

Lecture 21 :

Module 5: An Introduction to Statistical Learning Theory:

We consider the problem of supervised learning in this module, but all the concepts generalize to other forms of learning (e.g. unsupervised learning).

A formal framework for learning theory:

(1) Data is assumed to be generated

iid according to a distribution

$(x, y) \sim P(x, y)$. The input domain is denoted by X ($x \in X$, e.g. $X = \mathbb{R}^p$), the label domain is denoted by Y (i.e. $y \in Y$, e.g. $Y \in \{0, 1\}$), and

D is a distribution over $X \times Y$.

2: The learner output, the learner's

goal is to find a predictive relation

$h: X \rightarrow Y$ which has low error /

high accuracy. Formally speaking, for

any function $h: X \rightarrow Y$ we define

its error as:

error / loss \swarrow \searrow a function from X to Y

$$L_D(h) = \Pr \{ h(x) \neq y \mid (x,y) \sim D \}$$

$D =$ distribution of the data

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

$$\mathbb{1} \{ A \} = \begin{cases} 1 & \text{if } A \text{ is true} \\ 0 & \text{if } A \text{ is false} \end{cases}$$

Ideally, we would like to find a predictor

h that has the smallest error; i.e.

our "gold standard" is to solve the following problem:

$$\begin{aligned} \underset{h}{\text{Minimize}} \quad & \underbrace{L_D(h)}_{\substack{P_{(x,y) \sim D} \{h(x) \neq y\} \\ = E_{(x,y) \sim D} [\mathbb{1}_{\{h(x) \neq y\}}]}} \end{aligned} \quad (1)$$

However, this task is impossible since D is unknown.

(3) The learner's input: Training Data!

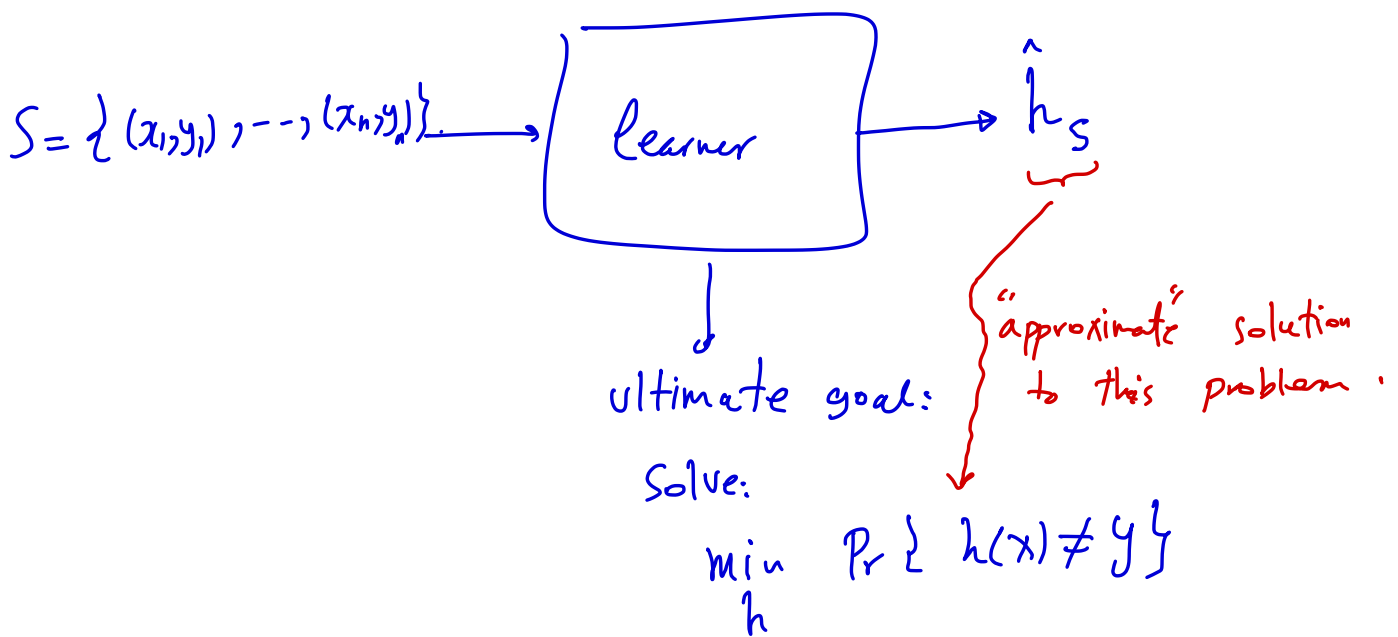
The only information that the learner has about the data distribution is a set of training samples:

$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

Where each (x_i, y_i) is generated i.i.d.

from the distribution D . Hence, the learner has to find "approximate" solutions to problem (1) using the training data.

