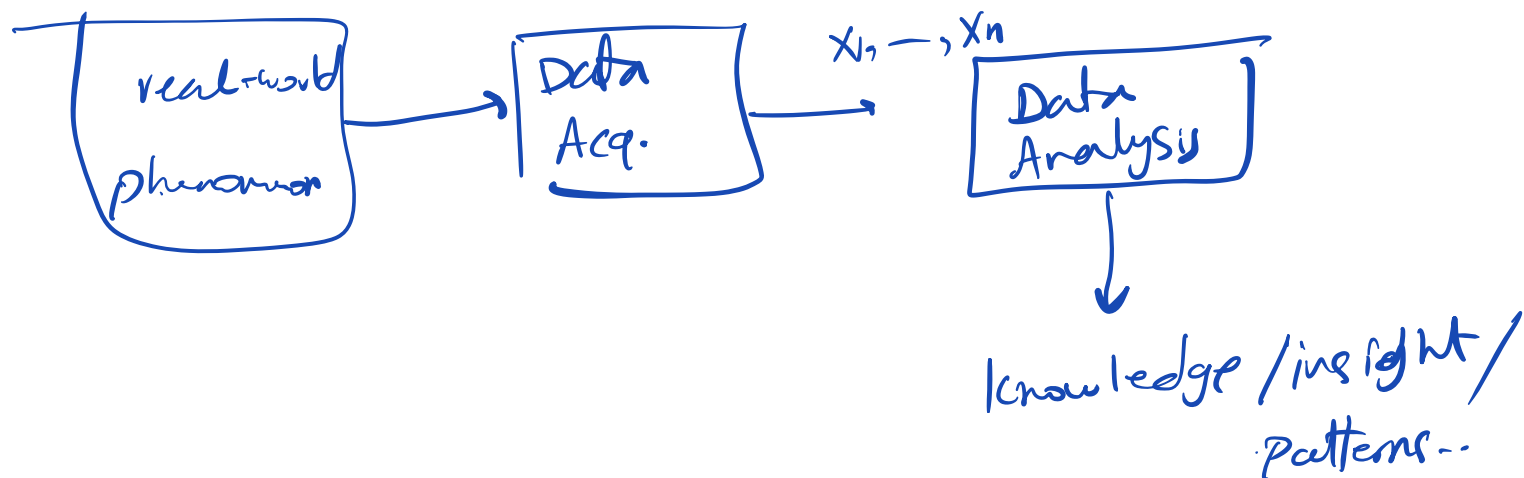


## Lecture 13:



- Modules 1, 2 :
- learned to estimate properties/parameters connected to the underlying phenomenon
  - hypothesis testing: verify different hypotheses connected with the underlying phenomenon.

~~pro~~ 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:

$\left\{ \begin{array}{l} x_i : \text{could be an image} \\ y_i : \text{could be an object inside} \\ \quad \text{the image} \end{array} \right.$

$\left\{ \begin{array}{l} x_i : \text{could be (the text of) an email} \\ y_i : \text{could be whether/not the email} \\ \quad \text{is spam.} \end{array} \right.$

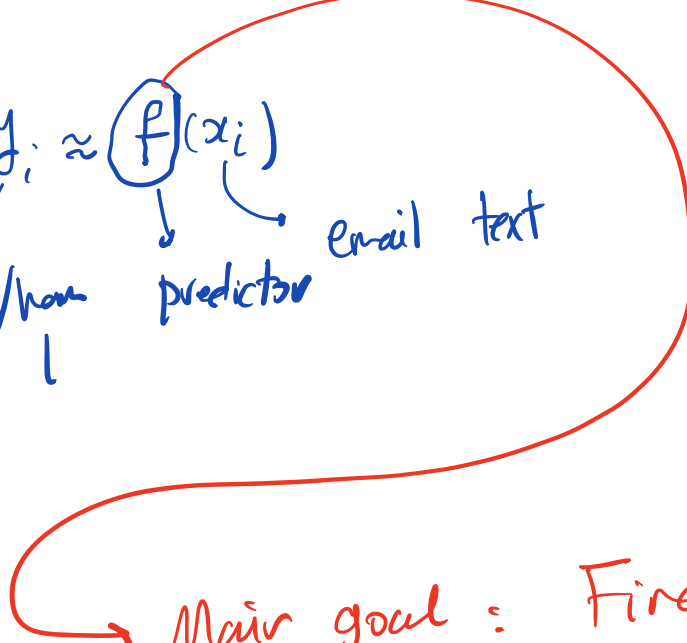
$\left\{ \begin{array}{l} x_i: \text{a vector representing a patient's profile, } x_i = (\text{age, height, weight, BP, ...}) \\ y_i: \text{Cancer / not cancer} \end{array} \right.$

