

Module 5: An Introduction to the Statistical Learning Theory

We will consider the supervised learning problem, but the framework can also be extended to unsupervised learning.

A formal framework for learning theory:

(1) Data is assumed to be generated according to a distribution $(x, y) \sim D(x, y)$. The input domain is denoted by X ($x \in X$), the label domain is denoted Y (i.e. $y \in Y$), and D is a distribution over $X \times Y$.

(2) The learner's output: we are

looking for predictive relations

$h: X \rightarrow Y$ with low prediction error. Formally speaking, for each function $h: X \rightarrow Y$ we define its error as:

$$\begin{aligned} \text{prediction error} & \quad \Pr_D(h) = \Pr_{(x,y) \sim D} \{ h(x) \neq y \} \\ \text{generalization error} & \\ \text{true error} & \quad \text{Distribution of data} \\ & = E_{(x,y) \sim D} [\prod \{ h(x) \neq y \}] \end{aligned}$$

Ideally, we'd like to find a predictor h that has the

Smallest true error, i.e. our gold standard is to solve

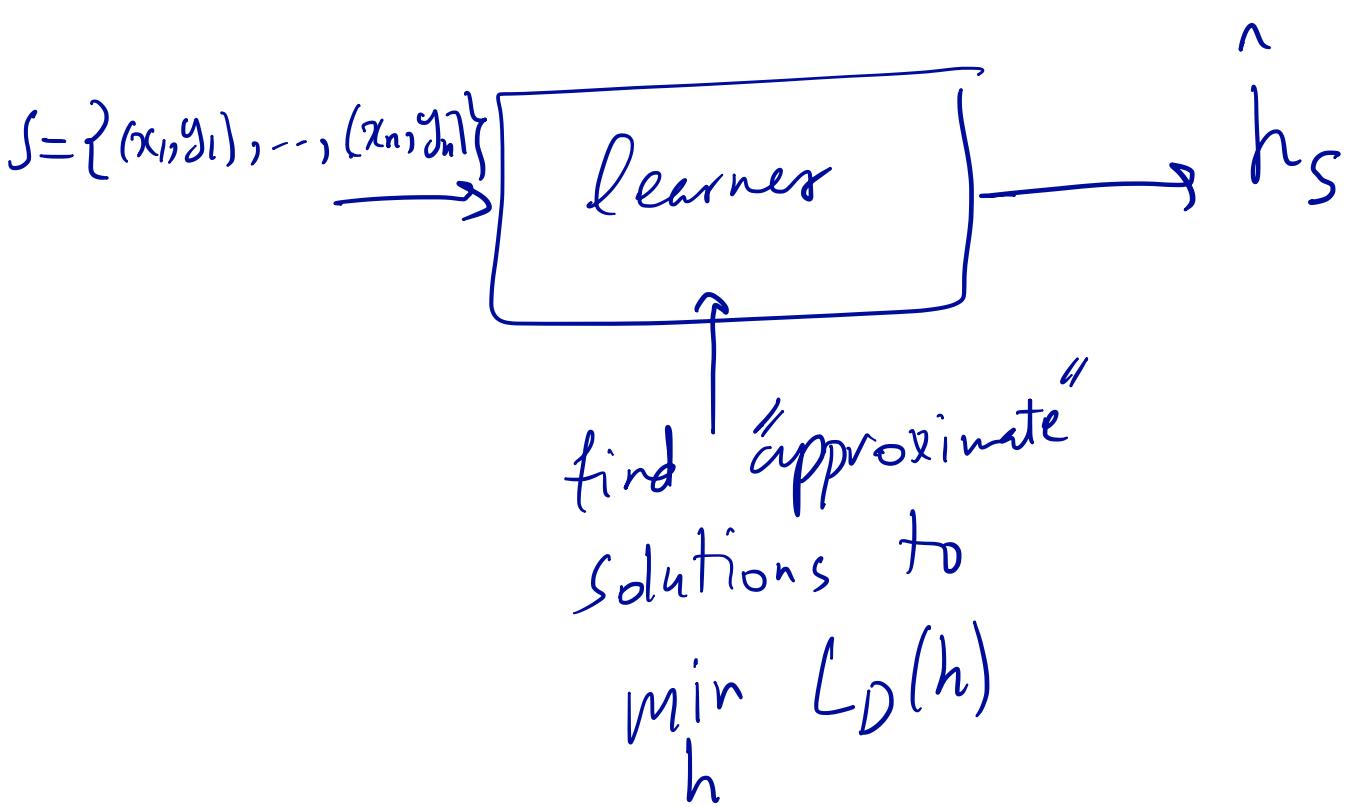
$$\underset{h}{\text{minimize}} \quad L_D(h) \quad (1)$$

$$L_D(h) = E_{(x,y) \sim D} [\#\{h(x) \neq y\}]$$

However, this task is impossible since D is unknown.

(3) the learner's input: Training data! The only information that the learner has about the data distribution is a set of training samples (data points), $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ where each (x_i, y_i) is

Generated i.i.d. according to the distribution D . Hence, the learner has to find "approximate" solutions to problem (1) using the training data set.



$$(x_i, y_i) \stackrel{iid}{\sim} D$$

4) Empirical Risk Minimization (ERM):

$$L_D(h) \longrightarrow \text{find an unbiased estimator}$$

↑
Can't compute
exactly

$$L_D(h) = \mathbb{E}_{(x,y) \sim D} \left[\mathbb{1}_{\{h(x) \neq y\}} \right]$$

unbiased estimate

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

- - - - -

Note: $E_S [L_S(h)] = L_D(h)$

Empirical risk minimization (ERM) is the task of finding a predictor that minimizes the training error:

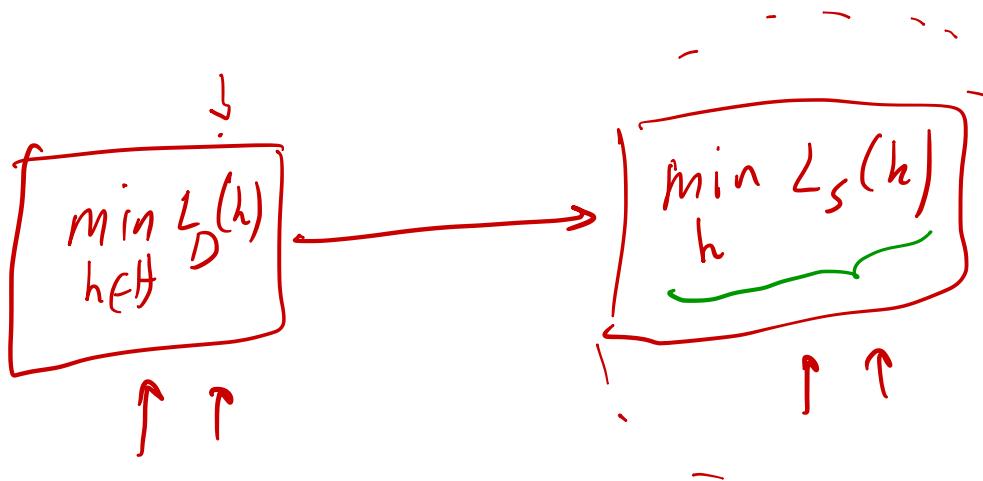
$$\min_h L_S(h) \quad (2)$$

$$\min_h L_D(h) \xrightarrow{\text{instead}} \min_h L_S(h)$$

we can't do

$$\begin{aligned}
 L_D(h) &\xrightarrow{\text{instead}} L_S(h) \\
 &= \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{h(x_i) \neq y_i\}
 \end{aligned}$$

Lecture 23:



Claim:

$$\min_h L_s(h) = 0$$

$$L_s(h) = \frac{1}{n} \sum_{i=1}^n \#\{h(x_i) \neq y_i\}$$

Define $\tilde{h} : X \rightarrow Y$ as follows:

$$\tilde{h}(x) = \begin{cases} y_i & \text{if } x = x_i \\ 0 & \text{if } x \in \{x_1, \dots, x_n\} \end{cases}$$

It's easy to see that

$$L_S(\tilde{h}) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{\tilde{h}(x_i) \neq y_i\} = 0$$

$$L_S(\tilde{h}) = 0 \rightarrow$$

$$\min_h L_S(h) = 0$$

5) Although ERM is natural,
it can fail miserably if
we are not careful: It
can "overfit" easily.

