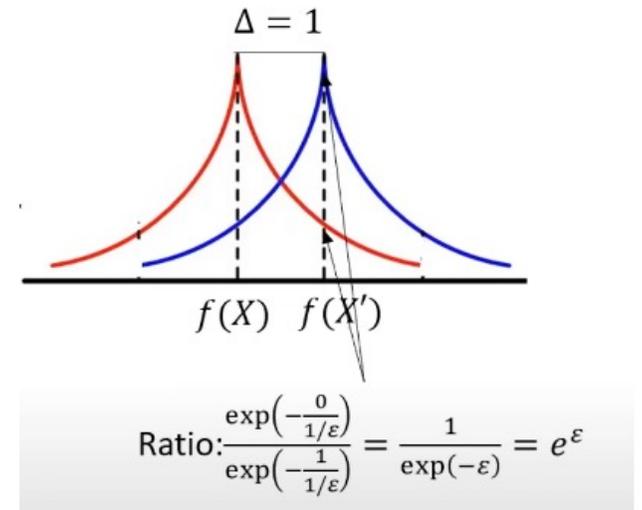
The Laplace Mechanism

- 拉普拉斯机制 (Laplace Mechanism)

$M(D) = f(D) + \text{Lap}\left(\frac{S(f)}{\epsilon}\right)$ $\text{Lap}\left(\frac{S(f)}{\epsilon}\right)$ 表示位置参数为 0, 尺度参数为 $\frac{S(f)}{\epsilon}$ 的拉普拉斯分布

Example:

- Given $X_1, \dots, X_n \in \{0,1\}$
- Goal: privately compute sum $f(X) = \sum_{i=1}^n X_i$
- Sensitivity of $f(X)$ is $\Delta = 1$
- Claim: $f(X) + Z$, where $Z \sim \text{Laplace}(\sigma = \Delta/\epsilon)$ is $(\epsilon, 0)$ - DP
- $\text{Laplace}(\sigma = \Delta/\epsilon) \propto \exp(-|x|/\sigma)$, two-sided exponential distribution



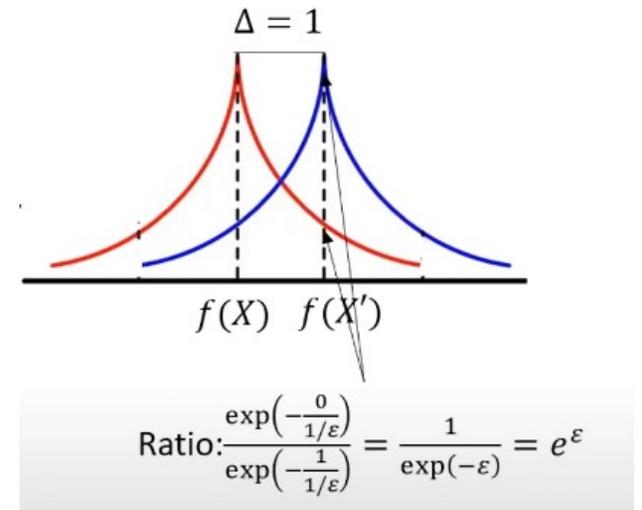
The Laplace Mechanism

- 拉普拉斯机制 (Laplace Mechanism)

$M(D) = f(D) + \text{Lap}\left(\frac{S(f)}{\epsilon}\right)$ $\text{Lap}\left(\frac{S(f)}{\epsilon}\right)$ 表示位置参数为 0, 尺度参数为 $\frac{S(f)}{\epsilon}$ 的拉普拉斯分布

Example: Counting queries

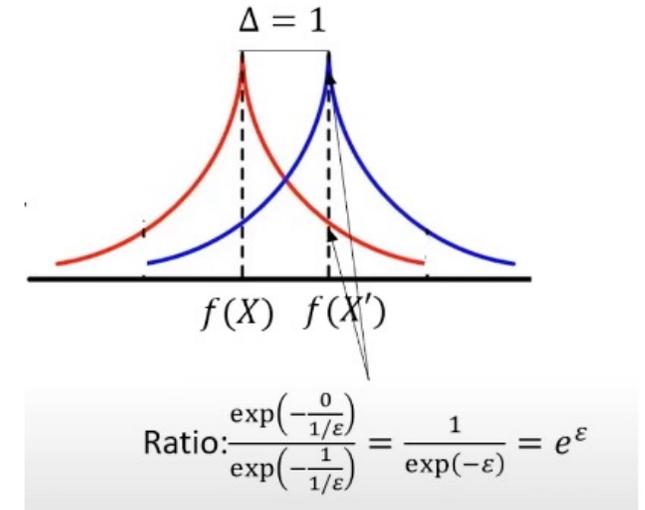
- Asking 100 people, how many smokers?
- E.g., 20 say yes
- Can guarantee $(1,0)$ – DP by adding Laplace $\left(1 = \sigma = \frac{\Delta}{\epsilon}\right)$ noise
 - Sampled outputs: 20.68, 19.24, 20.28, 19.83
- Stronger privacy: $(0.1,0)$ – DP by adding Laplace $\left(10 = \sigma = \frac{\Delta}{\epsilon}\right)$ noise
 - Sampled outputs: 22.45, 11.45, 2.4, 15.03, 29.47



