

Lecture 3

Recap:

- PAC learnability for finite hypothesis classes
- No-free lunch theorem
- Agnostic PAC learning
- Define Bayes-optimal predictor for 0/1 loss.

$$h^*(x) = \mathbb{1}(\eta(x) \geq 1/2)$$

$$\eta(x) = P_{\mathcal{D}}(y=1 | x=x)$$

$$R(h^*) = \mathbb{E}_{x \sim \mathcal{D}(x)} [\min\{\eta(x), 1 - \eta(x)\}]$$

Today

- Define Agnostic PAC
- Define Uniform Convergence (UC)
- UC \Rightarrow Agnostic PAC Learning
- VC results
- Concentration of measure

Agnostic PAC

Definition (Agnostic PAC Learnability).

A hypothesis \mathcal{H} is agnostic PAC learnable if for every $\epsilon, \delta \in (0, 1)$ \exists a function $n_{\mathcal{H}}(\epsilon, \delta)$ & a learning algorithm s.t. for every distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$, if the algorithm is run on $n \geq n_{\mathcal{H}}(\epsilon, \delta)$ samples drawn iid from \mathcal{D} , then the algorithm returns a hypothesis \hat{h} with $R(\hat{h}) \leq \min_{h \in \mathcal{H}} R(h) + \epsilon$, except with probability δ .

PAC learning assumes realizability.

Realizability : h^* (Bayes optimal predictor) $\in \mathcal{H}$.

\therefore Agnostic PAC more general than PAC learning

e.g. $\mathcal{H} = \{ h_w(x) : \mathbb{1}(w^T x > 0), w \in \mathbb{R}^d \}$

(linear classifiers)

Agnostic PAC for this \mathcal{H} asks we do as well as best possible linear classifier.

Uniform Convergence

Idea: ERM outputs $\hat{h} \in \mathcal{H}$, which has minimum risk over training set

Want: \hat{h} is close to minimizer of population risk ($\min_{h \in \mathcal{H}} R(h)$).

Suffices: All empirical risks of all members of \mathcal{H} are close to their population risk.

$$\tilde{h} = \operatorname{argmin}_{h \in \mathcal{H}} R(h) \quad (\text{Note that } \tilde{h} \neq h^*)$$

$$\begin{aligned} R(h_{s, \text{ERM}}) - R(\tilde{h}) &= R(h_{s, \text{ERM}}) - \hat{R}(h_{s, \text{ERM}}) \\ &\quad + \hat{R}(h_{s, \text{ERM}}) - \hat{R}(\tilde{h}) \\ &\leq 0 \quad + \hat{R}(\tilde{h}) - R(\tilde{h}) \end{aligned}$$

trickier (pointing to the first two terms)
easy to bound (pointing to the last two terms)

$$\hat{R}(\tilde{h}) = \frac{1}{n} \sum_{i=1}^n \ell(\tilde{h}(x_i), y_i)$$

avg. of n i.i.d. random variable

$$\forall i \in [n], \quad \mathbb{E}_{(x_i, y_i) \sim \mathcal{D}} [\ell(\tilde{h}(x_i), y_i)] = R(\tilde{h})$$

$$P_{\mathcal{D}}[|\hat{R}(\tilde{h}) - R(\tilde{h})| > \epsilon] \leq \delta.$$

Holds for any fixed hypothesis

Definition (Uniform convergence)

A hypothesis class \mathcal{H} has the uniform convergence property if for every $\epsilon, \delta \in (0, 1)$ there exists a function $n_{\mathcal{H}}^{uc}(\epsilon, \delta)$ such that for every dist. \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$, if S is a training set of $n \geq n_{\mathcal{H}}^{uc}(\epsilon, \delta)$ samples drawn iid from \mathcal{D} , then w.p. $1 - \delta$,

$$\forall h \in \mathcal{H}, |\hat{R}(h) - R(h)| \leq \epsilon.$$

Proposition (UC \Rightarrow Agnostic PAC learning)

If \mathcal{H} has the UC property with $n_{\mathcal{H}}^{\text{UC}}(\epsilon, \delta)$, then \mathcal{H} is Agnostic PAC learnable with sample complexity $n_{\mathcal{H}}(\epsilon, \delta) \leq n_{\mathcal{H}}^{\text{UC}}(\epsilon/2, \delta)$. Moreover, ERM is an algorithm which achieves this sample complexity.

