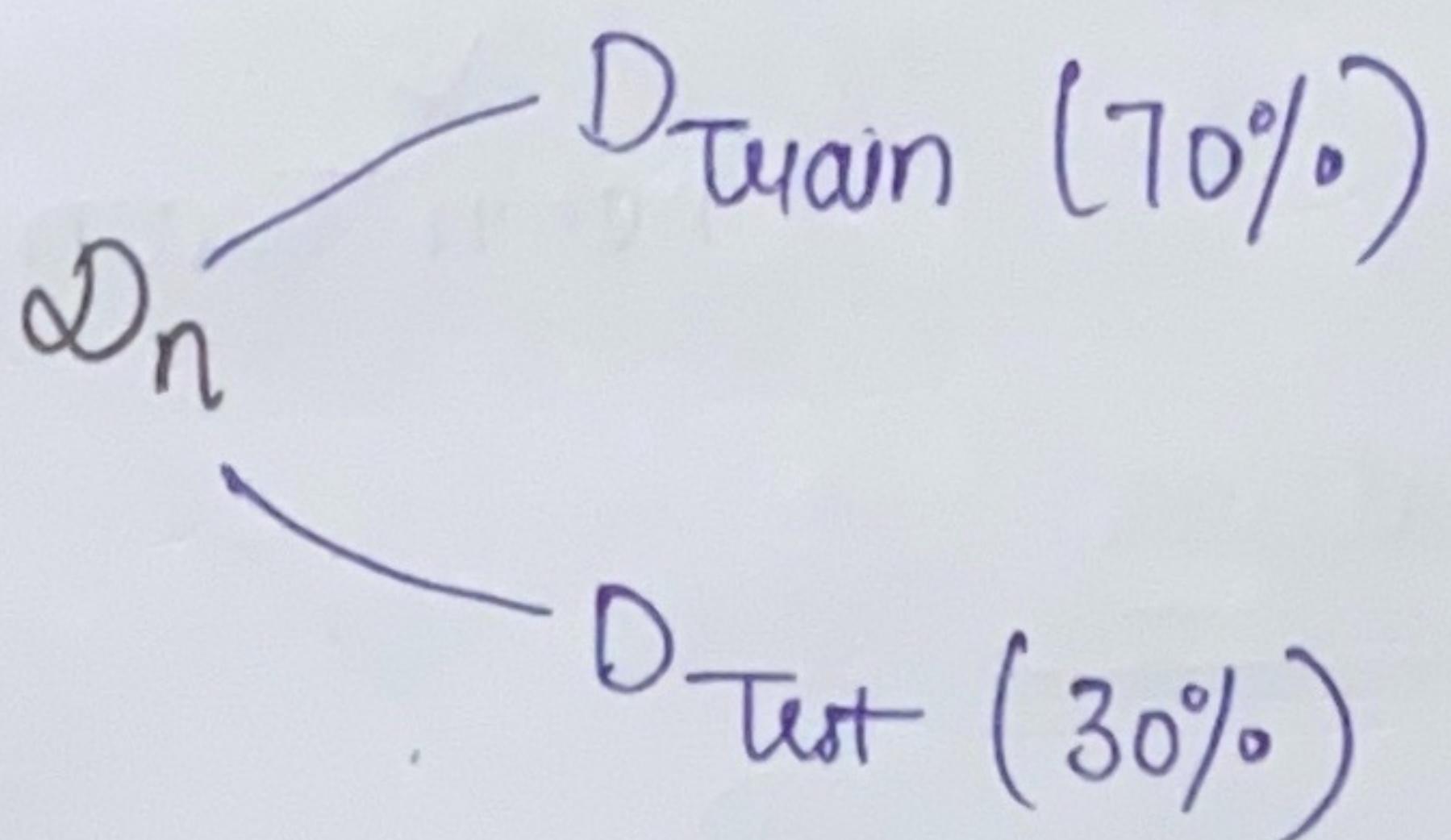


Ques How to determine 'K'



One idea

	Train	acc on $D_{\text{test}}$	no of correctly pts Total # pts
$K=1$	$D_{\text{train}}$	0.78	
$K=2$	"	0.82	
$K=3$	"	0.85	
$\vdots$		$\vdots$	

$$D_{\text{test}} = \left\{ (x_i, y_i) \mid x_i \in \mathbb{R}^d, y_i \in \{0, 1\} \right\}_{i=1}^{n_{\text{test}}}$$

