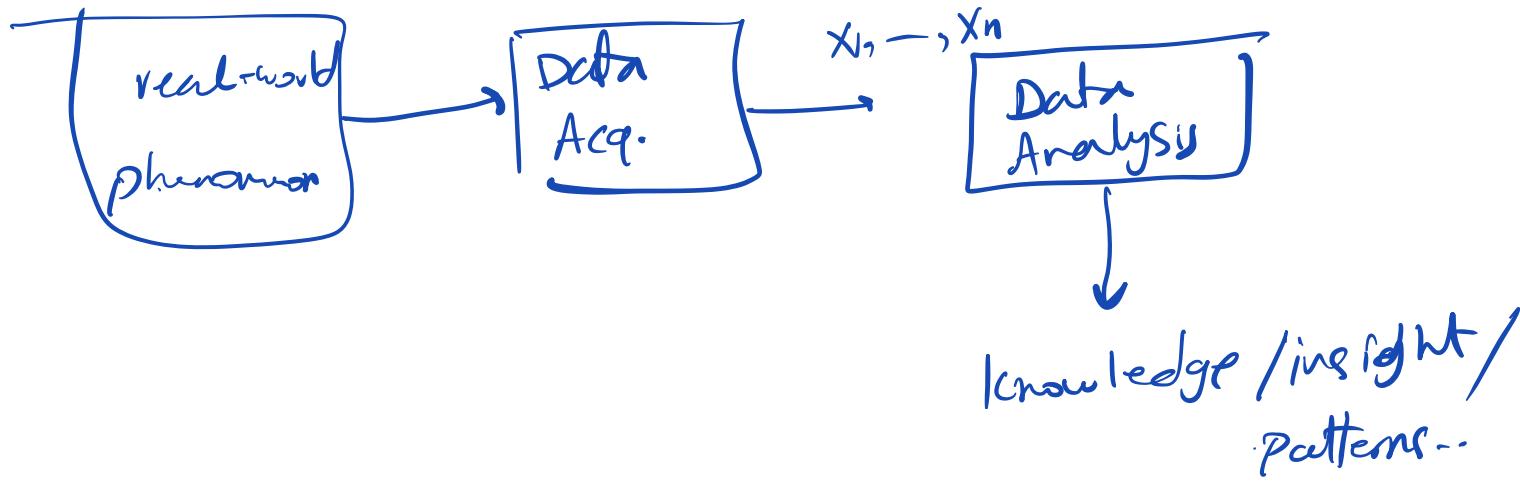


## Lecture 13:



Modules 1, 2 :

- learned to estimate properties/parameters connected to the underlying phenomenon
- hypothesis testing: verify different hypotheses connected with the underlying phenomenon.

~~Module 3:~~ **Supervised Learning:** "learn some 'predictive relation'" connected to the phenomenon.

## Supervised learning :

We assume that each data point  $x_i$  is associated with a "label" (which is a part of the data). I.e. each data point is of the form  $(x_i, y_i)$ , where  $x_i$  is the "feature vector" and " $y_i$ " is the label associated with  $x_i$ .

### Examples:

- {  $x_i$  : could be an image
- {  $y_i$  : could be an object inside the image
- {  $x_i$  : could be (the text of) an email
- {  $y_i$  : could be whether/not the email is spam.

$\left\{ \begin{array}{l} x_i: \text{a vector representing a patient's profile, } X_i = (\text{age}, \text{height}, \text{weight}, \text{...}) \\ y_i: \text{Cancer/ not cancer} \end{array} \right.$

---

Data :  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

Goal :

$y_i \approx f(x_i)$   
 spam/hom predictor  
 0 | 1  
 email text

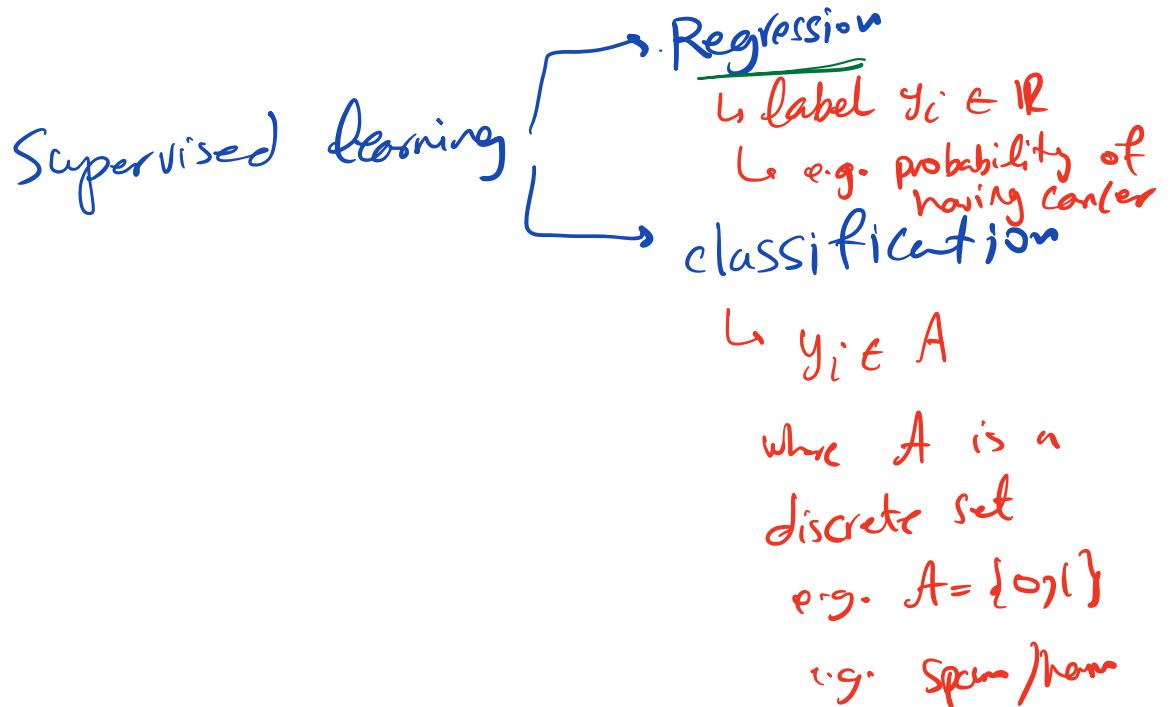
→ Main goal : Find the best predictor  $f$ .

Setting:

Data :  $\{(x_i, y_i)\}_{i=1}^n$

$x_i = (x_{i1}, x_{i2}, \dots, x_{ip}) \in \mathbb{R}^P$  ←  
 - image: pixel representation  
 - profile of the  $i$ -th patient

$y_i$  = label



## Regression:

- Let's start with the simplest setting (for the regression problem):
- 1-d regression

Data:  $\{(x_i, y_i)\}_{i=1}^n$

$$x_i \in \mathbb{R}$$

$$y_i \in \mathbb{R}$$

$$(p=1)$$

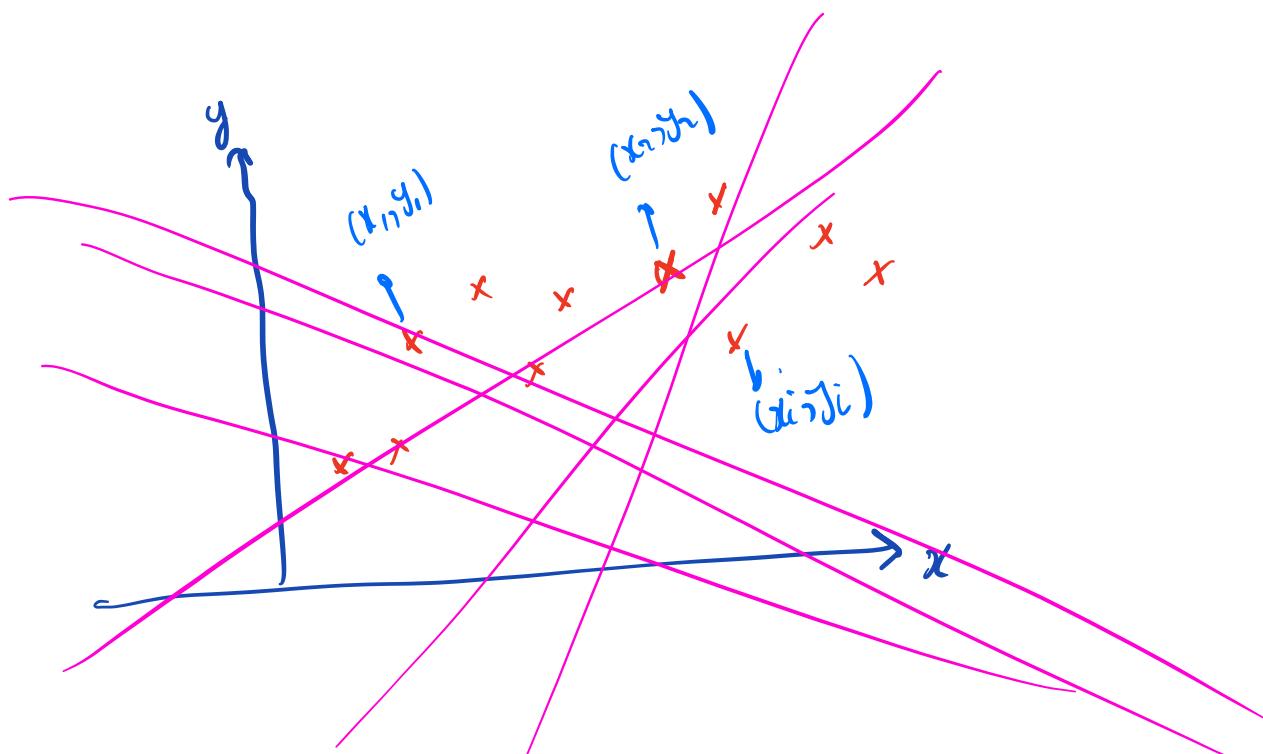
Goal: learn a "predictive" relation:

$$\hat{y} = \hat{f}(x)$$

Example:  $x$  = blood sugar at the age of 20

$y$  = probability of having Diabetes  
at the age of 40.

Data:  $\{(x_i, y_i)\}_{i=1}^n$



$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

$\underbrace{\quad}_{f(x)}$

Approach:

Goal: learn  $y = f(x)$

- Search for the best "candidate" for  $f$  inside a prescribed family of functions  $F$ .
- find the best  $f$  inside  $F$
- typically  $F$  is a parametric family  $\rightarrow f(x; \beta)$ .  
e.g. linear functions

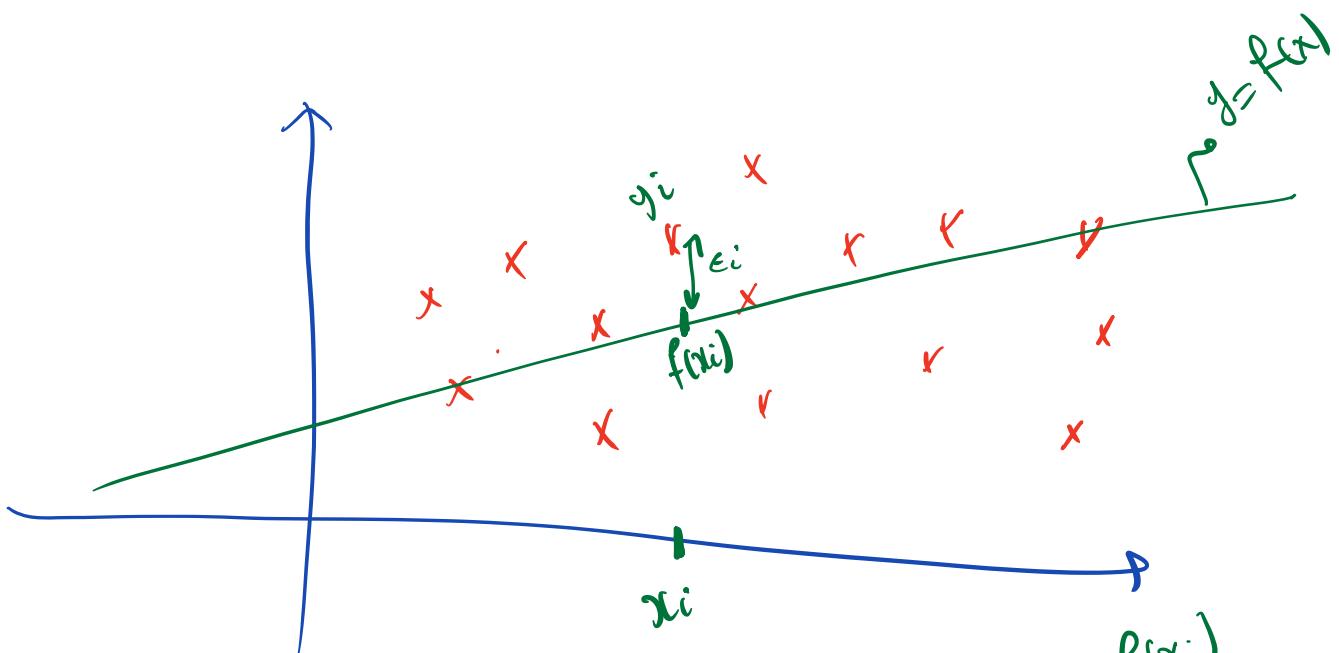
- Let's consider the family of linear functions as our first step. If  $x \in \mathbb{R}$  then this family is described by two parameters  $\hat{\beta}_0, \hat{\beta}_1$

$$f(x; \hat{\beta}_0, \hat{\beta}_1) = \hat{\beta}_0 + \hat{\beta}_1 x$$

$$\underbrace{\beta_2}_{\beta}$$

$$\underbrace{\downarrow}_{\beta_0} \quad \underbrace{\downarrow}_{\beta_1}$$

Goals: find the best choice of  $\hat{\beta}_0, \hat{\beta}_1$  from data.



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

error =  $\varepsilon_i := y_i - \underbrace{f(x_i)}_{\hat{\beta}_0 + \hat{\beta}_1 x_i}$

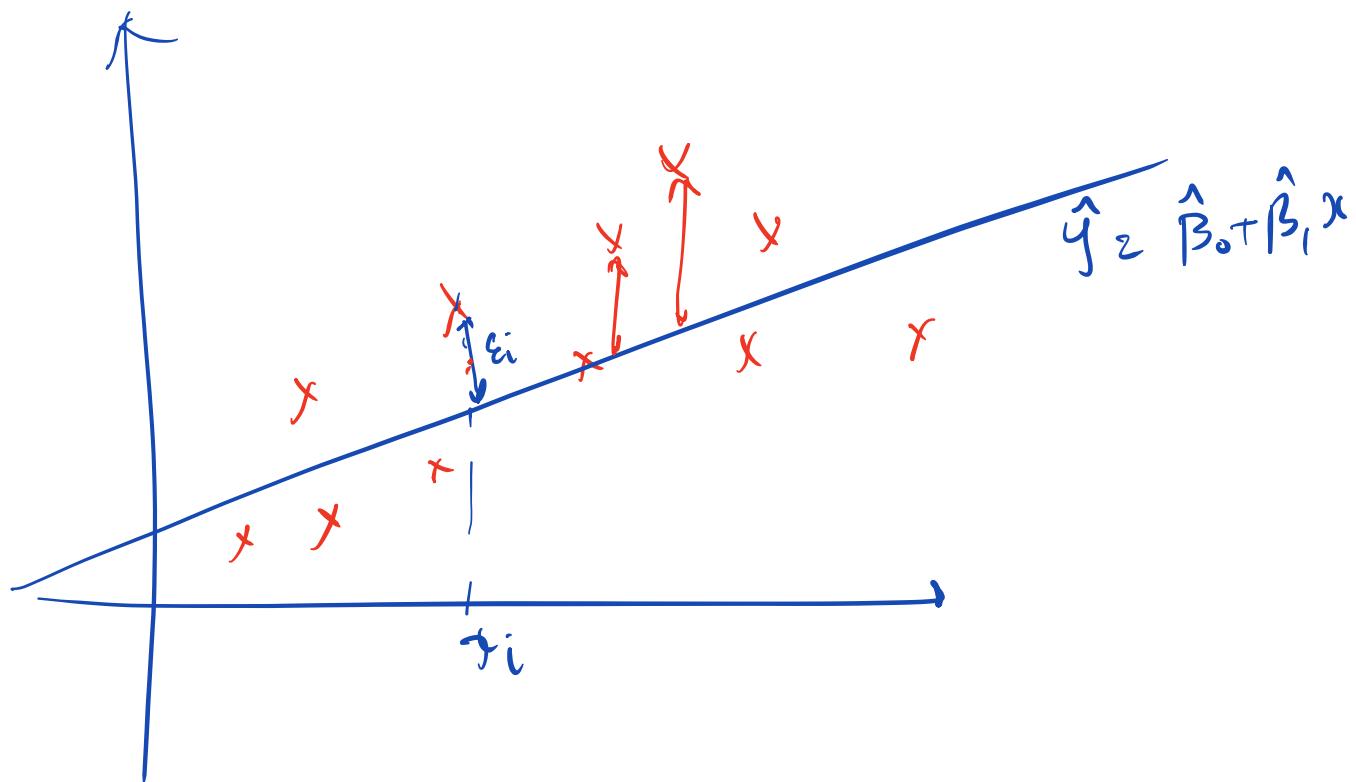
$$y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + \varepsilon_i$$

the error term models what we miss when we use this simple linear model. The true relation between the input and output is probably not linear as many other factors can be involved

(there may also be noise in  
the label)

Question:

What is the fundamental procedure  
principle behind learning/choosing  
the parameters  $\hat{\beta}_0, \hat{\beta}_1$ ?



Idea 1: find  $\hat{\beta}_0, \hat{\beta}_1$  s.t.

sum of  
the errors

$$\sum_{i=1}^n |\epsilon_i|$$

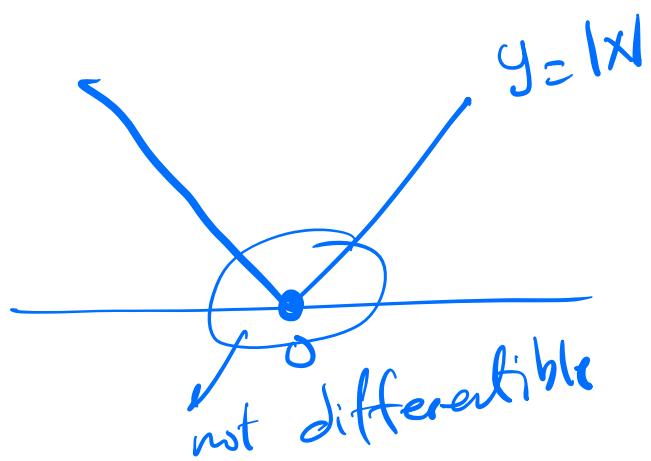
$$\epsilon_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)$$

is minimized

$$\min_{\hat{\beta}_0, \hat{\beta}_1} \sum_{i=1}^n |y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i|$$

↓

- In order to solve the above optimization problem, we'll have to take the derivative w.r.t.  $\hat{\beta}_0, \hat{\beta}_1$  and set it to zero.
- The problem with the above objective is that it's not differentiable.



- Idea 2: To make things differentiable  
let's consider the "square" of  
the errors:

minimize  $\hat{\beta}_0, \hat{\beta}_1$   $\sum_{i=1}^n \left( y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \right)^2$

$\underbrace{\qquad\qquad\qquad}_{\varepsilon_i^2}$

## Lecture 14:

supervised learning

Regression  $\rightarrow y \in \mathbb{R}$

Classification

Data :  $\{(x_i, y_i)\}_{i=1}^n$

$$x_i \in \mathbb{R}^P$$

$$y_i \in \mathbb{R}$$

$$\rightarrow y_i \approx f(x_i)$$

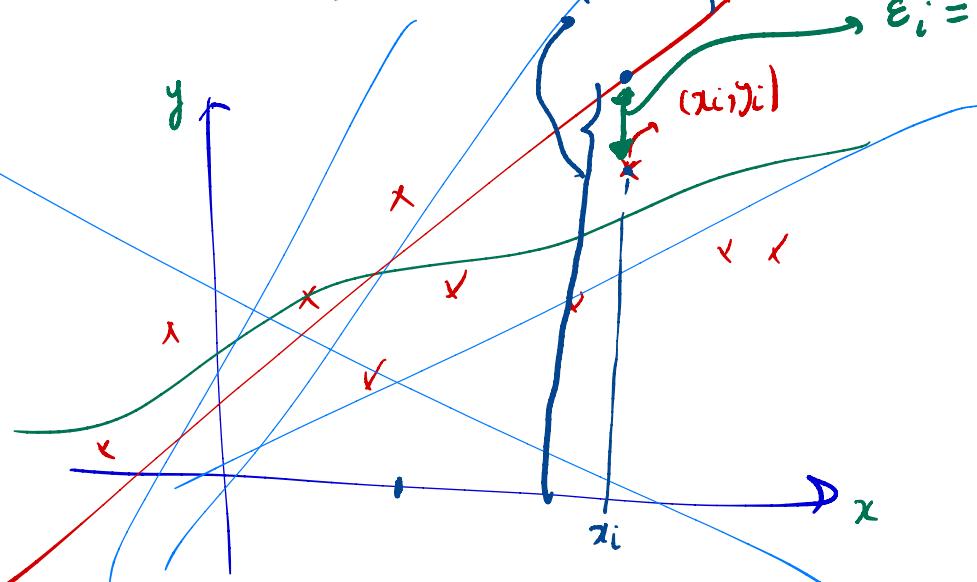
will be used for prediction on newly  
unseen data points

$$P=1$$

$$\begin{cases} x_i \in \mathbb{R} \\ y_i \in \mathbb{R} \end{cases}$$

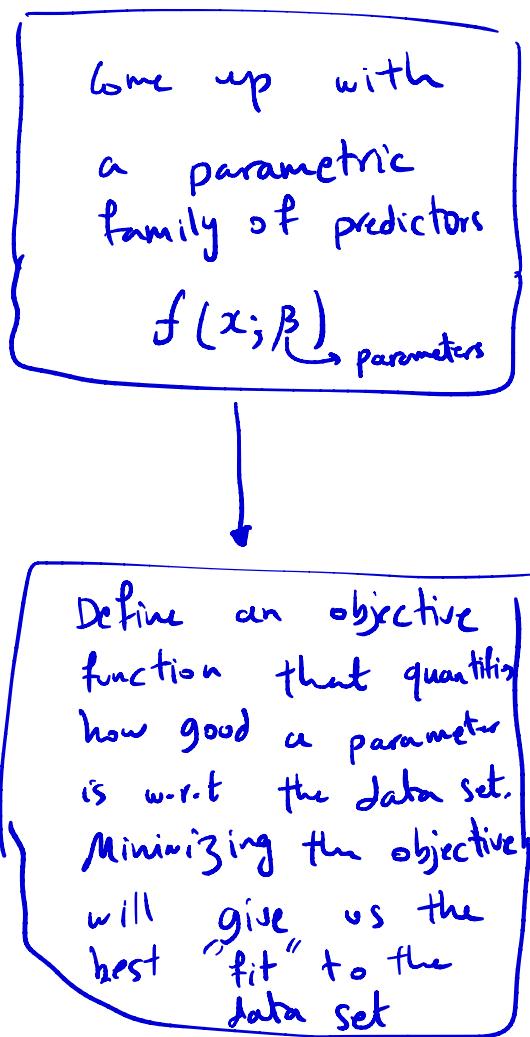
$$f(x; \beta)$$

$$\epsilon_i = y_i - f(x_i; \beta)$$



The simplest parametric family:  $y = f(x; \beta) = \hat{\beta}_0 + \hat{\beta}_1 x$

our general  
methodology:



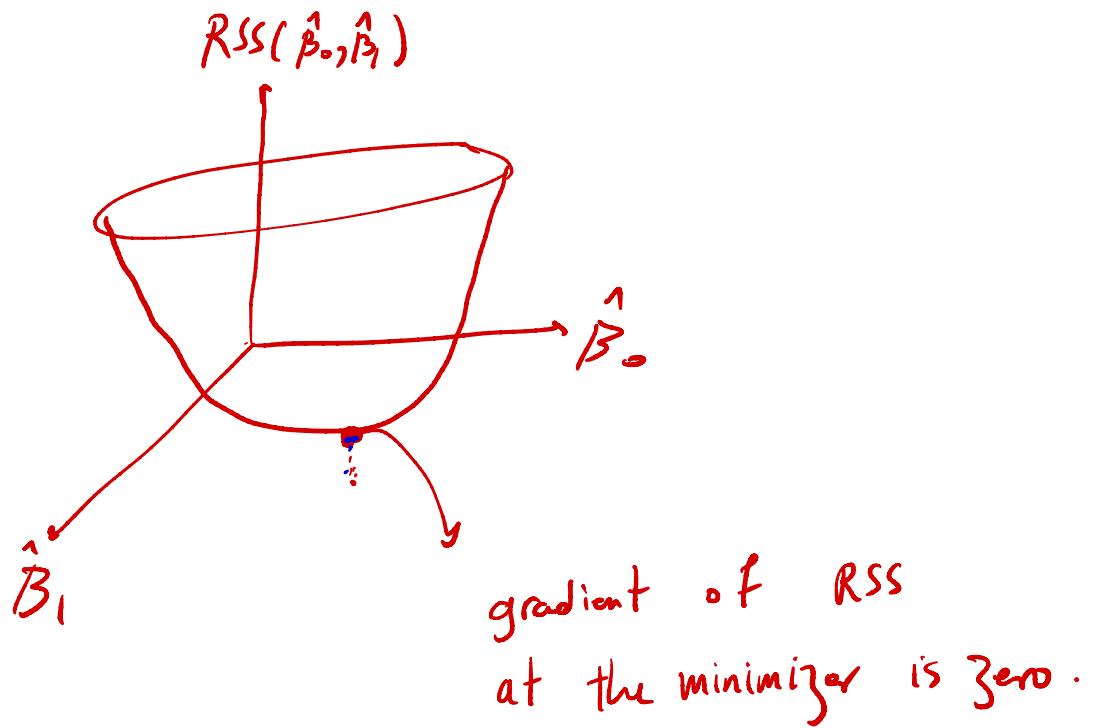
---

objective for the linear model:

$$\min_{\hat{\beta}_0, \hat{\beta}_1} \sum_{i=1}^n \epsilon_i^2 = \min_{\hat{\beta}_0, \hat{\beta}_1} \sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2$$

$\underbrace{\qquad\qquad\qquad}_{\text{RSS } (\hat{\beta}_0, \hat{\beta}_1)}$

Residual      Sum of      Squares



These two equations  
Should hold at  
the minimizer.

$$\left. \begin{aligned} \frac{d \text{RSS}(\hat{\beta}_0, \hat{\beta}_1)}{d \hat{\beta}_0} &= 0 \\ \frac{d \text{RSS}(\hat{\beta}_0, \hat{\beta}_1)}{d \hat{\beta}_1} &= 0 \end{aligned} \right\}$$

$$\text{RSS}(\hat{\beta}_0, \hat{\beta}_1) = \sum_{i=1}^n \left( y_i - \hat{\beta}_0 + \hat{\beta}_1 x_i \right)^2$$

$$\frac{d}{d\hat{\beta}_0} RSS(\hat{\beta}_0, \hat{\beta}_1) \quad (1)$$

$$= -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$

$$\frac{d}{d\hat{\beta}_1} RSS(\hat{\beta}_0, \hat{\beta}_1) \quad (2)$$

$$= -2 \sum_{i=1}^n x_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$


---

$$(1) \rightarrow \sum_{i=1}^n y_i - n \hat{\beta}_0 - \hat{\beta}_1 \sum_{i=1}^n x_i = 0$$

$\cancel{n} \div n$

$$\left( \frac{\sum_{i=1}^n y_i}{n} \right) - \hat{\beta}_0 - \hat{\beta}_1 \left( \frac{\sum_{i=1}^n x_i}{n} \right) = 0$$

$$\bar{y} - \hat{\beta}_0 - \hat{\beta}_1 \bar{x} = 0$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}$$

(3)  
(equivalent to (1))

$$(2) \Rightarrow \sum_{i=1}^n x_i y_i - \hat{\beta}_0 \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i^2 = 0$$

$$\div n \Rightarrow \frac{\sum x_i y_i}{n} - \underbrace{\hat{\beta}_0 \bar{x}} - \hat{\beta}_1 \frac{\sum x_i^2}{n} = 0$$

use (3)

$$\Rightarrow \frac{\sum x_i y_i}{n} - (\bar{y} - \hat{\beta}_1 \bar{x}) \bar{x} - \hat{\beta}_1 \frac{\sum x_i^2}{n} = 0$$

$$\frac{\sum x_i y_i}{n} - \bar{x} \bar{y} = \hat{\beta}_1 \left( \frac{\sum x_i^2}{n} - \bar{x}^2 \right)$$

Two relations (exercise)

$$\frac{\sum x_i y_i}{n} - \bar{x} \bar{y} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \leftarrow$$

$$\frac{\sum x_i^2}{n} - \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$(3) \rightarrow \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\Rightarrow \hat{\beta}_0 = \bar{y} - \bar{x} \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

**Question:** How good are the above estimates  $\hat{\beta}_0, \hat{\beta}_1$ ?

To evaluate the performance of our estimates of

$\hat{\beta}_0, \hat{\beta}_1$ , let's assume that the data is generated according to the following probabilistic model:

$$y = \underbrace{\beta_0}_{\text{true}} + \underbrace{\beta_1}_{\sim} x + \varepsilon$$

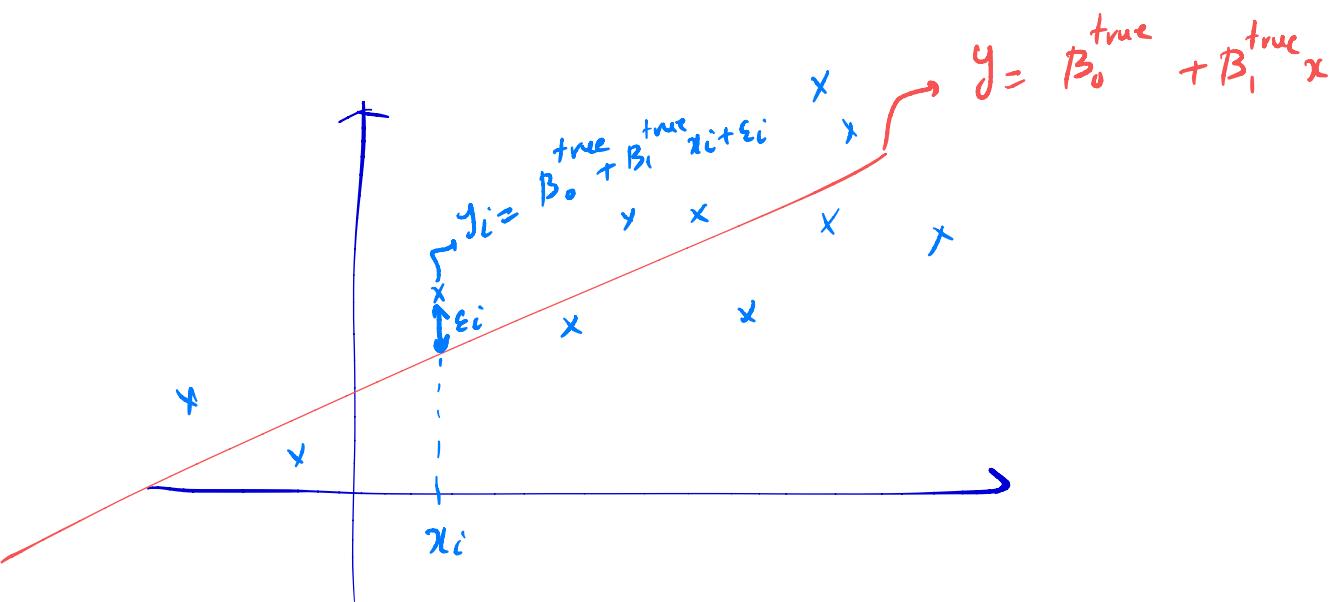
noise  
 zero-mean  
 independent  
 $\text{var}[\varepsilon] = 6^2$

Now, assume that we generate data points

$\{(x_i, y_i)\}_{i=1}^n$  from the above model. Let  $\hat{\beta}_0, \hat{\beta}_1$  be the solutions of the regression problem.

as  $n \rightarrow \infty$

$\hat{\beta}_0 \xrightarrow{\sim} \beta_0^{\text{true}}$	$\hat{\beta}_0$ is an estimator of: $\beta_0^{\text{true}}$
$\hat{\beta}_1 \xrightarrow{\sim} \beta_1^{\text{true}}$	$\hat{\beta}_1$ is an estimator of: $\beta_1^{\text{true}}$



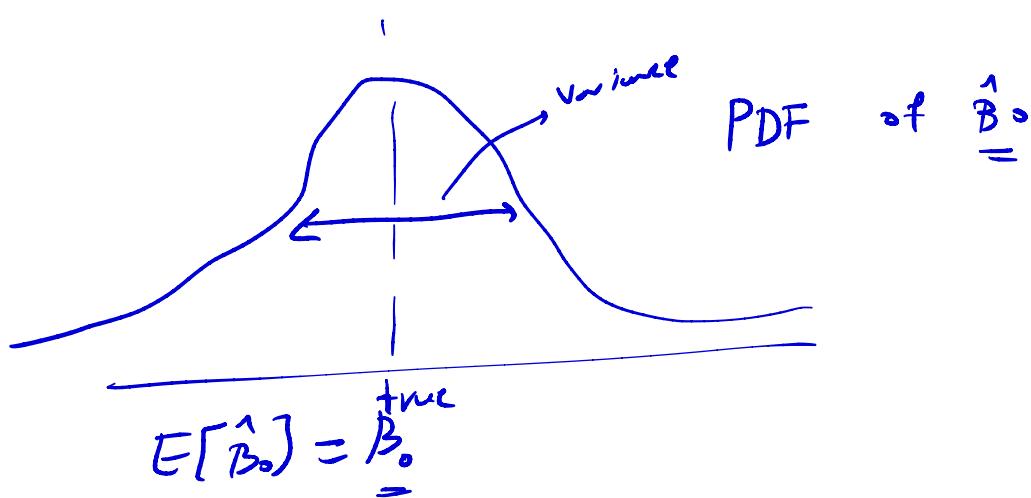
$$\hat{B}_0 \longrightarrow B_0^{\text{true}}$$

$$\hat{B}_1 \longrightarrow B_1^{\text{true}}$$

- In order to quantify the performance of the estimators  $\hat{B}_0$ ,  $\hat{B}_1$ , we need to compute their expectation (Bias) and Variance.
- The expectation and the variance of  $\hat{B}_0$  (and  $\hat{B}_1$ ) will tell us how close  $\hat{B}_0$  (and  $\hat{B}_1$ ) are w.r.t their true value  $B_0^{\text{true}}$  (and  $B_1^{\text{true}}$ )
- Note that, because of the fact that the noises  $\epsilon_i$ 's are random variables, our

data set  $\{(x_i, y_i)\}$  is also random.

This means that  $\hat{\beta}_0, \hat{\beta}_1$ , which depend on the data, are random variables as well.



