CS 499/599: Machine Learning Security 03.07: (Differential) Privacy

Mon/Wed 12:00 – 1:50 pm

Sanghyun Hong

sanghyun.hong@oregonstate.edu





Notice

- Due dates (in Mar.)
 - 9th: Final project presentation
 - 11 min presentation + 3-5 min Q&A (strict)
 - Presentation *MUST* cover:
 - 1-2 slides on your research *motivation* and *goals*
 - 1-2 slides on your *ideas* (how do you plan to achieve your goals)
 - 1 slides on your *experimental design*
 - 2-3 slides on your *most interesting results*
 - 1 slides on your *conclusions* and *next steps*
 - 14th: Final exam (online, 24 hrs., unlimited trials)
 - 14th: Final project report (Template is on the website)
 - 16th: HW4 deadline (HW 1-3 late submissions are available; HW4 won't have late submissions)
- Sign-up (on Canvas)
 - Scribe lecture note [only 1 slots remain; today!]



In-class Presentation (Akshith Gunasekaran) – Red Teaming Language Models (LMs) with LMs

Topics for Today

- Privacy Attacks and Defenses
 - Non-ML: Data anonymization
 - Membership inference
 - Attacks: Yeom et al. and Shokri et al.
 - Defensive techniques
 - Model inversion
 - Attacks: Fredrikson et al. and Carlini et al.
 - Defensive techniques
 - Model extraction
 - Attacks: Tramer et al. and Jagielski et al.
 - Defensive techniques
 - Differential Privacy
 - DP-SGD
 - DP-SGD in Practice



What Can We Do To Reduce Privacy Risks of ML?

Abadi et al., Deep Learning with Differential Privacy

- ϵ -Differential Privacy
 - A randomized algorithm $M: D \to R$ with domain D and a range R satisfies ϵ -differential privacy if for any two adjacent inputs $d, d' \in D$ and any subset of outputs $S \subset R$ it holds

$$\Pr[\mathcal{M}(d) \in S] \le e^{\varepsilon} \Pr[\mathcal{M}(d') \in S]$$

• (ϵ, δ) -Differential Privacy

 $\Pr[\mathcal{M}(d) \in S] \le e^{\varepsilon} \Pr[\mathcal{M}(d') \in S] + \delta$

- δ : Represent some catastrophic failure cases [Link, Link]
- $\delta < 1/|d|$, where |d| is the number of samples in a database



Revisit'ed – Differential Privacy

• (ϵ, δ) -Differential Privacy [Conceptually]

 $\Pr[\mathcal{M}(d) \in S] \le e^{\varepsilon} \Pr[\mathcal{M}(d') \in S] + \delta$

- You have two databases d, d' differ by one item
- You make the same query M to each and have results M(d) and M(d')
- You ensure the distinguishability between the two under a measure ϵ
 - ϵ is large: those two are distinguishable, less private
 - ϵ is small: the two outputs are similar, more private
- You also ensure the catastrophic failure probability δ



Revisit'ed - Differential Privacy

• (ϵ, δ) -Differential Privacy

 $\Pr[\mathcal{M}(d) \in S] \le e^{\varepsilon} \Pr[\mathcal{M}(d') \in S] + \delta$

• Mechanism for (ϵ, δ) -DP: Gaussian noise

 $\mathcal{M}(d) \stackrel{\Delta}{=} f(d) + \mathcal{N}(0, S_f^2 \cdot \sigma^2)$

- M(d): (ϵ, δ) -DP query output on d
- f(d): non (ϵ, δ) -DP (original) query output on d
- $N(0, S_f^2 \cdot \sigma^2)$: Gaussian normal distribution with mean 0 and the std. of $S_f^2 \cdot \sigma^2$

Post-hoc: Set the Goal ϵ and Calibrate the noise $S_f^2 \cdot \sigma^2$!



How Do We Use DP for ML?

