

Lecture 21: minimum norm solutions

Plan for rest of semester.

- 21 GD prefers low/minimum norm solutions
 - 22 project
 - 23 concentration
 - 24 deep network generalization
 - 25
 - 26
 - 27 double descent
 - 28 ????
- ~~Kolmogorov Arnold~~

Goal of course: we can get low test error.

$$R(\hat{f}) - R(g) = \underbrace{R(\hat{f}) - \hat{R}(\hat{f})}_{\text{error}} + \underbrace{\hat{R}(\hat{f}) - \hat{R}(f)}_{\text{opt}} + \underbrace{\hat{R}(f) - R(f)}_{\text{gap}} + \underbrace{R(f) - R(g)}_{\text{generalization}}$$

$\hat{f} \in \mathcal{F}$ a/g
 $f \in \mathcal{F}$ not inside
 g not outside

"rethinking generalization: ..." Zhang et al.
 "..." Tomasek - Neyshabur - Srebro

One solution to issue:

* \mathcal{F} is adapted to data.

(our approach: \mathcal{F} is norm bounded, & the norm we need depends on data).

In optimization:

① Near-initialization ("moderately convex" / "neural tangent kernel" ²⁰⁰⁹)

Adaptivity to low norm: $\frac{1}{2\epsilon} \|w_{\epsilon} - z\|^2 + \hat{R}(w_{\epsilon}) \leq \frac{1}{2} \|w_0 - z\|^2 + \hat{R}(z)$

② Far from initialization

* norm preservation

* next two lectures

Low-norm predictors

Seen it before: $\frac{1}{2} \|w_k - z\|^2 + \hat{R}(w_k) \leq \frac{1}{2} \|w_0 - z\|^2 + \hat{R}(z)$ (convex + smooth, add step for full version near init)

- Two topics we've missed:
 - * convergence to the minimum norm solution
 - * Far from initialization (& is that practical)

Recall appeared in linear regression:

$$\lim_{\lambda \downarrow 0} \operatorname{argmin}_w \frac{1}{2} \|Xw - y\|^2 + \frac{\lambda}{2} \|w\|^2 = \lim_{\lambda \downarrow 0} (X^T X + \lambda I)^{-1} X^T y = X^+ y$$

} ordinary least squares

Also GD on $\frac{1}{2} \|Xw - y\|_2^2 \xrightarrow[\text{iterations} \rightarrow \infty]{\text{step size to 0}}$ $X^+ y$ $X^T (Xw - y)$

Good news: GD gets min norm solution.

Question: (a) how about classification losses, (b) how about nonlinear

Remark (how to prove min norm)
 Consider opt cond
 ∂f
 min $\frac{1}{2} \|w\|^2$
 s.t. $X^T X w = X^T y$

Realizability/separability assumptions

Suppose first our model perfectly labels the data

Linear Regression version:
 min $\frac{1}{2} \|w\|_2^2$
 s.t. $Xw = y$
 & assume feasible

Classification version:
 min $\frac{1}{2} \|w\|_2^2$
 s.t. $y_i \odot (Xw)_i \geq 1$ [i.e., $\forall i, y_i x_i^T w \geq 1$]
 & assume feasible

Remark: why? without this assumption, norms seem to blow up.

non linear version:
 min $\frac{1}{2} \|w\|^2$
 s.t. $y_i f(x_i; w) \geq 1 \forall i$.

this "margin" is maximized

We'll consider two cases: ① Linear case + exponential loss (similar to logistic/cross-entropy) ② L-homogeneous.

Linear case

Theorem (A). Consider GF on $\mathcal{L}(w) = \sum_i \exp(-y_i x_i^T w)$ [$w_0 = 0, w_t = -\nabla \mathcal{L}(w_t)$].

