

CSC2412: Private Gradient Descent & Empirical Risk Minimization

Sasho Nikolov

Empirical Risk Minimization

Learning: Reminder

- Known data universe \mathcal{X} and an unknown probability distribution D on \mathcal{X}
- Known concept class C and an unknown concept $c \in C$
- We get a dataset $X = \{(x_1, c(x_1)), \dots, (x_n, c(x_n))\}$, where each x_i is an independent sample from D .

Goal: Learn c from X .

Binary loss \rightarrow

$$l(c', (x, y)) = \begin{cases} 1 & c'(x) \neq y \\ 0 & c'(x) = y \end{cases}$$

data pt
label

$$\mathbb{P}_{x \sim D}(c'(x) \neq c(x)) = \mathbb{E}_{x \sim D} \left[\underbrace{1 \{c'(x) \neq c(x)\}}_{= l(c', (x, c(x)))} \right]$$

$$\stackrel{=}{=} L_{D,c}(c') = \mathbb{E}_{x \sim D} [l(c', (x, c(x)))] = \mathbb{P}_{x \sim D}(c'(x) \neq c(x))$$

We want an algorithm \mathcal{M} that outputs some $c' \in \mathcal{C}$ and satisfies

$$\mathbb{P}(L_{D,c}(\mathcal{M}(X)) \leq \alpha) \geq 1 - \beta.$$

Agnostic learning

Maybe no concept gives 100% correct labels. \rightarrow agnostic setting
(as opposed to realizable)

Generally, we have a distribution D on $\mathcal{X} \times \{-1, +1\}$. distribution on labeled examples

$$L_D(c) = \mathbb{E}_{(x,y) \sim D}[\ell(c, (x,y))] = \mathbb{P}_{(x,y) \sim D}(c(x) \neq y)$$

D is unknown but we are given iid samples $X = \{(x_1, y_1), \dots, (x_n, y_n)\}$.
 $(x_i, y_i) \sim_{\text{iid}} D$

We want an algorithm \mathcal{M} that outputs some $c' \in C$ and satisfies

$$\mathbb{P}(L_D(\mathcal{M}(X)) \leq \underbrace{\min_{c \in C} L_D(c)}_{\text{best possible loss achievable by } C} + \alpha) \geq 1 - \beta.$$

best possible loss
achievable by C

Empirical risk minimization, again

Issue: We want to find $\arg \min_{c \in C} L_D(c)$, but we do not know D .

Solution: Instead we solve $\arg \min_{c \in C} L_X(c)$, where

$$L_X(c) = \frac{1}{n} \sum_{i=1}^n \ell(c, (x_i, y_i)). \quad = \frac{\sum_{i=1}^n \mathbb{1}\{c(x_i) \neq y_i\}}{n}$$

is the empirical error.

for binary loss

Theorem (Uniform convergence)

Suppose that $n \geq \frac{\ln(|C|/\beta)}{2\alpha^2}$. Then, with probability $\geq 1 - \beta$,

~~$\max_{c \in C} L_D(c) - L_X(c) \leq \alpha$~~

$$\max_{c \in C} L_D(c) - L_X(c) \leq \alpha$$

Example: Linear Separators

- $x = [0, 1]^d$ unit cube in \mathbb{R}^d

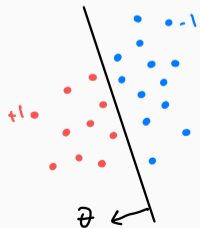
$$\langle u, v \rangle = \sum_{i=1}^d u_i v_i \quad \text{sign}(z) = \begin{cases} +1 & z > 0 \\ 0 & z = 0 \\ -1 & z < 0 \end{cases}$$

- \underline{C} is all functions of the type $c_\theta(x) = \text{sign}(\langle x, \theta \rangle + \theta_0)$ for $\theta \in \mathbb{R}^d, \theta_0 \in \mathbb{R}$.