---

We'll now compute the expectation of  $\hat{\beta}_0, \hat{\beta}_1$  and show that they are unbiased estimators of  $\beta_0^{\text{true}}$  and  $\beta_1^{\text{true}}$ , respectively.

$$E[\hat{\beta}_1] \stackrel{?}{=} \beta_1^{\text{true}}$$

↓

$$E[\hat{\beta}_1] = E\left[ \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

note:

$$y_i = \beta_0^{\text{true}} + \beta_1^{\text{true}} x_i + \varepsilon_i$$

↓  
zero-mean  
independent  
of everything else

Hence:

$$y_i = \beta_0^{\text{true}} + \beta_1^{\text{true}} x_i + \varepsilon_i$$

↓

$$\bar{y} = \beta_0^{\text{true}} + \beta_1^{\text{true}} \bar{x} + \bar{\varepsilon}$$
(1)

$(\bar{\varepsilon} = \frac{\varepsilon_1 + \dots + \varepsilon_n}{n})$

$$E[\hat{\beta}_1] = E\left[ \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \right]$$

(1)

$$= E\left[ \frac{\sum_i (x_i - \bar{x}) (\beta_1^{\text{true}}(x_i - \bar{x}) + \varepsilon_i - \bar{\varepsilon})}{\sum_i (x_i - \bar{x})^2} \right]$$

$$\begin{aligned}
 & \left( y_i - \bar{y} \right) \\
 & \downarrow \\
 & \left( B_0^{\text{true}} + B_1^{\text{true}} x_i + \varepsilon_i \right) - \left( B_0^{\text{true}} + B_1^{\text{true}} \bar{x} + \bar{\varepsilon} \right) \\
 & = B_1^{\text{true}} (x_i - \bar{x}) + \varepsilon_i - \bar{\varepsilon}
 \end{aligned}$$

$$E[\hat{B}_1] = E \left[ \frac{B_1^{\text{true}} \sum_i (x_i - \bar{x})^2 + \sum (\varepsilon_i - \bar{\varepsilon})(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} \right]$$

$$= B_1^{\text{true}} + E \left[ \sum_{i=1}^n \frac{(\varepsilon_i - \bar{\varepsilon})(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} \right]$$

noises  
 were independent  
 of everything  
 else

$$= B_1^{\text{true}} + \sum_{i=1}^n E \left[ \frac{(\varepsilon_i - \bar{\varepsilon})(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} \right]$$

$$= B_1^{\text{true}} + \sum_{i=1}^n \underbrace{E[\varepsilon_i - \bar{\varepsilon}]}_{E[\varepsilon_i] - E[\bar{\varepsilon}]} \underbrace{E \left[ \frac{(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} \right]}_{0}$$

$$\Rightarrow E[\hat{\beta}_1] = \beta_1^{\text{true}}$$

In general:

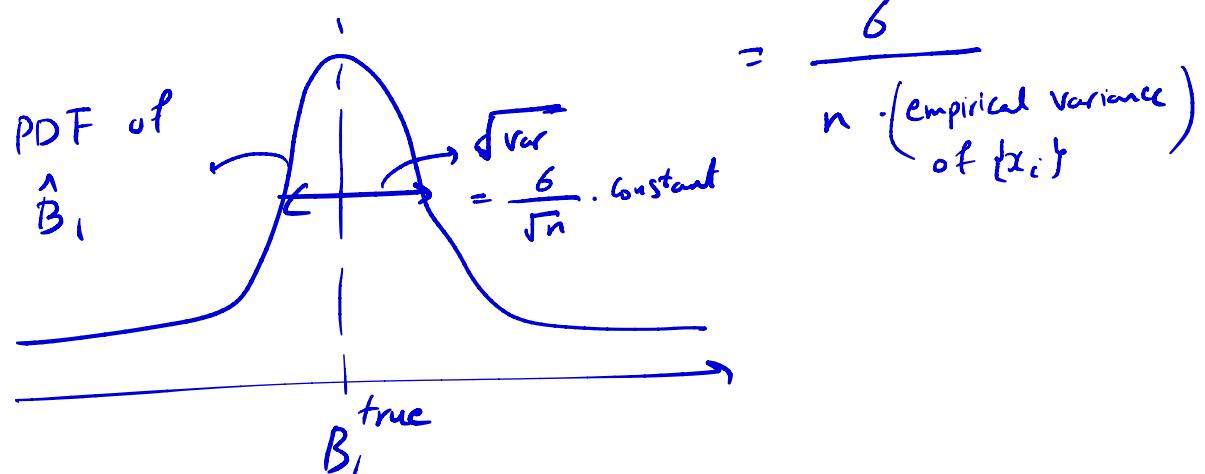
$$(1) \quad E[\hat{\beta}_1] = \beta_1^{\text{true}}$$

$$(2) \quad E[\hat{\beta}_0] = \beta_0^{\text{true}}$$

$$(3) \quad \text{Var}(\hat{\beta}_0) = E[(\hat{\beta}_0 - \beta_0^{\text{true}})^2]$$

$$(\text{Var}(\varepsilon_i) = \sigma^2) \quad = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)$$

$$(4) \quad \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sigma^2}{n \underbrace{\left( \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right)}_{\text{empirical variance of } \{x_i\}}}$$



So far we've only been talking about the 1-dim case ( $p=1$ ). Let's now consider the general case:

$$\{ (x_i, y_i) \}_{i=1}^n \quad \begin{array}{l} x_i \in \mathbb{R}^p, y_i \in \mathbb{R} \\ \text{---} \\ x_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix} \in \mathbb{R}^p, y_i \in \mathbb{R} \end{array}$$

simplest parametric family (linear):

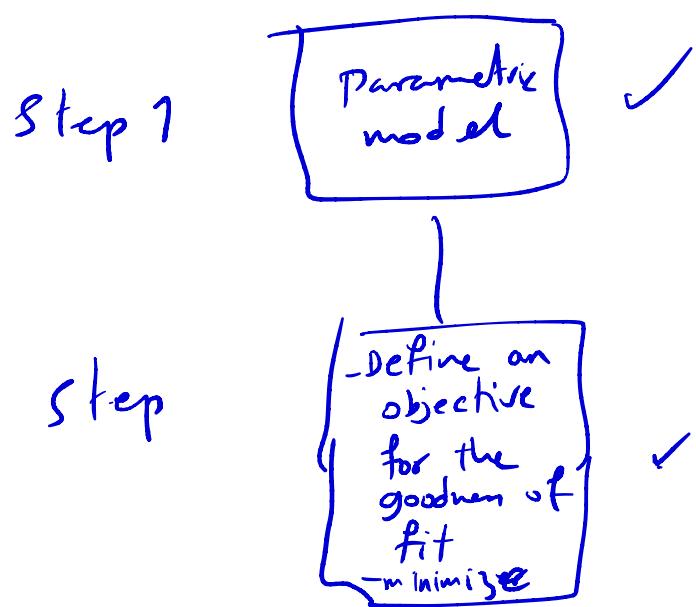
$$y_i \approx \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \cdots + \hat{\beta}_p x_{ip}$$

$\underbrace{\qquad\qquad\qquad}_{f(x_i, \vec{\beta})}$

## Lecture 15

P-dimensional linear regression:

$$\text{Data: } \{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^P$$



Step 1:

$$x_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix} \in \mathbb{R}^P \rightarrow y_i \approx \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip}$$

P+1 parameters  
(to be learned)  
from data

Step: 2:

$$RSS = \sum_{i=1}^n (\text{error}_i)^2$$

$$= \sum_{i=1}^n \left( y_i - \underbrace{\left( \hat{B}_0 + \hat{B}_1 x_{i1} + \dots + \hat{B}_p x_{ip} \right)}_{\substack{\text{true label} \\ \text{prediction}}} \right)^2$$

to find  $\{\hat{B}_j\}_{j=0}^p$  we'll have to  
minimize RSS.

$$RSS = \sum_{i=1}^n (y_i - \vec{x}_i^T \vec{B})^2$$

$$\hat{B} = \begin{pmatrix} \hat{B}_0 \\ \hat{B}_1 \\ \vdots \\ \hat{B}_p \end{pmatrix} \in \mathbb{R}^{p+1}, \vec{x}_i = \begin{pmatrix} 1 \\ x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix} \in \mathbb{R}^{p+1}$$

$$RSS = \|\vec{Y} - X \vec{B}\|_2^2$$

$$\boxed{\vec{Y} = \begin{bmatrix} 1, x_{i1}, x_{i2}, \dots, x_{ip} \end{bmatrix}^T, \vec{B} = \begin{pmatrix} \hat{B}_0 \\ \hat{B}_1 \\ \vdots \\ \hat{B}_p \end{pmatrix}}$$

$\hat{B}_0 + \hat{B}_1 x_{i1} + \dots + \hat{B}_p x_{ip}$   
 $\text{data}_i$

$$\vec{Y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}_{n \times 1}$$

$$\vec{X} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1p} \\ 1 & x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{np} \end{pmatrix}_{n \times (p+1)}$$

$$= \begin{pmatrix} \tilde{x}_1^T \\ \tilde{x}_2^T \\ \vdots \\ \tilde{x}_n^T \end{pmatrix}$$

$$\| \vec{z} \|_2^2 = \sum_{i=1}^n z_i^2$$

$$\vec{z} = \begin{pmatrix} z_1 \\ \vdots \\ z_n \end{pmatrix}$$

$$\underset{\vec{\beta}}{\text{minimize}} \quad RSS(\vec{\beta})$$

$$\underset{\vec{\beta}}{\text{minimize}} \quad \underbrace{\| \vec{Y} - \vec{X} \vec{\beta} \|_2^2}_{RSS(\vec{\beta})}$$

$$\Rightarrow \underbrace{\nabla \text{RSS}(\vec{B})}_{\downarrow} = 0$$

$$-2 \underset{=}{} \vec{X}^T (\vec{y} - \vec{X} \vec{B}) = 0$$

$$\Rightarrow \vec{X}^T \vec{X} \vec{B} = \vec{X}^T \vec{y}$$

$$\Rightarrow \vec{B} = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{y}$$

$$\begin{pmatrix} \hat{B}_0 \\ \hat{B}_1 \\ \vdots \\ \hat{B}_p \end{pmatrix}$$

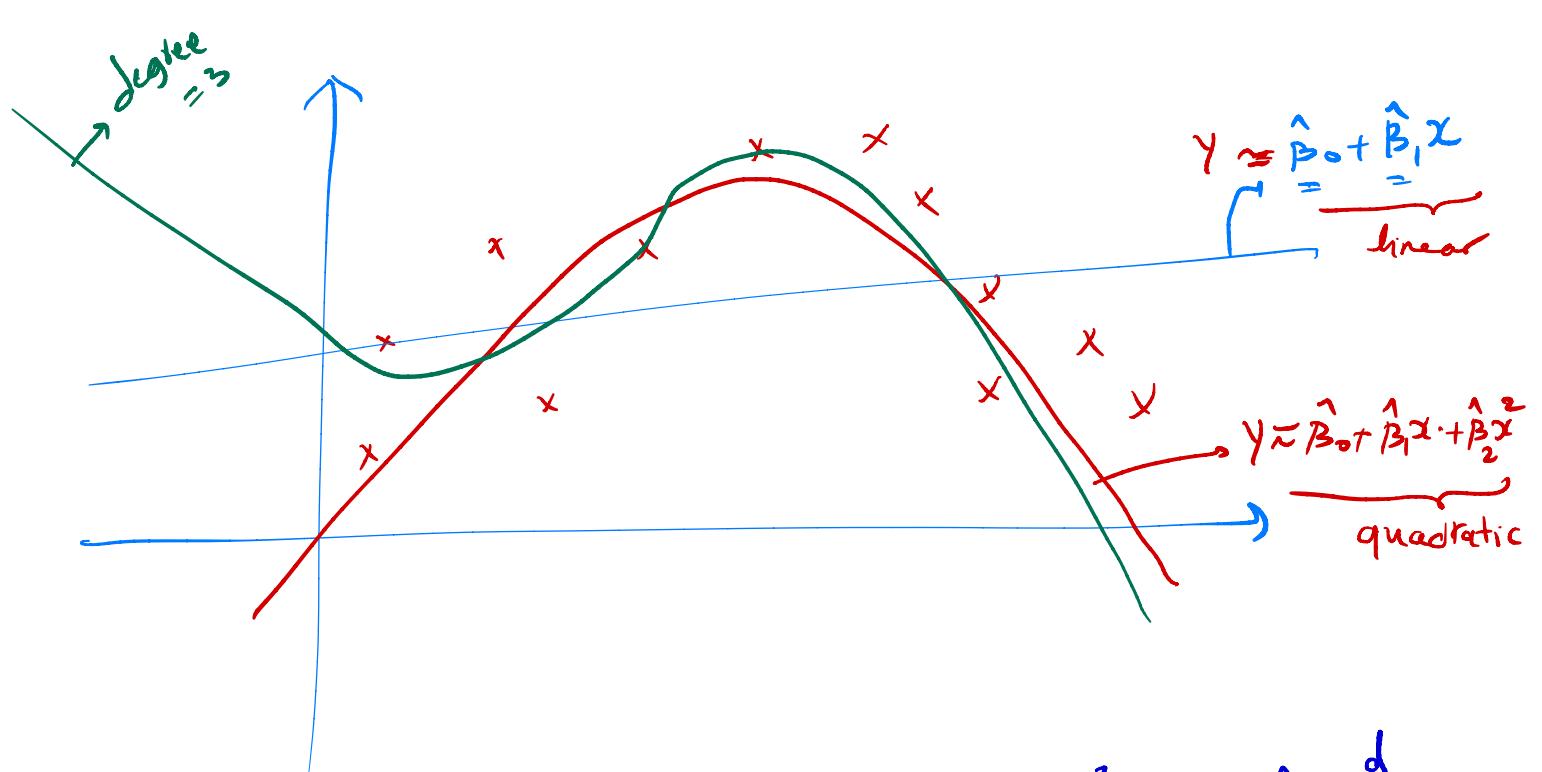
$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

$$X = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1p} \\ 1 & x_{21} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{i1} & \cdots & x_{ip} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{pmatrix}$$

i-th

- So far, we've only talked about linear models.  
 Let's now see how we can construct more complex (non-linear) parametric models..

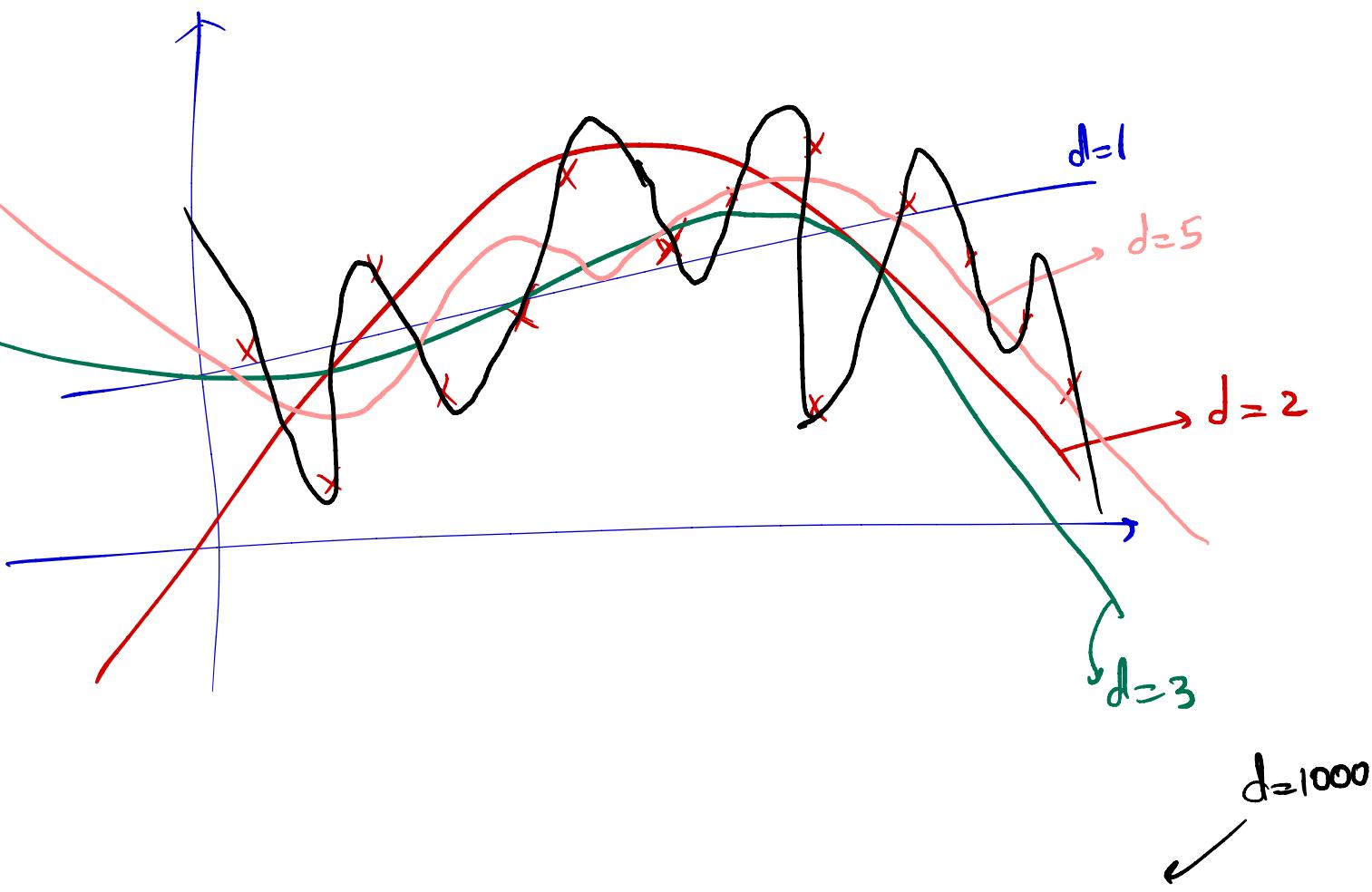
Let's assume  $p=1$ .



polynomial :  $y = p(x) = \hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 x^2 + \dots + \hat{\beta}_d x^d$

One way to construct non-linear parametric models is by using polynomial models:

$$\left\{ \begin{array}{l} \text{degree of the model: } d \\ \text{parameters: } \hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_d \quad (\text{d+1 parameters}) \end{array} \right.$$



$$y \approx \hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 x^2 + \dots + \hat{\beta}_{1000} x^{1000}$$

Question: Given a polynomial family of degree  $d$ , how do we find the best parameters (the best fit)?

} Data:  $\{(x_i, y_i)\}_{i=1}^n$      $x_i \in \mathbb{R}$ ,  $y_i \in \mathbb{R}$   
 } Step 1: Parametric model: polynomials of degree  $d$ :  $y = \hat{\beta}_0 + \hat{\beta}_1 x + \dots + \hat{\beta}_d x^d$

Step: Define an objective for the goodness of fit:

$$\xrightarrow{\text{objective}} \text{RSS} = \sum_{i=1}^n \text{error}_i^2$$

$$= \sum_{i=1}^n \left( y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{\beta}_2 x_i^2 + \dots + \hat{\beta}_d x_i^d) \right)^2$$

this objective is analogous to the multidimensional linear regression:

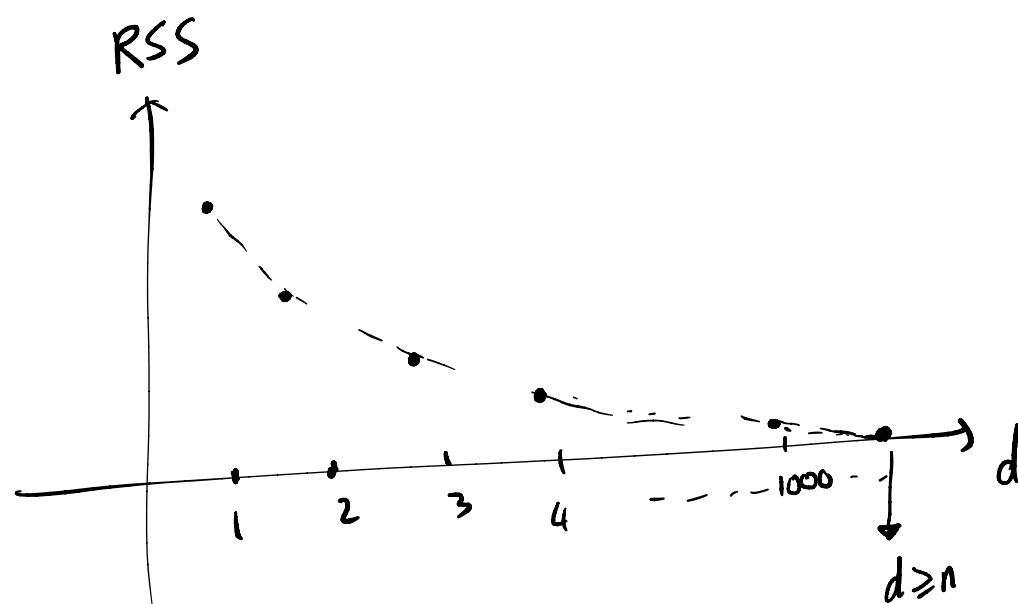
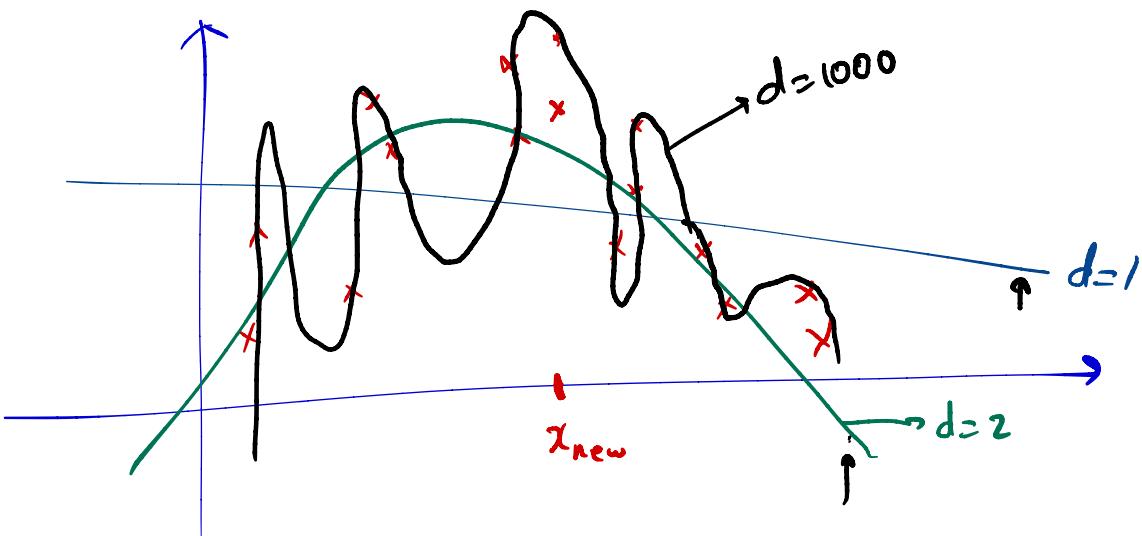
$$\underline{x_i} \in \mathbb{R} \rightarrow \tilde{x}_i = \begin{pmatrix} 1 \\ x_i \\ x_i^2 \\ \vdots \\ x_i^d \end{pmatrix}, \quad \vec{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \vdots \\ \hat{\beta}_d \end{pmatrix}$$

new data point #i

$$\text{RSS} = \sum_{i=1}^n (y_i - \tilde{x}_i^T \vec{\beta})^2$$

$$\Rightarrow \vec{\beta} = (\vec{x}^T \vec{x})^{-1} (\vec{x}^T \vec{y})$$

$$\vec{x} = \begin{pmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^d \\ 1 & x_2 & x_2^2 & \cdots & x_2^d \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^d \end{pmatrix}_{n \times (d+1)} \quad \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$



- Remember that the main goal behind learning predictive models was to be able to predict well on new and unseen data points.

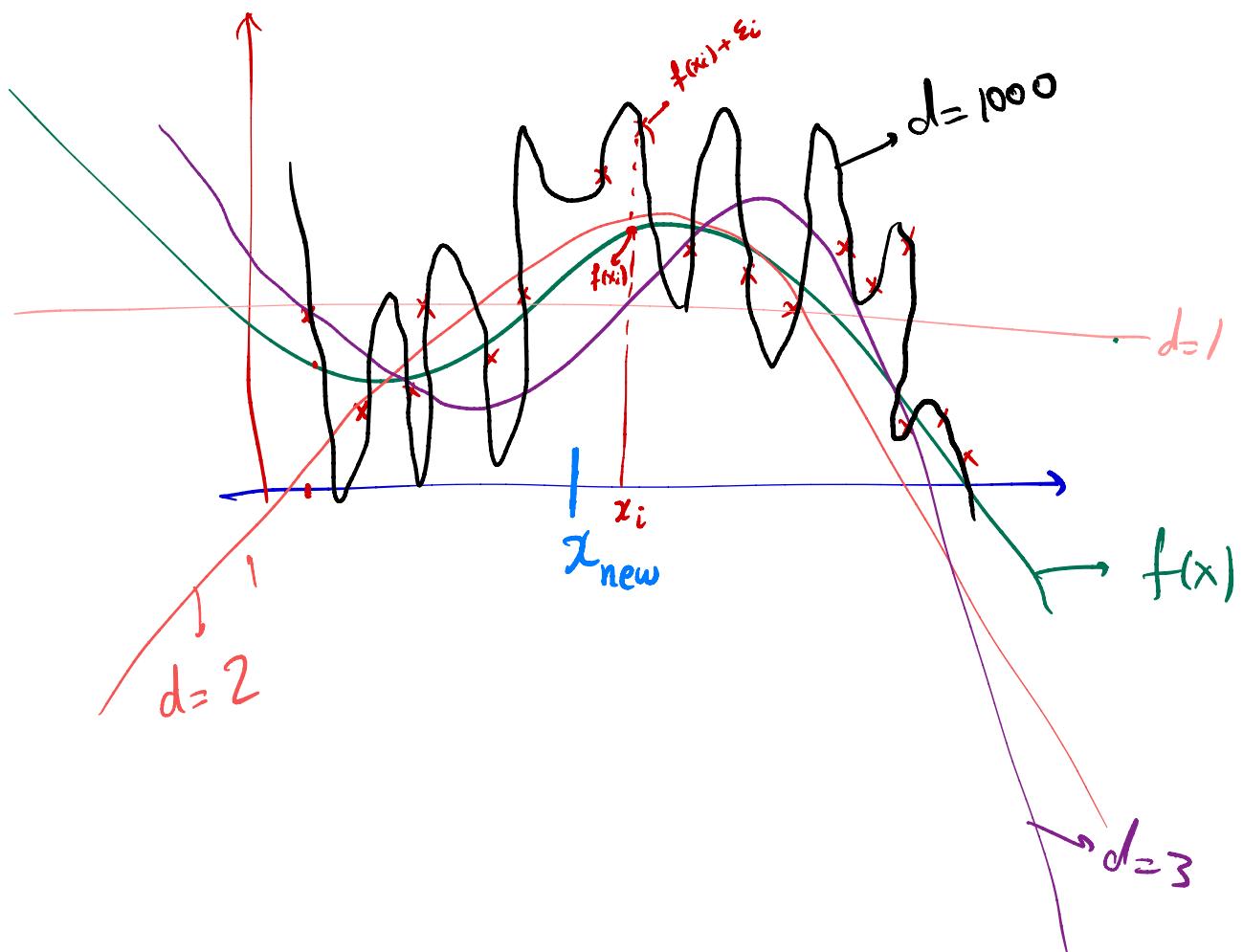
- Consider the following model for data:

how the data is generated

$$y = f(x) + \varepsilon$$

noise

e.g.  $f(x) = x^3 + 3x^2 - 4x + 5$



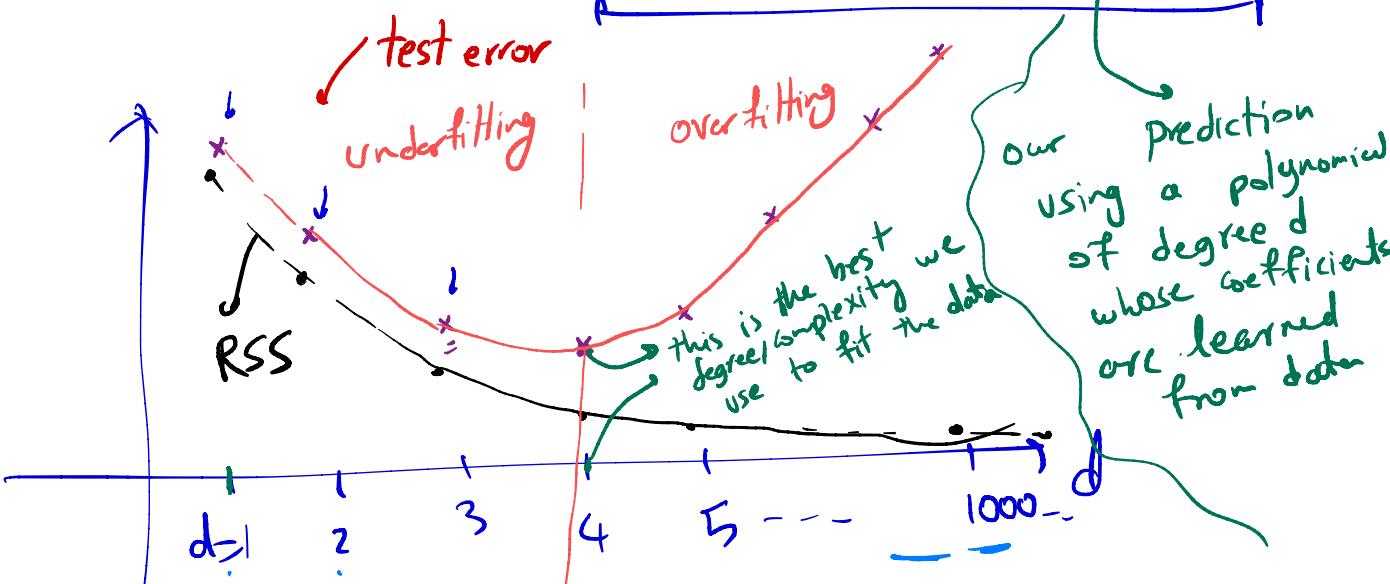
- On a new data point  $x_{\text{new}}$

$$\rightarrow y_{\text{new}} = \underline{f(x_{\text{new}})} + \varepsilon$$

w

the model learned from data

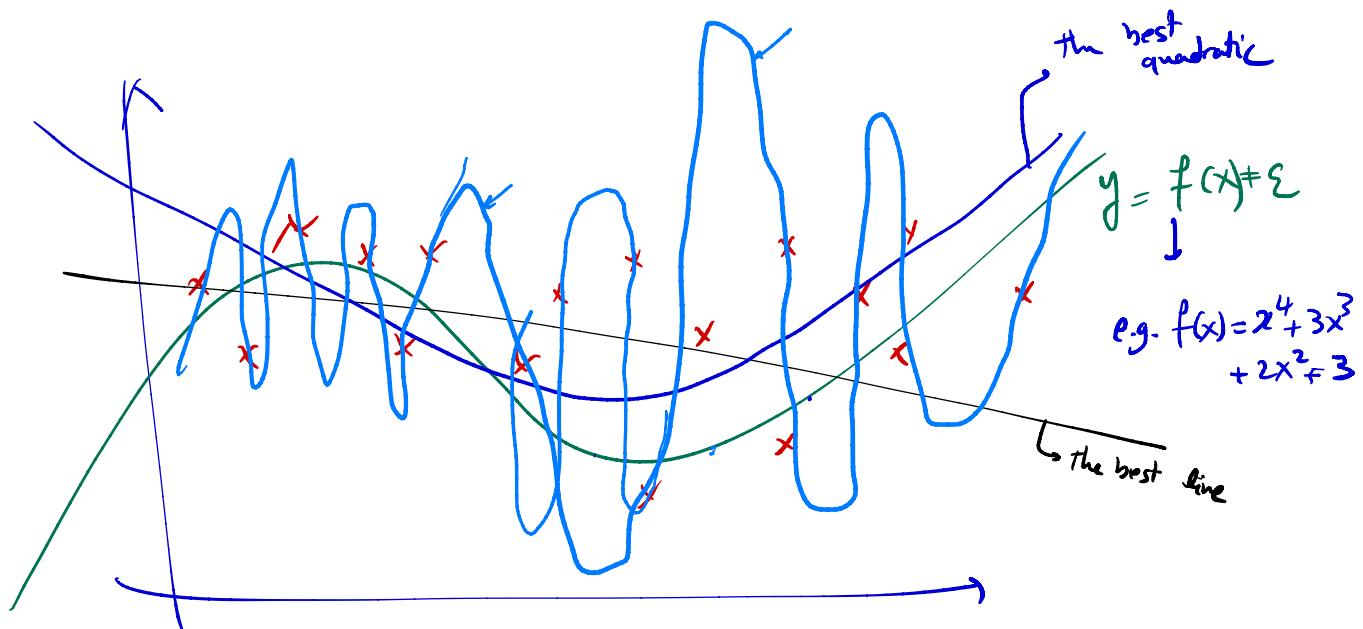
$$\rightarrow \text{Test error} = E \left[ \left( y_{\text{new}} - \underbrace{P_d(x_{\text{new}}, \vec{\beta})}_{y_{\text{new}} = f(x_{\text{new}}) + \varepsilon} \right)^2 \right]$$



we have two regions depending on the complexity (degree) of the parametric model that we're fitting to the data:

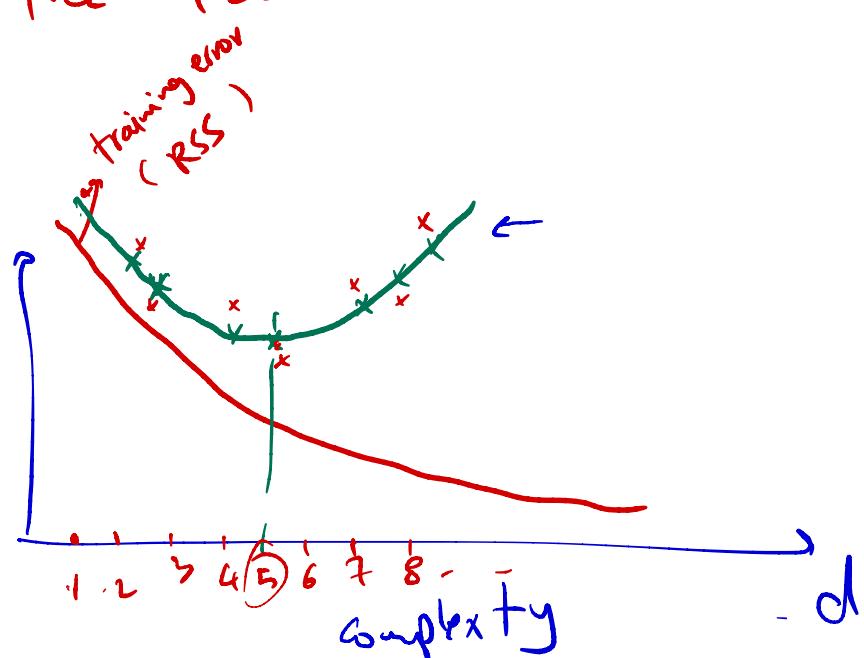
- Underfitting:
  - the model is not sufficiently complex to explain the data (high bias)
  - the trained model does not change much if we add new training data points (low variance)
- Overfitting:
  - model is too complex for training the data (low bias)
  - the trained model will change significantly if we add new data points (high variance)

## Lecture 16 :



the best  
fit / model with  
degree 100

Question : How can we estimate the  
the test error / curve ?

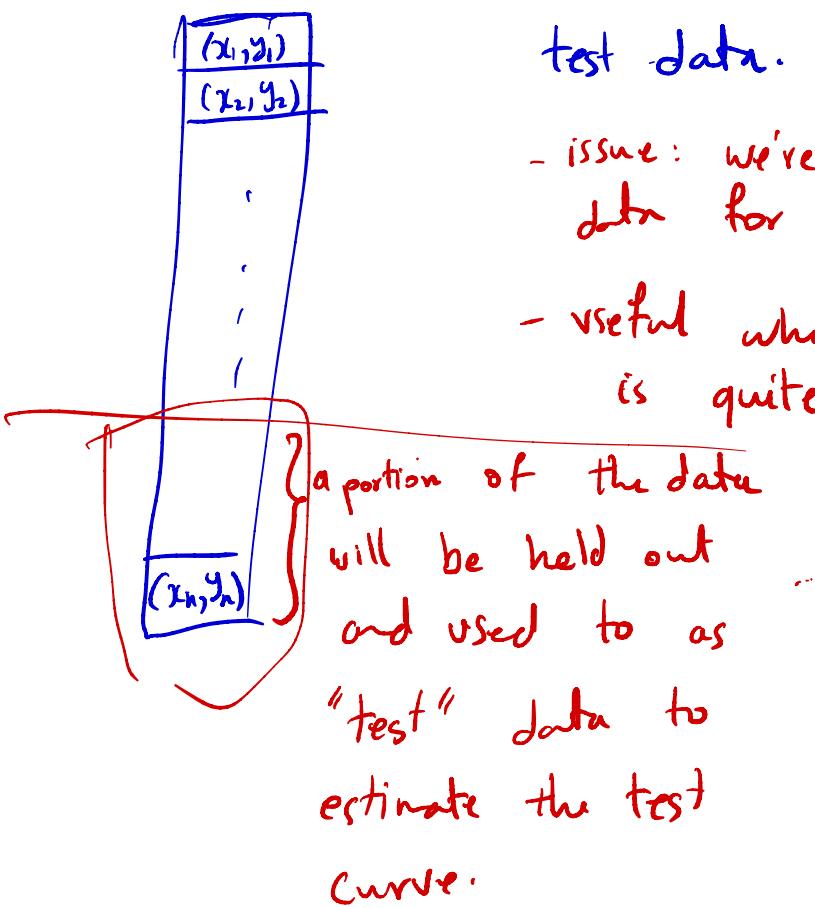


Two methods / answers :

(1) validation / hold-out data:

Randomly choose a (given) portion of the data and use it as test data. Use the rest for training.

