

Lecture 17

- Private Machine Learning

- DP SGD

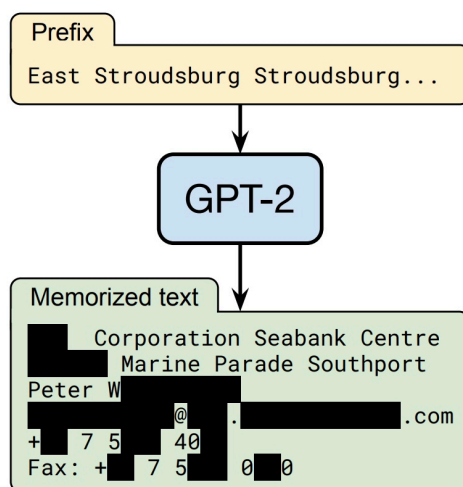
- Privacy Analysis

- DP SGD in code.

Announcement: Release HW3 this week
includes = ① Written Component
② Programming component.

Recitation (in-person)

"Memorization" Attack.



Extracting Training Data from Large Language Models

Nicholas Carlini¹ Florian Tramèr² Eric Wallace³ Matthew Jagielski⁴
Ariel Herbert-Voss^{5,6} Katherine Lee¹ Adam Roberts¹ Tom Brown⁵
Dawn Song³ Úlfar Erlingsson⁷ Alina Oprea⁴ Colin Raffel¹
¹Google ²Stanford ³UC Berkeley ⁴Northeastern University ⁵OpenAI ⁶Harvard ⁷Apple

Private SGD. (DP-SGD)

$$\text{Private SGD} \left(\overset{\text{Loss}}{\downarrow} L(\cdot) = \frac{1}{n} \sum_{i=1}^n \ell(\cdot; x_i), \overset{\text{feasible set}}{\downarrow} C, \overset{\text{learning rate}}{\downarrow} \eta, \overset{\text{Noise rate}}{\swarrow} \beta \right) =$$

$$\text{Init: } w_0 \in C$$

For $t=1, \dots, T$:

Random subsample $B_t \subseteq \{1, \dots, n\}$
"mini-batch"

$$g_t = \frac{1}{|B_t|} \sum_{i \in B_t} \nabla_w \ell(w_{t-1}; x_i)$$

$$\tilde{g}_t = g_t + N(0, \beta^2 I_d) \quad \text{: Gaussian Mechanism.}$$

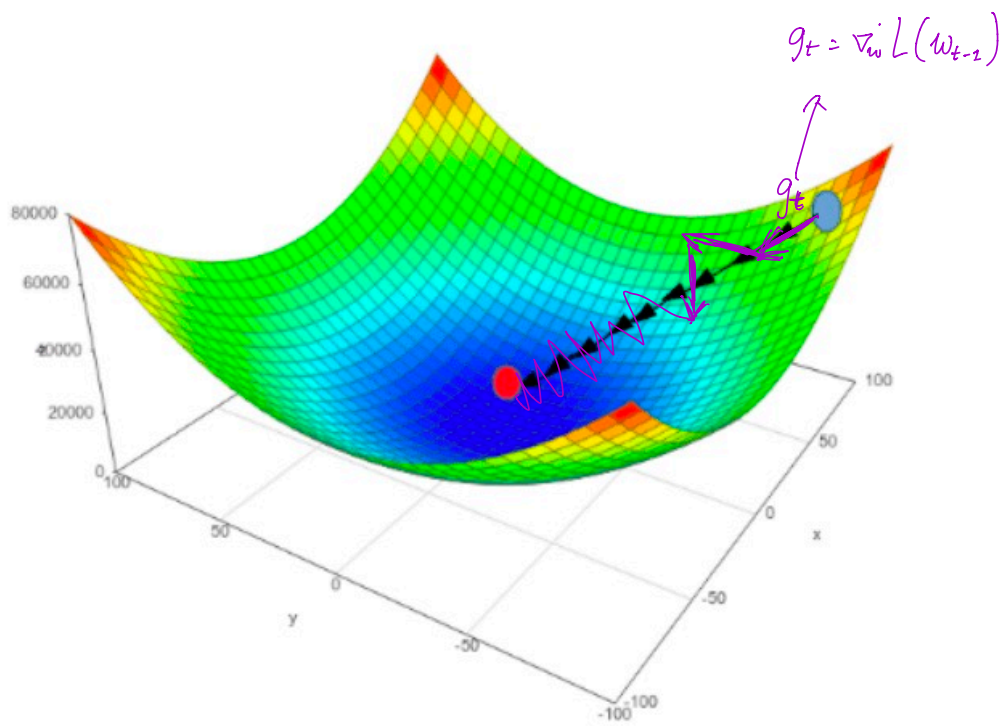
$$u_t = w_{t-1} - \eta \cdot \tilde{g}_t$$