The Laplace Mechanism

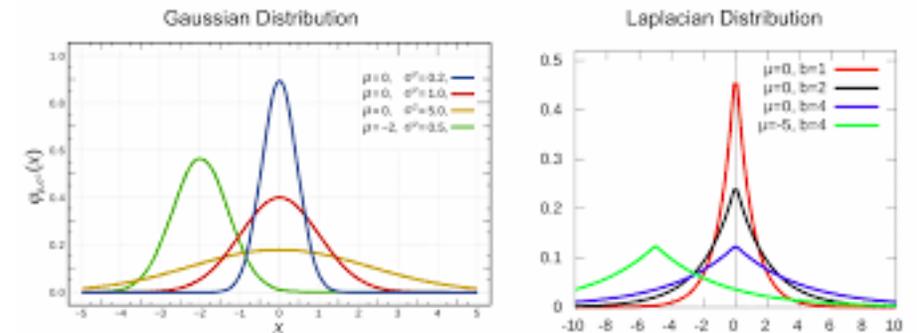
Given a function $f(X): \mathcal{X}^n \rightarrow \mathbb{R}^d$

- Let $\Delta_1^f = \max_{X, X' \text{ differ in one entry}} \|f(X) - f(X')\|_1$ be l_1 - sensitivity of f : how much can the function change by modifying one data point
- Theorem: the Laplace Mechanism $f(X) + \text{Lap}(\Delta_1^f / \epsilon)^{\otimes d}$ is $(\epsilon, 0)$ -DP
 - Adding Laplace noise to each coordinate, proportional to the l_1 - sensitivity



