# Modeling Privacy Erosion: Differential Privacy Dynamics in Machine Learning

#### Reza Shokri

Data Privacy and Trustworthy ML Research Lab National University of Singapore



reza@comp.nus.edu.sg 💟 @rzshokri





Based on joint work with: Rishav Chourasia\*, and Jiayuan Ye\*

https://arxiv.org/abs/2102.05855

# Privacy Risks in Machine Learning



<u>How to prevent this leakage?</u> Secure multi-party computation, homomorphic encryption, trusted hardware, ...

## Privacy Risks in Machine Learning

<u>What is leakage?</u> Inferring information about members of X, beyond what can be learned about its underlying distribution



[Shokri, Stronati, Song, Shmatikov] Membership Inference Attacks against Machine Learning Models, SP'17 [Nasr, Shokri, Houmansadr] Comprehensive Privacy Analysis of Deep Learning: Passive and Active Whitebox Inference Attacks against Centralized and Federated Learning, SP'19

# How to Quantify the Leakage?

- Indistinguishability game: Can an adversary distinguish between two models that are trained on two neighboring datasets (only one includes data point x)?
  - <u>Membership inference</u>: Given a model, can an adversary infer whether data point x is part of its training set?



[Shokri, Stronati, Song, Shmatikov] Membership Inference Attacks against Machine Learning Models, SP'17

### Tool: ML Privacy Meter



ML Privacy Meter is a Python library (ml\_privacy\_meter) that enables quantifying the privacy risks of machine learning models. <u>https://github.com/privacytrustlab/ml\_privacy\_meter</u>



# Privacy Risks in Machine Learning

<u>What is leakage?</u> Inferring information about members of X, beyond what can be learned about its underlying distribution



[Shokri, Stronati, Song, Shmatikov] Membership Inference Attacks against Machine Learning Models, SP'17 [Nasr, Shokri, Houmansadr] Comprehensive Privacy Analysis of Deep Learning: Passive and Active Whitebox Inference Attacks against Centralized and Federated Learning, SP'19

# **Differential Privacy**

- A randomized algorithm  $\mathscr{A}$  satisfies  $(\epsilon, \delta)$ -DP, if for any two neighboring datasets D, D', and all sets S

 $\Pr[\mathscr{A}(D) \in S] \leq e^{\epsilon} \Pr[\mathscr{A}(D') \in S] + \delta$ 

# Renyi Differential Privacy

• A randomized algorithm  $\mathscr{A}$  satisfies  $(\alpha, \epsilon)$ -Renyi DP, for  $\alpha > 1$ , if for any two neighboring datasets D, D',

 $R_{\alpha}(\mathcal{A}(D) \, \big| \, \mathcal{A}(D')) \leq \epsilon$ 

•  $R_{\alpha}(P \parallel Q)$  is the  $\alpha$ -Renyi divergence of P with respect to Q

$$R_{\alpha}(P \| Q) = \frac{1}{\alpha - 1} \log \mathbb{E}_{z \sim Q} \left[ \left( \frac{P(z)}{Q(z)} \right)^{\alpha} \right]$$

• 
$$(\alpha, \epsilon)$$
-RDP is  $(\epsilon + \frac{\log 1/\delta}{\alpha - 1}, \delta)$ -DP for any  $\delta$ 

• Composition of  $(\alpha, \epsilon_1)$ -RDP and  $(\alpha, \epsilon_2)$ -RDP is  $(\alpha, \epsilon_1 + \epsilon_2)$ -RDP

[Mironov] Rényi differential privacy. CSF 2017

# Learning with Differential Privacy





# Learning with Differential Privacy





# Learning with Differential Privacy





#### Observation 1:

This analysis accounts for privacy loss of all iterations,

even if the only observables are the model parameters at the final iteration K

#### Observation 2:

#### Privacy loss increases with a linear rate



## Differential Privacy Dynamics

- Assume that adversary observes the model parameters at iteration K, and the state of the algorithm is **private** throughout the training
- How does privacy loss change over time?
  - What is the difference between  $R_{\alpha}(\mathscr{A}_{K-1}(D) | \mathscr{A}_{K-1}(D'))$ and  $R_{\alpha}(\mathscr{A}_{K}(D) | \mathscr{A}_{K}(D'))$  for various K?

