

Lecture 19: GD descent guarantees; positive homogeneity

Recall:

Clarke differential $\partial f(w) = \text{conv} \left\{ \lim_{i \rightarrow \infty} \nabla f(w_i) : \begin{matrix} w_i \rightarrow w \\ \nabla f(w_i) \text{ exists} \\ \lim_i \nabla f(w_i) \end{matrix} \right\}$

e.g. $\partial f(0) = [-1, +1]$.

- key properties:
- ① always nonempty if f locally Lipschitz (& agrees with subdifferential if f convex)
 - ② flow property: $\dot{w}_t \in -\partial f(w_t)$, \Rightarrow for almost every $t \geq 0$, $\dot{w}_t = -\overline{\partial} f(w_t)$ (minimum norm element) and $\frac{d}{dt} f(w_t) = -\|\overline{\partial} f(w_t)\|_2^2$.
 - ③ Pytorch does not implement the Clarke differential

Descent guarantees

Consider $\dot{w}_t \in -\partial \hat{R}(w_t)$. [Recall $\hat{R} = \mathcal{L} \circ H$]

Proposition (stationary point). $\inf_{s \leq t} \|\overline{\partial} \hat{R}(w_s)\|_2^2 \leq \frac{1}{t} (\hat{R}(w_0) - \hat{R}(w_t))$.

Remark. GD on β -smooth \hat{R} , $\eta \leq \frac{1}{\beta} \Rightarrow \min_{s \leq t} \|\nabla \hat{R}(w_s)\|_2^2 \leq \frac{1}{2t\eta} (\hat{R}(w_0) - \hat{R}(w_t))$. Note $t\eta$ is morally equivalent to t .

Proof. $\hat{R}(w_t) - \hat{R}(w_0) \stackrel{FTC}{=} \int_0^t \frac{d}{ds} \hat{R}(w_s) ds \stackrel{a.e.}{=} \int_0^t -\|\overline{\partial} \hat{R}(w_s)\|_2^2 ds \leq \int_0^t \inf_{j \leq k} \|\overline{\partial} \hat{R}(w_j)\|_2^2 ds = -t \cdot \inf_{j \leq k} \|\overline{\partial} \hat{R}(w_j)\|_2^2$

Rem. $\hat{R} = \mathcal{L} \circ H \Rightarrow \|\overline{\partial} \hat{R}(w_s)\|_2^2 \geq \sigma_{\min}(\partial H(w_0))^2 \cdot \|\partial \mathcal{L}(H(w_s))\|_2^2$. \Rightarrow if $\sigma_{\min}(\partial H(w_0))$ bounded below (by positive constant) then pass near a stationary point of \mathcal{L} .

Prop. Consider $w \mapsto \max\{|w_1| - |w_2|, -100\}$.

Prop. If \hat{R} is λ -sc $\Rightarrow \hat{R}(w_t) - \inf_v \hat{R}(v) \leq (\hat{R}(w_0) - \inf_v \hat{R}(v)) \exp(-2t\lambda)$.

Remark. For β -smooth GD, we had $\exp(-t\eta\lambda)$, $\eta \leq \frac{1}{\beta}$.

Proof. For almost all $s \geq 0$ $\frac{d}{ds} [\hat{R}(w_s) - \inf_v \hat{R}(v)] = -\|\overline{\partial} \hat{R}(w_s)\|_2^2 \leq -2\lambda (\hat{R}(w_s) - \inf_v \hat{R}(v))$. [Recall Grönwall's inequality $\dot{u}_t \leq -\lambda u_t \Rightarrow u_t \leq u_0 \exp(-\int_0^t \lambda ds)$.]
 almost everywhere version

therefore, by Grönwall's inequality $\hat{R}(w_t) - \inf_v \hat{R}(v) \leq (\hat{R}(w_0) - \inf_v \hat{R}(v)) \exp(-2t\lambda)$

Similar curious theorem:

Theorem. Consider $\dot{w}_t \in -\partial \hat{R}(w_t)$, $\dot{v}_t \in -\partial \hat{R}(v_t)$, \hat{R} is λ -sc and w_0 & v_0 arbitrary (i.e., $w_0 \neq v_0$ is possible). Then $\frac{1}{2} \|w_t - v_t\|_2^2 \leq \frac{1}{2} \|w_0 - v_0\|_2^2 \exp(-2t\lambda)$.

