

## Lecture 2: Supervised Learning Setup (Mathematical Formalization)

- (our data)  
- Given labeled examples, find the right prediction of an unlabel example

- Setup: Given a data set D

$$D = \{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n)\}$$

$(\vec{x}_i, y_i) \sim P$  (unknown distribution)

$\vec{x}_i$  → number

$$\vec{x}_i \in \mathbb{R}^d$$

(e.g. spam filter)

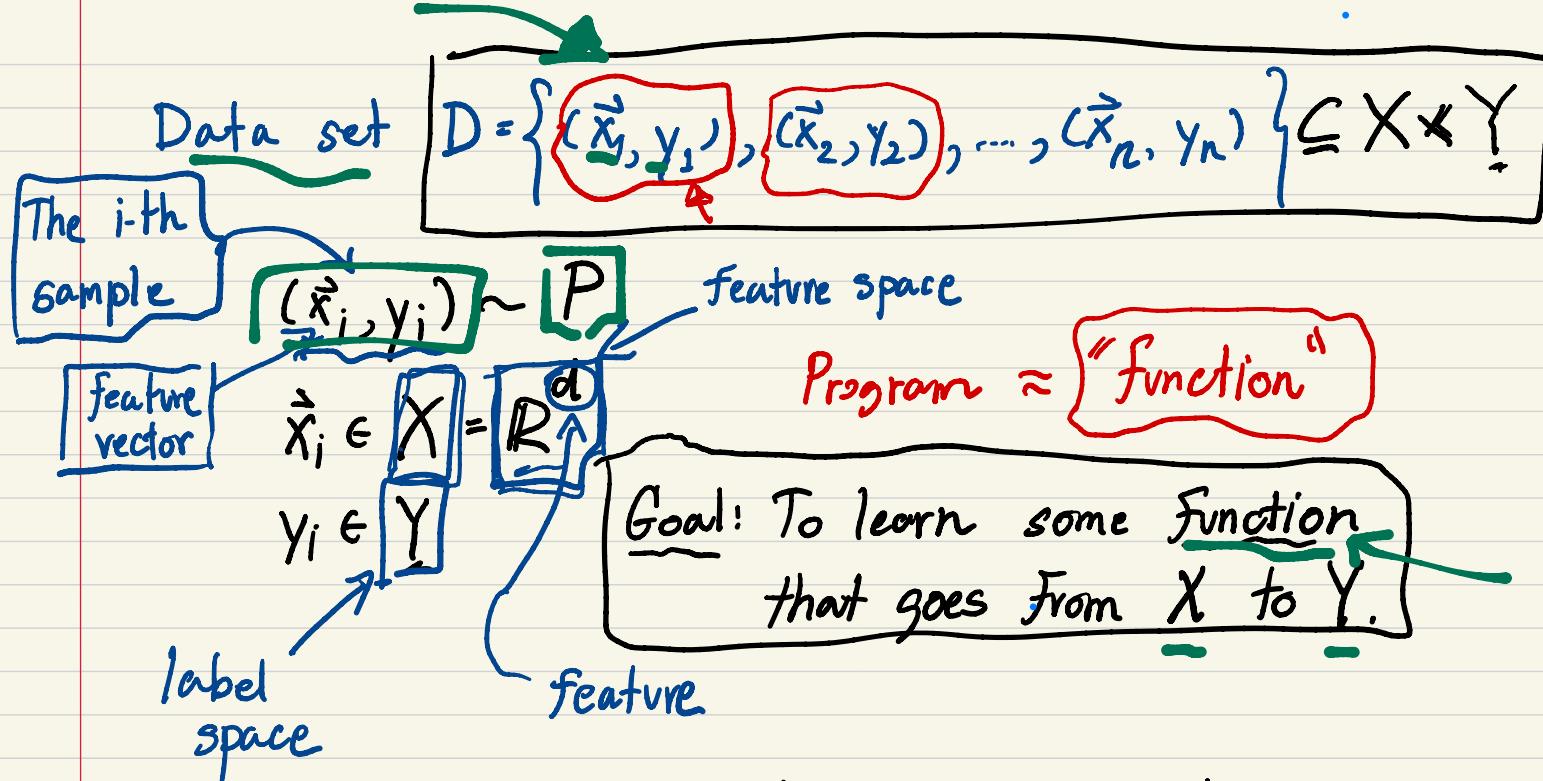
Binary classification:  $y_i \in \{0, 1\}$

