

# Classification & Regression (k-NN)

Ques How classification works? { Amazon food review }  
 → +ve/-ve

$x_i \rightarrow$  Text → Vector  
 (360K+)  
 [ Bow / tfidf / w2v ]

meal would eg  
 MNIST (PCA / tSNE)

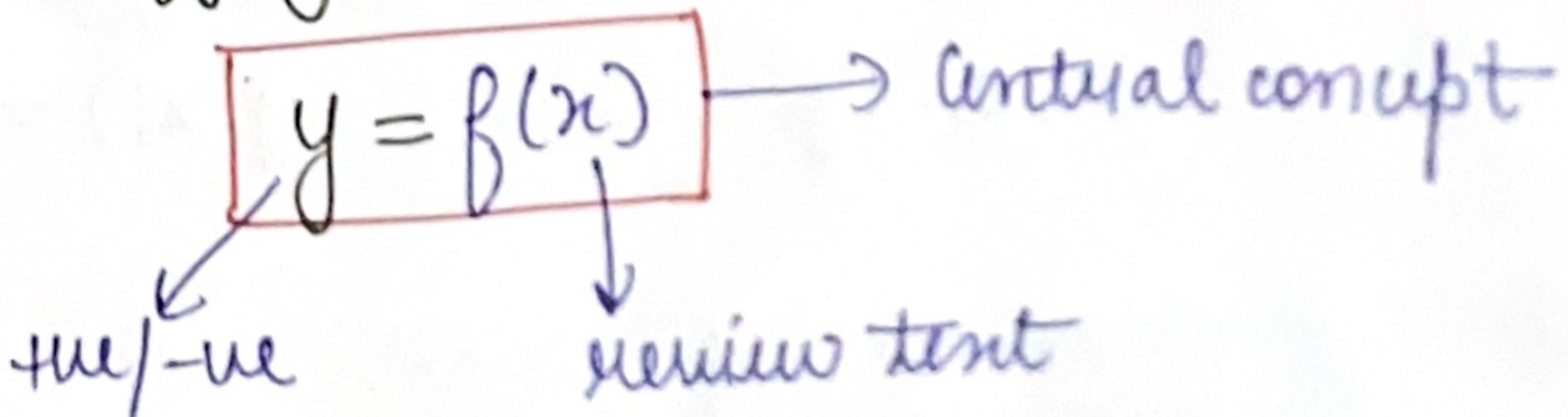
## classification

364K+  $x_i \rightarrow (x_i^{(1)}, +ve/-ve)$

work { Given a new review determine / predict if review is -ve/+ve.

$x_q \rightarrow$  +ve/-ve

## classifying

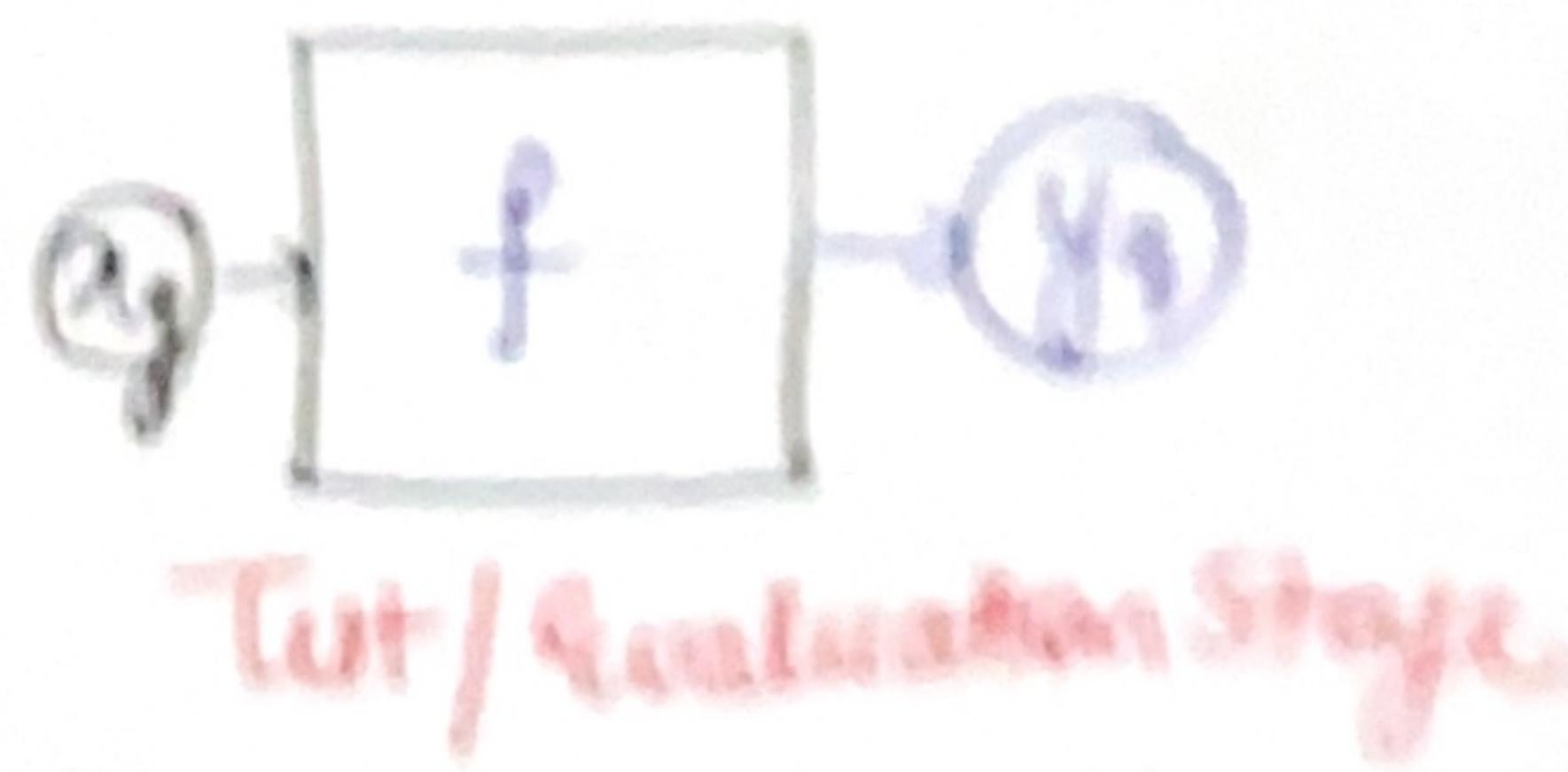
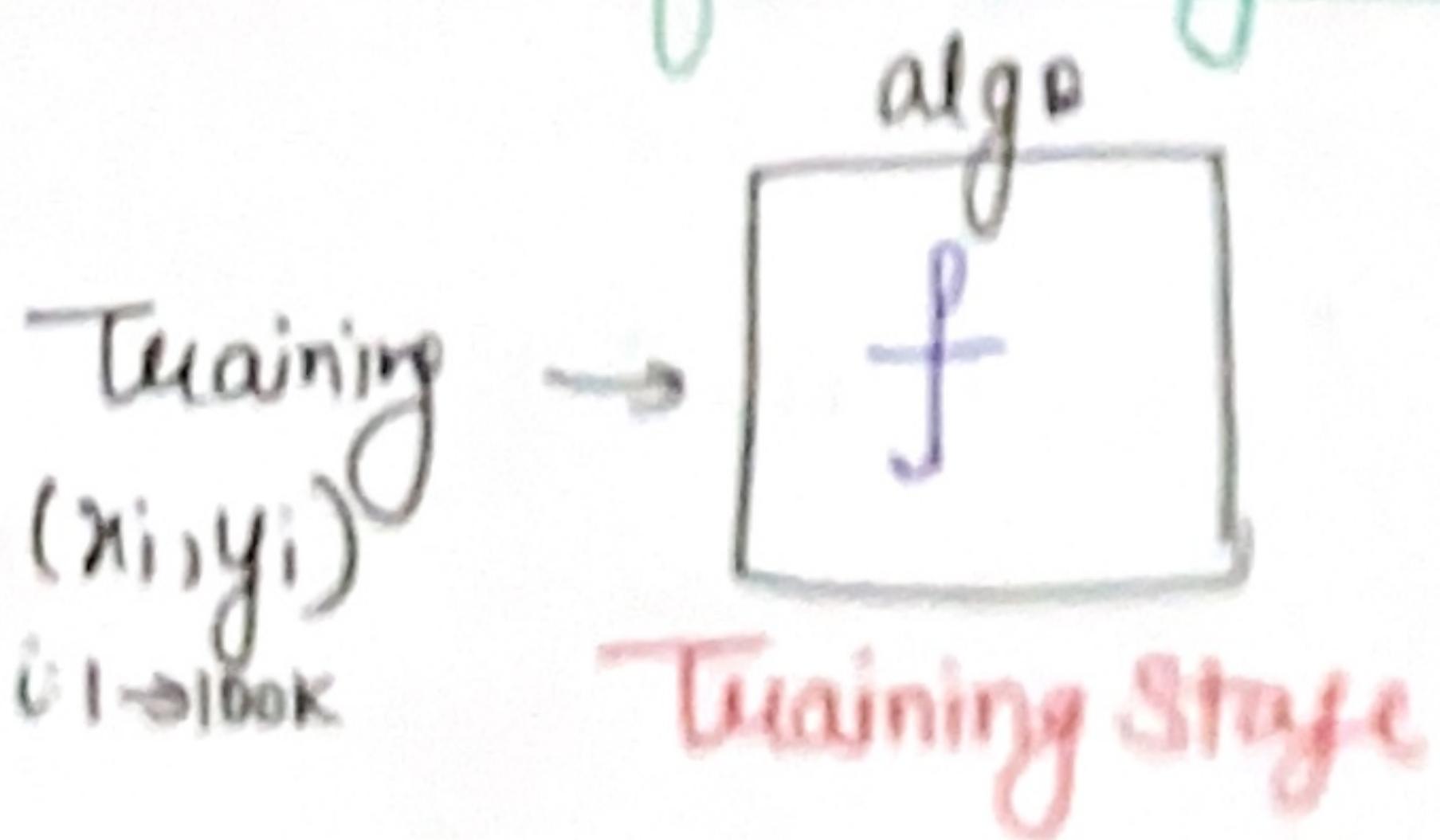


