

Lecture 13: opt: convex \mathcal{L} & Lip H ; smooth $\mathcal{L} \circ H$.

Outline for near initialization proofs:

* Recall standard convex opt settings & rates

* ^{convex} Lipschitz & Bounded $\Rightarrow \frac{1}{\sqrt{t}}$ to global

* Smooth $\Rightarrow \frac{1}{t}$ to local opt

* smooth & convex $\Rightarrow \frac{1}{t}$ to global opt

* smooth & strongly convex $\Rightarrow e^{-t}$ to global opt

instead consider $\nabla(\mathcal{L} \circ H)$
left column column

strategy: make minimal assumptions
on H so that we
get these rates

+ worst case terms for H .

Remark: recovers, extends, & unifies
proofs from the literature

Lemma (from last time: Lipschitz & bounded).

Given convex $L: \mathbb{R}^n \rightarrow \mathbb{R}$, feature mapping $H: \mathbb{R}^n \rightarrow \mathbb{R}^m$, closed convex $U \subseteq \mathbb{R}^m$, and:

$$\begin{aligned} \epsilon &:= \max_{z \in U} \left\| H(z) - \left[H(w_0) + \partial H(w_0)^T (z - w_0) \right] \right\|_{\infty} \\ &\quad \text{(linearization error)} \\ A &:= \max_{z \in U} \left\| \partial L(H(w_0)) \right\|_2 \quad \text{(Lipschitz)} \\ B &:= \max_{z \in U} \left\| \partial H(w_0) \right\|_{2, \infty} \quad \text{(features Lipschitz)} \end{aligned}$$

Then, choosing $\eta := \frac{\eta_0}{\sqrt{t}}$ & arbitrary $z \in U$,

$$\begin{aligned} \frac{1}{2\eta_0\sqrt{t}} \|w_t - z\|^2 + \min_{S \subseteq U} (L \circ H)(w_t) &\leq \frac{1}{2\eta_0\sqrt{t}} \|w_0 - z\|^2 + \frac{1}{t} \sum_{s=1}^t (L \circ H)(w_s) \\ &\leq \frac{1}{2\eta_0\sqrt{t}} \|w_0 - z\|^2 + \underbrace{(L \circ H)(z)}_R + \frac{A\epsilon}{2\eta_0} + \frac{\eta_0 A^2 B^2}{2\sqrt{t}} \end{aligned}$$

(interp: tuning error $\leq \underbrace{\epsilon}_{x_{opt}}$ tuning error + $\frac{1}{\sqrt{t}}$)

Examples: $F(x; V) = \sum_i a_i \sigma(w_i^T x) = a^T \sigma(Vx)$

$$H(w) = \begin{bmatrix} \sigma_1 F(x; w) \\ \vdots \\ \sigma_n F(x; w) \end{bmatrix}$$

$$\left\| \partial H(w) \right\|_{2, \infty}^2 = \left\| \begin{bmatrix} \sigma_1 \partial F(x; w) \\ \vdots \\ \sigma_n \partial F(x; w) \end{bmatrix} \right\|_{2, \infty}^2$$

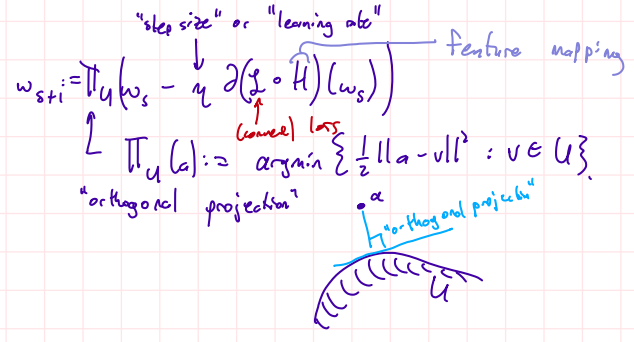
$$\begin{aligned} &= \max_i \left\| \partial F(x; w) \right\|_2^2 \\ &\leq \max_i \sum_j a_j^2 \left\| \sigma'(w_j^T x) x \right\|^2 \leq \max_i \sum_j \frac{1}{m} = 1 \end{aligned}$$

$$\epsilon = \frac{O(B^{4/3} + B \ln(\frac{1}{\eta_0})^{1/3})}{m^{1/6}}; \text{ choose } m = B^3 t^3 \ln(\frac{1}{\eta_0})^2$$

where $B = \sup_{a \in U, b \in U} \|a - b\|$

$$\Rightarrow \epsilon \leq \frac{1}{\sqrt{t}}$$

$$\Rightarrow \min_{S \subseteq U} (L \circ H)(w_s) \leq (L \circ H)(z) + O\left(\frac{1}{\sqrt{t}}\right)$$



Examples:

* $H = \text{identity}$: $\epsilon = 0$,

$$H(w) = w \Rightarrow \partial H(w) = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\left\| \partial H(w) \right\|_{2, \infty} = \left\| \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right\|_{\infty} = 1$$

* Linear: $H(w) = Zw = \begin{bmatrix} -y_1 x_1^T \\ \vdots \\ -y_n x_n^T \end{bmatrix} w$, $\partial H(w) = \begin{bmatrix} y_1 x_1 \\ \vdots \\ y_n x_n \end{bmatrix}$

$$\Rightarrow \left\| \partial H(w) \right\|_{2, \infty} = \max_i \|y_i x_i\| \leq 1$$