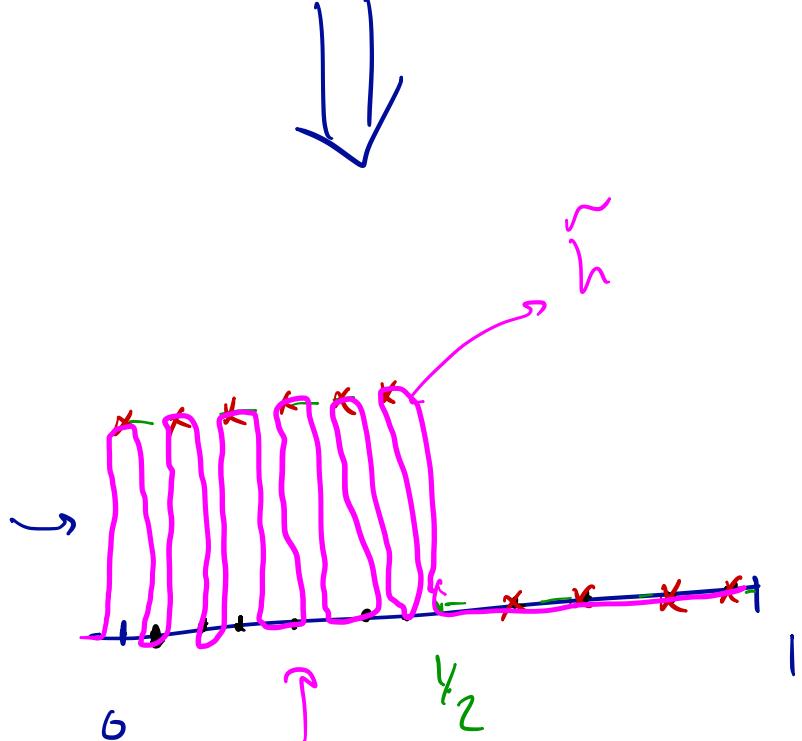
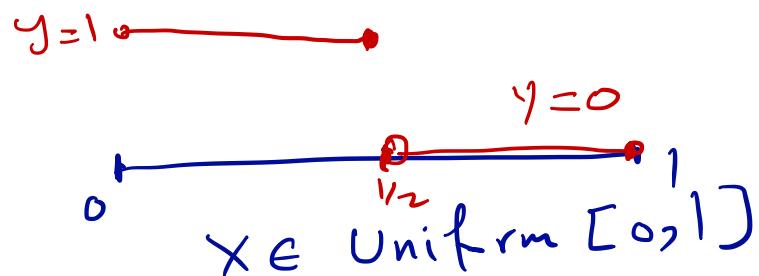
The minimizer of

$$\min_h L_s(h) \quad (2)$$

is always zero.

Example:

Distribution of data:



Prediction
on unseen
data here
is very bad

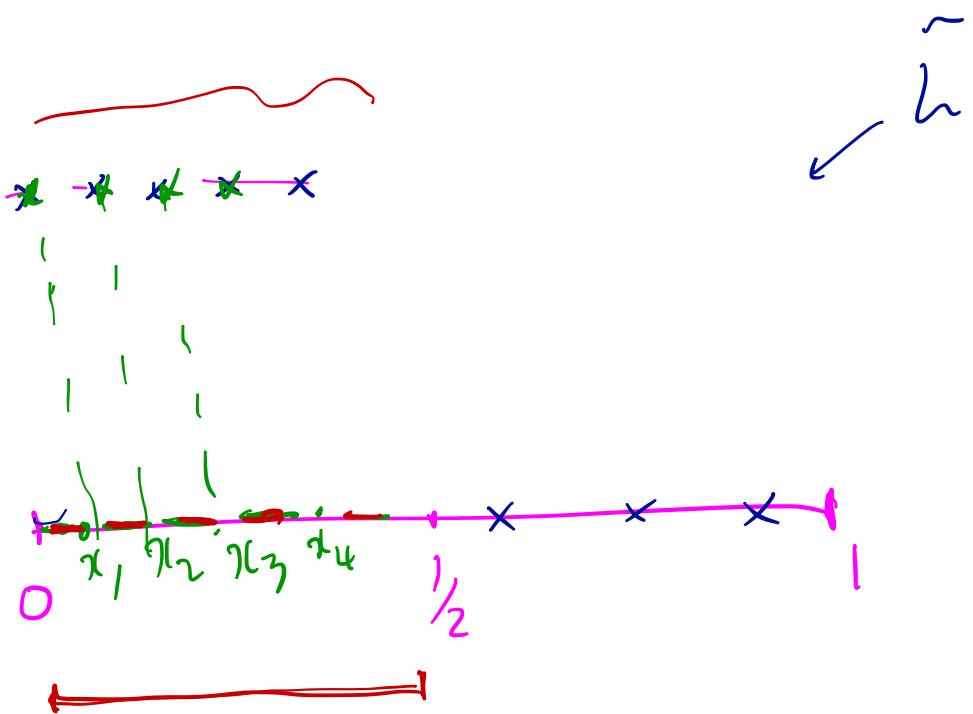
$$\min_h L_S(h)$$

$$L_S(\hat{h}) = 0$$

$$L_D(\hat{h}) \approx 0.5$$

$$\min_h L_D(h)$$

$$\min_h L_S(h) \quad f \hat{h}$$



$$\tilde{h} = \begin{cases} h(x_i) = y_i \\ h(x) = 0 \end{cases}$$

$x \in \{x_1, \dots, x_n\}$

$$\left\{ \begin{array}{l} \zeta_s(\tilde{h}) = 0 \\ \neg D^{**} = \bigcup_{(x,y) \sim D} \neg [h(x) \neq y] \end{array} \right.$$

$$\equiv \not\models$$

6) we need to search for conditions under which there is a guarantee that ERM does not overfit, namely conditions ~~under~~ which when the ERM predictor has a good performance with respect to the training data, it is also likely to perform well over the underlying (true) distribution.

$$\min_h L_S(h)$$

$$\min_{h \in H} L_S(h)$$

restrict the class of functions (complexity)

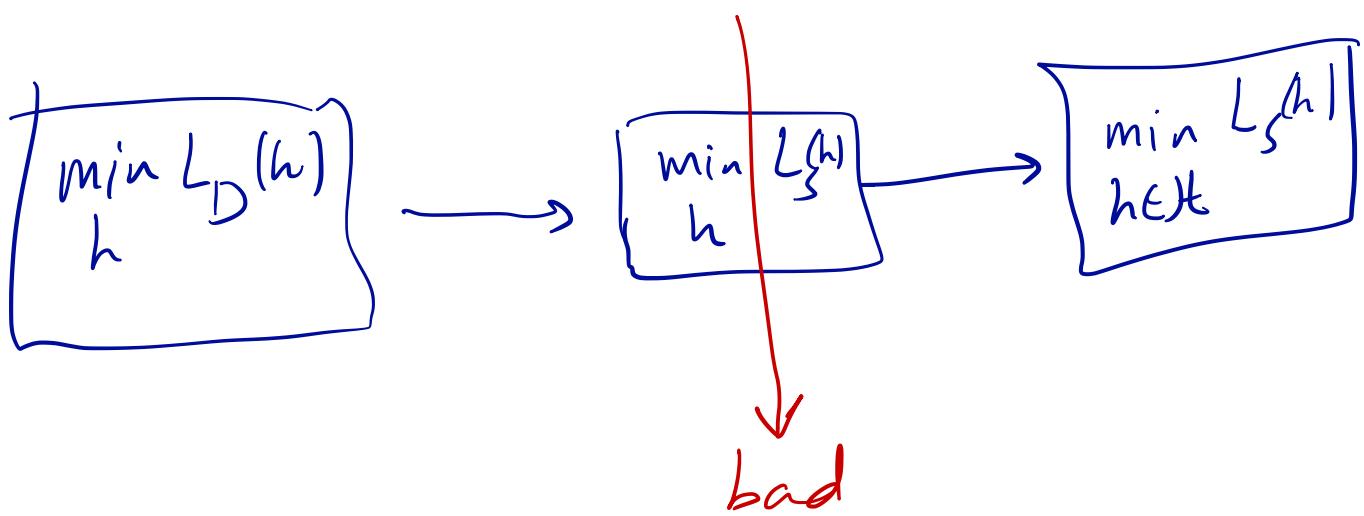


↓ prevent overfitting

A common solution is to apply ERM over a restricted set of functions (classifiers). Formally the learner should choose in advance (before seeing the data) a set of functions (classifiers). This set is called "the predictor class" or "the function class" or "the hypothesis class" and is denoted by H . Each $h \in H$ is a function from X to Y .