Proof

Let S be a sample of size $n \geq n_{\mathcal{H}}^{\text{UC}}(\epsilon/2, \delta)$.

By definition,

$$\forall h \in \mathcal{H}, |\hat{R}(h) - R(h)| \leq \epsilon/2.$$

Consider ERM $h_{S, \text{ERM}}$ & let $\tilde{h} = \underset{h \in \mathcal{H}}{\text{argmin}} R(h)$.

$$\begin{aligned} R(h_{S, \text{ERM}}) - R(\tilde{h}) &= R(h_{S, \text{ERM}}) - \hat{R}(h_{S, \text{ERM}}) \\ &\leq \epsilon/2 \\ &\leq 0 + \hat{R}(h_{S, \text{ERM}}) - \hat{R}(\tilde{h}) \\ &\leq \epsilon/2 + \hat{R}(\tilde{h}) - R(\tilde{h}) \\ &\leq \epsilon/2 + \epsilon/2 \\ &= \epsilon. \end{aligned}$$

$$R(h_{S, \text{ERM}}) - R(\tilde{h}) \leq \epsilon.$$

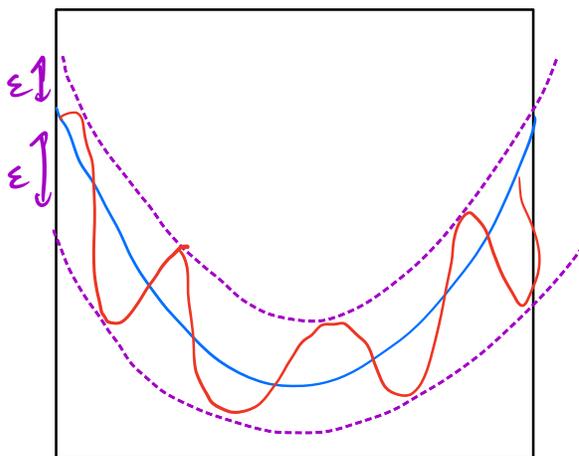
•

Consider: \mathcal{H} is parameterized by single parameter
($d \in \mathbb{R}$)

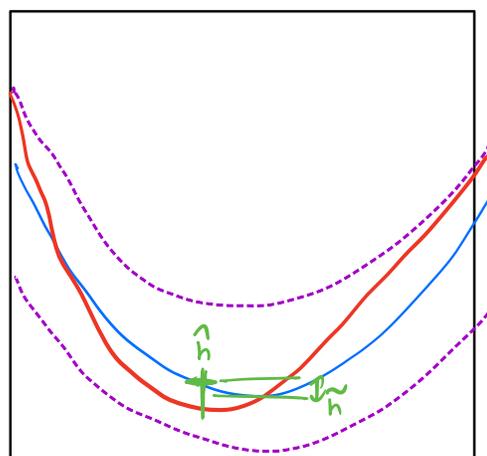
— : $R(h)$

--- : VC bound

— : $\hat{R}_s(h)$



— d →



— d →

Theorem (Agnostic PAC for finite classes)

Let \mathcal{H} be a class $|\mathcal{H}| < \infty$. Then \mathcal{H} is agnostic - PAC learnable with

$$n_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{2 \log(2|\mathcal{H}|/\delta)}{\epsilon^2} \right\rceil.$$

Proof (via Uniform convergence)

We will show:

(1) For any fixed $h \in \mathcal{H}$ and $\epsilon > 0$,

$$\Pr[|\hat{R}(h) - R(h)| \leq \epsilon] \geq 1 - 2e^{-2n\epsilon^2}.$$

(2) For any $\epsilon > 0$,

$$\Pr[\forall h \in \mathcal{H}, |\hat{R}(h) - R(h)| \leq \epsilon] \geq 1 - 2|\mathcal{H}|e^{-2n\epsilon^2}.$$

