

PAC Learning:

PAC = Probably Approximately Correct

Definition (PAC) : A function class H is called PAC learnable if for every $\epsilon, \delta \in (0, 1]$, there exists a number $n_{\epsilon, \delta}$ such that the following holds:

There exist a learning algorithm that for any distribution D over $X \times Y$, by using a set S of $n_{\epsilon, \delta}$ data points, we obtain a predictor function h_S^* such that with probability $1 - \delta$ we have

$$\min_{h \in H} L_D(h) \leq L_D(h_S^*) \leq \min_{h \in H} L_D(h) + \epsilon.$$

$\mathcal{H} \rightarrow \left\{ \begin{array}{l} \text{Algorithm } \rightarrow h^*_S \\ n(\epsilon, \delta) \end{array} \right. \Rightarrow \mathcal{H} \text{ is PAC learnable}$

Let's start with the simplest possible hypothesis class and see what PAC-learning means with that class.

$$\mathcal{H} = \{h_1, h_2, \dots, h_m\}$$

$$\min_{h \in \mathcal{H}} L_D(h) \rightarrow \min_{h \in \mathcal{H}} L_S(h)$$



h^*_S is one of the functions $h \in \mathcal{H}$ with smallest empirical error.

$$\min_{h \in \mathcal{H}} L_S(h) = \min \{L_S(h_1), \dots, L_S(h_m)\}$$

What we'd like to find in this case
is a number $n(\epsilon, \delta)$ such that
for my distribution D over data
we have:

with probability $1-\delta$:

$$L_D(h^*_S) \leq \min_{h \in H} L_D(h) + \epsilon$$

In other words we're asking how large
the number of training data points
should be s.t. we have w.p. $1-\delta$
that $L_D(h^*_S) \leq \min_{h \in H} L_D(h) + \epsilon$.

In order to find $n(\epsilon, \delta)$ we need
to answer two main questions.

(1) Given a fixed function $h: x \rightarrow y$,
 how many training data point should
 we have such that:
 with probability $1-\delta$:

$$\left| \underbrace{L_S(h)}_{\downarrow} - L_D(h) \right| < \epsilon$$

$$\left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{h(x_i) \neq y_i\} - E_D[\mathbb{1}\{h(x) \neq y\}] \right| < \epsilon$$

what should be $|S| = n_0(\epsilon, \delta)$?

To answer this question we'll use
 a very important probabilistic tool -
 which is call the Hoeffding's
 inequality. This inequality has very
 useful in a variety of applications
 in data science (both in terms of theory

and algorithm design).

Hoeffding's Inequality:

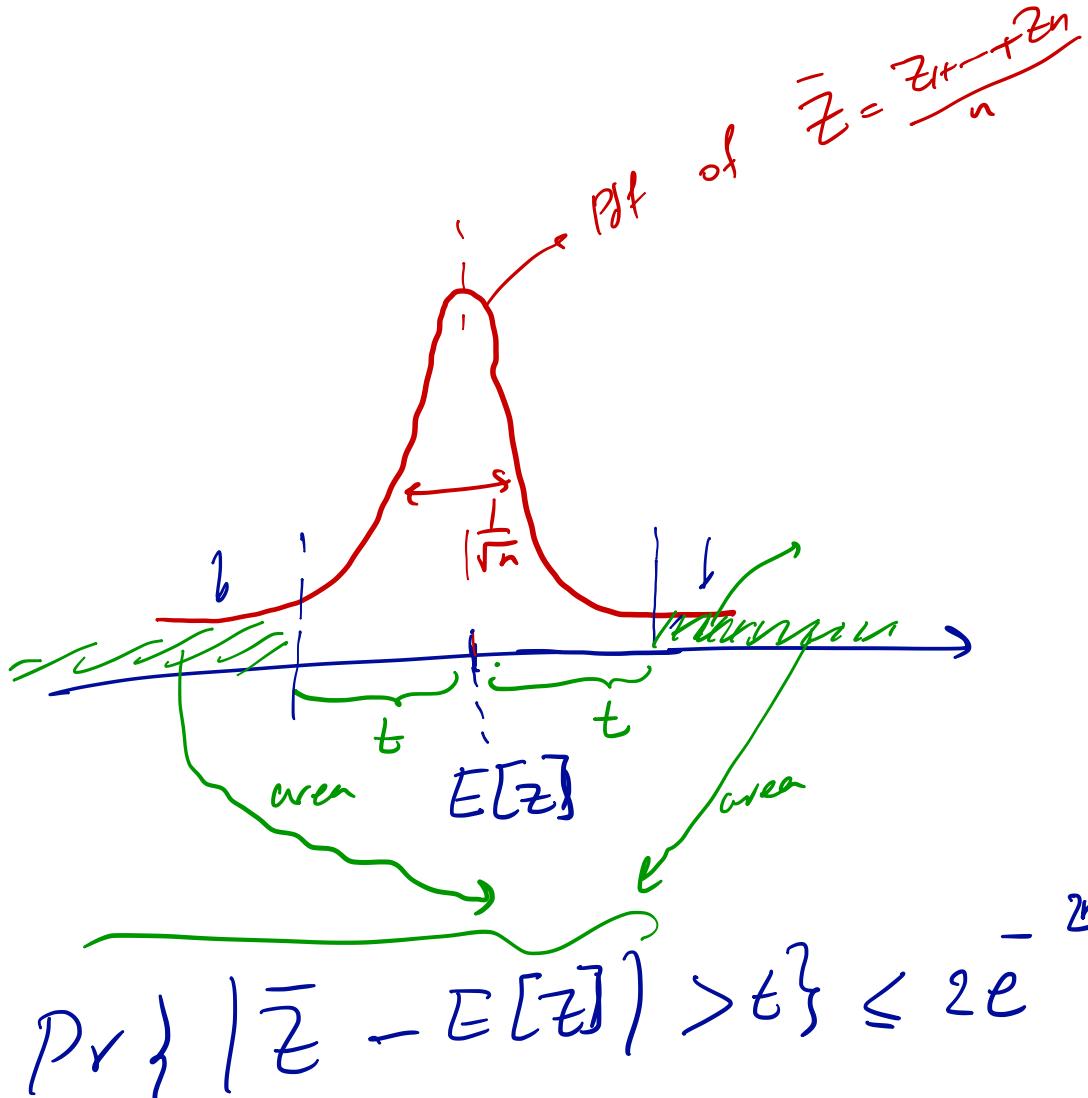
Let Z_1, Z_2, \dots, Z_n be n iid random variables that take value in the unit interval (i.e. $Z_i \in [0, 1]$).

Then:

$$\Pr \left\{ \left| \frac{\sum_{i=1}^n Z_i}{n} - E[Z_i] \right| \geq t \right\} \leq 2e^{-2nt^2}$$

$$Z_1, Z_2, \dots, Z_n \stackrel{iid}{\sim} Z \rightarrow \delta^2 = \text{var}(Z)$$

$$\frac{Z_1 + Z_2 + \dots + Z_n}{n} \sim E[Z] + \frac{1}{\sqrt{n}} N(0, \delta^2)$$



e.g. $n=1000$, $t=0.5$

$$\frac{-2 \cdot 1000 \cdot (0.5)^2}{e} = \frac{-500}{e} \rightarrow 0$$

$H \rightarrow$ PAC learnable:

(1) $n_0(\epsilon, \delta)$

(2) learning algorithm that takes as input a set of training data points S , and outputs h^* .

for any data distribution D , as

long as $|S| \geq n_0(\epsilon, \delta)$ then

with probability $1 - \delta$ we

hence:

$$L_D(h_s^*) \leq \min_{h \in H} L_D(h) + \epsilon$$

Example:

$$\mathcal{H} = \{ h_1, \dots, h_m \}$$

(1) specify what $n_{\epsilon, \delta}$ is.