& suppose $\exists u, \|u\|_2 = 1, \exists \gamma > 0$, s.t. $\forall i, y_i x_i^T u \geq \gamma$ (separability assumption).

need $\|w\| \leq$ Then $\frac{\min_i y_i x_i^T w_t}{\|w_t\|} = \min_i y_i x_i^T \left(\frac{w_t}{\|w_t\|}\right) \geq \gamma - O\left(\frac{\ln n}{\ln t}\right)$

Remark: $\frac{1}{\ln t}$ is slow; can be "accelerated", but in practice slow $\ln t$ seems red.

Lemma. Consider setting of (A). Then $\mathcal{L}(w_t) \leq \frac{n}{t} + \frac{(\ln t)^2}{2t\gamma^2}$, $\|w_t\| \geq \ln(t\gamma^2) - \ln\left(t + \frac{\ln t}{n}\right)$.

Remark. Note solutions unbounded & $\|w_t\| \rightarrow \infty$

Proof. Recall $\forall z \in \mathbb{R}^d, \frac{1}{2t} \|w_t - z\|^2 + \mathcal{L}(w_t) \leq \frac{1}{2t} \|z\|^2 + \mathcal{L}(z)$.

Now cleverly guess $z = u \frac{\ln t}{\gamma}$ $y_i x_i^T u \frac{\ln t}{\gamma} \geq \gamma \cdot \frac{\ln t}{\gamma} = \ln t$

$$\Rightarrow \mathcal{L}(w_t) \leq \frac{1}{2t} \frac{(\ln t)^2}{\gamma^2} + \sum \exp(-y_i x_i^T z) \leq \frac{(\ln t)^2}{2t\gamma^2} + \frac{n}{t}$$

(Remark: if we included dropped $\frac{1}{2t} \|w_t - z\|^2$, then can show $\|w_t\| = O\left(\frac{\ln t}{\gamma}\right)$.)

for norm lower bound: $\|w_t\| = -\ln \exp(-\|w_t\|) \geq -\ln \sum \exp(-\|w_t\|) \stackrel{\text{uses } \|x_i\| \leq 1}{\geq} -\ln \sum_i \exp(-y_i x_i^T w_t) \geq -\ln \left[\frac{(\ln t)^2}{2t\gamma^2} + \frac{n}{t} \right]$

Proof of (A).

$$\min_i y_i x_i^T w_t = -\ln \min_i \exp(y_i x_i^T w_t) \geq -\ln \sum_i \exp(-y_i x_i^T w_t)$$

$$= -\ln \sum \exp(-y_i x_i^T w_t) + \ln \sum \exp(-y_i x_i^T w_0) - \ln \sum \exp(-y_i x_i^T w_0)$$

$$\stackrel{\text{FTC}}{=} \int_0^t \frac{d}{ds} -\ln \sum \exp(-y_i x_i^T w_s) ds = \int_0^t \left\langle \frac{-\nabla \mathcal{L}(w_s)}{\sum \exp(-y_i x_i^T w_s)}, w_s \right\rangle ds \geq \gamma$$

$$= \int_0^t \frac{\|\nabla \mathcal{L}(w_s)\|^2}{\sum \exp(-y_i x_i^T w_s)} ds \geq \int_0^t \|\nabla \mathcal{L}(w_s)\| \cdot \frac{\sum \exp(-y_i x_i^T w_s) \langle y_i x_i, u \rangle}{\sum \exp(-y_i x_i^T w_s)} ds$$

$$\geq \gamma \int_0^t \|\nabla \mathcal{L}(w_s)\| ds \geq \gamma \left\| \int_0^t \nabla \mathcal{L}(w_s) ds \right\| = \gamma \|w_t - w_0\| = \gamma \|w_t\|$$

$$\Rightarrow \frac{\min_i y_i x_i^T w_t}{\|w_t\|} \geq \gamma \left(\frac{\|w_t\|}{\|w_t\|} \right) - O\left(\frac{\ln n}{\ln t}\right) \stackrel{\text{lemma}}{\geq} \gamma - O\left(\frac{\ln n}{\ln t}\right)$$