

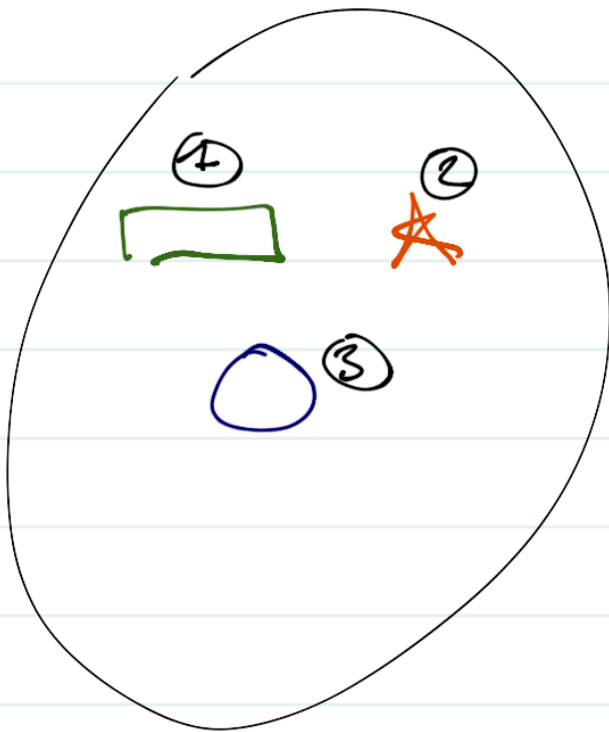
Lecture 1

TRAINING SET | datapoints $x_1, \dots, x_n \in \mathbb{R}^d$
 $y_1, \dots, y_n \in \{-1, +1\}$

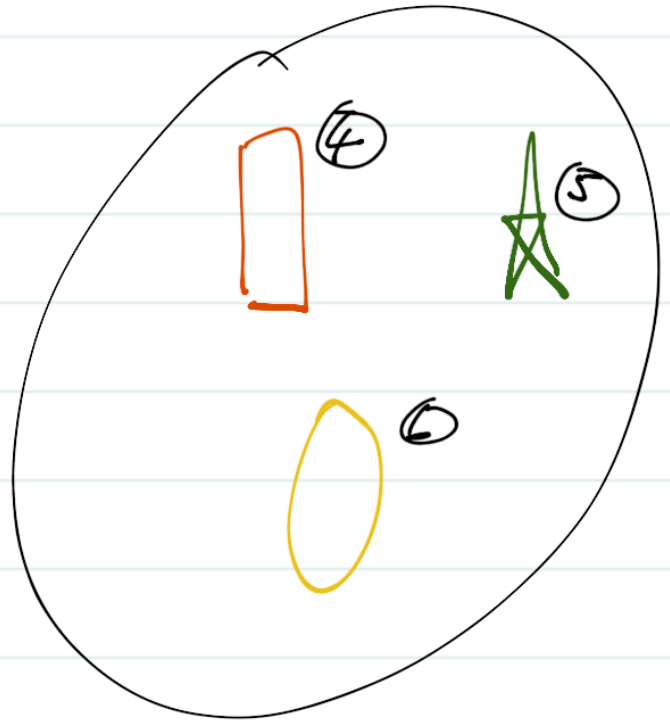
UNSEEN DATA | new datapoint $x \rightarrow$ what is y ?

Learn the dependency

$$x \longrightarrow y$$



+ 1



- 1

①

②

③

④

⑤

⑥

Color

green

red

blue

red

green

yellow

Height

small

small

small

large

large

large

Width

large

small

small

small

small

small

$$\textcircled{1} \quad x_1 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \in \mathbb{R}^3$$

green	0	height	small	0
red	1		large	1
blue	2			
yellow	3	width	small	0
			large	1