(e.g. face detection)

Multi-class classification:  $y_i \in \{0, 1, 2, \dots, k\}$

(e.g. house price prediction)

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



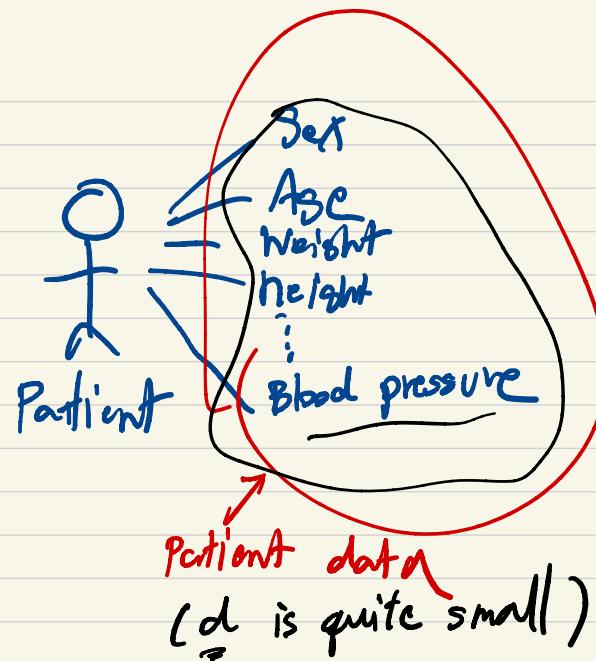
Assumption: Each  $(\vec{x}_i, y_i) \in D$  is i.i.d.  
(independent and identically distributed)

## Examples of Feature Vectors

Patient Data

$$\vec{x}_i = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \\ x_d \end{bmatrix} = \begin{bmatrix} 0 \\ 25 \\ 69 \\ 170 \\ \vdots \\ 110 \end{bmatrix} \begin{array}{l} \text{male / female} \\ \text{age in years} \\ \text{weight in kg} \\ \text{height in cm} \\ \text{Blood pressure} \end{array}$$

$$\vec{x}_i \in \mathbb{R}^{d_x}$$



(B.o.W.)

Text document using Bag-of-Word representation

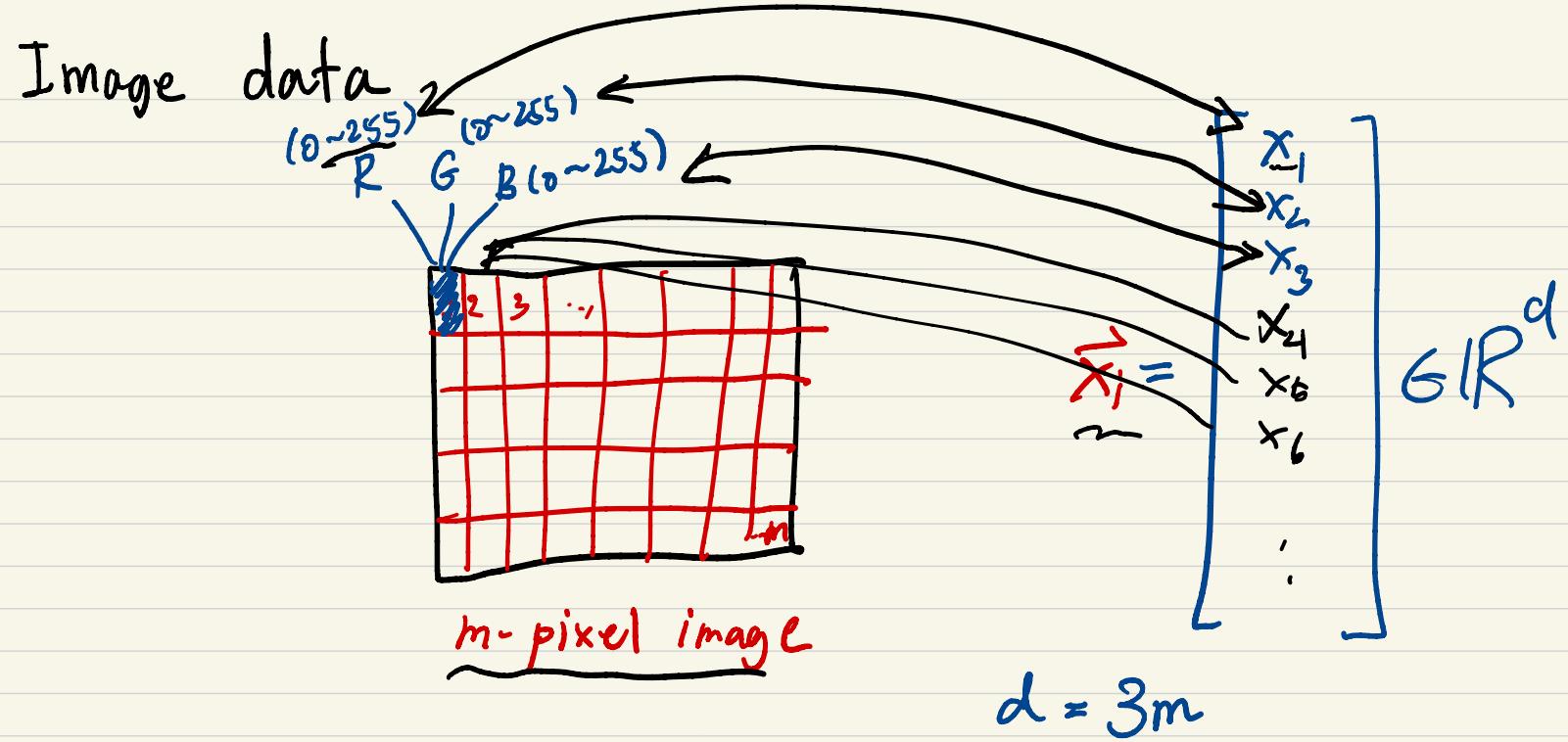
"zebra Ant Ant"  
 "Ant zebra Ant Ant"