And the restricted ERM problem ($ERM_{\mathcal{H}}$) is :

minimize $L_S(h)$ -
 $h \in \mathcal{H}$



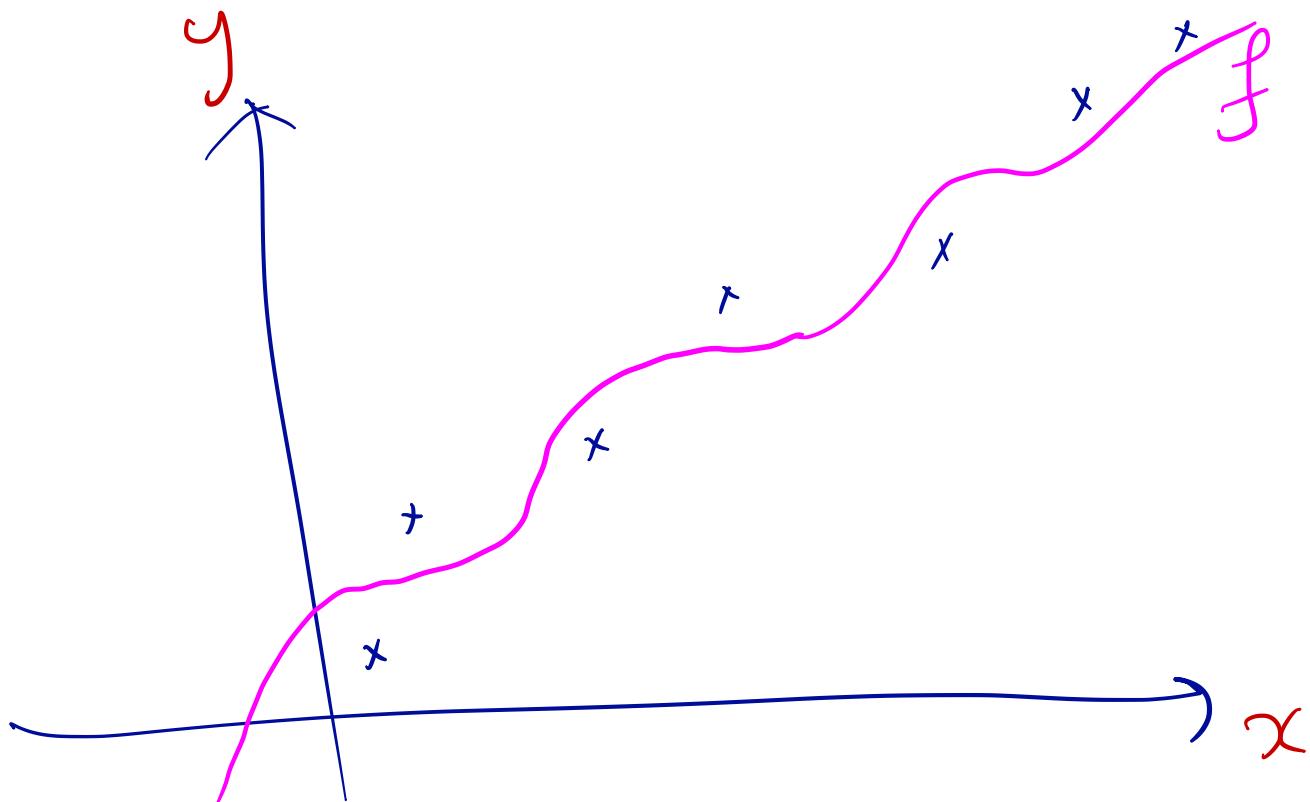
So we are biasing the learner towards a particular set of predictors.

To explain things further, let's look at an example about regression (even though it is not classification, it will help us understand the concepts better).

$$y = f(x) + \epsilon \xrightarrow{\text{independent gaussian}} N(0, \sigma^2)$$

We don't know f and we'd like to estimate the predictive relation between x and y using training data.

$$S = \left\{ (x_i, y_i) \right\}_{i=1}^n$$



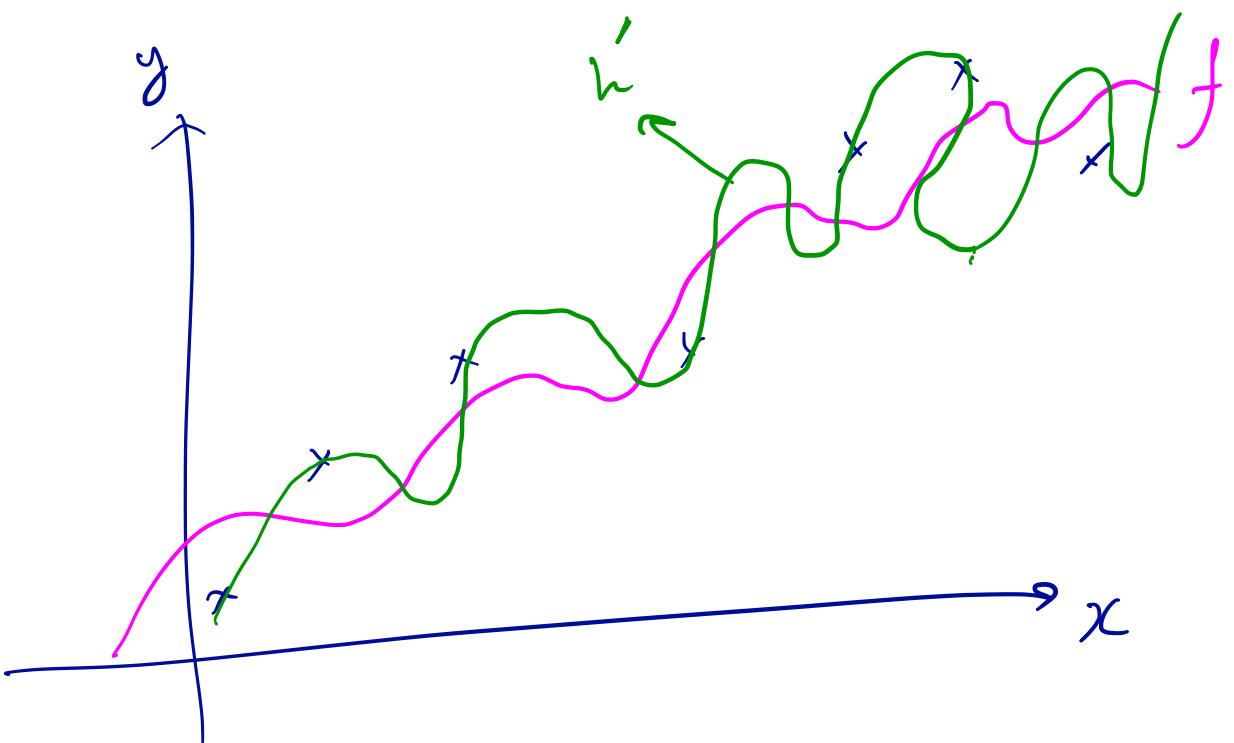
for any $h: \mathcal{X} \rightarrow \mathcal{Y}$

