

Lecture 6

①

D Empirical Processes

Consider X_1, \dots, X_n generated iid from (\mathcal{X}, Σ)

Let $\mu_n(A) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i \in A)$, $A \in \Sigma$.

Then, $E[\mu_n(A)] = \mu(A)$.

Furthermore, by Hoeffding, $P(|\mu_n(A) - \mu(A)| > \epsilon) \leq 2e^{-2n\epsilon^2}$.

Thus, by Borel-Cantelli lemma, $\mu_n(A) \xrightarrow{a.s.} \mu(A) \forall A$.

Does $\mu_n \rightarrow \mu$ in total variation distance?

No. If μ is continuous, $\mu(\text{supp}(\mu_n)) = 0 \forall n$.

$$d_{TV}(\mu_n, \mu) = \sup_A |\mu_n(A) - \mu(A)|$$

→ For continuous μ , $d_{TV}(\mu_n, \mu) = 1 \forall n$.

We shall be interested in $\boxed{\sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)|}$

(i) If $\mathcal{A} = \{(-\infty, x], x \in \mathbb{R}\}$,

$$\sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| = \sup_x \underbrace{|F_n(x) - F(x)|}_{\text{Kolmogorov-Smirnov}}$$

$d_{KS}(\mu_n, \mu) \equiv$ Kolmogorov-Smirnov
distance b/w μ_n and μ

Theorem (Glivenko-Cantelli)

$$d_{KS}(\mu_n, \mu) \rightarrow 0 \text{ a.s.}$$

Proof. Suffices to show: $P(d_{KS}(M_n, M) \geq \delta_n + \epsilon) \leq e^{-n\epsilon^2}$ for $\delta_n \rightarrow 0$. (2)

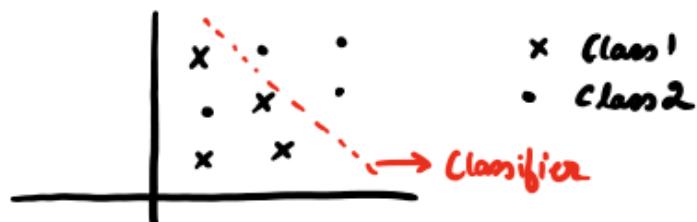
Note that for $\Delta_n(X_1, \dots, X_n) = d_{KS}(M_n, M)$, by changing only one of X_i we change $|F_n(x) - F(x)|$ for all $x \geq x_i$ by at most $\frac{1}{n}$. Thus, Δ_n satisfies BDP with $(\frac{1}{n}, \dots, \frac{1}{n})$.

By McDiarmid's inequality, we get

$$P(d_{KS}(M_n, M) \geq E[d_{KS}(M_n, M)] + \epsilon) \leq e^{-\frac{n\epsilon^2}{2}}.$$

We shall bound $E[d_{KS}(M_n, M)]$.

(II) Empirical risk minimization (ERM)



Given a family of classifiers \mathcal{F} (e.g. all linear/affine classifiers), find f^* which performs roughly as well as the best classifier in \mathcal{F} .

Bayesian formulation: $P(x^*) = \theta P_1(x^*) + \bar{\theta} P_2(x^*)$

Find $i \in \{1, 2\}$ by observing X_1, \dots, X_n .
(MAP is called naive Bayes)

A classifier is specified by an "acceptance" region A where you declare P_i .

Practice: θ, P_1, P_2 are unknown.

(3)

$$\begin{aligned}
 \text{ERM procedure: Select } f_n^* &= \underset{f \in \mathcal{F}}{\operatorname{argmin}} \mu_n(f(x) \neq y) \\
 &= \underset{f \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}(f(x_i) \neq y_i) \\
 &\equiv \underset{A \in \mathcal{A}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n [\mathbb{1}(x_i \notin A, y_i = 1) \\
 &\quad + \mathbb{1}(x_i \in A, y_i \neq 1)]
 \end{aligned}$$

Vapnik and Chervonenkis showed that

$$\left(\mu / f_n^*(x) \neq y - \min_{f \in \mathcal{F}} \mu / f(x) \neq y \right) \rightarrow 0 \text{ a.s.}$$

Proof: (a) McDiarmid for conc. (b) Bound for $\mathbb{E} \left[\max_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| \right]$

For A , denote

$$\Delta_n(A) := \sup_{A' \in \mathcal{A}} |\mu_n(A') - \mu(A')|$$

Our goal is to bound $\mathbb{E}[\Delta_n(A)]$.

(I) Method 1: Using Vapnik Chervonenkis Theory

$$\mathbb{E}[\Delta_n(A)] = \mathbb{E} \left[\sup_{A' \in \mathcal{A}} \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{x_i \in A'\}} - \mu(A') \right| \right]$$

Step 1. Symmetrization

Let X'_1, \dots, X'_n be an indep. copy of X_1, \dots, X_n .

Then, conditioned on $X = (x_1, \dots, x_n)$,

$$\begin{aligned} & \sup_{A \in \mathcal{A}} \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{x_i \in A\}} - \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X'_i \in A\}} \right] \right| \\ & \leq \sup_{A \in \mathcal{A}} \mathbb{E} \left[\left| \frac{1}{n} \sum_{i=1}^n (\mathbb{1}_{\{x_i \in A\}} - \mathbb{1}_{\{X'_i \in A\}}) \right| \right] \\ & \quad (\text{Jensen's inequality}) \\ & \leq \mathbb{E} \left[\sup_{A \in \mathcal{A}} \left| \frac{1}{n} \sum_{i=1}^n (\mathbb{1}_{\{x_i \in A\}} - \mathbb{1}_{\{X'_i \in A\}}) \right| \right] \end{aligned}$$

Therefore,

$$\mathbb{E}[\Delta_n(A)] \leq \mathbb{E} \left[\sup_{A \in \mathcal{A}} \frac{1}{n} \left| \sum_{i=1}^n (\mathbb{1}_{\{x_i \in A\}} - \mathbb{1}_{\{X'_i \in A\}}) \right| \right]$$

Step 2: Rademacher Sums

Let $\sigma_1, \dots, \sigma_n$ be iid unif $\{-1, +1\}$.