$$y_q = f(x_q)$$

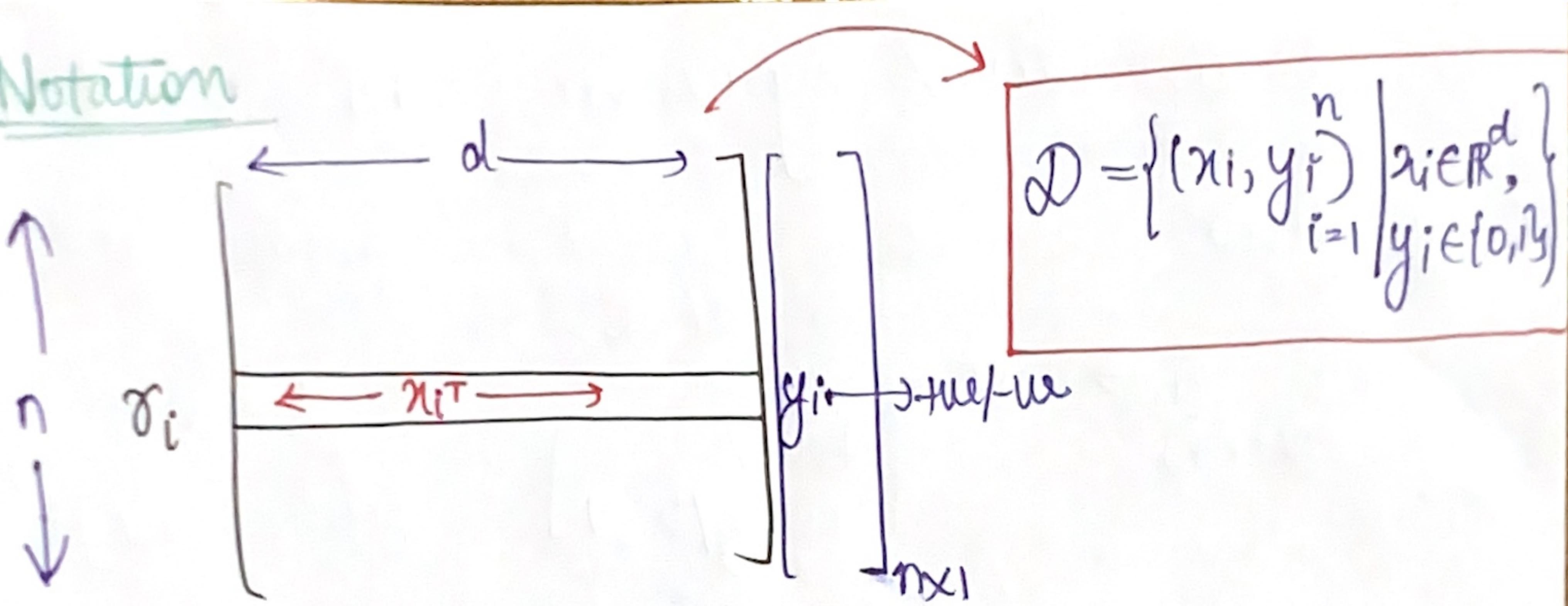
↑  
Query review

+ve/-ve

## Classification algorithm work



## Notation



$x_i \rightarrow \text{vector}(x_i)$   
 $\hookrightarrow \text{column vector}$   
 $\rightarrow \text{class(+ve/-ve)}$   
 $y = f(x)$

$$\mathcal{D}_n = \left\{ (x_i, y_i)_{i=1}^n \mid x_i \in \mathbb{R}^d, y_i \in \{0, 1\} \right\}$$

$\downarrow$  such that  
 $\swarrow$  -ve       $\searrow$  +ve

$x_i \rightarrow x_i$   
 $\hookrightarrow \{0, 1\} y_i$

classification:  
 $f(x_i) = y_i$

## Classification vs Regression

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \mid x_i \in \mathbb{R}^d, y_i \in \{0, 1\}\}$$

$y_i \in \{0, 1\}$  ← Amazon food review  
 $\downarrow$  -ve       $\downarrow$  +ve      ↗ 2 classes  
 ↗ 2 class - classification (binary)

## MNIST:

$$y_i \in \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

↗ 10-class / multi-class classification

What if  $y_i \in R \rightarrow \text{Regression}$   
 $\hookrightarrow y_i$  is no more part of small finite set of classes

$i=1 \rightarrow 10K$

$x_i = (\text{weight}, \text{age}, \text{gender}, \text{place})$   
 $\hookrightarrow y_i = \text{height}$   $f(x_i) = y_i$   
 $\hookrightarrow \text{real no}$   
 $(182.6\text{cm}, 152.7\text{cm})$

$y_i \in \{0, 1\} \rightarrow \text{Classification}$   
 $y_i \in R \rightarrow \text{Regression}$

K-Nearest Neighbors (Geometric)

2D - toy dataset

(binary classification)



$x$ : the data pt  
 $\alpha$ : -ve class data pt

$$\mathcal{D} = \{(x_i, y_i) | i \in \mathbb{R}, y_i \in \{0, 1\}\}$$

geom close to  $x_q$ :

conclude  $x_q \rightarrow \text{blue}$   
 $(+ve)$

given  
 $x_q \rightarrow y_q \rightarrow \text{final } \{0, 1\}$

① Find k-nearest points to  $x_q$  in  $\mathcal{D}$   
Let  $k=3$   $[x_1, x_2, x_3]$  - 3 neighbours to  $x_q$

②  $\underline{k=\text{odd}}$   $\{y_1, y_2, y_3\} \rightarrow$  majority vote / note