---

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

Goal :

$y_i \approx \underbrace{f}_{\text{predictor}}(\underbrace{x_i}_{\text{email text}})$   
 spam/not-spam  
 0 1


 Main goal : Find the best predictor  $f$ .

---

Setting :

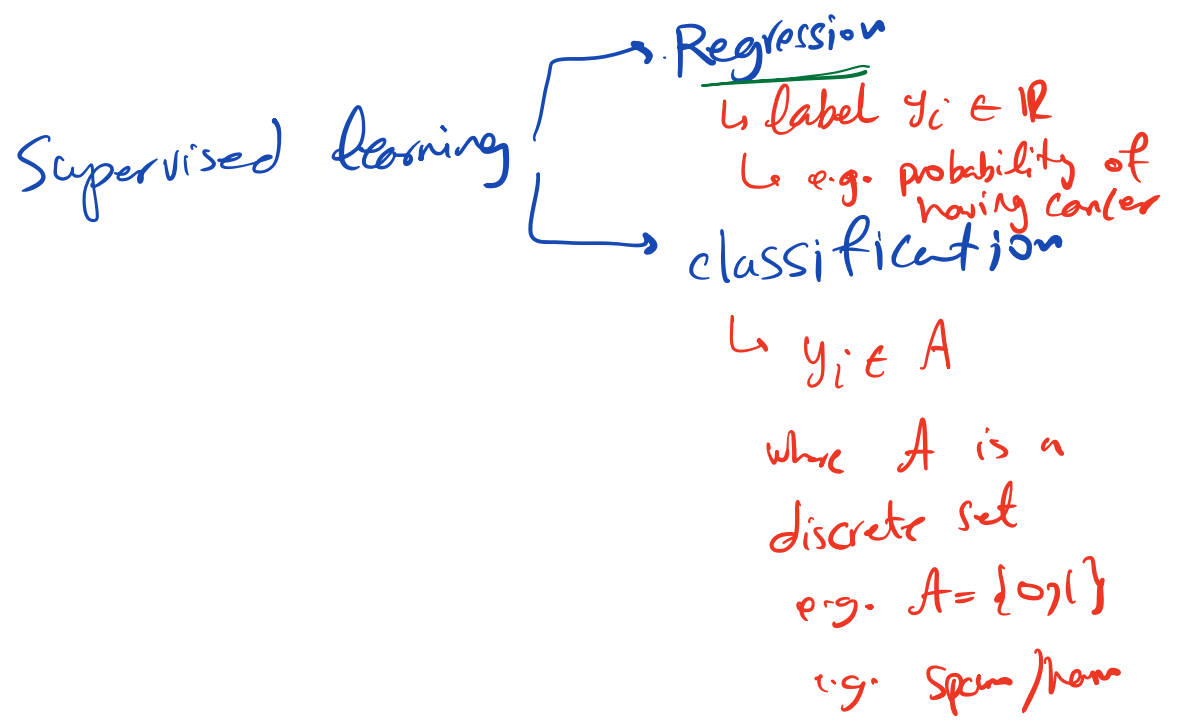
Data :  $\{(x_i, y_i)\}_{i=1}^n$

$x_i = (x_{i1}, x_{i2}, \dots, x_{ip}) \in \mathbb{R}^p$

- image: pixel representation

- profile of the  $i$ -th patient

$y_i = \text{label}$



## Regression:

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

- 1-d regression