Then, R.H.S. above equals

$$\mathbb{E} \left[\underbrace{\sup_{A \in \mathcal{A}} \frac{1}{n} \left| \sum_{i=1}^n \sigma_i (\mathbb{1}_{\{x_i \in A\}} - \mathbb{1}_{\{X'_i \in A\}}) \right|}_{n^{-1} \cdot Z(X_1, \dots, X_n, X'_1, \dots, X'_n)} \right]$$

$$\begin{aligned} & \mathbb{E}[Z | X = x, X' = x'] \\ & = \mathbb{E} \left[\sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \sigma_i (\mathbb{1}_{\{x_i \in A\}} - \mathbb{1}_{\{X'_i \in A\}}) \right| \right] \quad (1) \end{aligned}$$

Recall: $\mathbb{E} \left[\max_{1 \leq i \leq n} X_i \right] \leq \sigma \sqrt{2 \log n}$ if
 $X_i \in \mathcal{G}(\sigma^2)$

If we apply this bound naively to (1), we will be hit by $|\mathcal{A}|$, which may even be ∞ . (5)

Step 3: Shatter coefficients

Definition (Shatter coefficient, VC dimension)

The n^{th} shatter coefficient of \mathcal{A} , $S_{\mathcal{A}}(n)$, is

$$\text{given by } S_{\mathcal{A}}(n) = \sup_{x_1, \dots, x_n} |\{ \{x_1, \dots, x_n\} \subseteq A, A \in \mathcal{A} \}|.$$

$x_1, \dots, x_n \leftarrow \text{denoted } V(\mathcal{A})$

The VC dimension of \mathcal{A} is the largest n such that $S_{\mathcal{A}}(n) = 2^n$. (Sauer: $S_{\mathcal{A}}(n) \leq (n+1)^{V(\mathcal{A})}$)

How does it apply to our quantity in (1)?

Note that the maximum number of distinct sequences

$$b(A) = (\mathbb{1}_{\{x_1 \in A\}}, \dots, \mathbb{1}_{\{x_n \in A\}}, \mathbb{1}_{\{x'_1 \in A\}}, \dots, \mathbb{1}_{\{x'_n \in A\}}) \\ \in \{0, 1\}^{2n}$$

that we see for a fixed (x, x') and different $A \in \mathcal{A}$ is bounded above by $S_{\mathcal{A}}(2n)$.

Thus, there exists $\hat{\mathcal{A}}$ with $|\hat{\mathcal{A}}| \leq S_{\mathcal{A}}(2n)$ such that the right-side of (1) equals

$$\mathbb{E} \left[\sup_{A \in \hat{\mathcal{A}}} \left| \sum_{i=1}^n \sigma_i (\mathbb{1}_{\{x_i \in A\}} - \mathbb{1}_{\{x'_i \in A\}}) \right| \right]$$

$$\leq \sqrt{2\sigma \log 2S_A(2n)} \quad (2)$$

(6)

if

$$\sum_{i=1}^n \sigma_i (\mathbb{1}_{\{\underline{x}_i \in A\}} - \mathbb{E}\mathbb{1}_{\{\underline{x}_i \in A\}}) \in \mathcal{G}(\sigma). \quad (3)$$

By Hoeffding's lemma, (3) holds with
 $\sigma = n$.

Therefore, (2) and (1) imply

$$\mathbb{E}[\Delta_n(A)] \leq \sqrt{\frac{2}{n} \log 2S_A(2n)}$$

Finally, note that $S_A(2n) \leq S_A(n)^2$
 to get

$$\mathbb{E}[\Delta_n(A)] \leq \sqrt{\frac{4 \log S_A(n) + 2 \log 2}{n}}$$

Example.

For $A = \{(-\infty, \theta), \theta \in \mathbb{R}\}$, for any fixed sequence $(\underline{x}_1, \dots, \underline{x}_n) \in \mathbb{R}^n$,

$$|\{A \cap \{\underline{x}_1, \dots, \underline{x}_n\}, A \in A\}| \leq n+1$$

since if we arrange $\underline{x}_1, \dots, \underline{x}_n$ in an increasing order, the only possible $b(\underline{x})$ are $(0, \dots, 0), (1, 0, \dots, 0), (1, 1, 0, \dots, 0), \dots, (1, \dots, 1)$

Thus, for this case,

(7)

$$\begin{aligned}\mathbb{E}[\Delta_n(A)] &= \mathbb{E}[d_{KS}(\mu_n, \mu)] \\ &\leq \sqrt{\frac{4 \log(n+1) + 2 \log 2}{n}}\end{aligned}$$

Combining this with Mc Diarmid, we get

$$P(d_{KS}(\mu_n, \mu) > \sqrt{\frac{8 \log n}{n}} + \epsilon) \leq e^{-\frac{n\epsilon^2}{2}}.$$

By using the Borel-Cantelli lemma :

$$d_{KS}(\mu_n, \mu) \rightarrow 0 \text{ a.s. } (\text{Gliivenko-Cantelli})$$