

Lecture 15: smooth & strongly convex optimization

Optimization near initialization so far

| | rates | setting |
|----------------------------|---|--|
| [Lipschitz & convex | $1/\sqrt{\epsilon}$ | ReLU & smooth activation |
| [Smooth | $1/\epsilon$ | "realizable" & $\ k'\ \leq L$ (e.g., logistic) (nonclassical) |
| [Smooth & convex | $\min_{\text{set}} \ \nabla \hat{R}\ ^2 \leq O(\epsilon)$ | $\hat{R} = 2 \circ H$ smooth |
| [smooth & strongly convex | $1/\epsilon$ | $\hat{R} = 1 \circ H$ smooth |
| | $\exp(-O(\epsilon))$ | $\hat{R} = 1 \circ H$ smooth] |

Recall from last time as a consequence of β -smoothness of \hat{R}

$$\|\partial \mathcal{L}(w) - \partial \mathcal{L}(w_0)\| \leq \beta \|w - w_0\|$$

$$\Rightarrow \text{if } \eta \leq \frac{1}{\beta}, \text{ then } \hat{R}(w_{s+1}) = \hat{R}(w_s) - \eta \partial \hat{R}(w_s) \leq \hat{R}(w_s) - \frac{\eta}{2} \|\partial \hat{R}(w_s)\|^2$$

$$\Rightarrow \left\{ \begin{array}{l} \text{recurse } \Rightarrow \min_{z \in \mathcal{S}} \|\partial \hat{R}(z)\|^2 \leq \frac{1}{\epsilon} \sum_{s=0}^{\epsilon-1} \|\partial \hat{R}(w_s)\|^2 \leq \frac{2}{\epsilon \eta} (\hat{R}(w_0) - \hat{R}(w_\epsilon)) \\ \text{recurse k "min lemma"} \Rightarrow \|w_\epsilon - z\|^2 + 2\epsilon \eta \hat{R}(z) \leq \|w_0 - z\|^2 + 2\epsilon \eta \hat{R}(z) \end{array} \right.$$

Remark. today we'll get $\hat{R}(w_t) \leq \hat{R}(w_0) \exp(-\mathcal{O}(t))$
 $\Rightarrow \|w_t - z\|^2 = \|w_0 - z\|^2 \exp(-\mathcal{O}(t))$ "contraction towards global optimum".
 (common setting).

Process:

$$\begin{aligned} \hat{R}(w_{s+1}) &\leq \hat{R}(w_s) - \frac{\eta}{2} \|\partial \hat{R}(w_s)\|^2 = \hat{R}(w_s) - \frac{\eta}{2} \|\partial \mathcal{L}(\mathcal{H}(w_s))\|^2 \\ &\leq \hat{R}(w_s) - \frac{\eta}{2} \|\partial \mathcal{L}(w_s)\|_2^2 \|\partial \mathcal{H}(w_s)\|_2^2 \\ &= \hat{R}(w_s) - \frac{\eta}{2} \underbrace{\sigma_n(\partial \mathcal{H})^2}_{\lambda_{\min}(\partial \mathcal{H}^T \partial \mathcal{H})} \|\partial \mathcal{L}(w_s)\|_2^2 \end{aligned}$$

$\mathbb{R}^{n \times n}$ "kernel gram matrix"

$$(\partial \mathcal{H}^T \partial \mathcal{H})_{ij} = \partial F(x_i; w_s)^T \partial F(x_j; w_s)$$

not w_0 !