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

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

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

$$(p = 1)$$

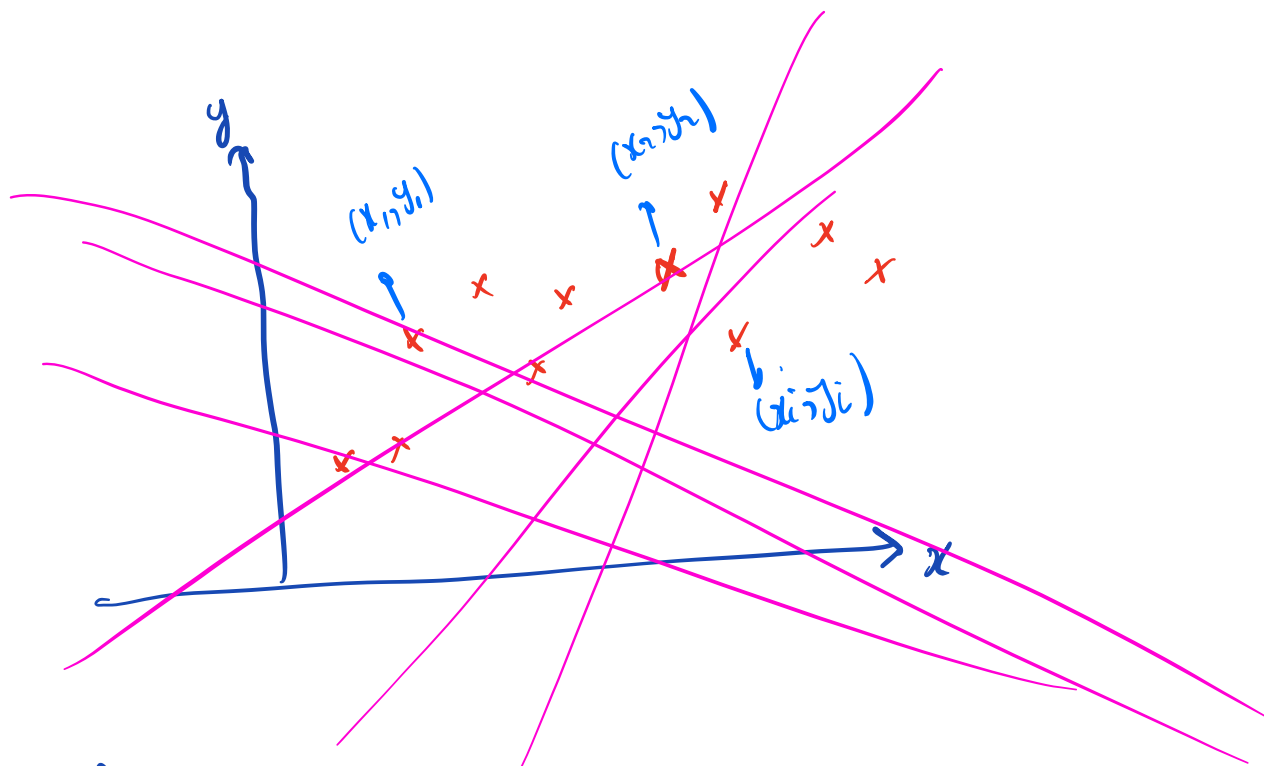
Goal: learn a 'predictive' relation:

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

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

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

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



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

$\underbrace{\hspace{10em}}_{f(x)}$

Approach:

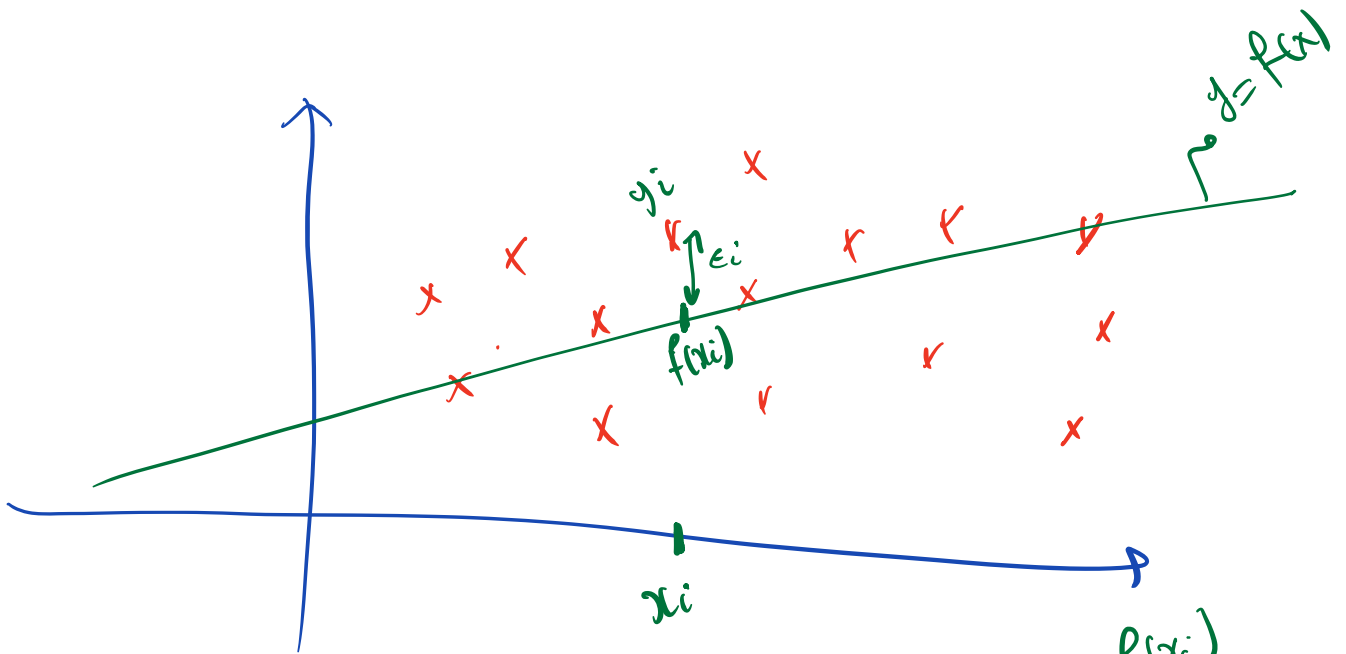
Goal: learn  $y = f(x)$

- 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  $\hat{\beta}_0, \hat{\beta}_1$