$$w_t = \underset{m \in C}{\text{argmin}} \|w - u_t\|_2$$

$$\text{Output: } \frac{1}{T} \sum_{t=1}^T w_t$$

or

$$w_T$$



Privacy Proof

Proof idea: • Think of releasing w_1, w_2, \dots, w_T .

- Suffices to release the update between iterates

$$w_0 \left[\xleftarrow{\text{update}} \rightarrow \right] w_1 \left[\xleftarrow{\text{update}} \rightarrow \right] w_2 \dots \dots w_T$$

- Suffices to release the sequence of gradient estimates

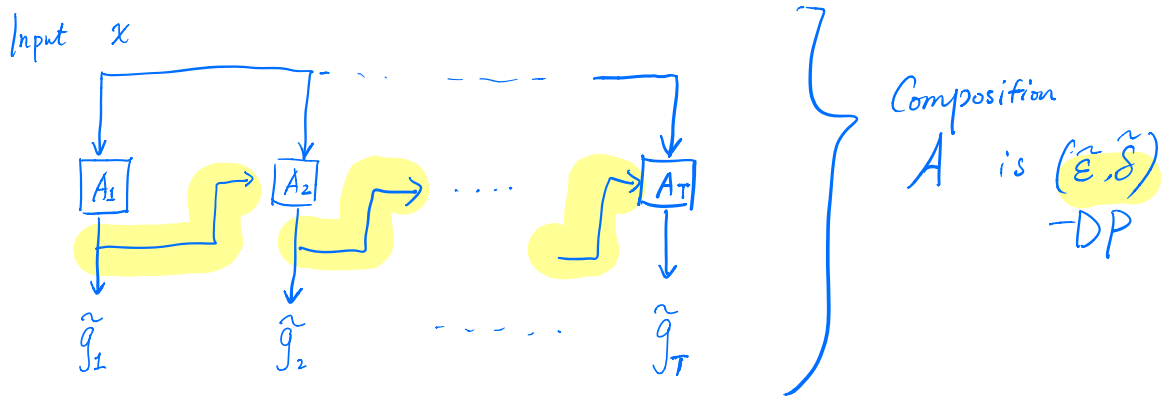
$$\tilde{g}_1, \tilde{g}_2, \dots, \tilde{g}_T$$

The output is a post-processing

Show releasing $(\tilde{g}_1, \tilde{g}_2, \dots, \tilde{g}_T)$ satisfies DP.

- ① Each step (releasing \tilde{g}_t) satisfies (ϵ, δ) -DP
- ② Adaptive composition across T steps

Adaptive Composition.



Suppose each step A_1, \dots, A_T is (ϵ, δ) -DP.

What are the values of $\tilde{\epsilon}$ and $\tilde{\delta}$?

- (Basic) Composition: $\tilde{\epsilon} = T\epsilon$, $\tilde{\delta} = T\delta$.
- Advanced Composition: $\tilde{\epsilon} = \epsilon \cdot \sqrt{2T \ln(\frac{1}{\delta'})} + T \cdot \epsilon \frac{e^\epsilon - 1}{e^\epsilon + 1}$
 $\tilde{\delta} = T\delta + \delta'$

If $\epsilon < \frac{1}{\sqrt{T}}$

$(\epsilon \sqrt{T})^2$ is "smaller" than $\epsilon \sqrt{T}$

Then $\tilde{\epsilon}$ is in the order of $\epsilon \cdot \sqrt{T \ln(\frac{1}{\delta'})}$

$\ll \epsilon \cdot T$
for large T .