$$x_q \rightarrow y_q$$

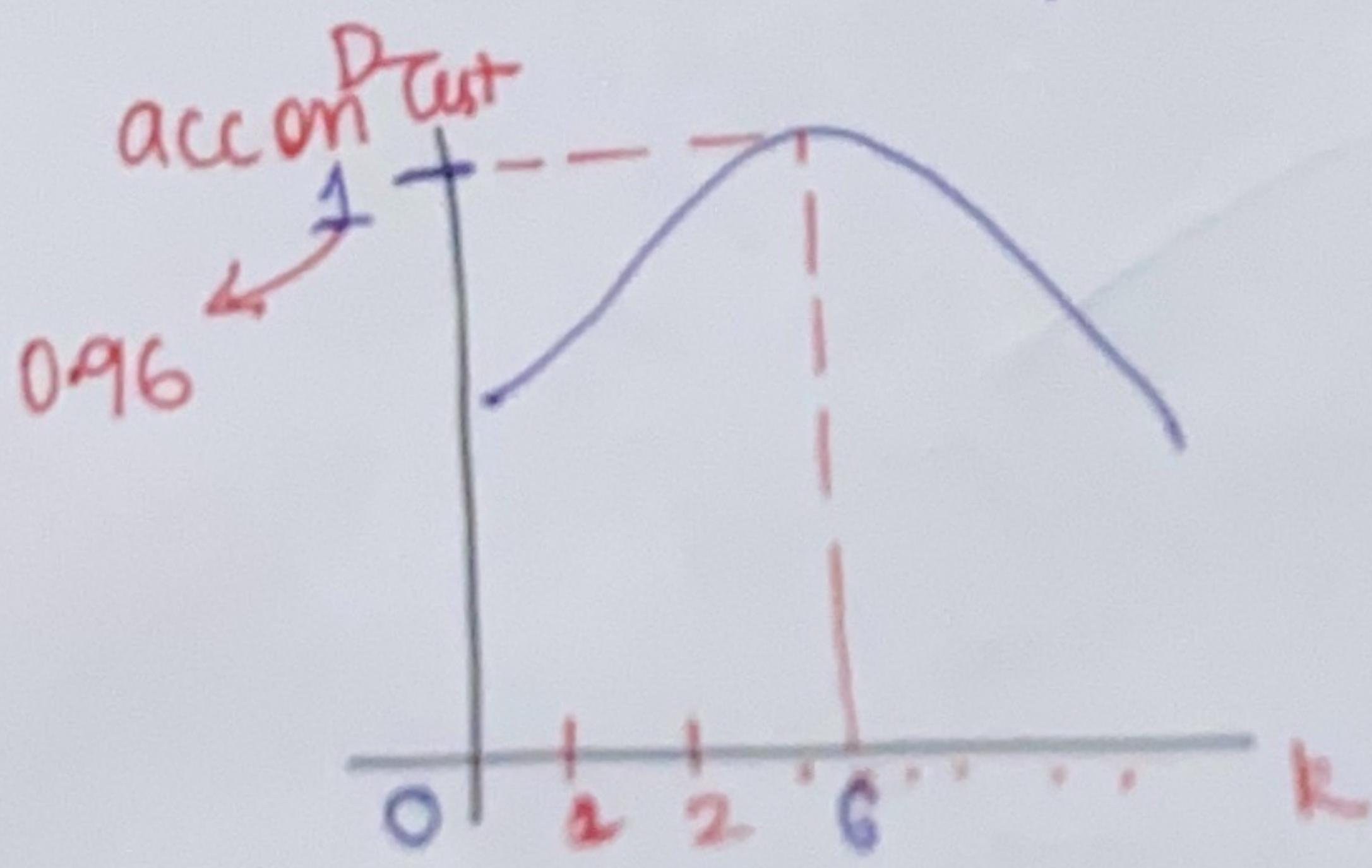
$$x_q = x_1$$

$$x_q = x_2$$

$$x_q = x_3$$

$\vdots$

$$x_q = x_{n_{\text{test}}}$$



(Typically)

$K=6$  gives me  
but accuracy on  
 $D_{\text{test}}$  taken using  
 $D_{\text{train}}$  as training  
data.