↓

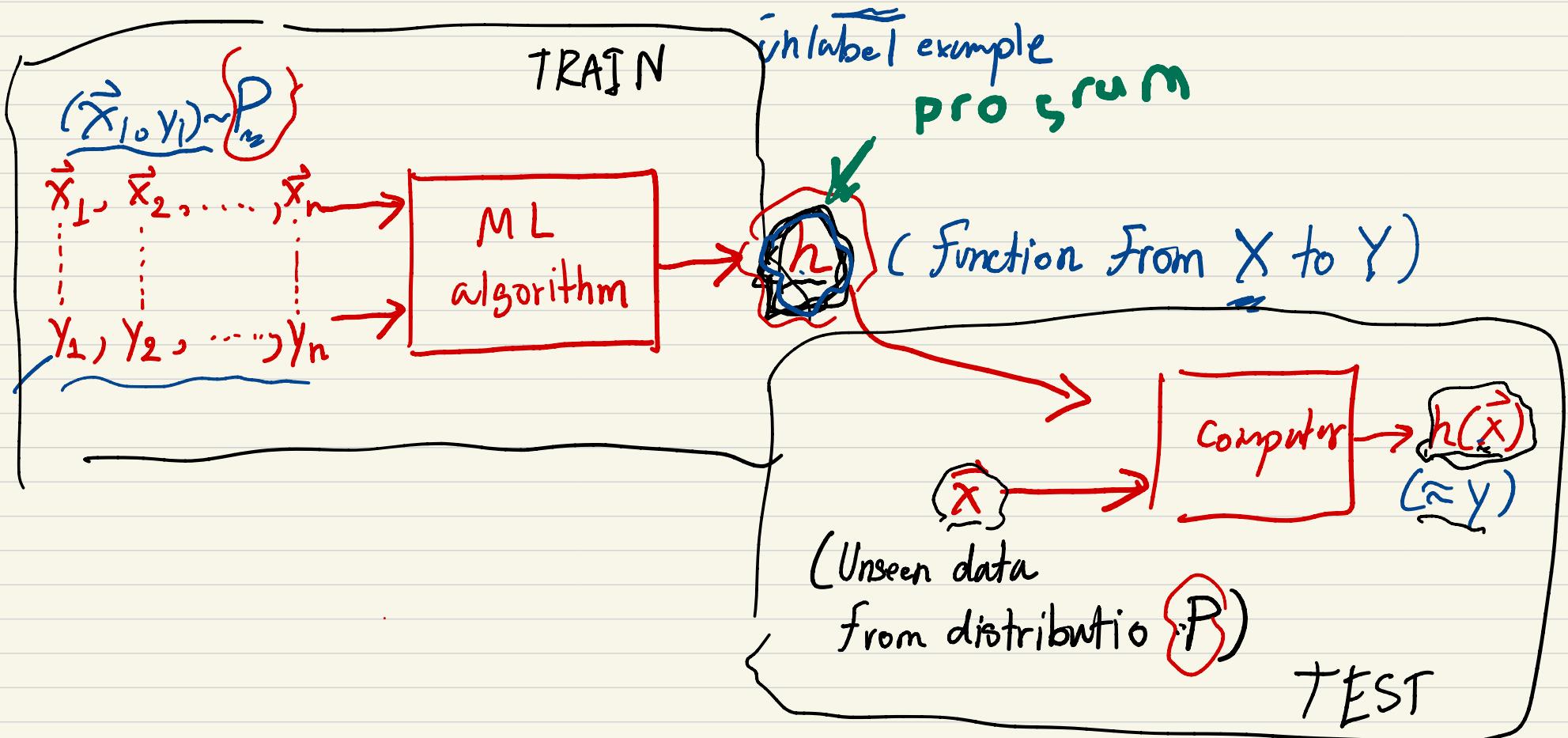
$$\vec{x}_i = \begin{bmatrix} 3 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \begin{array}{l} \dots \text{Ant} \\ \dots \text{zebra} \end{array}$$