- issue: we're not using all the available data for training
- useful when the size of the data is quite large



(2) k-fold cross-validation:

partition / Divide the data randomly into  $k$  parts:

$$D = D_1 \cup D_2 \cup D_3 \dots \cup D_k$$

for each  $i=1, \dots, k$

- train on  $D \setminus D_i$
- test the resulting model on  $D_i$

$\Rightarrow \text{test}_i = \text{test error of the trained model on } D_i$

$$\text{test-error} \approx \frac{1}{k} (\text{test}_1 + \text{test}_2 + \dots + \text{test}_k)$$

---

$\Rightarrow$  both of the methods are used to estimate the test curve

- once the "best complexity" is found (minimum in the estimated test curve),  $d^{\text{best}}$  we will use the whole data set to train a model with the "best complexity"
- cross validation is done for each complexity to estimate the test error.



# Classification:

Regression :

label  $y_i \in \mathbb{R}$

Classification:

label  $\in$  Discrete set

$y_i \in \{0, 1\}$

$\in \{\text{Spam, not-spam}\}$

$\in \{\text{Sunny, Snow, rain}\}$

Example:

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

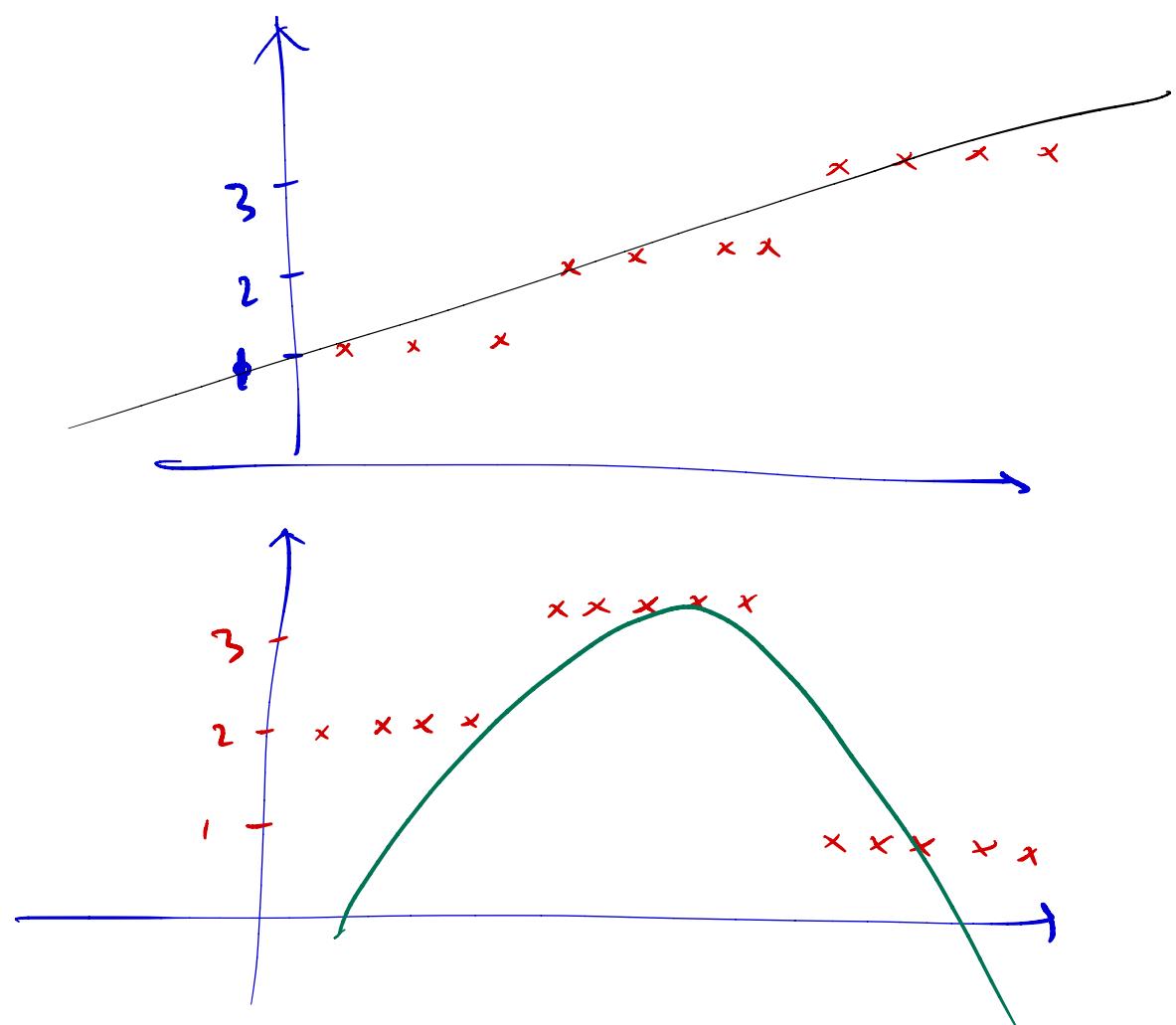
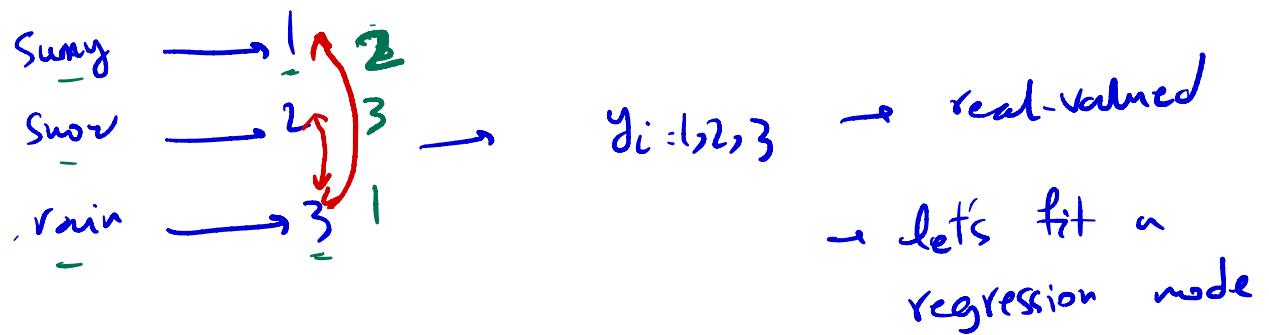
$x_i$  = weather data for day  $i$  ((clouds, --))

$y_i = \{\text{sunny, snowy, rainy}\}$

$$\underline{y_i} \approx f(\underline{x_i})$$

(remember, we want to be accurate at test time.)

As a naive approach, let's try model/solve  
a classification problem using regression:



Example:

input  $x$  (input data)  
weather data

$$\left. \begin{array}{l} \Pr\{\text{Sunny} \mid x\} = 0.2 \\ \Pr\{\text{Snow} \mid x\} = 0.5 \\ \Pr\{\text{Rain} \mid x\} = 0.3 \end{array} \right\}$$

- what we really need to learn/estimate  
is the "conditional probabilities":

$$\Pr\{\text{class}_k \mid x\} \quad \text{for } k=1,\dots,K$$

Bayes Optimal Classifier:

---

Assume for the moment that all the  
conditional probabilities are known

$$\Pr\{\text{class}_k \mid x\} \quad \text{known for } k=1,\dots,K.$$

How should we (optimally) predict the

class associated to input  $x$ ?

Bayes optimal classifier:

$$\begin{aligned}\hat{y}_{\text{Bayes}}(x) &= \text{predicted class}(x) \\ &= \underset{k}{\operatorname{argmax}} \Pr\{\text{class}_k | x\}\end{aligned}$$

(break ties arbitrarily)

Theorem: The Bayes optimal classifier has the smallest classification error among all the possible classifiers.

This means: for every other classifier

$\hat{y}_{\text{other}}(x)$ , we have:

$$\Pr_{(x,y)}\{\hat{y}_{\text{other}}(x) \neq y\} \geq \Pr_{(x,y)}\{\hat{y}_{\text{Bayes}}(x) \neq y\}$$

- The fundamental problem with the Bayes classifier is that it assumes the knowledge of the conditional probabilities (which not true in practice)
- All the methods for classification in ML aim at estimating the conditional Probabilities.
- We'll cover several methods to estimate these conditional probabilities from data.

## Logistic Regression:

For now, assume that we have two classes (the binary classification problem):

Class<sub>0</sub>  $\rightarrow y_i = 0$

$x_i \in \mathbb{R}$

Class<sub>1</sub>  $\rightarrow y_i = 1$

Goal: given  $x \rightarrow \underbrace{\left\{ \begin{array}{l} \Pr\{0|x\} \\ \Pr\{1|x\} \end{array} \right\}}$

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

Step 1:  
parametric model

↓  
Step 2:  
fit the  
params to data

Step: we're going to design a parametric model for estimation

$$\underbrace{\Pr\{0|x\}}_{\in [0,1]}, \quad \underbrace{\Pr\{1|x\}}_{\in [0,1]}$$

Let's begin with linear models.

What is a linear model in the binary classification setting?

$$\Pr\{0|x\} = f(x; \underline{\text{parameters}})$$

$$\Pr\{1|x\} = 1 - f(x; \underline{\text{parameters}})$$

let's come up with simplest candidates for the parametric model.

- Perhaps the simplest model would be something like

$$\begin{cases} \Pr\{0|x\} = \beta_0 + \beta_1 x \\ \Pr\{1|x\} = 1 - (\beta_0 + \beta_1 x) \end{cases}$$

But, one can not enforce the constraint that the probabilities are between 0 and 1 using the above model (i.e. it's not possible to ensure that  $\beta_0 + \beta_1 x \in [0, 1]$ ).

Instead, we consider the following model:

$$\Pr\{0|x\} = \frac{e^{B_0 + B_1 x}}{1 + e^{B_0 + B_1 x}}$$

$$\Pr\{1|x\} = 1 - \Pr\{0|x\}$$

$$= \frac{1}{1 + e^{B_0 + B_1 x}}$$

It's easy to see that  $\frac{e^{B_0 + B_1 x}}{1 + e^{B_0 + B_1 x}} \in [0,1]$

# Lecture 17:

Step 1: Devise a parametric model:

$$\Pr\{0|x\} = \frac{e^{B_0 + B_1 x}}{1 + e^{B_0 + B_1 x}} \quad \leftarrow$$

$$\Pr\{1|x\} = 1 - \Pr\{0|x\}$$

Step 2:

Define an objective for the "goodness of fit" → minimize the objective → find  $\hat{B}_0, \hat{B}_1$

$$\Pr\{y_i | x_i; B_0, B_1\} \quad \begin{cases} y_i = 0 \rightarrow \frac{e^{B_0 + B_1 x_i}}{1 + e^{B_0 + B_1 x_i}} \\ y_i = 1 \rightarrow \frac{1}{1 + e^{B_0 + B_1 x_i}} \end{cases}$$

$$\Pr\{x_i, y_i | B_0, B_1\} = \Pr\{x_i\} \cdot \Pr\{y_i | x_i, B_0, B_1\}$$

$$P(A, B) = P(A) \cdot P(B | A)$$

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$

$x_i \quad y_i \quad x_i \quad y_i \quad x_i$

$$P(A, B) = P(A) \cdot P(B | A)$$

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$

$x_i \quad y_i \quad x_i \quad y_i \quad x_i$

as we learned in module 1, we can be used to estimate parameters of the distribution that generates the data

mle:

$$\begin{aligned} \text{lik}(\beta_0, \beta_1) &= \prod_{i=1}^n \Pr\{x_i, y_i \mid \beta_0, \beta_1\} \\ &= \prod_{i=1}^n \Pr\{x_i\} \Pr\{y_i \mid x_i; \beta_0, \beta_1\} \\ &= \prod_{i=1}^n \Pr\{x_i\} \cdot \underbrace{\prod_{i=1}^n \Pr\{y_i \mid x_i; \beta_0, \beta_1\}}_{\text{objective to be optimized}} \end{aligned}$$

$$\begin{aligned} \arg \max_{\beta_0, \beta_1} \text{lik}(\beta_0, \beta_1) &\quad \text{the same value for any choice of } \beta_0, \beta_1 \\ &= \arg \max_{\beta_0, \beta_1} \underbrace{\prod_{i=1}^n \Pr\{x_i\}}_C \cdot \underbrace{\prod_{i=1}^n \Pr\{y_i \mid x_i; \beta_0, \beta_1\}}_{} \\ &= \bigcup_{\beta_0, \beta_1} \arg \max_{\beta_0, \beta_1} \underbrace{\prod_{i=1}^n \Pr\{y_i \mid x_i; \beta_0, \beta_1\}}_{} \\ &= \arg \max_{\beta_0, \beta_1} \sum_{i=1}^n \log \Pr\{y_i \mid x_i; \beta_0, \beta_1\} \\ &= \arg \max_{\beta_0, \beta_1} \sum_{i: y_i=0} \log \left( \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \right) + \sum_{i: y_i=1} \log \left( \frac{1}{1 + e^{\beta_0 + \beta_1 x_i}} \right) \end{aligned}$$

$$= \underset{B_0, B_1}{\operatorname{argmax}} \left\{ \sum_{i:y_i=0} (B_0 + B_1 x_i) - \sum_{i=1}^n \log (1 + e^{B_0 + B_1 x_i}) \right\}$$

$$\underline{g(B_0, B_1)}$$

↳ convex in terms of  
 $B_0, B_1$

$$\left\{ \begin{array}{l} \frac{\partial g}{\partial B_0} = 0 \\ \frac{\partial g}{\partial B_1} = 0 \end{array} \right.$$

- two equations in terms  
of the two unknowns  $B_0, B_1$ .

$$\Rightarrow$$

- no closed-form solution.  
⇒ we need to solve it  
numerically / iteratively

- the solution is unique.

$\Rightarrow \hat{B}_0, \hat{B}_1 \leftarrow$  the best fit  
to the data.

given a new  
data point  $x$

$$\hat{P}_{\text{o}}\{0|x\} = \frac{e^{\hat{B}_0 + \hat{B}_1 x}}{1 + e^{\hat{B}_0 + \hat{B}_1 x}}$$

$$\hat{P}_{\text{i}}\{1|x\} = \frac{1}{1 + e^{\hat{B}_0 + \hat{B}_1 x}}$$

$$\rightarrow \hat{y}_{\text{logistic-reg}}(x) = \underset{y \in \{0,1\}}{\operatorname{argmax}} \hat{P}_{\text{y|x}}$$

- logistic-regression for  $\underline{x} \in \mathbb{R}^P \rightarrow \underline{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_P \end{pmatrix} \in \mathbb{R}^P$

- binary labels:  $y \in \{0, 1\}$

$$\Pr\{y=0|x\} = \frac{e^{B_0 + B_1 x_1 + \dots + B_p x_p}}{1 + e^{B_0 + B_1 x_1 + \dots + B_p x_p}}$$

$$\Pr\{y=1|x\} = 1 - \Pr\{y=0|x\}$$

- general label set:

$y \in \{0, 1, 2\}$

$$\Pr\{y=i|x\} = \frac{e^{\hat{B}_0^i + \hat{B}_1^i x_1 + \dots + \hat{B}_p^i x_p}}{\sum_{j=0,1,2} e^{\hat{B}_0^j + \hat{B}_1^j x_1 + \dots + \hat{B}_p^j x_p}}$$

logistic regression is typically used for settings with binary labels ( $y \in \{0, 1\}$ ), and it does not perform well when we have multiple classes.

## - Linear / Quadratic Discriminant Analysis:

### - Linear Discriminant Analysis (LDA):

$$\Pr\{y \mid x\}$$

logistic regression directly models this conditional probability.

LDA/QDA try to estimate the conditional probabilities through a different route  $\rightarrow$  Bayes rule.

$$\Pr\{y=k \mid x=x\} = \frac{\text{Bayes rule}}{\Pr\{x=x \mid y=k\} \cdot \Pr\{y=k\}}$$

LDA/QDA aim at estimating these probabilities.

$$\Pr\{x=x\}$$

$$\Pr\{x=x\} = \sum_{k=1}^K \underbrace{\Pr\{y=k\}}_{\Pr\{y=k\mid x=x\}} \underbrace{\Pr\{x=x \mid y=k\}}_{\Pr\{x=x\mid y=k\}}$$

Assume K classes:

$$\Pr\{y=k \mid x=x\} =$$

$$\frac{f_k(x) = \text{gaussian}(\mu_k) e^{-\frac{(x-\mu_k)^2}{2\sigma^2}}}{\Pr\{x=x\}}$$

$$\Pr\{x=x \mid y=k\}$$

$$\Pr\{y=k\}$$

$$\Pr\{x=x\}$$

$$\Pr\{y=k\}$$

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

$$\text{example } (x_i \in \mathbb{R}, y_i \in \{-, +\})$$

$$\Pr\{x=x \mid y=+\}$$



$$\Pr\{y=+\} \stackrel{\text{estimate}}{=} \frac{\sum_{i: y_i=+} 1}{n}$$

$$\Pr \{ y = k \} \approx \frac{\# \{ i : y_i = k \}}{n}$$

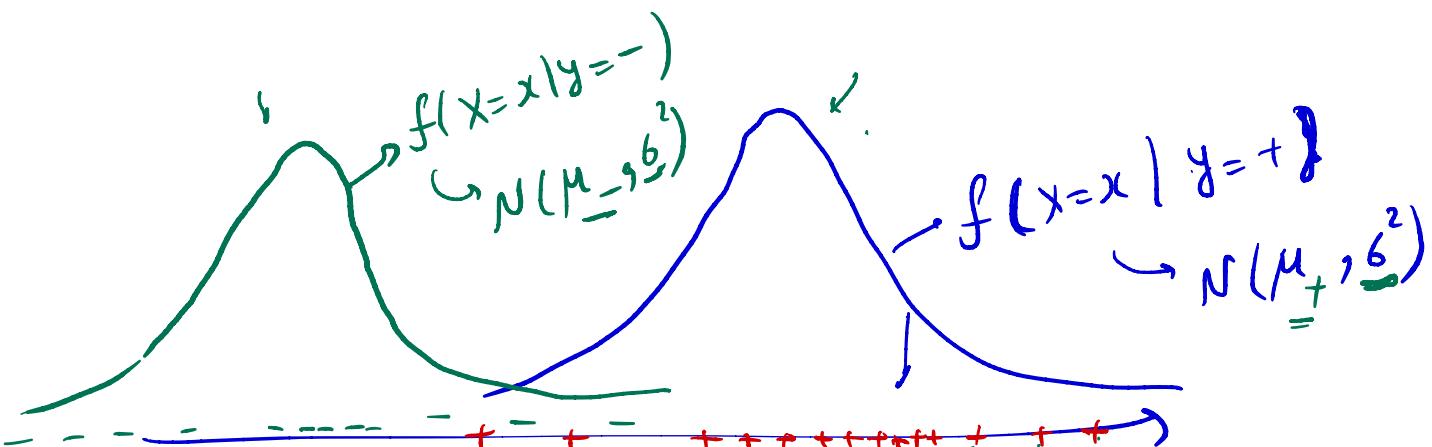
↓

Unbiased estimate of  
 $\Pr \{ y = k \}$

$$f \{ x = x \mid y = k \} \stackrel{\text{LDA}}{=} f_k(x)$$

$$= \frac{1}{\sqrt{2\pi}^6} e^{-\frac{(x-\mu_k)^2}{2\sigma^2}}$$

LDA "assumes" / "models" that data, given a class  $k$ , is generated according to the gaussian distribution.



In LDA, we assume that the variance of the data  $x$  is the same over all the classes :

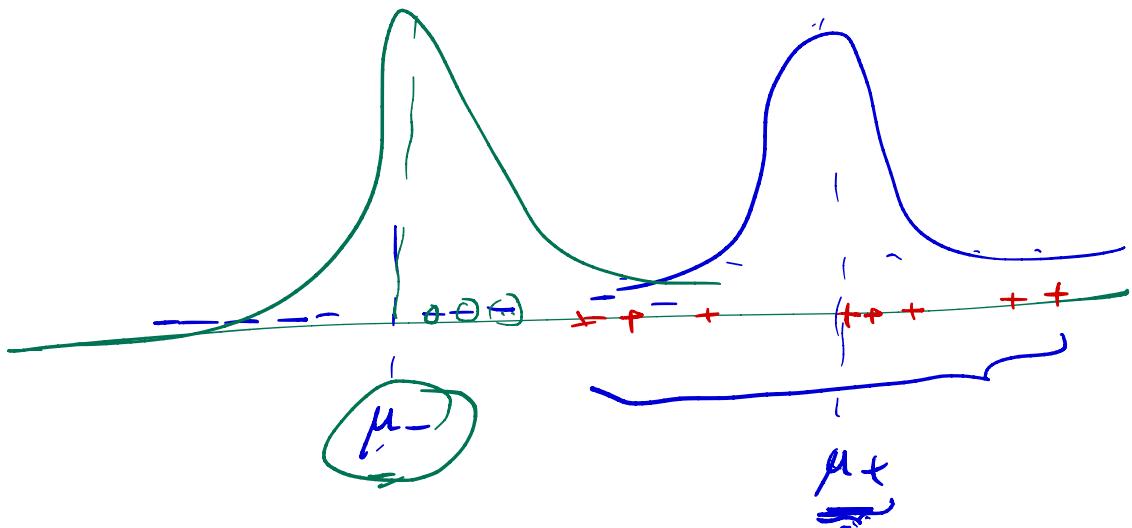
For a gaussian, we need to estimate two parameters :  $\mu_k, \sigma$ .

$$\hat{\mu}_k \xrightarrow{\text{mean of the density}} \approx \frac{\sum_{i: y_i=k} x_i}{n_k} \rightarrow \text{Sample mean}$$

$$n_k = \#\{ i : y_i=k \}$$

Let's now estimate  $\sigma$ :

$$\hat{\sigma}^2 = \frac{1}{n-k} \sum_{k=1}^K \underbrace{\sum_{i: y_i=k} (x_i - \hat{\mu}_k)^2}_{\text{variance across class } k}$$



\* LDA uses the same  $\sigma^2$  for each class,  $\Rightarrow$  LDA requires to learn less parameters than the case where we assign different  $\sigma^2$ 's to the classes.

$\Rightarrow$  QDA  $\rightarrow$  uses different  $\sigma^2$ 's for the classes.

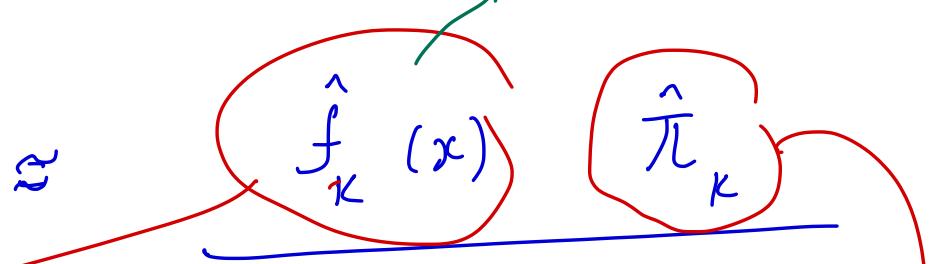
LDA :

*MR*

## Lecture 18 :

$$\text{LDA} \longrightarrow \Pr\{y=k \mid x=x\}$$

$N(\mu_k, \hat{\sigma}^2)$



$$\Pr\{x=x\}$$

$$\hat{\pi}_k = \frac{\#\{i: y_i=k\}}{n} \quad (\text{an estimate for } \Pr\{y=k\})$$

$$\frac{1}{\sqrt{2\pi}\hat{\sigma}} e^{-\frac{(x-\hat{\mu}_k)^2}{2\hat{\sigma}^2}}$$

(an estimate for  $\Pr\{x=x \mid y=k\}$ )

At prediction time, given a new input

$x:$

$$\hat{y}(x) = \operatorname{argmax}_k \Pr\{y=k \mid x=x\}$$

$$= \arg \max_K \frac{\hat{f}_K(x) - \hat{\pi}_K}{\text{(P.s.) } x - \bar{x}}$$

the same for every class  $K$

$$\begin{aligned} &= \arg \max_K \hat{f}_K(x) \hat{\pi}_K \\ &= \arg \max_K \left\{ \log \hat{f}_K(x) + \log \hat{\pi}_K \right\} \\ &= \arg \max_K \left\{ \log \frac{1}{\sqrt{2\pi}\hat{\sigma}} e^{-\frac{(x-\hat{\mu}_K)^2}{2\hat{\sigma}^2}} + \log \hat{\pi}_K \right\} \\ &= \arg \max_K \left\{ \log \frac{1}{\sqrt{2\pi}\hat{\sigma}} - \frac{(x-\hat{\mu}_K)^2}{2\hat{\sigma}^2} + \log \hat{\pi}_K \right\} \\ \text{Because of} \\ \text{the assumption} \\ \text{that the variances} \\ \text{were the same} \\ \text{over all} \\ \text{the classes,} &= \arg \max_K \left\{ \log \frac{1}{\sqrt{2\pi}\hat{\sigma}} + \frac{x^2}{2\hat{\sigma}^2} - \frac{\hat{\mu}_K^2}{2\hat{\sigma}^2} + \frac{x\hat{\mu}_K}{\hat{\sigma}^2} + \log \hat{\pi}_K \right\} \\ &\quad \text{do not depend} \\ &\quad \text{on the class } K \\ &= \arg \max_K \left\{ -\frac{\hat{\mu}_K^2}{2\hat{\sigma}^2} + \frac{x\hat{\mu}_K}{\hat{\sigma}^2} + \log \hat{\pi}_K \right\} \\ &\quad \text{linear in } x \end{aligned}$$

- Since the decision rule is linearly dependent on  $x$ , the classifier is called "linear" discriminant analysis.

Extension to the multi-dimensional setting:

$$\text{Data} = \left\{ (x_i, y_i) \right\}_{i=1}^n \quad \begin{cases} y_i \in \{1, \dots, k\} \\ x_i \in \mathbb{R}^p \end{cases}$$

$$\Pr\{y=k \mid x=x\} = \frac{\Pr\{x=x \mid y=k\} \cdot \Pr\{y=k\}}{\Pr\{x=x\}}$$

$f_k(x) \sim N(\tilde{\mu}_k, \tilde{\Sigma})$       the same over all classes  
 $\pi_{ik}$        $\Pr\{y=k\}$

After

# Multi-variate Gaussian Distribution

$$N(\mu, \Sigma)$$

↓  
 $\in \mathbb{R}^p$   
 p x p matrix

$$f(x | \mu, \Sigma)$$

↑ mean  
 ↑ covariance matrix

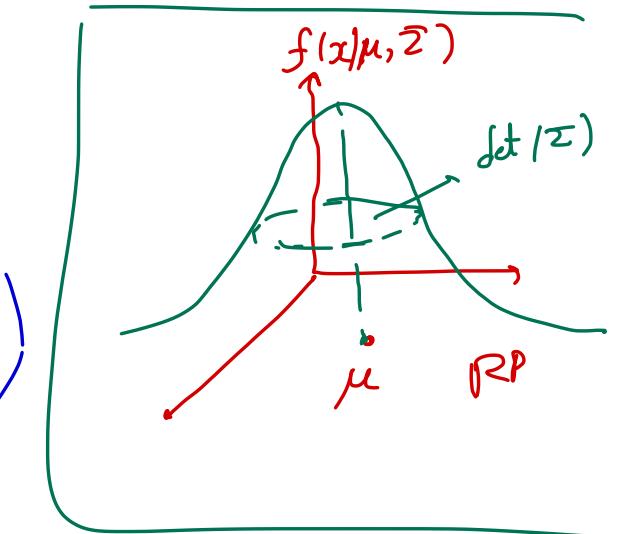
$$= \frac{1}{(2\pi)^{p/2} \cdot \det(\Sigma)^{1/2}} \exp \left( -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right)$$

Transpose

properties:

$$X \sim N(\mu, \Sigma)$$

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix} \in \mathbb{R}^p$$



$$E[X] = \mu, \quad E[(X-\mu)(X-\mu)^T] = \Sigma$$

$$\forall r, s \in \{1, \dots, p\}$$

$$E[(X_r - \mu_r)(X_s - \mu_s)] = \sum_{r,s}$$

$$X = \begin{pmatrix} X_1 \\ \vdots \\ X_r \\ \vdots \\ X_s \\ \vdots \\ X_p \end{pmatrix} \quad \begin{aligned} E[X_r] &= \mu_r \\ E[X_s] &= \mu_s \\ E[(X_r - \mu_r)(X_s - \mu_s)] & \end{aligned}$$

$$= \sum_{r,s} \underbrace{\begin{pmatrix} & & & s \\ & & & | \\ & & - & - \sum_{rs} \\ r & - & - & \end{pmatrix}}_{\text{matrix } \Sigma}^{p \times p}$$

$$\pi_k = \Pr\{Y=k\} \xrightarrow{\text{estimate}} \hat{\pi}_k = \frac{\#\{i : y_i=k\}}{\text{total number of data points}}^{\overbrace{n_k}}$$

$$f_k(x) = N(\mu_k, \Sigma)$$

$\hat{\mu}_k = \frac{\sum_{i:y_i=k} x_i}{n_k}$   
 $n_k \rightarrow \text{number of data points with label } k$

$$\hat{\Sigma} = r \begin{pmatrix} s & & & & \\ \vdots & \ddots & \ddots & \ddots & \\ \vdots & & \hat{\Sigma}_{rs} & & \\ \vdots & & & \ddots & \ddots \end{pmatrix}_{P \times P}$$

$$\hat{\Sigma}_{rs} = \frac{\sum_{k=1}^K \sum_{i:y_i=k} (x_{i,r} - \hat{\mu}_{k,r})(x_{i,s} - \hat{\mu}_{k,s})}{n-K}$$

we estimate  $\hat{\mu}_k$  and  $\hat{\Sigma}$  and  $\hat{\pi}_k$   
from data

⇒ an estimate for  
the conditional probabilities

$$- \Pr \left\{ y=k \mid X=x \right\} = \frac{\Pr \{ X=x \mid y=k \} \Pr \{ y=k \}}{\Pr \{ X=x \}}$$

estimated  
 $\hat{f}(x | \hat{\mu}_k, \hat{\Sigma})$

$$\frac{\hat{f}(x | \hat{\mu}_k, \hat{\Sigma})}{\Pr \{ X=x \}} \hat{\pi}_k$$

z.p.

- Prediction of . The label at a new data point  $x$ :

$$\hat{y}(x) = \underset{k}{\operatorname{argmax}} \left\{ \hat{f}(x | \hat{\mu}_k, \hat{\Sigma}) \hat{\pi}_k \right\}$$

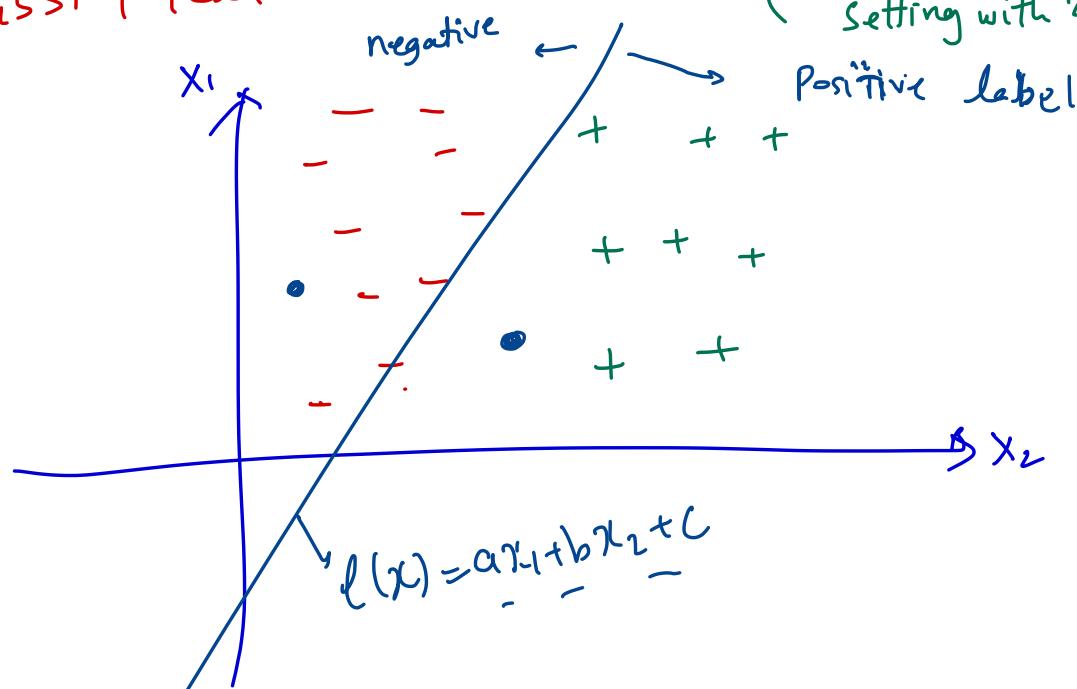
$$= \underset{k}{\operatorname{argmax}} \left\{ x^T \hat{\Sigma}^{-1} \hat{\mu}_k - \frac{1}{2} \hat{\mu}_k^T \hat{\Sigma}^{-1} \hat{\mu}_k + \log \hat{\pi}_k \right\}$$

$\ell_k(x)$

$\ell_k(x) \rightarrow \text{linear in } x$

This will lead to what we call linear classification boundaries.