(2) Algorithm: ERM It:

$$h_S^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} L_S(h)$$

what is $n_{\epsilon, \delta}$, st.

w.p. $1 - \delta$

$$L_D(h_S^*) \leq \min_{h \in \mathcal{H}} L_D(h) + \epsilon$$

Two questions;

(1) Given a fixed function h ,
how many data points, S , do
we need such that

W.P. 1-8

$$\left| L_S(h) - L_D(h) \right| < \varepsilon$$

$$\frac{1}{|S|} \sum_{i=1}^{|S|} \mathbb{1}\{h(x_i) \neq y_i\}$$

\Updownarrow equivalent

$$\Pr \left\{ \left| L_S(h) - L_D(h) \right| \geq \varepsilon \right\} \leq \delta$$



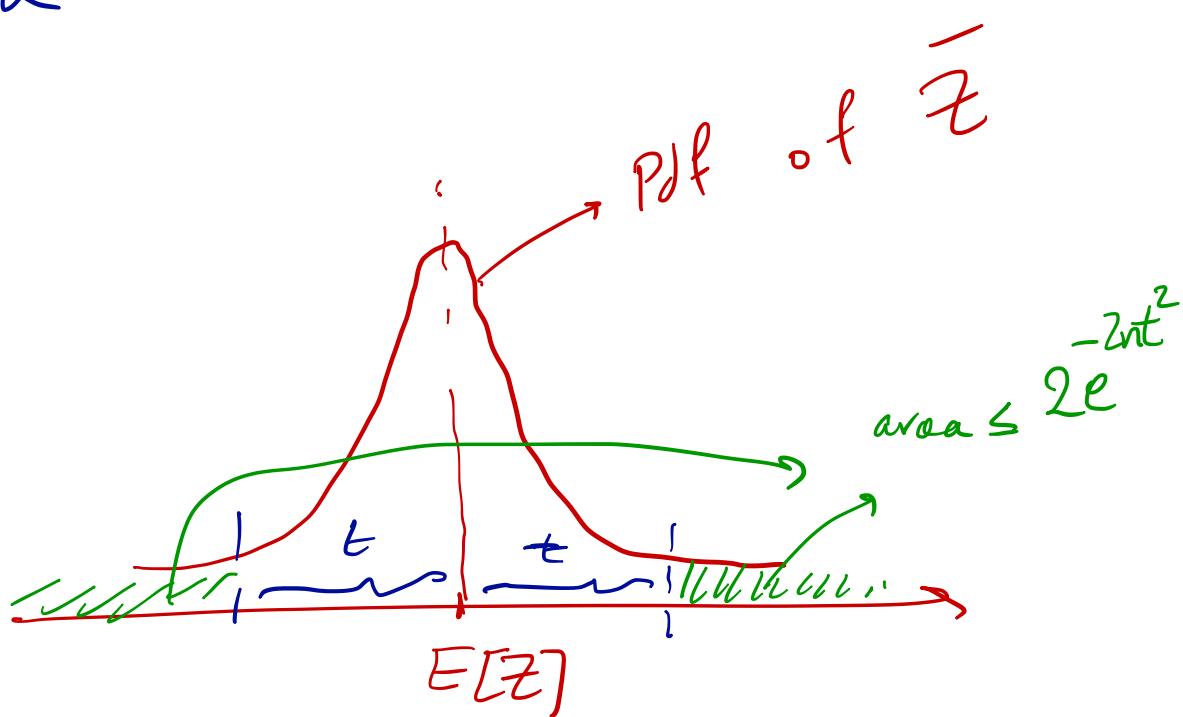
Recall that Hoeffding's inequality was given as follows:

For any iid random variables

Z_1, Z_2, \dots, Z_n , s.t. $Z_i \in [0, 1]$,

We have

$$\Pr \left\{ \left| \overline{\frac{1}{n} \sum_{i=1}^n Z_i} - \bar{E}[Z] \right| > t \right\} \leq 2e^{-2nt^2}$$



$$L_S(h) = \frac{1}{n} \sum_{i=1}^n \underbrace{\mathbb{1}\{h(x_i) \neq y_i\}}_{z_i}$$

