

Sample complexity of agnostic PAC learning

- Recall the notion of ϵ -approximation

Definition: Let D be a probability distribution over X . Let $H \subseteq \mathcal{P}(X)$.

A finite set (multiset) $S \subseteq X$ is an ϵ -approximation for D, H

if for all $h \in H$,

$$\left| D(h) - \frac{|S \cap h|}{|S|} \right| \leq \epsilon$$

Lemma: Let D be a distribution over $X \times Y$. Let $H \subseteq Y^X$.

Let $S \subseteq X \times Y$ be an ϵ -approximation

for $D, G^c(H)$ then

$$1) \forall h \in H: |\text{err}_D(h) - \widehat{\text{err}}_S(h)| \leq \epsilon$$

$$2) \text{err}_D(\text{ERM}(S)) \leq \min_{h \in H} \text{err}_D(h) + 2\epsilon$$

Proof:

$$\bullet \text{err}_D(h) = \Pr[h(x) \neq y]$$

$$= \mathbb{D}\left(\{ (x, y) \in X \times Y : h(x) \neq y \}\right)$$

$$= \mathbb{D}(g^c(h))$$

$$\bullet \widehat{\text{err}}_S(h) = \frac{|\{ (x, y) \in S : h(x) \neq y \}|}{|S|}$$

↳ call

$$= \frac{|S \cap g^{-1}(h)|}{|S|}$$

- If S is ϵ -approximation for $D, G^c(H)$ then

$$\forall h \in H: \left| D(g^c(h)) - \frac{|S \cap g^c(h)|}{|S|} \right| \leq \epsilon$$

- Therefore

$$\forall h \in H: |\text{err}_D(h) - \widehat{\text{err}}_S(h)| \leq \epsilon$$

This proves part 1).

Consider FERM algorithm

Lemma 1.1.1 (ERM)

Since ERM minimizes $\widehat{\text{err}}_S(\cdot)$,

- $\forall h \in H: \widehat{\text{err}}_S(\text{ERM}(S)) \leq \widehat{\text{err}}_S(h)$
- Use part 1) twice:

For all $h \in H$:

$$\text{err}_D(\text{ERM}(S)) \leq \widehat{\text{err}}_S(\text{ERM}(S)) + \epsilon$$

$$\leq \widehat{\text{err}}_S(h) + \epsilon$$

$$\leq \text{err}_D(h) + \epsilon + \epsilon$$

$$= \text{err}_D(h) + 2\epsilon$$

Now we can conclude part 1)

Here we use the following:

- $\text{err}_D(\text{ERM}(S)) \leq \text{err}_D(h) + 2\epsilon$

- Therefore

$$\text{err}_D(\text{ERM}(S)) \leq \min_{h \in H} \text{err}_D(h) + 2\epsilon$$



- We want to prove that if S is an i.i.d. from D then with high probability S is an ϵ -approximation.

- We will need two tools

from probability theory

1) McDiarmid's inequality

2) Massart's lemma

McDiarmid's inequality:

Let X_1, X_2, \dots, X_n be independent random variables such that $X_i \in A_i$.

Let $f: A_1 \times A_2 \times \dots \times A_n \rightarrow \mathbb{R}$ be a function such that for any i ,

for any $x_1, x_2, \dots, x_n, x'_i$

$$\left| f(x_1, x_2, \dots, x_n) - f(x_1, x_2, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n) \right| \leq c_i$$

(bounded difference property).

Then

$$\Pr \left[\left| f(x_1, \dots, x_n) - \mathbb{E}[f(x_1, \dots, x_n)] \right| > \varepsilon \right] \leq 2e^{-\frac{\varepsilon^2}{\sum_{i=1}^n c_i^2}}$$

- We omit the proof
(The proof relies on martingales.)
- The inequality is a generalization of Hoeffding's inequality.
- We can recover Hoeffding's inequality if we take $f(x_1, \dots, x_n) = x_1 + x_2 + \dots + x_n$

Massart's lemma:

... .. variables

Let X_1, X_2, \dots, X_n be random variables
such that for all $i, s > 0$,

$$\mathbb{E} \left[e^{sX_i} \right] \leq e^{s^2 \sigma^2 / 2}$$

1) Then

$$\mathbb{E} \left[\max_{i=1..n} X_i \right] \leq \sigma \sqrt{2 \ln n}$$

2) If additionally for all $i, s > 0$,

$$\mathbb{E} \left[e^{-sX_i} \right] \leq e^{s^2 \sigma^2 / 2}$$

then

$$\mathbb{E} \left[\max_{i=1..n} |X_i| \right] \leq \sigma \sqrt{2 \ln(2n)}$$

Proof:

\rightarrow follows from by applying 1)

• \hookrightarrow follows from ...

to $\underbrace{X_1, \dots, X_n, -X_1, \dots, -X_n}_{2n \text{ variables}}$

• It suffices to prove 1)

$$e^{s \mathbb{E}[\max_i X_i]} \leq \mathbb{E}\left[e^{s \max_i X_i}\right]$$

$$= \mathbb{E}\left[\max_i e^{s X_i}\right]$$

$$\leq \sum_{i=1}^n \mathbb{E}\left[e^{s X_i}\right]$$

$$\leq n \cdot \frac{s^2 \sigma^2}{2}$$

$$\leq n \epsilon$$

Take logarithm and divide
by $s > 0$:

$$\mathbb{E} \left[\max_i X_i \right] \leq \frac{\ln n}{s} + \frac{s\sigma^2}{2}$$

• Find $s > 0$ that minimizes
right-hand side.

• Optimal $s = \frac{\sqrt{2 \ln n}}{\sigma}$

• $\mathbb{E} \left[\max_i X_i \right] \leq \sigma \sqrt{\frac{1}{2} \ln n} + \sigma \sqrt{\frac{1}{2} \ln n}$
 $= \sigma \sqrt{2 \ln n}$





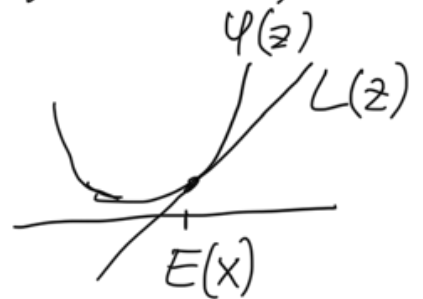
Lemma: Let X be a random variable. Let $\varphi: \mathbb{R} \rightarrow \mathbb{R}$ be a convex function. Then,

$$\varphi(\mathbb{E}[X]) \leq \mathbb{E}[\varphi(X)]$$

Proof: $L(z) = \varphi(\mathbb{E}[X]) + (z - \mathbb{E}[X])\varphi'(\mathbb{E}[X])$

$$\forall z \in \mathbb{R} \quad L(z) \leq \varphi(z)$$

$$\mathbb{E}[L(X)] \leq \mathbb{E}[\varphi(X)]$$



Lemma: Let \mathcal{D} be a distribution over X . Let $H \subseteq \mathcal{P}(X)$.

Let S be a i.i.d. sample from \mathcal{D} of size m . Then,

$$\mathbb{E} \left[\sup_{h \in H} \left| D(h) - \frac{|S_{nh}|}{|S|} \right| \right] \leq \sqrt{\frac{2 \ln(2 \pi_H(2m))}{m}}.$$

Proof:

• Let $S = (x_1, x_2, \dots, x_m)$

$$\frac{|S_{nh}|}{|S|} = \frac{1}{m} \sum_{i=1}^m \mathbb{1}[x_i \in h]$$

• Let $\gamma_1, \gamma_2, \dots, \gamma_m$ be an independent ghost sample.

$$D(h) = \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \mathbb{1}[\gamma_i \in h] \right]$$

• Let

$$P_m(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}[X_i \in h]$$

$$Q_m(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}[Y_i \in h]$$

• We need to upper bound

$$\mathbb{E} \left[\sup_{h \in H} |P_m(h) - \mathbb{E}[Q_m(h)]| \right]$$

$$= \mathbb{E} \left[\sup_{h \in H} | \mathbb{E}[P_m(h) - Q_m(h) \mid X_1, \dots, X_m] | \right]$$

$$\leq \mathbb{E} \left[\sup_{h \in H} \mathbb{E} \left[|P_m(h) - Q_m(h)| \mid X_1, \dots, X_m \right] \right]$$

$$\leq \mathbb{E} \left[\mathbb{E} \left[\sup |P_m(h) - Q_m(h)| \mid X_1, \dots, X_m \right] \right]$$

$$\left[\sup_{h \in H} |P_n(h) - Q_n(h)| \right]$$

$$= \mathbb{E} \left[\sup_{h \in H} |P_n(h) - Q_n(h)| \right]$$

$$= \frac{1}{m} \mathbb{E} \left[\sup_{h \in H} \left| \sum_{i=1}^m (\mathbb{1}[X_i \in h] - \mathbb{1}[Y_i \in h]) \right| \right]$$

• Let us introduce Rademacher variables $\sigma_1, \sigma_2, \dots, \sigma_m$.