(let's consider a binary classification setting with 2-d data)



LDA:  $\operatorname{argmax} \{ \ell_+(x), \ell_-(x) \}$

↓

$\left\{ \begin{array}{l} \text{if } l_+(x) > l_-(x) \Rightarrow \hat{y}(x) = + \\ \text{if } l_+(x) \leq l_-(x) \Rightarrow \hat{y}(x) = - \end{array} \right.$

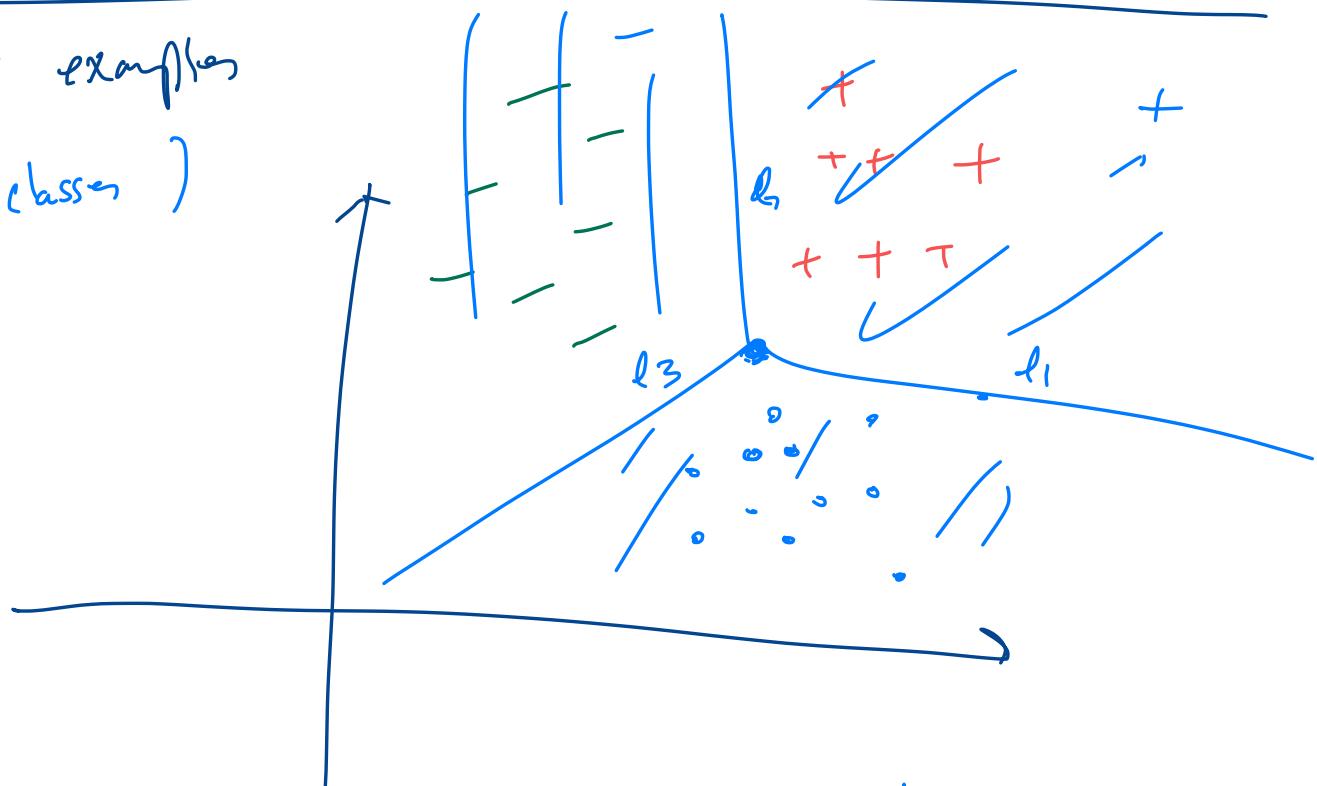
↓

$$l_+(x) - l_-(x) \geq 0 \quad \hat{y}(x) = +$$

$$\triangleq l(x) \quad ax_1 + bx_2 + c$$

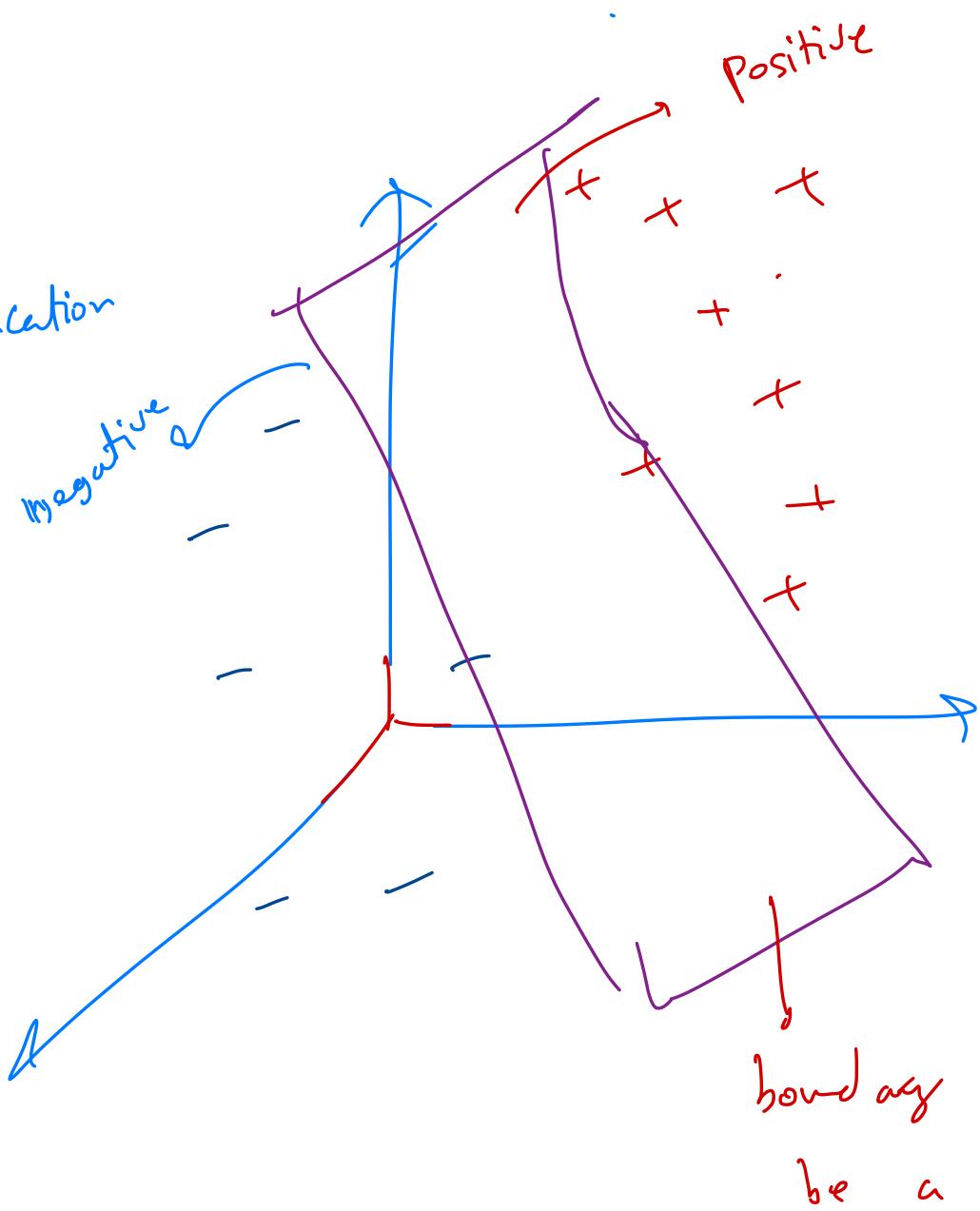
$$\hat{y}(x) = -$$

other examples  
(3 classes)

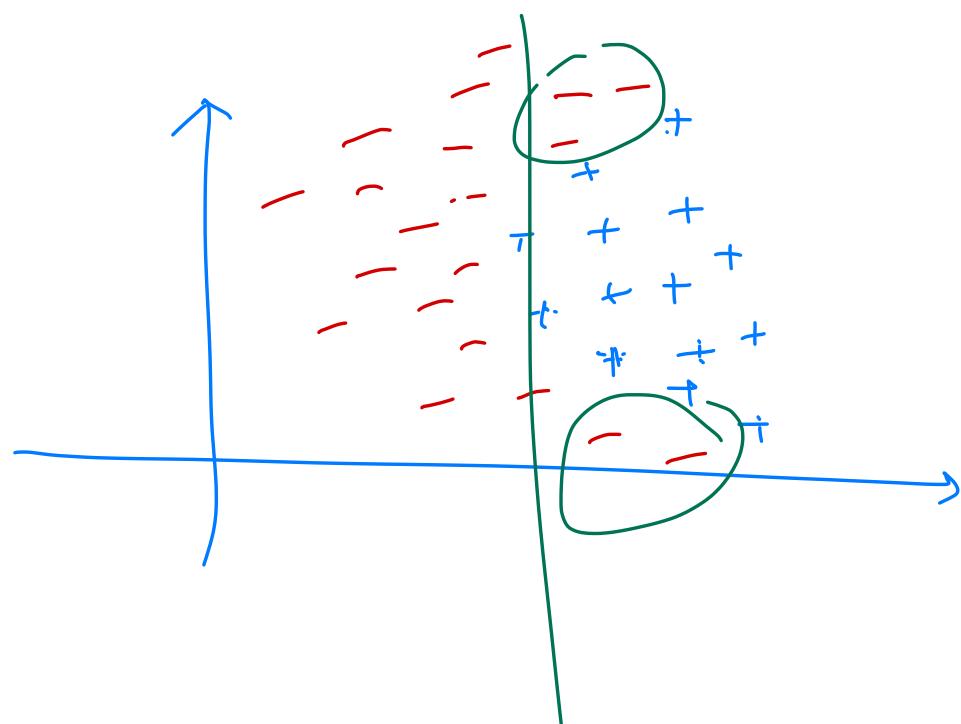


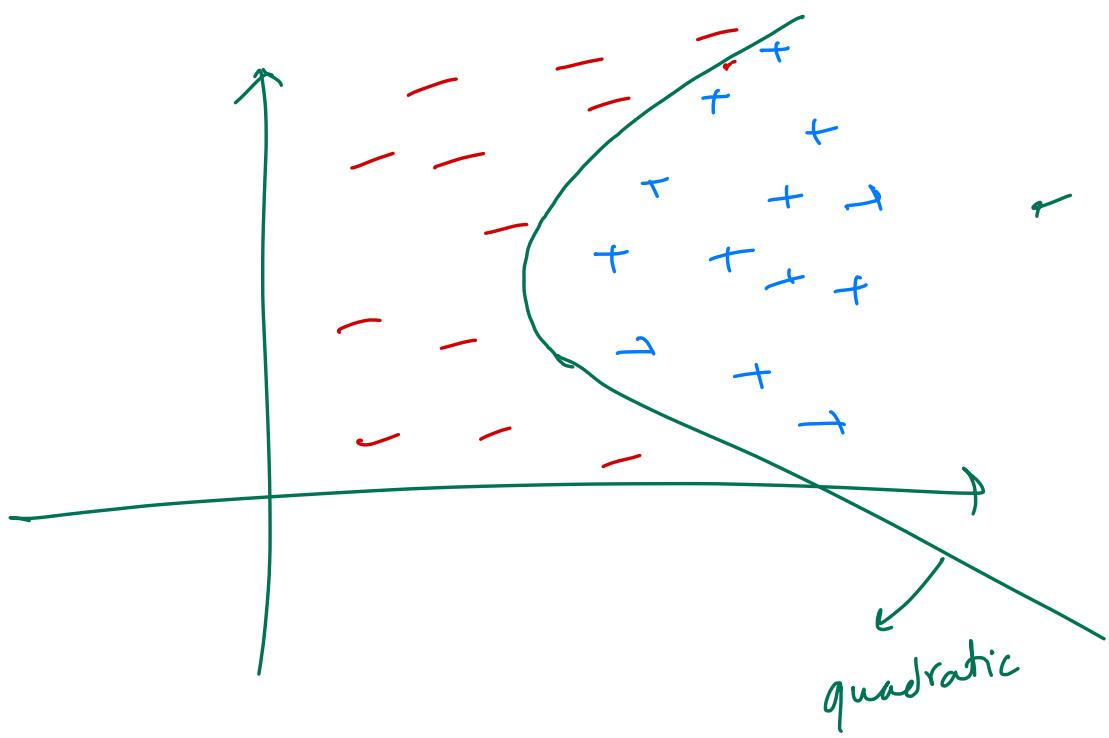
$$\hat{y}(x) = \arg \max \{ l_+(x), l_-(x), l_*(x) \}$$

$\mathbb{H}$   
3-dim  
binary  
classification



In what cases can LDA be ineffective?





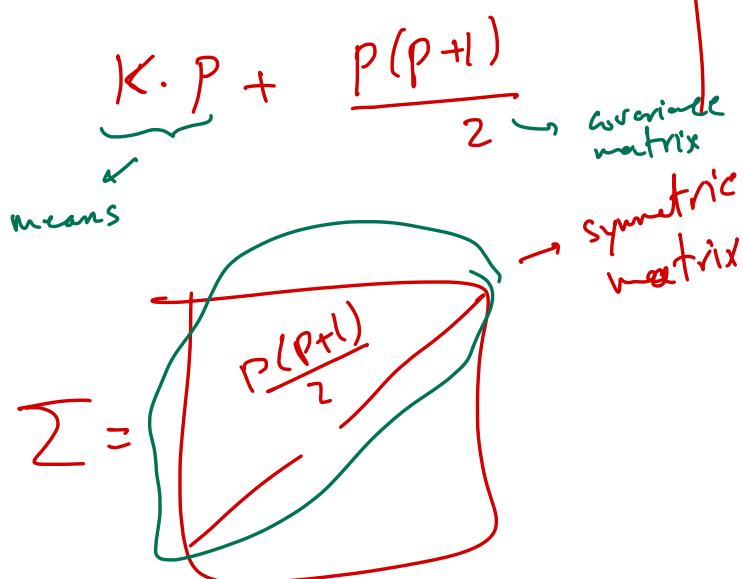
## Quadratic Discriminant Analysis:

- Similar to LDA except that the variances can change across the classes.
  - for each class  $k$  :  $f_k(x) = \mathcal{N}(\mu_k, \Sigma_k)$
- learned using  
 data from  
 class  $.k$ .
- quadratic term
- $$\hat{y}_{QDA}(x) = \underset{k}{\operatorname{argmax}} \left\{ -\frac{1}{2} (x - \hat{\mu}_k)^T \hat{\Sigma}_k^{-1} (x - \hat{\mu}_k) - \frac{1}{2} \log \det(\hat{\Sigma}_k) + \log \hat{\pi}_k \right\}$$

- Why QDA is a more complex class than LDA ?

LDA

→ Learn the mean vector per class + the covariance matrix.

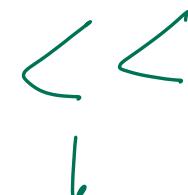


QDA

- 1 mean. vector per class
- 1 covariance matrix per class

$$k \cdot P + k \frac{P(P+1)}{2}$$

$$k_p + \frac{P(P+1)}{2}$$



$$kP + \frac{kP(P+1)}{2}$$

when  $p$  is large

LDA

- less parameters
- less flexibility / complexity

might underfit

QDA

might, overfit

LDA is a better bet than QDA if

there are relatively few training data

points. In contrast, QDA is a better

choice if the number of training

data points is large or if

the assumption of a common covariance

matrix for all the  $k$  classes is off.

(in this case LDA)  
will underfit

## Module 4 : Unsupervised learning:

Supervised learning:  $\left\{ \begin{matrix} (x_i, y_i) \\ = = \end{matrix} \right\}_{i=1}^n \rightarrow y_i \approx \tilde{f}(x_i)$

unsupervised learning  $\left\{ x_i \right\}_{i=1}^n$

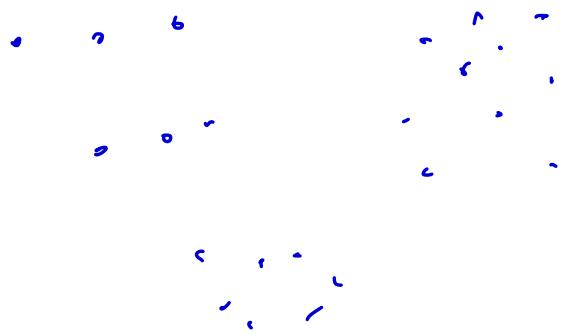
## Lecture 19.

### Module 4: Unsupervised Learning:

- There is no label (in contrast with supervised ML)

Data:  $x_1, x_2, \dots, x_n$

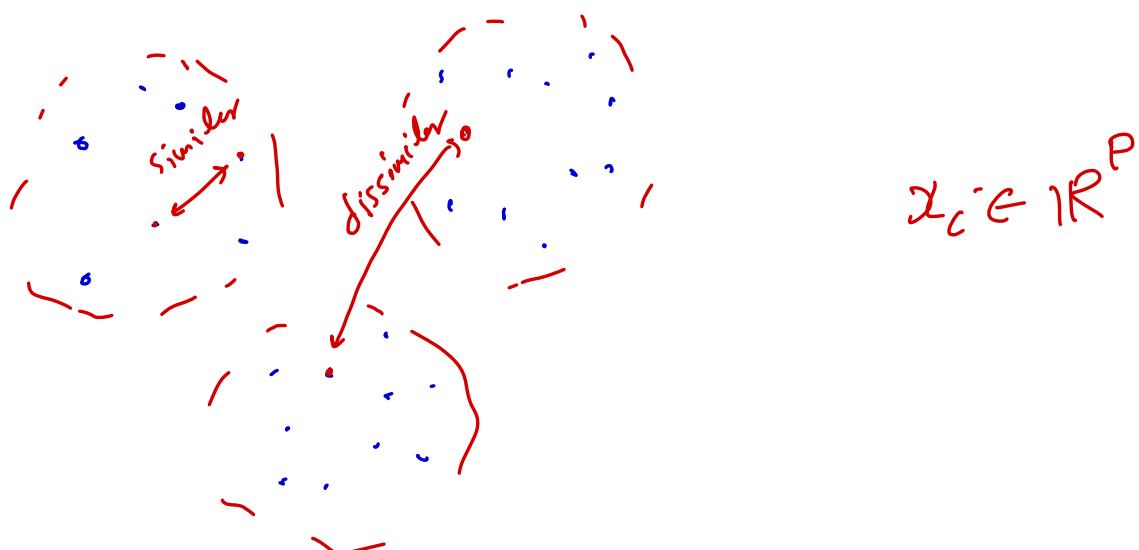
Goal: Discover informative patterns, structures, subgroups, etc within data.



In this course, we'll consider two specific instances of unsupervised learning:

- clustering -
- Dimensionality Reduction

clustering: Partition data into distinct sub-groups such that data points within each group are similar to each other and points from different groups are "dissimilar".



"similarity" will be quantified by a distance function:

$$x \leftarrow d(x, y) \rightarrow y$$

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix} \in \mathbb{R}^P, y = \begin{pmatrix} y_1 \\ \vdots \\ y_p \end{pmatrix} \in \mathbb{R}^P$$

Example:

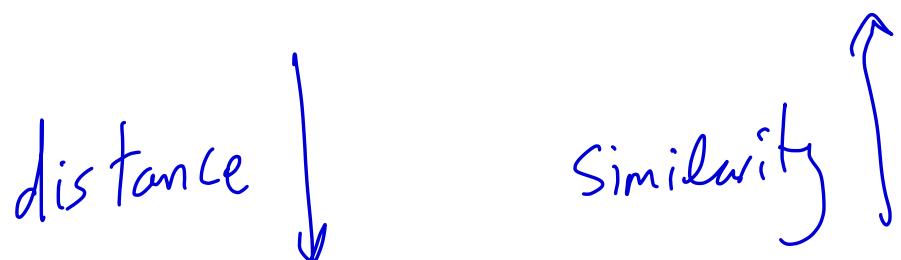
$$(1) \quad d(x, y) = \left( \sum_{i=1}^P (x_i - y_i)^2 \right)^{\frac{1}{2}} = \|x - y\|_2$$

- the Euclidean distance  
- the L<sub>2</sub> distance

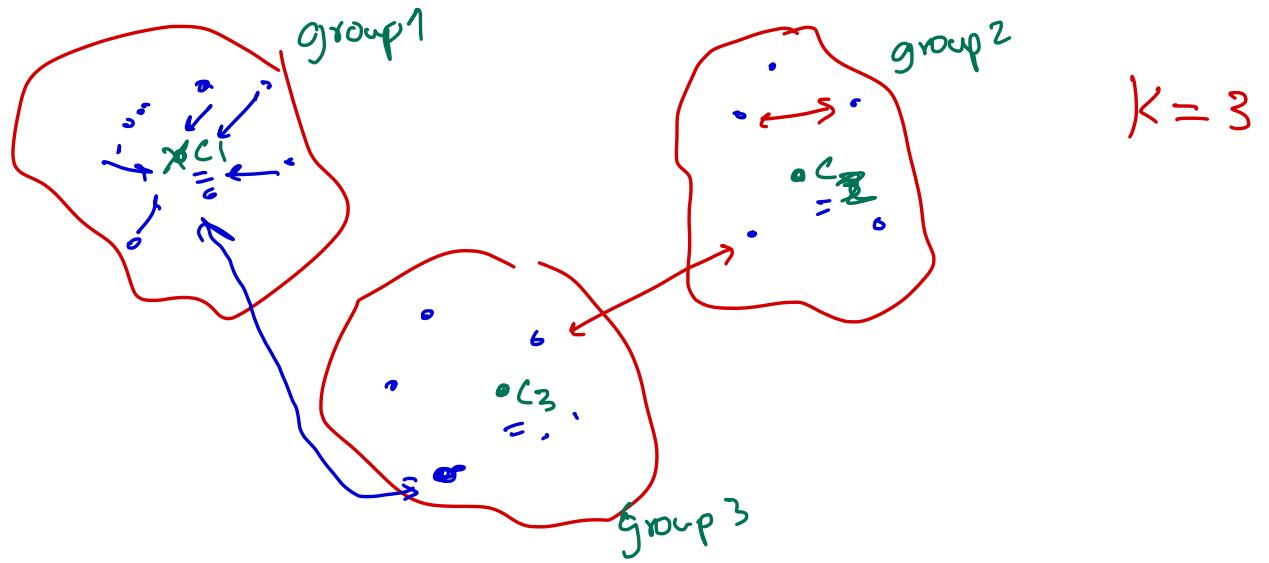
$$(2) d(x, y) = \sum_{i=1}^p |x_i - y_i| = \|x - y\|_1$$

-  $L_1$  distance

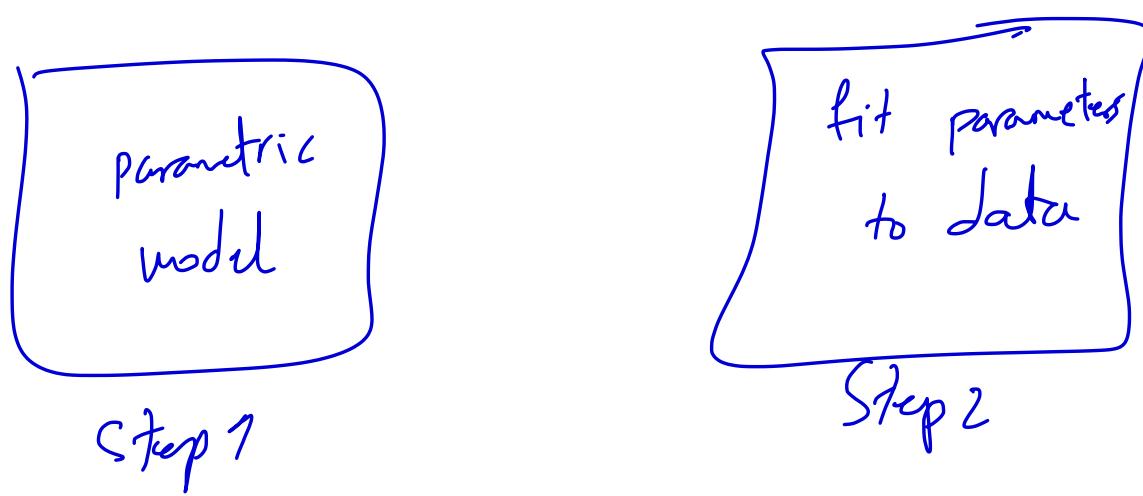
- the Manhattan distance



clustering: partition data into  $K$  subgroups ( $K$  is given) such that points in each group are similar (i.e. smaller distance) and points in different groups are dissimilar (relatively larger distance).



Remember our ~~our~~ general approach:



Step 1 (parametric modeling):

How do we represent a group?

- Assign to group  $i$  a "center",  $c_i \in \mathbb{R}^P$ , such that points in group  $i$  are close to the center  $c_i$  and points outside group  $i$  are far from  $c_i$ .

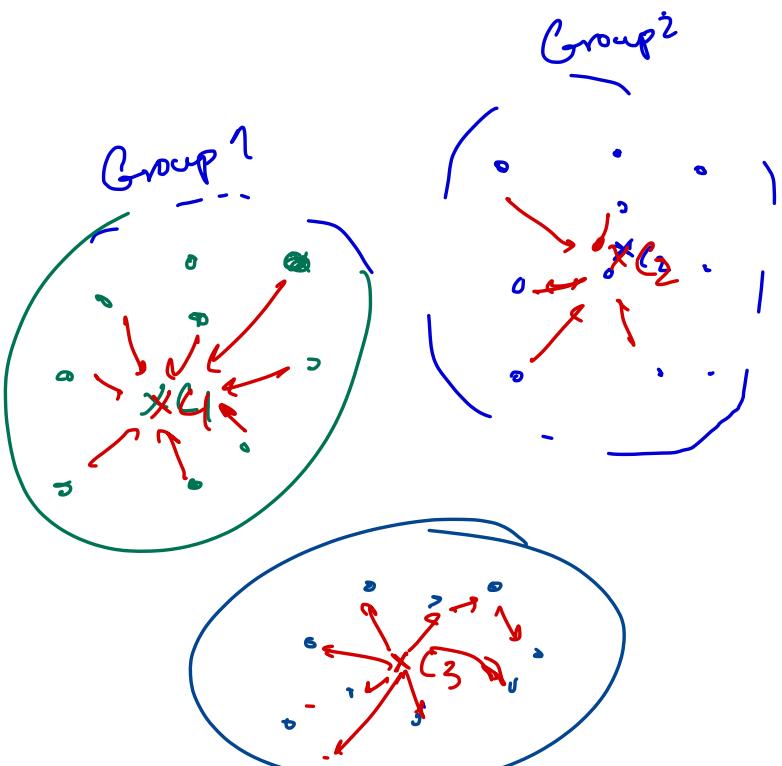
$\Rightarrow$  parameters of the clustering problem

$$\text{are } c_1, c_2, \dots, c_k \in \mathbb{R}^P$$

Step 2:

$$k = 3$$

$$c_1, c_2, c_3$$



$k$ -means objective

$$\min_{c_1, \dots, c_k} \sum_{j=1}^k \sum_{i \in \text{Group}_j} \|x_i - c_j\|^2 \quad (\text{k-means})$$

internal distance of Group $j$

$k$ -means clustering: solve ( $k$ -means) and find  $c_1, \dots, c_k$  along with the groups.

- Difficulty: Finding the exact solution of (k-means) is very difficult (it's an NP-hard problem)
- we'll describe an algorithm that finds an approximate (but sub-optimal) solution for the k-means problem which works well in practice!

## The K-means Algorithm:

---

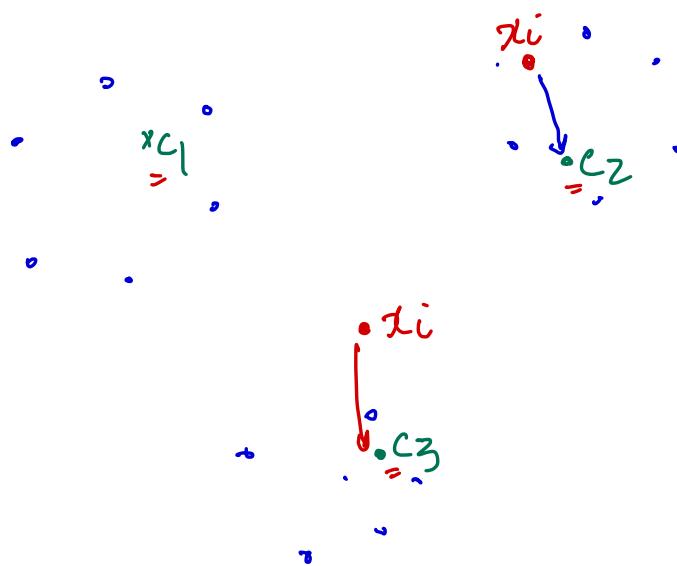
the k-means algorithm is based on two simple principles:

$$\begin{aligned}
 & \left( \sum_{j=1}^K \left( \sum_{x_i \in \text{Group}_j} \|x_i - c_j\|^2 \right) \right) \\
 & = \sum_{i=1}^n \|x_i - c(x_i)\|^2
 \end{aligned}$$

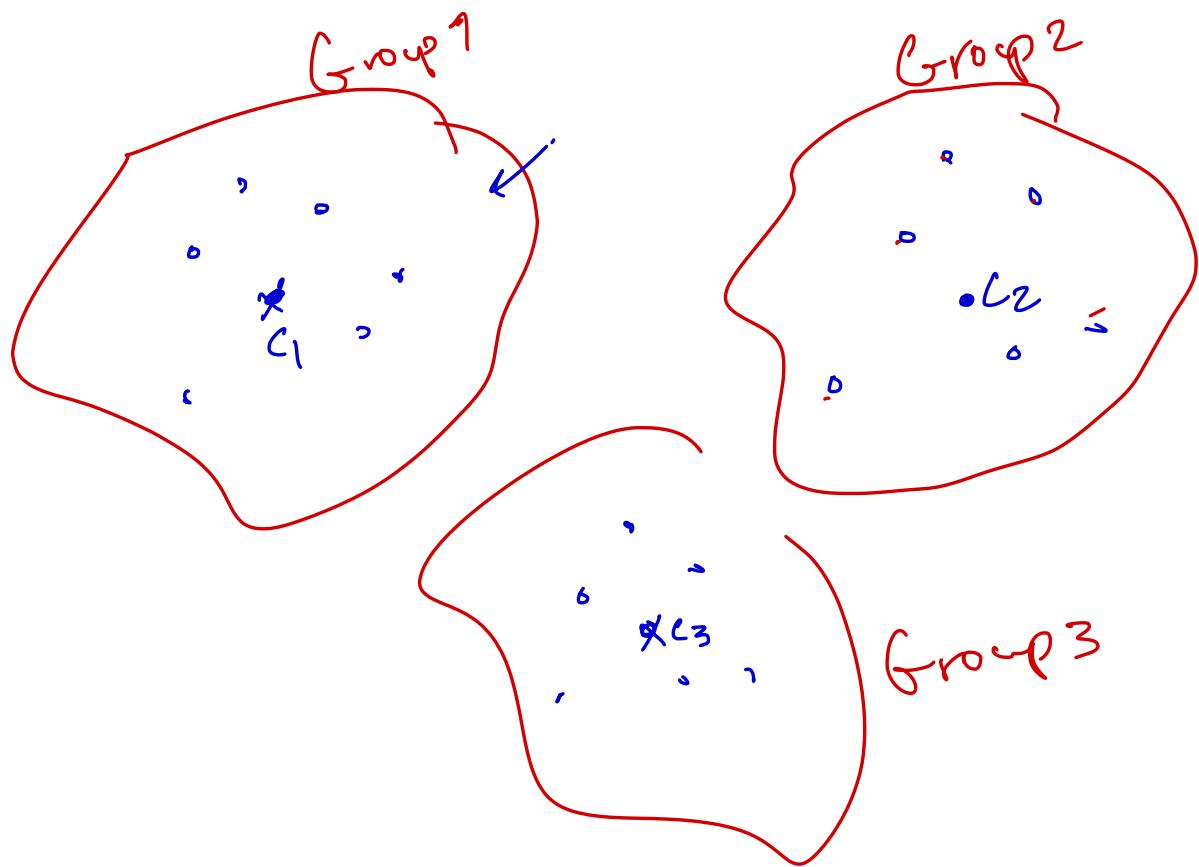
\$c(x\_i)\$ is the center of the group that \$x\_i\$ belongs to.

principle 1: Given centers \$c\_1, \dots, c\_K\$, the center (group) that we assign to data point \$x\_i\$ is the one that's closest to \$x\_i\$.

$$c(x_i) = \underset{j \in \{1, \dots, K\}}{\operatorname{argmin}} \|x_i - c_j\|$$



Principle 2: Let's assume that we fix the groups:



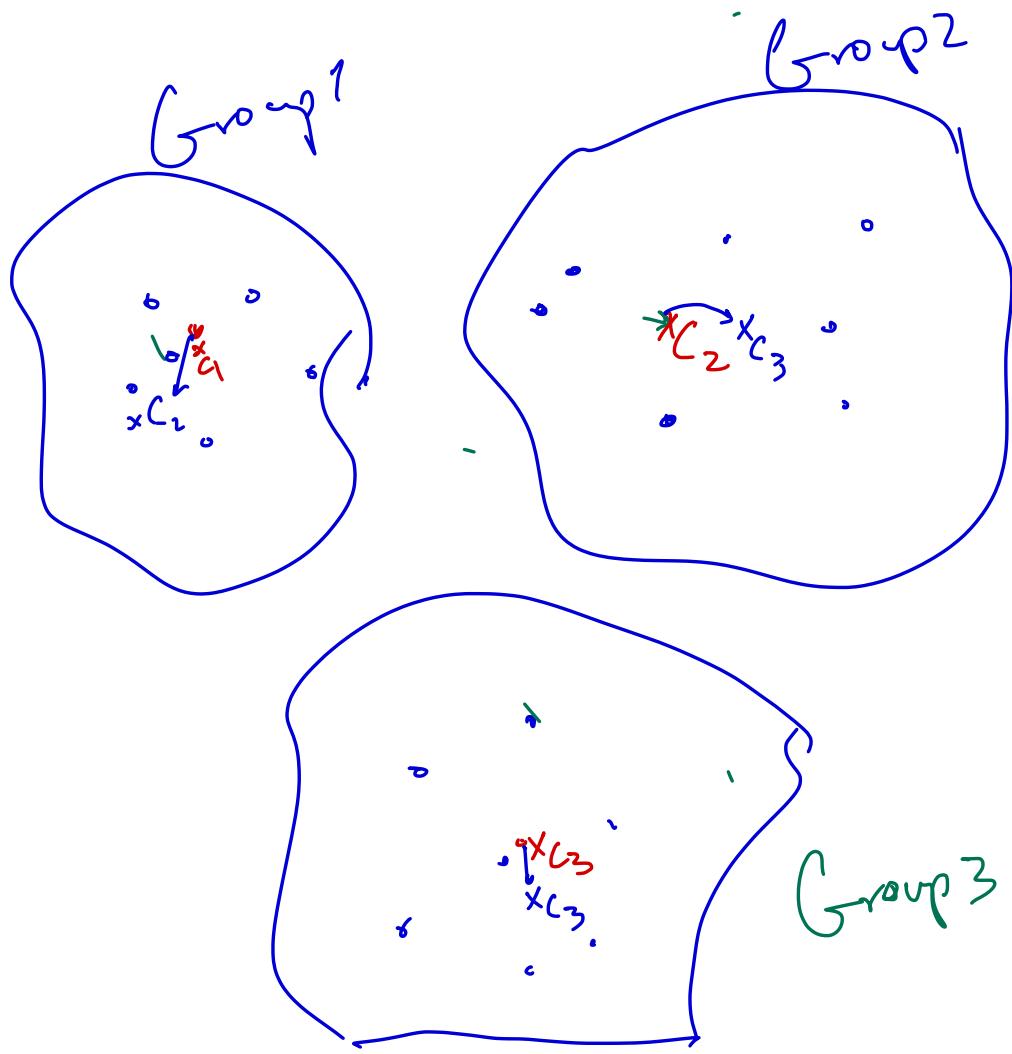
$$c_j = \frac{\sum_{x_i \in \text{group } j} x_i}{N_j}$$

$N_j = \# \text{ points in group } j$

let's try to solve the k-means problem

$$\min_{c_i \in \mathbb{R}^p} \sum_{i=1}^n \|x_i - c_i\|_2^2$$

$$\rightarrow c_1 = \frac{\sum x_i}{n}$$



$k$ -means algorithm:

- Choose  $\{ \underline{c_1^0, c_2^0, \dots, c_k^0} \}$  according to some (arbitrary) choice.
- For  $t=1, 2, 3, \dots$
- (i) Define Group  $j$  to be the set of all the data points that are closest to  $c_j^{t-1}$ .  $\rightarrow$  (principle 1)

$$(2) \text{ update } c_j^t = \frac{\sum_{x_i \in \text{Group}_j} x_i}{N_j} \rightarrow (\text{principle})$$

(3) if  $c_j^t = c_j^{t-1}$  for all  $j \in \{1 \rightarrow k\}$   
 we will stop!