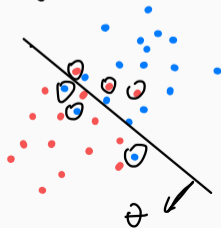
For convenience, replace x by $(x, 1) \in [0, 1]^{d+1}$ and θ, θ_0 by $(\theta, \theta_0) \in \mathbb{R}^{d+1}$

$$c_\theta(x) = \text{sign}(\langle x, \theta \rangle)$$

Realizable



Agnostic



$$c_\theta(x) = \begin{cases} +1 & \text{if } x \text{ "below" the plane} \\ -1 & \text{if above} \end{cases}$$

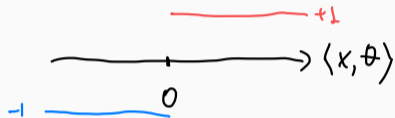
$$\langle x, \theta \rangle + \theta_0 = \left\langle \begin{pmatrix} x \\ 1 \end{pmatrix}, \begin{pmatrix} \theta \\ \theta_0 \end{pmatrix} \right\rangle$$

From now, will ignore θ_0

Finding best separator is generally computationally hard 6

Logistic Regression

Sign for sigmoid: given θ and x predict $\begin{cases} +1 & \text{w/ prob. } \frac{1}{1+e^{-\langle x, \theta \rangle}} \\ -1 & \text{w/ prob. } \frac{1}{1+e^{\langle x, \theta \rangle}} \end{cases}$



Logistic loss

$$\ell(\theta, (x, y)) = \log \left(\frac{1}{\mathbb{P}(\text{predict } y \text{ from } \langle x, \theta \rangle)} \right) = \log(1 + e^{-y \cdot \langle x, \theta \rangle}).$$

\downarrow
 c_θ

Logistic regression: Given $X = \{(x_1, y_1), \dots, (x_n, y_n)\}$ solve

can be efficiently
solved

$$\arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i \cdot \langle x_i, \theta \rangle})$$

Empirical loss
Connection to population
loss in another class

(Private) Gradient Descent

Convex loss

The function $L_X(\theta) = \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i \cdot \langle x_i, \theta \rangle})$ is convex in θ .

$$\forall \theta, \theta' \in \mathbb{R}^{d+1} : L_X\left(\frac{\theta + \theta'}{2}\right) \leq \frac{1}{2} L_X(\theta) + \frac{1}{2} L_X(\theta')$$

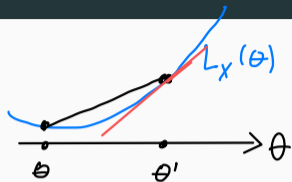
$$\Downarrow \quad \forall \lambda \in [0, 1] \quad L_X(\lambda\theta + (1-\lambda)\theta') \leq \lambda L_X(\theta) + (1-\lambda)L_X(\theta')$$

$$\forall \theta, \theta' : L_X(\theta) \geq L_X(\theta') + \langle \nabla L_X(\theta'), \theta - \theta' \rangle$$

$$\nabla L_X(\theta) = \begin{pmatrix} \frac{\partial L_X(\theta)}{\partial \theta_0} \\ \vdots \\ \frac{\partial L_X(\theta)}{\partial \theta_d} \end{pmatrix}$$

Convex functions can be minimized efficiently.

- for non-convex, it's complicated



Gradient descent

Logistic regression

$$\nabla L_x(\theta) = \begin{pmatrix} \frac{\partial L_x(\theta)}{\partial \theta_0} \\ \vdots \\ \frac{\partial L_x(\theta)}{\partial \theta_d} \end{pmatrix}$$

$$\theta_0 = 0$$

for $t = 1 \dots T - 1$ do

$$\tilde{\theta}_t = \theta_{t-1} - \eta (\nabla L_x(\theta_{t-1}) + z_t)$$

$$\theta_t = \tilde{\theta}_t / \max\{1, \|\tilde{\theta}_t\|_2 / R\}$$

end for

output $\frac{1}{T} \sum_{t=0}^{T-1} \theta_t$

$$\arg \min_{\theta \in B_2^{d+1}(R)} \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i \cdot \langle x_i, \theta \rangle})$$