$y_1, y_2, y_3 \rightarrow$  Majority vote

$+, +, + \rightarrow \oplus y_q = +ve (\text{blue})$

$+, +, - \rightarrow y_q = +ve$

if  $k=4$

$+, +, -, - \rightarrow \oplus \text{ or } \ominus$

$x_q \rightarrow k-\text{NN of } x_q$

$x_1, x_2, \dots, x_k$

$\downarrow$

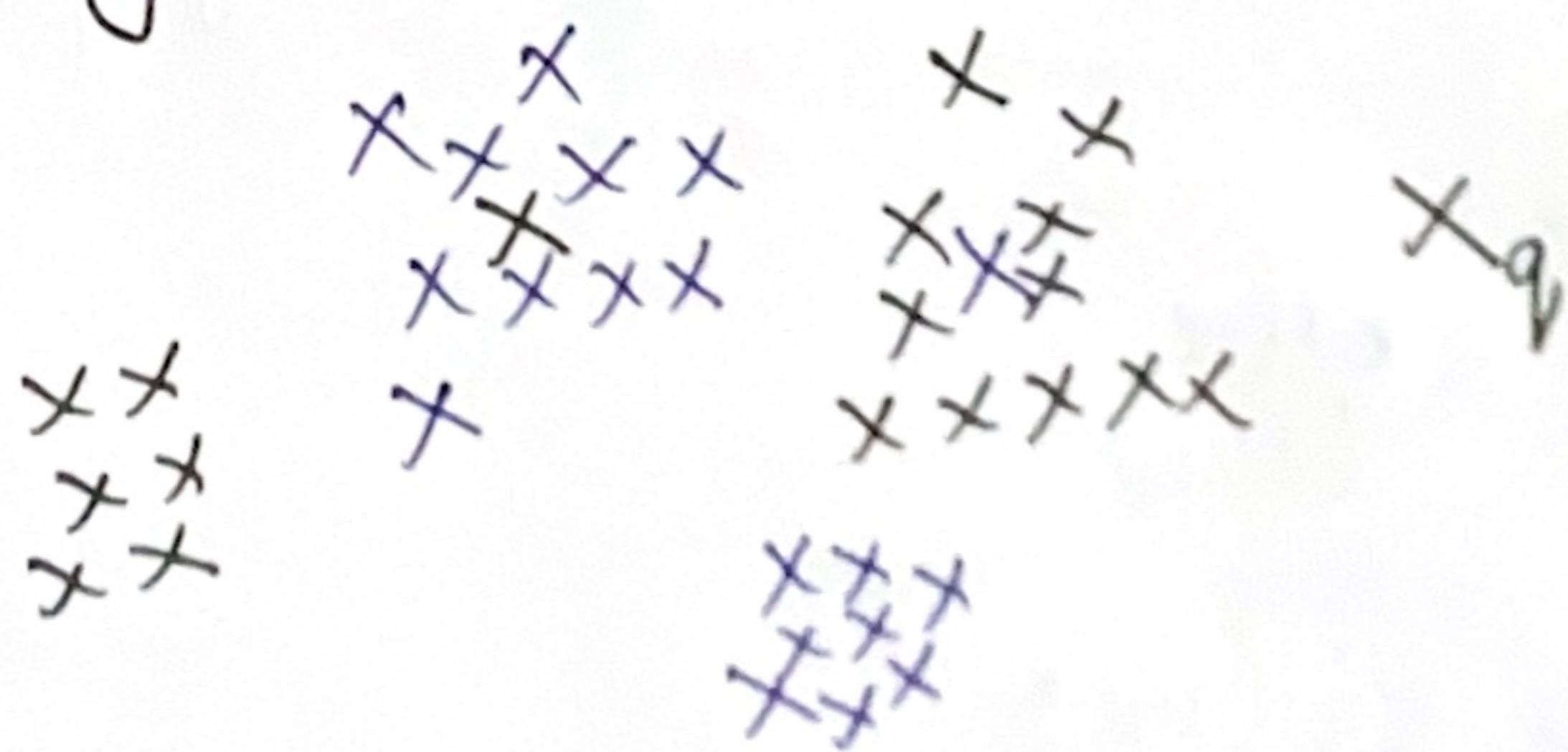
$y_1, y_2, \dots, y_k$

Majority vote

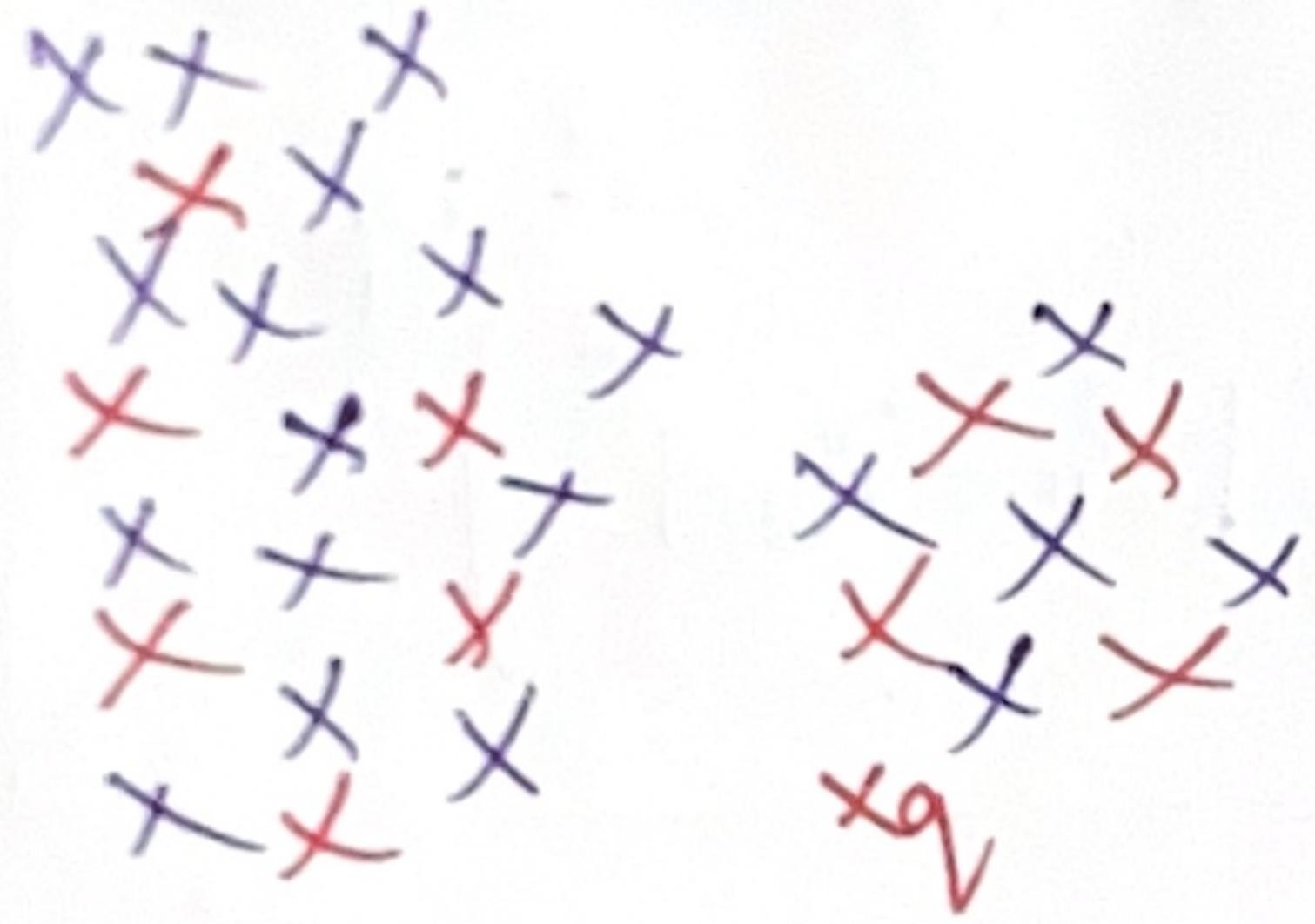
$\textcircled{y}_q$

Failure of K-Nearest Neighbour (geometric)

2D toy dataset

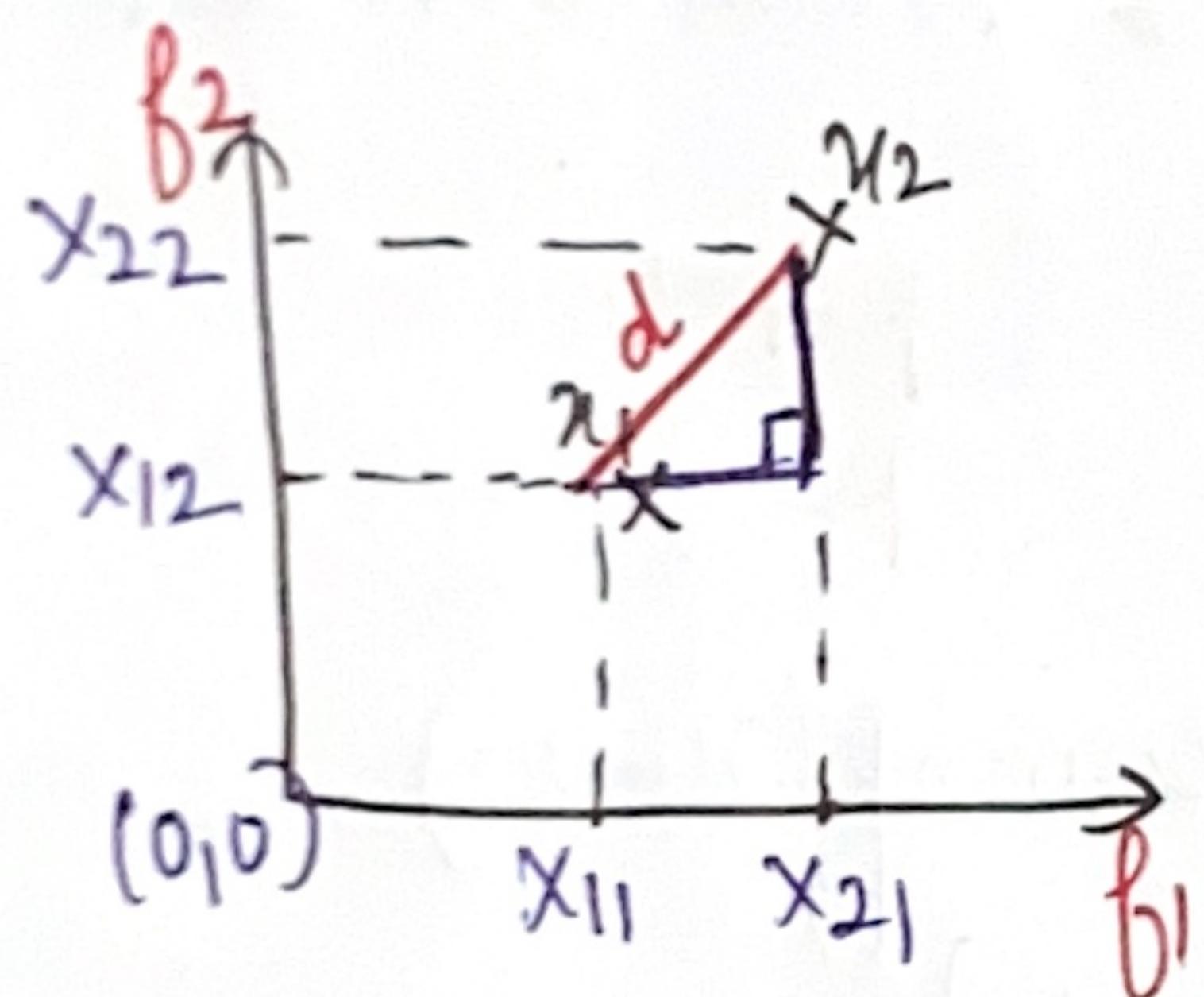


$x_q$  is very far away from other points



Data is very jumbled  
KNN algorithm makes  
↓  
does not

### Distance Measure



$$x_1 = (x_{11}, x_{12})$$

$$x_2 = (x_{21}, x_{22})$$

$d = \text{length of shortest line from } x_1 \text{ to } x_2$

Euclidean distance  $\leftarrow d = \sqrt{(x_{21} - x_{11})^2 + (x_{22} - x_{12})^2} = \|\mathbf{x}_1 - \mathbf{x}_2\|$

$$\mathbf{x}_i \in \mathbb{R}^d, \mathbf{x}_2 \in \mathbb{R}^d$$

Eucl. dist.  $\|\mathbf{x}_1 - \mathbf{x}_2\|_2 = \left( \sum_{i=1}^d (x_{1i} - x_{2i})^2 \right)^{1/2}$