Using  $D_{\text{train}}$  & 6-NN on Amazon

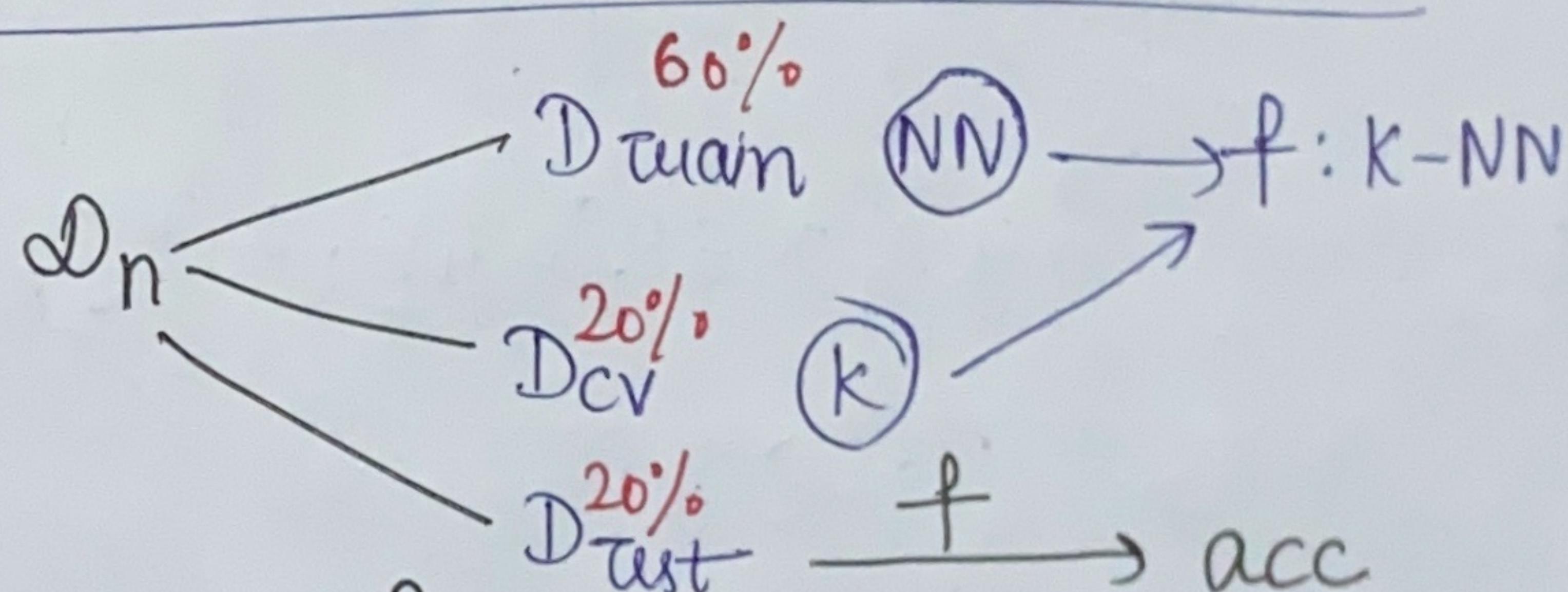
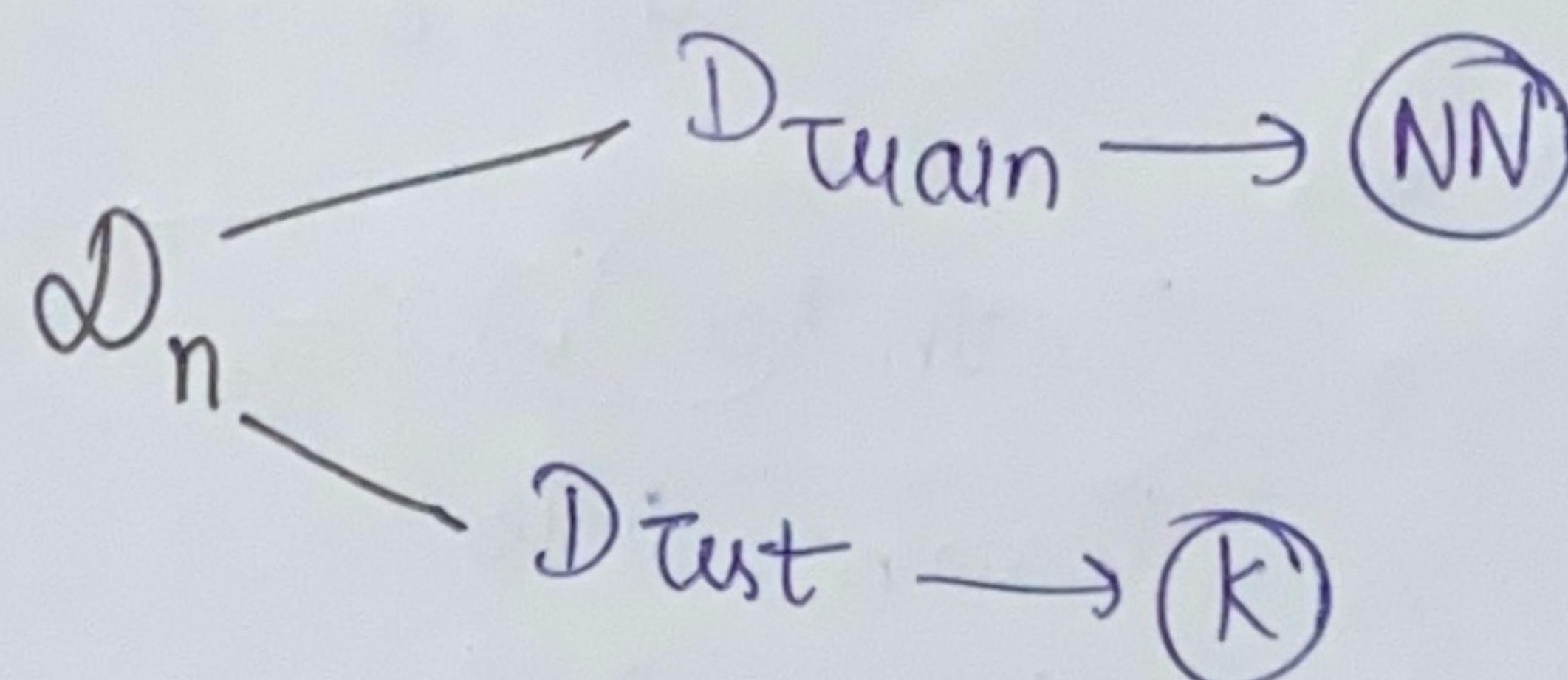
food review dataset I get an accuracy of 96%

