





Membership and Property Inference Attacks Against Machine Learning

Emiliano De Cristofaro https://emilianodc.com



Most papers on privacy attacks in ML focus on inferring:



Most papers on privacy attacks in ML focus on inferring:

1. Inclusion of a data point in the training set (aka "membership inference")



Most papers on privacy attacks in ML focus on inferring:

- Inclusion of a data point in the training set (aka "membership inference")
- 2. What class representatives look like



Adversary wants to test whether data of a target victim has been used to train a model

Adversary wants to test whether data of a target victim has been used to train a model

Serious problem if inclusion in training set is privacy-sensitive

Adversary wants to test whether data of a target victim has been used to train a model

Serious problem if inclusion in training set is privacy-sensitive

E.g., main task is: predict whether a smoker gets cancer

Adversary wants to test whether data of a target victim has been used to train a model

Serious problem if inclusion in training set is privacy-sensitive E.g., main task is: predict whether a smoker gets cancer [Shokri et al., S&P'17] show it for discriminative models

Adversary wants to test whether data of a target victim has been used to train a model

Serious problem if inclusion in training set is privacy-sensitive E.g., main task is: predict whether a smoker gets cancer [Shokri et al., S&P'17] show it for discriminative models [Hayes et al. PETS'19] for generative models (later in the talk)

Adversary wants to test whether data of a target victim has been used to train a model

Serious problem if inclusion in training set is privacy-sensitive E.g., main task is: predict whether a smoker gets cancer [Shokri et al., S&P'17] show it for discriminative models

[Hayes et al. PETS'19] for generative models (later in the talk)

Membership inference is a very active research area, not only in machine learning...

Membership inference is a very active research area, not only in machine learning...

Membership inference is a very active research area, not only in machine learning...

Given f(data), infer if $x \in$ data (e.g., f is aggregation)

Membership inference is a very active research area, not only in machine learning...

Given f(data), infer if $x \in$ data (e.g., f is aggregation) [Homer et al., Science'13] for genomic data [Pyrgelis et al., NDSS'18] for mobility data

Membership inference is a very active research area, not only in machine learning...

Given f(data), infer if $x \in$ data (e.g., f is aggregation) [Homer et al., Science'13] for genomic data [Pyrgelis et al., NDSS'18] for mobility data

Well-understood problem, besides the more obvious leakage

Membership inference is a very active research area, not only in machine learning...

Given f(data), infer if $x \in$ data (e.g., f is aggregation) [Homer et al., Science'13] for genomic data [Pyrgelis et al., NDSS'18] for mobility data

Well-understood problem, besides the more obvious leakage

- Establish wrongdoing
- Assess protection, e.g., from differentially private defenses

Prior work shows how infer properties of an entire class, e.g.:

Prior work shows how infer properties of an entire class, e.g.: Model Inversion [Fredrikson et al. CCS'15]

Prior work shows how infer properties of an entire class, e.g.: Model Inversion [Fredrikson et al. CCS'15] GAN attacks [Hitaji et al. CCS'17]

Prior work shows how infer properties of an entire class, e.g.: Model Inversion [Fredrikson et al. CCS'15] GAN attacks [Hitaji et al. CCS'17]

E.g.: given a gender classifier, infer what a male looks like

Prior work shows how infer properties of an entire class, e.g.: Model Inversion [Fredrikson et al. CCS'15] GAN attacks [Hitaji et al. CCS'17]

E.g.: given a gender classifier, infer what a male looks like

But...any useful machine learning model does reveal something about the population from which the training data was sampled

Prior work shows how infer properties of an entire class, e.g.: Model Inversion [Fredrikson et al. CCS'15]

GAN attacks [Hitaji et al. CCS'17]

E.g.: given a gender classifier, infer what a male looks like

But...any useful machine learning model does reveal something about the population from which the training data was sampled

Privacy leakage != Adv learns something about training data

Prior work shows how infer properties of an entire class, Model Inversion [Fredrikson et al. CCS'15] GAN attacks [Hitaji et al. CCS'17]

E.g.: given a gender classifier, infer what a male looks like

But...any useful machine learning model does reveal something about the population from which the training data was sampled

Privacy leakage != Adv learns something about training data







... but not of the whole class?



... but not of the whole class?



... but not of the whole class?

