# Private Gradient Descent

Recap:

- (Projected) Gradient Descent
  Privacy , Convergence.

- Practical Aspects of Private Deep Learning (w/ slides)

  [ HW2    Due ]

# Projected Gradient Descent (PGD)

Constraint Set

learning rate

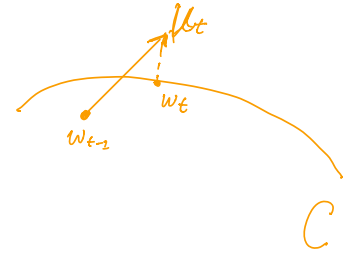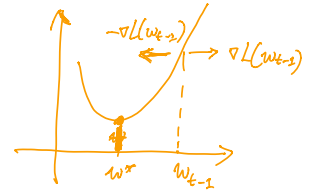$PGD(L, C, \eta):$

→ Init: $w_0 \in C$ arbitrary

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

$g_t = \nabla L(w_{t-1})$ ← Gradient

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

$w_t \leftarrow \overline{\Pi}_C(u_t).$

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

# Private SGD.

Private SGD $(L, C, \eta)$ :

    *learning rate* ↓ $\eta$

→ Init: $w_0 \in C$   arbitrary

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

*Also · Subsampling a minibatch* ⟶

$I_t \leftarrow unif\left(\{1,\ldots,n\}\right)$ ; $g_t = \nabla \ell(w_{t-1}; x_{I_t})$

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

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

$w_t \leftarrow \overline{\Pi}_C (u_t).$

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

# Privacy Amplification

- Keep $I_t$ secret
- Use their randomness.

In general: $A : \mathcal{X} \longmapsto Y$ is $(\varepsilon, \delta)$-DP.

<span style="color:red">Computation at each step $\nearrow$</span>

$\curvearrowleft$ take one data point

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

$$\begin{cases} I \leftarrow \text{unif} (\{1, \ldots, n\}) \\ \\ \text{Return } A(\mathcal{X}_I) \end{cases}$$

- $A'$ is $(\varepsilon', \delta')$-DP where

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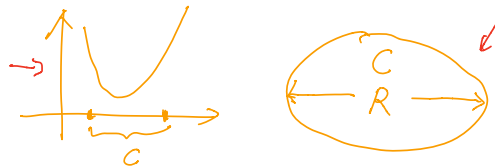
$$\delta' = \frac{\delta}{n}$$

Can generalize to subsample of size $\boxed{m \leq n}$.

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

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

# Convergence / Optimality.

**Theorem.** Let $L : C \longrightarrow R$ be convex and $\underline{G\text{-Lipschitz}}$

$C \subseteq R^d$ be a closed and convex set with $\underline{\text{diameter } R}$

$$w^* \in \underset{w \in C}{\arg\min} \; L(w)$$

(Part a)

- For regular PGD, set $\eta = \dfrac{R}{G\sqrt{T}}$, then $L(\hat{w}) - L(w^*) \leq \boxed{\dfrac{RG}{\sqrt{T}}}$   $\downarrow 0$   $T \nearrow \infty$

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

For theory: $T \approx n^2$

Practice: Trial-&-error.

(so is learning rate)

"Cost of privacy" Gap: $\dfrac{\sqrt{d}}{n\varepsilon}$   $\leftarrow$ "tight" in the worst-case

Gap for EM: $\dfrac{d}{n\varepsilon}$

Proof   (for   regular PGD).

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

Claim.  ( Measure of Progress).

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

2 Key Quantities

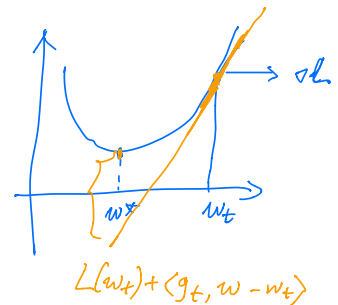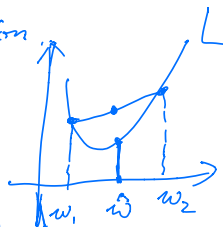Excess Risk       Distance to $w^*$

Reduction on Squared distances.

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

By Jensen Inequality. for Convex function.

$$L(\hat{w}) \leq \underbrace{\frac{1}{T}\sum_t L(w_t)}_{\text{Compare w/}} \qquad \frac{1}{T}\left(T \cdot L(w^*)\right)$$



$$L(w_t) + \langle g_t, w - w_t\rangle$$

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

$$\leq \frac{\eta}{2} \cdot \max_t \|g_t\|^2 + \frac{1}{2\eta T}\left(\|w_1 - w^*\|_2^2 - \|w_{T+1} - w^*\|^2\right)$$

$$\leq \frac{\eta}{2} \cdot G^2 + \frac{1}{2\eta T}\left(\|w_1 - w^*\|_2^2\right)$$

$$\leq \underbrace{\frac{\eta}{2} G^2 + \frac{R^2}{2\eta T}}_{\text{Equalize}} \quad \underset{\text{Set } \eta}{=} \quad \frac{GR}{\sqrt{T}}$$

$$= \frac{R}{G} \cdot \frac{1}{\sqrt{T}}$$

# Noisy / Private PGD.

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

## "New" Progress Claim.

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

Proof.
$$\mathbb{E}\left[ L(w_t) - L(w^*) \right] \leq \mathbb{E}_t\left[ \langle \eta g_t, w_t - w^* \rangle \right]$$
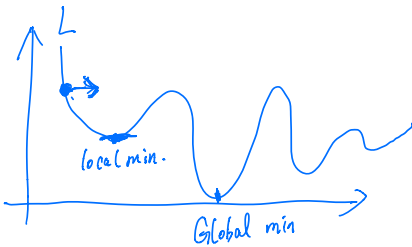
$$= \mathbb{E}_{w_t}\left[ \langle \eta \, \mathbb{E}[\tilde{g}_t | w_t], w_t - w^* \rangle \right]$$

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

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

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

# What about Noncovex Case?


local min.
Global min


$w_t$

Smoothness.
(Lipschitz Gradient)

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

$$L(w') \leq L(w) + \nabla L(w)^\top (w' - w) + \frac{\beta}{2} \|w - w'\|^2$$

Can Show: $w_1, \ldots, w_T$

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

$$\longrightarrow \frac{\sqrt{d}}{n\varepsilon} \sqrt{\ln(\tfrac{1}{\delta})} \quad (\text{DP})$$