$$L_D(h) = E_{(x,y) \sim D} [(h(x) - y)^2]$$

$$\min_h L_D(h) = \sigma^2$$

↳ minimizer $\rightarrow f$

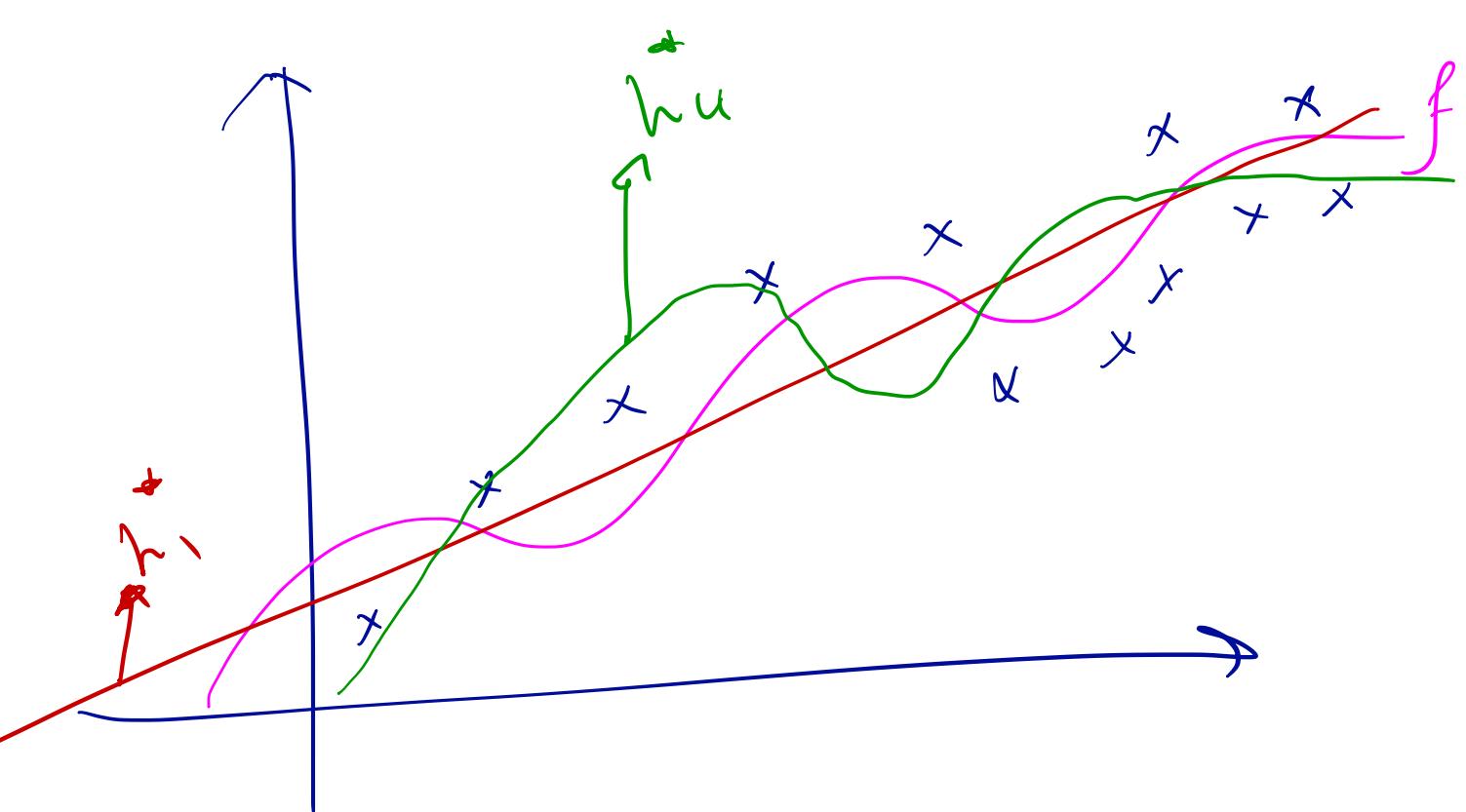
$$L_D(f) = E[(y - f(x))^2] = \sigma^2$$



$$\left\{ \begin{array}{l} \min_h L_D(h) = b^2 \\ \min_h L_S(h) = 0 \rightarrow L_S(h') = 0 \end{array} \right.$$

Let's now consider restricted
 \mathcal{E}_{R^n} objectives:

$$\min_{h \in \mathcal{E}} L_S(h)$$



$\text{ERM}_t \rightarrow \min_{h \in \mathcal{H}} L_s(h)$

e.g. when $\mathcal{H}_1 = \{ \text{polynomials of degree 1} \}$

$\hookrightarrow \min_{h \in \mathcal{H}_1} L_s(h) \rightarrow h_1^*$

$\mathcal{H}_4 = \{ \text{polynomials of degree 4} \}$

$\min_{h \in \mathcal{H}_4} L_s(h) \rightarrow h_4^*$

These are two important points:

(1) since we are restricting our function class, our best bet would be to achieve

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

which is larger than

$$\min_{h \in \mathcal{H}} L_D(h) \geq \underbrace{\min_{h \in \mathcal{H}} L_D(h)}_{\text{---}}$$

$$\min L_D(h) \neq 0$$

Hence, if the class \mathcal{H} is not rich enough we

may have (underfitting)

$$\min_{h \in \mathcal{H}} L_D(h) \gg \min_h L_D(h)$$

$\Rightarrow \mathcal{H}$ should be sufficiently complex

(2) The minimizer of

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

is not the same as

the minimizer of

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

In other words if we denote

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

↳ this is computable

then it is possible that

$$L_D(h^*) \gg \min_{h \in \mathcal{H}} L_D(h)$$

↑
overfitting

To avoid over fitting in \mathcal{H}

We should make sure
that we have sufficiently
many training data points.

h_S^* → minimizer of $L_S(h)$ within the class \mathcal{H} .

Since $h \in \mathcal{H} \Rightarrow L_D(h) \geq \min_{h \in \mathcal{H}} L_D(h)$

$$L_D(h^*) \gg \min_{h \in \mathcal{H}} L_D(h)$$

↑
Overfitting

two potential Potential problems:

$$(1) \quad \min_{h \in \mathcal{H}} L_D(h) \gg \min_{h \in \mathcal{H}} L_D(h)$$

↑
underfitting

we need
to increase
the complexity of \mathcal{H}

solution of ERM $_{\mathcal{H}}$

$$(2) \quad L_D(h_S^*) \gg \min_{h \in \mathcal{H}} L_D(h)$$

↑
overfitting

we need to
increase the # of training
data points

Lecture 24:

$$\min_h L_D(h) \rightarrow \min_h L_S(h)$$

$$\tilde{h} = \begin{cases} y_i & x = x_i \\ 0 & x \notin \{x_1, \dots, x_n\} \end{cases}$$

$$L_S(\tilde{h}) = 0$$

$$L_D(\tilde{h}) \text{ very large}$$

— — — —

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

$$h \in \mathcal{H}$$

1. Since we're restricting our search to \mathcal{H}

$$\min_{h \in \mathcal{H}} L_D(h) \geq \min_h L_D(h)$$

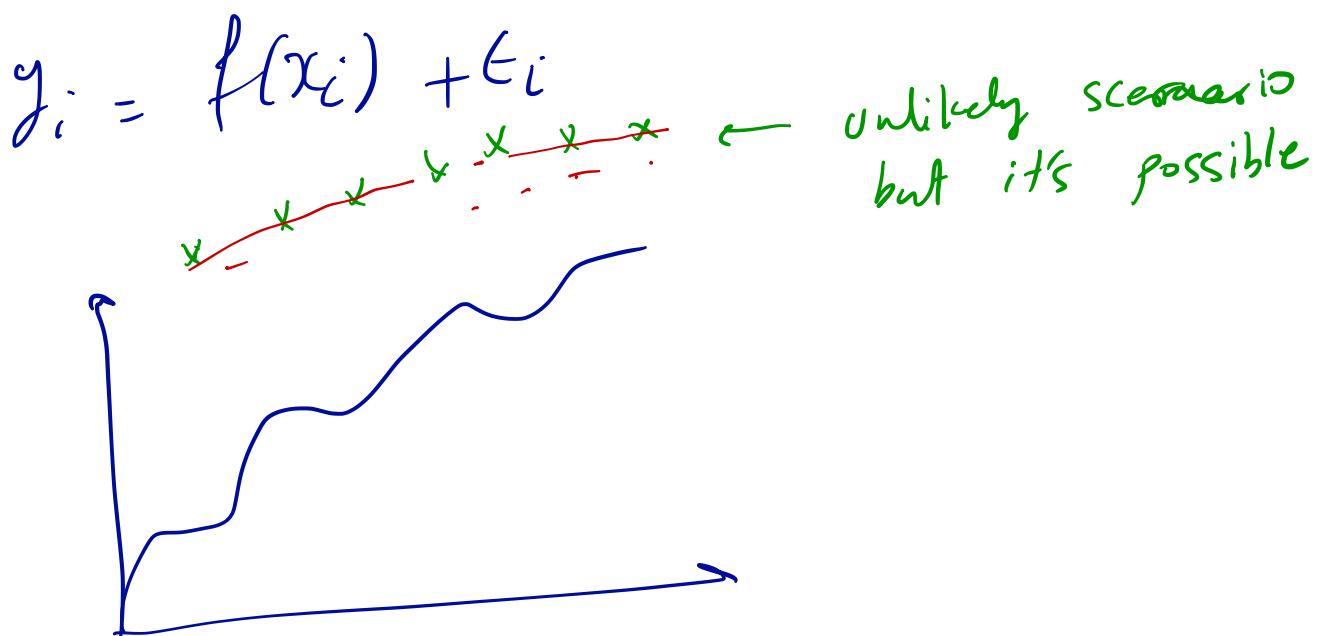
2. what we're actually optimizing is

$$\min_{h \in \mathcal{H}} L_S(h) \Rightarrow \text{Let's denote the minimizer by } h^*_S.$$

$$L_D(h^*) \geq \min_{h \in \mathcal{H}} L_D(h)$$

↳ if $L_D(h^*)$ is much larger than $\min_{h \in \mathcal{H}} L_D(h)$, then this is called overfitting.

3. All the guarantees and bounds are probabilistic (they'll hold with high probability excluding the unlikely scenarios).