$\|\mathbf{x}_1 - \mathbf{x}_2\|_2 \rightarrow L_2 \text{ norm of vector}$

$\|\mathbf{x}_1\|_2 = \text{dist of } \mathbf{x}_1 \text{ from origin} = \left( \sum_{i=1}^d x_{1i}^2 \right)^{1/2}$

### Manhattan dist.

$L_1 \text{ norm of vector } \|\mathbf{x}_1 - \mathbf{x}_2\|_1 = \sum_{i=1}^d |x_{1i} - x_{2i}|$

absolute

$$\|\mathbf{x}_1\|_1 = \sum_{i=1}^d |x_{1i}|$$

$L_p$  norms  $\rightarrow$  Minkowski distance

$$\|x_1 - x_2\|_p = \left( \sum_{i=1}^d |x_{1i} - x_{2i}|^p \right)^{1/p}$$

$p=2 \rightarrow$  Minkowski dist  $\rightarrow$  Euclidean dist

$p=1 \rightarrow$  "  $\rightarrow$  Manhattan dist

✓  $\|x_i\|_p = \left( \sum_{i=1}^d |x_{ii}|^p \right)^{1/p}$

$p \neq 0$   
 $p > 0$

$L_p$  norm

$$\text{Eucl dist}(x_1, x_2) = L_2 \text{ norm of } (x_1 - x_2)$$
$$= \|x_1 - x_2\|$$

dust( $x_1, x_2$ ) [b/w two points]

norm  $\rightarrow$  vector [b/w k vectors]

Hamming dist (boolean vector)

$x_1, x_2 \rightarrow$  boolean vector  $\rightarrow$  Binary B/w

$$x_1 = [0, 1, 1, 0, 1, 0, 0, \dots]$$

$$x_2 = [1, 0, 1, 0, 1, 0, 1, \dots]$$

Hamming dist( $x_1, x_2$ ) = # locations / dimensions  
where binary vectors differ

↳ 3

## Strings

$\pi_1 = abc\underset{|}{a}def\underset{|}{g}hik \leftarrow \text{gene code}$   
 $\pi_2 = ac\underset{|}{b}ade\underset{|}{g}f\underset{|}{hik}$  eg AGTC

$$\text{hamming dist } (\pi_1, \pi_2) = 4$$

## Cosine similarity & cosine distance

$\pi_1, \pi_2$

Similarity  
↓  
↑

distance  
↑ inc  
↓ dec

(opposite)