Numeric Example.

$$\varepsilon = \frac{1}{1000}, \quad \delta$$

$$T = 500,$$

$$\text{Basic Composition: } \tilde{\varepsilon} = 0.5, \quad \tilde{\delta} = T\delta$$

$$\text{Advanced Composition: } \tilde{\varepsilon} \leq 0.1, \quad \tilde{\delta} = 10^{-6} + T\delta$$

Privacy Proof

Proof idea: • Think of releasing w_1, w_2, \dots, w_T .

- Suffices to release the update between iterates

$$w_0 \left[\xleftarrow{\text{update}} \rightarrow \right] w_1 \left[\xleftarrow{\text{update}} \rightarrow \right] w_2 \dots \dots w_T$$

- Suffices to release the sequence of gradient estimates

$$\tilde{g}_1, \tilde{g}_2, \dots, \tilde{g}_T$$

The output is a post-processing

Show releasing $(\tilde{g}_1, \tilde{g}_2, \dots, \tilde{g}_T)$ satisfies DP.

This step. \rightarrow ① Each step (releasing \tilde{g}_t) satisfies (ϵ, δ) -DP

② Adaptive composition across T steps ✓

Multivariate Gaussian Mechanism

e.g.
average
gradient

$$f: \mathcal{X}^n \mapsto \mathbb{R}^d$$

$$\Delta_2(f) = \max_{x, x' \text{ neighbors}} \|f(x) - f(x')\|_2$$

L₂ sensitivity.

$$A(x) = f(x) + N\left(0, \underbrace{\frac{2 \Delta_2^2 \log(2/\delta)}{\epsilon^2}}_{\text{per-coordinate variance. } \zeta^2} I_{\text{dim}}\right)$$

$$f(x) = \begin{pmatrix} f(x)_1 \\ f(x)_2 \\ \vdots \\ f(x)_d \end{pmatrix} + \begin{pmatrix} N(0, \zeta^2) \\ N(0, \zeta^2) \\ \vdots \\ N(0, \zeta^2) \end{pmatrix}$$

“fresh” independent noise

Theorem: $\forall \epsilon \leq 1, \delta > 0$

$A(\cdot)$ satisfies (ϵ, δ) -DP.

What is Δ_2 ?

$$f: \mathcal{X}^n \mapsto \mathbb{R}^d$$

$$\Delta_2(f) = \max_{\substack{x, x' \\ \text{neighbors}}} \|f(x) - f(x')\|_2$$

Average gradient

$$f \leftarrow \frac{1}{|B_t|} \sum_{i \in B_t} \nabla_w l(w; x_i)$$

Suppose $|B_t|=1$.

$$f \leftarrow \nabla_w l(w; x_i)$$

$$\Delta_2 = \max_{x_i, x_i'} \|\nabla_w l(w; x_i) - \nabla_w l(w; x_i')\|_2$$

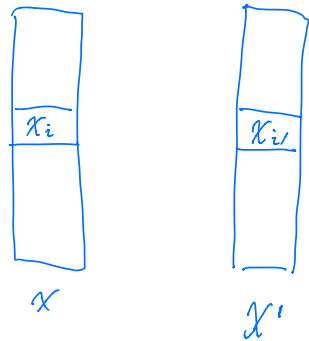
In theory, we make assumption on l

so that $\Delta_2 \leq G$

Lipschitzness

\Rightarrow Changing w a bit does not change $l(w; x_i)$ too much for $x_i \in \mathcal{X}$.

Privacy Amplification by Sub-sampling.



→ first sample minibatch B_t of size 1.

→ Compute gradient using B_t .

In general: suppose $A: \mathcal{X} \mapsto \mathcal{Y}$ is (ϵ, δ) -DP.

↑
an algo that takes input a data set of size 1.

(e.g., add gaussian noise to the gradient of one example)

Consider $A': \mathcal{X}^n \mapsto \mathcal{Y}$

Random index $\rightarrow I \leftarrow \text{unif}\{1, \dots, n\}$
Return $A(x_I)$

A' is (ϵ', δ') -DP where

$$\epsilon' = \ln\left(1 + \frac{e^\epsilon - 1}{n}\right) \approx \frac{\epsilon}{n} \quad \text{for } \epsilon \leq 1$$

$$\delta' = \frac{\delta}{n}$$

Can generalize to $|B_t| > 1$.

$$\epsilon' \approx \frac{|B_t|}{n} \epsilon, \quad \delta' \approx \frac{|B_t|}{n} \cdot \delta$$

Wrapping up the privacy proof.

