

Lecture 3 VC Theory for Generalization Error

He Li Machine Learning Course

2017/10/10

1. Simple Classifiers

$(x_i, y_i)_{i=1}^n$, $x_i \in \mathbb{R}^d$, $x = (x^{(1)}, \dots, x^{(d)})$ instance, $y \in \{\pm 1\}$ label

- Decision Tree: hypothesis space $\mathcal{F} = \{\text{all decision tree, depth} \leq \alpha\}$
- Linear Classifier: hypothesis space $\mathcal{F} = \{(\omega, b) : \omega \in \mathbb{R}^d, b \in \mathbb{R}, \|\omega\| = 1\}$, then classifier is:

$$f(x) = \text{sgn}(\omega^T x + b)$$

Empirical Risk Minimization

Given \mathcal{F} , training data $(x_i, y_i)_{i=1}^n$. **Empirical Risk Minimization(ERM)** is an algorithm, which finds $f \in \mathcal{F}$, s.t.

$$\hat{f} = \underset{f \in \mathcal{F}}{\text{argmin}} \frac{1}{n} \sum_{i=1}^n I[y_i \neq f(x_i)]$$

- Empirical Error: \hat{f} 's performance on training data. Denote

$$\mathbb{P}_S(y_i \neq f(x_i)) = \frac{1}{n} \sum_{i=1}^n I[y_i \neq \hat{f}(x_i)]$$

where $S = (x_i, y_i)_{i=1}^n$

- Generalization Error: \hat{f} 's performance on test data. Denote

$$\mathbb{P}_D(y_i \neq f(x_i)) = \mathbb{E}_{(x,y) \sim D} \{I(y \neq \hat{f}(x))\}$$

A small empirical error cannot indicate a small generalization error (cannot use Chernoff bound), since on training data, $z_i = I[y_i \neq \hat{f}(x_i)]$ are not independent.

Finite hypothesis space

Consider \mathcal{F} is finite, $|\mathcal{F}| < \infty$. ERM learns $\hat{f} \in \mathcal{F}$, then

$$P \{\mathbb{P}_D(y_i \neq f(x_i)) - \mathbb{P}_S(y_i \neq f(x_i)) \geq \epsilon\} \leq |\mathcal{F}| e^{-2n\epsilon^2}$$

called **union bound**. If we fix $f \in \mathcal{F}$, $P \{\mathbb{P}_D(y_i \neq f(x_i)) - \mathbb{P}_S(y_i \neq f(x_i)) \geq \epsilon\}$ satisfies Chernoff bound.

$$P \{\mathbb{P}_D(y_i \neq f(x_i)) - \mathbb{P}_S(y_i \neq f(x_i)) \geq \epsilon\} \leq e^{-2n\epsilon^2}$$

Then, by union bound

$$P \{\exists f \in \mathcal{F}, \mathbb{P}_D(y_i \neq f(x_i)) - \mathbb{P}_S(y_i \neq f(x_i)) \geq \epsilon\} \leq |\mathcal{F}| e^{-2n\epsilon^2}$$

Then we have our conclusion.

2. VC bound

Uniform Law of Large Numbers

Denote $z_i = (x_i, y_i)$, $\phi_f(z_i) = I(y_i \neq f(x_i))$. Consider:

$$\mathbb{P} \left(\sup_{\phi \in \Phi} \left| \frac{1}{n} \sum_i \phi(z_i) - \mathbb{E}[\phi(z)] \right| \geq \epsilon \right)$$

Step I (Double Sample Trick)

Proposition: $X_1, \dots, X_n, X_{n+1}, \dots, X_{2n}$ are i.i.d Bernoulli R.V.. $\mathbb{E}(X) = p$. Denote $\nu_1 = \frac{1}{n} \sum_{i=1}^n X_i$, $\nu_2 = \frac{1}{n} \sum_{i=n+1}^{2n} X_i$. If $n \geq \frac{\ln 2}{\epsilon^2}$, $\epsilon > 0$, then

$$\frac{1}{2} \mathbb{P}(|\nu_1 - p| \geq 2\epsilon) \leq \mathbb{P}(|\nu_1 - \nu_2| \geq \epsilon) \leq 2\mathbb{P}\left(|\nu_1 - p| \geq \frac{\epsilon}{2}\right)$$

proof

right inequality:

$$|\nu_1 - \nu_2| \geq \epsilon \implies |\nu_1 - p| \geq \frac{\epsilon}{2} \text{ or } |\nu_2 - p| \geq \frac{\epsilon}{2}$$

