

Lecture 14: Smooth optimization

So far:

* Setting $\hat{R} = \mathcal{L} \circ H$ where \mathcal{L} convex & Lipschitz, H Lipschitz (kind of) & "linearizable".

This time: $\mathcal{L} \circ H$ is smooth, meaning $\|\partial(\mathcal{L} \circ H)(a) - \partial(\mathcal{L} \circ H)(b)\| \leq \beta \|a-b\|$.

After this, one more near initialization setting: \mathcal{L} strongly convex

$$\begin{aligned} & \langle \partial(\mathcal{L} \circ H)(w_s), z - w_s \rangle \\ & \quad \circlearrowleft \partial H \partial \mathcal{L} \\ & \langle \partial \mathcal{L}(H(w_s)), \partial H(w_s)^T (z - w_s) \rangle \\ & = \langle \partial \mathcal{L}(H(w_s)), H(z) - H(w_s) \rangle \\ & \quad + \langle \partial \mathcal{L}(H(w_s)), \underbrace{H(w_s) - H(z) + \partial H(w_s)^T (z - w_s)}_{\text{Small near initialization}} \rangle \end{aligned}$$

Why smooth:

- * last iterate "for free"
- * guaranteed improvement in each iteration

- * mimics what we got from gradient flow
- * gives local optimality guarantees with no further assumptions

$$\underbrace{\|\partial \hat{R}(a) - \partial \hat{R}(b)\| \leq \beta \|a-b\|}_{\text{classical opt}} \quad \Rightarrow \quad \underbrace{\left| \hat{R}(a) - \hat{R}(b) - \langle \nabla \hat{R}(b), b-a \rangle \right| \leq \frac{\beta}{2} \|a-b\|^2}_{\text{only use this version in proofs}}$$

if $\hat{R} = \mathcal{L} \circ H$, then

$$\begin{aligned} \|\partial(\mathcal{L} \circ H)(a) - \partial(\mathcal{L} \circ H)(b)\| &= \|\partial H(a) \partial \mathcal{L}(H(a)) - \partial H(b) \partial \mathcal{L}(H(b))\| \\ &= \|\partial H(a) \partial \mathcal{L}(H(a)) - \partial H(a) \partial \mathcal{L}(H(b)) + \partial H(a) \partial \mathcal{L}(H(b)) - \partial H(b) \partial \mathcal{L}(H(b))\| \\ &\leq \|\partial H(a) [\partial \mathcal{L}(H(a)) - \partial \mathcal{L}(H(b))]\| + \|[\partial H(a) - \partial H(b)] \partial \mathcal{L}(H(b))\| \\ &\leq \underbrace{\|\partial H(a)\|_{2 \rightarrow 2}}_{\text{smooth}} \underbrace{\|\partial \mathcal{L}(H(a)) - \partial \mathcal{L}(H(b))\|}_{\text{smooth}} + \underbrace{\|\partial \mathcal{L}(H(b))\|_2}_{\text{smooth}} \cdot \underbrace{\|\partial H(a) - \partial H(b)\|_{2 \rightarrow 2}}_{\text{smooth}} \end{aligned}$$

hard to satisfy (for now).

$$\begin{aligned} \|A_x\| &\leq \|A\|_{2 \rightarrow 2} \|x\|, \quad \text{where } \|A\|_{2 \rightarrow 2} = \max_{\|v\|_2 \leq 1} \|Av\|_2 \\ &= \max_{v \neq 0} \frac{\|Av\|_2}{\|v\|_2} \\ \|A\|_{p \rightarrow q} &= \max_{v \neq 0} \frac{\|Av\|_q}{\|v\|_p} \end{aligned}$$

Consequences to GD: ① local optimality, ② last iterate guarantee when convex.

① Local optimality (without projection)

Recall GD: $w_{s+1} := w_s - \eta \nabla \hat{R}(w_s)$

By smoothness:
$$\begin{aligned} \hat{R}(w_{s+1}) &= \hat{R}(w_s - \eta \nabla \hat{R}(w_s)) \leq \hat{R}(w_s) + \langle \nabla \hat{R}(w_s), w_s - \eta \nabla \hat{R}(w_s) - w_s \rangle \\ &= \hat{R}(w_s) - \eta \left(1 - \frac{\eta \beta}{2}\right) \|\nabla \hat{R}(w_s)\|^2 \end{aligned}$$

Note $\hat{R}(w_{s+1}) \leq \hat{R}(w_s)$ when $\eta \leq \frac{2}{\beta}$

Theorem Suppose \hat{R} is β -smooth, not necessarily convex.

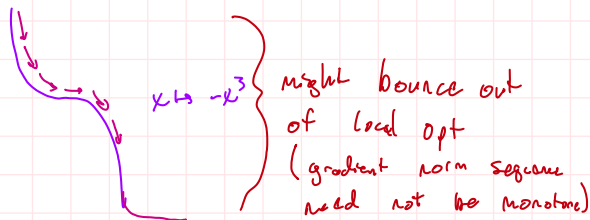
① If $\eta \in [0, \frac{2}{\beta}] \Rightarrow \hat{R}(w_s)$ is nonincreasing