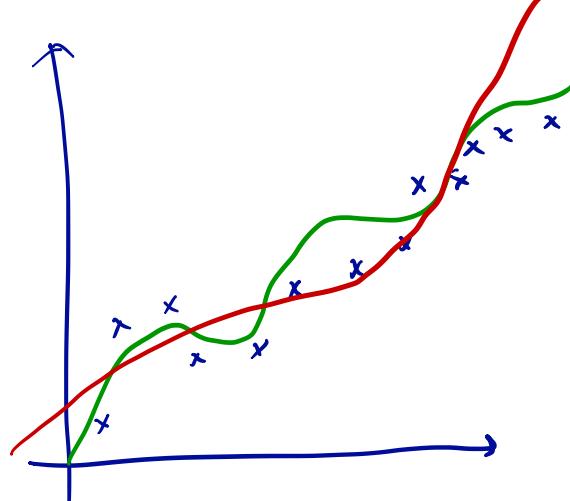
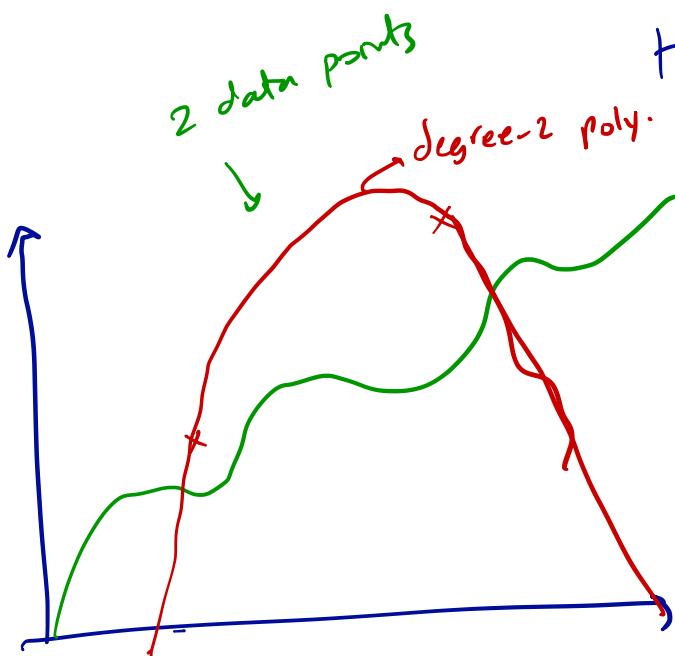
$$L_D \rightarrow E_D$$

$$L_S \rightarrow \frac{1}{n} \sum_i$$

$$L_S \rightarrow L_D$$

$$\min_{h \in H} L_S(h)$$

in order to ensure that overfitting does not take place, we need to have a sufficient number of training data points.



$$y_i = f(x_i) + \epsilon_i$$

$$H = \{ \text{polynomials of degree } 3 \}$$

\Rightarrow overfitting can be avoided provided that we have sufficiently many training data points.

Hence, the main question is:

How many training data points do we need to make sure that overfitting does not happen?

↪ # of data points that we need will depends on:

$$H, \epsilon, \delta$$

What does "avoiding overfitting" mean?

$$h_s^* = \underset{h \in H}{\operatorname{argmin}} L_s(h)$$

what we can do computationally

$$\min_{h \in H} L_D(h)$$



the best error that we can hope for.

overfitting is avoided when: $L_D(h_s^*)$ is close to $\min_{h \in H} L_D(h)$

Overfitting is avoided when

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

Small value
↓

With probability $1-\delta$ we have:

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

provided that
sufficiently
many data
points are
available.

if we have more than no data points, then with high probability the following is true: with probability at least $1-\delta$

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

Now depends on $\epsilon, \delta, \mathcal{H}$.

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:

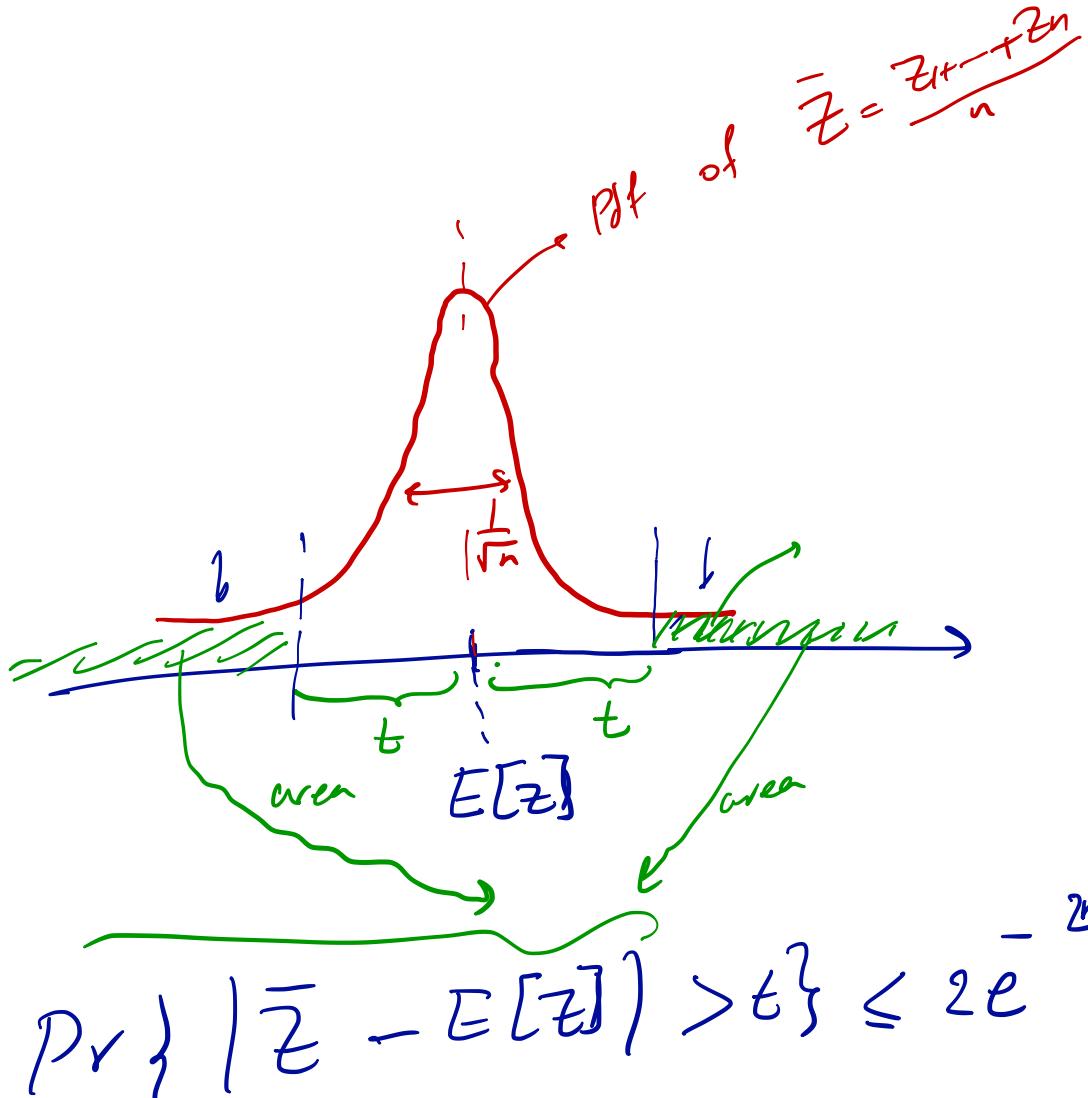
Let Z_1, Z_2, \dots, Z_n be n iid random variables that take value in the unit interval (i.e. $Z_i \in [0, 1]$).