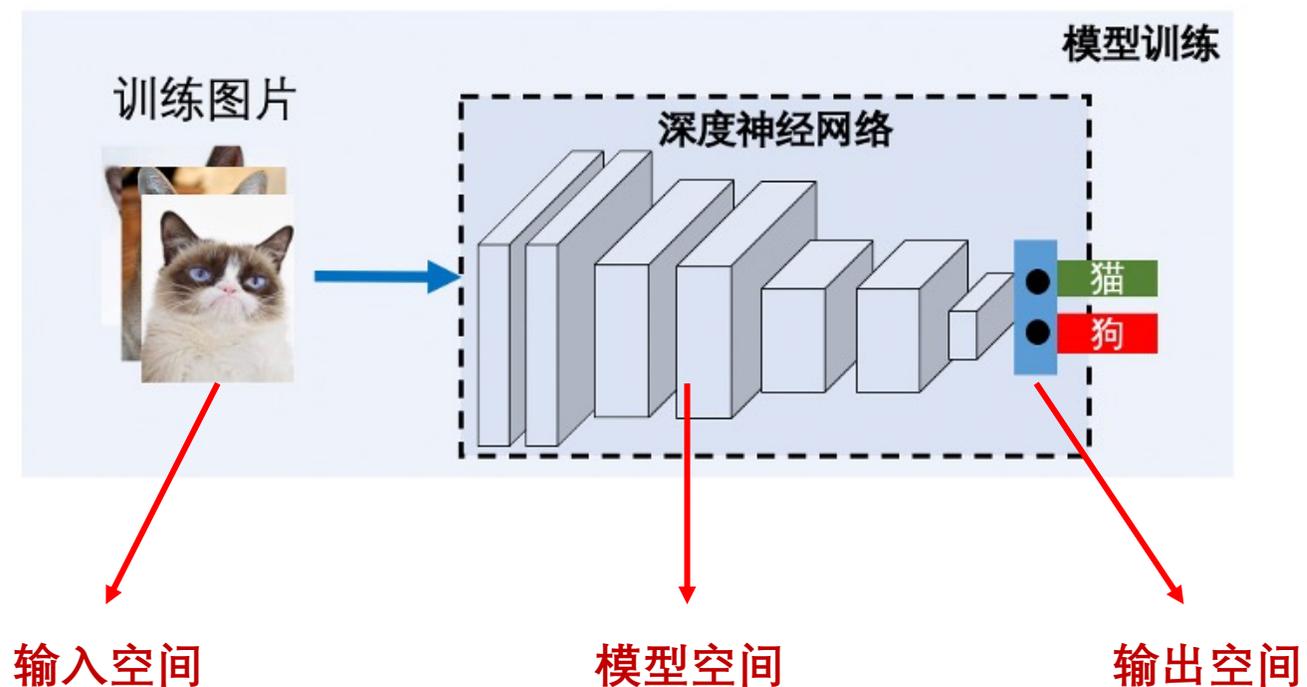
Laplace vs. Gaussian

- Both add noise to $f(X)$
- Gaussian often adds less noise than Laplace
 - As sensitivity is based on l_2 - sensitivity $< l_1$ - sensitivity
- But Gaussian gives (ϵ, δ) -DP rather than Laplace's $(\epsilon, 0)$ -DP
 - Weaker privacy
- Most of the time, (ϵ, δ) -DP is good enough
 - Necessary if you're doing a lot of queries on the same datasets
 - "Advanced composition"