Some important (practical) points about the K-means algorithm:

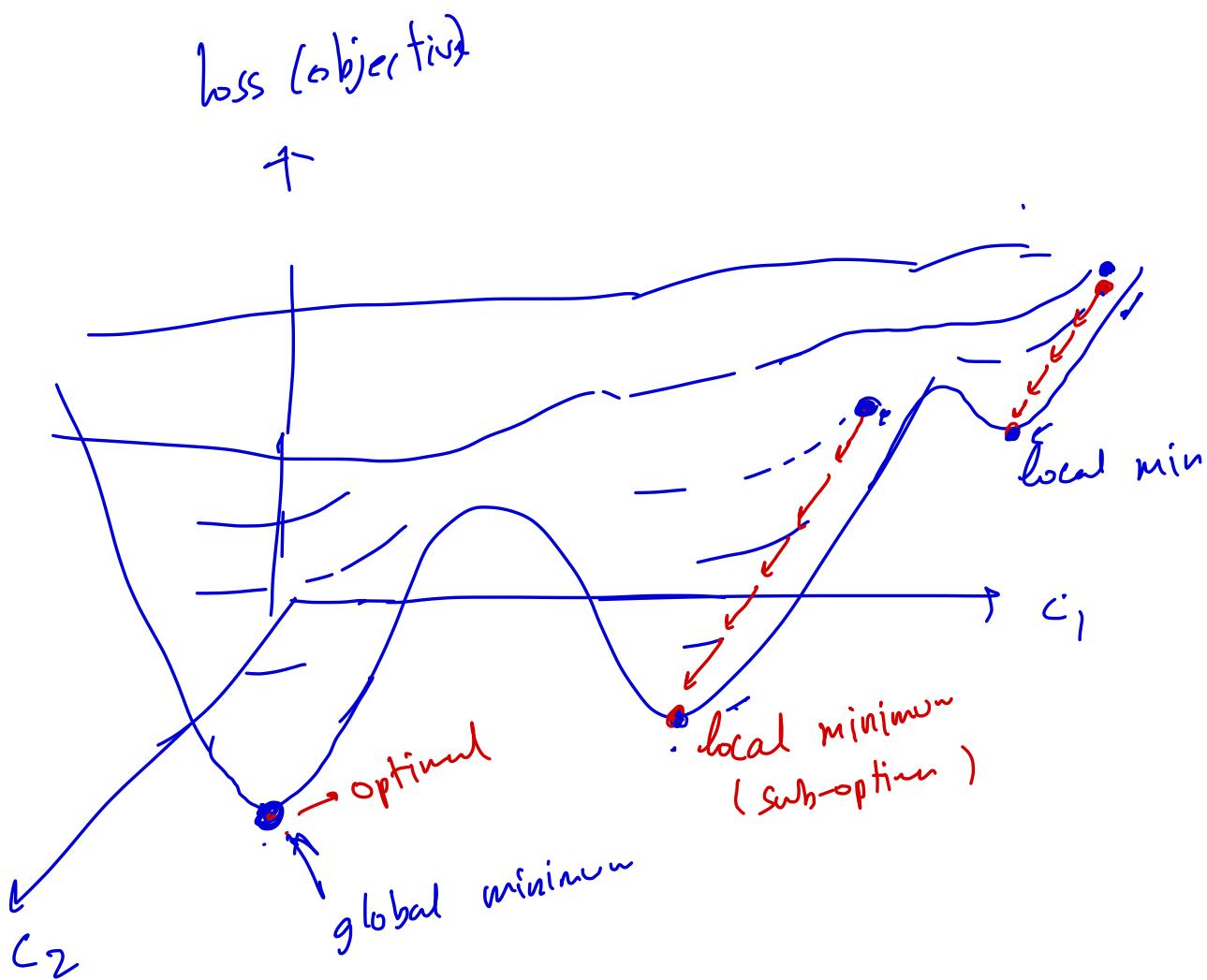
- Let  
 $L_t \triangleq L(\{c_1^t, \dots, c_k^t\}) = \sum_{i=1}^n \|x_i - c^t(x_i)\|^2$   
 objective computed  
 for  $\{c_1^t, \dots, c_k^t\}$

- then it can be proven that

$$L_t \leq L_{t-1}$$

- Also, it can be proven that

after a "finite" number of iterations,  
the k-means algorithm will stop  
(Converge.)



- The quality of the final set of centers that we obtain from the k-means algorithm depends on the initial centers.
- In practice, we often choose 10-20-50 (random) different initializations and run the k-means algorithm for each initialization, and choose the set of

Centers that has the best loss.

Centers that has the best loss.

## Lecture 20 :

unsupervised learning

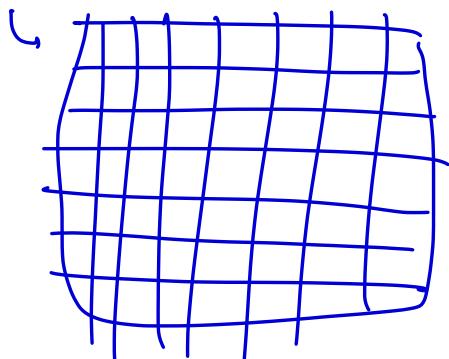
clustering

dimensionality reduction

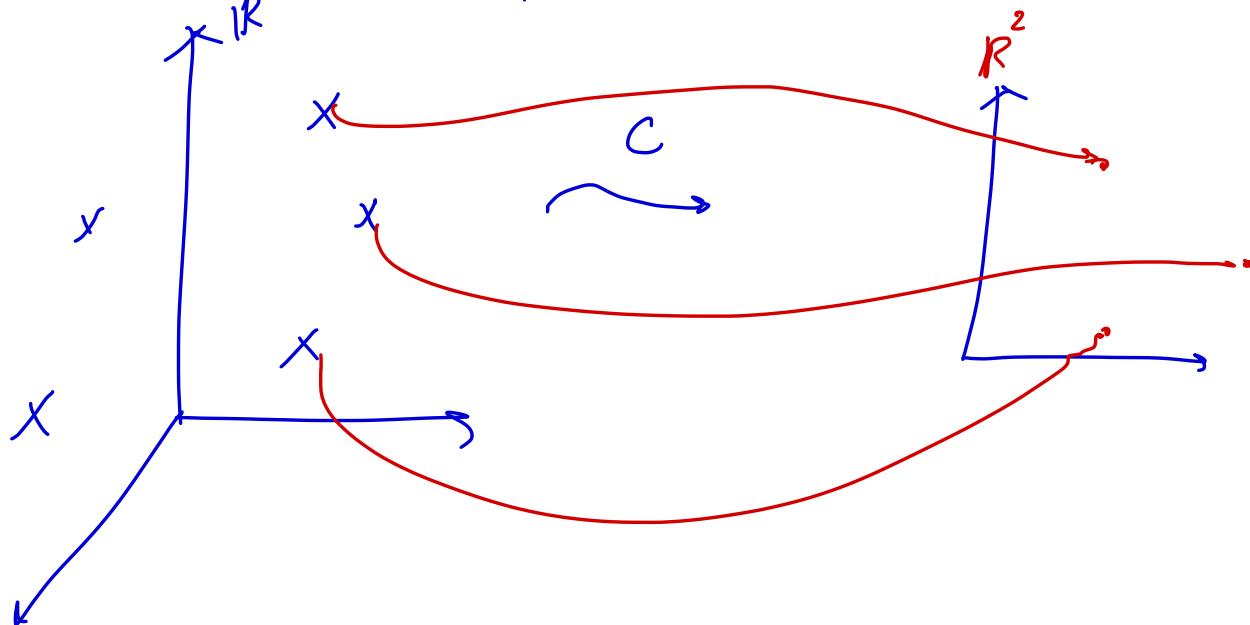
PCA

$\{x_1, x_2, \dots, x_n\}$

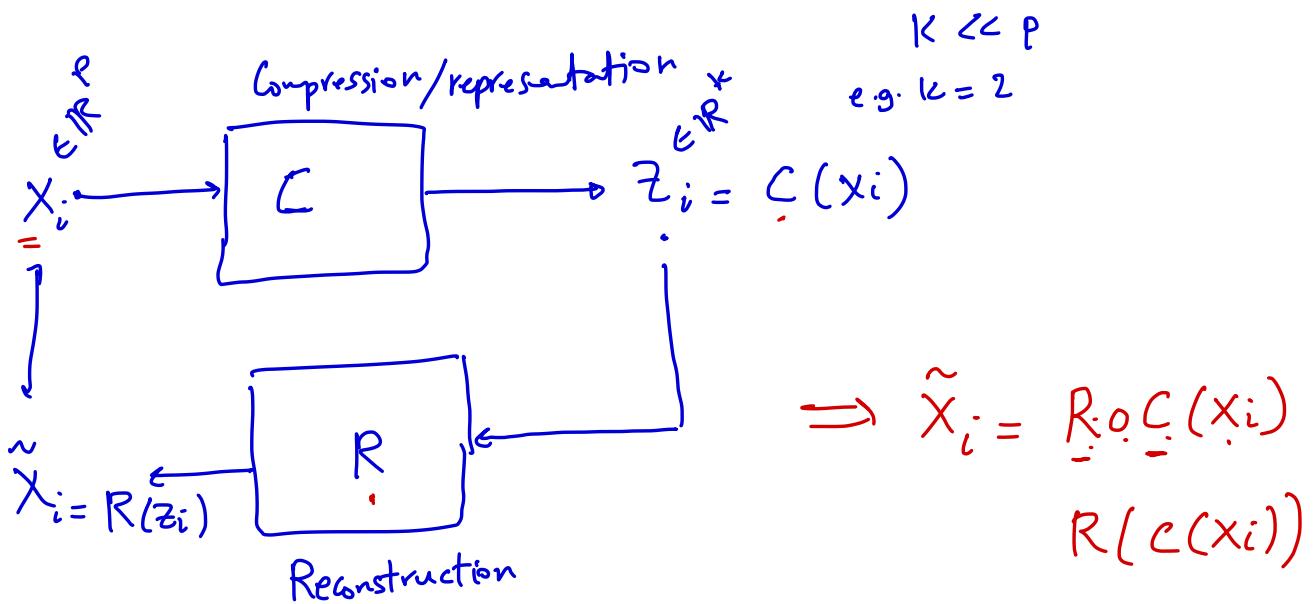
image



$x \in \mathbb{R}^P \rightarrow P: \# \text{ pixels}$



- $x_1, x_2, \dots, x_n \in \mathbb{R}^P$  ( $P$  is typically large)
- we'd like to represent the data in a low-dimensional space (e.g.  $\mathbb{R}^k$ )



our goal is to design  $C, R$  such that  $x_i$  and  $\hat{x}_i$  are as close as possible.

objective:

$$\begin{aligned} & \text{minimize}_{C, R} \sum_{i=1}^n \|x_i - \hat{x}_i\|_2^2 \\ &= \text{minimize}_{C, R} \sum_{i=1}^n \|x_i - R \circ C(x_i)\|_2^2 \end{aligned}$$

$$x_i = \begin{pmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{pmatrix} \in \mathbb{R}^p$$

$$x_i \xrightarrow{C} \begin{pmatrix} 0 \\ 0 \end{pmatrix} = z_i$$

As usual, let's start with the simplest class of mappings which are linear mappings.

$$z_i = C(C(x_i)) = C x_i$$

$C$  matrix  $C = [ ]_{k \times p}$

$$\text{e.g. } k=2 \rightarrow C = \left[ \begin{array}{c} \quad \\ \quad \end{array} \right]_{2 \times p}$$

$$= \left[ \frac{U_1^T}{U_2^T} \right]$$

$$\Rightarrow z_i = c x_i = \begin{bmatrix} U_1^T x_i \\ U_2^T x_i \end{bmatrix}$$

Also,  $R$  = linear mapping

$$\tilde{x}_i = R z_i = \underbrace{R}_{\in \mathbb{R}_{p \times k}} \underbrace{c}_{\in \mathbb{R}^k} \underbrace{x_i}_{\in \mathbb{R}^p}$$

objective:  $\min_{R, C} \sum_{i=1}^n \|x_i - R \cdot c\|_2^2$  (PCA)

$$C: \mathbb{R}^p \rightarrow \mathbb{R}^k$$

$C$  = matrix  $\in \mathbb{R}^{p \times k}$

This problem is called  
the Principal Component Analysis (PCA)

In order to solve this problem,  
we need to review an important  
tool in linear algebra:

(this tool is important for many  
other branches of Data Science)

# The Singular Value Decomposition:

Theorem: Every symmetric matrix,  $A_{p \times p}$ , can be written as

$$A_{p \times p} = U_{p \times p}^T \Lambda_{p \times p} U_{p \times p}$$

- Where  $U U^T = I_{p \times p}$   $\xrightarrow{\text{identity matrix}}$

$$\Lambda = \text{diagonal: } \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{pmatrix}$$

$$- A = U^T \cdot \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{pmatrix} U$$

$$- \lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_n$$

-  $\lambda_i$  = eigenvalues of  $A$ .

- Each row of  $U$  is an "eigenvector" of the matrix  $A$ .

Example:

Assume  $A$  is  $2 \times 2$  matrix.

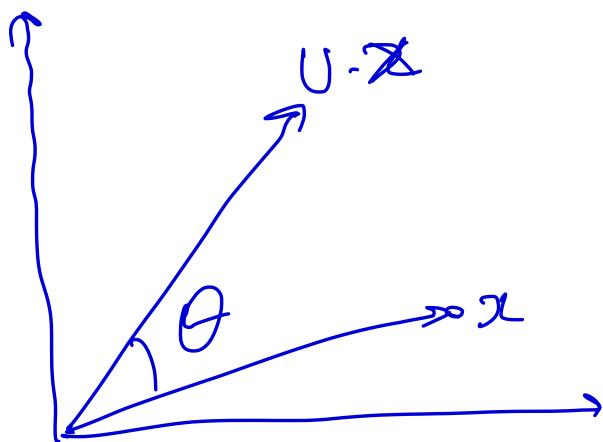
$$A = \begin{pmatrix} a & c \\ c & b \end{pmatrix}$$

$\brace{rescaling}$

$$A = \underbrace{U}_{\text{Rotation}(\theta)} \cdot \underbrace{\begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}}_{\text{Scaling}} \cdot \underbrace{U^T}_{\text{Rotation}(-\theta)}$$

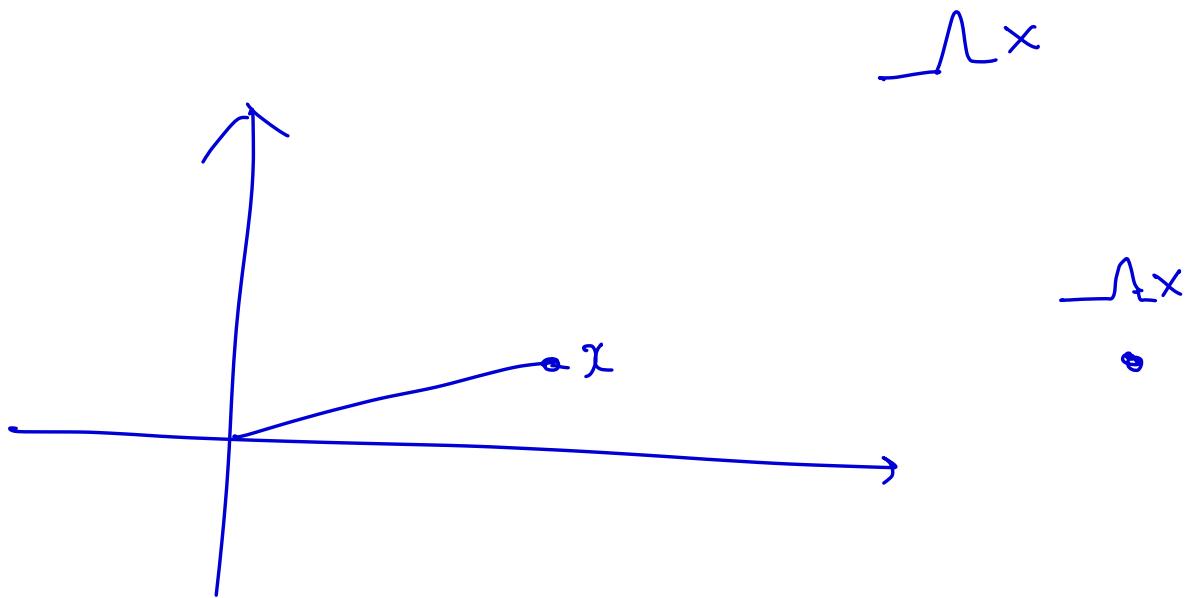
$$U \cdot U^T = I \rightarrow U = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}$$

$\brace{}$



$$A = \begin{pmatrix} 4 = \lambda_1 > 0 & 0 \\ 0 & 1 = \lambda_2 > 0 \end{pmatrix}$$

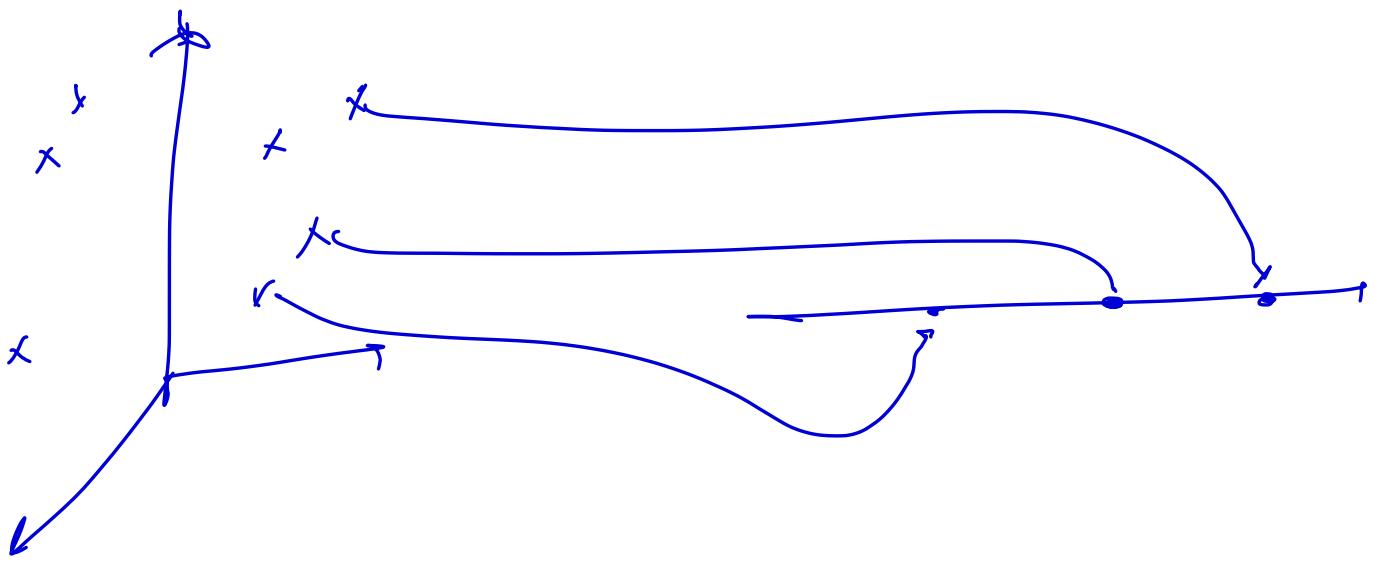
? Re Scaling



$A = \text{Rotation}(-\theta), \text{Rescaling} \cdot \text{Rotation}(\theta)$

let's go back to our problem.

let's assume for simplicity that  $k=1$  ; we're looking for the best 1-dimensional representation of the data.

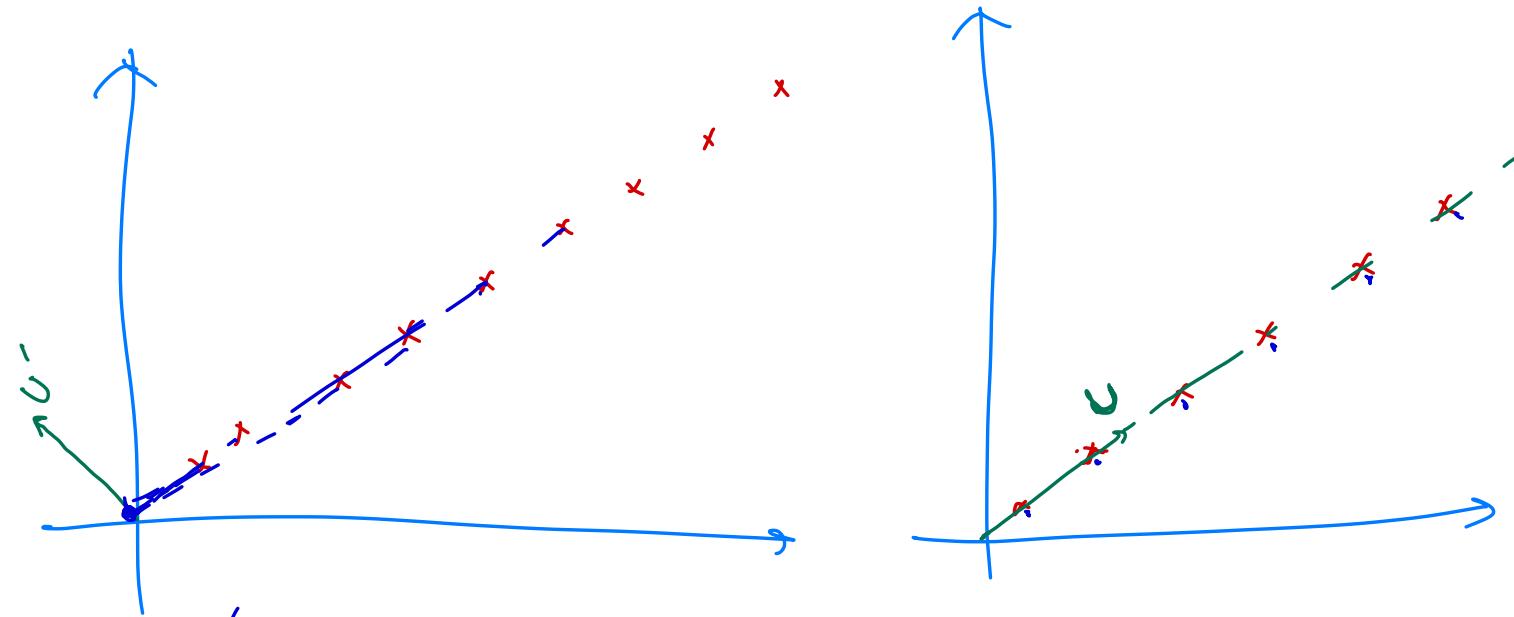


$$C = [ \quad ]_{k \times p} \longrightarrow C = U^T$$

$$C \cdot x_i = \underset{\substack{= \\ \text{vector}}}{U^T} x_i$$

When  $k=1$ , our goal is to find a single vector  $U^T$  such that the representation  $z_i = U^T x_i$  is "the most informative" representation.

Example:

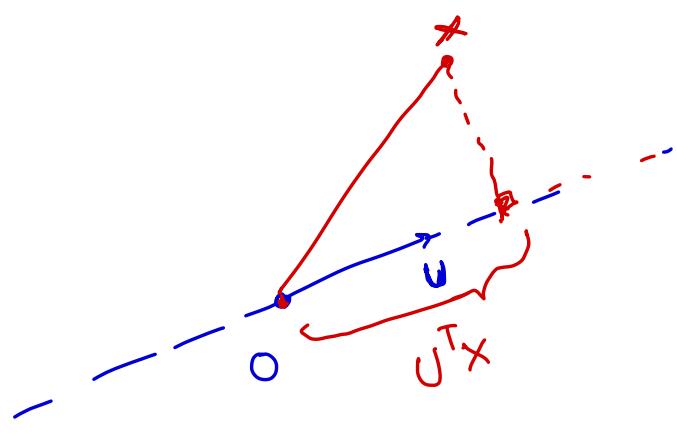


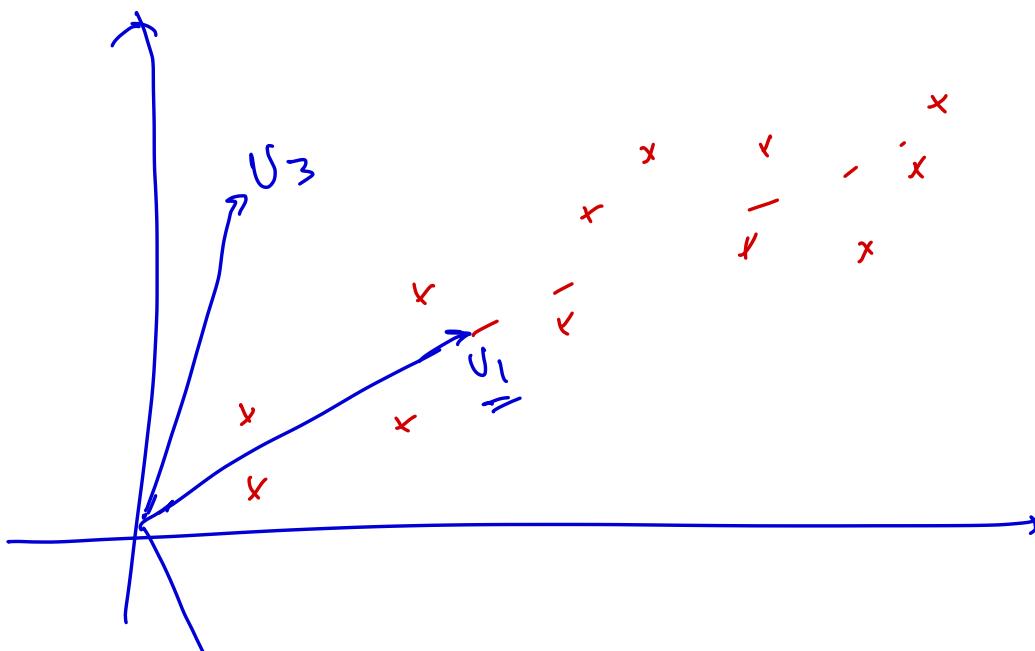
$$U'^T X_i \approx 0$$

$$U^T X_i$$

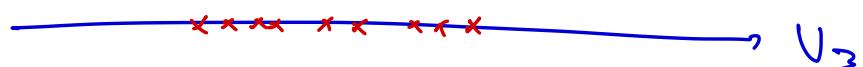
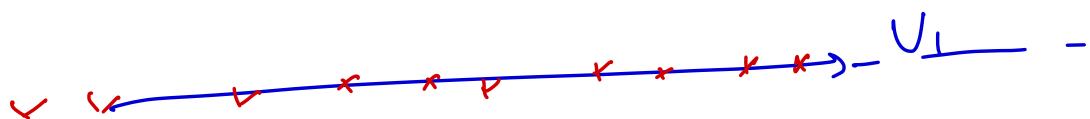
$U$  is a better direction than  $U'$  in terms of capture more information about the data

Assume  $\|U\|_2 = 1$





$$U_1 > U_3 > U_2$$



We're looking for a direction  $U$  along which data has the most "variation".

- Mathematically, we're looking for a direction  $u$  such that the "variance" of  $u^T x$  is the most.
- Assume that the data comes from a distribution  $X \sim P(x)$ . Also assume for simplicity that  $E[x] = 0$ .

$$\begin{aligned}
 & \max_{\substack{u \in \mathbb{R}^{1 \times p} \\ \|u\|_2=1}} \text{Var}(u^T x) \\
 &= \max_{\substack{u : \|u\|_2=1}} E[(u^T x)^2] \\
 &= \max_{\substack{u : \|u\|_2=1}} E[u^T x x^T u] \\
 &= \max_{\substack{u : \|u\|_2=1}} u^T \underbrace{E[x x^T]}_{\text{Covariance matrix of the data}} u
 \end{aligned}$$

$$\left( \begin{array}{l} x \in \mathbb{R}^p \\ E[x x^T] \\ = \sum_{p \times p} \end{array} \right)$$

$$\text{Let } \Sigma \triangleq E[X X^T]$$

$$= \max_u u^T \Sigma u$$

$u:$

$$\|u\|_2 = 1$$

The solution of the problem above can be found from the singular value decomposition of  $\Sigma$ .

$$\Sigma = \begin{pmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_p \end{pmatrix}$$

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p$$

$$\text{SVD} \rightarrow \Sigma = U \Sigma^T U^T$$

It can be shown the the first row of the matrix  $U$ , is the maximizer of the following problem:

$$\max_u u^T \Sigma u$$

$$u: \|u\|_2 =$$

$$U = \begin{pmatrix} U_1^T \\ U_2^T \\ \vdots \\ U_p^T \end{pmatrix}_{p \times p}$$

- $U_1^T$  is the solution to the problem above.
- $U_1^T$  is the principal component of the data.

$-v_1^T$  is the vector along which the data has the most "variation".

$$\Sigma = E[X X^T]$$

- We know that the data

$$x_1, x_2, \dots, x_n \stackrel{\text{iid distribution}}{\sim} p(x)$$

in practice:

- We cannot compute the matrix  $\Sigma$  because we do not know the distribution of the data

$$\Sigma \approx \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

  
unbiased estimator

## Algorithm (PCA) : ( $k=1$ )

- Input:  $X_1, X_2, \dots, X_n \in \mathbb{R}^P$
- Compute the empirical covariance matrix of data:  $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n X_i X_i^T$
- Compute the Singular value decomposition of  $\hat{\Sigma}$ :  
$$\hat{\Sigma} = \hat{U}^T \Lambda \hat{U}$$
- Use the first row of  $U$ ,  $U_1$ , as the principal component of the data:

$$Z_i = U_1^T X_i .$$

---

\* let's now find the best 2-dimensional representation of the data.

$$\rightarrow \max_{U_1, U_2} U_1^T \Sigma U_1 + U_2^T \Sigma U_2$$
$$\|U_1\|, \|U_2\| = 1$$

(orthogonal)  $U_1^T U_2 = 0$

If can be shown that  $u_1, u_2$  are the first two rows of the matrix  $U$ .

- \* In general if we want to find a matrix  $C = \begin{bmatrix} \frac{u_1^T}{\|u_1\|} \\ \frac{u_2^T}{\|u_2\|} \\ \vdots \\ \frac{u_k^T}{\|u_k\|} \end{bmatrix}_{k \times n}$  then
  - we need to solve the following optimization problem:

*Cumulative variance ←  
of the data along  
the directions*

$u_1, \dots, u_k$

$$\sum_{i=1}^k u_i^T \Sigma u_i$$

$$\text{s.t. } \|u_i\|=1$$

$$\text{and } u_i^T u_j = 0 \text{ for } i \neq j$$

& it can be shown that the solution to the above optimization problem is the first  $k$  rows of the matrix  $U$ .

Algorithm PCA for general  $k$ :

- Input :  $x_1, \dots, x_n \in \mathbb{R}^P$

- Compute:  $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$

- Compute: (SVD)  $\hat{\Sigma} = U \Sigma U^T$

-  $C = \begin{bmatrix} U_1 \\ U_2 \\ \vdots \\ U_k \end{bmatrix}$  where  $U_1, \dots, U_k$  are the first  $k$  rows of the matrix  $U$ .

If can be proven that PCA gives the optimal linear representation of the data (it solves equation (PCA) exactly).

## Lecture 21 :

### Module 5: An Introduction to statistical Learning Theory:

We consider the problem of supervised learning in this module, but all the concept 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.  $Y \in \{0, 1\}^F$ ), 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:

a function from  $X$  to  $Y$

$$\text{error / loss } \leftarrow L_D(h) = \Pr_{(x,y) \sim D} \{ h(x) \neq y \}$$

$D$  = 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:

$$\underset{h}{\text{minimize}} \quad L_D(h) \quad (1)$$
$$= E_{(x,y) \sim D} [\mathbb{1}\{h(x) \neq y\}]$$

$\mathbb{P}\{h(x) \neq y\}$   
 $(x,y) \sim D$

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:

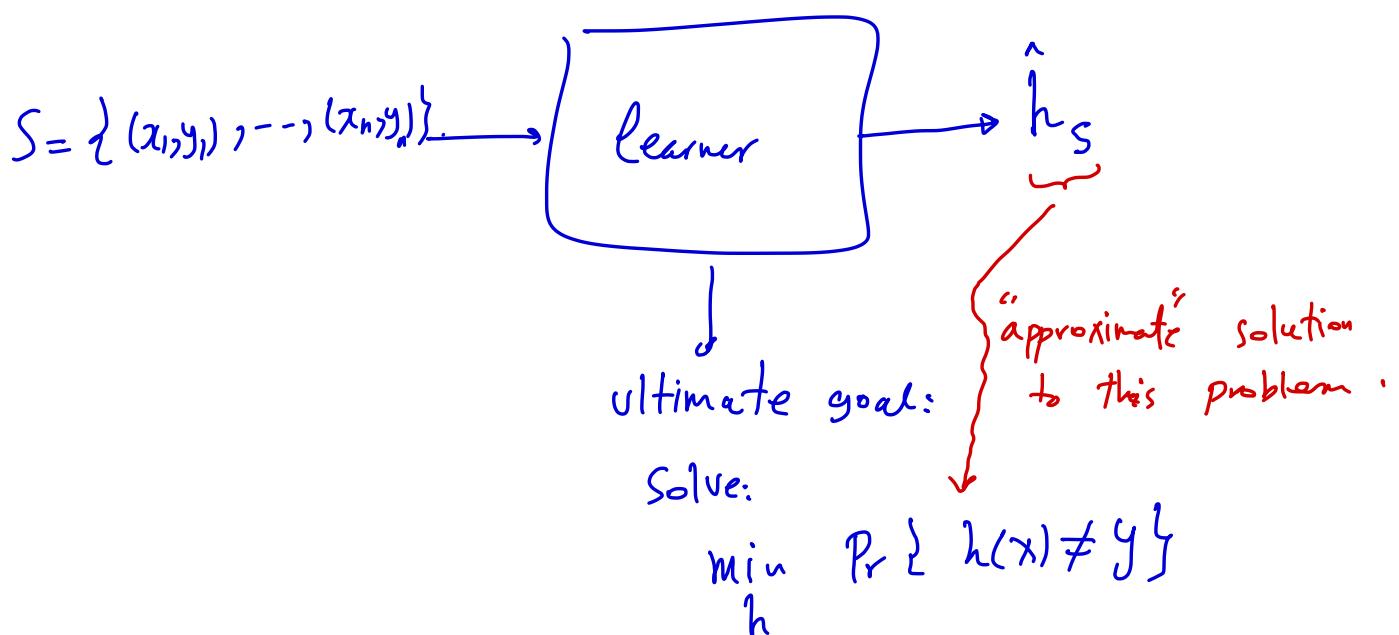
$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

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  $D$



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

$$\min_h E \left[ \# \{ h(x) \neq y \} \right]_{(x,y) \sim D}$$

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

#### (4) Empirical Risk Minimization (ERM):

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

↓ unbiased estimate

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

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

$$\min_h L_S(h)$$

instead of solving  $\rightarrow$

$$\boxed{\min_h L_D(h)}$$

ERM solves  $\rightarrow$

$$\boxed{\min_h L_S(h)}$$

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

→ overfitting

Example: Consider the following distribution on data:

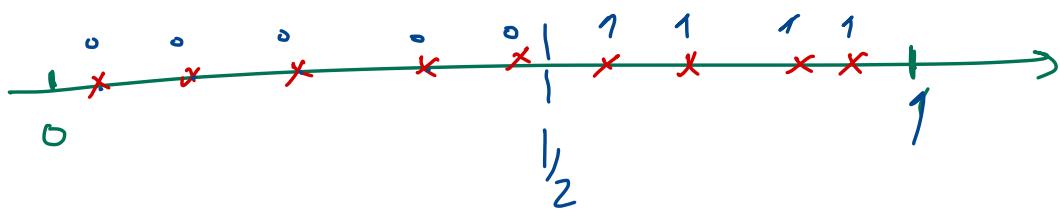
$$X = [0, 1]$$

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



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

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



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

$$\min_h \underbrace{L_S(h)}_{=} =$$

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

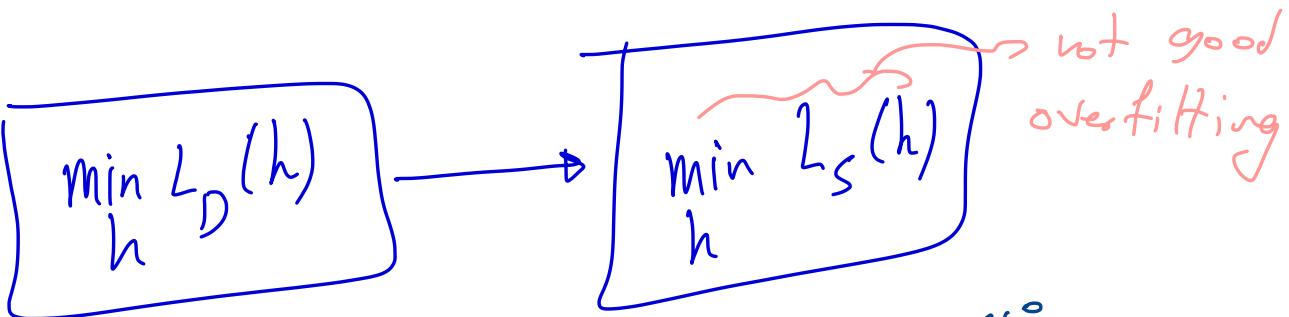
$$\hat{h}_S = \begin{cases} y_i & \cdot \text{ if } x = x_i \\ 0 & x \neq \{x_1, \dots, x_n\} \end{cases}$$

Proof of the claim:

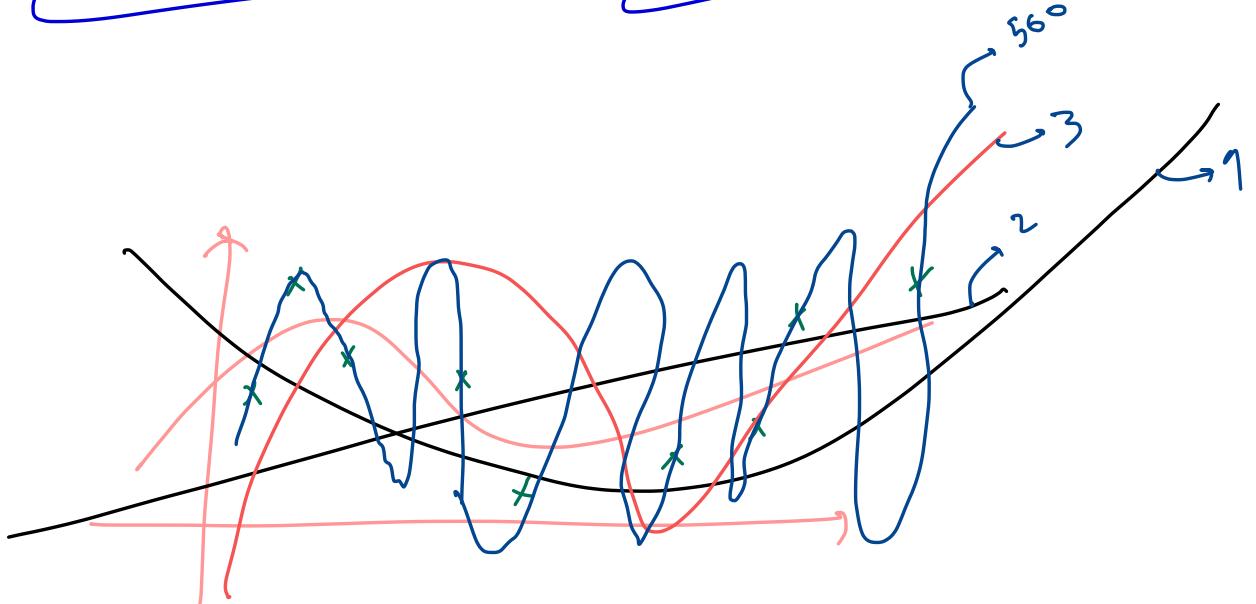
Define  $\hat{h}_S(x) = \begin{cases} y_i & \text{if } x = x_i \\ 0 & \text{if } x \notin \{x_1, \dots, x_n\} \end{cases}$

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

$$= 0$$



Recall  
(Regression)



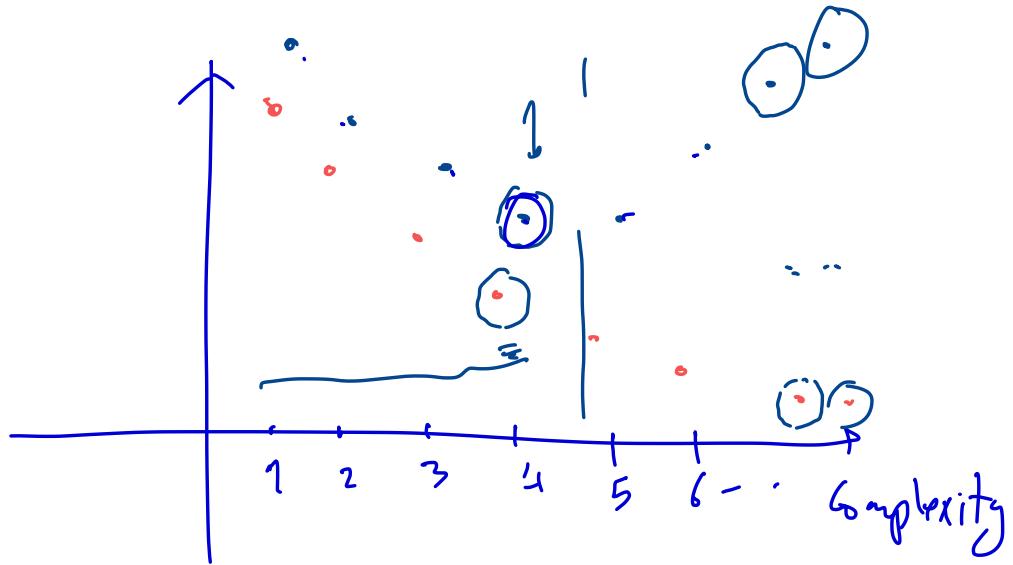
(5) Although ERM is natural, it can miserably if we are not careful:

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.