obj perform well on "unseen" pts



not in train & test

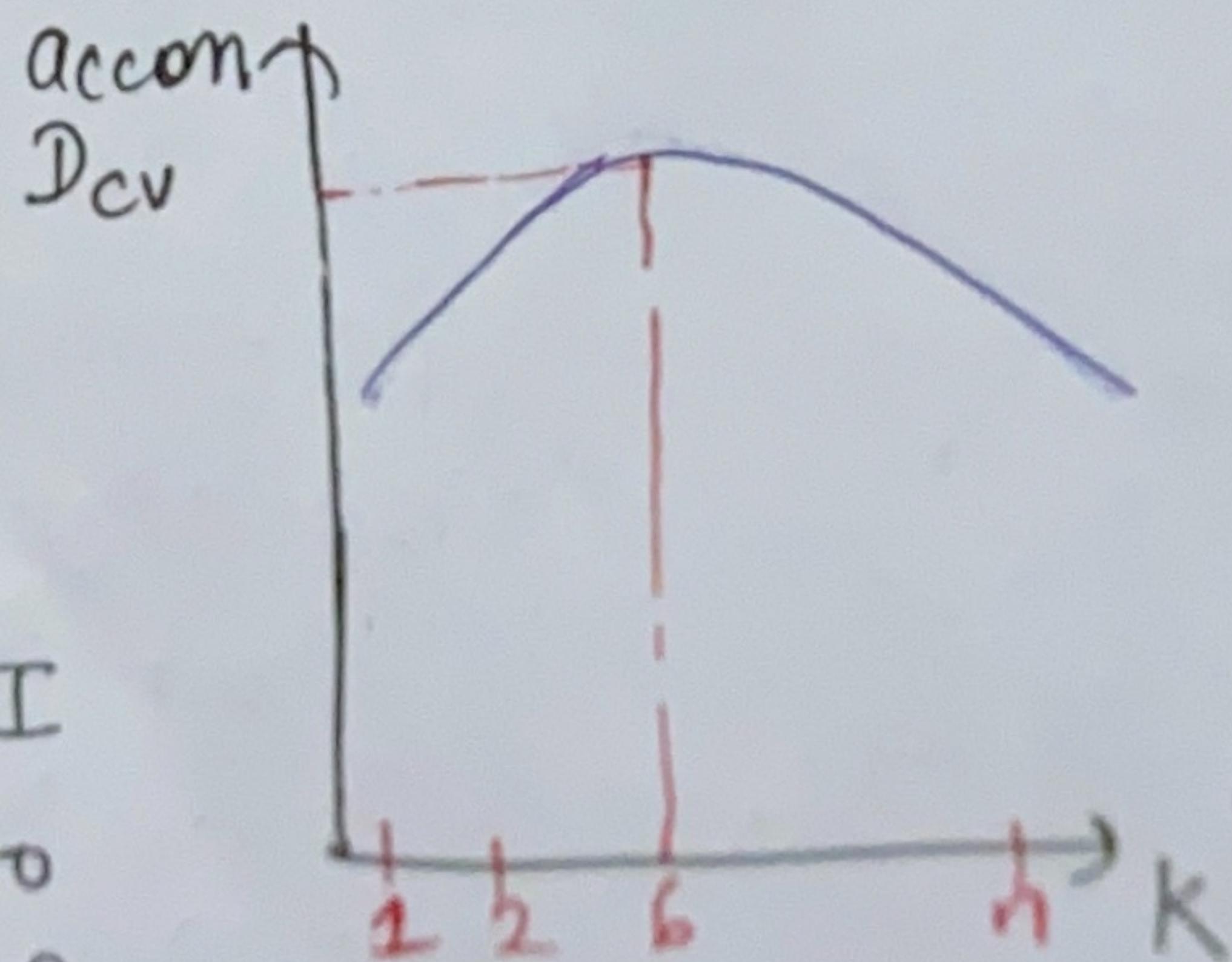
## Cross Validation



{ can I now say 6-NN has  
acc of 93% on unseen  
data? }

(unseen data)

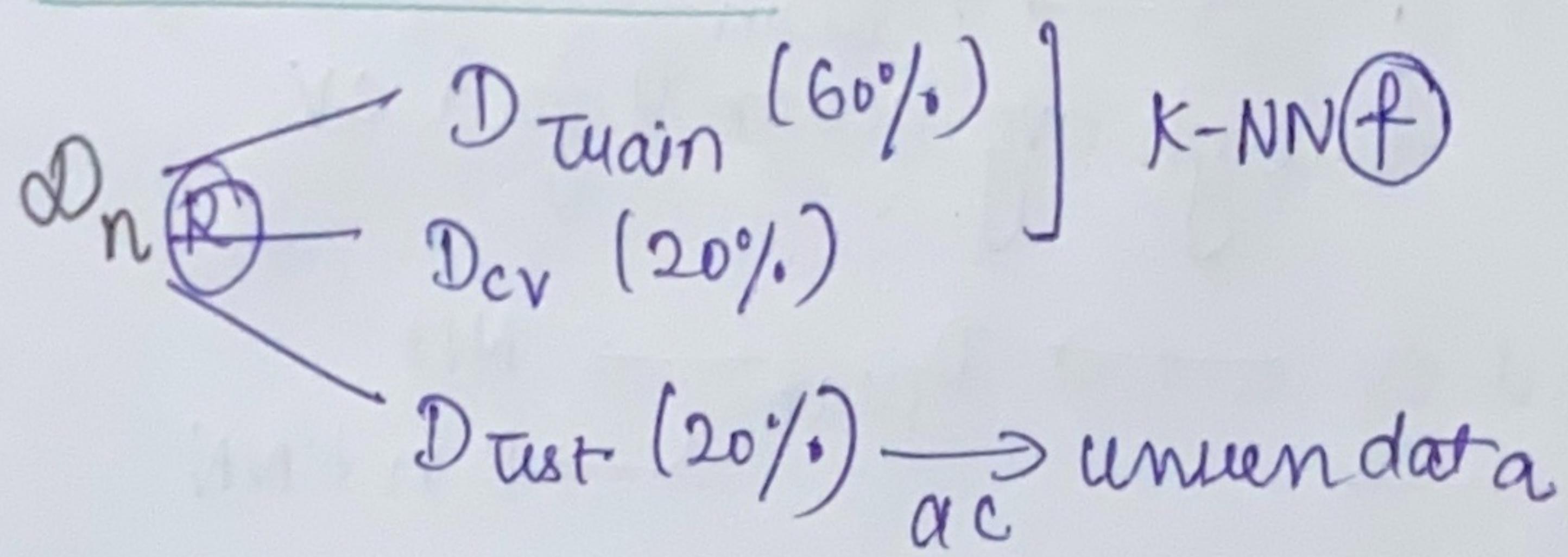
→  $x_i$  in  $D_{\text{test}}$  among  
acc 93%



using  $D_{\text{train}}$  as  
training data I  
found 6-NN to  
have an acc of  
93% on unseen data