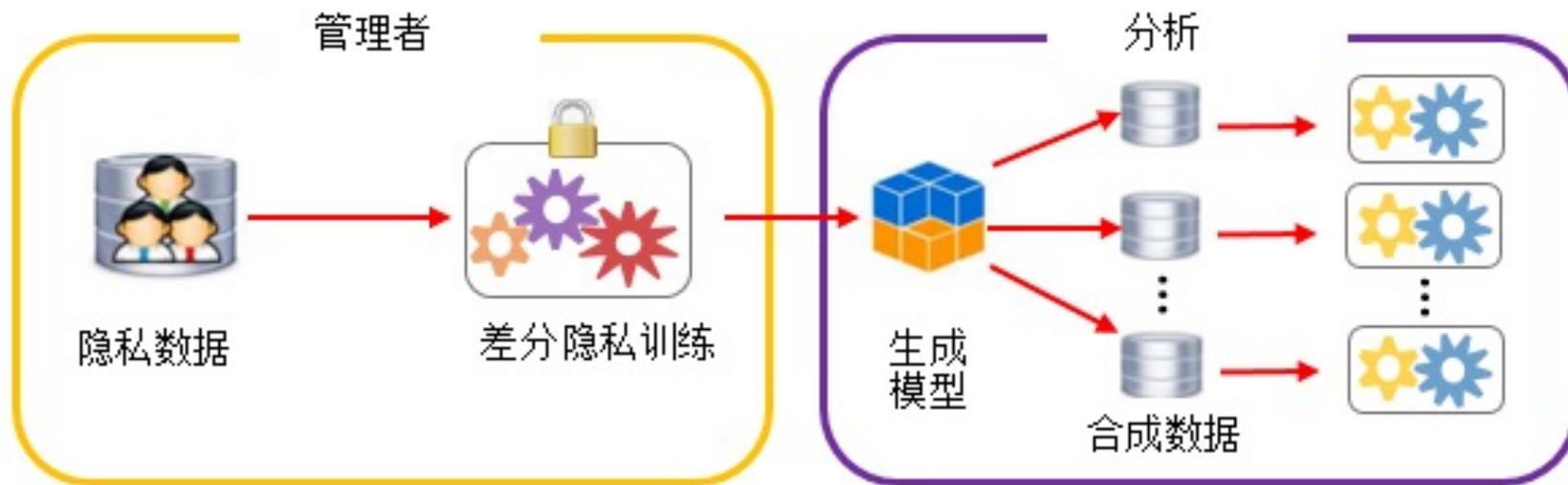
DP + Deep Learning

□ 问题：在哪里添加噪声？



输入空间DP

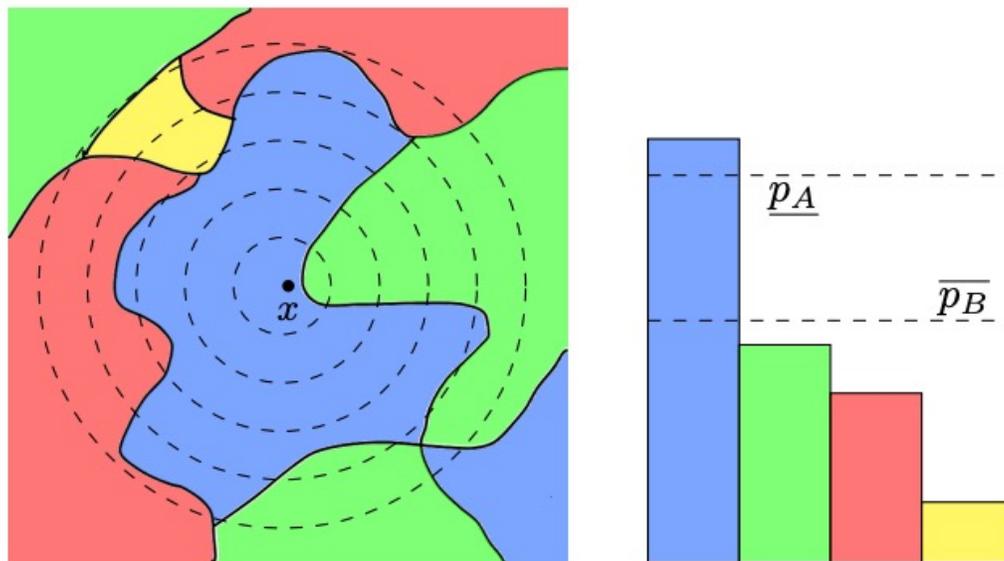
□ 差分隐私预处理训练数据



dp-GAN pipeline

输入空间DP

□ 随机平滑 Randomized Smoothing



用随机噪声填充输入空间，
得到对抗鲁棒性边界

随机平滑：可验证对抗防御

□ 差分隐私平滑模型参数：DP-SGD算法

Algorithm 5.1 Differentially Private SGD (DP-SGD) [Abadi et al., 2016]

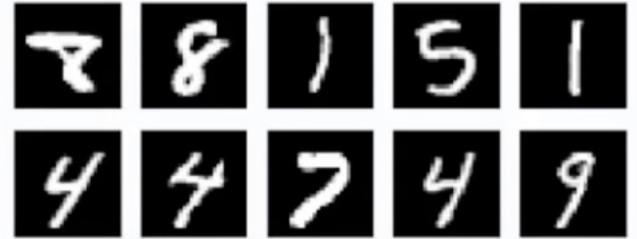
输入: 样本 $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, 损失函数 $\mathcal{L}(\theta) = \frac{1}{n} \sum_i \mathcal{L}(\theta, \mathbf{x}_i)$ 。超参数: 学习率 η_t , 噪声参数 σ , 分组大小 L , 梯度约束范数 C

输出: θ_T , 同时利用隐私统计方法计算总体的隐私损失 (ϵ, δ)

- 1: 随机初始化模型 θ_0
 - 2: **for** $t \in [T]$ **do**
 - 3: 以概率 L/n 随机采取一组样本 L_t
 - 4: **计算梯度:** 对每一个样本 $i \in L_t$, 计算 $g_t(\mathbf{x}_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, \mathbf{x}_i)$
 - 5: **裁剪梯度:** $\bar{g}_t(\mathbf{x}_i) \leftarrow g_t(\mathbf{x}_i) / \max(1, \frac{\|g_t(\mathbf{x}_i)\|_2}{C})$
 - 6: **噪声添加:** $\bar{g}_t \leftarrow \frac{1}{L} (\sum_i \bar{g}_t(\mathbf{x}_i) + \mathcal{N}(0, \sigma^2 C^2 I))$
 - 7: **梯度下降:** $\theta_{t+1} \leftarrow \theta_t - \eta_t \bar{g}_t$
-

DP-SGD性能

- MNIST: black and white image classification
 - Canonical “easy” ML task
- Non-private test accuracy: $\approx 100\%$
- Private (ϵ from 1 to 3): **98% - 99%**
 - [Tramer-Boneh, '21]



DP-SGD性能

- CIFAR-10: Low resolution images
 - Same size as MNIST, but harder
- Non-privately: 98%+
- Privately ($\epsilon = 3$): 69%
 - [Tramer-Boneh, '21]
 - Much worse!
- Recent results: 73.5% for $\epsilon = 4$ and 82.5% for $\epsilon = 8$
 - [De-Berrada-Hayes-Smith-Balle, '22], [Klaue-Ziller-Rueckert-Hammernik-Kaissis, '22]

