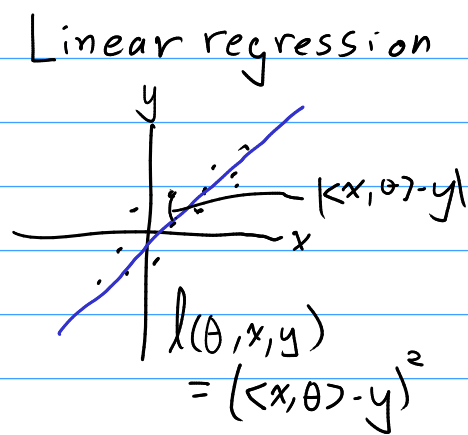


# Lecture 13

## Differentially Private ERM

### A Quick Primer on Non-Private Machine Learning

**Formulation**  
Dataset  $D$  of  $(x, y)$  pairs  
loss function:  $l(\theta, x, y)$



$$J(\theta, D) = \sum_{i=1}^n l(\theta, x_i, y_i)$$

### Empirical Risk Minimization (ERM)

Goal: minimize  $J(\theta, D)$

Given: Dataset  $D$  of  $(x, y)$   
"Parameter space"  $C$

Ideal  $\theta^* = \operatorname{argmin}_{\theta \in C} J(\theta, D)$

Algo will output  $\hat{\theta}$ .

Expected excess empirical risk

$$f_{\hat{\theta}}(x) = \langle \hat{\theta}, x \rangle$$

$$E[J(\hat{\theta}, D) - J(\theta^*, D)]$$

↑ over algo's randomness

Restrictions:

- Diameter of  $C$  is bounded
- $l(\cdot, x_i, y_i)$  convex, L-Lipschitz  $\forall (x_i, y_i) \in X$

## Terminology

- Gradient:  $l(\theta): \mathbb{R}^d \rightarrow \mathbb{R}$   
Gradient at  $\bar{\theta}$  is  $\nabla l(\bar{\theta}) \in \mathbb{R}^d$  where  $i$ th coord is  $\frac{\partial}{\partial \theta_i} (l(\bar{\theta}))$

- Diameter of  $C$ ,  $\|C\|_2$   
Max dist between 2 pts in  $C$

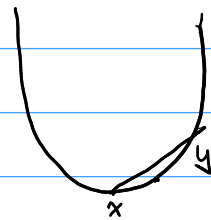
-  $l: C \rightarrow \mathbb{R}$  is convex if  $\forall x, y \in C, \forall t \in [0, 1]$

$$f(tx + (1-t)y) \leq tf(x) + (1-t)f(y)$$

-  $l: C \rightarrow \mathbb{R}$  is  $L$ -Lipschitz if  $\forall x, y \in C$

$$|l(x) - l(y)| \leq L \|x - y\|_2$$

$$\Rightarrow \|\nabla l\|_2 \leq L$$



## Loss Functions

$$x, \theta \in \mathbb{R}^d$$

1. Linear regression:  $y \in \mathbb{R}$ ,  $\ell(\theta, x, y) = (\langle x, \theta \rangle - y)^2$
2. Logistic Regression:  $y \in \{\pm 1\}$ ,  $\ell(\theta, x, y) = \log(1 + e^{-y \langle x, \theta \rangle})$
3. (Geometric) Median:  $\ell(\theta, x, y) = \|\theta - x\|_2$ , 1-Lipschitz.  
Mean:  $\ell(\theta, x, y) = \|\theta - x\|_2^2$   
 $\downarrow$   
 $\| \cdot \|_2$
4. Support Vector Machine (SVM):  $y \in \{\pm 1\}$ ,  $\ell(\theta, x, y) = \max(0, 1 - y \langle x, \theta \rangle)$



## Optimization

$$\theta^* = \arg \min_{\theta} \mathcal{J}(\theta, D)$$

Gradient descent

Assume non-priv. optimizer exists

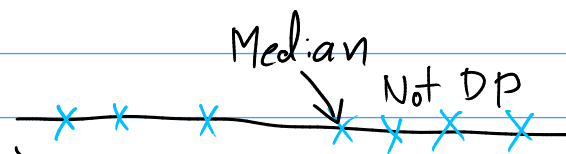
## Generalization

Uniform convergence + ERM  $\rightarrow$  generalization

## Privacy Considerations

$\hat{\theta}$  which minimizes  $\mathcal{L}(\cdot, D)$

$\hat{\theta}$  is DP wrt  $D$

Median: 

SVM: 

- Output perturbation
- Objective "
- Gradient "

- Input perturbation!  $\leftrightarrow$  LDP

# Output Perturbation

Chaudhuri, Monteleoni, Sarwate '11 (CMS'11)

