

DP-GD for Generalized Linear Models

Today

- Generalized Linear Models
- DP-GD, again
- Convergence analysis

Tools for DP optimization

- Better composition
- Convergence analysis
- Privacy amplification by subsampling
- DP-FTRL

GLMs

Prediction task: $(x_1, y_1), \dots, (x_n, y_n)$

$$x_i \in \mathbb{R}^d$$

y_i discrete or $\in \mathbb{R}$.

Modeling assumption: Our prediction depends on parameters θ through $\langle x_i, \theta \rangle$

loss function: $l(\theta; x_i, y_i) = l(\langle x_i, \theta \rangle, y_i)$

Example 1 Linear regression, $y_i \in \mathbb{R}$

Predict $\hat{y}_i = \langle x_i, \theta \rangle$

Squared loss: $l(\theta; x_i, y_i) = (\langle x_i, \theta \rangle - y_i)^2$

Example 2 Logistic regression, $y_i \in \{0, 1\}$

Predict $\hat{y}_i = \frac{1}{1 + e^{-\langle x_i, \theta \rangle}}$

cross-entropy loss: $l(\theta; x_i, y_i) = -y_i \ln(\hat{y}_i) - (1 - y_i) \ln(1 - \hat{y}_i)$

Claim For some scalar a ,

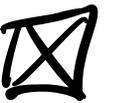
$$\nabla_{\theta} l(\langle x_i, \theta \rangle, y_i) = a x_i$$

Proof Let $z = \langle x_i, \theta \rangle$. For any $j \in [d]$

$$\frac{\partial}{\partial \theta_j} l(\langle x_i, \theta \rangle, y_i) = \frac{\partial l}{\partial z} \frac{\partial z}{\partial \theta_j}$$

indep of j

$\hookrightarrow x_{ij}$



Algorithm 1 (Noisy Gradient Descent)

Input: data $(x_1, y_1), \dots, (x_n, y_n)$; loss function l ;
num iterations T ;
noise scale σ ; learning rate η

Output: $\tilde{\theta} \in \mathbb{R}^d$

① $\theta_0 \leftarrow \vec{0}$

② For $t=1, \dots, T$

③ $\tilde{g}_t \leftarrow \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} l(\theta_t; x_i, y_i) + z_t, z_t \sim \mathcal{N}(0, \sigma^2)$

④ $\theta_{t+1} \leftarrow \theta_t - \eta \tilde{g}_t$

⑤ Return $\tilde{\theta} = \frac{1}{T} \sum_{t=1}^T \theta_t$

Privacy?

Claim If $\forall \theta \in \Theta$ and $(x, y), \|\nabla_{\theta} l(\theta; x, y)\|_2 \leq L$,
then Alg 1 is (ϵ, δ) -DP for

$$\sigma = \frac{2L\sqrt{T} \sqrt{\log(1/\delta)}}{n\epsilon}$$

How can we be sure gradient is not too big? Clipping!

Def A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is L-Lipschitz
if $\forall x, y \quad |f(x) - f(y)| \leq L \cdot \|x - y\|_2$

Fact If $\nabla f(x)$ exists, then $\|\nabla f(x)\|_2 \leq L$.

Accuracy? Need to control effect of noise.

Convergence for GLMs

Theorem (Song, Steinke, Thakkar, Thakurta 2021)

Suppose loss function $\ell(\langle x, \theta \rangle, y)$ is convex in its first parameter and L -Lipschitz over $\theta \forall x, y$.

$$\text{Let } \mathcal{L}(\theta; X, Y) = \frac{1}{n} \sum_{i=1}^n \ell(\langle \theta, x_i \rangle, y_i).$$

Let $\theta^* = \arg\min_{\theta \in D} \mathcal{L}(\theta; D)$ and M the projector to the eigenspace of $\sum_{i=1}^n x_i x_i^T$.

For Alg 1 with $\sigma = \frac{2L\sqrt{T} \sqrt{\log 1/\delta}}{n\epsilon}$, $T = n^2 \epsilon^{-2}$,

& some learning rate, we have

$$\mathbb{E}[\mathcal{L}(\hat{\theta}; X, Y)] - \mathcal{L}(\theta^*; X, Y) \leq \frac{L \|\theta^*\|_M \sqrt{1 + 2\text{rank}(M) \log 1/\delta}}{\epsilon n}$$

where $\|v\|_M = v^T M v$.

(Walk through this)

$$\hat{\theta}^* = \frac{\sqrt{\|\theta^*\|_M^2}}{\sqrt{T(L^2 + \text{rank}(M)\sigma^2)}} \\ f(x) = \frac{A}{x} + Bx, \quad x > 0 \\ \Rightarrow x^2 = \frac{A}{B}$$

Fact If f is convex, ^{& differentiable} then for any x, y ,
 $f(x) \geq f(y) + \langle \nabla f(y), x - y \rangle$

Idea: in \mathbb{R}^d , Taylor's thm

$$f(x) \approx f(a) + f'(a)(x-a)$$



Jensen's ineq

Proof $\mathcal{L}(\bar{\theta}) \stackrel{\downarrow}{\leq} \frac{1}{T} \sum_{t=1}^T \mathcal{L}(\theta_t)$

$$\begin{aligned} \mathcal{L}(\bar{\theta}) - \mathcal{L}(\theta^*) &\leq \frac{1}{T} \sum_{t=1}^T (\mathcal{L}(\theta_t) - \mathcal{L}(\theta^*)) \\ &\leq \frac{1}{T} \sum_{t=1}^T \langle \nabla \mathcal{L}(\theta_t), \theta_t - \theta^* \rangle \end{aligned}$$

Potential function

$$\Psi_t(\theta) = \mathbb{E}_{z_1, \dots, z_t} [\|\theta - \theta^*\|_m^2]$$

apply to RV.

$$\Psi_t(\theta_{t+1}) = \mathbb{E}_{z_{1:t}} [\|\theta_{t+1} - \theta^*\|_M^2]$$

$$= \mathbb{E}_{z_{1:t}} [\|\theta_t - \eta(g_t + z_t) - \theta^*\|_M^2]$$

$$= \mathbb{E}_{z_{1:t}} [\|\theta_t - \theta^*\|_M^2 - 2\eta \langle \theta_t - \theta^*, g_t + z_t \rangle + \eta^2 \|g_t + z_t\|_M^2]$$

$$= \underbrace{\Psi_t(\theta_t)}_{(i)} - 2\eta \underbrace{\mathbb{E}_{z_{1:t}} [\langle \theta_t - \theta^*, g_t + z_t \rangle]}_{(ii)} + \eta^2 \underbrace{\mathbb{E}_{z_{1:t}} [\|g_t + z_t\|_M^2]}_{(iii)}$$

$$i) \quad \Psi_t(\theta_t) = \Psi_{t-1}(\theta_t)$$

$$ii) \quad \mathbb{E}_{z_{1:t}} [\langle \theta_t - \theta^*, g_t + z_t \rangle] = \mathbb{E}_{\substack{z_{1:t-1} \\ z_t}} [\langle \theta_t - \theta^*, g_t \rangle + \langle \theta_t - \theta^*, z_t \rangle]$$
$$= \mathbb{E}_{z_{1:t-1}} [\langle \theta_t - \theta^*, g_t \rangle]$$

$$\text{iii) } \mathbb{E}_{z_{1:t}} \left[\|g_t + z_t\|_M^2 \right] \leq L^2 + \text{rank}(M) \sigma^2$$