Proof. (similar)

Remark. $\hat{R} = \mathcal{L} \circ H$, $\|\overline{\partial} \hat{R}(w_s)\|_2^2 \geq \sigma_{\min}(\partial H(w_s))^2 \cdot \|\partial \mathcal{L}(H(w_s))\|_2^2 \geq \sigma_{\min}(\partial H(w_s))^2 2\lambda (\hat{R}(w_s) - \inf_v \hat{R}(v))$.

Proposition. also in last lecture Given reference point $z \in \mathbb{R}^d$, and convex \hat{R} , $\frac{1}{2t} \|w_t - z\|_2^2 + \hat{R}(w_t) \leq \frac{1}{2t} \|w_0 - z\|_2^2 + \hat{R}(z)$.

Rem. For GD, this guarantee had a constraint set.

Pf. $\frac{1}{2} \|w_t - z\|_2^2 - \frac{1}{2} \|w_0 - z\|_2^2 = \int_0^t \frac{d}{ds} \frac{1}{2} \|w_s - z\|_2^2 ds = \int_0^t \langle \overline{\partial} \hat{R}(w_s), z - w_s \rangle ds \leq \int_0^t [\hat{R}(z) - \hat{R}(w_s)] ds \leq \int_0^t [\hat{R}(z) - \hat{R}(w_t)] ds$

Rem. $\hat{R} = \mathcal{L} \circ H$, then $\langle \overline{\partial} \hat{R}(w_s), z - w_s \rangle = \langle \overline{\partial} \mathcal{L}(H(w_s)), \partial H(w_s)^T (z - w_s) \rangle = \langle \overline{\partial} \mathcal{L}(H(w_s)), H(z) - H(w_s) \rangle$ (controlled via convexity of \mathcal{L})
 $- \langle \overline{\partial} \mathcal{L}(H(w_s)), H(z) - H(w_s) - \partial H(w_s)^T (z - w_s) \rangle$ (controlled via large width).

Positive homogeneity.

Definition. f is L -positive-homogeneous if $\forall c \geq 0, f(cx) = c^L f(x)$.

Remarks. ReLU $\sigma(z) = \max\{0, z\}$ is 1-homogeneous

$$\left(\overset{\forall c \geq 0}{\sigma(cz)} = \max\{0, cz\} = c \cdot \max\{0, z\} = c \sigma(z). \right)$$

ReLU network is L -homogeneous in parameters:

$$\begin{aligned} F(x; cW) &= cW_L \sigma(cW_{L-1} \sigma(\dots \sigma(cW_1 x) \dots)) \\ &= c^L W_L \sigma(W_{L-1} \dots W_1 x) = c^L F(x; W). \end{aligned}$$

Attention layer is not positive homogeneous.

Norms are 1-homogeneous: $\|cx\| = |c| \cdot \|x\|$.

Homogeneous polynomial: all terms have same degree (common degree is L).

Relationship to gradients.

Proposition. Suppose g is L -positive homogeneous & locally Lipschitz.

$\forall w, \forall s \in \partial g(w)$ then $\langle s, w \rangle = L \cdot g(w)$.

Examples. $\langle \partial \frac{1}{2} \|w\|^2, w \rangle = \langle w, w \rangle = \|w\|^2$.

should = $2 \cdot \frac{1}{2} \|w\|^2 = \|w\|^2$

$$\text{ReLU: } \sigma'(r)r = \begin{cases} r < 0 & 0 \cdot r = 0 = \sigma(r) \\ r = 0 & ? \cdot 0 = 0 = \sigma(r) \\ r > 0 & 1 \cdot r = r = \sigma(r) \end{cases} = \sigma(r).$$

ReLU network