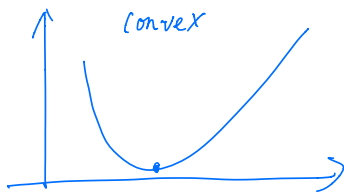
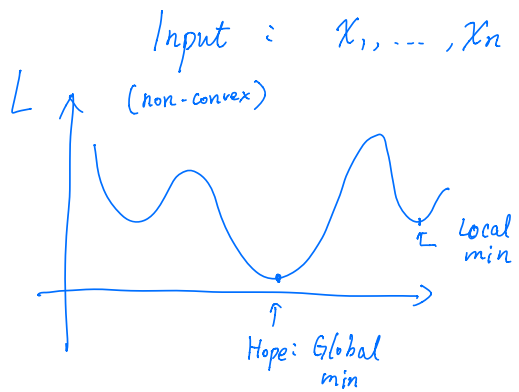


Private Gradient Descent

- Recap: Private ERM Problem
Exponential Mechanism
- (Projected) Gradient Descent
Privacy, Convergence.

HW2 Due on Weds

Formulation.



feasible set of models/parameters

$$l: \overset{\downarrow}{C} \times X \mapsto \mathbb{R}$$

$l(w, x)$ measures "loss"

$$L: C \mapsto \mathbb{R}$$

$$L(w) = \frac{1}{n} \sum_{i=1}^n l(w, x_i)$$

$$\Pi_C(w) = \arg \min_{w' \in C} \|w - w'\|_2$$

"Projection"

Goal: output \hat{w} s.t.

$$L(\hat{w}) \approx \min_{w \in C} L(w)$$

Projected Gradient Descent (PGD)

PGD (L, C, η):

→ Init: $w_0 \in C$ arbitrary

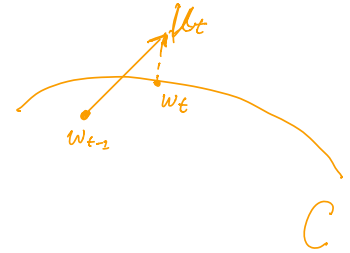
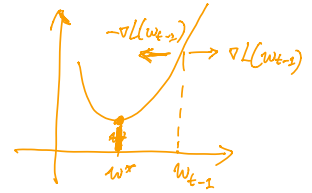
For $t = 1, \dots, T$:

$$g_t = \nabla L(w_{t-1}) \leftarrow \text{Gradient}$$

$$u_t \leftarrow w_{t-1} - \eta \cdot g_t$$

$$w_t \leftarrow \Pi_C(u_t)$$

→ Output $\hat{w} = \frac{1}{T} \sum_{t=1}^T w_t$



Robustness to noise in gradient estimation. (g_t)

→ For privacy: $\tilde{g}_t = g_t + \underbrace{N(0, \delta^2 I_d)}_{\substack{d\text{-dim} \\ \text{i.i.d. Gaussian}}}$

$$\mathbb{E}[\tilde{g}_t] = g_t = \nabla L(w_{t-1})$$

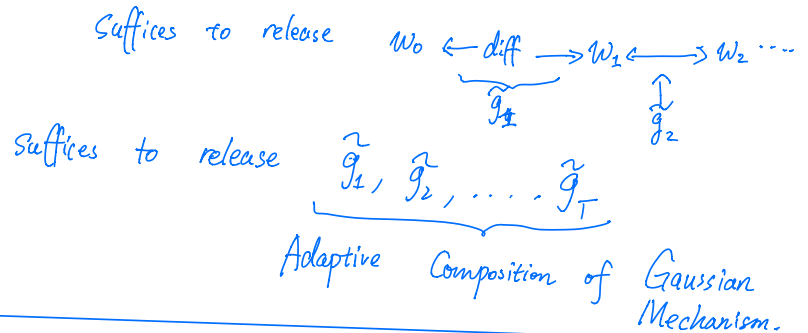
→ For efficiency: $\tilde{g}_t = \nabla L(w_{t-1}, \mathcal{K}_{I_t})$, $I_t \leftarrow \text{unif}\{1, \dots, n\}$

$$\mathbb{E}[\tilde{g}_t] = g_t = \nabla L(w_{t-1})$$

↓
"Unbiased"

Noisy/Private PGD (outputs $\hat{w} = \frac{1}{T} \sum_{t=1}^T w_t$).

- Proof idea for privacy: think of releasing w_1, \dots, w_T



Lemma. If l is G -Lipschitz for every x . ($\|\nabla l(w; x)\| \leq G$) on C ,
 $(|l(w, x) - l(w', x)| \leq G \|w - w'\|)$.

then Noisy PGD is (ϵ, δ) -DP when $\left[\epsilon \geq \frac{2G}{n} \cdot \frac{\sqrt{2T \ln(1/\delta)}}{\epsilon} \right]$

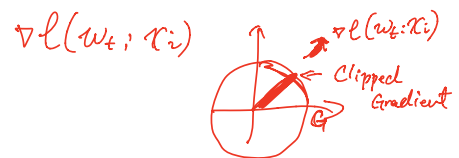
Proof Ideas:

- ① Overall l_2 -sensitivity of all T gradients: $\frac{2G}{n} \cdot \sqrt{T}$
 $\nabla L(w_t) = \frac{1}{n} \sum_{i=1}^n \nabla l(w_t; x_i)$

- ② Apply Advanced Composition to T adaptive Gaussian Mech
 loose bound

Better techniques = $\left[\begin{array}{l} \text{Renyi DP accountant (Tensorflow)} \\ \text{Gaussian DP "-----"} \\ \text{(zero) - Concentrated DP "-----"} \end{array} \right]$

NB: Enforce low sensitivity by "Gradient Clipping" on



Private SGD.