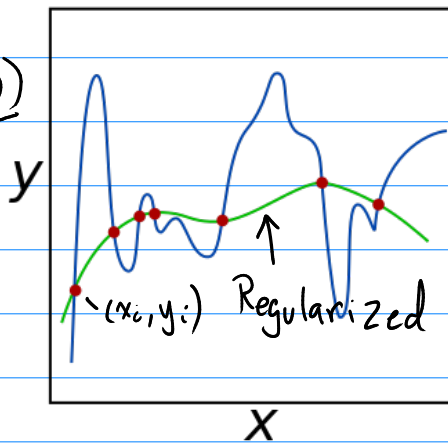
Regularization

Polynomial Regression.  $(x_i, y_i) \in \mathbb{R}^2$ .  $\theta \in \mathbb{R}^{k+1}$

$$l(\theta, x, y) = \left( \sum_{j=0}^k \theta_j x^j - y \right)^2 \quad \hat{y} = \sum \theta_j x^j$$

$$\mathcal{L}(\theta, D) = \sum_{k=n} l(\theta, x_i, y_i) + \lambda N(\theta)$$

$$N(\theta) = \sum \theta_j^2$$



- Poor gen.  
- Sensitive.

**Lemma 1** (Corollary 8 in [CMS11]). If  $N$  is differentiable and 1-strongly convex, and  $l$  is convex and 1-Lipschitz, then ~~for all  $D$~~  the  $\ell_2$ -sensitivity of  $(\arg \min_{\theta \in \mathcal{C}} \mathcal{L}(\theta, D))$  is at most  $\frac{2n}{\lambda}$ .

$$\hat{\theta} = \arg \min_{\theta \in \mathcal{C}} \mathcal{L}(\theta, D)$$

$$\hat{\theta} = \bar{\theta} + b \leftarrow \text{noise}$$

$$b \sim \mathcal{N}\left(0, O\left(\frac{n^2 \log(1/\delta)}{\lambda^2 \epsilon^2}\right) \cdot \text{Id}_{d \times d}\right)$$

( $\epsilon, \delta$ )-DP

## Objective Perturbation

Before:  $f(\theta, D) = \sum \ell(\theta, x_i, y_i) + \frac{\lambda}{2} \|\theta\|_2^2$

Random perturb  $f$

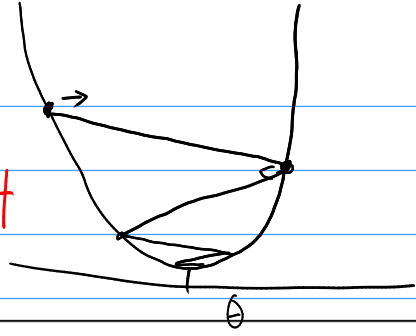
$$f^{\text{Priv}}(\theta, D) = \sum \ell(\theta, x_i, y_i) + \frac{\lambda}{2} \|\theta\|_2^2 + \langle b, \theta \rangle$$

$b$  is Gamma or Gaussian

1. Perturb  $f$  to get  $f^{\text{Priv}}$
2. Solve  $f^{\text{Priv}}$  non-privately

Drawbacks:

- Exact optimization req'd
  - ↳ Numerical prec Mironov '12
  - ↳ Iterative methods (!)
- ⇒ Approx optimum.
- Restrictive assns.



# Gradient Perturbation Non-Private Gradient Descent

## Algorithm 1: Projected Gradient Descent

```

Set  $\theta_0 \in C$  arbitrarily
for  $t = 1$  to  $T$  do
  | Compute  $\theta_t = \Pi_C(\theta_{t-1} - \eta(t) \nabla \mathcal{L}(\theta_{t-1}, D))$ 
end
return  $\hat{\theta} = \theta_T$ 

```

learning rate  
projection onto  $C$

**Theorem 2** (Theorem 2 from [SZ13]). Let  $F$  be a convex function and let  $\theta^* = \arg \min_{\theta \in C} F(\theta)$ . Let  $\theta_0$  be an arbitrary point in  $C$ , and  $\theta_{t+1} = \Pi_C(\theta_t - \eta(t) G_t(\theta_t))$ , where  $\mathbf{E}[G_t(\theta_t)] = \nabla F(\theta_t)$  and  $\mathbf{E}[\|G_t(\theta_t)\|_2^2] \leq G^2$ , and the learning rate function  $\eta(t) = \frac{\|C\|_2}{G\sqrt{t}}$ . Then for any  $T > 0$ , we have the following:

$$\mathbf{E}[F(\theta_T) - F(\theta^*)] = O\left(\frac{\|C\|_2 G \log T}{\sqrt{T}}\right).$$

$$\nabla f(\theta_{t+1}, D) = G_t(\theta_t)$$

$$G_t = n \cdot L$$

$$O\left(\frac{\|C\|_2 n L}{\sqrt{T}}\right)$$

$$T \geq \Omega\left(\frac{n^2 L^2 \|C\|_2^2}{\alpha^2}\right) \rightarrow \leq \alpha.$$