- Revisit'ed Stochastic Gradient Descent (SGD)
 - 1. At each step t, it takes a mini-batch L_t
 - 2. Computes the loss $\mathcal{L}(\theta)$ over the samples in L_t , w.r.t. the label y
 - 3. Computes the gradients g_t of $\mathcal{L}(\theta)$
 - 4. Update the model parameters θ towards the direction of reducing the loss



Make an SGD Step (ϵ, δ)-DP

- Stochastic Gradient Descent (SGD)
 - 1. At each step t, it takes a mini-batch L_t
 - 2. Computes the loss $\mathcal{L}(\theta)$ over the samples in L_t , w.r.t. the label y
 - 3. Computes the gradients g_t of $\mathcal{L}(\theta)$
 - 4. Clip (scale) the gradients to 1/C, where C > 1
 - 5. Add Gaussian random noise $N(0, \sigma^2 C^2 \mathbf{I})$ to g_t
 - 6. Update the model parameters θ towards the direction of reducing the loss



Make the Whole SGD Process (ϵ, δ)-DP

- Stochastic Gradient Descent (SGD)
 - SGD iteratively computes the (ϵ , δ)-DP step T times
 - Problem: how do we compute the total privacy leakage ϵ_{tot} over T iterations?
- Privacy accounting with moment accountant
 - Key intuition: DP has the composition property
 - Suppose the two mechanism M_1 and M_2 satisfies $(\varepsilon_1, \delta_1)$ and $(\varepsilon_2, \delta_2)$ -DP the composition of those mechanisms $M_3 = M_2(M_1)$ satisfies $(\varepsilon_1 + \varepsilon_2, \delta_1 + \delta_2)$ -DP
 - If each step t satisfies (ε , δ)-DP, the total SGD process satisfies (ε T, δ T)-DP
 - Moment accountant: tracking the total privacy leakage εT over T iterations



Putting All Together

DP-Stochastic Gradient Descent (DP-SGD)

Algorithm 1 Differentially private SGD (Outline) // we train a model θ with the privacy budget ε_{budget} **Input:** Examples $\{x_1, \ldots, x_N\}$, loss function $\mathcal{L}(\theta)$ = $\frac{1}{N}\sum_{i}\mathcal{L}(\theta, x_{i})$. Parameters: learning rate η_{t} , noise scale σ , group size L, gradient norm bound C. **Initialize** θ_0 randomly // iterate over T mini-batches for $t \in [T]$ do Take a random sample L_t with sampling probability L/N// compute the gradient Compute gradient For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$ **Clip** gradient // clip the magnitude of the gradients $\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max\left(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C}\right)$ Add noise // add Gaussian random noise to the gradients $\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} \left(\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}) \right)$ Descent $\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$ // compute the privacy cost (leakage) up to t iterations $\varepsilon, \delta \leftarrow$ compute the privacy cost (leakage) so far If $\varepsilon > \varepsilon_{buget}$: then break; // if the cost is over the budget, then stop training **Output** θ_T and compute the overall privacy cost (ε, δ) using a privacy accounting method.



- Setup
 - Datasets: MNIST | CIFAR-10/100
 - Models:
 - MNIST: 2-layer feedforward NN on 60-dim. PCA projected inputs
 - CIFAR-10/100: A CNN with 2 conv. layers and 2 fully-connected layers
 - Metrics:
 - Classification accuracy
 - Privacy cost (ε_{budget})



- Impact of Noise
 - Dataset, Models: MNIST, 2-layer feedforward NN
 - Setup: 60-dim PCA projected inputs | Clipping threshold (C): 4 | Noise (σ): 8, 4, 2 (from the left)
 - Summary:
 - On MNIST, DP-SGD offers reasonable acc. under various privacy costs (clean: 98.3%)
 - The accuracy of private models decreases as we decrease the privacy cost



- Impact of Noise
 - Dataset, Models: MNIST, 2-layer feedforward NN
 - Setup: 60-dim PCA projected inputs | Clipping threshold (C): 4 | Noise (σ): 8, 4, 2 (from the left)
 - ummary:
 - On MNIST, DP-SGD offers reasonable acc. under various privacy costs (clean: 98.3%)
 - The accuracy of private models decreases as we decrease the privacy cost