Therefore,

$$\begin{aligned} \{|\nu_1 - \nu_2| \geq \epsilon\} &\subset \{|\nu_1 - p| \geq \frac{\epsilon}{2}\} \cup \{|\nu_2 - p| \geq \frac{\epsilon}{2}\} \\ \mathbb{P}(|\nu_1 - \nu_2| \geq \epsilon) &\leq \mathbb{P}\left(|\nu_1 - p| \geq \frac{\epsilon}{2} \cup |\nu_2 - p| \geq \frac{\epsilon}{2}\right) \\ &\leq \mathbb{P}\left(|\nu_1 - p| \geq \frac{\epsilon}{2}\right) + \mathbb{P}\left(|\nu_2 - p| \geq \frac{\epsilon}{2}\right) \\ &= 2\mathbb{P}\left(|\nu_1 - p| \geq \frac{\epsilon}{2}\right) \end{aligned}$$

Similarly, left inequality:

$$|\nu_1 - p| \geq 2\epsilon, |\nu_2 - p| \leq \epsilon \implies |\nu_1 - \nu_2| \geq \epsilon$$

□

Lemma (Homework)

$$\begin{aligned} &\frac{1}{2} \mathbb{P} \left(\sup_{\phi \in \Phi} \left| \frac{1}{n} \sum_{i=1}^n \phi(z_i) - \mathbb{E}[\phi(z)] \right| \geq 2\epsilon \right) \\ &\leq \mathbb{P} \left(\sup_{\phi \in \Phi} \left| \frac{1}{n} \sum_{i=1}^n \phi(z_i) - \frac{1}{n} \sum_{i=n+1}^{2n} \phi(z_i) \right| \geq \epsilon \right) \\ &\leq 2\mathbb{P} \left(\sup_{\phi \in \Phi} \left| \frac{1}{n} \sum_{i=1}^n \phi(z_i) - \mathbb{E}[\phi(z)] \right| \geq \frac{\epsilon}{2} \right) \end{aligned}$$

Step II (Symmetrization)

Denote

$$N^\Phi(z_1, \dots, z_n) = |\{(\phi(z_1), \dots, \phi(z_n)), \phi \in \Phi\}|$$

$$N^\Phi(n) = \max_{z_1 \dots z_n} N^\Phi(z_1 \dots z_n)$$

$N^\Phi(n)$ is growth function.

Lemma

$$\mathbb{P} \left(\sup_{\phi \in \Phi} |\nu_1(z) - \nu_2(z)| \geq \epsilon \right) \leq \mathbb{E} [N^\Phi(z_1, \dots, z_n)] 2e^{-2n\epsilon^2} \leq N^\Phi(2n) 2e^{-2n\epsilon^2}$$

where $\nu_1(z) = \frac{1}{n} \sum_{i=1}^n \phi(z_i)$, $\nu_2(z) = \frac{1}{n} \sum_{i=n+1}^{2n} \phi(z_i)$.

Draw a set, permutation, fix set (draw without replacement)

$$\mathbb{P}_{\text{permutation}} \left(\sup_{\phi \in \Phi} |\nu_1(z) - \nu_2(z)| \geq \epsilon \right) \leq N^\Phi(z_1, \dots, z_n) 2e^{-2n\epsilon^2}$$

Take expectation,

$$\mathbb{E}_{\text{set}} \left\{ \mathbb{P}_{\text{permutation}} \left(\sup_{\phi \in \Phi} |\nu_1(z) - \nu_2(z)| \geq \epsilon \right) \right\} \leq \mathbb{E} [N^\Phi(z_1, \dots, z_n)] 2e^{-2n\epsilon^2}$$