

Lecture 24

Last time

- ▷ Controlling the VC dimension
- ▷ Generalization for linear models.

Today

- ▷ Linear Models continued
- ▷ Convexity & Stability

Linear models continued

Lemma: Consider $X = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$, and $\mathcal{H} = \{a \mapsto \langle w, a \rangle \mid \|w\|_2 \leq 1\}$.

Then,

$$\bar{\mathcal{R}}(\mathcal{H}(X)) \leq \sqrt{\sum_{i=1}^n \|x_i\|_2^2},$$

thus,

$$\bar{\mathcal{R}}_n(\mathcal{H}) \leq \frac{\mathbb{E} \max_i \|x_i\|_2}{\sqrt{n}}. \quad \text{Dimension free.}$$

Proof: Expanding

$$\begin{aligned}\bar{\mathcal{R}}(\mathcal{H}(X)) &= \mathbb{E}_\varepsilon \sup_{\|w\|_2 \leq 1} \sum_{i=1}^n \varepsilon_i \langle w, x_i \rangle \\ &= \mathbb{E}_\varepsilon \sup_{\|w\|_2 \leq 1} \langle w, \sum \varepsilon_i x_i \rangle \\ &= \mathbb{E}_\varepsilon \|\sum \varepsilon_i x_i\|_2\end{aligned}$$

$$\text{Jensen's} \rightarrow \leq \sqrt{\mathbb{E}_{\varepsilon} \|\sum \varepsilon_i x_i\|_2^2}$$

$\varepsilon_i \text{ iid}$ $\rightarrow = \sqrt{\sum \|x_i\|_2^2}.$

Recall

$$\bar{R}_n(\mathcal{H}) = \mathbb{E}_{\mathbf{x}} \bar{R}\left(\frac{1}{n} \mathcal{H}(\mathbf{x})\right)$$

$$\leq \frac{1}{n} \mathbb{E} \sqrt{\sum \|x_i\|_2^2}$$

$\text{Jensen's} \rightarrow \leq \frac{1}{n} \mathbb{E}_{\mathbf{x}} \max \|x_i\|_2.$

□

Lastly, let's do the same with the ℓ_1 norm.

Lemma: Consider $X = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$, and

$\mathcal{H} = \{a \mapsto \langle w, a \rangle \mid \|w\|_1 \leq 1\}.$
Then,

$$\bar{R}(\mathcal{H}(X)) \leq \sqrt{2n \log(2d)} \max \|x_i\|_{\infty},$$

thus,

$$\bar{R}_n(\mathcal{H}) \leq \sqrt{\frac{2 \log 2d}{n}} \max_{\|x_i\|_\infty} \|x_i\|_\infty$$

Almost dim
free.

Proof: Same strategy as before

$$\begin{aligned}\bar{R}(\mathcal{H}(x)) &= \mathbb{E}_{\varepsilon} \sup_{\|w\|_1 \leq 1} \langle w, \sum \varepsilon_i x_i \rangle \\ &\stackrel{(why?)}{\downarrow} = \mathbb{E}_{\varepsilon} \|\sum \varepsilon_i x_i\|_\infty \\ &= \mathbb{E}_{\varepsilon} \sup_{j \in [d]} \left| \sum \varepsilon_i x_{ij} \right| \\ &\stackrel{\text{Massart's lemma.}}{\leq} \sqrt{2 \log(2d)} \sup_{x_i} \sqrt{\sum_j x_{ij}^2} \\ &\leq \sqrt{2n \log(2d)} \max_{\|x_i\|_\infty} \|x_i\|_\infty.\end{aligned}$$

□

Generalization via convexity and stability

Suppose we wanted to minimize

$$(*) \quad \min_{\theta \in \Theta} F(\theta) \text{ with } F(\theta) = \mathbb{E}_{x \sim P} f(\theta, x).$$

and we had an algorithm A s.t. given an iid sample $X = (x_1, \dots, x_n)$ it outputs $\hat{\theta} = A(X)$ an estimate of a minimizer of (\star) . Our goal next is to establish generalization guarantees for $A(X)$ (not necessarily uniformly).

The main idea is to use stability. Define

$$X^i = (x_1, \dots, x_{i-1}, \underset{\uparrow}{x'}, x_{i+1}, \dots, x_n)$$

where $x' \sim P$ and independent of X .

Theorem: We have that

$$\begin{aligned} \mathbb{E}_X [F(A(X))] - \frac{1}{n} \sum_{i=1}^n f(A(X), x_i) \\ = \mathbb{E}_{\substack{(X, x') \sim \text{iid } P \\ i \sim \text{Unif}(\mathbb{N})}} [f(A(X^i), x_i) - f(A(X), x_i)]. \end{aligned}$$

Proof: The proof is simple, for every i we have

$$\mathbb{E}_X F(A(X)) = \mathbb{E}_{X, x'} [f(A(X), x')] = \mathbb{E}_{X, x'} [f(A(X^i), x_i)].$$

Also

$$\mathbb{E}_X \left[\frac{1}{n} \sum_i f(A(x), x_i) \right] = \mathbb{E}_{x,i} f(A(x), x_i).$$

Substituting these two gives (10). \square

Intuition: If the value of the loss is stable to changing one element of the training set, then the method generalize well. \rightarrow

Def: We say that an algorithm $A(\cdot)$ is **replace-one-stable** with rate $\epsilon(n)$ if

$$\mathbb{E}_{\substack{x, x' \sim P \\ i \sim \text{Unif}[m]}} [f(A(x^i), x_i) - f(A(x), x_i)] \leq \epsilon(n). \quad \rightarrow$$

Assumption (2): We will focus on the case $f(\cdot, x)$ δ -Lipschitz convex for all x and \mathbb{H} is a convex closed set. \rightarrow

Just as we did for ridge regression we consider a regularized estimate

for

$$(A) A(\theta) := \underset{\theta \in \Theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n f(\theta, x_i) + \frac{\lambda}{2} \|\theta\|_2^2$$

In turn, this objective exhibits nice properties.