These are i.i.d. variables independent of $X_1, \dots, X_m, Y_1, \dots, Y_m$ such that

$$\Pr[\sigma_i = +1] = \Pr[\sigma_i = -1] = \frac{1}{2}$$

• Note that the distribution of $\left(\mathbb{1}[X_i \in h] - \mathbb{1}[Y_i \in h] \right)$ is the same as the distribution of $\sigma_i \left(\mathbb{1}[X_i \in h] - \mathbb{1}[Y_i \in h] \right)$.

$$\begin{aligned} & \bullet \mathbb{E} \left[\sup_{h \in H} \left| \sum_{i=1}^m \left(\mathbb{1}[X_i \in h] - \mathbb{1}[Y_i \in h] \right) \right| \right] \\ &= \mathbb{E} \left[\sup_{h \in H} \left| \sum_{i=1}^m \sigma_i \left(\mathbb{1}[X_i \in h] - \mathbb{1}[Y_i \in h] \right) \right| \right] \end{aligned}$$

• Let $U = (X_1, \dots, X_m, Y_1, \dots, Y_m)$ be sample of size $2m$.

- Condition on U

- There are at most $\pi_H(2m)$ possible behaviors of H on U .

- After conditioning on U

$$\mathbb{E} \left[\sup_{h \in H} \left| \sum_{i=1}^m \sigma_i (\mathbb{1}[x_i \in h] - \mathbb{1}[\gamma_i \in h]) \right| \right]$$

$$= \mathbb{E} \left[\max_{h \in \pi_H(U)} \left| \sum_{i=1}^m \sigma_i (\mathbb{1}[x_i \in h] - \mathbb{1}[\gamma_i \in h]) \right| \right]$$

- For any $h \in H$, conditioned on U ,

$$Z_h = \sum_{i=1}^m \sigma_i (\mathbb{1}[x_i \in h] - \mathbb{1}[\gamma_i \in h])$$

is random variable such that

$$1) \mathbb{E}[z_n] = 0$$

$$2) z_n = \sum_{i=1}^m z_{n,i} \quad \text{where}$$

$$z_{n,i} \in [-1, 1] \quad \text{and} \quad \mathbb{E}[z_{n,i}] = 0$$

3) By Hoeffding's lemma, $s \in \mathbb{R}$,

$$\mathbb{E}[e^{s z_{n,i}}] \leq e^{s^2/2}$$

$$\mathbb{E}[e^{s z_n}] = \prod_{i=1}^m \mathbb{E}[e^{s z_{n,i}}] \leq e^{m s^2/2}$$

• By Massart's lemma

$$\mathbb{P} \left[\left| \sum_{i=1}^m z_{n,i} \right| \geq t \right] \leq 2 e^{-t^2/2m}$$

$$\mathbb{E} \left[\max_{h \in \Pi_H(\mathcal{U})} \left| \sum_{i=1}^m \sigma_i (\| [x_i \in h] - \mathbb{1}[y_i \in h] \|) \right| \right]$$

$$= \mathbb{E} \left[\max_{h \in \Pi_H(\mathcal{U})} |z_h| \right]$$

$$\leq \sqrt{2m \ln(2\pi_H(2m))}$$



Theorem: Let \mathcal{D} be a probability distribution over X . Let $\delta \in (0, 1)$. Let $H \subseteq \mathcal{P}(X)$. Let S be an i.i.d. sample from \mathcal{D} of size m . Then with probability at least $1 - \delta$, S is an ϵ -approximation for

H_1, D where

$$\varepsilon = \sqrt{\frac{2 \ln(2T_H(2m))}{m}} + \sqrt{\frac{\ln(7/5)}{2m}}$$

Proof:

• Let $Z = \sup_{h \in H} \left| D(h) - \frac{|S \cap h|}{|S|} \right|$

• Z is a random variable that is a function of $S = (X_1, X_2, \dots, X_m)$

• We can apply McDiarmid's inequality to Z .