v. similar

$$\text{cos sim}(\pi_1, \pi_2) = +1$$

$$1 - \text{Cos Sim}(\pi_1, \pi_2) = \text{Cos dist}(\pi_1, \pi_2)$$

v. diff

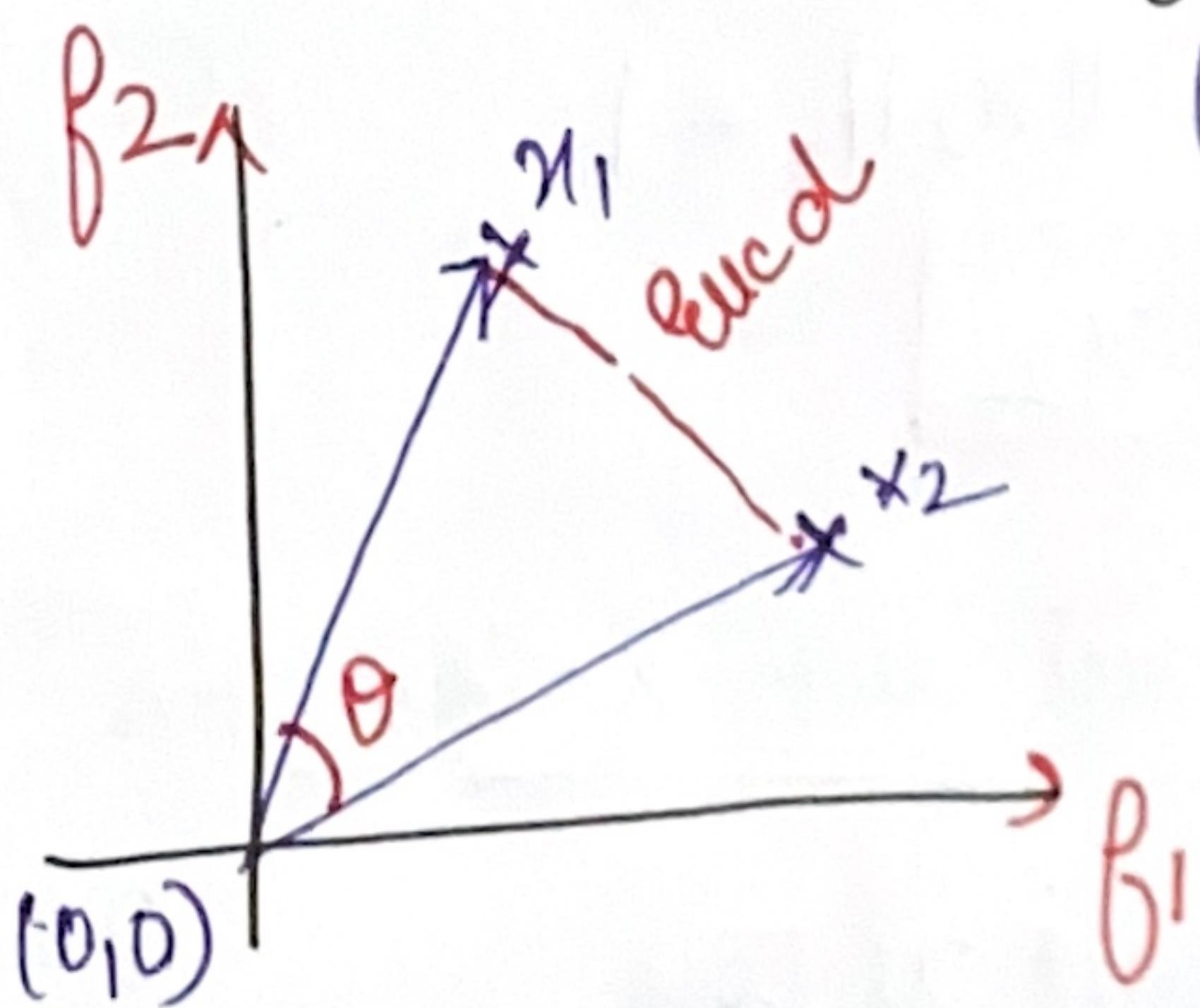
$$\text{cos sim}(\pi_1, \pi_2) = -1$$

$$\text{det } [-1, +1]$$

$$\text{Cos Sim} = \cos \theta$$

$(\pi_1, \pi_2)$

$\theta$ : angle b/w  $\pi_1$  &  $\pi_2$



$$\text{cos dist} = 1 - \cos \theta$$

$(\pi_1, \pi_2)$  Euclidist

$$d_{13} > d_{12}$$

$$\text{cos dist}_{13} < \text{cos dist}_{12}$$

$$\text{cos dist} = 1 - 1 = 0$$

$$\text{Cos Sim}(\pi_1, \pi_2) = \cos 0^\circ$$

$$\text{Cos Sim}(\pi_1, \pi_3) = 1$$

$$\theta_{\pi_1, \pi_3} = 0^\circ$$

$$\cos 0^\circ = 1$$

$\pi_1$

$\pi_2$

$\pi_3$

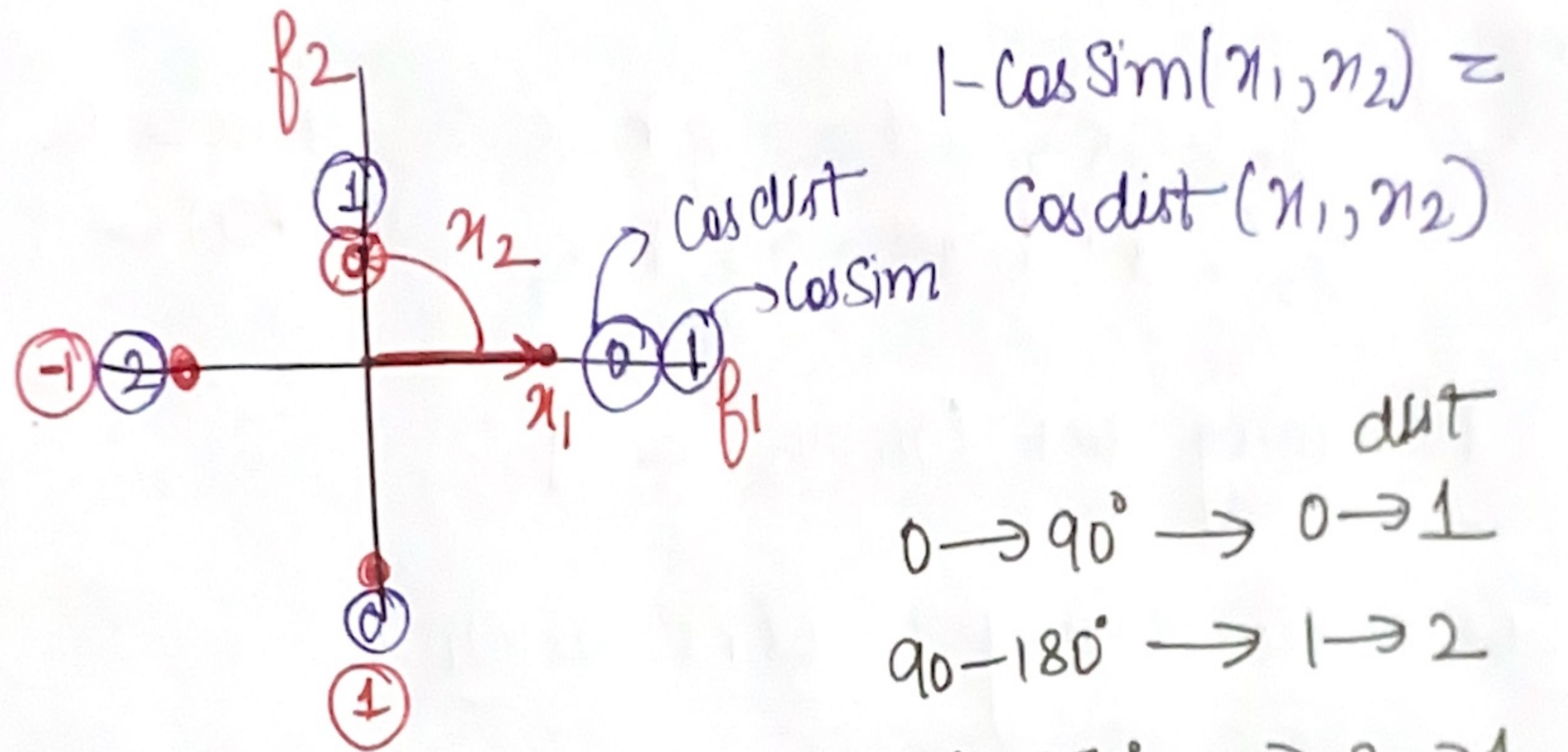
$$\text{cos dist} = 1 - 1 = 0$$

$$\text{Cos Sim}(\pi_1, \pi_2) = \cos 0^\circ$$

$$\text{Cos Sim}(\pi_1, \pi_3) = 1$$

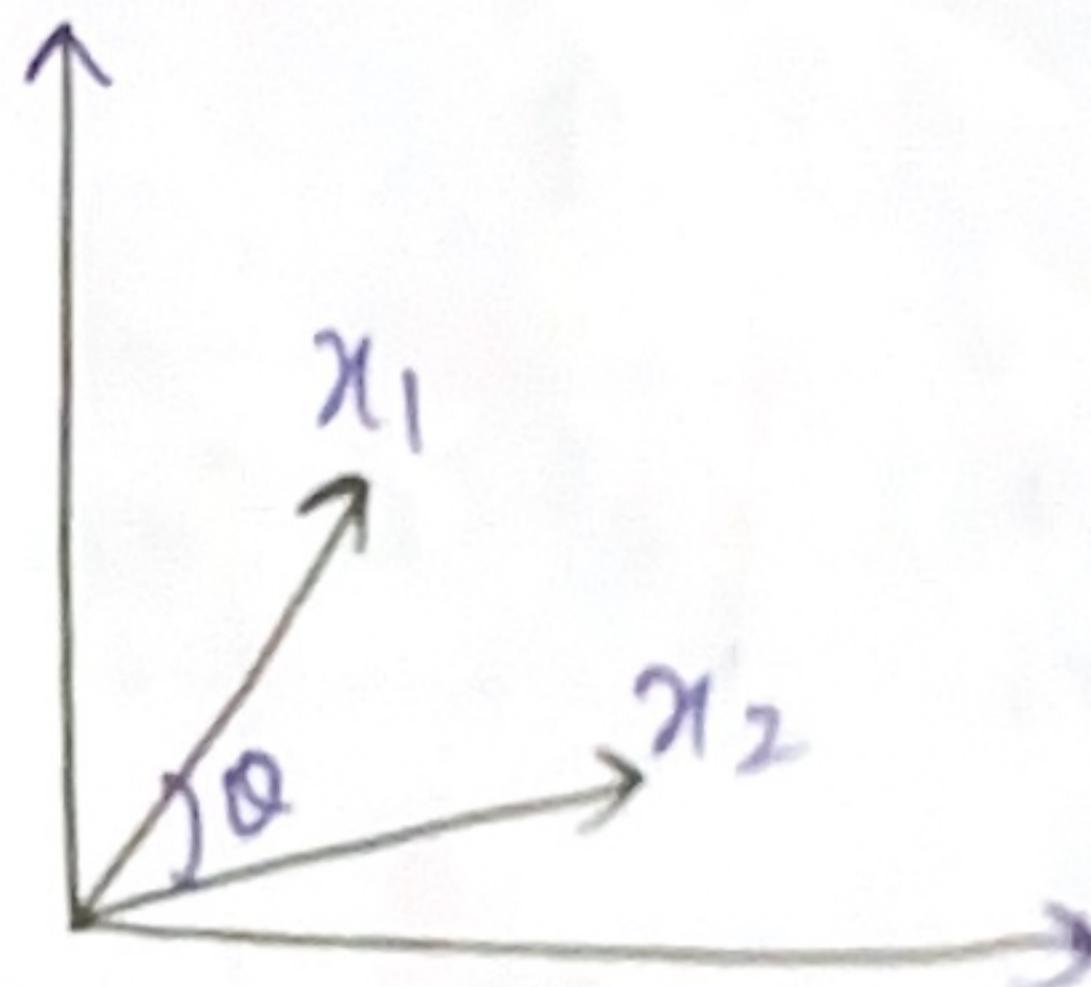
$$\theta_{\pi_1, \pi_3} = 0^\circ$$

$$\cos 0^\circ = 1$$



$$\cos(\theta) = \frac{\underline{x_1} \cdot \underline{x_2}}{\|\underline{x_1}\|_2 \|\underline{x_2}\|_2}$$