$$x_1, \dots, x_n \in \mathbb{R}$$

$$n=6 \quad d=3$$

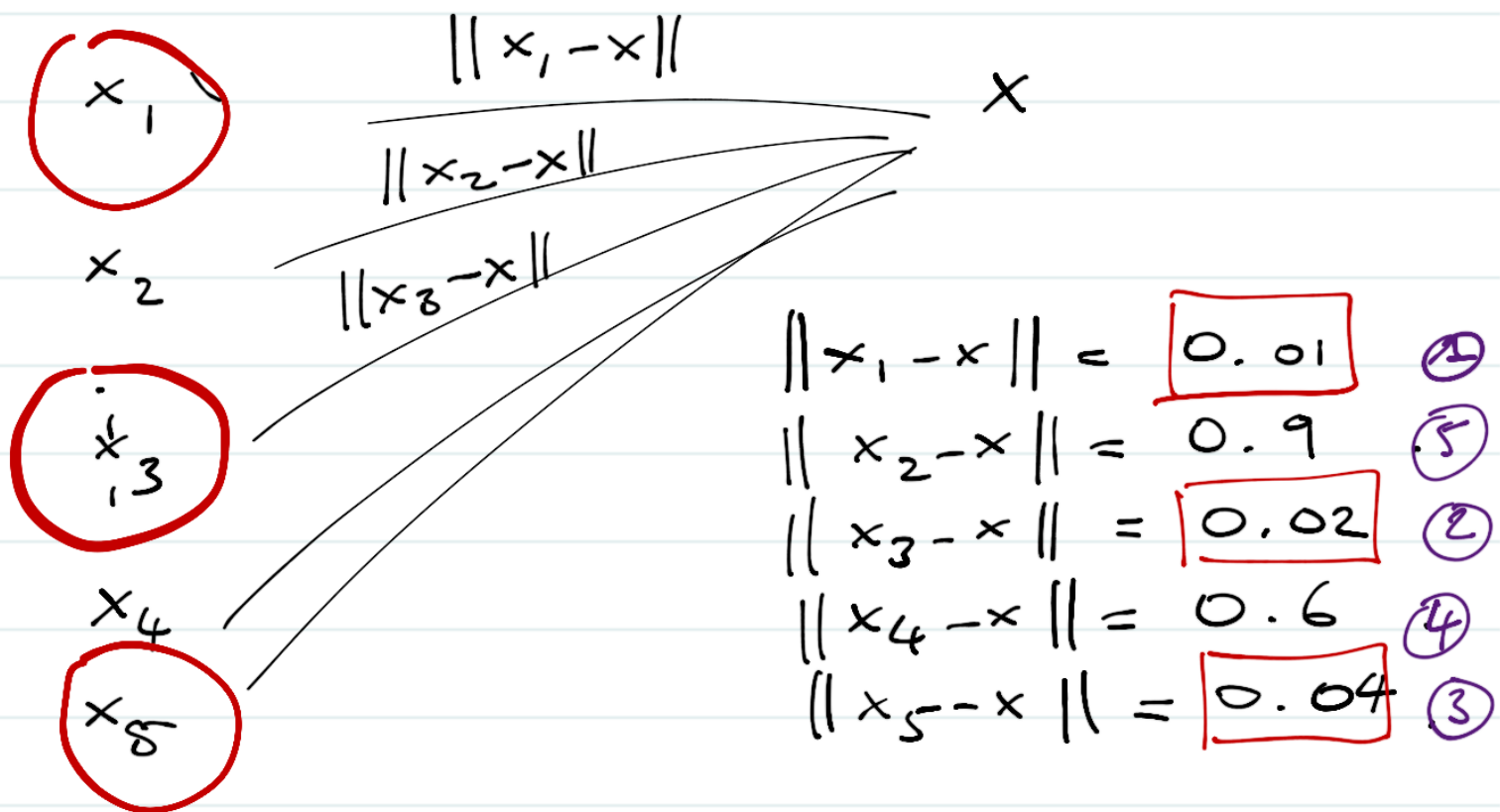
$$P(y = \pm 1 \mid x, \mathcal{D})$$

\uparrow new datapoint \uparrow dataset $\{(x_1, y_1), \dots, (x_n, y_n)\}$

$$= \frac{1}{K} \sum_{i \in \text{NN}_K(x, \mathcal{D})} \mathbb{1}(y_i = \pm 1)$$