• Suppose $S = (X_1, \dots, X_m)$,

$$S^i = (X_1, \dots, X_{i-1}, X_i', X_{i+1}, \dots, X_m)$$

and
$$z' = \sup_{h \in H} \left| D(h) - \frac{|S'nh|}{|S'|} \right|$$

• For all $h \in H$:

$$\left| \left| D(h) - \frac{|S'nh|}{|S'|} \right| - \left| D(h) - \frac{|Snh|}{|S|} \right| \right|$$

$$\leq \left| \frac{|S'nh|}{|S'|} - \frac{|Snh|}{|S|} \right|$$

$$= \frac{1}{m} \underbrace{\left| |S'nh| - |Snh| \right|}_{\leq 1}$$

$$\leq \frac{1}{m}$$

• Therefore

$$\sup_{h \in H} \left| D(h) - \frac{|S_{nh}|}{|S|} \right| - \sup_{h \in H} \left| D(h) - \frac{|S'_{nh}|}{|S'|} \right| \leq \frac{1}{m}$$

- So Z satisfies McDiarmid's theorem with $c_i = \frac{1}{m}$.

- Therefore, with probability at least $1 - \delta$,

$$|Z - \mathbb{E}[Z]| \leq \sqrt{\frac{\ln(2/\delta)}{2m}}$$

- By previous lemma

$$\mathbb{E}[Z] \leq \sqrt{\frac{2 \ln(2\pi_H(2m))}{m}}$$

- Thus, w.p. at least $1 - \delta$,

$$Z \leq \mathbb{E}[Z] + \sqrt{\frac{\ln(2/\delta)}{2m}}$$

$$\leq \sqrt{\frac{2 \ln(2\pi_H(2m))}{m}} + \sqrt{\frac{\ln(2/\delta)}{2m}}$$



Corollary: Let \mathcal{D} be a distribution over X . Let $H \subseteq \mathcal{P}(X)$ and let $d = VC(H)$. Let $\epsilon, \delta \in (0, 1)$. Let S be an i.i.d. sample from \mathcal{D} of size m . With probability at least $1 - \delta$, S is an ϵ -approximation provided

$$m \geq \max \left(\frac{2 \ln(2/\delta)}{\epsilon^2}, \frac{4 \ln 2}{\epsilon^2}, \frac{16d \ln \left(\left(\frac{32e}{\epsilon^2} \right)^2 \right)}{\epsilon^2} \right)$$

Proof:

• It suffices to show that

$$\sqrt{\frac{2 \ln(2\pi_H(2m))}{m}} + \sqrt{\frac{\ln(2/\delta)}{2m}} \leq \varepsilon.$$

• Clearly

$$\sqrt{\frac{\ln(2/\delta)}{2m}} \leq \varepsilon/2$$

• It remains to show that

$$\sqrt{\frac{2 \ln(2\pi_H(2m))}{m}} \leq \varepsilon/2$$

• Since $2m \geq d$, $\pi_H(2m) \leq \left(\frac{2me}{1}\right)^d$

and

$$\sqrt{\frac{2 \ln(2\pi H(2m))}{m}}$$

$$\leq \sqrt{\frac{2 \ln\left(2\left(\frac{2me}{d}\right)^d\right)}{m}}$$

$$= \sqrt{\frac{2 \ln 2 + 2d \ln\left(\frac{2me}{d}\right)}{m}}$$

• Squaring both sides

$$\frac{2 \ln 2 + 2d \ln\left(\frac{2me}{d}\right)}{m} \leq \frac{\epsilon^2}{4}$$

• $\frac{2 \ln 2}{m} \leq \frac{\epsilon^2}{8}$ • It suffices

to prove

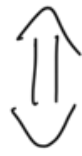
$$\frac{2d \ln\left(\frac{2me}{d}\right)}{m} \leq \frac{\epsilon^2}{8}$$

• Left-hand side is decreasing in m . It suffices to verify it for

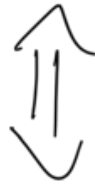
$$m = \frac{16d \ln\left(\left(\frac{32e}{\epsilon^2}\right)^2\right)}{\epsilon^2}$$