Def: A function $g: \mathbb{R}^d \rightarrow \mathbb{R} \cup \infty$ is α -strongly convex if $g - \frac{\alpha}{2} \|\cdot\|^2$ is convex.

The objective in (A) is λ -strongly convex by construction.

Lemma (*) Proven in Nonlinear I: Suppose g is α -strongly convex. Then it has a unique minimizer θ^* and

$$\frac{\alpha}{2} \|\theta - \theta^*\|^2 \leq g(\theta) - g(\theta^*) \quad \forall \theta.$$

Thus, $A(\cdot)$ is well defined. Next, we establish a generalization guarantee.

Theorem: Suppose Assumption (*) holds. Then, algorithm (A) is

replace -one- stable with rate $\frac{2\rho^2}{\lambda n}$.
 Thus,

$$\mathbb{E}_X [F(A(X)) - \frac{1}{n} \sum_{i=1}^n f(A(X), x_i)] \leq \frac{2\rho^2}{\lambda n}.$$

Proof: Notice that by ρ -Lipschitzness

$$f(A(X^i), x_i) - f(A(X), x_i) \leq \rho \|A(X^i) - A(X)\|_2.$$

So it suffices to show that

$$\|A(X^i) - A(X)\| \leq \frac{2\rho}{\lambda n}.$$

Define the function

$$f_X(\theta) := \frac{1}{n} \sum_{i=1}^n f(\theta, x_i) + \frac{\lambda}{2} \|\theta\|_2^2.$$

By Lemma (V) we have that for each i :

$$\begin{aligned} & \frac{\lambda}{2} \|A(X^i) - A(X)\|_2^2 \\ & \leq f_X(A(X^i)) - f_X(A(X)) \\ & = \frac{1}{n} \sum_{x \in X} f(A(X^i), x) + \frac{\lambda}{2} \|A(X^i)\|_2^2 \\ & \quad - \frac{1}{n} \sum_{x \in X} f(A(X), x) - \frac{\lambda}{2} \|A(X)\|_2^2 \\ & = \frac{1}{n} \sum_{x \in X} f(A(X^i), x) + \frac{\lambda}{2} \|A(X^i)\|_2^2. \end{aligned}$$

Notice the index set

$$- \frac{1}{n} \sum_{x \in X} f(A(x), x) - \frac{\lambda}{2} \|A(x)\|_2^2.$$

$$+ \frac{f(A(x^i), x_i) - f(A(x), x_i)}{n}$$

$$+ \frac{f(A(x), x') - f(A(x^i), x')}{n}$$

$$\leq f_{x^i}(A(x)) - f_{x^i}(A(x^i)) \\ + \frac{2\rho}{n} \|A(x) - A(x^i)\|_2$$

$$\leq -\frac{\lambda}{2} \|A(x^i) - A(x)\|_2^2 + \frac{2\rho}{n} \|A(x) - A(x^i)\|_2$$

$A(x^i)$ minimizes f_{x^i} so this uses Lemma (v).

Rearranging gives

$$\|A(x^i) - A(x)\|_2^2 \leq \frac{2\rho}{\lambda n},$$

as we wanted. □

We can use this to prove an excess risk bound.

Corollary: Let θ^* be any minimizer of

(★). Then, we have

$$\mathbb{E}_x F(A(x)) \leq \min_{\theta \in \Theta} F(\theta) + \frac{\lambda}{2} \|\theta^*\|_2^2 + \frac{2\rho}{\lambda n}.$$

Thus, if we set $\lambda = \sqrt{\frac{4\rho^2}{n\|\theta^*\|^2}}$,

$$\mathbb{E}_x F(A(x)) \leq \min_{\theta \in \Theta} F(\theta) + \frac{2\rho\|\theta^*\|}{\sqrt{n}}.$$

Proof: Adding and subtracting

$$\mathbb{E}_x F(A(x))$$

$$= \mathbb{E}_x \left[\frac{1}{n} \sum_{x \in X} f(A(x), x) \right]$$

$$+ \mathbb{E}_x \left[F(A(s)) - \frac{1}{n} \sum_{x \in X} f(A(x), x) \right]$$

$$\leq \mathbb{E}_x \left[\frac{1}{n} \sum_{x \in X} f(A(x), x) \right] + \frac{2\rho^2}{\lambda n}.$$

Notice that by definition

$$\mathbb{E}_x \left[\frac{1}{n} \sum_{x \in X} f(A(x), x) \right] \leq \mathbb{E}_x [f_x(A(x))]$$

$$\begin{aligned} &\leq \mathbb{E}_{\mathbf{x}} [f_{\mathbf{x}}(\boldsymbol{\theta}^*)] \\ &= F(\boldsymbol{\theta}^*) + \frac{\lambda}{2} \|\boldsymbol{\theta}^*\|_2^2. \end{aligned}$$

Combining these two bounds yield the first stated bound. The second one follows by optimizing λ . \square