→ generalisation accuracy 7% = generalisation error

# K-fold Cross Validation

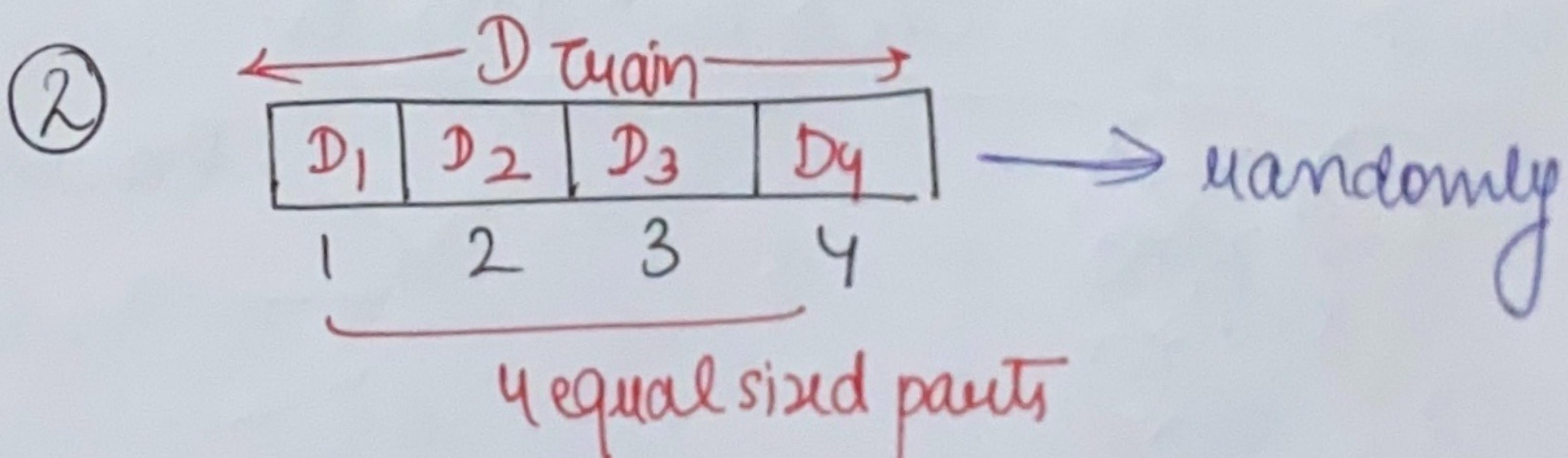
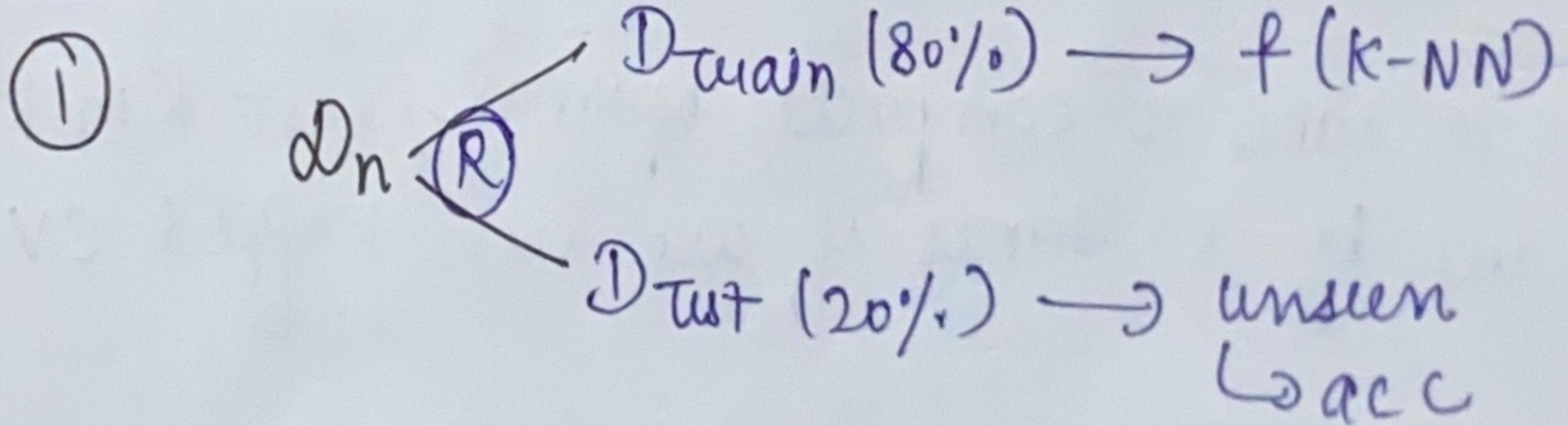


## Problem

only using 60% of data to compute NN  
 20%  $\rightarrow$  CV  
 20% Test (unseen)