$\text{NN}_K(x, \mathcal{D})$ the set of
nearest-neighbors of x in \mathcal{D}

$\|x - x_i\|_2$ distance between
 x and x_i



$$\underline{k=3} \quad \mathcal{N}_3(x, \emptyset) = \{1, 3, 5\}$$

$$P(y = +1 \mid x, \emptyset)$$

$$= \frac{1}{3} \sum_{i \in \mathcal{N}_3(x, \emptyset)} \mathbb{1}(y = +1)$$

\uparrow
 $\{1, 3, 5\}$

$$x_1 \quad y_1 = +1$$

$$x_3 \quad y_3 = +1$$

$$x_5 \quad y_5 = -1$$

$$p(y = +1 | x, \infty)$$

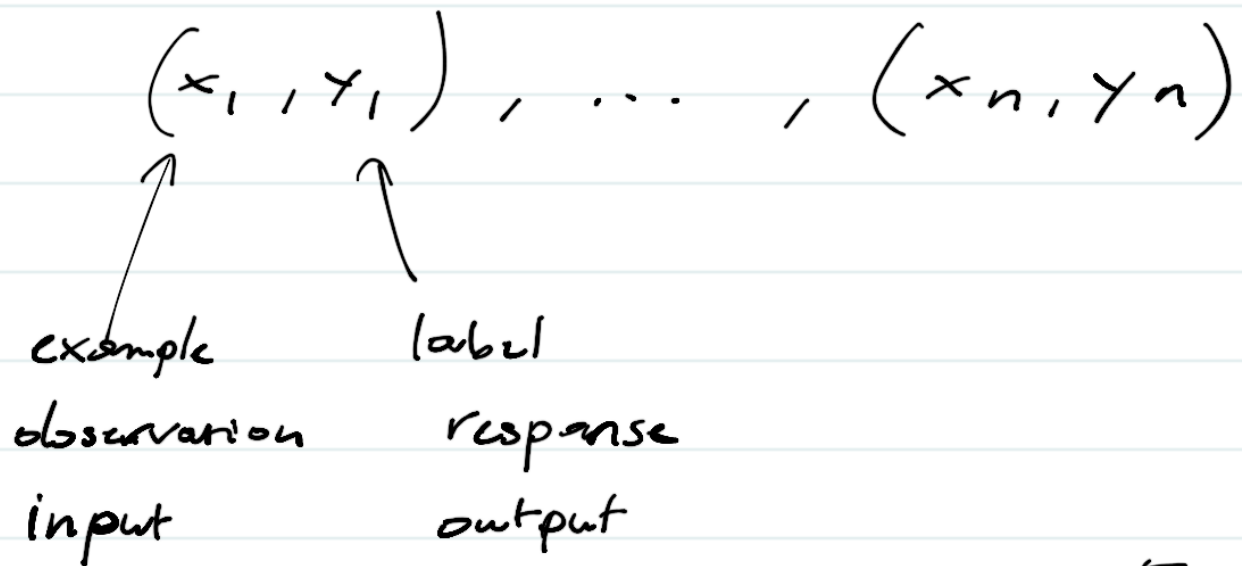
$$= \frac{1}{3} \sum_{i \in \{1, 3, 5\}} \mathbb{1}(y_i = +1)$$

$$= \frac{1}{3} \times \{1 + 1 + 0\}$$

$$= \underline{\underline{\frac{2}{3}}}$$

Binary supervised classification

TRAINING SET



$$x_i \in \mathbb{R}^d \quad x_i = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} \begin{matrix} \uparrow \\ \downarrow \end{matrix} d$$

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

1 - Nearest - neighbor

$n \longrightarrow +\infty$ large sample
training set

$d \ll n$ asymptotic setting

error of 1-NN
number of mistakes

is

2x best error possible

Error = 1%

1-NN error = 2%

15%

30%

• \textcircled{K} wrt n

$$\frac{K}{n} \xrightarrow{n \rightarrow +\infty} 0$$

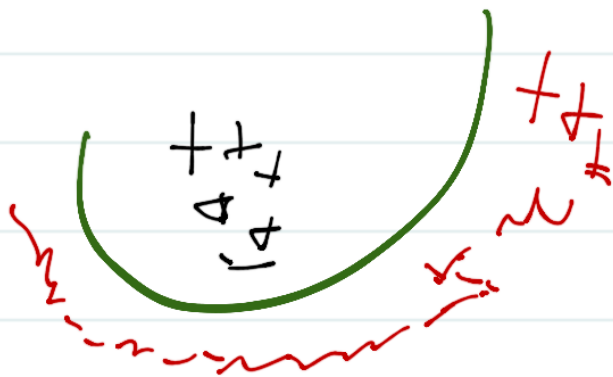
$$K \ll n$$

$$K = \frac{\sqrt{n}}{2}$$

• K -NN with $K = o(n)$

as $n \rightarrow +\infty$

K -NN achieves
the best possible error



Nearest-Neighbors $n \gg 1$
 $d \gg 1$

Approximate NN search

$$x \in \mathbb{R}^d \quad x_i \in \mathbb{R}^d$$

$$\cdot \quad \|x - x_i\|_2 \quad \left. \begin{array}{l} O(d) \text{ time} \\ O(1) \text{ space} \end{array} \right\}$$

