

Lecture 12: OPTIMIZATION

Ann.
* typed notes Sunday

First-order methods only (only use (Clarke) gradients).
 * low per-iteration complexity
 (the high per-iter complexity methods, e.g., Newton, give high accuracy, but that's irrelevant for problems with uncertainty.)
 * Seen to have favorable bias (towards predictors which generalize well).

Rough plan

- * GD & GF near initialization
- * homogeneity ($\sigma(cx) = c\sigma(x)$) & margin maximization & feature locality
- * Maybe (?): mean field
- * Maybe not:
 - * landscape
 - * SGD
 - * SDE
 - * Adam

Precise plan

Near initialization, we'll consider (projected) GD
 eventually we'll pick $S \supset \mathbb{R}^n$ (important)
 $w_0 \in S$, thereafter $w_{t+1} := w_t - \eta g_t$ ← "gradient"
 ↑ step size
 where $g_t = \nabla \hat{R}(w_t)$,
 and $\hat{R} = \mathcal{L} \circ H$ ← $\mathbb{R}^n \rightarrow \mathbb{R}^n$ (training set baked in)
 ↑ convex loss $\mathbb{R}^n \rightarrow \mathbb{R}$

Examples

① linear (rank by rank): $\mathcal{L}(\hat{z}) = \frac{1}{n} \sum_i \ell(\hat{z}_i) \in \mathbb{R}^n$

$$H(w) = \begin{bmatrix} -y_1 x_1^T \\ \vdots \\ -y_n x_n^T \end{bmatrix} w \in \mathbb{R}^{n \times n}$$

← $\ln(\text{kernel}(\hat{z}))$

$$\partial \hat{R} = \partial H \partial \mathcal{L}(H \cdot) = \frac{1}{n} \begin{bmatrix} \ell'(\hat{z}_1) \\ \vdots \\ \ell'(\hat{z}_n) \end{bmatrix}$$

← $\begin{bmatrix} | & | \\ y_1 x_1 & \dots & y_n x_n \\ | & | \end{bmatrix} \in \mathbb{R}^{n \times n}$

② 2-layer network, two inner layers

\mathcal{L} & $\partial \mathcal{L}$ as before

$$H: \mathbb{R}^n \rightarrow \mathbb{R}^n; H(w) = \begin{bmatrix} y_1 F(x_1; w) \\ \vdots \\ y_n F(x_n; w) \end{bmatrix}$$

where $F(x; w) = \sum_j a_j \sigma(v_j^T x)$
 $\uparrow \pm \frac{1}{\sqrt{m}}$

$$\partial H(w) = \begin{bmatrix} | & | \\ \text{vec}(\sum_j a_j \sigma'(v_j^T x) x e_j^T) & \vdots \\ | & | \end{bmatrix} \in \mathbb{R}^{n \times n}$$

← \mathbb{R}^n

More precise plan
 * Gradient methods on $\mathcal{L} \circ H$ near init
 * \mathcal{L} Lipschitz, ∂H "bounded in a certain sense"
 * Smooth \mathcal{L} & H
 * Smooth \mathcal{L} & H , \mathcal{L} strongly convex (usually)
 * weak implicit bias (remove projections)
 * GF
 * Perceptron / Polyak - Łojasiewicz
 Homogeneity...

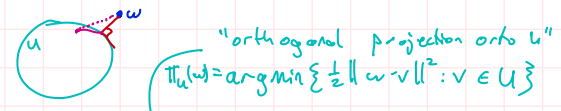
6D $w_{s,t} := w_s - \eta g_s$

"Mirror descent" / "Mirror" lemma

Let convex closed constraint U be given, & comparator $z \in U$. Then

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

where \uparrow is equality if projection never encountered.