University

- Impact of Hyper-parameter Choices
 - Dataset, Models: MNIST, 2-layer feedforward NN
 - Setup: 60-dim PCA projected inputs



- Impact of Noise
 - Dataset, Models: CIFAR-10, CNN
 - Setup: Clipping threshold (C): 3 | Noise (σ): 6
 - Summary:

Oregon State

- On CIFAR-10, DP-SGD offers reasonable acc. under various privacy costs (clean: 80%)
- The accuracy of private models decreases as we decrease the privacy cost



Topics for Today

• Privacy Attacks and Defenses

- Non-ML: Data anonymization
- Membership inference
 - Attacks: Yeom *et al.* and Shokri *et al.*
 - Defensive techniques
- Model inversion
 - Attacks: Fredrikson et al. and Carlini et al.
 - Defensive techniques
- Model extraction
 - Attacks: Tramer *et al*. and Jagielski *et al*.
 - Defensive techniques
- Differential Privacy
 - DP-SGD
 - DP-SGD in Practice



What Does It Mean by Epsilon = 2/4/6 in CIFAR-10?

Jayaraman et al., Evaluating Differentially Private Machine Learning in Practice

Empirical Evaluations of Privacy Risks in DP-Models

- Setup
 - Datasets: Purchase-100 | CIFAR-100 (on 50-dim PCA projected inputs)
 - Models: Logistic regressions | 2-layer feedforward NNs
 - Privacy Attacks:
 - Membership inference: Yeom *et al*. and Shokri *et al*.
 - DP-SGD:
 - Set the clipping norm (C) to 1
 - Set the prob. of catastrophic failures (δ) to $10^{-5} < 1/|N|$ (N~60k in MNIST and 50k in CIFAR)
 - Set the batch size to 200
 - Set the learning rate to 0.01 for Adam optimizer
 - Vary ε from 0.01 to 1000
 - Compare (ϵ, δ) -DP with other DP-mechanisms: AC, CDP, zCDP, and RDP
 - Run 5-times and measure the (TPR FPR) and accuracy loss on average



Evaluation on CIFAR-100, LRs

- Summary
 - Yeom et al. and Shokri et al. are weak privacy attacks
 - In other words, (ϵ, δ) -DP theoretically offers very strong privacy bounds
 - If a DP-mechanism offers stronger bound, the acc. of models decrease accordingly



Evaluation on CIFAR-100, NNs

- Summary
 - Yeom et al. and Shokri et al. are weak privacy attacks
 - In other words, (ϵ, δ) -DP theoretically offers very strong privacy bounds
 - If a DP-mechanism offers stronger bound, the acc. of models decrease accordingly
 - Compared to LRs, NNs leak more in higher privacy budgets



Evaluation on MI Predictions: LRs vs. NNs

- Summary
 - Yeom et al. and Shokri et al. are weak privacy attacks
 - In other words, (ϵ, δ) -DP theoretically offers very strong privacy bounds
 - If a DP-mechanism offers stronger bound, the acc. of models decrease accordingly
 - Compared to LRs, NNs leak more in higher privacy budgets
 - Predictions (TPRs and FPRs) are more consistent in LRs than NNs in CIFAR-100



Recap: Privacy!

- Privacy Attacks and Defenses
 - Non-ML: Data anonymization
 - Membership inference
 - Attacks: Yeom et al. and Shokri et al.
 - Defensive techniques
 - Model inversion
 - Attacks: Fredrikson et al. and Carlini et al.
 - Defensive techniques
 - Model extraction
 - Attacks: Tramer et al. and Jagielski et al.
 - Defensive techniques
 - Differential Privacy
 - DP-SGD
 - DP-SGD in Practice



Thank You!

Mon/Wed 12:00 – 1:50 pm

Sanghyun Hong

https://secure-ai.systems/courses/MLSec/W22