$T \cdot n$  Grad. comps  
 $\approx n^3$ ,

## Algorithm 2: Stochastic Projected Gradient Descent

```

Set  $\theta_0 \in C$  arbitrarily
for  $t = 1$  to  $T$  do
  | Select  $i \in [n]$  uniformly at random
  | Compute  $\theta_t = \Pi_C(\theta_{t-1} - \eta(t) (n \cdot \nabla \ell(\theta_{t-1}, x_i, y_i)))$ 
end
return  $\hat{\theta} = \theta_T$ 

```

$\nabla f(\theta_t, D)$

$$G_t(\theta_t) = n \nabla \ell(\theta_t, x_i, y_i)$$

$$\mathbf{E}[G_t(\theta_t)] = \nabla f(\theta_t, D)$$

$$\mathbf{E}[\|G_t(\theta_t)\|_2^2] \leq n^2 L^2 \rightarrow G_t = nL$$

$Tn \approx n^2$  comps. vs  $n^3$



# Private Stochastic Gradient Descent

**Algorithm 3:** Private Stochastic Projected Gradient Descent

Define  $\sigma^2 \leftarrow \frac{32L^2 n^2 \log(n/\delta) \log(1/\delta)}{\varepsilon^2}$

Set  $\theta_0 \in \mathcal{C}$  arbitrarily

**for**  $t = 1$  to  $n^2$  **do**

    Select  $i \in [n]$  uniformly at random

    Compute  $\theta_t = \Pi_{\mathcal{C}}(\theta_{t-1} - \eta(t)(n \cdot \nabla \ell(\theta_{t-1}, x_i, y_i) + b_{t-1}))$ , where  $b_{t-1} \sim N(0, \sigma^2 \cdot I_{d \times d})$

**end**

**return**  $\hat{\theta} = \theta_T$

## Utility

$$G_t(\theta_t) = n \cdot \nabla \ell(\theta_t, x_i, y_i) + b_t$$

$$E[G_t(\theta_t)] = \nabla \mathcal{L}(\theta_t, \mathcal{D}) + \theta$$

$$\begin{aligned} E[\|G_t\|_2^2] &= n^2 E[\|\nabla \ell\|_2^2] + 2n E[\langle \nabla \ell, b_t \rangle] + E[\|b_t\|_2^2] \\ &\leq n^2 L^2 + 0 + d\sigma^2 \end{aligned}$$

$$\begin{aligned} E[\mathcal{L}(\theta_T, \mathcal{D}) - \mathcal{L}(\theta^*, \mathcal{D})] &\leq \tilde{O}\left(\frac{\|\mathcal{C}\|_2 \sqrt{n^2 L^2 + d\sigma^2}}{\sqrt{n^2}}\right) \\ &= \tilde{O}\left(\frac{\|\mathcal{C}\|_2 L \sqrt{d \log(1/\delta)}}{\varepsilon}\right) \end{aligned}$$

# Privacy

1. Gaussian Mech
2. Amp. by subsamp.
3. Adv. comp.

1.  $n \cdot \nabla \ell(\theta_{t-1}, x_i, y_i) \quad \|\nabla \ell\|_2 \leq L$

$l_2$ -sens  $\leq 2nL$

Assume  $i$  is fixed

$\forall t=1$  to  $n^2$

G.M.  $\Rightarrow$  priv loss RV for  $G_t(\theta_t) \leq \frac{\epsilon}{2\sqrt{\log(1/\delta)}}$  w.p  $1 - \frac{\delta}{2}$ . ✓

2. Lemma: If  $A$  is  $\epsilon' < 1$  DP, if executed on a <sup>size  $\gamma n$</sup>  subsample ~~from~~ a dataset of size  $n$ , then result is  $2\gamma\epsilon'$ -DP.

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

$\forall t=1$  to  $n^2$

$\epsilon' = \frac{\epsilon}{2\sqrt{\log(1/\delta)}} \rightarrow \frac{\epsilon}{n\sqrt{\log(1/\delta)}} \text{-DP w.p } 1 - \frac{\delta}{2}$ .

3.  $n^2$  iterations

Basic comp  $(\epsilon n, \frac{\delta}{2})$ -DP

Advanced  $(\epsilon, \delta)$ -DP

---

### Algorithm 3: Private Stochastic Projected Gradient

---

Define  $\sigma^2 \leftarrow \frac{32L^2 n^2 \log(n/\delta) \log(1/\delta)}{\epsilon^2}$

Set  $\theta_0 \in \mathcal{C}$  arbitrarily

for  $t = 1$  to  $n^2$  do

    | Select  $i \in [n]$  uniformly at random

    | Compute  $\theta_t = \Pi_{\mathcal{C}}(\theta_{t-1} - \eta(t)(n \cdot \nabla \ell(\theta_{t-1}, x_i, y_i)))$

end

return  $\hat{\theta} = \theta_T$

---