$$\cdot \quad \|x - x_i\|_2 \quad \text{for all } i = \{1, \dots, n\}$$

$$O(nd)$$

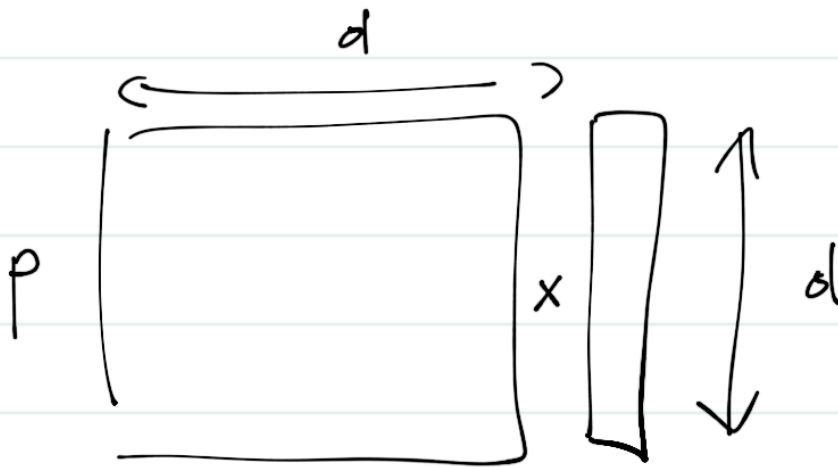
$$\tilde{d}_i = \|x - x_i\|_2$$

Hashing techniques

$$x \in \mathbb{R}^d \quad Ax$$

$$A_{ij} \stackrel{\text{iid}}{\sim} \text{er}(0, 1) \quad \forall i, j \in \{1, \dots, p\}$$

$$x \in \{1, \dots, d\}$$



$$p \ll d$$

$$\tilde{x} \in \mathbb{R}^p$$

$$d = 10^6$$

$$p = 10^2$$

Issues with K-NN :

