

Lecture 17: Nonsmooth & margins.

Announcements

* HW2 tomorrow (no late homework!)

Recap: nonsmoothness:

① "Clarke differential"

$$\partial f(w) := \text{conv} \left(\left\{ \lim_{w_s \rightarrow w} \nabla f(w_s) : \begin{array}{l} \nabla f(w_s) \text{ exists} \\ \lim \nabla f(w_s) \text{ exists} \end{array} \right\} \right)$$

key properties (A) exists everywhere if f locally lip
(B) has a reasonable analog to ∇f .
 caveat: I don't know it re-l.

② Positive homogeneity:

f is L -positive-homogeneous
 $\forall c \geq 0 \quad f(c \cdot x) = c^L f(x)$.

key properties: (A) reasonable abstraction of ReLU
(B) Gradient interaction:
 $\forall v \in \partial f(w) : \langle v, w \rangle = L \cdot f(w)$.

Plan for this week:

- Ⓘ one more Clarke consequence
 - Ⓙ "strong" implicit regularization property of GD.
 - ⓪ linear case
 - ⓫ multilayer L -homogeneous (with Clarke flow). ("beyond NTK").
 - Ⓚ Optimization topics I neglected.
-

Proposition ("norm preservation", Du-Hu-Lee, '18)
Proof "Suggested" Zwei.

Suppose $w(t)$ given by Clarke flow of

$$\hat{R}(w) = \frac{1}{n} \sum_k \mathcal{L}(y_n f(x_k; w)),$$

where f is 1-homogeneous in each layer,
and suppose the following chain rule:

$$\forall u \in \partial \hat{R}(w) \Rightarrow u = \frac{1}{n} \sum_k \underbrace{\mathcal{L}'(y_n f(x_k; w))}_{\mathcal{L}'_k(w)} y_n \uparrow u_k \in \partial f(x_k; w)$$

Then $\forall t \geq 0, \forall i, j,$

$$\frac{1}{2} \|W_i(t)\|^2 - \frac{1}{2} \|W_i(0)\|^2 = \frac{1}{2} \|W_j(t)\|^2 - \frac{1}{2} \|W_j(0)\|^2$$

Proof.

$$\begin{aligned} & \frac{1}{2} \|W_i(t)\|^2 - \frac{1}{2} \|W_i(0)\|^2 \\ &= \int_0^t \frac{d}{ds} \frac{1}{2} \|W_i(s)\|^2 ds = \int_0^t \langle W_i(s), \dot{W}_i(s) \rangle ds \end{aligned}$$

$$= \int_0^t \left\langle W_i(s), \frac{1}{n} \sum_k \mathcal{L}'_k(w(s)) u_{k,i}(s) \right\rangle ds$$

$$= \int_0^t \frac{1}{n} \sum_k \mathcal{L}'_k(w(s)) \langle W_i(s), u_{k,i}(s) \rangle ds$$

$$= \int_0^t \frac{1}{n} \sum_k \mathcal{L}'_k(w(s)) f(x_k; w(s)) ds,$$

which is independent of i .

"Strong" implicit regularization.

Now we'll show GD \rightarrow minimum norm solutions (not just small norm).

Our settings

* exponential or logistic loss
 $\uparrow \ln(1 + \exp(-z))$

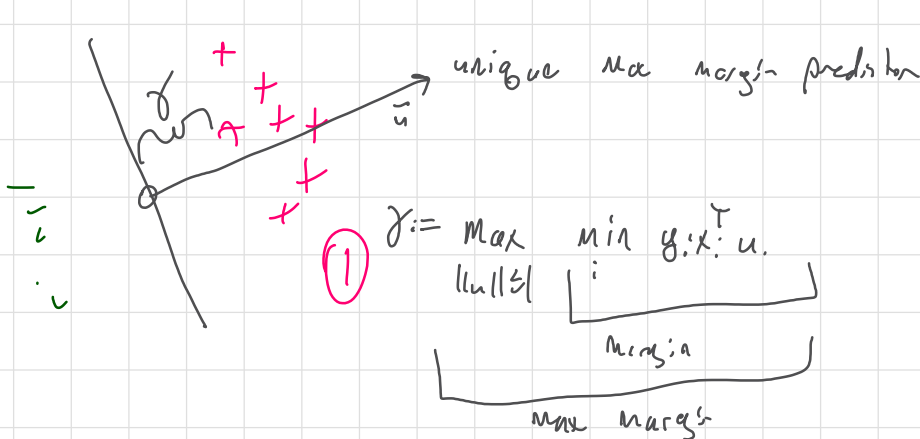
* GF/CF.