↳ l<sub>2</sub> norm



① If  $x_1$  &  $x_2$  are unit vectors

$$\|\underline{x_1}\|_2 = \|\underline{x_2}\|_2 = 1$$

$$\boxed{\cos(\theta) = \underline{x_1} \cdot \underline{x_2}}$$

Relationship b/w Eucl. dist & Cos-Sim

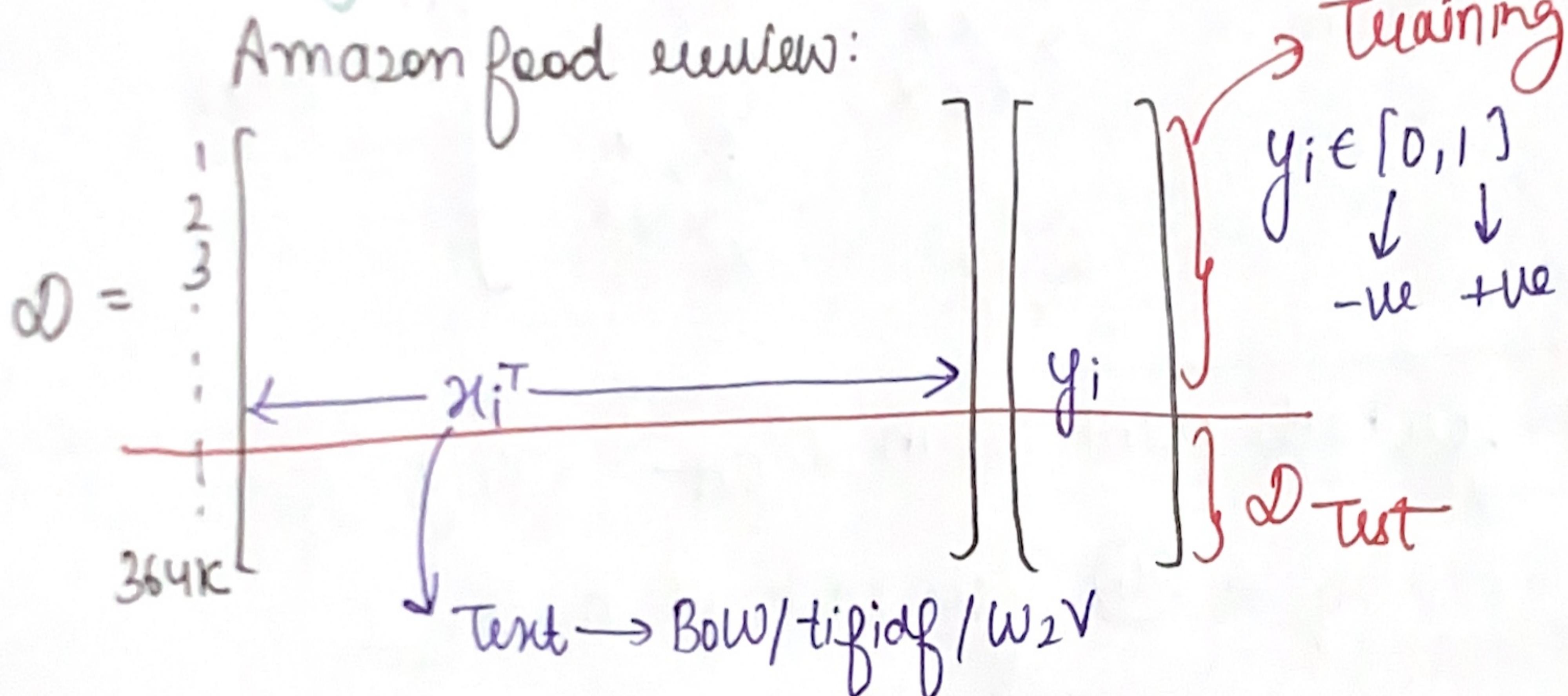
if  $x_1$  &  $x_2$  are unit vectors

$\theta = \text{angle b/w } x_1 \text{ & } x_2$

$$[\text{euc dist}(x_1, x_2)]^2 = 2(1 - \cos(\theta))$$

$$\hookrightarrow \boxed{2 \cos \text{dist}(x_1, x_2)}$$

# Measuring how good K-NN is



Problem: Given a new food review, what is its polarity (true/-ve)

$$x_q \rightarrow y_q$$

$$x_q \rightarrow \text{Text} \rightarrow x_q$$

"Measure" k-NN  $\rightarrow$  (k-NN + Majority vote)

$$D_n \xrightarrow{x_i} D_{train} \cup D_{test} = D_n$$

$$D_{train} \cap D_{test} = \emptyset$$

$$D_{train} \rightarrow n_1$$

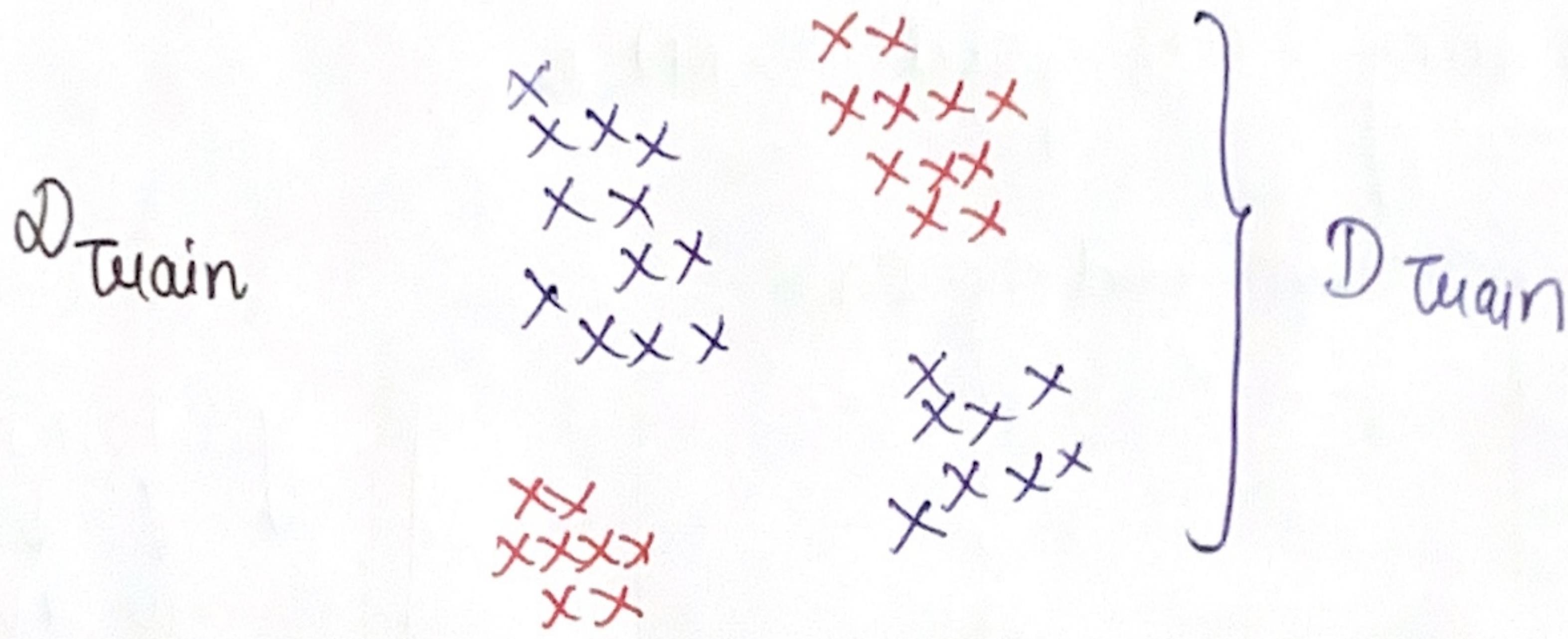
$$n_1 + n_2 = n$$

$$D_{test} \rightarrow n_2$$

$$\begin{aligned} \text{Split } D_n &\xrightarrow{\text{70\%}} D_{train} \\ &\xrightarrow{\text{30\%}} D_{test} \end{aligned} \quad \xrightarrow{\text{randomly}}$$

①  $D_{train} \rightarrow \text{k-NN}$   
 $(x_i, y_i)_{i=1}^{n_1}$

②  $D_{test} \rightarrow (x_i, y_i)_{i=1}^{n_2}$



for each point  $x_i$  in  $D_{\text{test}} = (x_i, y_i)_{i=1}^{n_2}$

$x_q = \textcircled{x}_i$  for each pt in  $D_{\text{test}}$

$\rightarrow x_q = \text{pt}$

$\rightarrow$  use  $D_{\text{train}}$  & KNN to predict  $y_q$

if  $y_q = y_{\text{pt}}$

$\text{cnt} + 1$

end

$\text{cnt} = \# \text{ pts}$  for which  $D_{\text{train}}$  + KNN gave a correct class label

$$\boxed{\text{Accuracy} = \frac{\text{cnt}}{n_2}}$$

$\checkmark$   
 $\# \text{ pts in } D_{\text{test}}$

$\#$  for which  $D_{\text{train}}$  + KNN gave a correct class label

$$0 \leq \text{Accuracy} \leq 1$$

$\text{Acc} = 0.91 \Rightarrow 91\% \text{ of times}$

$x_q \rightarrow y_q$

Conclude KNN on Amazon food review using  
 $D_{\text{train}}$  gives me an accuracy of 91%  
assumption

### Time / Evaluation time & Space Complexity

$$x_q \rightarrow y_q$$

$K$  is small  
5 or 10

Input:  $D_{\text{train}}, K, x_q \in \mathbb{R}^d$ ; output:  $y_q$

$$\text{KNN pts} = [ ]$$

$\rightarrow$  npts; d-dim  
 $\hookrightarrow \text{BOW}(10K)$

[ for each  $x_i$  in  $D_{\text{train}}$ :

$O(d)$  — compute  $d(x_i, x_q) \rightarrow d_i$

— keep the smallest  $K$ -distances  $(x_i, y_i, d_i)$

$$\text{cnt\_pos} = 0 \quad \text{cnt\_neg} = 0$$

for each  $x_i$  in KNN pts

if  $y_i$  is +ve

$$\text{cnt\_pos} += 1$$

else

$$\text{cnt\_neg} += 1$$

$O(1)$  } if  $\text{cnt\_pos} > \text{cnt\_neg}$   
return  $y_q = 1 \rightarrow +ve$   
else  
 $y_q = 0 \rightarrow -ve$