* Smooth & Lipschitz activations, e.g., $\sigma(z) = \frac{1}{1 + \exp(-z)}$

By some calculation: $\left\| \partial H(w_s) \right\| \leq 1$,

$$\text{and } \epsilon \leq \frac{\beta^2}{\sqrt{m}} \Rightarrow m = t B^4$$

Loose ends

Proof of lemma. We didn't quite finish; recall MD lemma

$$\|w_k - z\|^2 \leq \sum_{s \in S} \alpha_s \langle g_s, z - w_s \rangle + \sum_{s \in S} \alpha_s^2 \|g_s\|^2 + \|w_0 - z\|^2$$

\uparrow
 $\partial \mathcal{L}(H)(w_s)$

where $\|g_s\|_2 = \|\partial H(w_s) \partial \mathcal{L}(H(w_s))\|_2 = \left\| \sum_{i=1}^n \partial H(w_s)_i; \partial \mathcal{L}(H(w_s))_i \right\|$

$$\leq \sum_{i=1}^n |\partial \mathcal{L}(H(w_s))_i| \cdot \|\partial H(w_s)_i\|_2 \leq \|\partial \mathcal{L}(H(w_s))\|_2 \cdot \max_i \|\partial H(w_s)_i\|_2$$

$\underbrace{\hspace{10em}}_{\| \partial H(w_s) \|_{2, \infty}}$

Loose end #2: projection. When it appears: ^{in MD lemma} $\|w_{s+1} - z\|^2 = \|\Pi_U(w_s - \alpha g) - z\|^2$

$$\leq \|w_s - \alpha g - z\|^2$$
$$= \|w_s - z\|^2 + (\text{cross term } g \text{ goes away})$$

Proposition. If $U \subseteq \mathbb{R}^p$ is closed convex, $\forall a \in \mathbb{R}^d$, $\forall b \in U$, $\|\Pi_U(a) - b\| \leq \|a - b\|$.

Proof. By first-order optimality conditions applied to $\arg \min \{ \frac{1}{2} \|a - v\|^2 : v \in U \}$,

then $\forall b \in U$, $\langle a - \Pi_U(a), b - \Pi_U(a) \rangle \leq 0$

$$\|b - \Pi_U(a)\|^2 = \langle b - \Pi_U(a), b - \Pi_U(a) \rangle \leq \langle b - a, b - \Pi_U(a) \rangle$$
$$\leq \|b - a\| \cdot \|b - \Pi_U(a)\|$$

Refinement to lemma so that $m \ll t$ is possible

Lemma. Same setup as previous convex + Lipschitz lemma, except additionally $\|\partial \mathcal{L}(H(w_s))\| \leq C \cdot \mathcal{L}(H(w_s))$ (true with $C=1$ for logistic).

Then,

$$\frac{1}{2\epsilon t} \|w_\epsilon - z\|^2 + \frac{1}{t} \sum_{i=1}^t \left(1 - \underbrace{C\epsilon}_{\frac{1}{m}} - \underbrace{\frac{m}{2} AB}_{\frac{1}{2}}\right)^{\frac{1}{2}} (\mathcal{L} \circ H)(w_s) \leq \frac{1}{2\epsilon t} \|w_0 - z\|^2 + (\mathcal{L} \circ H)(z).$$

Remark: Suppose $\text{ReLU} \Rightarrow B^2 = 1$. Suppose logistic: $A > 1, C = 1$.

pick $m = 1$, want $\epsilon \leq \frac{1}{4}$; e.g. $m = 2^{24} B^9 \ln(m)^2$.

$$\Rightarrow \frac{2}{\epsilon} \|w_\epsilon - z\|^2 + \frac{1}{t} \sum_{i=1}^t (\mathcal{L} \circ H)(w_s) \leq \frac{2}{\epsilon} \|w_0 - z\|^2 + 4 (\mathcal{L} \circ H)(z).$$

(approach only makes sense if $(\mathcal{L} \circ H)(z) = O(\frac{1}{\epsilon})$;
whereby rhs is $O(\frac{1}{\epsilon})$.)

Proof comment: same proof as before, but use $\|\partial \mathcal{L}(H(w_s))\| \leq C \mathcal{L}(H(w_s))$ & ident terms.

Remark. Point of this page: $\frac{1}{\epsilon}$ rate and width $\ll t$ possible. //

Next time.

Suppose $\mathcal{L} \circ H = \hat{R}$ is β -smooth } (i.e., $\| \nabla \hat{R}(a) - \nabla \hat{R}(b) \| \leq \beta \|a - b\|$).

\Rightarrow GD with $\eta = \frac{1}{\beta}$ satisfies } if $\hat{R} = \mathcal{L} \circ H$ & some sets as before

min s.t $\| \nabla \hat{R}(w) \|^2 \leq 2\beta \left(\frac{f(w_0) - f(w_t)}{t} \right)$ } $\Rightarrow \frac{1}{t}$ rhs rate

interp: come close to local opt

very common ("optimality for nonconvex problems")