(3) For $n \geq \left\lceil \frac{\log(2|\mathcal{H}|/\delta)}{2\epsilon^2} \right\rceil$, with probability $1 - \delta$,

$$|\hat{R}(h) - R(h)| \leq \epsilon \quad \forall h \in \mathcal{H}.$$

(4) By UC, \mathcal{H} is agnostic PAC learnable

with $n \geq \left\lceil \frac{2 \log(2|\mathcal{H}|/\delta)}{\epsilon^2} \right\rceil$ samples.

Proof of (i)

Lemma (Hoeffding's inequality):

Let x_1, x_2, \dots, x_n be independent random variables such that $a_i \leq x_i \leq b_i$ for each $i \in [n]$. Then for any $\epsilon > 0$,

$$\begin{aligned} P_n \left[\left| \frac{1}{n} \sum_{i=1}^n x_i - \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n x_i \right] \right| \leq \epsilon \right] \\ \geq 1 - 2 \exp \left(- \frac{2 n^2 \epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2} \right) \end{aligned}$$

Given Hoeffding's, we prove (i).

Take each $x_i = \ell(h(x_i), y_i)$

Since $\ell(h(x_i), y_i) = \mathbb{1}(h(x_i) \neq y_i)$

$\overset{R(h)}{\underbrace{x_i}} \in \{0, 1\} \Rightarrow a_i = 0, b_i = 1.$

$$P_n \left[\left| \frac{1}{n} \sum_{i=1}^n \overset{R(h)}{\underbrace{x_i}} - \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \overset{R(h)}{\underbrace{x_i}} \right] \right| \leq \epsilon \right]$$

$$\geq 1 - 2 \exp(-2n\epsilon^2)$$

Proof of (2).

$$\begin{aligned} & P_n [\forall h \in \mathcal{H}, |\hat{R}(h) - R(h)| \leq \varepsilon] \\ &= 1 - P_n \left[\bigcup_{i=1}^{|\mathcal{H}|} |\hat{R}(h_i) - R(h_i)| > \varepsilon \right] \\ &\geq 1 - \sum_{i=1}^{|\mathcal{H}|} P_n [|\hat{R}(h_i) - R(h_i)| > \varepsilon] \\ &\geq 1 - |\mathcal{H}| (2 \exp(-2n \varepsilon^2)) \end{aligned}$$

Proof of (3).

$$\text{Set } n \geq \left\lceil \frac{\log(2|\mathcal{H}|/\delta)}{2\varepsilon^2} \right\rceil$$

failure prob. $\leq \delta$.

Proof of (4)

UC \Rightarrow Agnostic PAC



Concentration inequalities

Let $\{x_1, x_2, \dots\}$ be sequence of random variables
 $E[x_i] = \mu$, $\text{Var}(x_i) = \sigma^2 < \infty$. Let $\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i$.

How close is \bar{x}_n to μ ?

Central Limit Theorem

Thm (Central Limit Thm): As $n \rightarrow \infty$

$$(\bar{x}_n - \mu) \xrightarrow{\text{dist.}} \mathcal{N}(0, \sigma^2/n).$$

\bar{x}_n converges to μ at rate of $O\left(\frac{\sigma}{\sqrt{n}}\right)$.

Markov's bound.

Let x be a non-negative random variable.

Then for any $t > 0$, we have

$$P[x \geq t] \leq \frac{E[x]}{t}.$$

Proof.

$$E[x] = \sum_{i=0}^{\infty} i P[x=i]$$

$$\geq \sum_{i=t}^{\infty} i P[x=i] \quad (\text{non-negativity})$$

$$\geq t \cdot \sum_{i=t}^{\infty} P[x=i]$$

$$\Rightarrow P[x \geq t] \leq E[x]/t. \quad \blacktriangleright$$

Exercise:

Show Markov's bound is tight.