• We need to prove

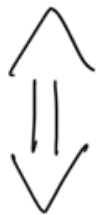
$$\frac{2d \ln\left(\frac{2me}{d}\right)}{m} \leq \frac{\epsilon^2}{8}$$



$$\frac{2d \ln \left(\frac{32ed}{\epsilon^2 d} \ln \left(\left(\frac{32e}{\epsilon^2} \right)^2 \right) \right)}{\frac{16d}{\epsilon^2} \ln \left(\left(\frac{32e}{\epsilon^2} \right)^2 \right)} \leq \frac{\epsilon^2}{8}$$

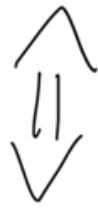


$$\frac{\ln \left(\frac{32e}{\epsilon^2} \ln \left(\left(\frac{32e}{\epsilon^2} \right)^2 \right) \right)}{\ln \left(\left(\frac{32e}{\epsilon^2} \right)^2 \right)} \leq 1$$

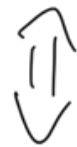


$$\frac{32e}{\epsilon^2} \ln \left(\left(\frac{32e}{\epsilon^2} \right)^2 \right) \leq \left(\frac{32e}{\epsilon^2} \right)^2$$

$$\varepsilon^4 \quad \left(\frac{\varepsilon^4}{\varepsilon^2} \right) \quad \left(\frac{\varepsilon^4}{\varepsilon^2} \right)$$



$$2 \ln \frac{32e}{\varepsilon^2} \leq \frac{32e}{\varepsilon^2}$$



$$2 \ln x \leq x$$

which is true for all $x > 0$.



Corollary: Let D be a distribution over $X \times Y$. Let $H \subseteq Y^X$. $d = VC(H)$

Let $\varepsilon, \delta \in (0, 1)$. Let S be

an i.i.d. sample from D

of size

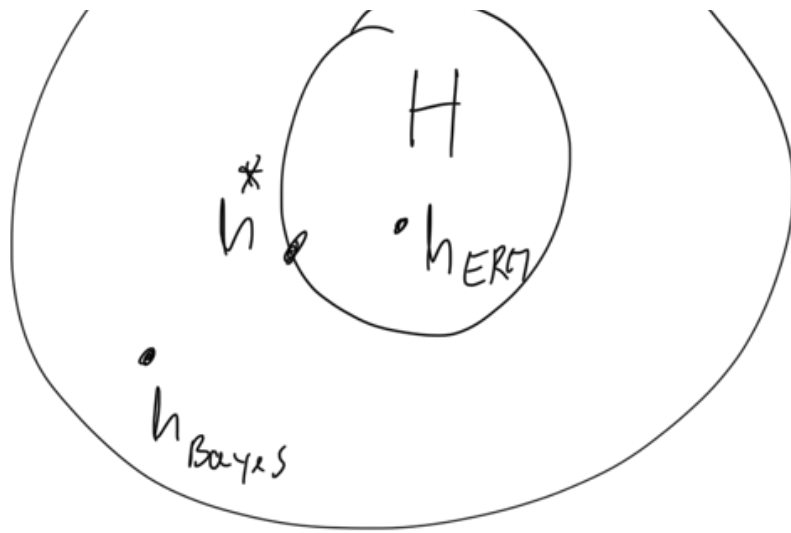
$$m \geq \max \left(\frac{8 \ln(2/\delta)}{\epsilon^2}, \frac{16 \ln 2}{\epsilon^2}, \frac{64d \ln \left(\left(\frac{128e}{\epsilon^2} \right)^2 \right)}{\epsilon^2} \right)$$

Then, with probability at least $1-\delta$,

$$\text{err}_D(\text{ERM}(S)) \leq \inf_{h \in H} \text{err}_D(h) + \epsilon.$$

Proof: Follows directly from previous corollary with $\epsilon \leftarrow \epsilon/2$ and lemmas above.





$$h^* = \underset{h \in H}{\operatorname{argmin}} \operatorname{err}_D(h)$$

$$\operatorname{err}_D(h_{\text{Bayes}}) - \operatorname{err}_D(h^*) = \text{approximation error}$$

$$\operatorname{err}_D(h_{\text{ERM}}) - \operatorname{err}_D(h^*) = \text{estimation error}$$