- choice of distance
- choice of k

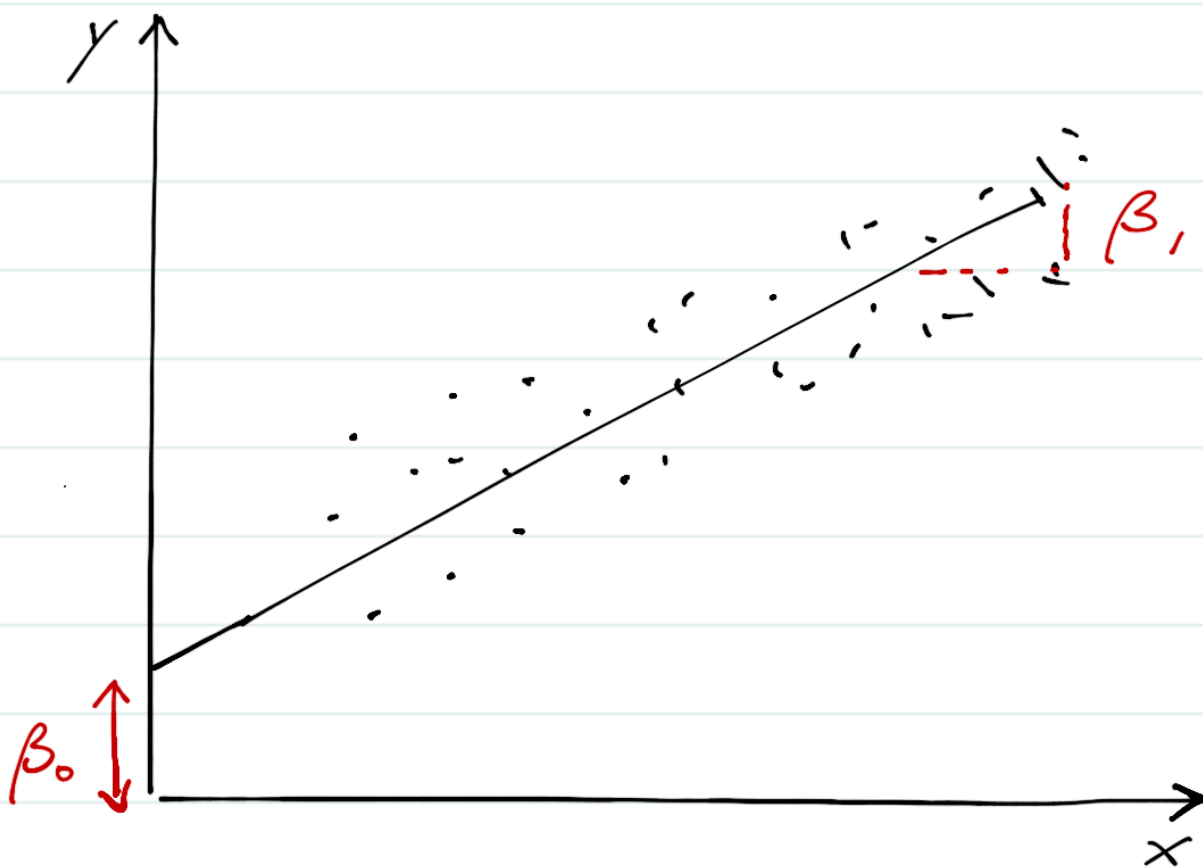
Linear models

$$y = x^T \beta + \beta_0$$

$$x \in \mathbb{R}^{d-1}$$

$$\begin{pmatrix} \beta_1 \\ \vdots \\ \beta_{d-1} \end{pmatrix} \in \mathbb{R}^{d-1}$$

$$\begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{d-1} \end{pmatrix} \in \mathbb{R}^d$$



$$\begin{aligned} \text{Min}_{\beta, \beta_0} \quad & \frac{1}{n} \sum_{i=1}^n \left(y_i - (x_i^T \beta + \beta_0) \right)^2 \\ & + \boxed{\frac{\lambda}{2} \|\beta\|_2^2} \end{aligned}$$

$\lambda = 0$ Ordinary Least-Squares

$\lambda > 0$ L_2 -regularized Least-Squares

$$y = x^T \beta + \beta_0 \leftarrow \begin{array}{l} \text{offset} \\ \text{intercept} \end{array}$$

$$x^T \beta = \begin{bmatrix} x^1 & x^2 & x^3 & x^4 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_4 \end{bmatrix}$$

$$= \sum_{j=1}^4 \beta_j x^j$$

$$\text{Min}_{\beta, \beta_0}$$

$$\frac{1}{n} \sum_{i=1}^n \left(y_i - x_i^T \beta \right)^2$$

$$+ \frac{\lambda}{2} \sum_{j=1}^4 \beta_j^2 \quad \text{REGULARIZATION}$$

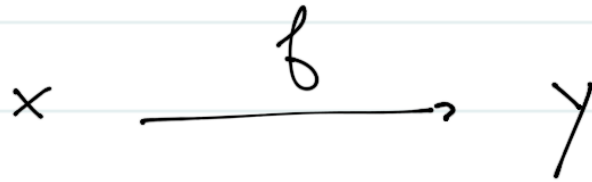
• $\|\beta\|_2^2$ L_2 reg^o penalty

• $\|\beta\|_1$ L_1 reg penalty

→ Sparsity
select features

LASSO Min $\frac{1}{n} \sum_{i=1}^n (y_i - (x_i^T \beta + \beta_0))^2$
 $+ \lambda \|\beta\|_1$

Model evaluation



Empirical risk

$$\frac{1}{n} \sum_{i \in \mathcal{D}_{\text{TRAIN}}} \ell(y_i, f(x_i))$$

↑
sum over training set

Training error

$$\frac{1}{n} \sum_{i \in \mathcal{D}_{\text{TRAIN}}} \text{Err}(y_i, f(x_i))$$

Test error

$$\frac{1}{m} \sum_{i \in \mathcal{D}_{\text{TEST}}} \text{Err}(y_i, f(x_i))$$