- For each step: sub-sampled Gaussian mechanism
- Apply Adaptive Composition.

DP-SGD (in Theory)

Init: $w_0 \in C$

For $t=1, \dots, T$:

Random subsample $B_t \subseteq \{1, \dots, n\}$
"mini-batch"

$$g_t = \frac{1}{|B_t|} \sum_{i \in B_t} \nabla_w \ell(w_{t-1}; x_i)$$

$$\tilde{g}_t = g_t + N(0, \beta^2 I_d)$$

$$u_t = w_{t-1} - \eta \cdot \tilde{g}_t$$

$$w_t = \operatorname{argmin}_{m \in C} \|w - u_t\|_2$$

Assume
 ℓ is Lipschitz
or
gradient $\nabla_w \ell(w; x_i)$
for every $w \in C$ & $x_i \in X$

DP-SGD (in practice)

For $t=1, \dots, T$

Sample minibatch $B_t \subseteq \{1, \dots, n\}$

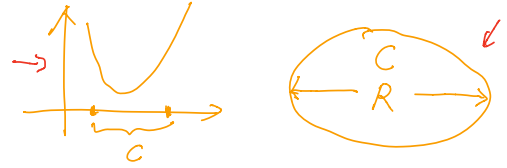
$$g_t = \frac{1}{|B_t|} \sum_{i \in B_t} \operatorname{Clip}(\nabla_w \ell(w; x_i), G)$$

$$\operatorname{Clip}(g, G) = g \min\left(1, \frac{G}{\|g\|_2}\right)$$

shrink
gradient
if too
large.

$$\tilde{g}_t = g_t + \text{Gaussian Noise.}$$

Convergence / Optimality.



Theorem. Let $L: C \rightarrow \mathbb{R}$ be convex and G -Lipschitz
 $C \subseteq \mathbb{R}^d$ be a closed and convex set
with diameter R

(Part a)

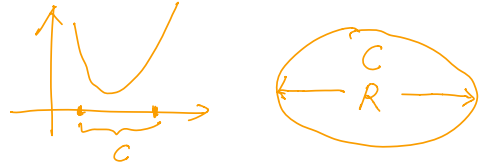
↓

$$w^* \in \operatorname{argmin}_{w \in C} L(w)$$

• For regular PGD, set $\eta = \frac{R}{G\sqrt{T}}$, then $L(\hat{w}) - L(w^*) \leq \frac{RG}{\sqrt{T}}$

• For noisy PGD, set η, T, δ^2 so that, $\mathbb{E}[L(\hat{w}) - L(w^*)] \leq \frac{RG\sqrt{d \ln(1/\delta)}}{n\epsilon}$

Convergence / Optimality.



Theorem. Let $L: C \rightarrow \mathbb{R}$ be convex and G -Lipschitz
 $C \subseteq \mathbb{R}^d$ be a closed and convex set
 with diameter R

(Part a)

↓

$$w^* \in \underset{w \in C}{\operatorname{argmin}} L(w)$$

• For regular PGD, set $\eta = \frac{R}{G\sqrt{T}}$, then $L(\hat{w}) - L(w^*) \leq \frac{RG}{\sqrt{T}}$ ↓ 0

• For noisy PGD, set η, T, δ^2 so that, $\mathbb{E}[L(\hat{w}) - L(w^*)] \leq O\left(\frac{RG\sqrt{d \ln(1/\delta)}}{n\epsilon}\right)$

"cost of privacy" Gap: $\frac{\sqrt{d}}{n\epsilon}$ ← "tight" in the worst-case

Gap for EM: $\frac{d}{n\epsilon}$

Proof (for regular PGD).

$$w^* = \operatorname{argmin}_{w \in C} L(w)$$

Claim. (Measure of Progress)

$$\underbrace{L(w_t) - L(w^*)}_{\text{Excess Risk}} \leq \frac{\eta \cdot \|g_t\|^2}{2} + \frac{1}{2\eta} \left(\underbrace{\|w_t - w^*\|^2}_{\text{Reduction on Squared distances}} - \underbrace{\|w_{t+1} - w^*\|^2}_{\text{Reduction on Squared distances}} \right)$$

