

Privacy Evaluation and Accuracy for Different ML Secure Models

Privacy Project

Anastasiya Merkusheva
Alessandro Stanghellini
Yannick Martin

University Basel

April 2024

Introduction

- ML in medicine
- Sensitive data
- Privacy Preserving Training techniques
- Membership Inference Attacks



- Texas100
- Technical Dataset
- 67330 records
- 6169 binary features (information about the patient, the causes of injury, the diagnosis)
- 101 classes which represent the most frequent medical procedures

PyVacy

- Privacy Algorithms for PyTorch (DPSGD)
- Not well maintained

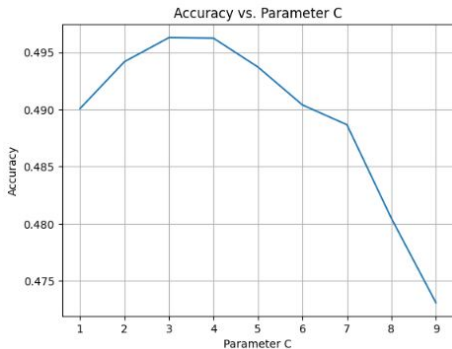
Opacus

- Privacy Algorithms for PyTorch (DPSGD)
- User-Friendly

Algorithm 3 Setting up Opacus

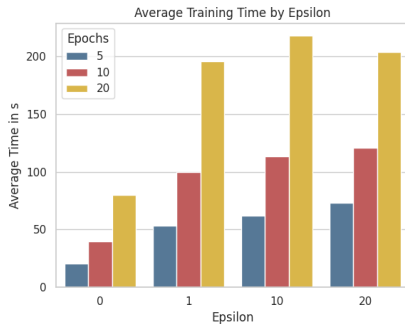
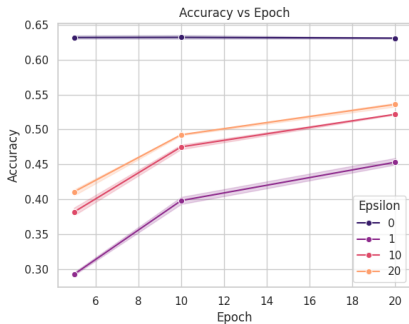
```
1: define your components(network, optimizer) as usual
2: Initialize Privacy Engine:
3: privacy_engine = PrivacyEngine()
4:
5: Make Network Private with Epsilon:
6: network, optimizer, trainloader =
7:     privacy_engine.make_private_with_epsilon(
8:         module=network,
9:         optimizer=optimizer,
10:        data_loader=trainloader,
11:        max_grad_norm=C,
12:        target_epsilon=epsilon,
13:        target_delta=Delta,
14:        epochs=epochs
15:    )
16: Now Start the training as usual
```

Hyper parameter finding & Results

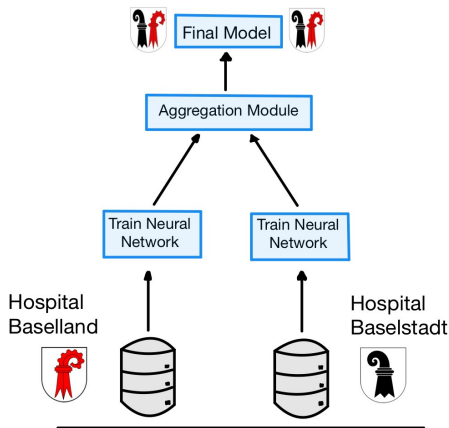


Model	Accuracy	Epsilon	Epochs
NN	0.63	-	10
Opacus	0.55	20	20
Opacus	0.53	10	20
Opacus	0.45	1	20
PyVacy	0.37	28	20
PyVacy	0.38	12	20
PyVacy	0.26	1.3	20

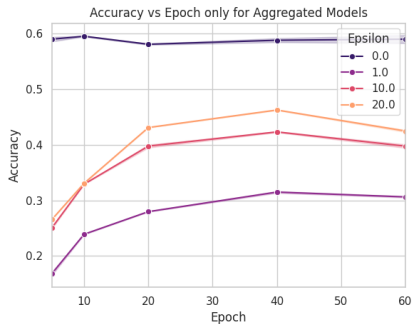
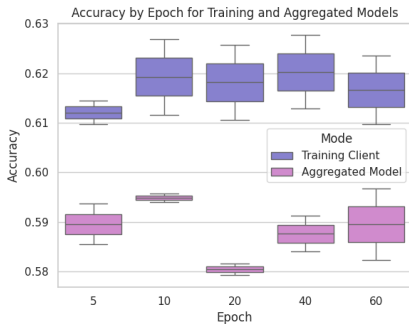
Results

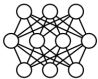



Federated Learning



Results from Federated Learning



Was  trained
on the example  ?