80% to compute NN

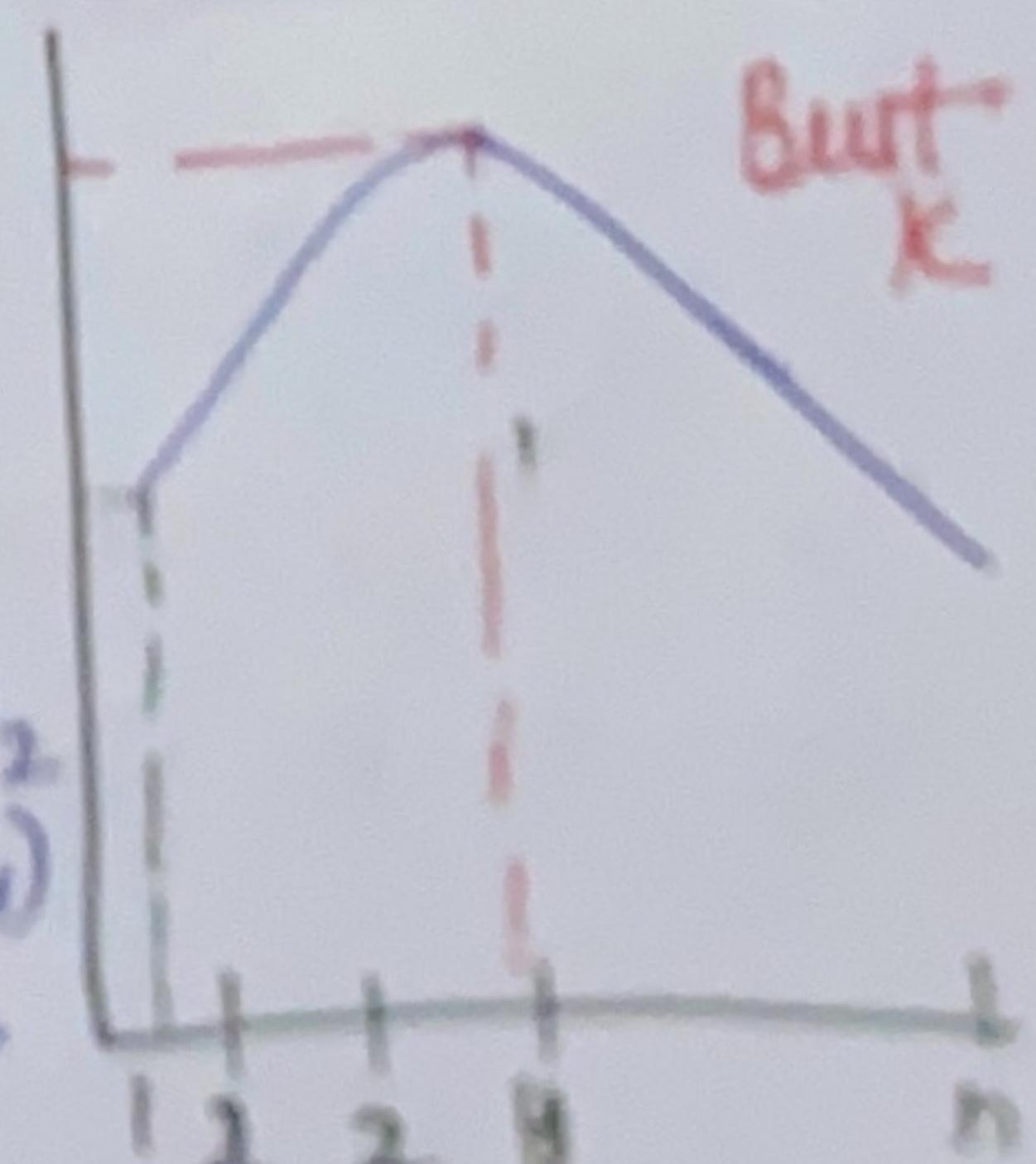
## Solution



③

	Train	CV	acc on CV
$k=1$	$D_1, D_2, D_3$	$D_4$	$a_4$
$k=1$	$D_1, D_2, D_4$	$D_3$	$a_3$
$k=1$	$D_1, D_3, D_4$	$D_2$	$a_2$
$k=1$	$D_2, D_3, D_4$	$D_1$	$a_1$

acc on CV dataset



4 fold cross validation times

again repeat

avg all of k-NN for k=1 on CV data

$f: 3\text{-NN} \rightarrow D_{\text{train}}$   
 $\downarrow$   
 through avg acc on 4 fold CV

$k'$  fold CV  $\rightarrow D_{\text{train}} \rightarrow$  NN  
 $\downarrow$   $\hookrightarrow k \text{ in } k\text{-NN}$