$$\vec{x}_i \in \mathbb{R}^{d_x}$$

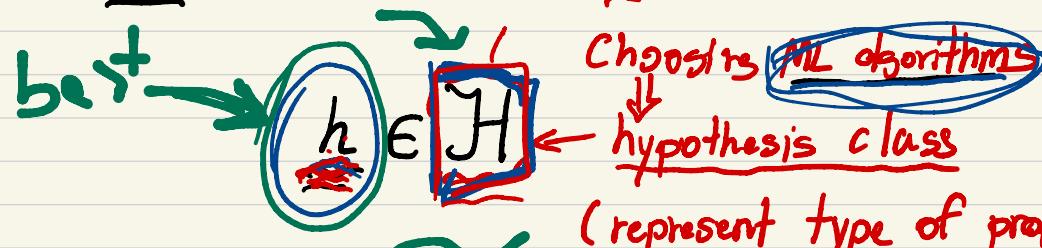
(In English dictionary  
 $d \approx 170\text{ k}$ ).



$$7\text{MP Image} \Rightarrow d = 3 \times 7M = \underline{\underline{21M}}$$



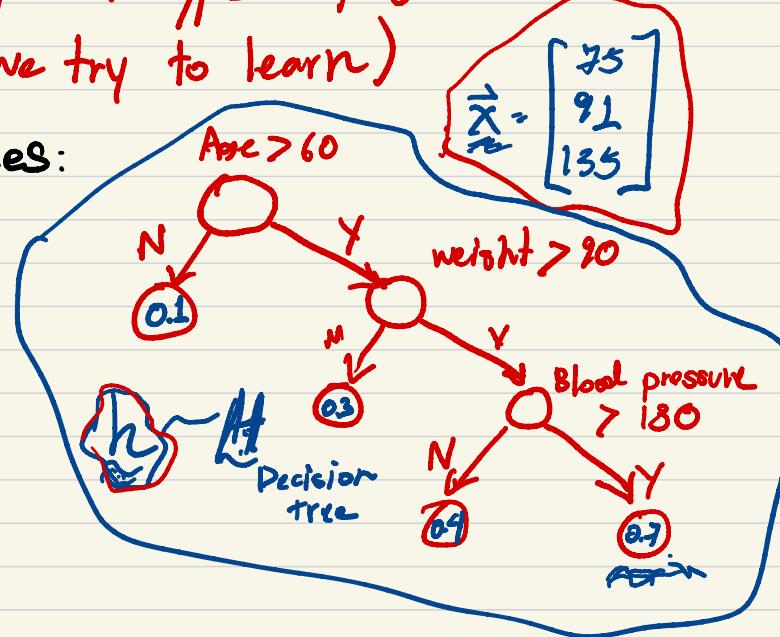
Goal: To learn function  $h$  from  $X$  to  $Y$ .



Automated job for  
defining  $H$ .  
(Unnecessary)

- Types of hypothesis classes:

- Decision tree
- Neural networks
- Linear classifier
- Etc.



## Learning the hypothesis

① Defining  $H$ : Choosing appropriate ML algorithm.