Why?


Curiosity!

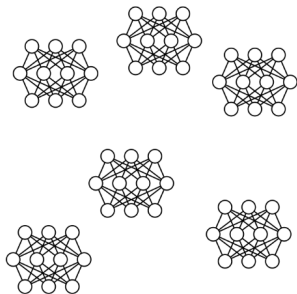
Reconnaissance!

Data Extraction!

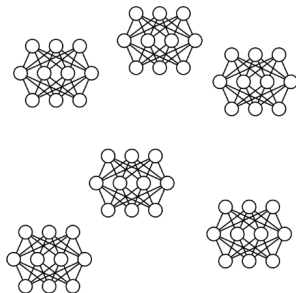
Auditing!

Attack via Population Overview

Models trained on :




Models not trained on :

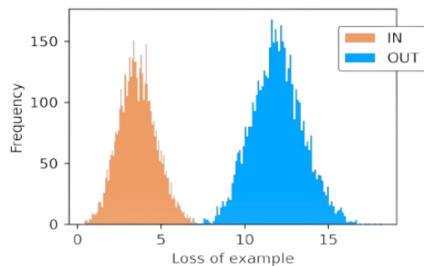


Membership Inference

$$A = \Pr(\text{Loss}_{\text{model}} (\text{image}) \mid \underbrace{L(\text{image})}_{\text{losses}}))$$

The distribution over **losses**
of models trained on 

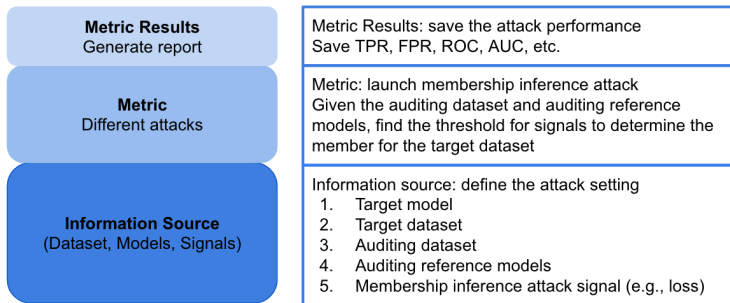
Membership Inference



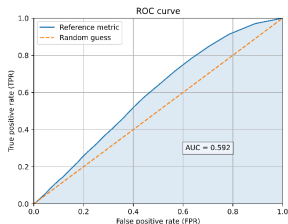
Privacy Meter: Attack via Population Data

- Privacy Meter [1] library for auditing data privacy in ML algorithms.
- Population attack uses direct statistical analysis of the target model, avoiding shadow models.
- Attack thresholds are tailored to each target model, ensuring robustness across varied datasets.
- Empirically, attackers approximate distribution by sampling records from the population data pool.

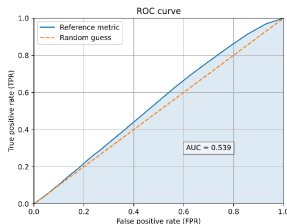
Privacy Meter: Attack via Population Data



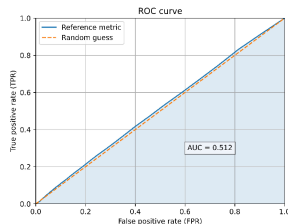
Privacy Meter: Attack via Population Data



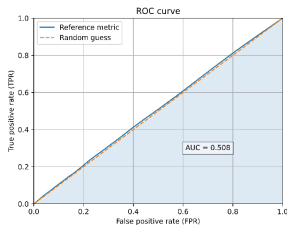
(a)



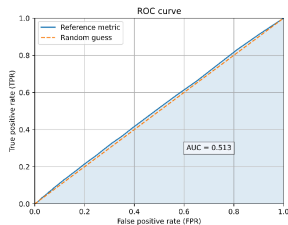
(b)



(c)



(d)



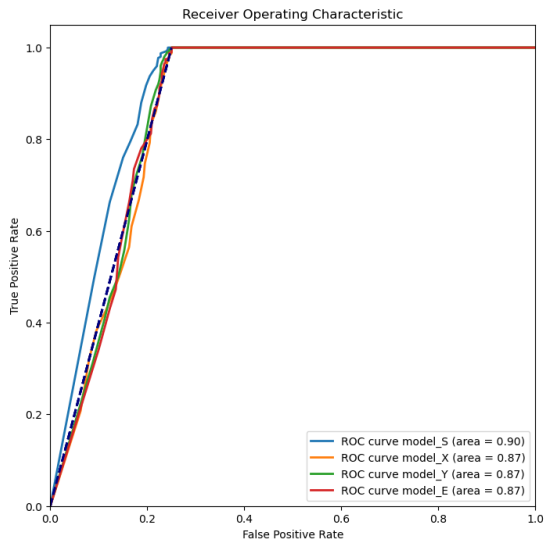
(e)