$\|\partial \mathcal{L}(w_s)\|_2^2 \leq \lambda (\hat{R}(w_s) - \mathcal{L}(\bar{y}))$
 strong convexity constant for \mathcal{L}

Looking ahead:

$$\hat{R}(w_{s+1}) - \mathcal{L}(\bar{y}) \leq \hat{R}(w_s) - \mathcal{L}(\bar{y}) - \frac{\eta}{2} \sigma_n(\partial \mathcal{H})^2 \lambda (\hat{R}(w_s) - \mathcal{L}(\bar{y}))$$

$\sigma_{\geq s} := \min_{r \geq s} \sigma_n(\partial \mathcal{H}(w_r))$

$$= [\hat{R}(w_s) - \mathcal{L}(\bar{y})] (1 - \eta \lambda \sigma_n(\partial \mathcal{H})^2)$$

$$\leq [\hat{R}(w_s) - \mathcal{L}(\bar{y})] \exp(-\eta \lambda \sigma_{\geq s}^2)$$

$$\Rightarrow \hat{R}(w_t) - \mathcal{L}(\bar{y}) \leq [\hat{R}(w_0) - \mathcal{L}(\bar{y})] \exp(-t \eta \lambda \sigma_{\geq t}^2)$$

near initialization
"classical term"

Remark. ① contraction $w_t \rightarrow \bar{w}$
 inverse condition number
 ① MD lemma / master lemma
 ② apply strong convexity of \mathcal{L} at end
 ③ $t \eta \lambda \leq \frac{t \lambda}{\beta}$ is classical; $\sigma_{\geq t}^2$ is to handle NTK

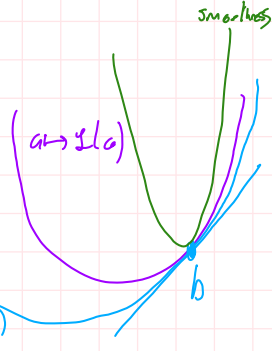
③

Strong convexity

Definition. f is λ -strongly-convex (λ -sc) if $\forall a, b$

$$f(b) \geq f(a) + \langle \partial f(b), b-a \rangle + \frac{\lambda}{2} \|b-a\|^2$$

first-order convexity data



Remark. If f is twice-diff,

$$f \text{ } \lambda\text{-sc} \Leftrightarrow \nabla^2 f \succeq \lambda I.$$

i.e., $Q_b(a) \leq f(a) \forall a,$

$$\& Q_b(b) = f(b)$$

$$\& Q_b(b) + \langle \partial Q_b(b), b-a \rangle = f(b) + \langle \partial f(b), b-a \rangle.$$

Example. ① If f is convex, then $f + \frac{\lambda}{2} \|\cdot\|^2$ is λ -sc.

② $\hat{R}(w) = \frac{1}{2} \|Xw - y\|^2$

is $\sigma_n(X)^2$ -sc.

③ note $\nabla^2 \hat{R} = X^T X$

all equalities

$$\begin{aligned} \text{④ } \frac{1}{2} \|Xa - y\|^2 &= \frac{1}{2} \|X(a-b+b) - y\|^2 \quad \langle X^T X(b-a), a-b \rangle \\ &= \frac{1}{2} \|Xb - y\|^2 + \langle Xb - y, X(a-b) \rangle \\ &\quad + \frac{1}{2} \|X(a-b)\|^2 \\ &\geq \hat{R}(b) + \langle \nabla \hat{R}(b), a-b \rangle + \frac{\sigma_{\min}(X)^2}{2} \|a-b\|^2 \end{aligned}$$

Proposition If f is λ -sc, $\forall w$ $\| \nabla f(w) \|^2 \geq 2\lambda (f(w) - \inf_{v \in \mathbb{R}^n} f(v))$

Proof. Recall $Q_b(a) := f(b) + \langle \partial f(b), a-b \rangle + \frac{\lambda}{2} \|a-b\|^2 \leq f(a)$.

Then $\inf_v f(v) \geq \inf_v Q_w(v)$ optimized at $a = b - \frac{1}{\lambda} \partial f(b)$

$$\begin{aligned} &= f(w) - \langle \partial f(w), \frac{1}{\lambda} \partial f(w) \rangle + \frac{\lambda}{2} \left\| \frac{1}{\lambda} \partial f(w) \right\|^2 \\ &= f(w) - \frac{1}{2\lambda} \|\partial f(w)\|^2 \end{aligned}$$

$\partial f(b) + \lambda(b-a) = 0 \Rightarrow a = b - \frac{1}{\lambda} \partial f(b)$