The set  $H$  defines the set of functions we could learn. (non-automated job)

② Learning the best  $h \in H$ :

The ML algorithm will find the best hypothesis  $h$  within the class (automated job)

Ideally, we wish to find  $\underset{h \in H}{\text{that makes the fewest mistakes within our data}}!!!$

## Loss Functions:

$\rightarrow L(h, D)$  bad? Good?

- Loss functions are the functions that evaluate if one function is better than another.
- We use loss functions to evaluate how bad any  $h \in H$  is on the data
  - The higher loss  $\rightarrow$  the worst hypothesis
  - Normally, we normalize the loss by the number of examples (average loss)

$h(x) = y$   
Zero  
 $\forall (x, y) \in D$

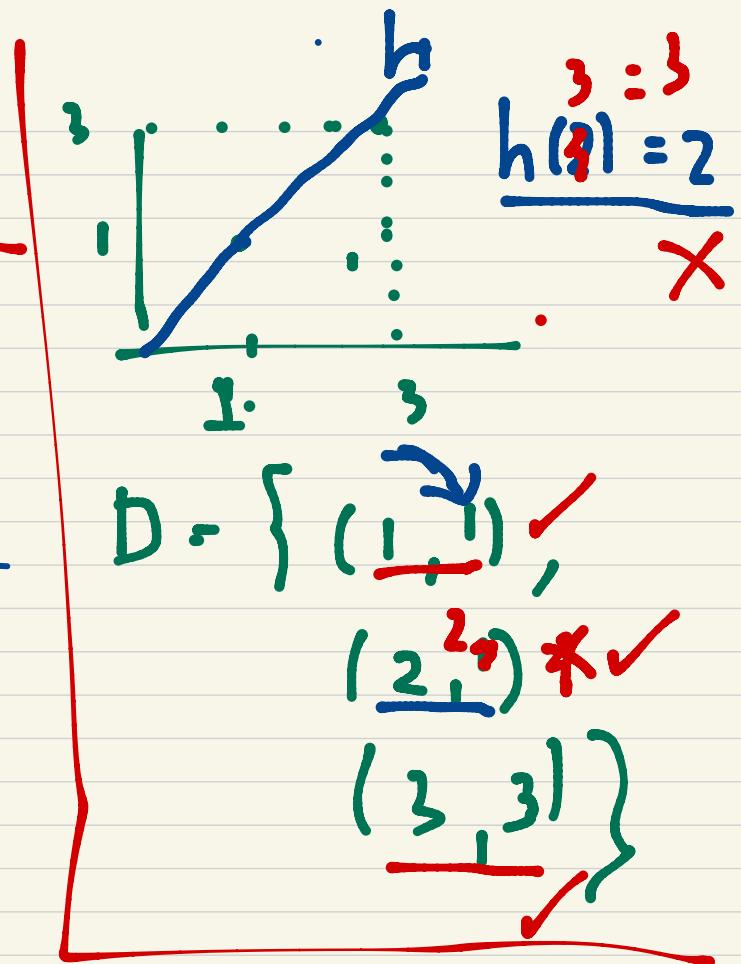
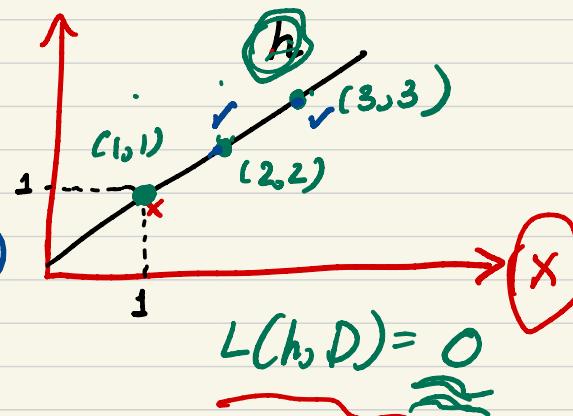
## Common Loss Functions:

0/1-loss:  $L(h, D) = \frac{1}{|D|} \sum_{\forall (\bar{x}, y) \in D} d(h(\bar{x}), y)$

$$d(h(\bar{x}), y) = \begin{cases} 0, & \text{if } h(\bar{x}) = y \\ 1, & \text{if } h(\bar{x}) \neq y \end{cases}$$

$$h(\vec{x}) = x$$

$$L(h, D) = \frac{1}{3} (1+0+0) = 0.33$$



$$L(h, D) = 0 \rightarrow 0$$