### Noisy Gradient Descent

#### Input: Dataset $\mathcal{D} = (\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_n)$ , loss function $\ell$ , learning rate $\eta$ , noise variance $\sigma^2$ , initial parameter vector $\theta_0$ . 1: for $k = 0, 1, \cdots, K - 1$ do 2: $g(\theta_k; \mathcal{D}) = \sum_{i=1}^n \nabla \ell(\theta_k; \mathbf{x}_i)$ 3: $\theta_{k+1} = \theta_k - \frac{\eta}{n} g(\theta_k; \mathcal{D}) + \sqrt{2\eta\sigma^2} \mathcal{N}(0, \mathbb{I}_d)$ 4: Output $\theta_K$



# Dynamics of RDP in Noisy GD

- Noisy GD is a discrete-time stochastic process
- Coupled stochastic processes: Let D, D' be two neighboring datasets. Let  $\{\Theta_k\}_{k\geq 0}$  and  $\{\Theta'_k\}_{k\geq 0}$  be the sequence of probability distributions over training iterations on D, D', respectively. We assume  $\Theta_0 = \Theta'_0$  are initial parameter distributions.
- Renyi divergence  $R_{\alpha}(\Theta_K | \Theta'_K)$  reflects the privacy loss at iteration K

# Dynamics of RDP in Noisy GD

• We trace the changes in privacy loss of this discrete-time stochastic process, with a continuous-time stochastic process, which matches the probability distributions at each iteration



# RDP for Noisy GD

**Theorem** (RDP for Noisy GD) The noisy GD algorithm with loss function  $\ell(\theta; \mathbf{x})$ , learning rate  $\eta$ , and noise variance  $\sigma^2$ , satisfies  $(\alpha, \varepsilon)$  Rényi differential privacy with

$$\varepsilon \ge \frac{\alpha S_g^2}{\lambda \sigma^2 n^2 (1 - \eta \beta)^2} (1 - e^{-\lambda \eta K/2}),$$

If

1) 
$$\theta_0 \sim \mathcal{N}\left(0, \frac{2\sigma^2}{\lambda} \mathbb{I}_d\right)$$
,

 loss function ℓ(θ; x) is λ-strongly convex and β-smooth, and total gradient g(θ; D) = ∑<sub>xi∈D</sub> ∇ℓ(θ; x<sub>i</sub>) has a finite sensitivity S<sub>g</sub>,
update step-size η < 1/β.</li>



### Lower Bound on RDP

**Theorem 5** (Lower Bound on Rényi DP for Noisy GD). There exist two neighboring datasets  $\mathcal{D}, \mathcal{D}' \in \mathcal{X}^n$ , a start distribution  $p_0$ , and a  $\beta$ -smooth loss function  $\ell(\theta; \mathbf{x})$  which has a total gradient  $g(\theta; \mathcal{D})$  with finite sensitivity  $S_g$ , such that for some constants  $a_1, a_2 > 0$ , and for for any  $K \in \mathbb{N}$ , the Rényi privacy loss of  $\mathcal{A}_{Noisy-GD}$  on  $\mathcal{D}, \mathcal{D}'$  with step-size  $\eta$  and noise variance  $\sigma^2$  is lower-bounded by

$$R_{\alpha}\left(\Theta_{\eta K} \| \Theta_{\eta K}'\right) \geq \frac{a_{1}\alpha}{\sigma^{2}n^{2}} \left(1 - e^{-a_{2}\eta K}\right).$$

### **Tightness Analysis**