— — — — —
Data distribution D



$$S = \{(x_1, y_1), \dots, (x_n, y_n)\} \quad \min_h \underbrace{E[\# \{ h(x) \neq y \}]}_{(x, y) \sim D}$$

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

(4) Empirical Risk Minimization (ERM):

true error $\rightarrow L_D(h) = E_{(x,y) \sim D} [\mathbb{1}\{h(x) \neq y\}]$

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

\downarrow unbiased estimate

ERM is the task of finding a predictor that minimizes the training error :

$$\min_h L_S(h)$$

instead of solving \rightarrow

$$\min_h L_D(h)$$

ERM solves \rightarrow

$$\min_h L_S(h)$$

claim:

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

overfitting

Example: Consider the following distribution on data:

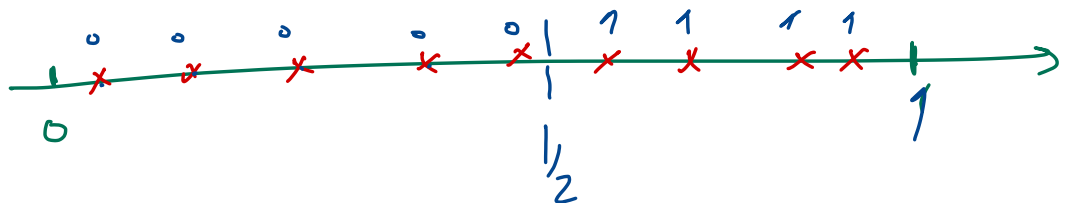
$$X = [0, 1]$$

$$Y = \{0, 1\}$$

→

$$x \sim \text{uniform}[0, 1]$$

$$y = \begin{cases} 0 & \text{if } x \leq 1/2 \\ 1 & \text{if } x > 1/2 \end{cases}$$



$$\min_h L_D(h) = \min_h \Pr\{h(x) \neq y\}$$

$$\min_h L_S(h) =$$

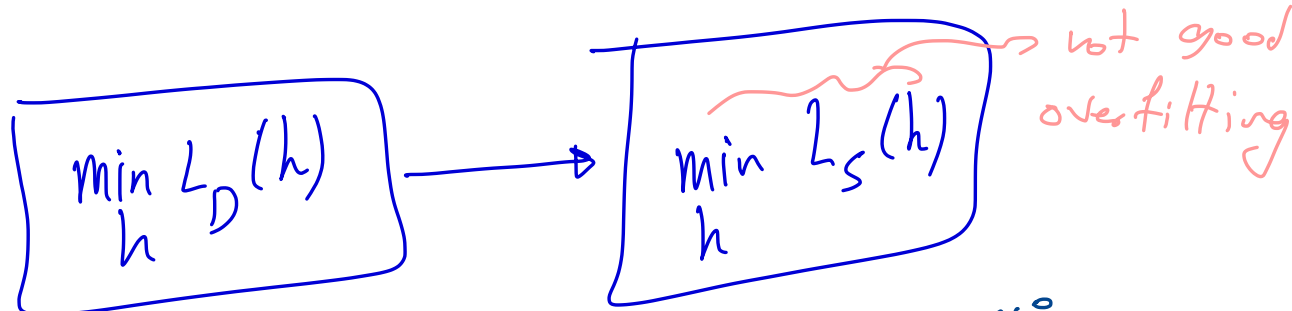
$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{h(x_i) \neq y_i\}$$

$$\hat{h}_S = \begin{cases} y_i & \text{if } x = x_i \\ 0 & x \neq \{x_1, \dots, x_n\} \end{cases}$$