Other settings

* Squared loss, GF/GD.
 (Some authors: Gunasekar, Srebro, Cohen, Arora;)

Today's plan: linear case.

Here, minimum norm solution = "maximum margin solution".



② "SVM form": $\min \frac{1}{2} \|v\|^2$ s.t. $y_i x_i^T v \geq 1 \forall i$

u is ① $\Leftrightarrow \frac{u}{\gamma}$ is optimal for ②

Connection to exp loss:

$$\min_i y_i x_i^T u = -\max_i -y_i x_i^T u \approx -\ln \sum \exp(-y_i x_i^T u)$$

empirical risk

In fact:

$$\gamma = \max_{\|u\| \leq 1} \min_i y_i x_i^T u = \max_{r \geq 0} \max_{\|u\| \leq r} \min_i -r \ln \sum \exp\left(\frac{-y_i x_i^T u}{r}\right)$$

Theorem (Telgarsky '13; Gunasekar et al. '18).

Let $w(t)$ be given by $\dot{w}(s) = \nabla_w \ln \sum \exp(-z_i^T w)$, with $w(0) = 0$. Then

$$\gamma \geq \frac{\min_i z_i^T w(t)}{\|w(t)\|} \geq \frac{-\ln \sum \exp(-z_i^T w(t))}{\|w(t)\|} \stackrel{\text{①}}{\geq} \tilde{\gamma}(t) \geq \gamma - \frac{\ln n}{t\gamma - (\ln n)/\gamma}$$

Remark.

Note $\nabla \ln \sum \exp(v) = \frac{\nabla \sum \exp(v)}{\sum \exp(v)}$

i.e., both point in same direction.

\Rightarrow Can show GF on both gives same path. ("time rescaling")

Proof. $\min_i z_i^T w(t) = -\max_i -z_i^T w(t) \geq -\ln \sum \exp(-z_i^T w(t))$

$$\tilde{\gamma}(t) = \frac{u(t) - u(0)t\gamma}{\|w(t)\|} = \frac{u(0)}{\|w(t)\|} + \frac{\int_0^t \dot{u}(s) ds}{\|w(t)\|} \geq \gamma$$

lose this step to \leftarrow homogeneous

Note: $\|\dot{w}(s)\| \geq \langle \dot{w}(s), u \rangle$

$$= \langle -\nabla \ln \sum \exp(-z^T w(s)), u \rangle = \frac{\sum \exp(-z^T w(s)) z^T u}{\sum \exp(-z^T w(s))} \geq \gamma \frac{\sum \exp(-z^T w(s))}{\sum \exp(-z^T w(s))} = \gamma$$

$$\Rightarrow \int_0^t \dot{u}(s) ds = \int_0^t \langle \nabla \ln \sum \exp(-), \dot{w}(s) \rangle = \int_0^t \|\dot{w}(s)\|^2 ds \geq \gamma \int_0^t \|\dot{w}(s)\| ds \geq \gamma \left\| \int_0^t \dot{w}(s) ds \right\| = \gamma \|w(t) - w(0)\| = \gamma \|w(t)\|$$

$$\|w(t)\| \gamma \geq \|w(t)\| \tilde{\gamma}(t) = -\ln \sum \exp(-z^T w(s)) = u(0) + \int_0^t \dot{u}(s) ds = u(0) + \int_0^t \|\dot{w}(s)\|^2 ds \geq \frac{u(0)}{-\ln n} + t\gamma^2$$