What is right  $k'$

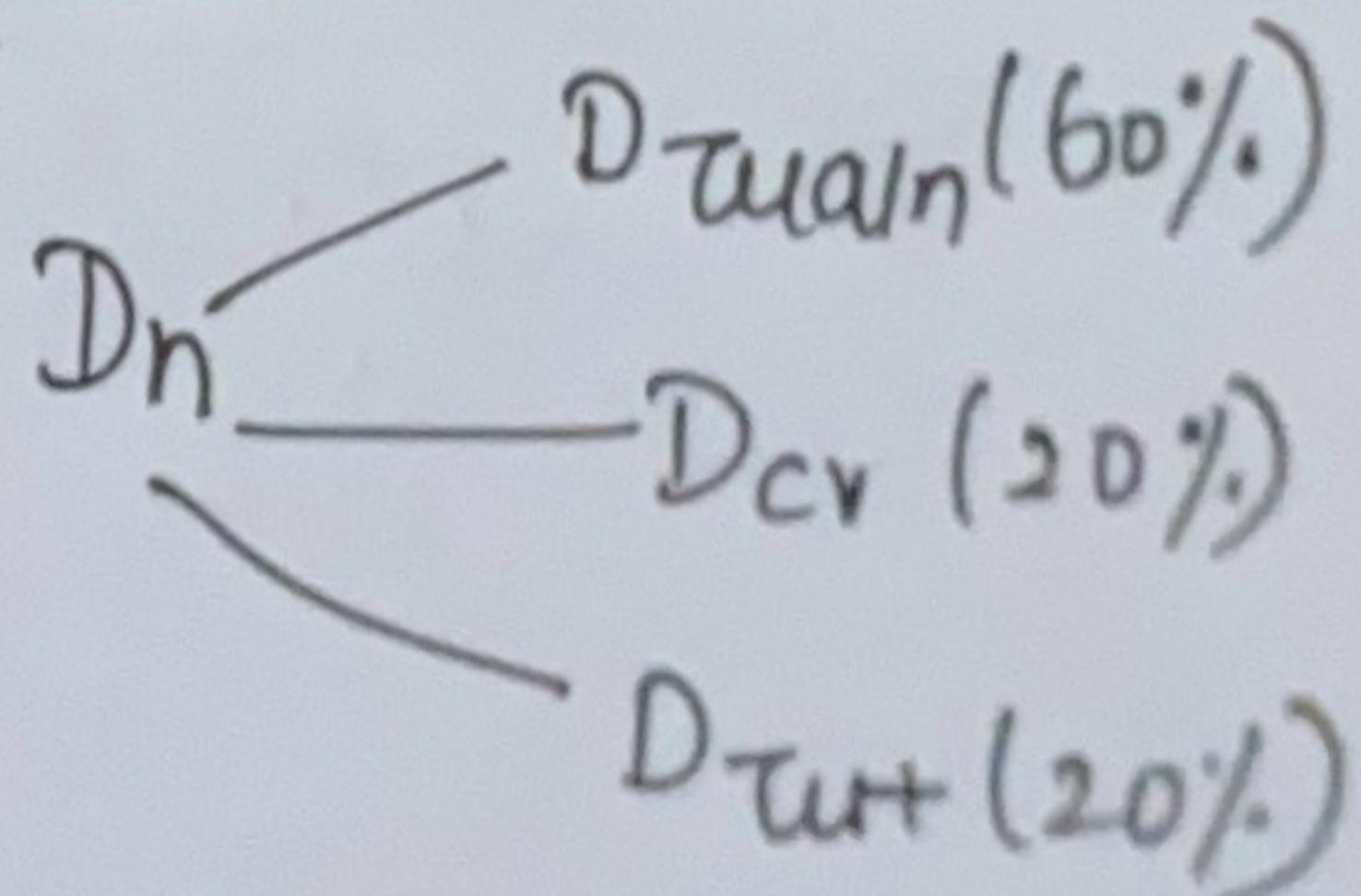
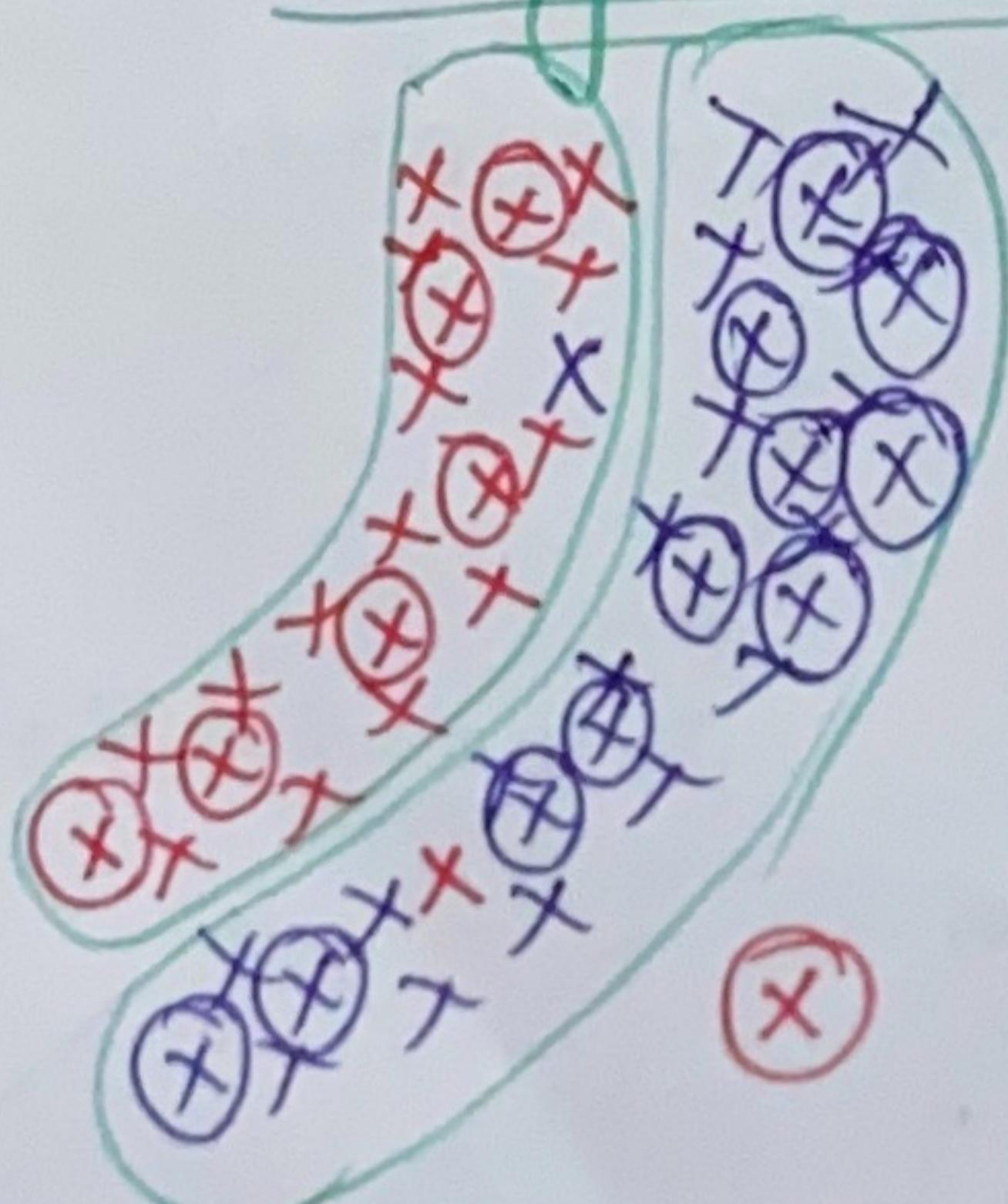
$k' = 4$   $k' = 10$   $k' = 100$

rule of thumb: 10 fold CV

limitation

time it takes to compute optimal/last  $k$  in  $k\text{-NN}$   
 increases by  $k'$  times if we use  $k'$  fold CV