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

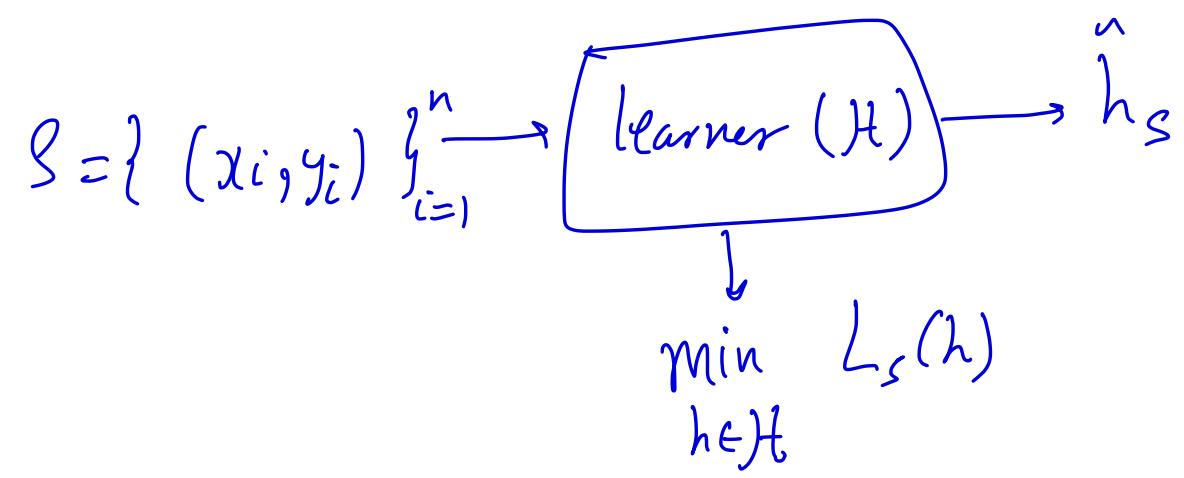
restrict the class of functions/predictors  
(restrict the complexity)

↓  
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.



# 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

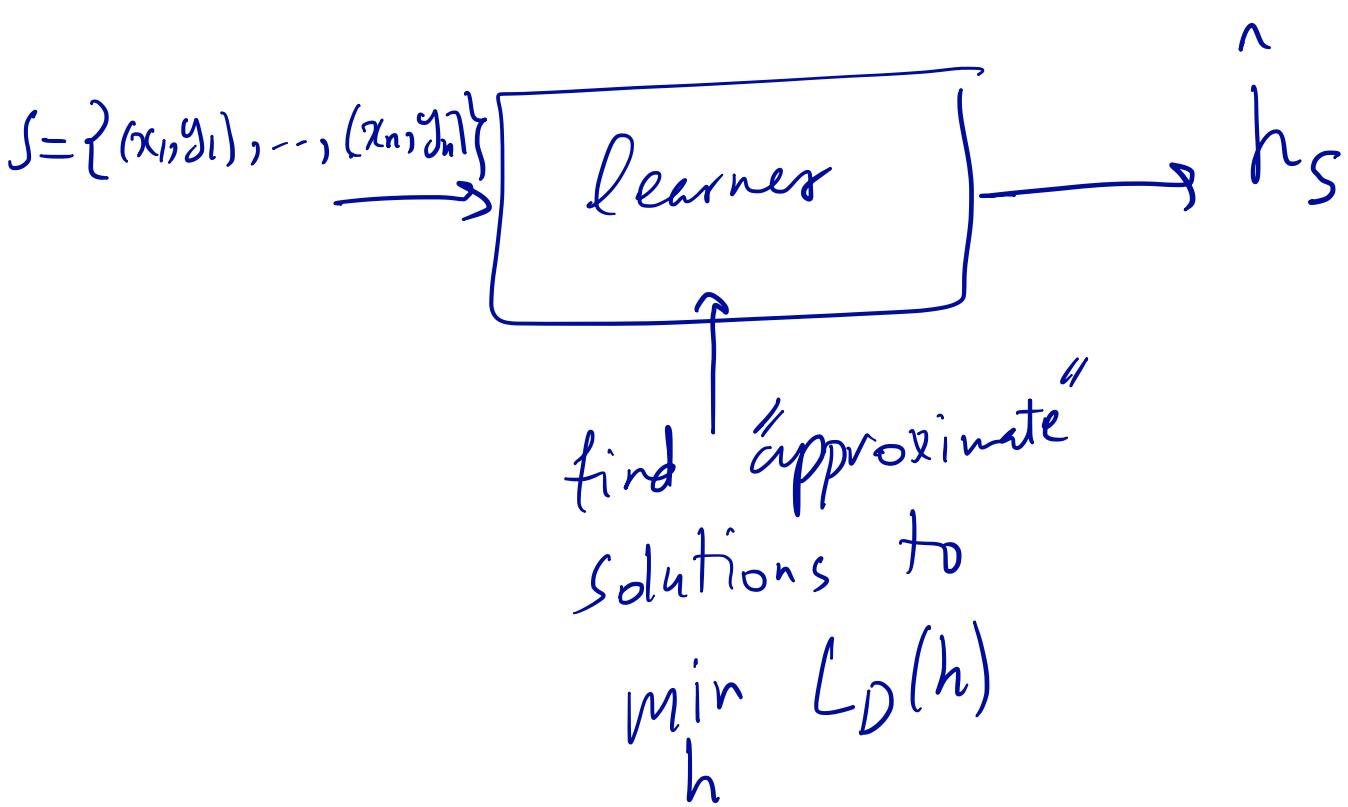
$$\underset{h}{\text{minimize}} \quad L_D(h) \quad (1)$$

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

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:

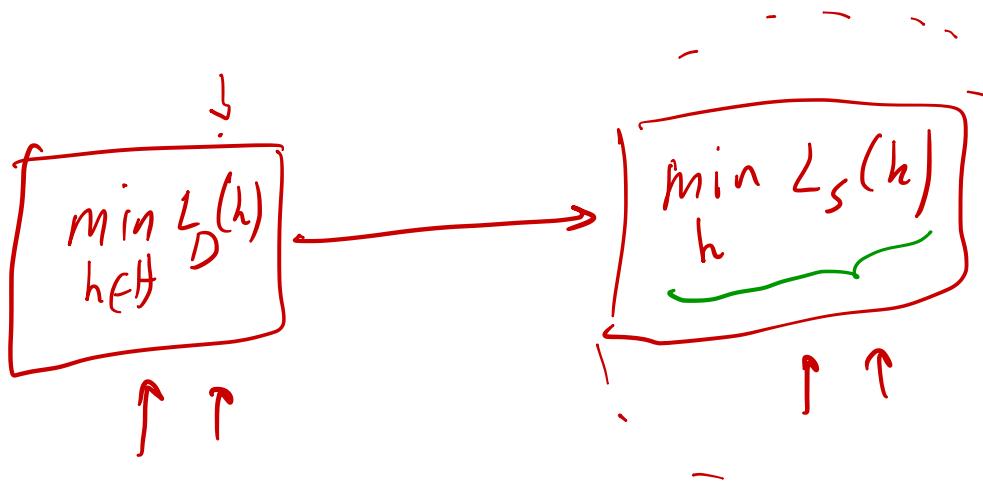
$$\min_h L_S(h) \quad (2)$$

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

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.

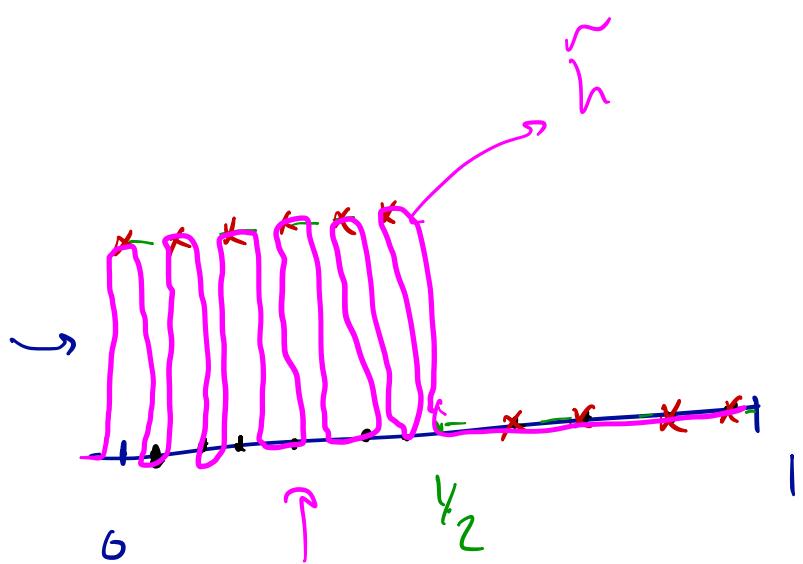
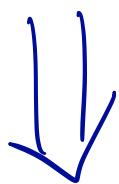
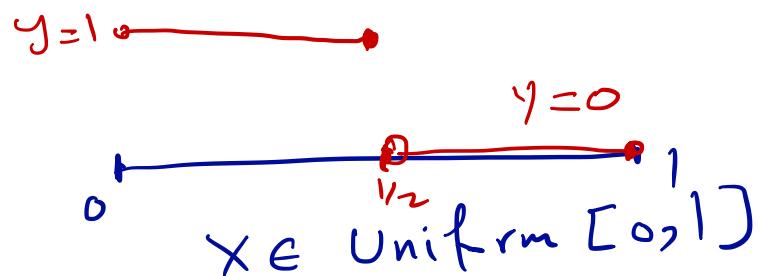
The minimizer of

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

is always zero.

Example:

Distribution of data:



$$\min_h L_S(h)$$

{  
on unseen  
data here  
is very bad}

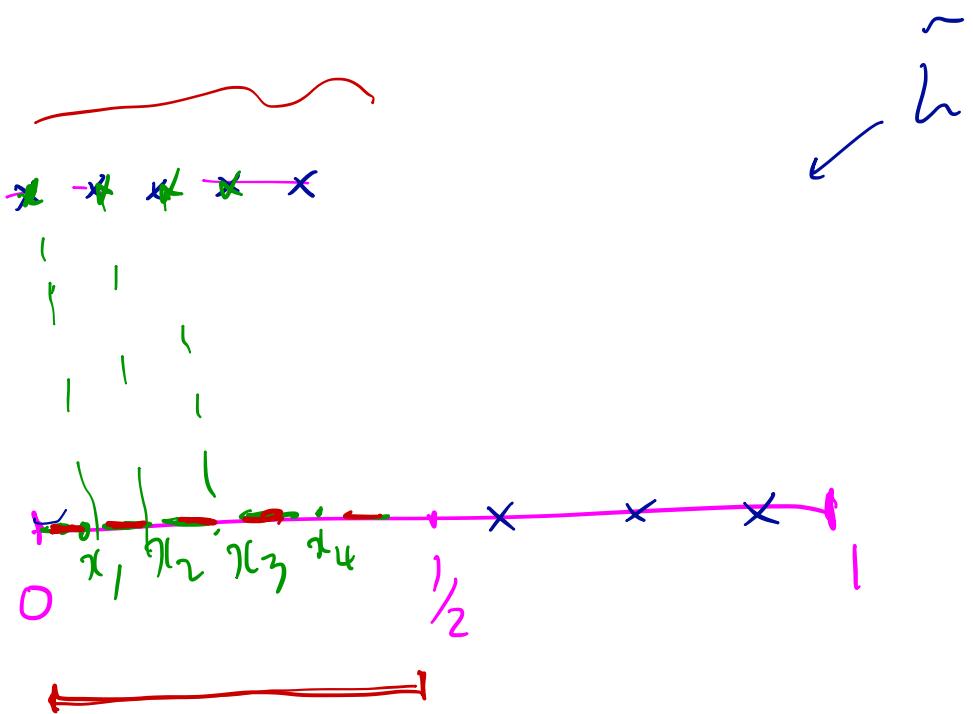
$$L_S(\tilde{h}) = 0$$

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

---

$$\min_h L_D(h)$$

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



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

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

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

$$\equiv \not\models$$

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)

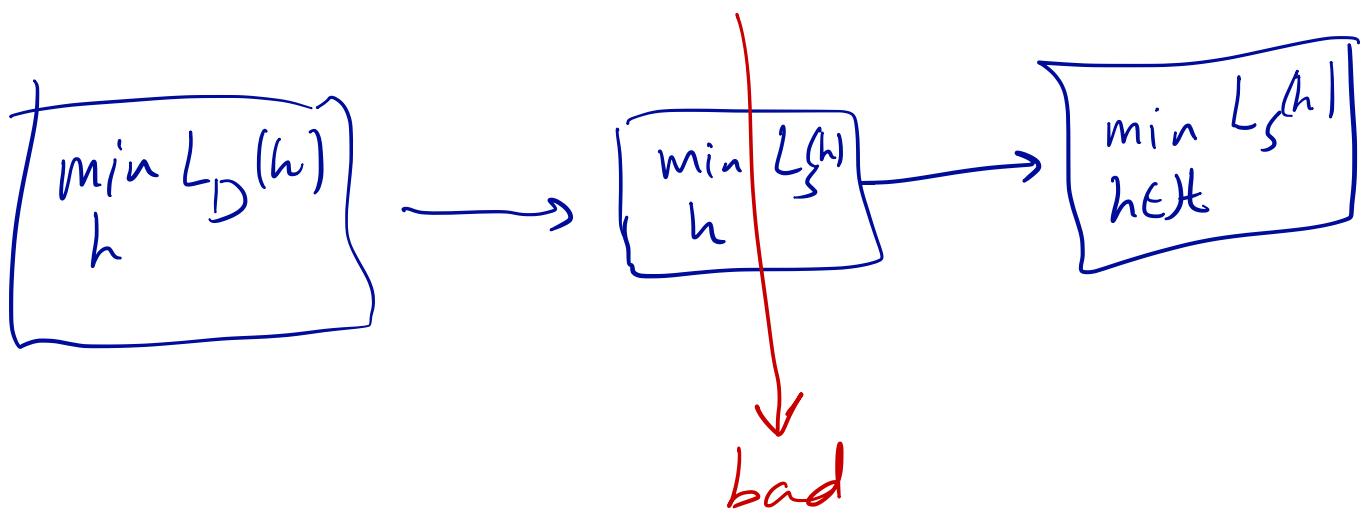


↓ prevent overfitting

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



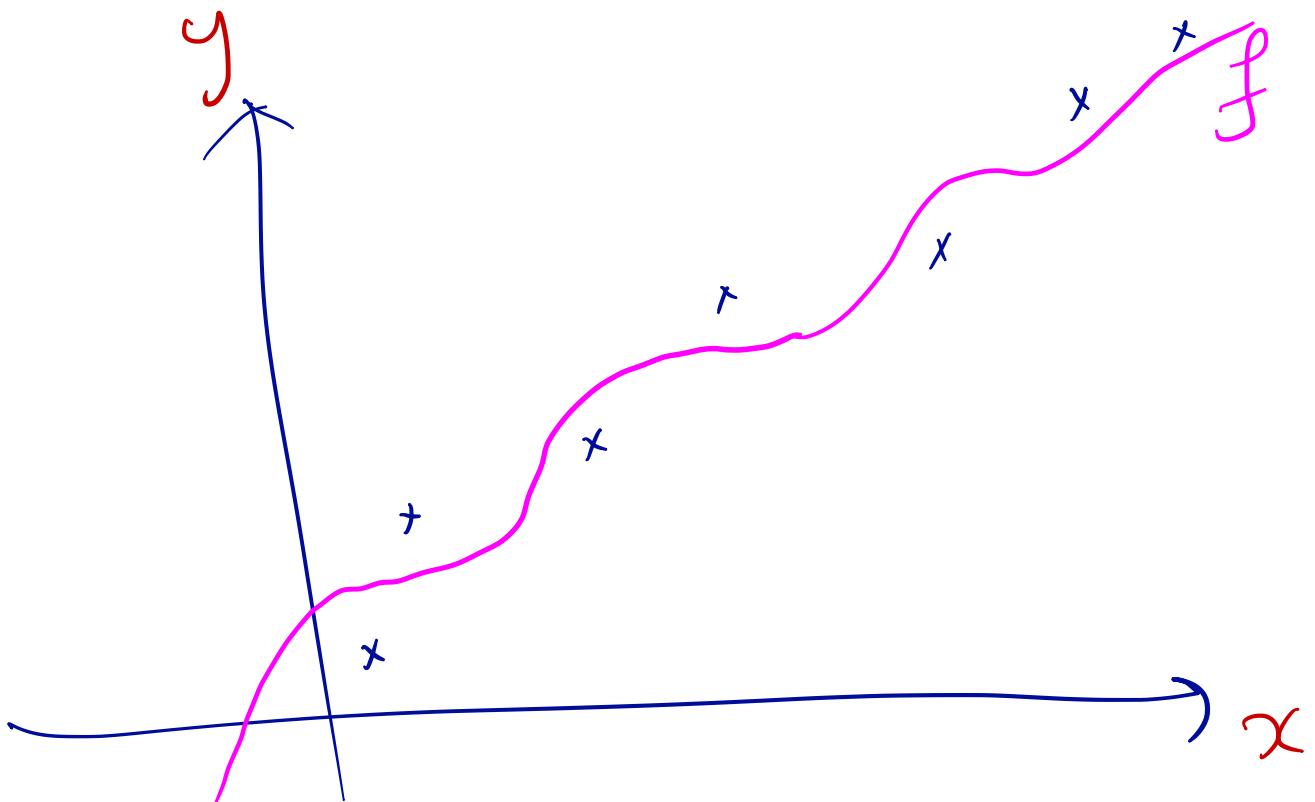
So we are biasing the learner towards a particular set of predictors.

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.

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



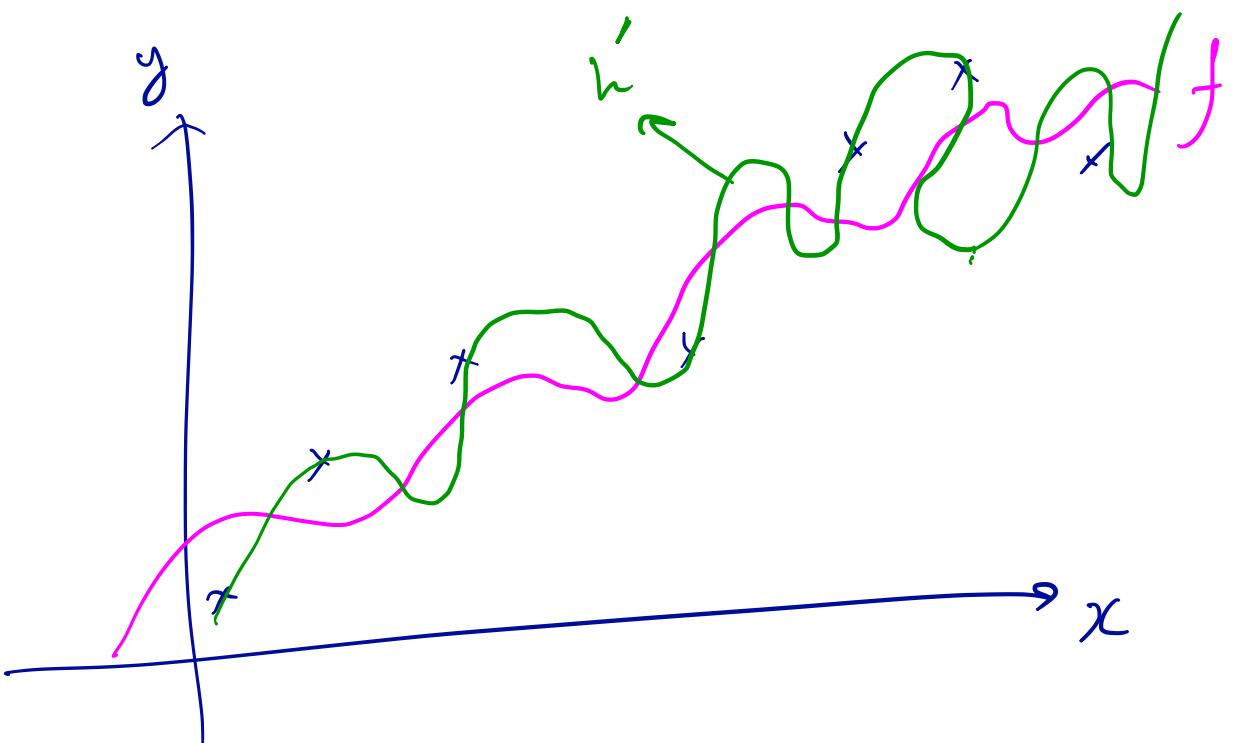
for any  $h: X \rightarrow Y$

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

$$\min_h L_D(h) = 6^2$$

↳ minimizer  $\rightarrow f$

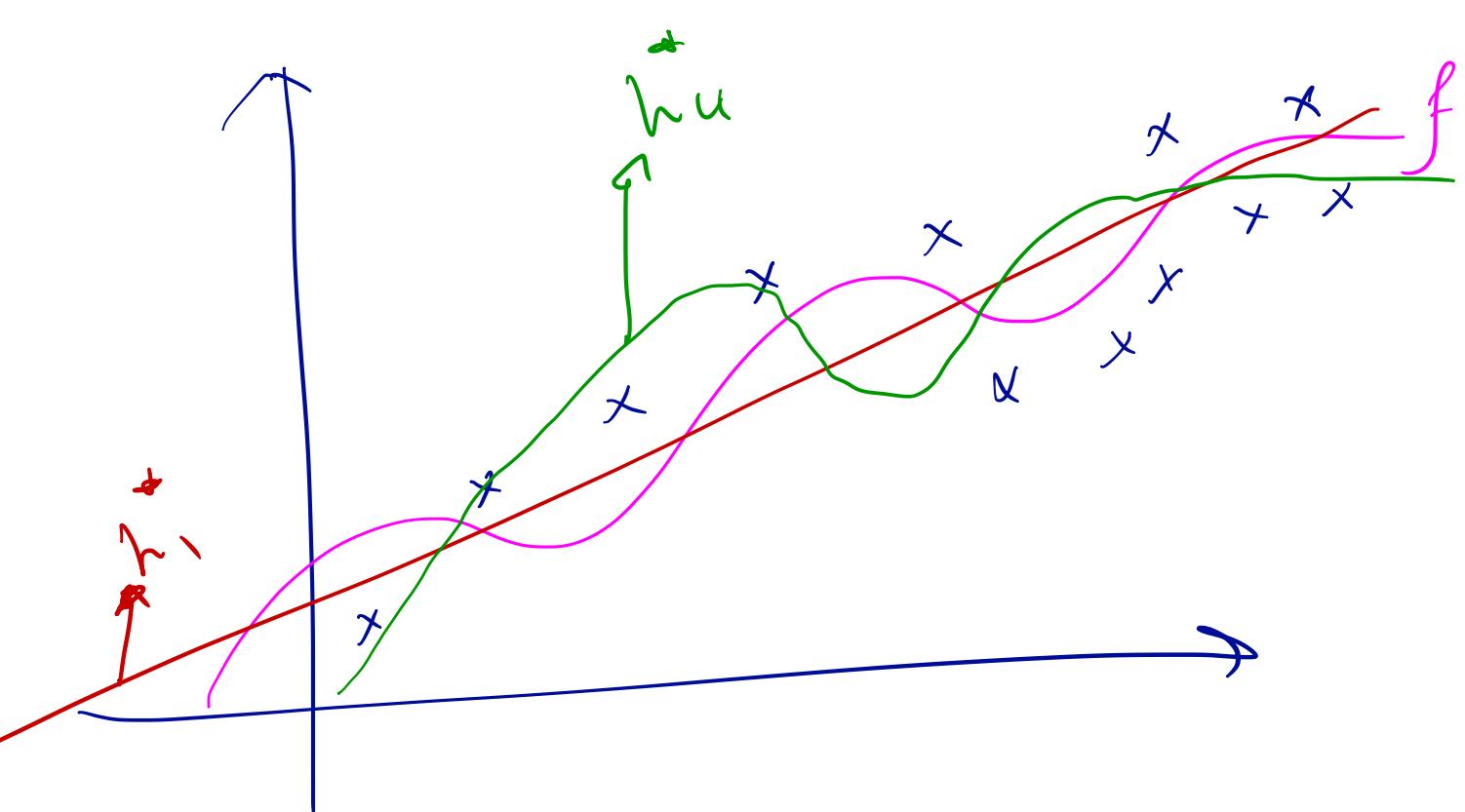
$$L_D(f) = E[(y - f(x))^2] = 6^2$$



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

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

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

↳ this is computable

then it is possible that

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

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

↑  
underfitting

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

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

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

↑  
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  $\mathcal{H}$

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

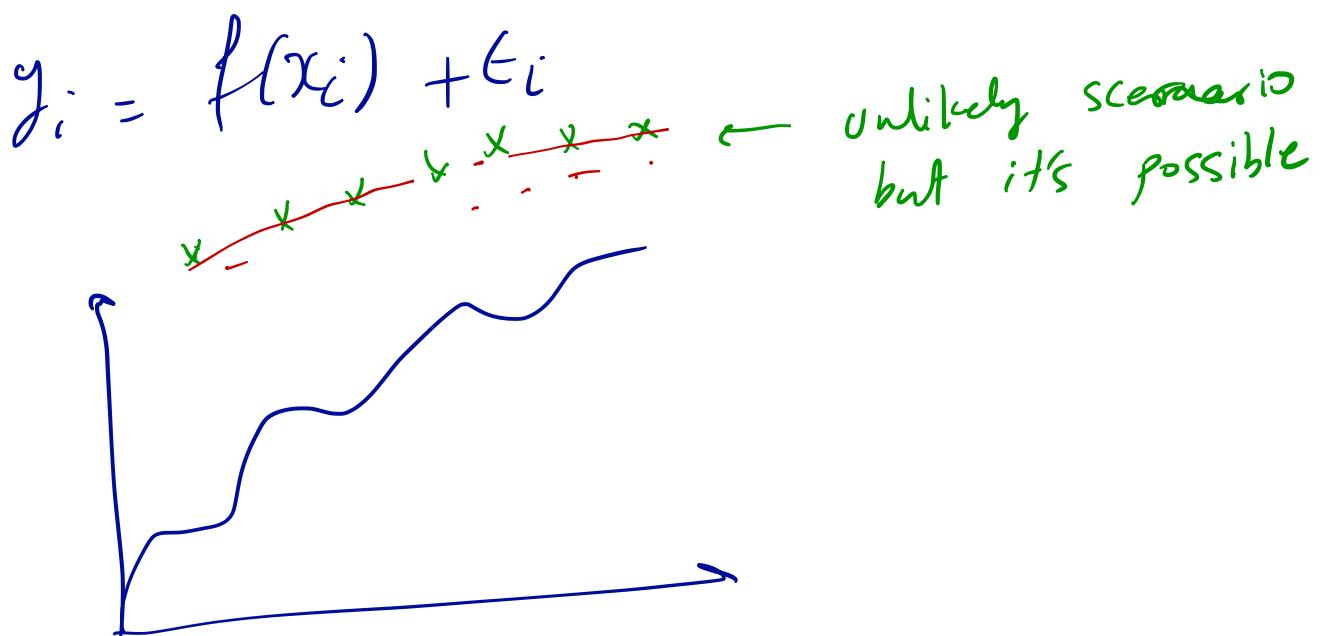
2. what we're actually optimizing is

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

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

↳ 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).



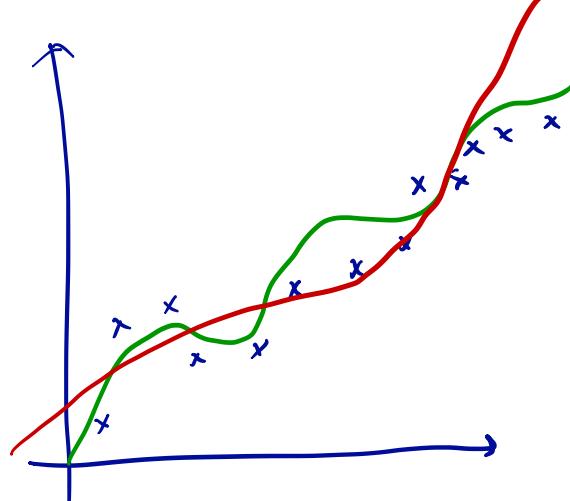
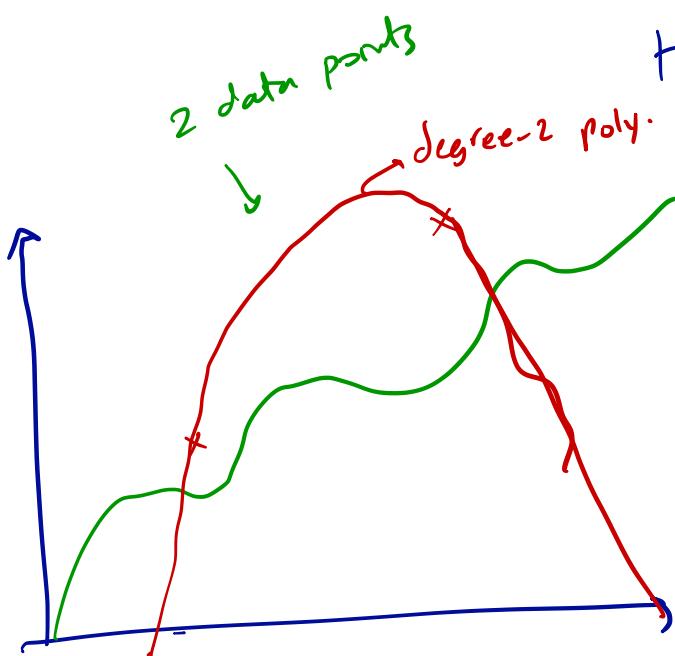
$$L_D \rightarrow E_D$$

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

$$L_S \rightarrow L_D$$

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

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  $1 - \delta$

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

$\epsilon$  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.

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

$$\begin{aligned}
 \min_{h \in \mathcal{H}} L_D(h) &\longrightarrow \min_{h \in \mathcal{H}} L_S(h) \\
 h^* & \downarrow \\
 h^* \text{ is one of the functions} \\
 h \in \mathcal{H} \text{ with smallest} \\
 \text{empirical error.} \\
 \min_{h \in \mathcal{H}} L_S(h) &= \min \{L_S(h_1), \dots, L_S(h_m)\}
 \end{aligned}$$

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:

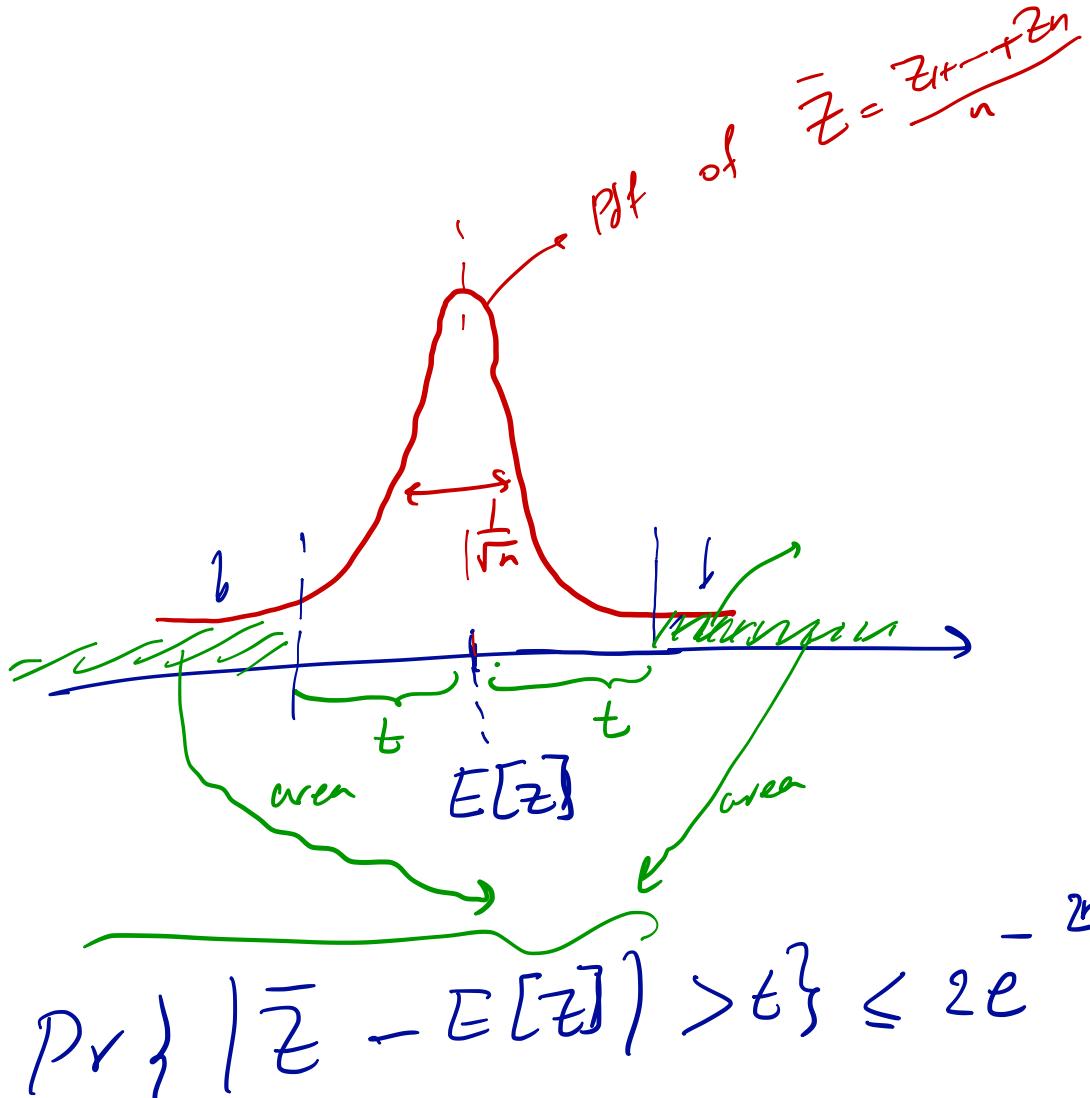
Let  $Z_1, Z_2, \dots, Z_n$  be  $n$  iid random variables that take value in the unit interval (i.e.  $Z_i \in [0, 1]$ ).