Time complexity =  $O(nd) + O(k) + O(k)$   
 $O(nd)$

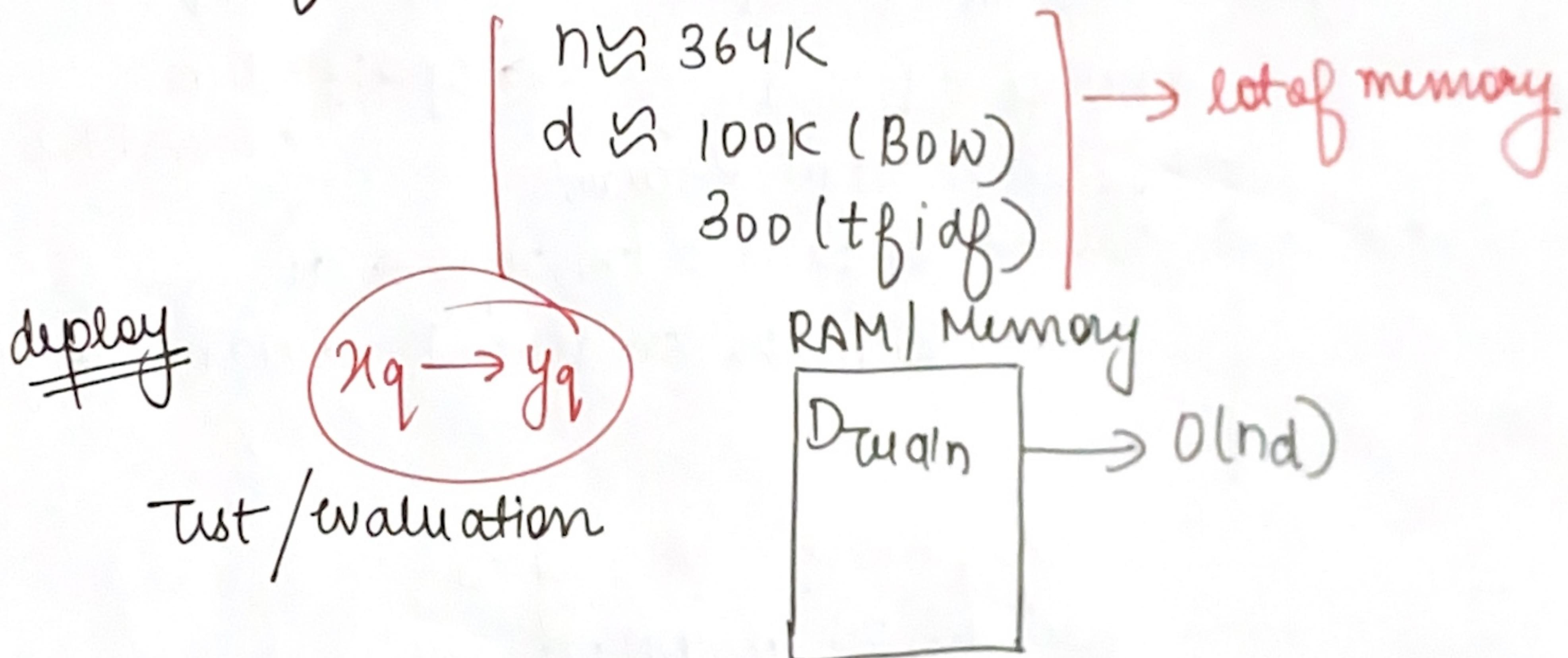
if  $d$  is small

if  $d \ll n$

$O(n)$

→ Time

### Amazon food review



Space complexity : space which is required to  
(evaluation) evaluate

$x_q \rightarrow y_q$   
 $O(nd)$

### Limitations of K-NN (Simple implementation)

(Amazon) fine food review: (real time system)

Time complexity  $O(nd)$

Space complexity  $O(nd)$

$n \approx 364K$   
 $d \approx 100K$   
36400 M  $\gg$  36GB

} a lot of space

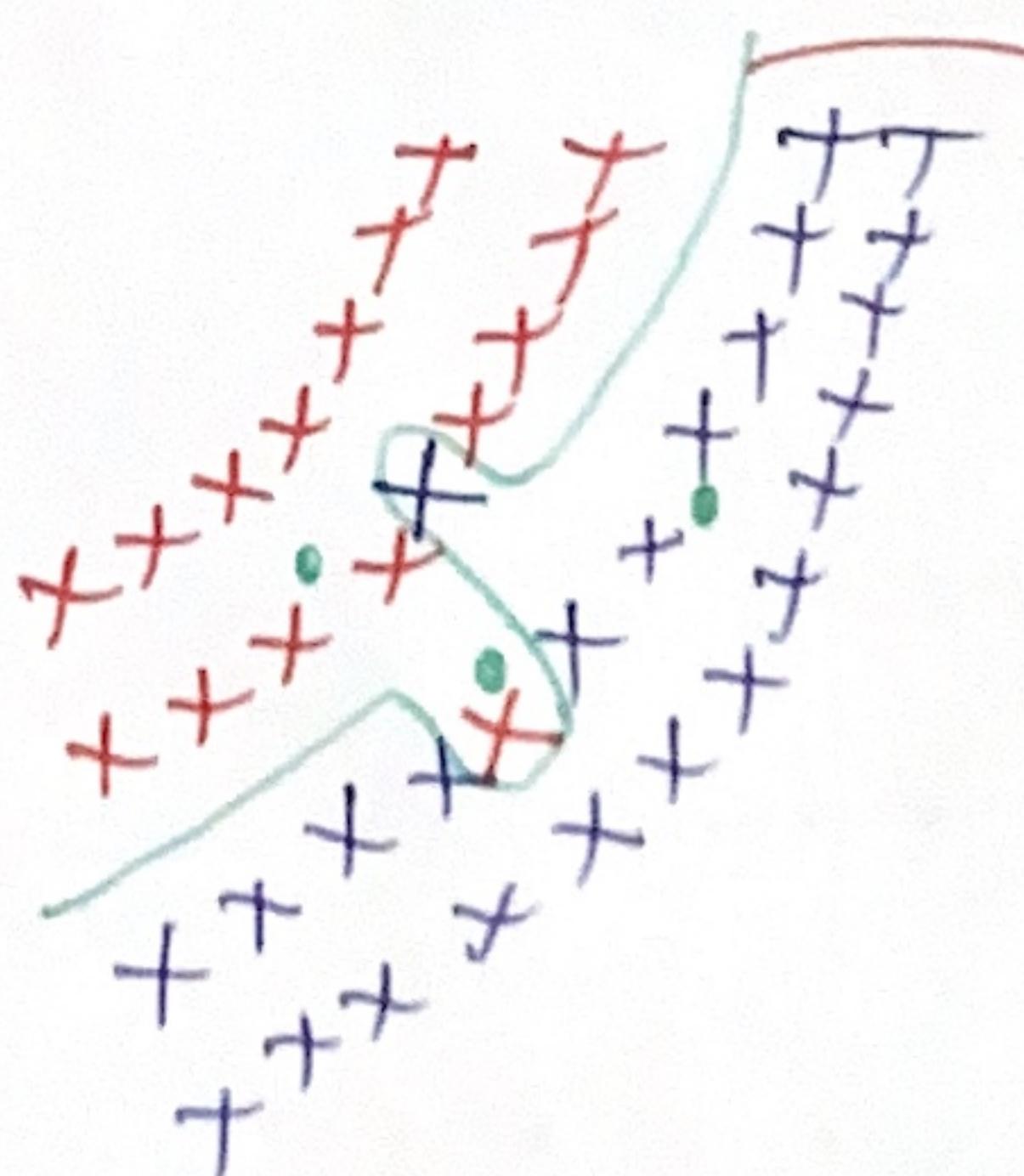
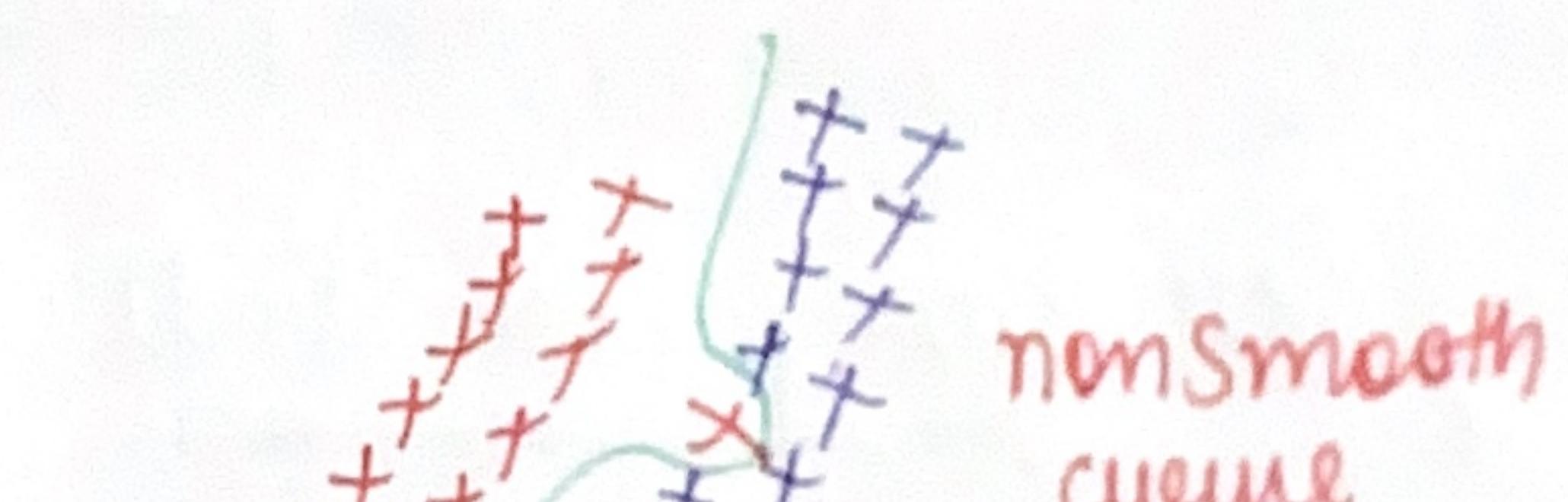
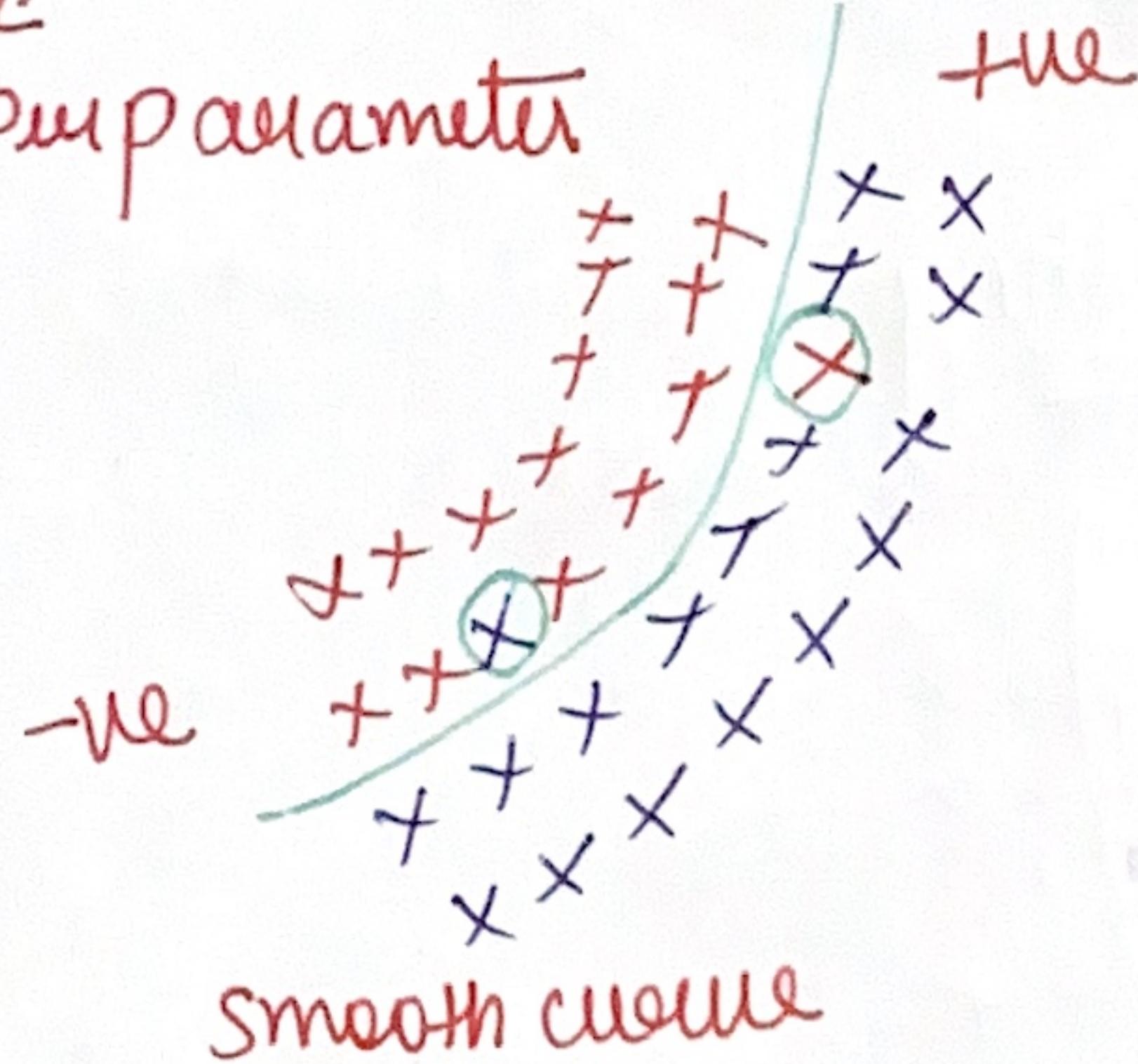
Time Compn: 86 Billion Computations