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

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



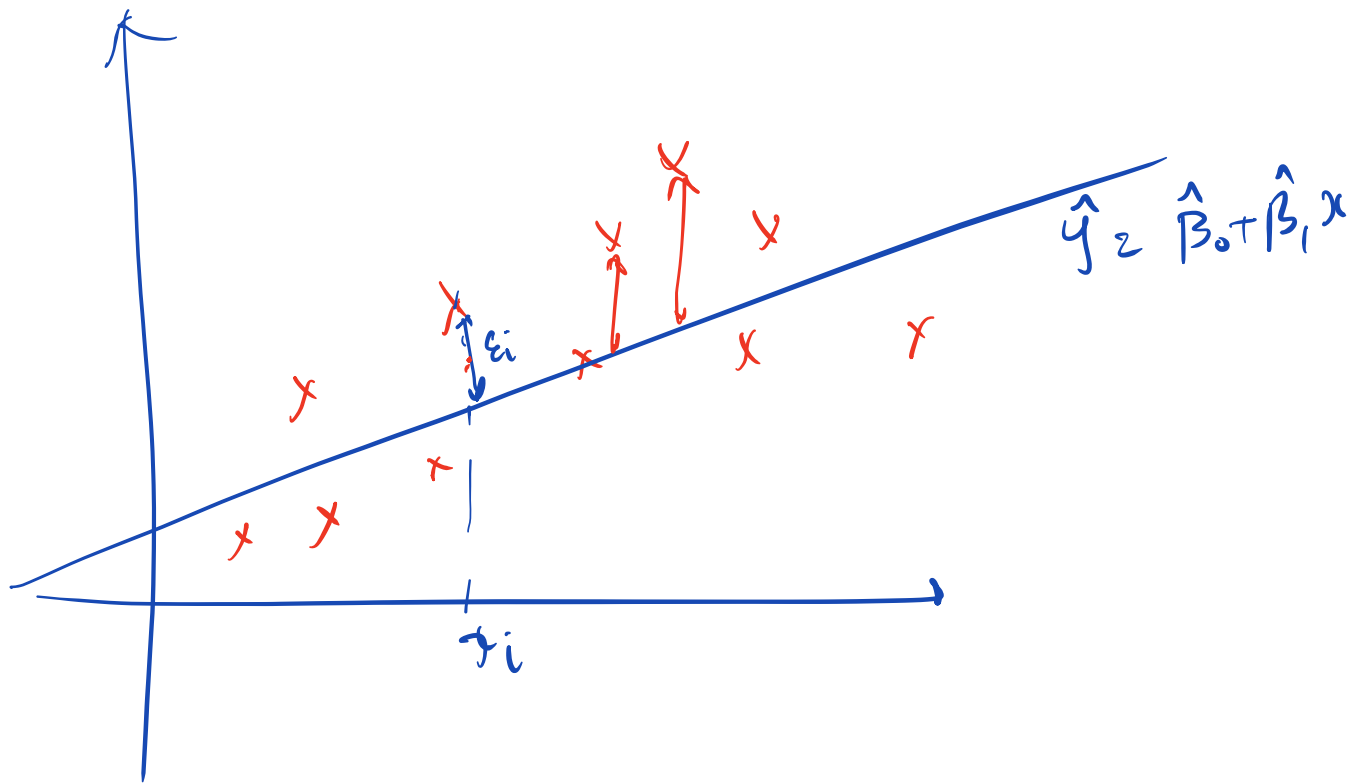
$$y_i = f(x_i) + \epsilon_i \quad \text{error} = \epsilon_i := y_i - \underbrace{f(x_i)}_{\hat{\beta}_0 + \hat{\beta}_1 x_i}$$

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

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  $\hat{\beta}_0, \hat{\beta}_1$  st.

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

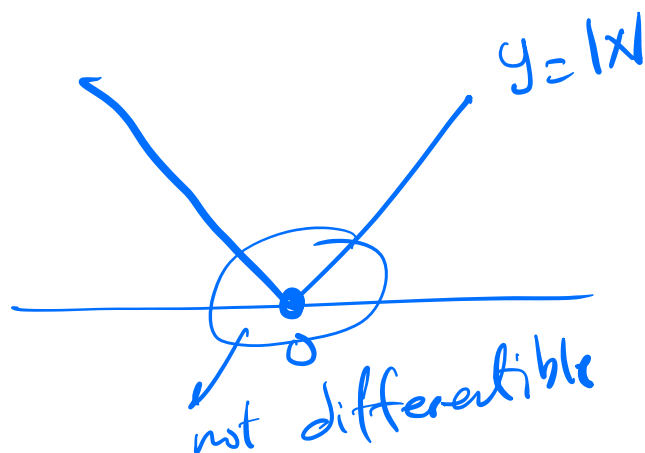
Sum of the errors  $\rightarrow \sum_{i=1}^n |\epsilon_i|$   $\rightarrow$  is minimize



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



- In order to solve the above optimization problem, we'll have to take the derivative wrt.  $\hat{\beta}_0, \hat{\beta}_1$  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:

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