Proof. For any z

$$\|w_{s,t} - z\|^2 = \|\Pi_U(w_s - \eta g_s) - z\|^2 \leq \|w_s - \eta g_s - z\|^2$$

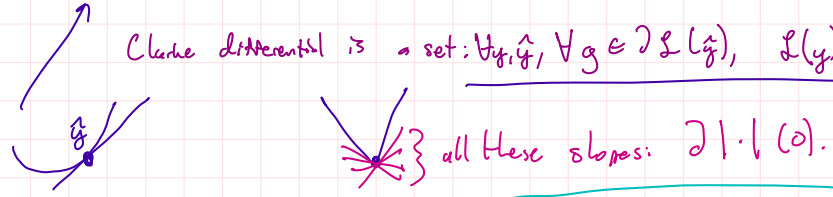
Π_U is contractive; will show it time

$$= \|w_s - z\|^2 + 2\eta \langle g_s, z - w_s \rangle + \eta^2 \|g_s\|^2$$

apply $\sum_{s \in \mathcal{S}}$ to both sides & cancel common term. //

Definition. If L is convex, then $\forall y, \hat{y}$ " $L(y) \geq L(\hat{y}) + \langle \partial L(\hat{y}), y - \hat{y} \rangle$ "

Clarke differential is a set: $\forall y, \hat{y}, \forall g \in \partial L(\hat{y}), L(y) \geq L(\hat{y}) + \langle g, y - \hat{y} \rangle$.



Corollary. U closed convex set, $z \in U$,
 $\varepsilon := \sup_{w \in U} \|H(z) - H(w) - \langle \partial H(w), z - w \rangle\|_{\infty}$

linearization error from near-init opt. lectures

$\sup_{w \in U} \|\partial L(H(w))\|_1 \leq A < \infty$

$\sup_{w \in U} \|\partial H(w)\|_{2, \infty} \leq B < \infty$

$$\frac{\|w_0 - z\|^2}{2\eta_0 \sqrt{\varepsilon}} + \frac{1}{\varepsilon} \sum_{s \in \mathcal{S}} \hat{R}(w_s) \leq \hat{R}(z) + \frac{\|w_0 - z\|^2}{2\eta_0 \sqrt{\varepsilon}} + \frac{A\varepsilon}{\eta_0} + \frac{A^2 B^2}{\eta_0 \sqrt{\varepsilon}}$$

$\hat{R}(w_s) \leq \frac{1}{\sqrt{\varepsilon}}$

Proofs recall from Lemma:

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

$$\langle \partial H(\partial L(H(w_s))), z - w_s \rangle = \langle \partial L(H(w_s)), \partial H^T(z - w_s) \rangle$$

$$= \langle \partial L(H(w_s)), H(z) - H(w_s) \rangle + \langle \partial L(H(w_s)), H(w_s) - H(z) + \partial H^T(z - w_s) \rangle$$

$\leq L(H(z)) - L(H(w_s))$ by convexity

$$\leq \|\partial L(H(w_s))\|_1 \cdot \varepsilon$$

$\|\partial H\|_{2, \infty}$ next time

cannot directly apply convexity

in general can't pick such a z ; more natural is $R(z) \leq \inf_{w \in U} R(w) + \frac{1}{\sqrt{\varepsilon}}$

$\Rightarrow \min_{s \in \mathcal{S}} \hat{R}(w_s) \leq \inf_{w \in U} \hat{R}(w) + O(\frac{1}{\sqrt{\varepsilon}})$

Example.

Recall shallow network, L logistic loss

$$\|\partial L(H(w))\|_1 = \frac{1}{2} \sum_i |l'(H(w)_i)| \leq \frac{1}{2} \sum_i 1 = 1 \leq A$$

$$\|\partial H(w)\|_{2, \infty}^2 = \max_i \|g_i \cdot \partial F(x_i; w)\|_F^2 = \max_i \sum_j \|a_j \sigma'(v_j^T x_i)\|_F^2$$

$$= \max_i \sum_j \sigma_j^2 \|x_i\|^2 \leq 1$$

$$\Rightarrow \frac{1}{\varepsilon} \sum_{s \in \mathcal{S}} \hat{R}(w_s) \leq \hat{R}(z) + \frac{\|w_0 - z\|^2}{\eta_0 \sqrt{\varepsilon}} + \frac{\varepsilon}{\eta_0} + \frac{1}{\eta_0 \sqrt{\varepsilon}}$$

Remark. $\frac{1}{\eta_0} \varepsilon = \frac{1}{\varepsilon} \eta_0 \Rightarrow \eta_0 = \varepsilon^2$; disaster!

later proofs we cover get much smaller widths.