Internet      Review  $\xrightarrow{1\text{ms}}$  +ve/-ve

(low latency)  $x_q \xrightarrow{\text{fast}} y_q$   
system

KNN Simple Implementation  $\rightarrow O(nd)$   $O(nd)$

$K$  in K-NN  
hyperparameter



curves -ve from +ve  
 $\hookrightarrow$  division surface

Decision Surface  
↓  
many query points  
1-NN ( $K=1$ )  
 $\rightarrow$  separate region using curves.

let  $K=1$

1-NN

let  $K=5$

Take majority vote

$\hookrightarrow$  Smooth curve

{ as  $K \uparrow$   
smoothening of curve  
increases

let say  $k=n$

$k=1, 2, \dots, n$

$n = \text{Total no of points}$

max value of  $k=n$

$$n_1 = +ve$$

$$n_1 + n_2 = n = 1000$$

$$n_2 = -ve$$

let's say  $n_1 > n_2$

$$n_1 = 600$$

$$n_2 = 400$$

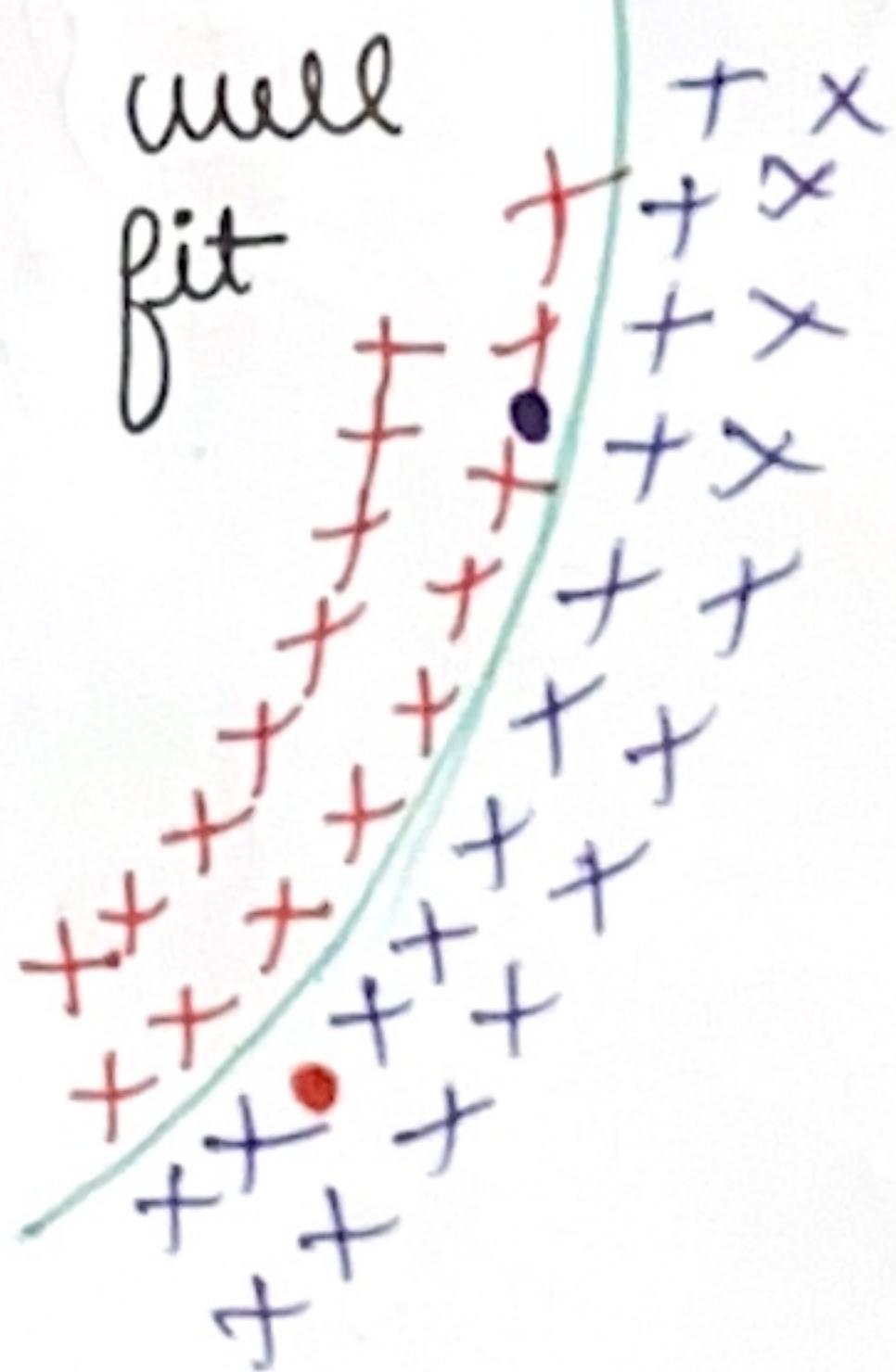
Everything will turn out +ve as majority points are positive...

## Overfitting and Underfitting (fundamental)

Over fitting



Smooth



every point belongs  
to the majority  
class

$k=1 \rightarrow k=5 \rightarrow k=n$

no-mistakes

smooth

fitting a function to data is **fitting**