Membership Inference Attack

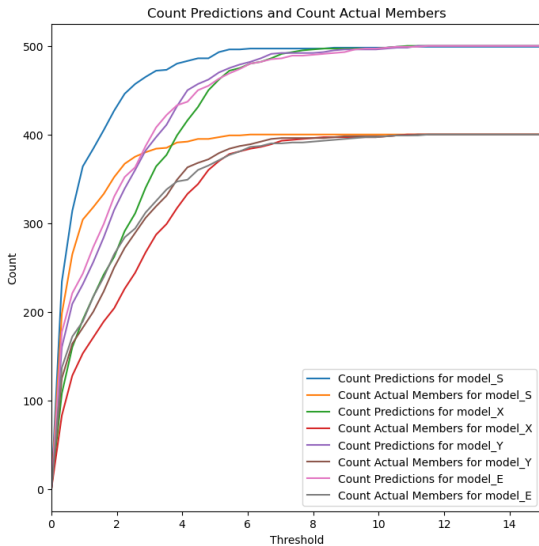
- Why can they be effective?
- Black-box and white-box attack

- Black-box attack using loss
- White-box attack using norm of the gradient of the loss function with respect to input point

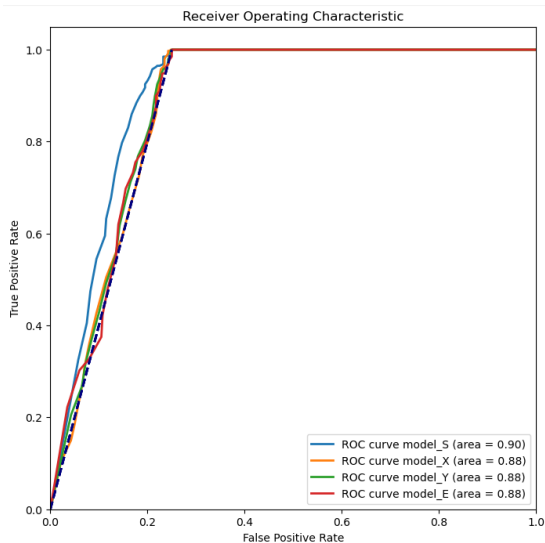
Black-box attack



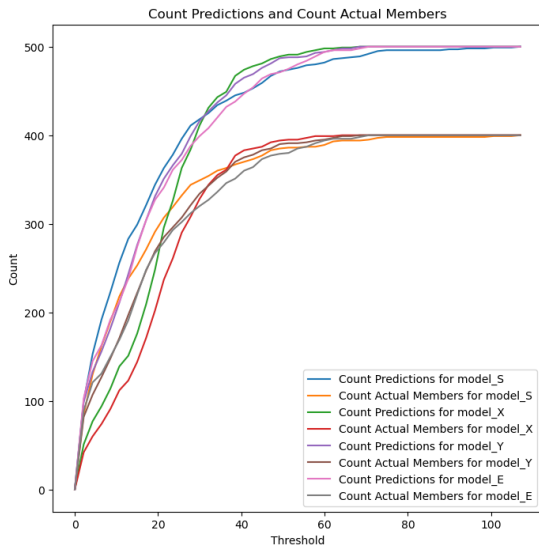
Black-box attack cont'd



White-box attack



White-box attack cont'd



Take-home messages

- We trained a model in private and non-private way
- Privacy-Meter investigation

Thank you for listening



Jiayuan Ye, Aadyaa Maddi, Sasi Kumar Murakonda, Vincent Bindschaedler, and Reza Shokri.

Enhanced membership inference attacks against machine learning models, 2022.

Privacy Meter: Attack via Population Dataset

- The hypothesis test with model-dependent attack threshold:

$$\text{If } \ell(\theta, x_z, y_z) \leq c_\alpha(\theta), \text{ reject } H_0$$

- Out world:

$$P_{out}(D, \theta, z) : D = D_0, \theta = \theta_0, z \sim \pi$$

- Empirical distribution approximation:

$$p_{\theta_0} = \{(\theta_0, z_i)\}_{i=1,2,\dots}, \text{ and } z_1, z_2, \dots \sim \pi$$

- Attack threshold for low false positive rate:

$$\frac{((\theta, z) \in p^{\theta_0}) : \ell(\theta, x_z, y_z) \leq c_\alpha(\theta_0))}{|p^{\theta_0}|} = \alpha$$