In a nutshell: given a gender classifier, infer race of people in Bob's photos

Agenda

1. Property Inference in Collaborative/Federated ML

2. Membership Inference against Generative Models



1. Property Inference in Collaborative/Federated ML

2. Membership Inference against Generative Models



1. Property Inference in Collaborative/Federated ML

2. Membership Inference against Generative Models

Luca Melis, Congzheng Song, Emiliano De Cristofaro, Vitaly Shmatikov. Exploiting Unintended Feature Leakage in Collaborative Learning. IEEE Symposium on Security & Privacy (S&P'19)

8

Deep Learning



• Map input x to layers of hidden representations h, then to output y

•
$$h_{l+1} = a(W_l \cdot h_l)$$
 with parameter W_l

- Train model to minimizes loss: $W = \operatorname{argmin}_W L(f(x), y)$
- Gradient descent on parameters:
 - Each iteration train on a batch
 - Update W based on gradient of L

Collaborative/Federated Learning



Collaborative

Federated

Algorithm 1 Parameter server with synchronized SGD	Algorithm 2 Federated learning with model averaging
Server executes:	Server executes:
Initialize θ_0	Initialize θ_0
for $t = 1$ to T do	$m \leftarrow max(C \cdot K, 1)$
for each client k do	for $t = 1$ to T do
$g_t^k \leftarrow \textbf{ClientUpdate}(heta_{t-1})$	$S_t \leftarrow (random set of m clients)$
end for	for each client $k \in S_t$ do
$ heta_t \leftarrow heta_{t-1} - \eta \sum_k g_t^k$	$\theta_t^k \leftarrow \text{ClientUpdate}(\theta_{t-1})$
end for	end for
ClientUpdate(θ): Select batch <i>b</i> from client's data	$ heta_t \leftarrow \sum_k rac{n^k}{n} heta_t^k$ end for

Select batch o from client's data **return** local gradients $\nabla L(b; \theta)$

ClientUpdate(θ):

for each local iteration do for each batch b in client's split do $\theta \leftarrow \theta - \eta \nabla L(b; \theta)$ end for end for **return** local model θ

Passive Property Inference Attack



Active Property Inference Attack



Dataset	Туре	Main Task	Inference Task	
LFW	Images	Gender/Smile/Age Eyewear/Race/Hair	Race/Eyewear	
FaceScrub	Images	Gender	Identity	
PIPA	Images	Age	Gender	
FourSquare	Locations	Gender	Membership	
Yelp-health	Text	Review Score	Membership Doctor specialty	
Yelp-author	Text	Review Score	Author	
CSI	Text	Sentiment	Membership Region/Gender/Veracity	

Property Inference on LFW



Multi-Party

Feature t-SNE projection





Passive vs Active Attack on FaceScrub

Main Task: $\blacktriangle / \bullet =$ female/male Inference Task: Blue points with the property (identity)



Inferring when a property occurs

Inferring when a property occurs

Batches with the property appear



Main task: Age / Two-party Inference task: people in the image are of the same gender (PIPA)

Inferring when a property occurs

Batches with the property appear

Participant with ID 1 joins training



Main task: Age / Two-party Inference task: people in the image are of the same gender (PIPA) Main task: Gender / Multi-Party Inference task: author identification

Defenses?

Defenses?

Selective gradient sharing

Dataset: Text reviews

Main Task: Sentiment classifier

Doesn't really work...

Property / % parameters shared	10%	50%	100%
Top region	0.84	0.86	0.93
Gender	0.90	0.91	0.93
Veracity	0.94	0.99	0.99

Defenses?

Selective gradient sharing Dataset: Text reviews Main Task: Sentiment classifier Doesn't really work...

Property / % parameters shared	10%	50%	100%
Top region	0.84	0.86	0.93
Gender	0.90	0.91	0.93
Veracity	0.94	0.99	0.99

Participant-level differential privacy

Hide participant's contributions

Only two mechanisms in the literature

Fail to converge for "few" participants

Agenda

1. Property Inference in Collaborative/Federated ML

2. Membership Inference against Generative Models

Agenda

1. Property Inference in Collaborative/Federated ML

2. Membership Inference against Generative Models

Machine Learning as a Service





Membership Inference in Generative Models

Membership Inference in Generative Models



Membership Inference in Generative Models



Jamie Hayes, Luca Melis, George Danezis, Emiliano De Cristofaro. LOGAN: Membership Inference Attacks Against Generative Models. PETS 2019.

Inference without predictions?

Use generative models!

Train GANs to learn the distribution and a prediction model at the same time

Inference without predictions?

Use generative models!

Train GANs to learn the distribution and a prediction model at the same time



White-Box Attack



Black-Box Attack



Datasets

Models

LFW

DR



automobile



CIFAR-10



airplane



bird





10)

deer



A. HEALTHY B. DISEASED Hemorrhages

Attacker Model: DCGAN

Target Model: DCGAN, DCGAN+VAE, BEGAN



White-Box Results



Black-Box Results



DR Dataset



DR Dataset



Defense? Differentially Private GAN*



*Triastcyn et al. "Generating differentially private datasets using GANs." arXiv 1803.03148

Thank you!







Thank you!