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

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

Lecture 25:

$H \rightarrow$ 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

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

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

↑ equivalent

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



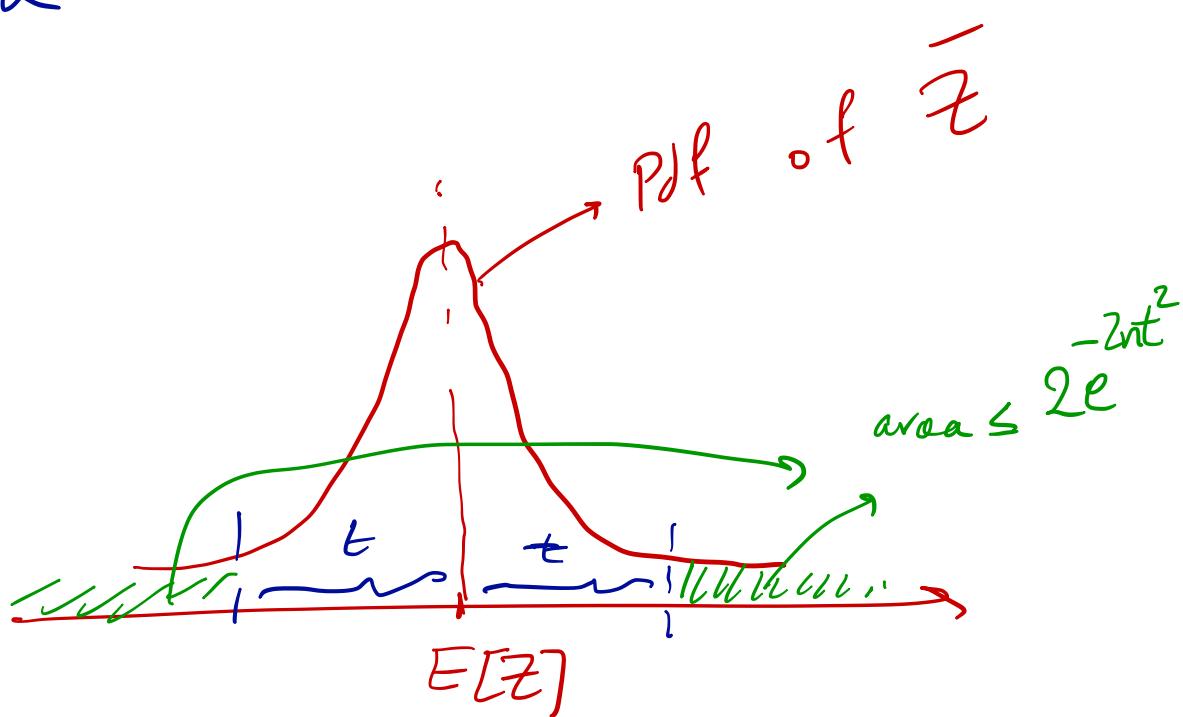
Recall that Hoeffding's inequality was given as follows:

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:

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

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

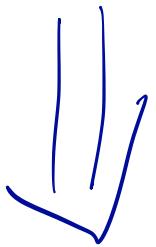
$$\Pr \{ A_1 \} + \Pr \{ A_2 \} + \dots + \Pr \{ A_m \} \leq \delta \quad (6)$$

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

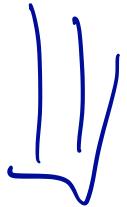
Now, to guarantee (6), it is sufficient to choose n large enough s.t. for every i

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

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

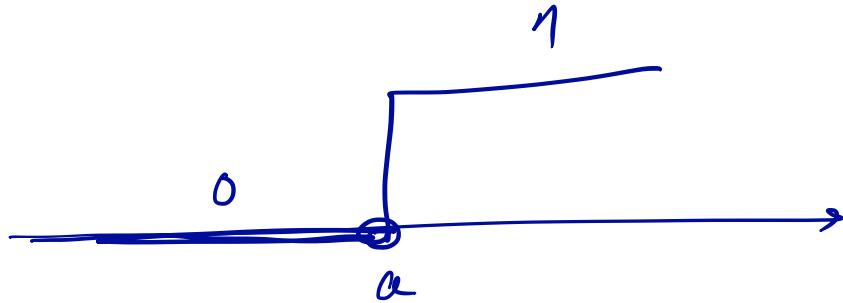
$$0 \leq L_D(h_S^*) - \min_{h \in \mathcal{H}} L_D(h) \leq 2\epsilon.$$

minimizer of

ERM over \mathcal{H} .

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

$$h_a(x) =$$



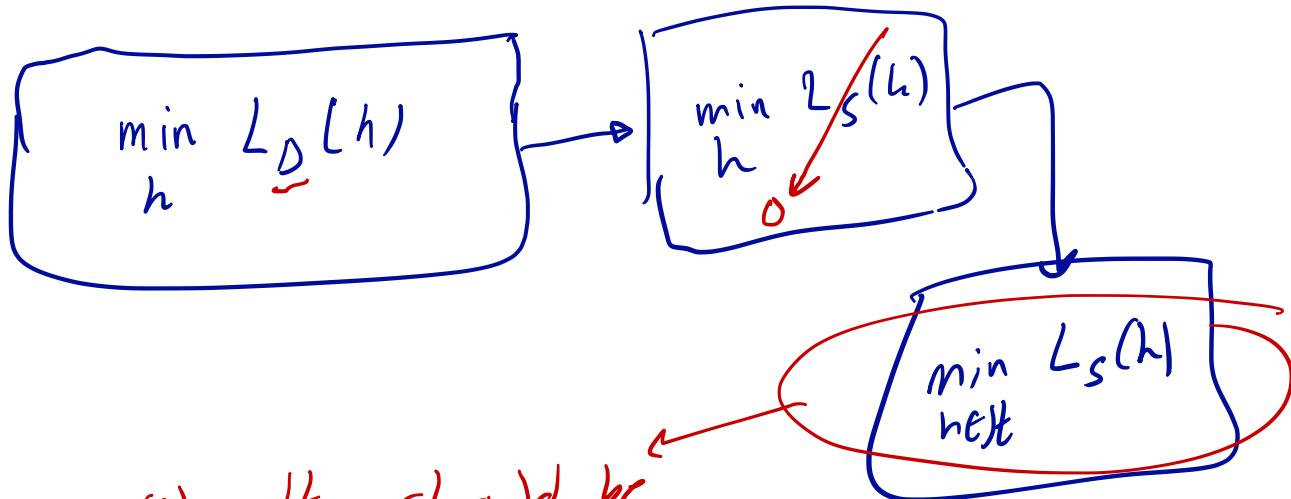
$$\mathcal{H} = \{ h_a(x), a \in [-\infty, \infty) \}$$

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.

Lecture 26:

- $(x, y) \sim D$ \rightarrow



(1) It should be sufficiently simple

(2) Avoid overfitting = have sufficiently many data points

PAC: H is PAC learnable:

(1) No $(\epsilon, \delta) \in \mathbb{N}$ training data

(2) learning alg. $S \rightarrow h_S^*$

If $|S| > n_{\epsilon, \delta}$:

w.p. $1 - \delta$:

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

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

$$- n_0(\epsilon, \delta) = \frac{1}{\epsilon^2} \log \frac{m}{\delta} \approx \frac{\log \frac{1}{\delta} + \overbrace{\log m}^{\text{up to constants}}}{\epsilon^2}$$

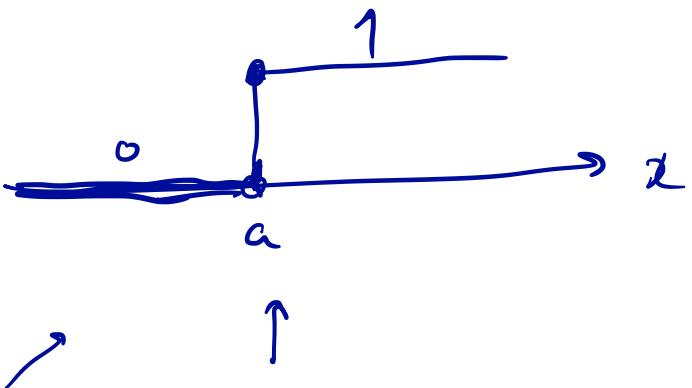
$$- \text{learning algorithm: } h_s^* = \operatorname{argmin}_{i=1, \dots, m} L_S(h_i)$$

\mathcal{H} is PAC-learnable with $n_0(\epsilon, \delta)$ samples and ERM-H as the learning algorithm.

what about infinite \mathcal{H} ?

$$\mathcal{H} = \{ h_a(x) : \mathbb{R} \rightarrow \{0, 1\}, a \in \mathbb{R} \}$$

$$h_a(x) \rightarrow$$



The fundamental theorem of learning theory:

For any function class H , if the number of training data point is larger than

$$n_0(\epsilon, \delta) = \frac{\log \frac{1}{\delta} + \sqrt{VC\text{-dim}(H)t}}{\epsilon^2}, \text{ then}$$

w.p. 1- δ :

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

where $h_s^* = \underset{h \in H}{\operatorname{argmin}} L_S(h)$.

→ $VC\text{-dim}(H)$ = Complexity of the class

→ if $VC\text{-dim}(H)$ is finite, then

It is PAC-learnable with

$$n_0(\epsilon, \delta) = \frac{\log \frac{1}{\delta} + VC\text{-dim}(H)}{\epsilon^2} \text{ and}$$

Term as the learning algorithm.