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

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



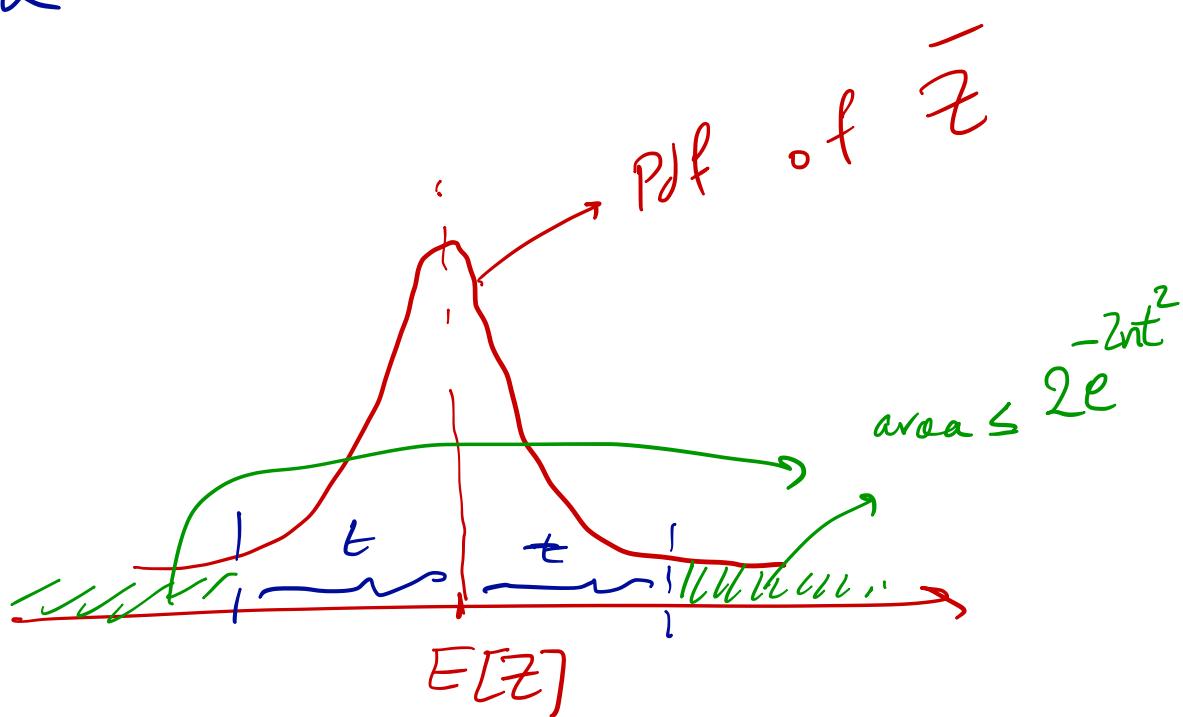
Recall that Hoeffding's inequality was given as follows:

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.

$$\Pr \left\{ \left| L_S(h) - L_D(h) \right| > \epsilon \right\} \leq \delta \quad (2)$$

using (1), and in order to guarantee

(2), we can let:

$$2e^{-2n\epsilon^2} \leq \delta \quad \left( \log \frac{1}{\delta} = -\log_2 \frac{1}{\delta} \right)$$
$$\Rightarrow n \geq \frac{1}{2\epsilon^2} \log \frac{1}{\delta}$$

So the final statement is:

Given any  $\epsilon, \delta$ , 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

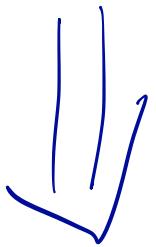
$$\Pr \{ A_1 \} + \Pr \{ A_2 \} + \dots + \Pr \{ A_m \} \leq \delta \quad (6)$$

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

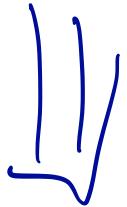
Now, to guarantee (6), it is sufficient to choose  $n$  large enough s.t. for every  $i$

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- $\epsilon$ -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

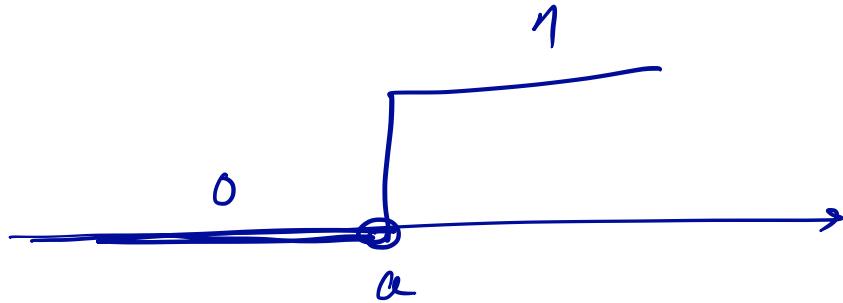
$$0 \leq L_D(h_S^*) - \min_{h \in \mathcal{H}} L_D(h) \leq 2\epsilon.$$

minimizer of

ERM over  $\mathcal{H}$ .

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

$$h_a(x) =$$



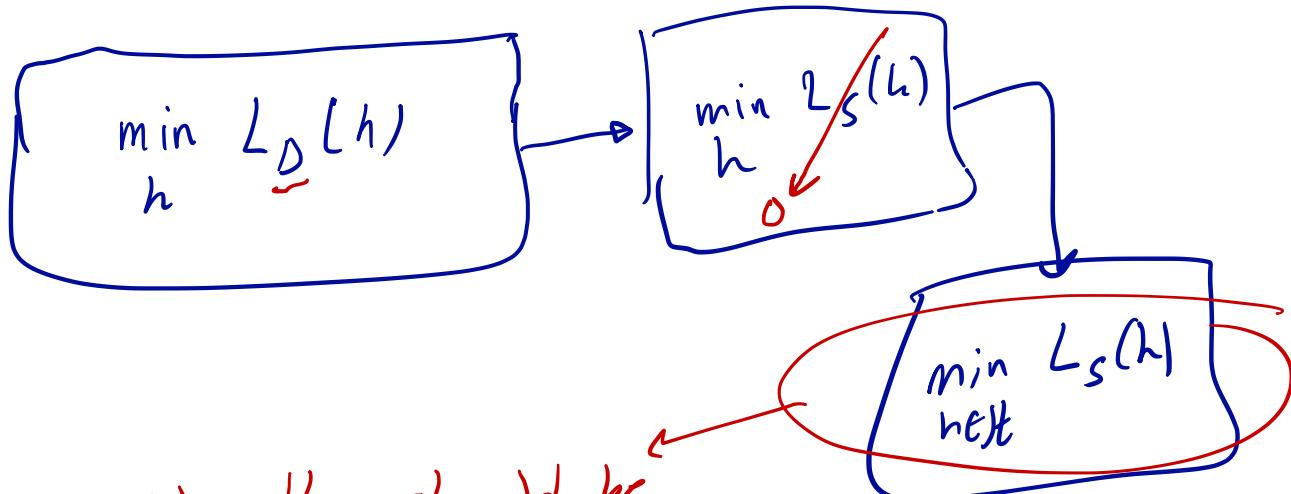
$$\mathcal{H} = \{ h_a(x), a \in [-\infty, \infty) \}$$

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)  $\exists n(\epsilon, \delta) \in \mathbb{N}$  training data

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

If  $|S| > n(\epsilon, \delta)$ :

w.p.  $1 - \delta$ :

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

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

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

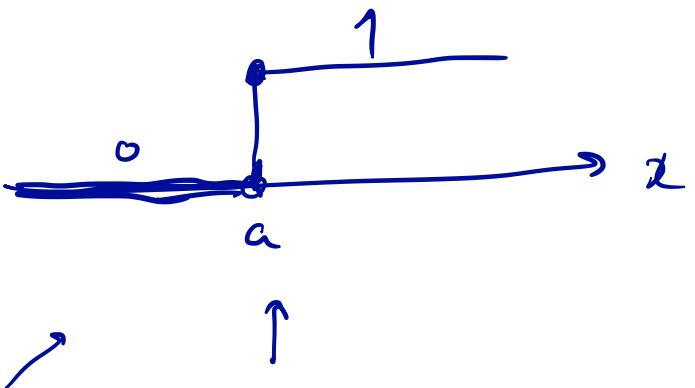
$$- \text{learning algorithm: } h_s^* = \operatorname{argmin}_{i=1, \dots, m} L_s(h_i)$$

$\mathcal{H}$  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- $\delta$ :

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

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-like problem and we lose at most a small value  $\epsilon$  with respect to the optimal accuracy.

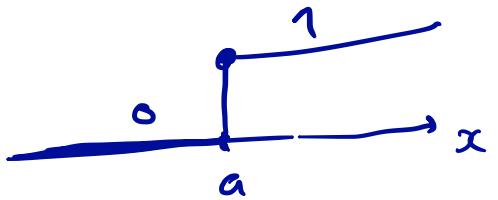
VC dimension:

1

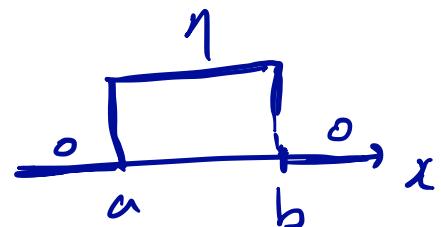
Vapnik-Chernovskii

Consider the following function classes:

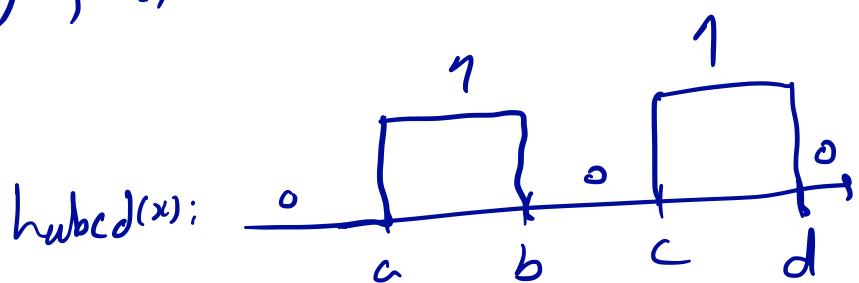
$$H_1 = \{ h_a(x), a \in \mathbb{R} \} \rightarrow$$



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

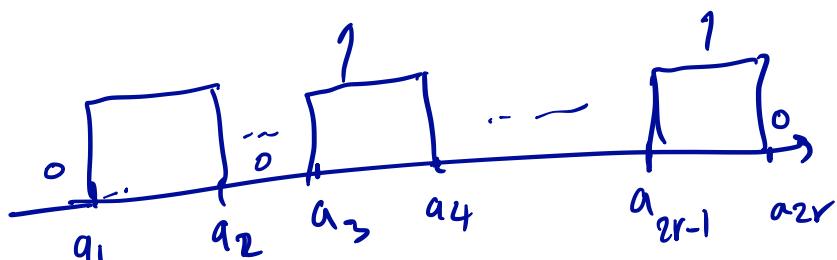


$$H_3 = \{ h_{abcd}(x), a, b, c, d \in \mathbb{R} \}$$



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

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



$H_1 < H_2 < H_3 < H_4 \rightarrow$  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  $C = \{x_1, x_2, \dots, x_k\} \subseteq X$ .

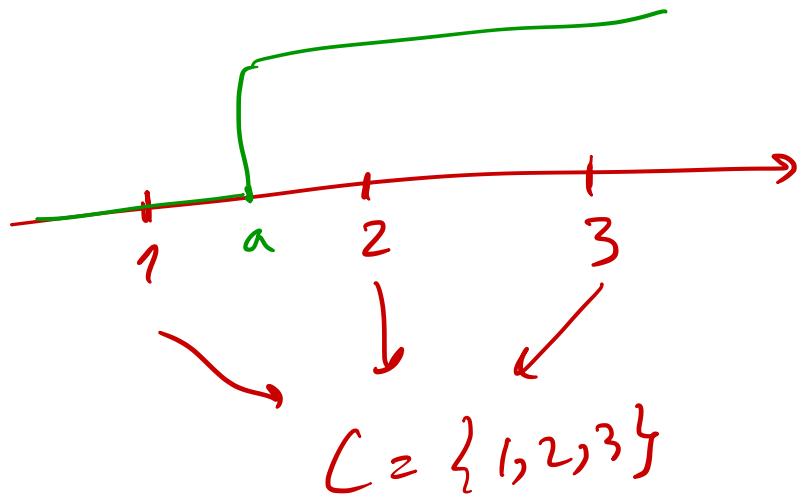
The restriction of  $\mathcal{H}$  to  $C$  is  
 the set of all the  $k$ -tuples

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

**Example:** if  $\mathcal{H} = \{h_a(x) , a \in \mathbb{R}\}$   
 e.g. if  $h_a(x) = \begin{cases} 1 & x > a \\ 0 & \text{otherwise} \end{cases}$

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

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

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

$2^k$

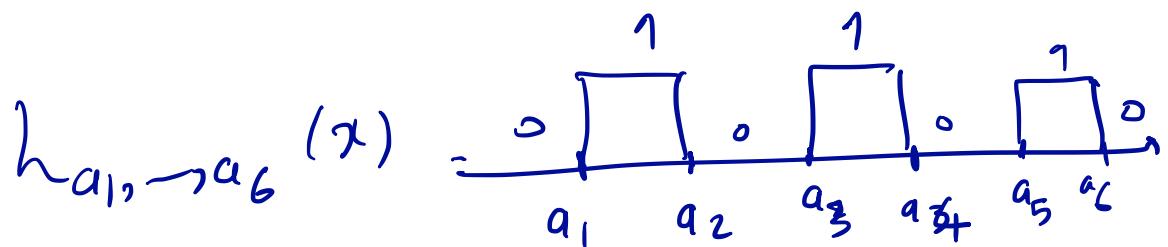
Definition (Shattering): We say that

a function class  $\mathcal{H}$  shatters a

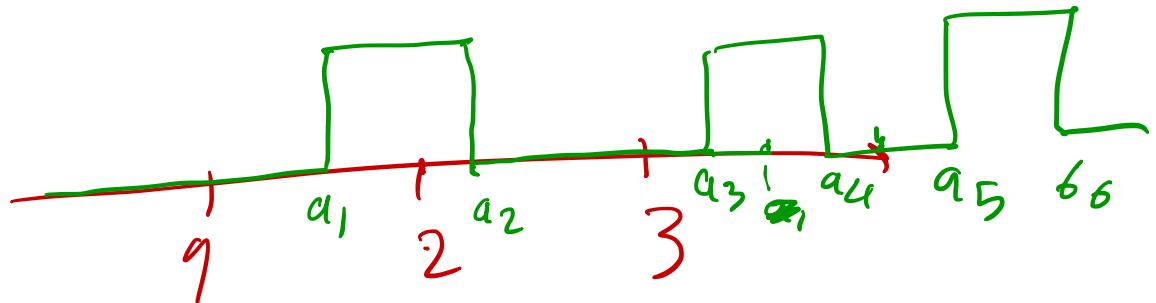
Set  $C$  if  $|\mathcal{H}_C| = 2^{|C|}$

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):

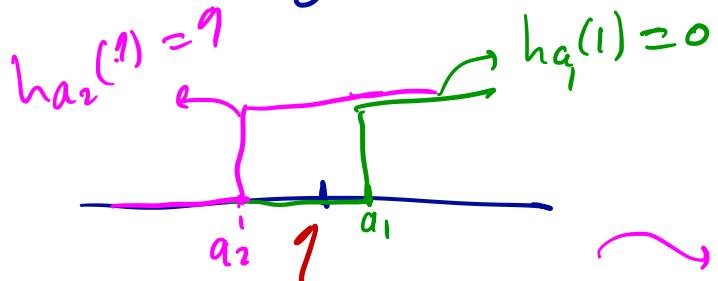
The VC-dimension of a function class  $\mathcal{H}$  is the largest number  $k$  such that exists a set  $C$  of size  $k$  which is shattered by  $\mathcal{H}$ .

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.

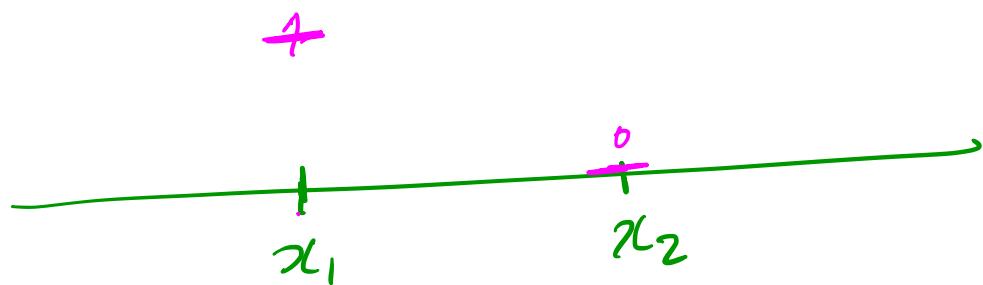


$$C = \{1\}$$

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

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

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

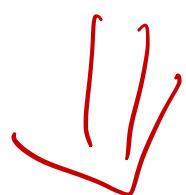


$$C = \{x_1, x_2\} \quad (\text{assume } x_1 < x_2)$$

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

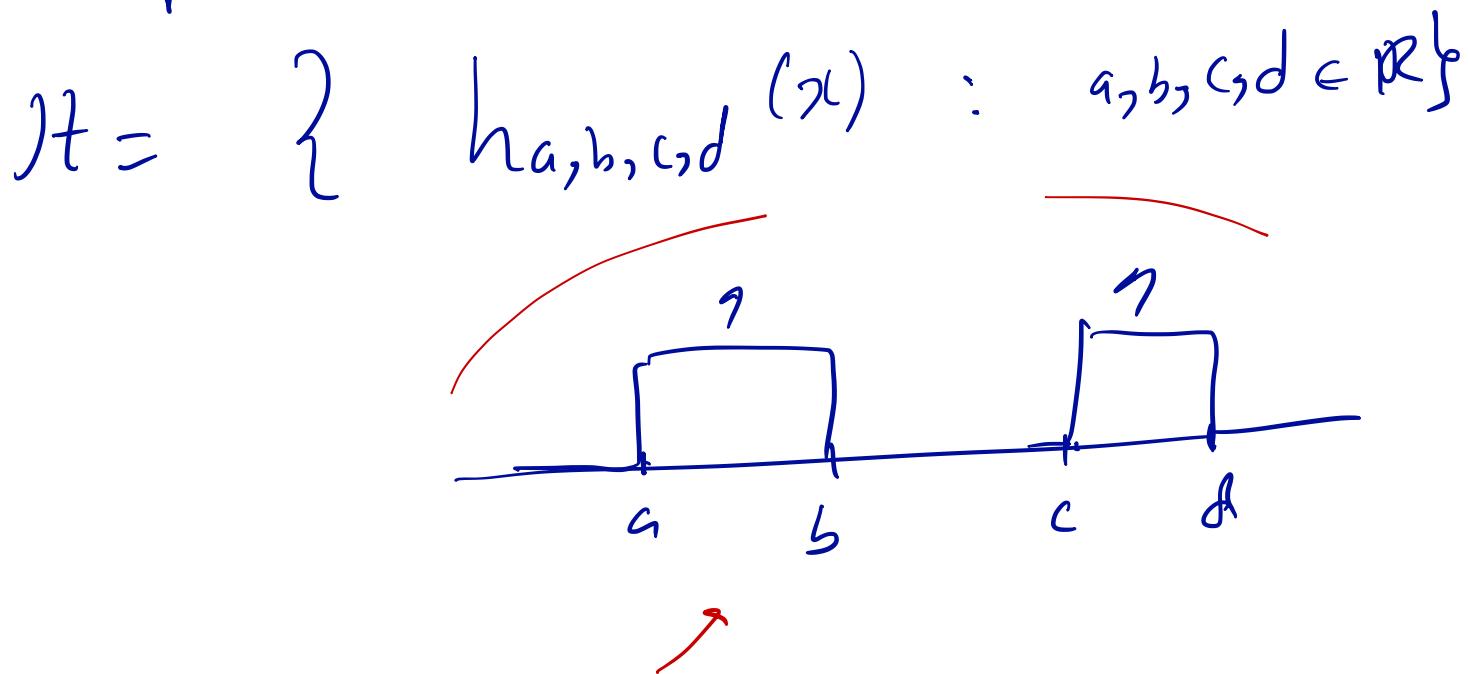
$\hookrightarrow$  there is no  $h_{\mathcal{H}}(x)$   
that produces this  
binary 2-tuple

$$\check{VC\_dim(\mathcal{H}) = }$$



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  $\mathcal{H}$ ?

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

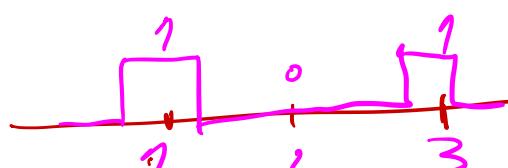
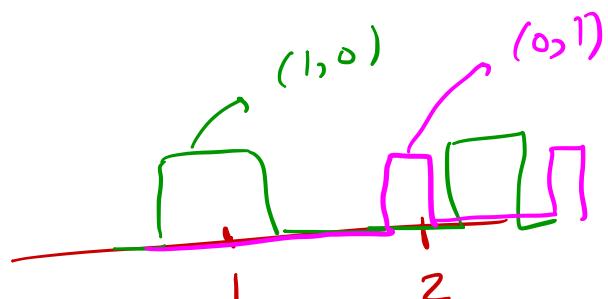
$k=2 \rightarrow \text{YES}$

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

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

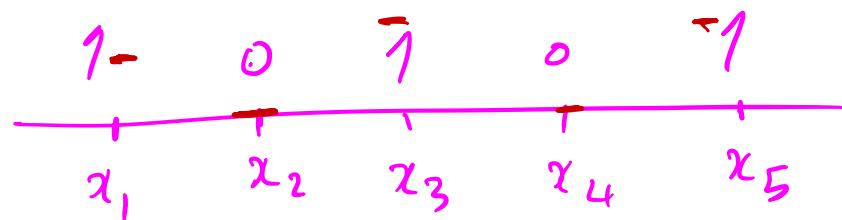
$k=5 \rightarrow \text{NO}$

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

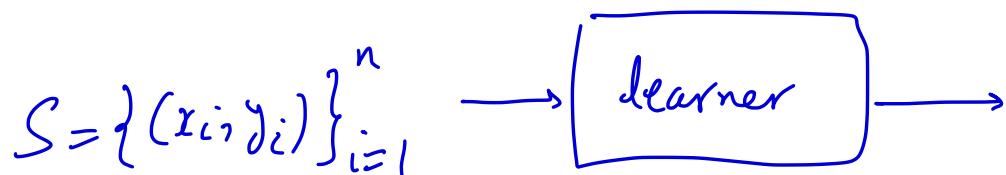


$k=5$

→  $h$  should have 3 jumps from 0 to 9.



## Lecture 23:



$$\min_h \Pr_{(x,y) \sim D} \{ h(x) \neq y \}$$

↓  
the first relaxation:

$$\min_h \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{h(x_i) \neq y_i\} \rightarrow \text{bad} = 0$$

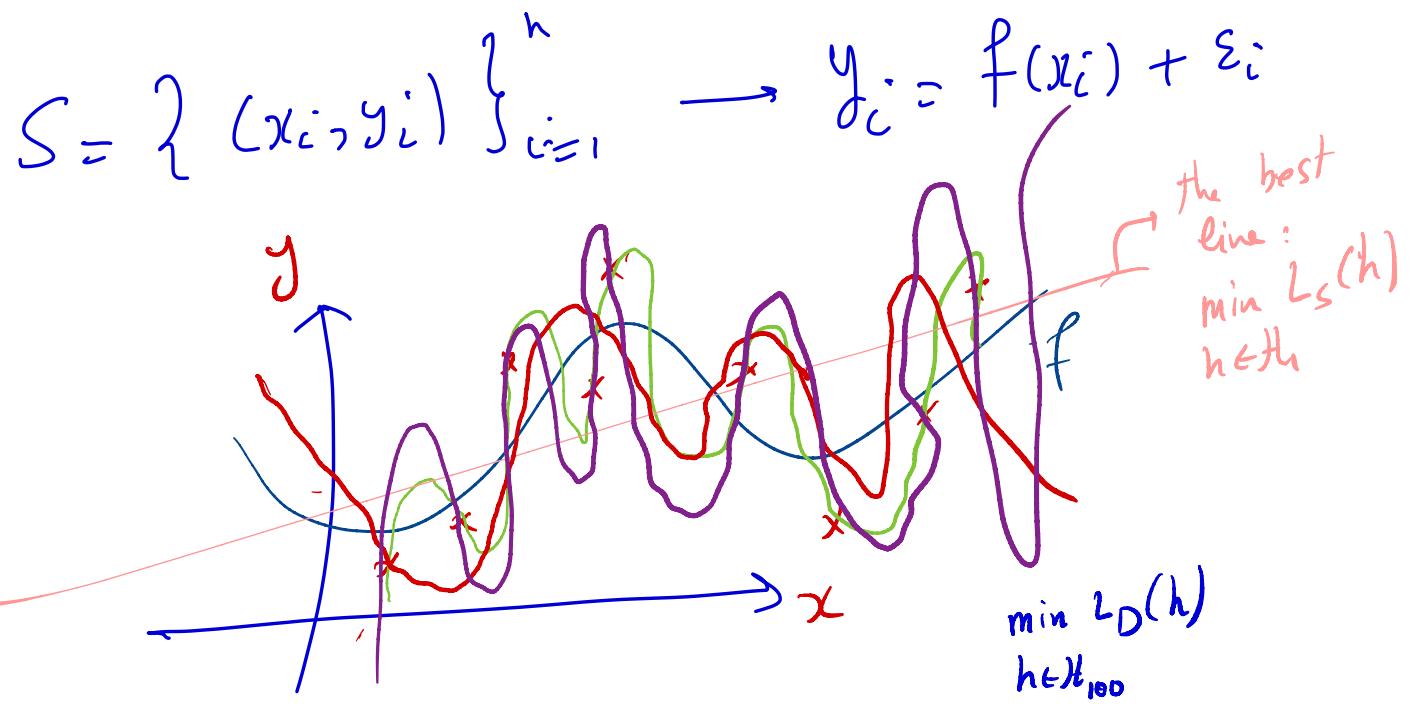
↓  
restricted ERM\_H

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

the predictor is inside the class H.

Example:

$$y = f(x) + \varepsilon \quad E[\varepsilon] = 0, \quad \text{var}[\varepsilon] = 6^2$$



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

$\min_h L_D(h) \Rightarrow$  the minimizer is  $f$

$$\begin{aligned} \min_{h \in H} L_D(h) &= E[(y - f(x))^2] \\ &= E[\varepsilon^2] = 6^2 \end{aligned}$$

training error

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

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

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

examples for  $\mathcal{H}$ :

$\mathcal{H}_1 = \left\{ \begin{array}{l} \text{all the linear functions} \\ \text{i.e. polynomials of degree 1} \end{array} \right\}$

↳ since  $\mathcal{H}_1$  is not rich enough, we will have underfitting

$\Rightarrow$  Hence the function class  $\mathcal{H}$  that we're going to use, should be sufficiently rich/complex.

$$H_6 = \{ \text{all the polynomials of degree 6} \}$$

⇒ the problem with using rich/complex function classes is: overfitting

In summary: we have 2 problems:

Problem 1:

if  $H$  is not rich enough, it'll underfit:

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

Solution: use a class  $H$  that's sufficiently rich.

Problem 2:

if  $H$  is complex, then it may suffer from overfitting:

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

*training error*

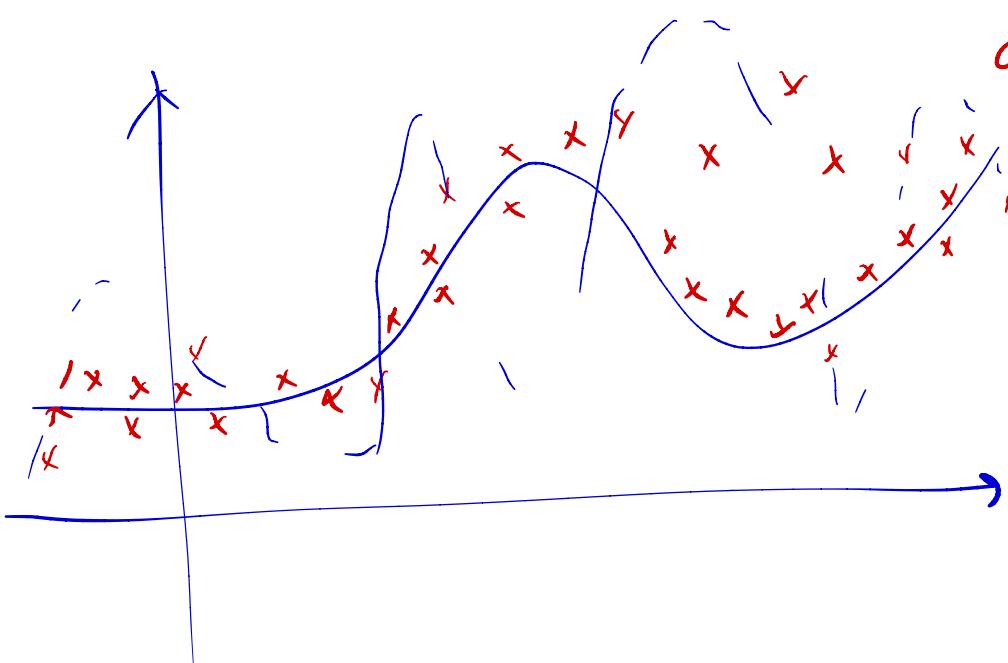
$$\min_{h \in \mathcal{H}} L_S(h) \rightarrow h_S^* \rightarrow L_S(h_S^*) \approx 0$$

*very rich/complex*

$$\min_h L_D(h) \ll L_D(h_S^*)$$

*true error of  
what we learned  
from fitting the  
function class  $\mathcal{H}$*

*to the training  
data S.*



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

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

$$\underbrace{\hspace{10em}}_{h \in \mathcal{H}} \quad \underbrace{\hspace{10em}}_{h \in \mathcal{H}}$$

$\Rightarrow$  to avoid overfitting we should have a sufficient number of training data points.

$$\text{for every } h \in \mathcal{H} \rightarrow \underbrace{L_S(h)}_{\approx L_D(h) \pm \epsilon}$$

$$\min_{h \in \mathcal{H}} \underbrace{L_S(h)}_{\approx \min_{h \in \mathcal{H}} L_D(h) \pm \epsilon} \approx \min_{h \in \mathcal{H}} L_D(h) \pm \epsilon$$

---

It seems that to avoid overfitting the number of training data points should be sufficiently large

$\hookrightarrow$  - depends on  $\cdot \mathcal{H}$

- depends on the accuracy  $\epsilon$

PAC learning (PAC = Probably Approximately Correct)

We say that a function class  $\mathcal{H}$  is PAC-learnable if for any  $\epsilon, \delta > 0$ , there exists a number  $n_0(\epsilon, \delta)$ , such that for any distribution of data,  $D$ , the following holds: If  $n \geq n_0(\epsilon, \delta)$ , then with probability  $1 - \delta$  we have that:

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

$$\underbrace{\min_{h \in \mathcal{H}} L_S(h)}_{\approx} \approx \underbrace{\min_{h \in \mathcal{H}} L_D(h)}_{\approx}$$

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:

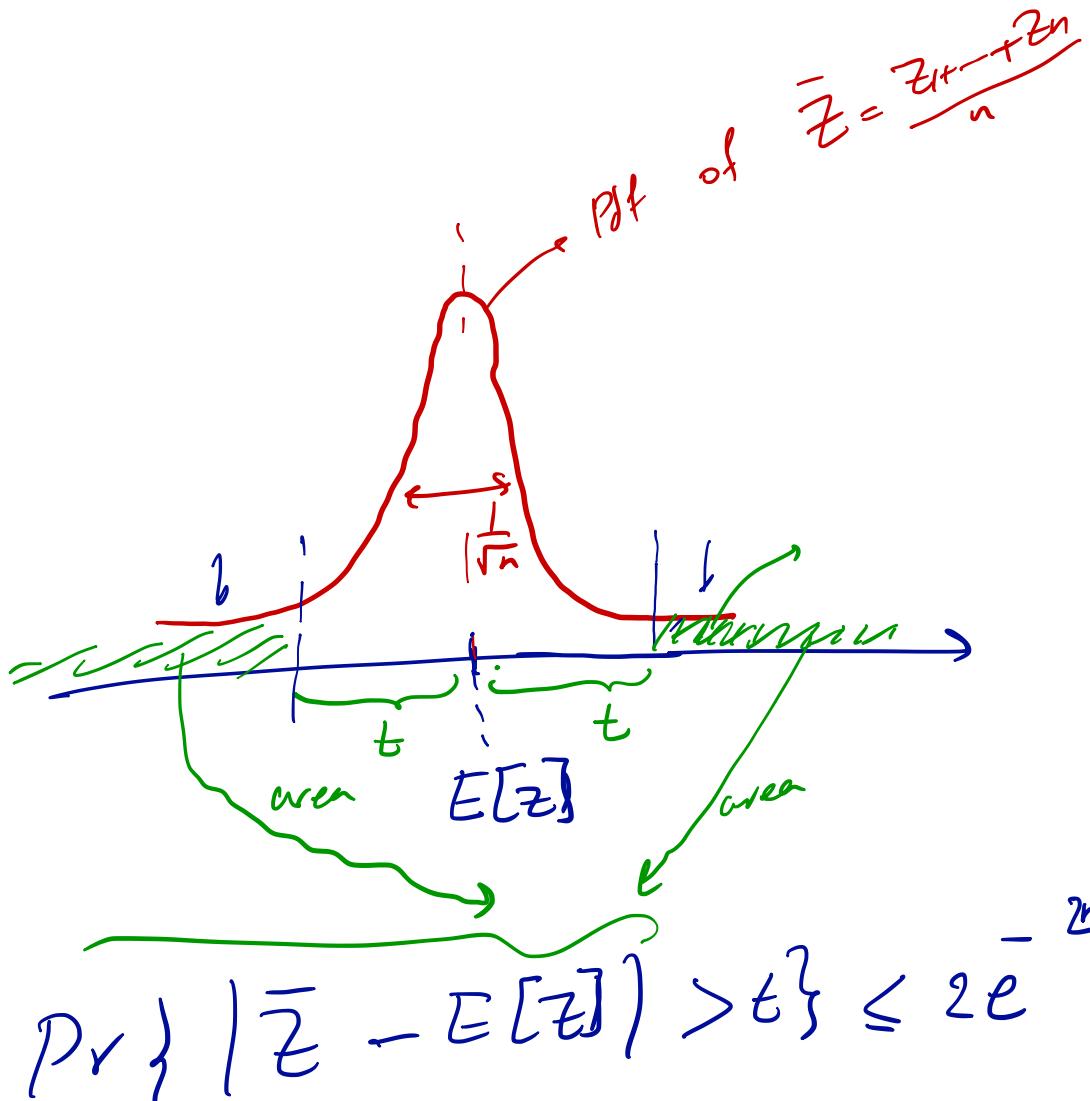
Let  $Z_1, Z_2, \dots, Z_n$  be  $n$  iid random variables that take value in the unit interval (i.e.  $Z_i \in [0, 1]$ ).