$(|S|=n)$

$$E_D(h) = E_{(x,y) \sim D} [\mathbb{1}\{h(x) \neq y\}]$$

$E[z]$

$$L_S(h) = \frac{1}{n} \sum_{i=1}^n z_i = \bar{z}$$

$$z_i = \begin{cases} 1 & \text{if } h(x_i) \neq y_i \\ 0 & \text{o.w.} \end{cases} \rightarrow z_i \in [0,1]$$

$$E[z_i] = E_{(x_i, y_i) \sim D} [\mathbb{1}\{h(x_i) \neq y_i\}]$$

$$= L_D(h)$$

Hoeffding:

$$\Pr \left\{ \left| L_S(h) - L_D(h) \right| > \epsilon \right\} \leq 2e^{-2n\epsilon^2}$$

for any $\epsilon > 0$

(1)

Recall that we're looking for the smallest n s.t.

$$\Pr \left\{ \left| L_S(h) - L_D(h) \right| > \epsilon \right\} \leq \delta$$

(2)

using (1), and in order to guarantee

(2), we can let:

$$2e^{-2n\epsilon^2} \leq \delta$$
$$n \geq \frac{-\log \delta}{2\epsilon^2}$$

$(\log \frac{1}{\delta} = -\log \delta)$

So the final statement is:

Given any ϵ, δ , as long as the number of training data points is larger than

$$n_1 = \frac{1}{2\epsilon^2} \log \frac{2}{\delta}, \text{ we have}$$

~~.....~~

$$\Pr_{(x,g) \in D} \{ |L_S(h) - L_D(h)| > \epsilon \} \leq \delta.$$

| for a fixed function L)

(2) Let's now assume that we have m functions h_1, h_2, \dots, h_m . What is the smallest value $n_0(\epsilon, \delta)$ such that

with probability $1 - \delta$: (3)

$$\forall i \in \{1, \dots, m\} : |L_S(h_i) - L_D(h_i)| < \epsilon.$$

To answer this question, we're going to write an equivalent formulation of relation (3):

Let A_i be the event that

$$|L_S(h_i) - L_D(h_i)| > \epsilon.$$

Then (3) is equivalent to:

$$\Pr \left\{ A_1 \cup A_2 \cup A_3 \cup \dots \cup A_m \right\} \leq \delta. \quad (4)$$

Remember that we are searching for the smallest value of n such that (4) holds.

$$\Pr \{ A_1 \cup A_2 \cup \dots \cup A_m \} \leq \delta$$

To answer this, we're going to use the Union bound:

$$\Pr \{ A \cup B \} \leq \Pr \{ A \} + \Pr \{ B \}$$

$$\Pr \{ A_1 \cup A_2 \cup \dots \cup A_m \} \leq \Pr \{ A_1 \} + \Pr \{ A_2 \} + \dots + \Pr \{ A_m \}$$

bad event # i : A_i :

$$|L_S(h_i) - L_D(h_i)| > \epsilon$$

We'd like to make sure that none of the bad events, A_i , would take place.

So in order to guarantee

$$\Pr \{ A_1 \cup A_2 \cup \dots \cup A_m \} \leq \delta \quad (5)$$

it is sufficient to guarantee that

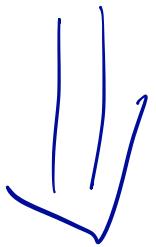
$$\Pr \{ A_1 \} + \Pr \{ A_2 \} + \dots + \Pr \{ A_m \} \leq \delta \quad (6)$$

Note that if (6) holds, then (5) would also hold as a result of the union bound.

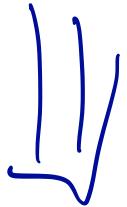
Now, to guarantee (6), it is sufficient to choose n large enough s.t. for every i

We have

$$\Pr \{ A_i \} \leq \frac{\delta}{m} \quad (7)$$



$$\Pr \{ A_1 \} + \Pr \{ A_2 \} + \dots + \Pr \{ A_m \} \leq \delta$$



$$\Pr \{ A_1 \cup A_2 \cup \dots \cup A_m \} \leq \delta$$

Now given our answer to question 1,

in order to guarantee that

$\Pr \{ A_i \} \leq \frac{\delta}{m}$ we need to

choose:

$$\text{if } n \geq n_0(\epsilon, \frac{\delta}{m}) = \underbrace{\frac{1}{2\epsilon^2} \log \frac{2m}{\delta}}_{n_0(\epsilon, \delta)}$$