$L_x(\theta)$

$$B_2^{d+1}(R) = \{\theta : \|\theta\|_2 \leq R\}$$

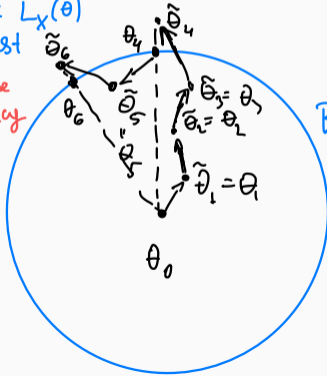
$$\sqrt{\sum \theta_i^2}$$

parameter

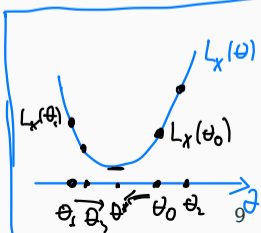
$-\nabla L_x(\theta)$: direction in which $L_x(\theta)$ decreases the fastest

Gauss noise for privacy

projection into $B_2^{d+1}(R)$



$B_2(R)$



Advanced composition - warmup

Publish k functions $f_1, \dots, f_k : \mathcal{X}^n \rightarrow \mathbb{R}^d$ with (ϵ, δ) -DP, where $\forall i \Delta_2 f_i \leq C$

Want to release $f_1(x), f_2(x), \dots, f_k(x)$. Can achieve noise / function

f_i could be adaptive: f_i depends on $f_1(x) + z_1, \dots, f_{i-1}(x) + z_{i-1}$

$$\approx \sqrt{k \log k}$$

1) Apply Gaussian mechanism k times, use composition: \rightarrow noise $\approx Ck \log k$

Release $f_i(x) + z_i$ $z_i \sim N(0, \frac{C^2}{\rho} I)$ for $\rho \approx \frac{\epsilon^2/k^2}{\log(k/\delta)}$, to achieve $(\frac{\epsilon}{k}, \frac{\delta}{k})$ -DP

By composition then $(f_1(x) + z_1, \dots, f_k(x) + z_k)$ is (ϵ, δ) -DP

2) $g(x) = \begin{pmatrix} f_1(x) \\ f_2(x) \\ \vdots \\ f_k(x) \end{pmatrix} \in \mathbb{R}^{dk}$; Use the Gaussian noise mechanism for g . \rightarrow noise $\approx C\sqrt{k}$

$$(\Delta_2 g)^2 = \max_{x \sim x'} \|g(x) - g(x')\|_2^2 = \max_{x \sim x'} \sum_{i=1}^k \|f_i(x) - f_i(x')\|_2^2 \leq k C^2$$

Advanced composition (for Gaussian noise)

$$f_i: \mathcal{X}^n \rightarrow \mathbb{R}^d$$

Suppose we release $Y_1 = f_1(X) + Z_1, \dots, Y_k = f_k(X) + Z_k$ where $f_i: \mathcal{X}^n \rightarrow \mathbb{R}^d$ depends also on Y_1, \dots, Y_{i-1} .

$$Z_i \stackrel{\text{e}}{\sim} \mathcal{N}\left(0, \frac{(\Delta_2 f)^2}{\rho} \cdot I\right)$$

Then the output (Y_1, \dots, Y_k) satisfies (ϵ, δ) -DP for $\epsilon = k\rho + \sqrt{2k\rho \ln(1/\delta)}$.

\Rightarrow for $\rho \approx \frac{\epsilon^2}{k \log(1/\delta)}$ to achieve (ϵ, δ) -DP

\Rightarrow noise per query f_i is $\approx \frac{\Delta_2 f \cdot \sqrt{k} \cdot \sqrt{\log(1/\delta)}}{\epsilon}$

Same as if queries were not adaptive

Sensitivity of gradients

$$X = \{(x_1, y_1), \dots, (x_n, y_n)\} \quad \nabla L_X(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla \log(1 + e^{-y_i \cdot \langle x_i, \theta \rangle})$$

$$X' = \{(x'_1, y'_1), \dots, (x'_n, y'_n)\}$$

Suppose $X = [-1, +1]^{d+1}$.