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

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

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

↑ equivalent

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



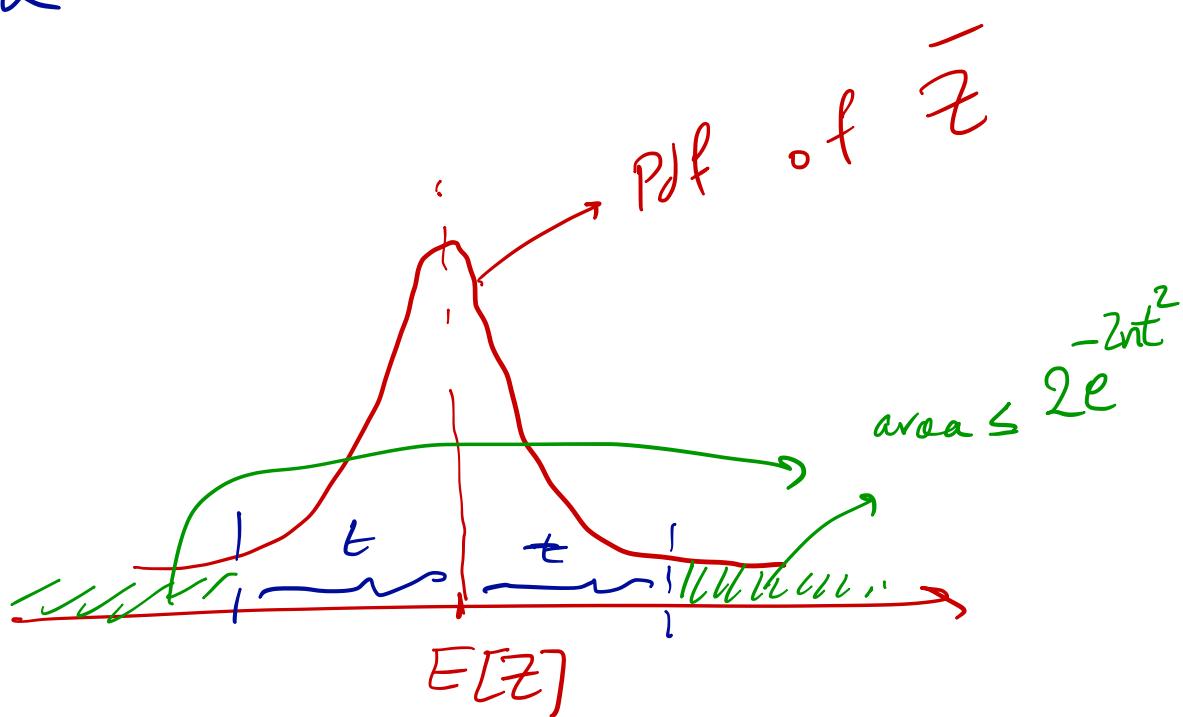
Recall that Hoeffding's inequality was given as follows:

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:

$$2e^{-2n\epsilon^2} \leq \delta$$
$$n \geq \frac{1}{2\epsilon^2} \log \frac{2}{\delta}$$

$\left( \log \frac{2}{\delta} = -\log_2 \delta \right)$

So the final statement is:

Given any  $\epsilon, \delta$ , 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

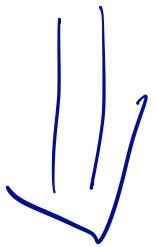
$$\Pr \{ A_1 \} + \Pr \{ A_2 \} + \dots + \Pr \{ A_m \} \leq \delta \quad (6)$$

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

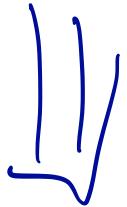
Now, to guarantee (6), it is sufficient to choose  $n$  large enough s.t. for every  $i$

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- $\epsilon$ -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

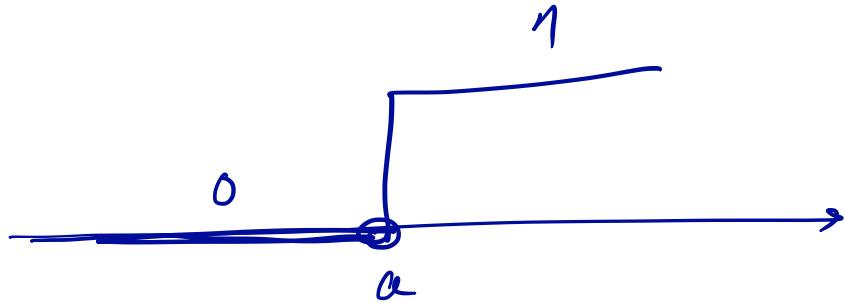
$$0 \leq L_D(h_S^*) - \min_{h \in \mathcal{H}} L_D(h) \leq 2\epsilon.$$

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 24:

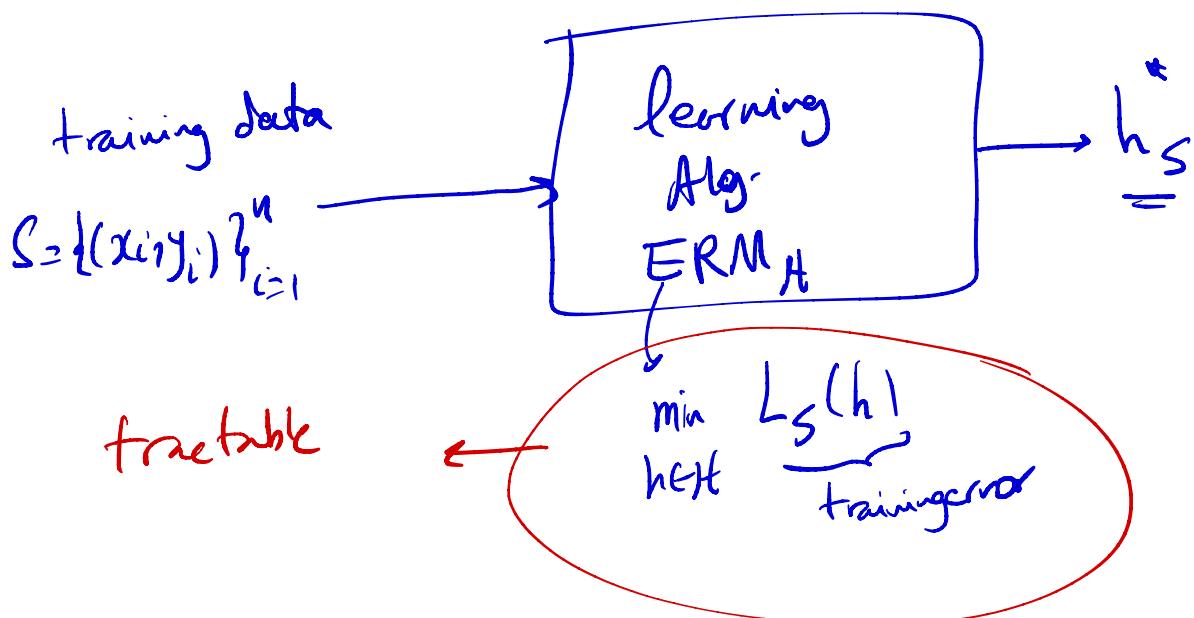
$$\mathcal{H} = \{h_1, \dots, h_m\}$$

PAC learning:

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

Then: For any data distribution, with probability  $1-\delta$ , the outcome of  $\text{ERM}_{\mathcal{H}}$ , which we denote by  $h_S^*$ , satisfies:

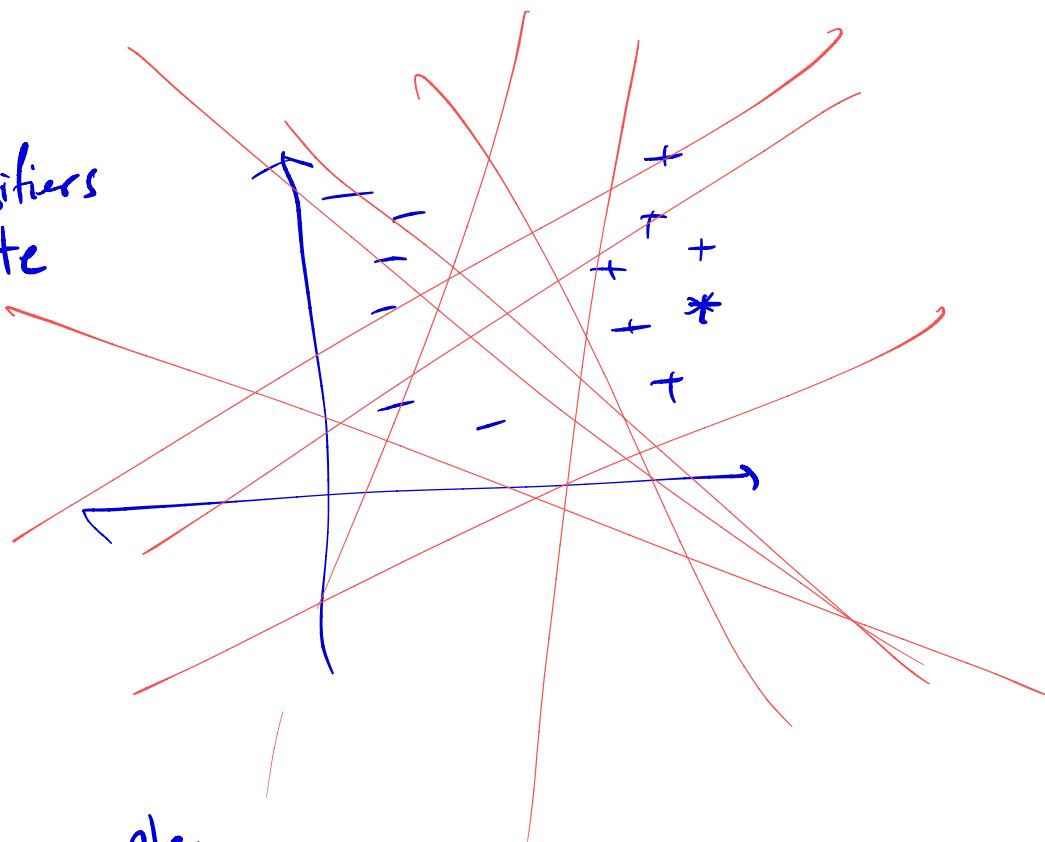
$$\min_{h \in \mathcal{H}} L_D(h) \leq L_D(h_S^*) \leq \min_{h \in \mathcal{H}} L_D(h) + \epsilon$$



So far, we've shown that classes  $H$  with finitely many functions are PAC learnable.

But in practice, all the function classes that we use are not finite.

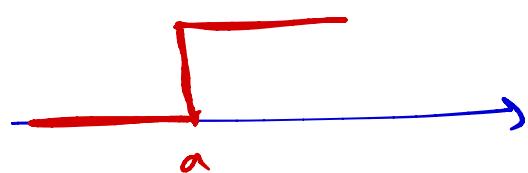
examples:  
linear classifiers  
with infinite



an other example:  
consider the following simple class

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

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



## The Fundamental Theorem of Learning Theory:

For any function class  $H$ , if the number of training data points is larger than  $n_0(\epsilon, \delta) = 4 \cdot \frac{\log \frac{1}{\delta} + (\text{VC-dim}(H))}{\epsilon^2}$ , then quantifying the "complexity" of  $H$

$H$  is PAC-learnable: For any distribution  $D$  on data, with  $1-\delta$ :

$$L_D(h_S^+) \leq \underbrace{\min_{h \in H} L_D(h)}_{\text{intractable}} + \epsilon$$

where  $h_S^+ = \underset{h \in H}{\operatorname{argmin}} \widehat{L}_S(h)$ .

As long as we have sufficiently many data points, we'll be fine with solving the "tractable" ERM<sub>H</sub> problem and we lose at most a small value  $\epsilon$  with respect to the best attainable accuracy.

⇒ Machine learning is possible as long as we have sufficient data.

$V_C$  - Dimension :

Vapnik Chernovskii

Example: Consider the following function

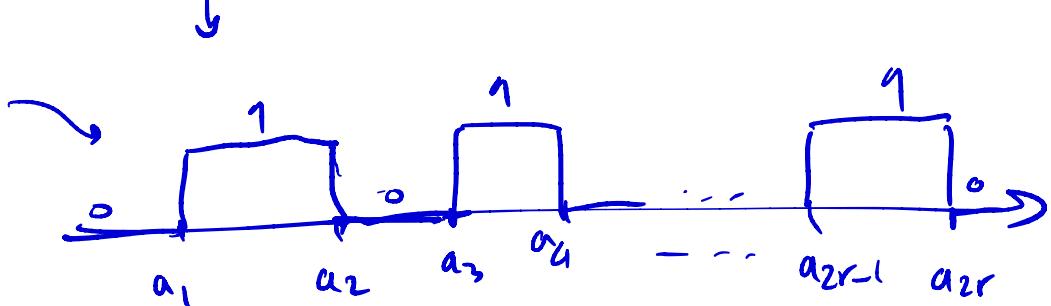
classes:

$$H_1 = \{ h_a(x), a \in \mathbb{R} \} \rightarrow \begin{array}{c} h_a(x) \\ \text{---} \quad | \\ \text{---} \end{array} \quad a$$

$$H_2 = \{ h_{a,b}(x), a, b \in \mathbb{R} \} \rightarrow \begin{array}{c} h_{a,b}(x) \\ \text{---} \quad | \\ \text{---} \quad | \\ a \quad b \end{array}$$

$$H_3 = \{ h_{abcd}(x), a, b, c, d \in \mathbb{R} \} \rightarrow \begin{array}{c} h_{abcd}(x) \\ \text{---} \quad | \\ \text{---} \quad | \\ a \quad b \quad c \quad d \end{array}$$

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



Intuitively, in terms of complexity  
we should have

$$H_1 < H_2 < H_3 < H_4$$

VC-dim is a precise mathematical notion to quantity "complexity".

In order to define VC-dim we need two other definitions (restriction, shattering).

Definition (restriction),

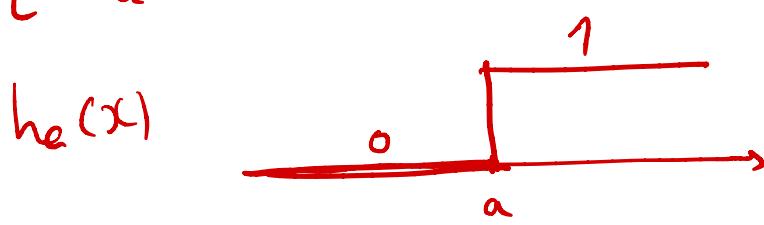
Let  $\mathcal{H}$  be a class of functions from a domain  $X$  to  $\{0, 1\}$  ( $\mathcal{H} : h : X \rightarrow \{0, 1\}$ ).

Let  $C = \{x_1, \dots, x_k\} \subseteq X$ . Then, the restriction of  $\mathcal{H}$  to  $C$  is the set of all the  $k$ -tuples of the following form:

$$\mathcal{H}_C = \left\{ \underbrace{(h(x_1), h(x_2), \dots, h(x_k))}_{k\text{-tuple}} ; h \in \mathcal{H} \right\}$$

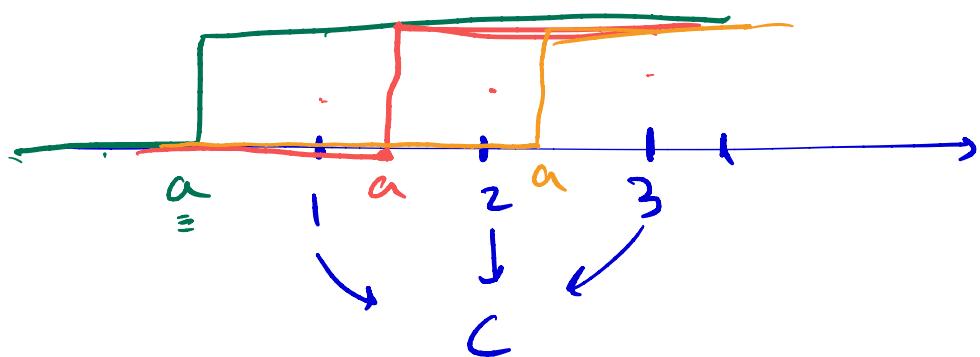
Example:

$$\mathcal{H} = \{h_a(x) ; a \in \mathbb{R}\}$$



Let  $C = \{1, 2, 3\}$ . What is  $H_C$ ?

To find  $H_C$  we need to go over all the functions  $h \in H$ .



$$\text{if } a \leq 1 \rightarrow (h_a(1), h_a(2), h_a(3)) \\ = (\underline{1}, 1, 1) \leftarrow$$

$$\text{if } 1 < a \leq 2 \rightarrow (h_a(1), h_a(2), h_a(3)) \\ = (\underline{0}, 1, 1) \leftarrow$$

$$\text{if } 2 < a \leq 3 \rightarrow (h_a(1), h_a(2), h_a(3)) \\ = (0, \underline{0}, 1) \leftarrow$$

$$\text{if } a > 3 \rightarrow (\underline{0}, 0, 0) \leftarrow$$

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

We note that if  $C = \{x_0 -> x_k\}$

then  $|\mathcal{H}_C| \leq 2^k$

↳ recall that there are  $2^k$  binary  $k$ -tuples.

$$1 \leq |\mathcal{H}_C| \leq \overbrace{2^k}^{\text{complexity} \uparrow}$$



**Definition (shattering):** We say that a function class  $\mathcal{H}$  shatters a set  $C$  of size  $k$  if  $|\mathcal{H}_C| = 2^k$ .

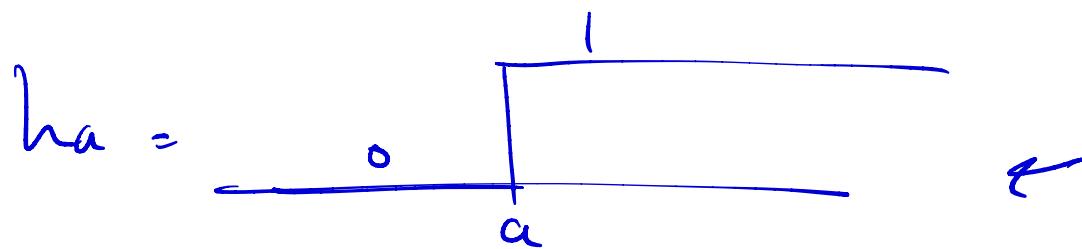
Definition (VC-dimension) :

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

---

Example:

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



In order to find the  $\text{vc-dim}(\mathcal{H})$  we need to find the largest  $K$  for which  $\mathcal{H}$  shatters a set  $C$  of size  $K$ :

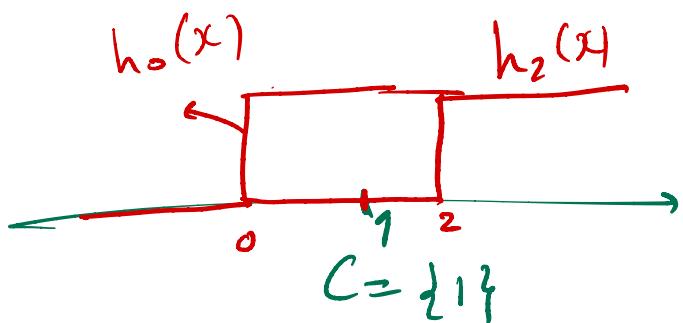
---

$K=1 \rightarrow$  Is there a set  $C_1$  that is shattered by  $\mathcal{H}$ ?

$$C = \{x_1\}$$

$$|C|=1, \quad H_C = \{(h(x_1)) ; h \in \mathcal{H}\}$$

$$C = \{x_1\} \quad = \{0, 1\}$$

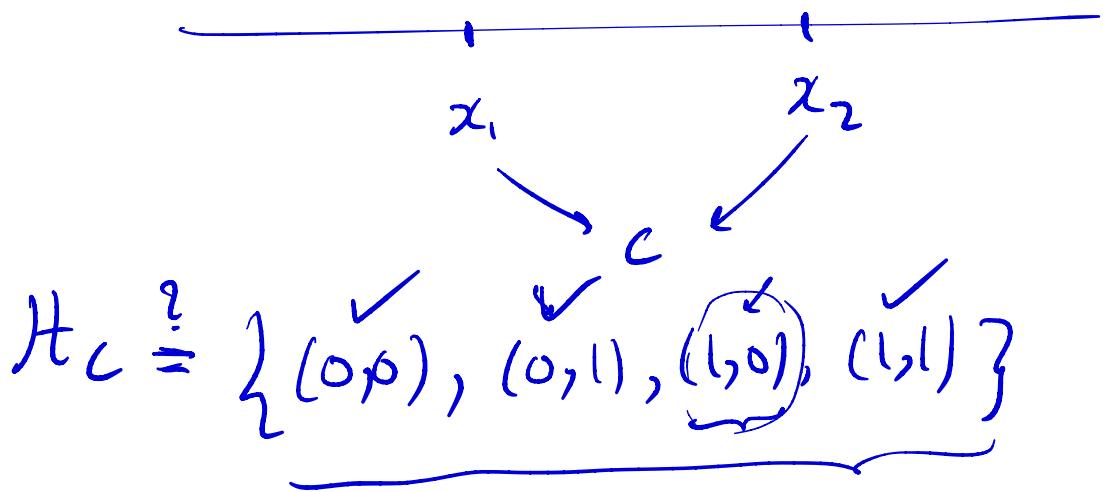


$$\begin{cases} h_0(1) = 1 \\ h_2(1) = 0 \end{cases} \rightarrow H_C = \{0, 1\}$$

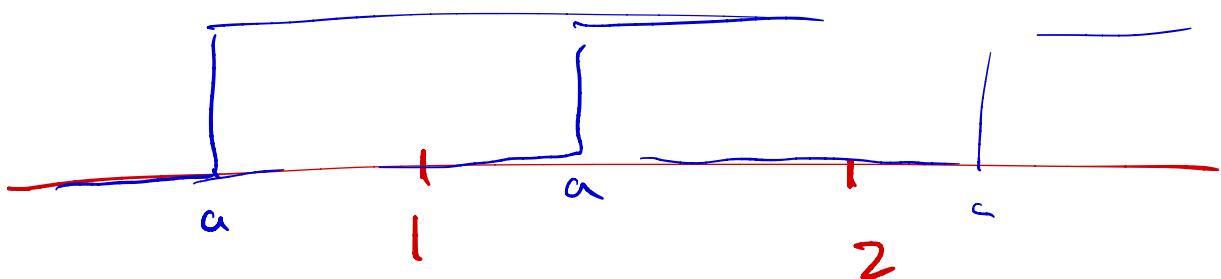
$\Rightarrow C = \{1\}$  is shattered by  $H$ .

Now consider  $k=2$ .

Is there a set  $C$  of size 2 that is shattered by  $H$ ?



let's take  $x_1=1$  and  $x_2=2$



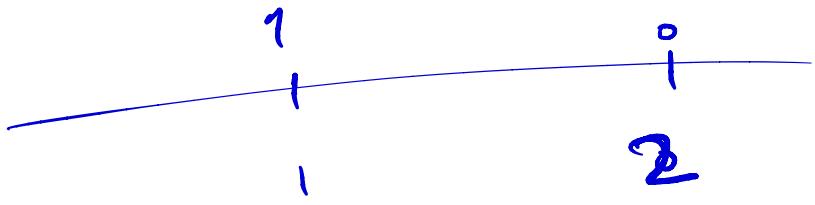
if  $a > 2 \Rightarrow (h_a(1), h_a(2)) = (0,0)$

if  $a < 1 \Rightarrow (h_a(1), h_a(2)) = (1,1)$

if  $1 < a \leq 2 \Rightarrow (h_a(1), h_a(2)) = (0,1)$

is there any function that generates  $(1,0)$ ?

→ No since all the functions in  $H$  are increasing.



We can show that for my set

$C = \{x_1, x_2\}$ ,  $x_1 < x_2$ , the 2-tuple  $(1, 0)$

can not be generated by any  $h \in H$ .

$\Rightarrow \nexists C$  of size 2 which is shattered by  $H$ .

$$k=1 \quad \checkmark$$

$$\rightarrow V_C - \dim(H) = 1.$$

$$k=2 \quad \times$$

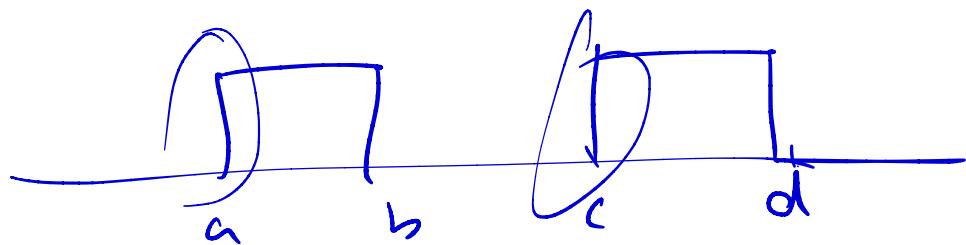
$$k=3 \quad \times$$

1

1

Example :

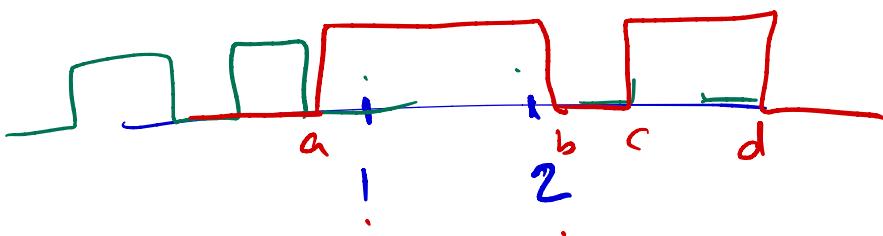
$$\mathcal{H} = \left\{ h_{a,b,c,d}(x) : a, b, c, d \in \mathbb{R} \right\}$$



$$k=1 \rightarrow \checkmark$$

of size 2

$k=2 \rightarrow$  is there a set  $C \supseteq$  that's shattered by  $\mathcal{H}$ ?

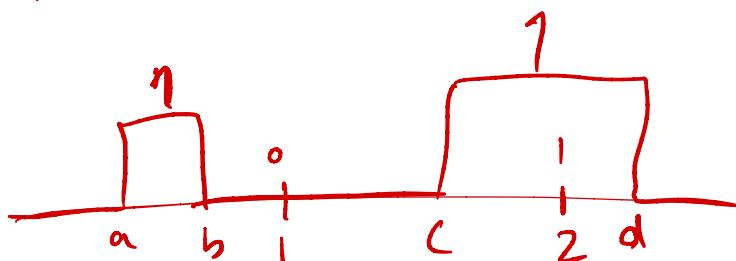


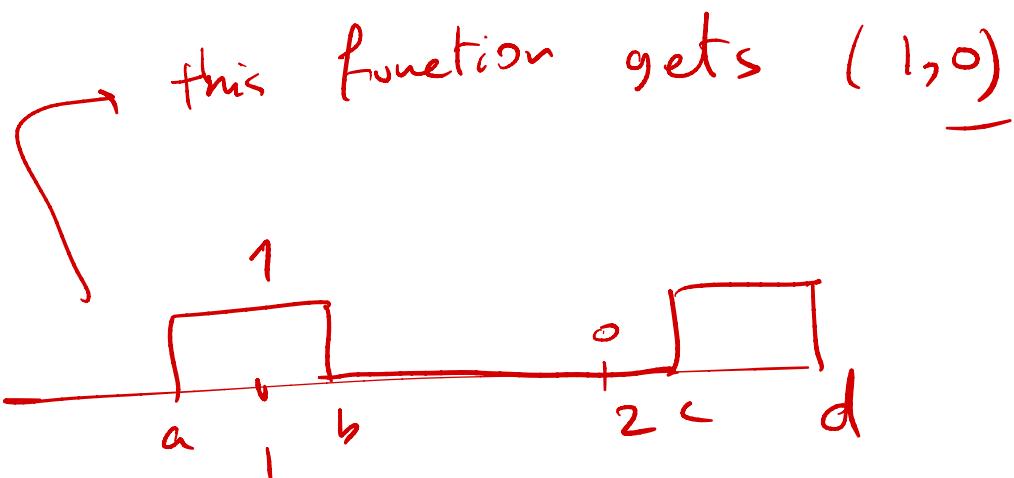
$$\mathcal{H}_C = \{(0,0), (0,1), (1,0), (1,1)\}$$

if  $d < 1$  then we get  $(0,0)$

if  $a < 1 < 2 < b$  then we get  $(1,1)$

$a, b < 1$  and  $c < 2 < d \rightsquigarrow \underline{(0,1)}$





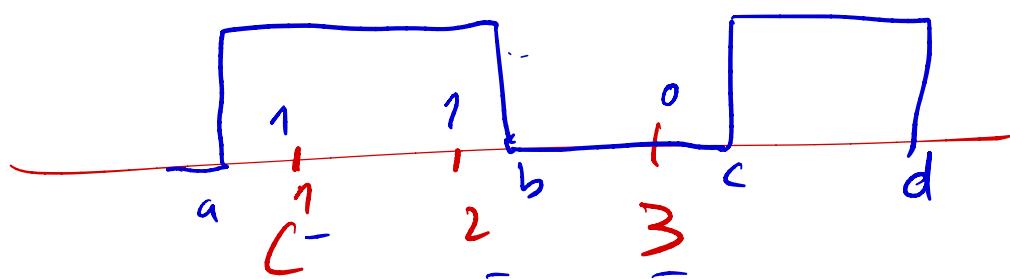
$\mathcal{H}$  shatters the set  $C = \{1, 2\}$ .

$k = 3$  ✓

$C = \{1, 2, 3\}$  will be shattered by  $\mathcal{H}$ .

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

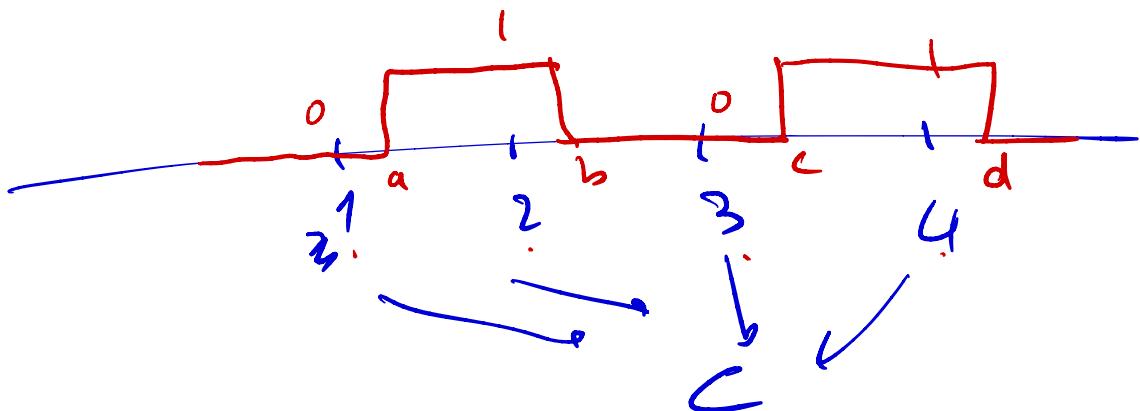
e.g.  $(\underline{1}, \underline{1}, \underline{0})$



$k = \underline{4}$ , the set  $C = \{1, 2, 3, 4\}$   
is shattered by  $H$ .

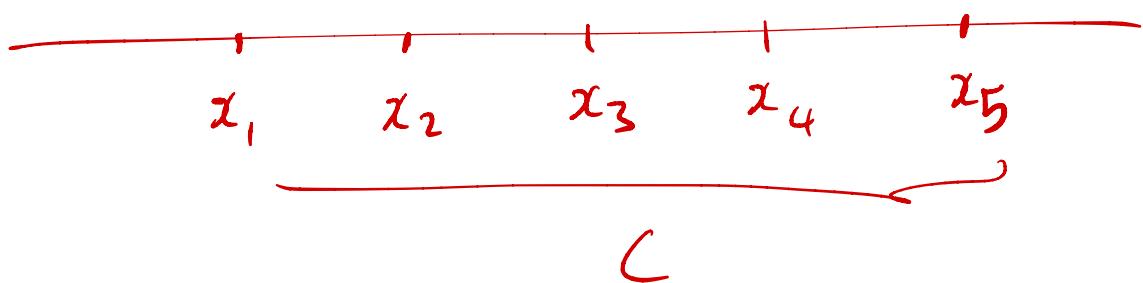
$H_C = \{ \text{all the } 16 \text{ binary 4-tups} \}$

e.g.  $(\underline{0}, \underline{1}, \underline{0}, \underline{1})$

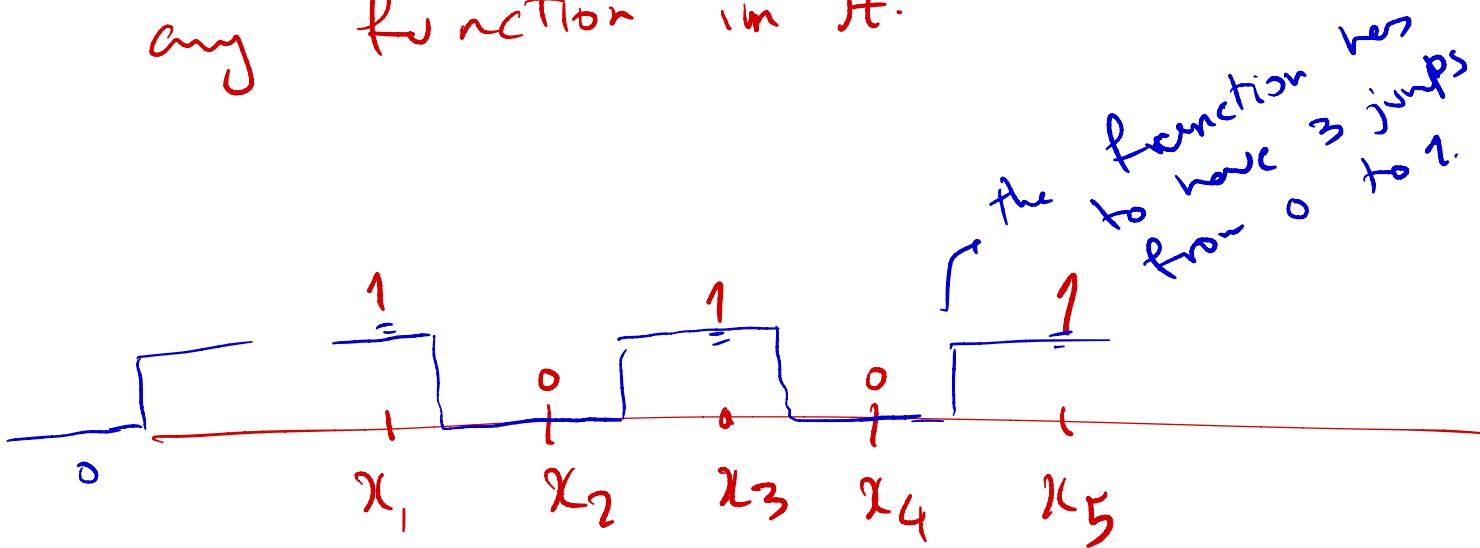


$H$  shatters a. set  $C$  of size 1 ✓  
2 ✓  
3 ✓  
4 ✓  
5 ? X

$k=5$  does  $H$  shatter a set  $C$  of size 5?



We claim that the 5-tuple  $(1, 0, 1, 0, 1)$  can not be generated by any function in  $H$ .



$$\Rightarrow \text{VC-dim}(H) = 4.$$