Important implication:

We can learn from finite data

Sets -

the best that we could do

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

intractable

As long as we have sufficiently many data points we'll be fine with solving the "tractable" ERM-H problem and we lose at most a small value ϵ with respect to the optimal accuracy.

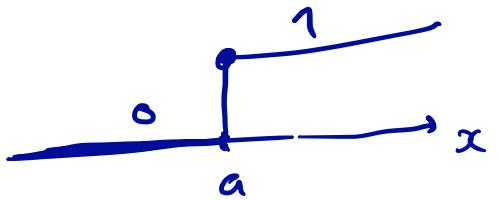
VC dimension:

1

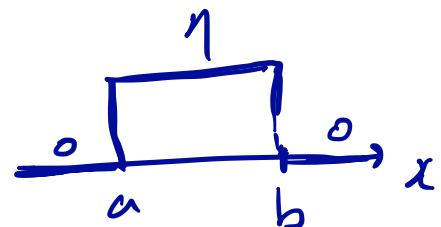
Vapnik-Chernovskii

Consider the following function classes:

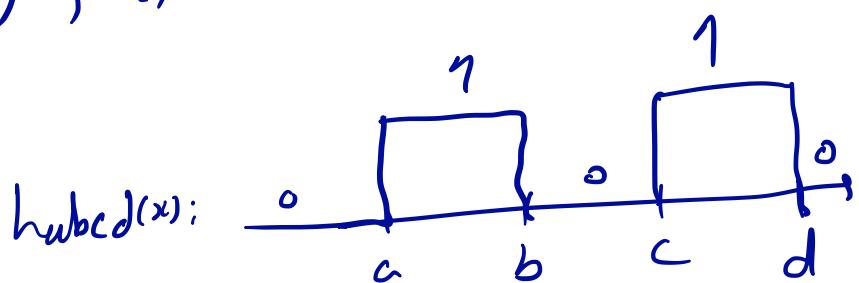
$$H_1 = \{ h_a(x), a \in \mathbb{R} \} \rightarrow$$



$$H_2 = \{ h_{ab}(x), a, b \in \mathbb{R} \} \rightarrow$$

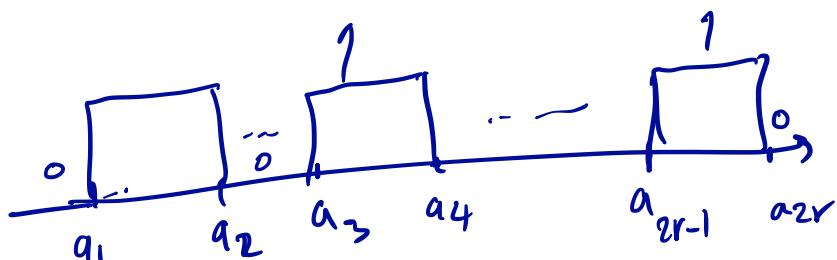


$$H_3 = \{ h_{abcd}(x), a, b, c, d \in \mathbb{R} \}$$



$$H_4 = \{ h_{a_1, a_2, \dots, a_{2r}}(x), a_1, \dots, a_{2r} \in \mathbb{R} \}$$

$$h_{a_1, a_2, \dots, a_{2r}}(x):$$



$H_1 < H_2 < H_3 < H_4 \rightarrow$ complexity

Definition: (restriction):
Let \mathcal{H} be a class of functions from
 X to $\{0, 1\}$ ($\forall h \in \mathcal{H}, h: X \rightarrow \{0, 1\}$).

Let $C = \{x_1, x_2, \dots, x_k\} \subseteq X$.

The restriction of \mathcal{H} to C is
the set of all the k -tuples

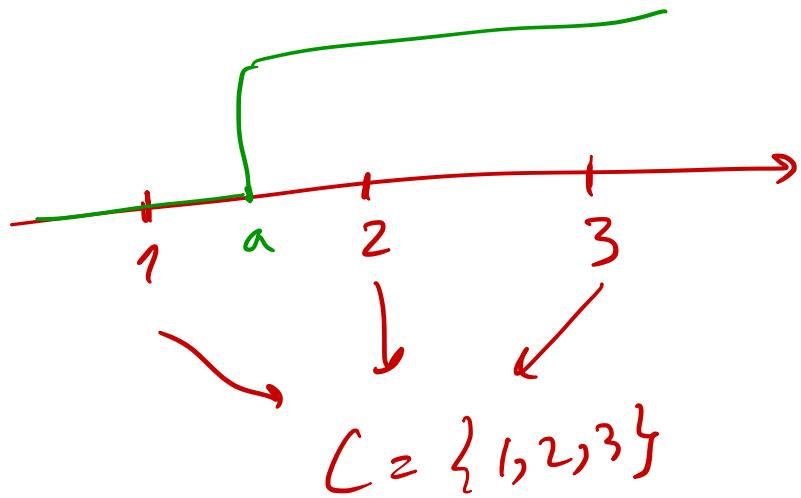
$$\mathcal{H}_C = \{ (h(x_1), h(x_2), \dots, h(x_k)) , h \in \mathcal{H} \}.$$

Example: e.g. if $\mathcal{H} = \{ h_a(x) , a \in \mathbb{R} \}$

$$h_a(x) = \begin{cases} 1 & \text{if } x > a \\ 0 & \text{otherwise} \end{cases}$$

Let $C = \{1, 2, 3\}$.

$$\mathcal{H}_C = \{ (h_a(1), h_a(2), h_a(3)) , \begin{matrix} h_a \in \mathcal{H} \\ a \in \mathbb{R} \end{matrix} \}$$



$$a < 1 : (h_a(1), h_a(2), h_a(3)) = (1, 1, 1)$$

$$a \in [1, 2) : \quad " \quad = (0, 1, 1)$$

$$a \in [2, 3) \quad " \quad = (0, 0, 1)$$

$$a > 3 \quad " \quad = (0, 0, 0)$$

$$\mathcal{H}_C = \{(1, 1, 1), (0, 1, 1), (0, 0, 1), (0, 0, 0)\}$$

$$|\mathcal{H}_C| = 4 < 2^{|C|} = 2^3 = 8$$

$\rightarrow \mathcal{H}$ is not sufficiently complex to produce all the possible binary 3-tuples on C .

Note :

If $C = \{x_1, \dots, x_k\}$

$$|\mathcal{H}_C| \leq 2^{|C|}$$

$$\mathcal{H}_C = \{(0, 0, \dots, 0), (1, 0, \dots, 0), (0, 1, 0, \dots, 0), \dots, (1, 1, \dots, 1)\}$$

2^k

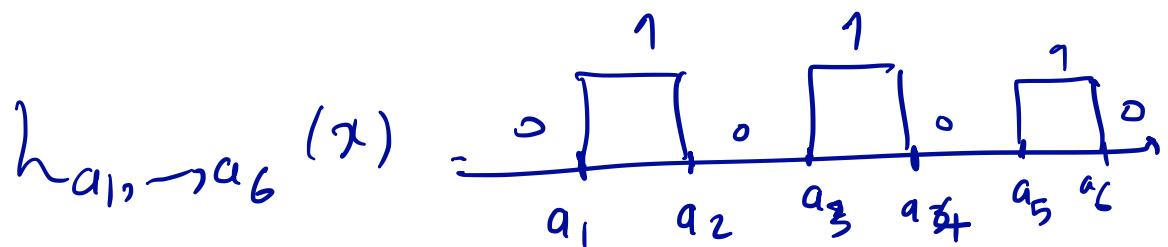
Definition (Shattering): We say that

a function class \mathcal{H} shatters a

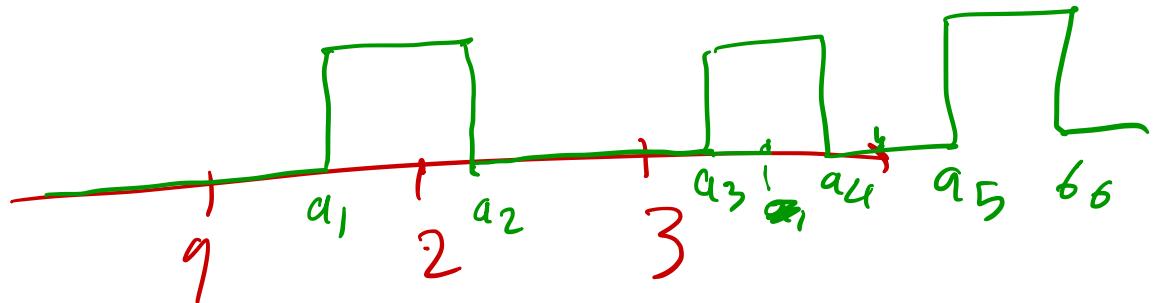
Set C if $|\mathcal{H}_C| = 2^{|C|}$

Example:

$$\mathcal{H} = \left\{ h_{a_1, a_2, a_3, a_4, a_5, a_6}(x) , a_1, \dots, a_6 \in \mathbb{R} \right\}$$



$$C = \{1, 2, 3\}$$



$$\mathcal{H}_C = \{ (0, 0, 0), (1, 0, 0), (0, 1, 0), (0, 0, 1), -\tau^{(1, 1, 1)} \}$$

\mathcal{H}_C is all the 8 possible binary 3-tuples on C.
It shatters C.