② If $\eta \in [0, \frac{2}{\beta}] \Rightarrow \min_{s \leq t} \|\nabla \hat{R}(w_s)\|^2 \leq \frac{1}{t} \sum_{s \leq t} \|\nabla \hat{R}(w_s)\|^2 \quad \textcircled{A}$

$$\leq \frac{1}{t} \left(\frac{2}{\eta(2-\eta\beta)} \right) (\hat{R}(w_0) - \inf_{v \in \mathbb{R}^p} \hat{R}(v))$$

[Note if $\eta = \frac{1}{\beta}$: $\min_{s \leq t} \|\nabla \hat{R}(w_s)\|^2 \leq \frac{2\beta}{t} (\hat{R}(w_0) - \hat{R}(w_t)) \leq \frac{2\beta}{t} (\hat{R}(w_0) - \underbrace{\inf_{v \in \mathbb{R}^p} \hat{R}(v)}_{\text{e.g., } \geq 0})$

Remark (interpretation). $\exists s \leq t$ s.t. $\|\nabla \hat{R}(w_s)\|^2 = O(\frac{1}{t})$



$\hat{R}(w_{s+1})$

$$= \hat{R}(w_s) - \eta \left(1 - \frac{\eta \beta}{2}\right) \|\nabla \hat{R}(w_s)\|^2$$

2

Proof of A. From before:

$$\frac{\eta}{2} (2 - \eta \beta) \|\nabla \hat{R}(w_s)\|^2 \leq \hat{R}(w_s) - \hat{R}(w_{s+1})$$

\rightarrow apply $\frac{1}{t} \sum_{s \leq t} (\cdot)$

$$\frac{\eta}{2} (2 - \eta \beta) \frac{1}{t} \sum_{s \leq t} \|\nabla \hat{R}(w_s)\|^2$$

$$\leq \hat{R}(w_0) - \hat{R}(w_t)$$

$$\Rightarrow \frac{1}{t} \sum_{s \leq t} \|\nabla \hat{R}(w_s)\|^2 \leq \frac{2}{\eta(2-\eta\beta)} (\hat{R}(w_0) - \hat{R}(w_t))$$

Theorem. Suppose \hat{R} is convex and β -smooth, & $\eta \leq \frac{1}{\beta}$.
 $\forall z \in \mathbb{R}^d$,

$$\frac{1}{2t\eta} \|w_t - z\|^2 + \hat{R}(w_t) \leq \frac{1}{2t\eta} \|w_t - z\|^2 + \frac{1}{t} \sum_{s=1}^t \hat{R}(w_s) \leq \frac{1}{2t\eta} \|w_0 - z\|^2 + \hat{R}(z).$$

Proof. Recall

$$\|w_t - z\|^2 + 2\eta \sum_{s=1}^t \hat{R}(w_s) \leq \|w_0 - z\|^2 + 2t\eta \hat{R}(z) + \sum_{s=1}^t \eta^2 \|\nabla \hat{R}(w_s)\|^2 \quad (\text{MD lemma})$$

$$\sum_{s=1}^t \|\nabla \hat{R}(w_s)\|^2 \leq \frac{2}{\eta(2-\eta\beta)} (\hat{R}(w_0) - \hat{R}(w_t)) \quad (\text{smoothness guarantee})$$

$$\leq \frac{2}{\eta} (\hat{R}(w_0) - \hat{R}(w_t))$$

combine: $\|w_t - z\|^2 + 2\eta \sum_{s=1}^t \hat{R}(w_s) \leq \|w_0 - z\|^2 + 2t\eta \hat{R}(z) + \underbrace{\eta^2 \left[\frac{2}{\eta} (\hat{R}(w_0) - \hat{R}(w_t)) \right]}_{2\eta}$

$$\Rightarrow \|w_0 - z\|^2 + 2\eta \sum_{s=1}^t \hat{R}(w_s) \leq \|w_0 - z\|^2 + 2t\eta \hat{R}(z).$$

Then divide by $2t\eta$.

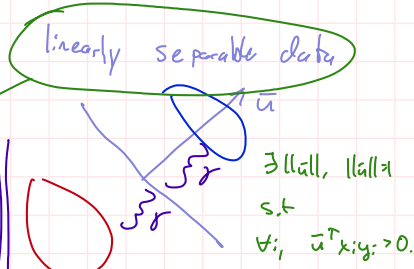
Example. logistic regression $l(z) = \ln(1 + \exp(-z))$.

prediction $x \mapsto w^T x$.

Define $\gamma := \min_i \bar{u}^T x_i \cdot y_i > 0$,
 and $z := \bar{u} \frac{\ln t}{\gamma}$, } \hat{R} has no minimum!

whereby $\hat{R}(z) = \frac{1}{n} \sum_i \ln(1 + \exp(-y_i x_i^T z))$
 $\leq \frac{1}{n} \sum_i \exp(-y_i x_i^T z)$
 $= \frac{1}{n} \sum_i \exp(-y_i x_i^T \bar{u} \frac{\ln t}{\gamma}) \leq \frac{1}{n} \sum_i \exp(-\ln t) = \frac{1}{t}.$

(Remark)
 2012-2021: DL gets zero training error (aka interpolation)
 2022: large language model training; not zero training error.



By extremely cute theorem with step size $\eta=1$, $w_0=0$

$$\frac{1}{2t} \|w_t - z\|^2 + \hat{R}(w_t) \leq \frac{1}{2t} \left\| \frac{\ln t}{\gamma} \right\|^2 + \hat{R}(z)$$

$$= \frac{1}{t} \left(\frac{(\ln t)^2}{2\gamma^2} + 1 \right).$$