$$\Delta_2 \nabla L_X(\theta) = \max_{X \sim X'} \|\nabla L_X(\theta) - \nabla L_{X'}(\theta)\|_2 = \max_{\substack{(x, y) \\ (x', y')}} \frac{1}{n} \|\nabla \log(1 + e^{-y \cdot \langle x, \theta \rangle}) - \nabla \log(1 + e^{-y' \cdot \langle x', \theta \rangle})\|_2$$

$$\leq \frac{2}{n} \cdot \max_{(x, y)} \|\nabla \log(1 + e^{-y \cdot \langle x, \theta \rangle})\|_2$$

$$\leq \frac{2\sqrt{d+1}}{n}$$

$$\nabla \log(1 + e^{-y \cdot \langle x, \theta \rangle}) = -\frac{1}{1 + e^{y \cdot \langle x, \theta \rangle}} yx$$

$$y \in \{\pm 1\}$$

$$x \in [-1, +1]^{d+1}$$

$$\|\nabla \log(1 + e^{-y \cdot \langle x, \theta \rangle})\|_2 = \left\| \frac{1}{1 + e^{y \cdot \langle x, \theta \rangle}} \cdot yx \right\|_2 \leq \|x\|_2 \leq \sqrt{d+1}$$



Private gradient descent

Think of $\nabla L_X(\theta_0), \nabla L_X(\theta_1), \dots;$
as the adaptively chosen functions $f_1, f_2, \dots;$

$$\theta_0^0 = 0$$

for $t = 1 \dots T - 1$ do

$$\tilde{\theta}_t^0 = \theta_{t-1}^0 - \eta(\nabla L_X(\theta_{t-1}^0)) + \underline{Z_{t,0}}$$

$$\theta_t^0 = \tilde{\theta}_t^0 / \max\{1, \|\tilde{\theta}_t^0\|_2 / R\}$$

end for

$$\text{output } \frac{1}{T} \sum_{t=0}^{T-1} \theta_t^0$$

\forall_t

$$Z_t \sim N\left(0, \frac{4d+1}{n^2 \rho} \mathbf{I}\right)$$

$$\rho \approx \frac{\epsilon^2}{T \log(1/\delta)}$$

$$\rho^2 \approx \frac{dT \cdot \log(1/\delta)}{\epsilon^2 n^2}$$

(ϵ, δ) -DP by advanced composition
+ post-processing

Accuracy analysis

Theorem

Suppose $\mathbb{E} \|\nabla L_X(\theta_t) + Z_t\|_2^2 \leq B^2$ for all t . For $\eta = \frac{R}{BT^{1/2}}$ we have

$$\mathbb{E} \left[L_X \left(\frac{1}{T} \sum_{t=0}^{T-1} \theta_t \right) \right] \leq \underbrace{\min_{\theta \in B_2^{d+1}(R)} L_X(\theta)}_{\substack{\downarrow \\ \text{optimal} \\ \text{value}}} + \frac{RB}{T^{1/2}}$$

↓ goes to 0
as $T \rightarrow \infty$

→ Proof in the notes (optional)

Plugging in

$$\mathbb{E} \|\nabla L_X(\theta_t) + Z_t\|_2^2 = \mathbb{E} \|\nabla L_X(\theta_t)\|_2^2 + \mathbb{E} \|Z_t\|_2^2 \leq \frac{d+1}{\epsilon^2} + \sigma^2(d+1) = B^2$$

$$\mathbb{E} \sum_{i=0}^d (z_t)_i^2 = \sigma^2(d+1)$$

$$\sigma^2 \approx \frac{dT \cdot \log(1/\delta)}{\epsilon^2 n^2}$$

$$B^2 \approx d + \frac{d^2 T \log(1/\delta)}{\epsilon^2 n^2}$$

For any θ $\|\nabla L_X(\theta)\|_2 \leq \sqrt{d+1}$

$$\text{error} = \frac{RB}{T^{1/2}} \approx \frac{R\sqrt{d}}{T^{1/2}} + \frac{Rd \sqrt{\log(1/\delta)}}{\epsilon n}$$

$$\begin{aligned} \|\nabla L_X(\theta)\|_2 &= \left\| \frac{1}{n} \sum \nabla \log(1 + e^{-y_i \langle x_i, \theta \rangle}) \right\|_2 \\ &\leq \frac{1}{n} \sum \|\nabla \log(1 + e^{-y_i \langle x_i, \theta \rangle})\|_2 \leq \sqrt{d+1} \end{aligned}$$