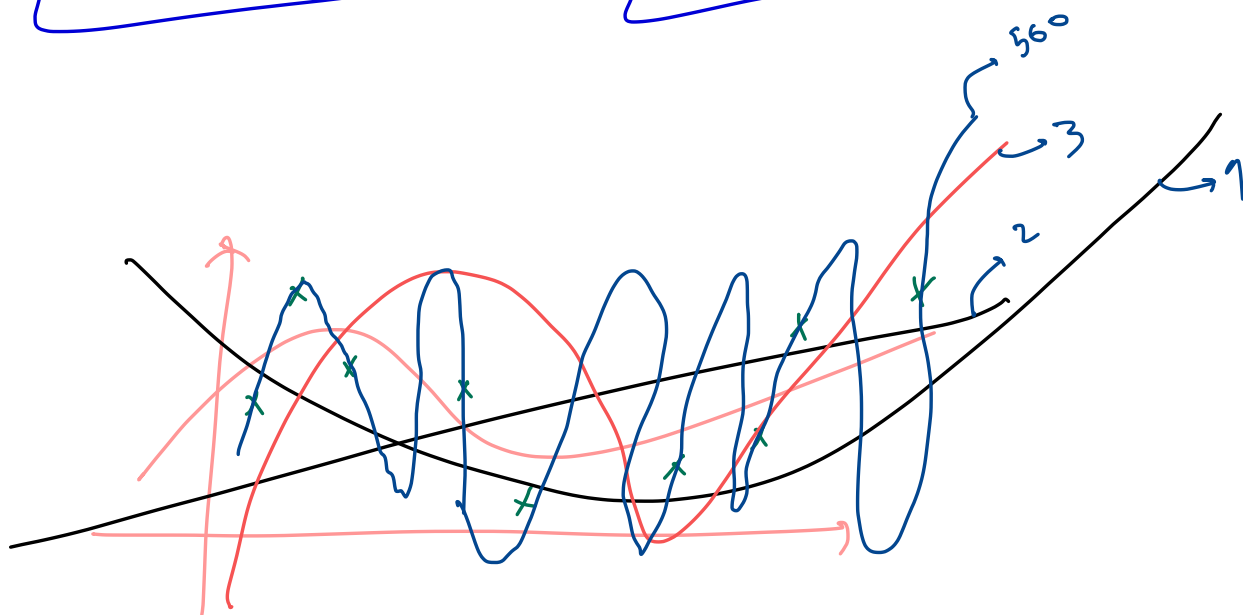
Proof of the claim:

$$\text{Define } \hat{h}_S(x) = \begin{cases} y_i & \text{if } x = x_i \\ 0 & \text{if } x \notin \{x_1, \dots, x_n\} \end{cases}$$

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

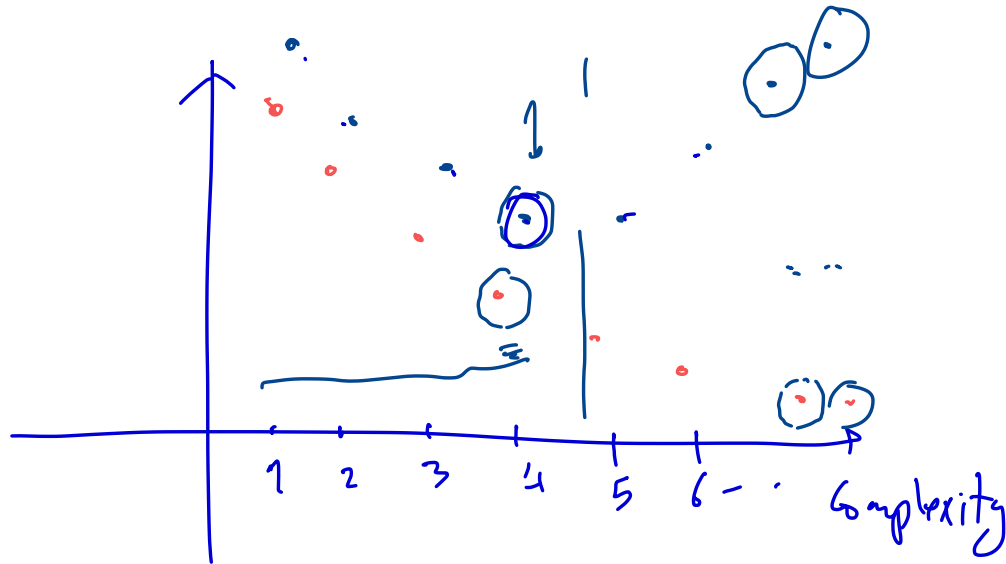


Recall
(Regression)



(5) Although ERM is natural, it can miserably if we are not careful:
It can overfit easily: the minimum of
$$\min_h L_S(h)$$
is always zero (we always overfit).

(6) We need to search for conditions under which there is a guarantee that ERM does not overfit; namely, conditions under which when ERM has good performance on training data, it also has good performance over the underlying distribution.



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

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

restrict the
class of functions/predictors
(restrict the complexity)

↓
prevents overfitting

The above problem is denoted by

$$\text{ERM}_{\mathcal{H}}$$

\mathcal{H} is the specific function class that
we use to fit the data.