2 Key Quantities

Excess Risk

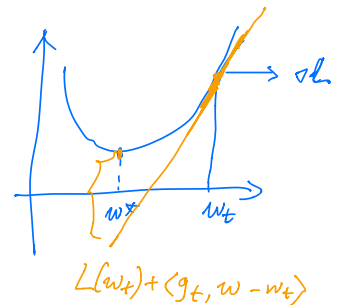
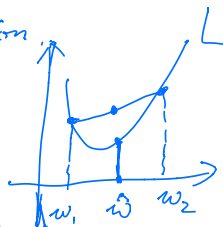
Distance to w^*

Proof for $\hat{w} = \frac{1}{T} \sum_{t=1}^T w_t$

By Jensen Inequality for convex function,

$$L(\hat{w}) \leq \frac{1}{T} \sum_t L(w_t)$$

Compare w/ $\frac{1}{T} (T \cdot L(w^*))$



$$L(\hat{w}) - L(w^*) \leq \frac{1}{T} \left(\sum_t (L(w_t) - L(w^*)) \right) \leftarrow \text{use "Progress Claim"}$$

$$\leq \frac{\eta}{2} \cdot \max_t \|g_t\|^2 + \frac{1}{2\eta T} \left(\|w_1 - w^*\|^2 - \|w_{T+1} - w^*\|^2 \right)$$

$$\leq \frac{\eta}{2} \cdot G^2 + \frac{1}{2\eta T} \left(\|w_1 - w^*\|^2 \right)$$

$$\leq \frac{\eta}{2} G^2 + \frac{R^2}{2\eta T} = \frac{GR}{\sqrt{T}}$$

Equalize

$$\stackrel{\text{Set } \eta}{=} \frac{R}{G} \cdot \frac{1}{\sqrt{T}}$$

Noisy / Private PGD.

$$\tilde{g}_t = \underline{g}_t + N(0, \beta^2 I)$$

"New" Progress Claim.

$$\mathbb{E}[L(w_t) - L(w^*)] \leq \frac{\eta}{2} \mathbb{E}[\|\tilde{g}_t\|^2] + \frac{1}{2\eta} \mathbb{E}[\|w_t - w^*\|^2 - \|w_{t+1} - w^*\|^2]$$

Proof. $\mathbb{E}[L(w_t) - L(w^*)] \leq \mathbb{E}[\langle \eta g_t, w_t - w^* \rangle]$

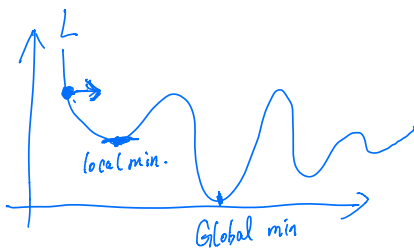
$$= \mathbb{E}[\langle \eta \mathbb{E}[\tilde{g}_t | w_t], w_t - w^* \rangle]$$

$$\Rightarrow \mathbb{E}[\langle \eta \tilde{g}_t, w_t - w^* \rangle]$$

$$\langle a, b \rangle = \|a\|_2^2 + \|b\|_2^2 - \|a - b\|_2^2$$

$$\mathbb{E}[\|\tilde{g}_t\|_2^2] \leq \|g_t\|_2^2 + \boxed{d\beta^2}$$

What about Nonconvex Case?



Smoothness.

(Lipschitz Gradient)

$$\|\nabla L(w) - \nabla L(w')\|_2 \leq \beta \|w - w'\|_2$$

$$L(w') \leq L(w) + \nabla L(w)^T (w' - w) + \frac{\beta}{2} \|w - w'\|_2^2.$$

Can Show: w_1, \dots, w_T

$$\frac{1}{T} \sum \|\nabla L(w_t)\|_2^2 \rightarrow O\left(\frac{1}{\sqrt{T}}\right). \quad (\text{non-DP})$$

$$\rightarrow \frac{\sqrt{d}}{n\epsilon} \sqrt{\ln\left(\frac{1}{\delta}\right)} \quad (\text{DP})$$

