

Lecture 18: end of margins & opt

Announcements

- * Any general opt questions, including "why didn't we cover ..." (at any point in lecture).
- Original (convex) opt lectures, some focus on "implicit regularization"
- * E.g., concretely, hw 2 p3.

Last two lectures: minimum norm solution (via margin characterization).

Last time: (linear models) $\gamma(w)$
 * hard margin: $\min_{\|w\|} y \cdot x^T w = \min_{\|w\|} y \cdot x^T \left(\frac{w}{\|w\|} \right)$
 * "smoothed"/"soft margin": $-\frac{\ln \sum \exp(-y \cdot x^T w)}{\|w\|^2}$ $\tilde{\gamma}(w)$
 * GF on $\ln \sum \exp$
 $\gamma(w(t)) \geq \tilde{\gamma}(w(t)) \geq \max_{\|u\|=1} \gamma(u) - O\left(\frac{1}{t}\right)$

This time: (L-homogeneous models)
 $f(x; w)$ L-homogeneous $f(x; cw) = c^L f(x; w) \quad c > 0$
 $\langle w, v \rangle = L f(x; w)$
 $z \in \partial f(x; w)$
 $m_i(w) = y_i f(x_i; w)$

* hard margin $\min_{\|w\|} m_i \left(\frac{w}{\|w\|} \right) = \frac{\min_i m_i(w)}{\|w\|^L}$ (control scaling)
 $\left[\begin{array}{l} \max_{\|u\|=1} \\ \max_{u \in \mathbb{R}^d} \end{array} \min_i m_i(u) \right]$ has maximum under various conditions (e.g., cont)

* "smoothed"/"soft" margin $-\frac{\ln \sum \exp(-m_i(w))}{\|w\|^L}$ $\tilde{\gamma}(w)$
 * CF on $\ln \sum \exp$
 we'll show $t \mapsto \tilde{\gamma}(w(t))$ non decreasing
 $\forall t \geq t_0, \tilde{\gamma}(w(t_0)) > 0$ (Kaifeng Yu, Jian Li) [proof by Zivi]
 very nice paper

Proposition: Suppose $\exists w \tilde{\gamma}(w) > 0$
 Then $\inf_{v \in \mathbb{R}^d} \sum \exp(-m_i(v)) = 0$
 (and it is not attained).

Proof. $0 < \frac{-\ln \sum \exp(-m_i(w))}{\|w\|^L}$
 $\Rightarrow \sum \exp(-m_i(w)) < 1$
 $\Rightarrow \forall i: m_i(w) > 0$. let $\epsilon = \min_i m_i(w) > 0$

$0 = \inf_{v \in \mathbb{R}^d} \sum \exp(-m_i(v)) \leq \liminf_{c > 0} \sum \exp(-m_i(c \cdot w))$
 $= \liminf_{c > 0} \sum \exp(-c^L \cdot m_i(w))$
 $\leq \liminf_{c > 0} \sum \exp\left(-\frac{c}{\epsilon} \cdot \epsilon\right) = 0$

Theorem (Lyu & Li '19, simplification by Zivi & Ji)
 [conditions giving chain rule]
 Suppose $w(t) \in \partial \ln \sum \exp(-m_i(w(t)))$ $\forall t \geq 0$ a.e.
 Suppose m_i L-homogeneous, $\exists t_0, \tilde{\gamma}(w(t_0)) > 0$.
 Then $\forall t \geq t_0, \tilde{\gamma}(w(t)) \geq \tilde{\gamma}(w(t_0))$.

Remark. Assuming all chain rules.

Proof. Note $\forall i: m_i(w(t_0)) > 0$.

Plan: show $\frac{d}{dt} \tilde{\gamma}(w(t)) \geq 0$.

define $u(t) = -\ln \sum \exp(-m_i(w(t)))$
 $v(t) = \|w(t)\|^L$

Note $\frac{d}{dt} \tilde{\gamma}(w(t)) = \frac{v(t)u'(t) - u(t)v'(t)}{v(t)^2}$.

Crucial fact: $\langle w(t), w'(t) \rangle = -\partial \ln \sum \exp(-m_i(w))$
 $= \left\langle w(t), \frac{\sum_i \exp(-m_i(w)) \partial m_i(w)}{\sum \exp(-m_i(w))} \right\rangle$ (chain rule)
 $= \sum_i \frac{\exp(-m_i(w))}{\sum_j \exp(-m_j(w))} \langle w(t), \partial m_i(w) \rangle$
 $= \sum_i \frac{\exp(-m_i(w))}{\sum_j \exp(-m_j(w))} L_i m_i(w) = -\ln \sum \exp(-m_i(w))$
 $\geq -\ln \sum \exp(-m_i(w))$

$\geq \left(\sum_i \frac{\exp(-m_i(w))}{\sum_j \exp(-m_j(w))} \right) u(t) \cdot L$
 $= L \cdot u(t)$. (lower bounded by what we'd get if $u(t)$ were homogeneous.)

Can use this to upper & lower bound $u(t), v(t)$
 & get $u(t)v'(t) - u'(t)v(t)$ "mechanically"

loss of global optimality:
 $u(t) = \langle \partial \ln \sum \exp(-m_i(w)), -\partial \ln \sum \exp(-m_i(w)) \rangle$
 $= \|u(t)\|^2$
 linear case: $\|u(t)\|^2 = \|u(t)\| \cdot \sup_{\|u\|=1} \langle u(t), u \rangle$
 homogeneous case: $\geq \|u(t)\| \cdot \langle u(t), u \rangle$
 $\geq \|u(t)\| \cdot L \cdot u(t) \cdot \frac{1}{\|u(t)\|} \geq \|u(t)\| \cdot \gamma$ (uses properties of linear hardmargin)

Remark (omitted topics).

* Escaping saddle points (typically by adding Gaussian noise to SGD) by Kong Ge & friends.

* Mean field optimization analysis (Song Mei, Andrea Montanari, & friends)
 Personally not sure of precise relationship to NTK. Answer is much different.

* Landscape analysis ("mode connectivity" by Ge et al; "all local optima are global" Tengyu Ma & friends)