Definition (VC-dimension):

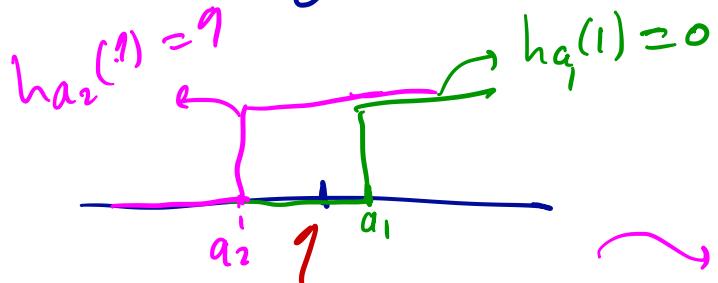
The VC-dimension of a function class \mathcal{H} is the largest number k such that exists a set C of size k which is shattered by \mathcal{H} .

Examples:

$$\mathcal{H} = \{ h_a(x) , a \in \mathbb{R} \} \rightarrow \begin{array}{c} \nearrow \\ \text{---} \\ \text{---} \end{array} \quad \begin{array}{c} \nearrow \\ \text{---} \\ \text{---} \\ \nearrow \\ a \end{array}$$

VC-dim(\mathcal{H}) =

Is there a set of size $?$ that's shattered by \mathcal{H} ? YES.

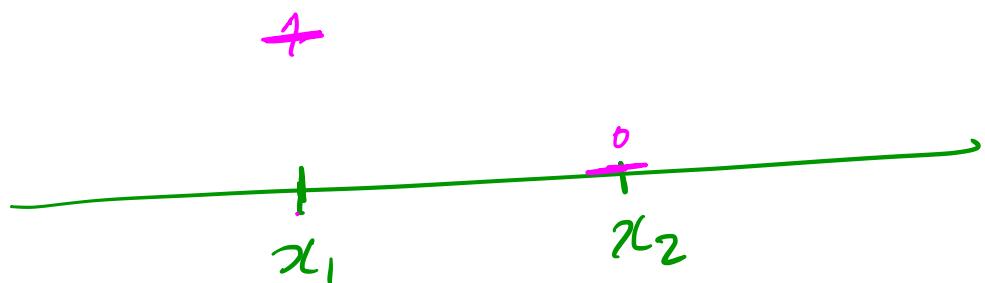


$$C = \{1\}$$

$$\mathcal{H}_C = \{0, 1\}$$

$\hookleftarrow \mathcal{H}$ shatters C .

Is there a set of size 2 which
is shattered by \mathcal{H} ? No

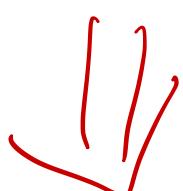


$$C = \{x_1, x_2\} \quad (\text{assume } x_1 < x_2)$$

$\hookrightarrow \mathcal{H}_C = \{(\circ, \emptyset), (\underline{1}, 0), (0, 1), (1, 1)\}$

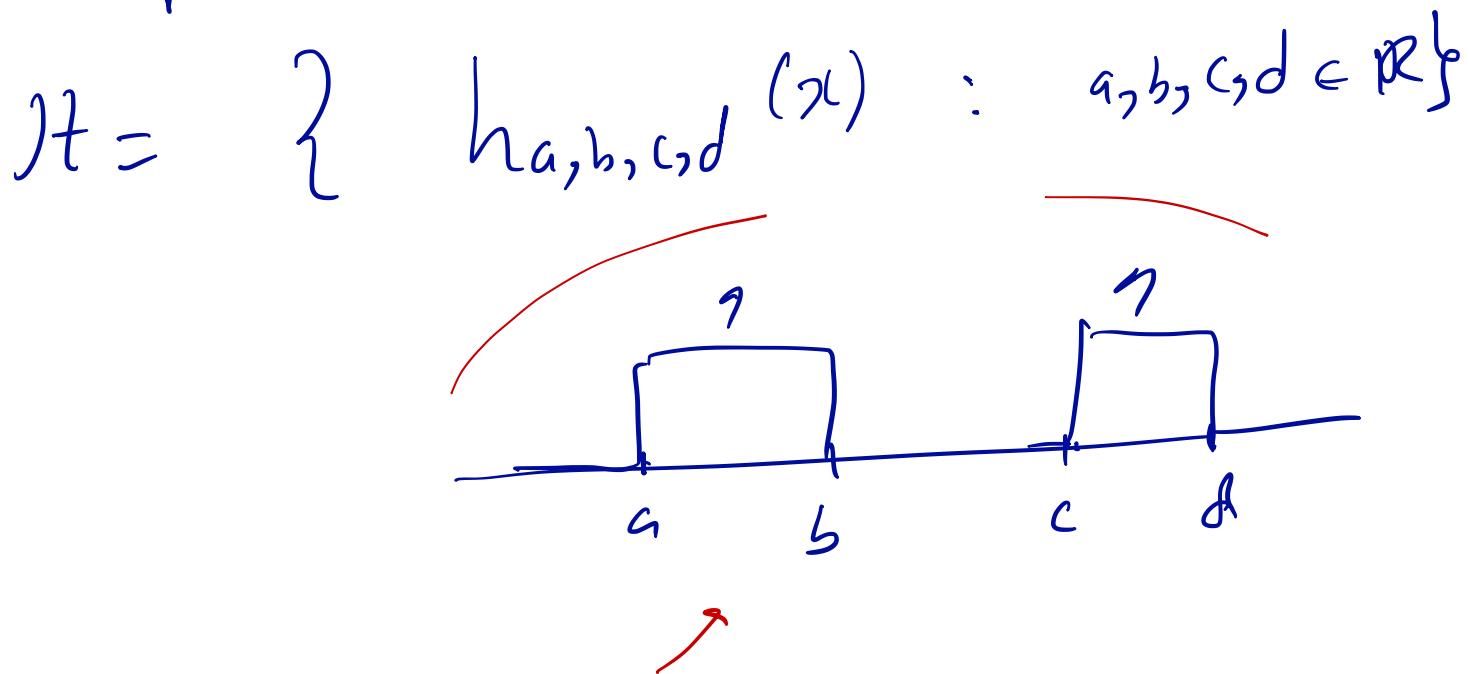
\hookrightarrow there is no $h_a(x)$
that produces this
binary 2-tuple

$$\check{VC_dim(\mathcal{H}) = }$$



there exists no set C
of size 2 which is shattered
by \mathcal{H}

Example:



is there a set C of size k
which is shattered by \mathcal{H} ?

$k=1 \rightarrow \text{YES}$

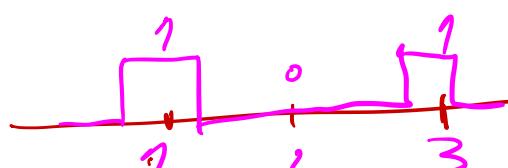
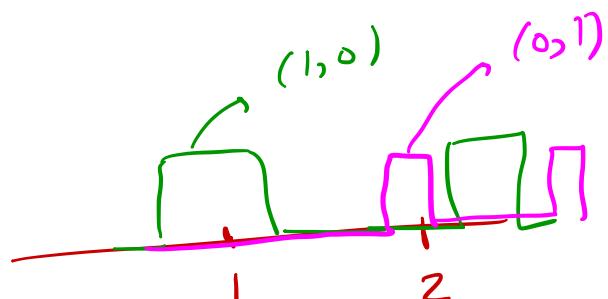
$k=2 \rightarrow \text{YES}$

$k=3 \rightarrow \text{YES}$

$k=4 \rightarrow \text{YES}$

$k=5 \rightarrow \text{NO}$

$$\mathcal{H}_C = \left\{ (0,0,0), (1,0,0), (0,1,0), (0,0,1), (1,0,1), (1,1,0), (0,1,1), (1,1,1) \right\}$$



$k=5$

→ h should have 3 jumps from 0 to 9.