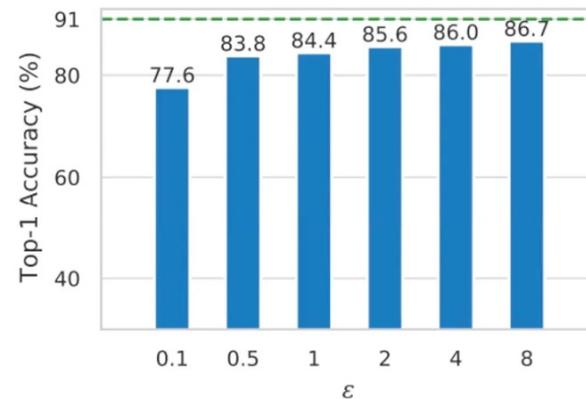


More Practical Solution?

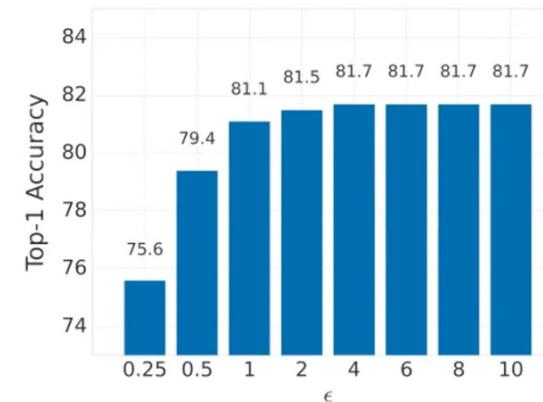
- 1: Training on public data like ImageNet
- 2: Finetuning with DP-SGD on private data

Using public data (ImageNet)

- [De, Berrada, Hayes, Smith, Balle, '22]
 - Pretrain with JFT-4B



- [Mehta, Thakurta, Kurakin, Cutkosky, '22]
 - Pretrain with JFT-4B



输出空间DP

□ 差分隐私扰动目标函数：多项式目标函数

- 回归模型
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_{i=1}^n \mathcal{L}(t_i, \mathbf{w})$$

- 根据Stone-Weierstrass 理论：*任意连续可微的函数可表示为：*

$$\mathcal{L}_D(\mathbf{w}) = \sum_{j=0}^J \sum_{\phi \in \Phi_j} \lambda_{\phi t_i} \sum_{t_i \in D} \phi(\mathbf{w})$$

□ 差分隐私扰动目标函数：多项式目标函数

Algorithm 5.2 函数机制 (Functional Mechanism) [Zhang et al., 2012]

输入: 数据集 D , 目标函数 $\mathcal{L}_D(\mathbf{w})$, 隐私预算 ϵ

输出: 差分隐私扰动后的模型参数 $\bar{\mathbf{w}}$

- 1: 令 $\Delta = 2 \max_t \sum_{j=1}^J \sum_{\phi \in \Phi_j} \|\lambda_{\phi t}\|_1$
 - 2: **for** $0 \leq j \leq J$ **do**
 - 3: **for** $\phi \in \Phi_j$ **do**
 - 4: 令 $\lambda_\phi = \sum_{t_i \in D} \lambda_{\phi t_i} + \text{Laplace}(\frac{\Delta}{\epsilon})$
 - 5: 令 $\bar{\mathcal{L}}_D(\mathbf{w}) = \sum_{j=1}^J \sum_{\phi \in \Phi_j} \lambda_\phi \phi(\mathbf{w})$
 - 6: 计算 $\bar{\mathbf{w}} = \arg \min_{\mathbf{w}} \bar{\mathcal{L}}_D(\mathbf{w})$
 - 7: 返回 $\bar{\mathbf{w}}$
-

输出空间DP

□ 差分隐私扰动目标函数 : cross-entropy

$$\tilde{f}_D(\omega) = \sum_{i=1}^{|D|} \sum_{l=1}^m \sum_{R=0}^{\infty} \frac{f_l^{(R)}(z_l)}{R!} (g_l(t_i, \omega) - z_l)^R$$

泰勒展开 Taylor Expansion

Remaining Challenges

□ Attack:

- Better Performance Metrics for MIA
- Attacking large-scale pretrained models

□ Defense:

- How to achieve both accuracy and privacy
- How to detect potential MIAs on the fly



谢谢！