Private SGD ($L, C, \overset{\text{learning rate}}{\eta}$) :

→ Init: $w_0 \in C$ arbitrary

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

Also . Subsampling
a minibatch

→ $I_t \leftarrow \text{unif}(\{1, \dots, n\})$; $g_t = \nabla \ell(w_{t-1}; x_{I_t})$

$$\tilde{g}_t = g_t + N(0, \delta^2 \text{Id})$$

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

$$w_t \leftarrow \Pi_C(u_t).$$

→ Output $\hat{w} = \frac{1}{T} \sum_{t=1}^T w_t$

SG Langevin
Dynamics.

SGD + Gaussian
Noise.

How to analyze Privacy?

Privacy Amplification

- Keep I_t secret
- Use their randomness.

In general: $A = \mathcal{X} \mapsto \mathcal{Y}$ is (ϵ, δ) -DP.
 \uparrow take one data point

• Consider: $A' = \mathcal{X}^n \mapsto \mathcal{Y}$ $\left\{ \begin{array}{l} I \leftarrow \text{unif}(\{1, \dots, n\}) \\ \text{Return } A(\mathcal{X}_I) \end{array} \right.$

• 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 subsample of size $m \leq n$.

$$\epsilon' \approx \frac{m}{n} \epsilon$$

$$\delta' \approx \frac{m}{n} \delta.$$

- Adding Gaussian noise w/ $\beta = \frac{2G}{\epsilon} \sqrt{2 \ln(1/\delta)}$ is (ϵ, δ) -DP.

$$\left[\nabla \ell(w_{t-1}; \mathcal{X}) + N(0, \beta^2 I) \right]$$

- Subsampling + Gaussian is (ϵ', δ') -DP

$$\epsilon' = \frac{\epsilon}{n}, \quad \delta' = \frac{\delta}{n}.$$

- Advanced Composition. over T iterations.

$$\left(\frac{\epsilon \sqrt{T}}{n}, \frac{\delta \cdot T}{n} \right) \text{-DP}$$

(ignoring constants).

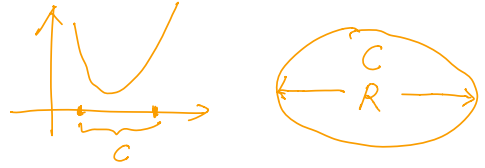
Set δ to be smaller than

$$\frac{1}{n}.$$

$$\text{Can } \delta = \frac{1}{n^5}$$

Can Run $T = n^2$ iterations!

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 $L(\hat{w}) - L(w^*) \leq \frac{RG}{\sqrt{T}}$ ↓ 0

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

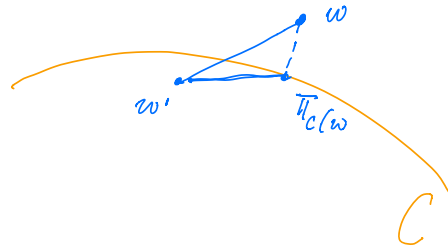
"cost of privacy" Gap: $\frac{\sqrt{d}}{n\epsilon}$ ← "tight" in the worst-case

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

Projection Lemma

Lemma If $C \subseteq \mathbb{R}^d$ is closed and convex.
Then for any $w \in \mathbb{R}^d$, $w' \in C$

$$\| \Pi_C(w) - w' \| \leq \| w - w' \|$$



Proof (for regular PGD).

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

Claim. (Measure of Progress)

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

2 Key Quantities

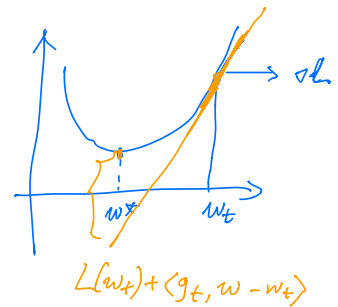
Excess Risk

Distance to w^*

Squared distances.

Proof. $L(w^*) \geq L(w_t) + \langle g_t, w^* - w_t \rangle$

$$L(w_t) - L(w^*) \leq \frac{1}{\eta} \langle \eta g_t, w_t - w^* \rangle$$



$$\forall a, b \in \mathbb{R}^d$$

$$\langle a, b \rangle = \frac{1}{2} (\|a\|^2 + \|b\|^2 - \|a-b\|^2)$$

$$\begin{aligned} L(w_t) - L(w^*) &\leq \frac{1}{2\eta} \left(\|\eta g_t\|^2 + \|w_t - w^*\|^2 - \|w_t - w^* - \eta g_t\|^2 \right) \\ &= \frac{\eta \cdot \|g_t\|^2}{2} + \frac{1}{2\eta} \left(\|w_t - w^*\|^2 - \|w_t - w^* - \eta g_t\|^2 \right) \\ &\leq \frac{\eta \cdot \|g_t\|^2}{2} + \frac{1}{2\eta} \left(\|w_t - w^*\|^2 - \|w_{t+1} - w^*\|^2 \right) \quad (\text{Projection}) \end{aligned}$$

Noisy / Private PGD.