Only applies to non-negative.

Can work with variance.

Chebyshev's inequality :

Let X be a random variable. Then for any $t > 0$.

$$P_X [|X - E(X)| \geq t] \leq \frac{\text{Var}[X]}{t^2}.$$

Proof.

$$P_X [|X - E(X)| \geq t] = P_X [|X - E(X)|^2 \geq t^2]$$

Apply Markov's inequality to the non-negative r.v. $(X - E(X))^2$,

$$\begin{aligned} P_X [(X - E(X))^2 \geq t^2] &\leq \frac{E [(X - E(X))^2]}{t^2} \\ &= \frac{\text{Var}[X]}{t^2}. \end{aligned}$$

Is this enough to get Hoeffding's?

Lemma (Hoeffding's inequality):

Let x_1, x_2, \dots, x_n be independent random variables such that $a_i \leq x_i \leq b_i$ for each $i \in [n]$. Then for any $\epsilon > 0$,

$$\begin{aligned} P \left[\left| \frac{1}{n} \sum_{i=1}^n x_i - \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n x_i \right] \right| \leq \epsilon \right] \\ \geq 1 - 2 \exp \left(- \frac{2 n^2 \epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2} \right) \end{aligned}$$

$$a_i = 0$$

$$b_i = 1$$

$$\mathbb{E}[x_i] = \mu \quad \forall i.$$

$$\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\text{Var}[\bar{x}_n] = \frac{1}{n^2} \sum_{i=1}^n \text{Var}[x_i]$$

$$\leq \frac{1}{n}$$

$$\mathbb{E}[\bar{x}_n] = \mu$$

For any $n > 0$.
 x_1, x_2 (independent),
 $\text{Var}[x_1 + x_2] =$
 $\text{Var}[x_1] + \text{Var}[x_2]$

$$\begin{aligned} & \text{Var}[x_i] \\ &= \mathbb{E}[x_i^2] - \mathbb{E}[x_i]^2 \\ &\leq 1 \end{aligned}$$

By Chebyshev,

$$\begin{aligned} P_n [|\bar{X}_n - \mu| \geq \varepsilon] &\leq \frac{\text{Var}[\bar{X}_n]}{\varepsilon^2} \\ &\leq \frac{1}{n\varepsilon^2}. \end{aligned}$$

But Hoeffding's inequality says,

$$P_n [|\bar{X}_n - \mu| \geq \varepsilon] \leq 2e^{-2n\varepsilon^2}.$$

Idea: No need to stop at second moment.

for any positive integer k ,

$$\begin{aligned} P_n [|X - E(X)| \geq t] &= P_n [|X - E(X)|^k \geq t^k] \\ &\leq \frac{E(|X - E(X)|^k)}{t^k}. \end{aligned}$$

Consider functions other than polynomial.

For $\lambda \geq 0$,

$$\begin{aligned} P_X [X - \mathbb{E}(X) \geq t] &= P_X [e^{\lambda(X - \mathbb{E}(X))} \geq e^{\lambda t}] \\ &\leq \frac{\mathbb{E}(e^{\lambda(X - \mathbb{E}(X))})}{e^{\lambda t}} \\ &\leq \inf_{\lambda \geq 0} \frac{\mathbb{E}(e^{\lambda(X - \mathbb{E}(X))})}{e^{\lambda t}} \end{aligned}$$

Chernoff-style bounds.

$$\begin{aligned} \log(P_X [X - \mathbb{E}(X) \geq t]) &\leq \\ &\inf_{\lambda \geq 0} (\log(\mathbb{E}(e^{\lambda(X - \mathbb{E}(X))})) - \lambda t). \end{aligned}$$

Moment-generating function.

For any r.v. X , the moment generating function

$$M_X(t) = \mathbb{E}(e^{tX}).$$

$$\text{Exercise: } \left. \frac{\partial^k M_X(t)}{\partial t^k} \right|_{t=0} = \mathbb{E}[X^k]$$