then

$$\Pr \{ A_i \}$$

$$= \Pr \{ |L_S(h_i) - L_D(h_i)| > \epsilon \} \leq \frac{\delta}{m}$$

Hence,

Statement: If the number of training data points, $|S|$, is larger

$$\text{then } n_0(\epsilon, \delta) = \frac{1}{2\epsilon^2} \log \frac{2m}{\delta}, \text{ then}$$

with probability $1 - \delta$ we have

$$\forall i \in \{1, \dots, m\} : |L_S(h_i) - L_D(h_i)| < \epsilon.$$

What we've shown is that we need

$$\overbrace{\frac{1}{2\epsilon^2} \lg \frac{2m}{\delta}}^{\text{no } (\epsilon, \delta)}$$

data points to
~~guarantee that for all the~~

m functions $h \in \mathcal{H}$ the value
of $L_S(h)$ is close to the

Value of $L_D(h)$:

w.p. $1-\delta$:

$\forall h \in \mathcal{H} : |L_S(h) - L_D(h)| < \epsilon.$

Now, let h_s^* be the minimizer of ERM over the class \mathcal{H} :

$$h_s^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} L_S(h) \quad \left(\begin{array}{l} L_S(h_s^*) \\ \text{is the} \\ \text{smallest} \\ \text{among } \mathcal{H} \end{array} \right)$$

Let $h_{\text{true}}^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} L_D(h)$

We will show that if the number of training data points is at least $N_0(\epsilon, \delta)$ then with $1 - \delta$ we have

$$L_D(h_s^*) - \min_{h \in \mathcal{H}} L_D(h) \leq 2\epsilon. \quad (8)$$

Hence, \mathcal{H} is PAC- ϵ -mable with $n_{\epsilon}(\epsilon, \delta)$ data points.

Let's see why (8) holds.

$$\begin{aligned}
 & L_D(h_S^*) - L_D(h_{\text{true}}^*) \\
 & \leq \underbrace{L_D(h_S^*) - L_S(h_S^*)}_{\leq 0} + \underbrace{L_S(h_S^*) - L_S(h_{\text{true}})}_{\text{h_S^* is the minimizer of } L_S(\cdot)} \\
 & = L_D(h_S^*) - L_S(h_{\text{true}}) + L_S(h_{\text{true}}^*) - L_D(h_{\text{true}}^*) \\
 & \leq \epsilon + \epsilon + \epsilon \leq 2\epsilon.
 \end{aligned}$$

Hence, if the number of training data points is

at least $n_0(\epsilon, \delta) = \frac{1}{2\epsilon^2} \log \frac{2m}{\delta}$

then w.p. $1-\delta$ we have

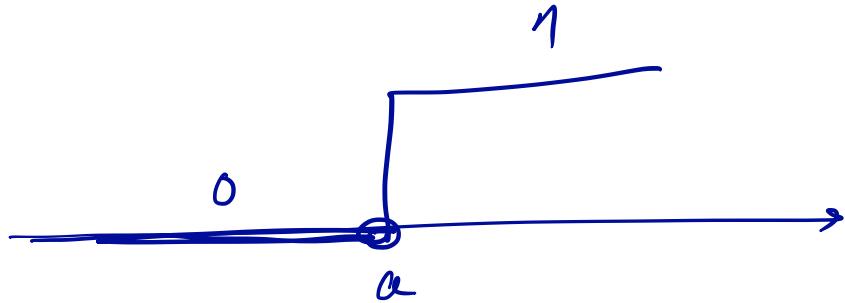
$$0 \leq L_D(h_S^*) - \min_{h \in \mathcal{H}} L_D(h) \leq 2\epsilon.$$

minimizer of

ERM over \mathcal{H} .

When $\mathcal{H} = \{h_1, \dots, h_m\}$.

$$h_a(x) =$$



$$\mathcal{H} = \{ h_a(x), a \in [-\infty, \infty) \}$$

infinitely many functions

One important consequence of PAC-learnability is that it's sufficient to work with a ~~finite~~ finite data set and we don't lose anything in terms of generalization.