Visualizing Train, CV & test datasets:



$\circlearrowleft$ : in datapoint in  $D_{\text{train}}$   
 $\times$ : the " " "

- ( $\circlearrowleft$ ): -ve pt from  $D_{\text{cv}}$
- ( $\times$ ): +ve pt from  $D_{\text{cv}}$

- ①  $D_{\text{train}}$  &  $D_{\text{cv}}$  do not overlap perfectly
- ② If there are many +ve pts from  $D_{\text{train}}$  in region then it is likely to find many +ve pts from

③ If there are very few true pts in region from  $D_{\text{train}}$  then it is unlikely to find true pts from  $D_{\text{cv}}$  in that region

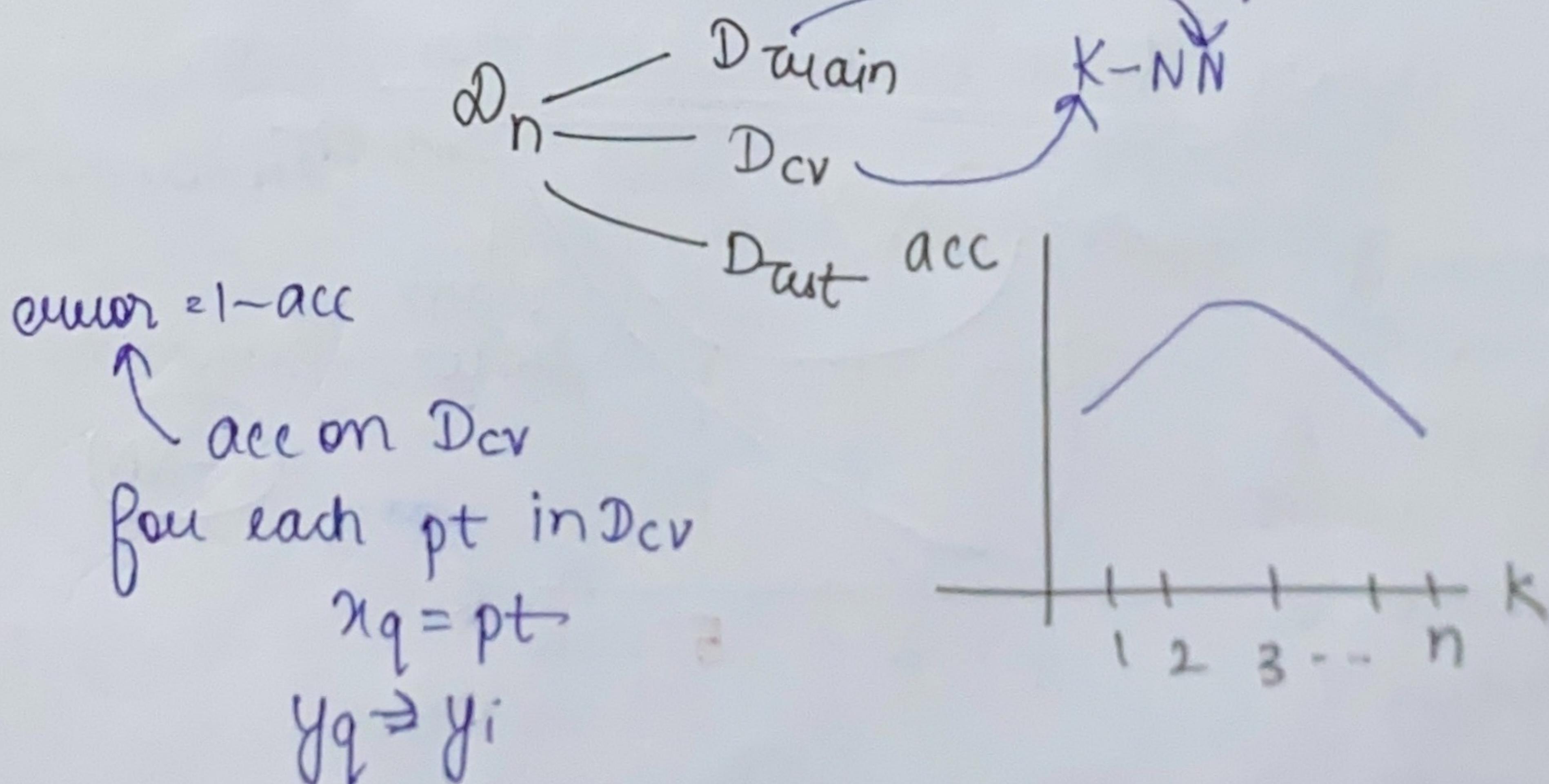
### Overshooting vs Underfitting

$k$ -fold CV or  $D_{\text{cv}}$  → best  $k \rightarrow$  neither overfit  
nor underfit

Plus How we can be sure that we are not underfitting or overfitting → plot

$$\max \uparrow \text{acc} = \frac{\# \text{ correctly class pts}}{\text{total } \# \text{ pts}} = 0.93$$

$$\min \downarrow \text{error} = 1 - \text{acc} = 0.07\% = 7\%$$



## Training error

- for each  $x_i$  in  $D_{\text{train}}$
- find 2 nearest  
neigh to  $x_i$  from

$D_{\text{Train}}$

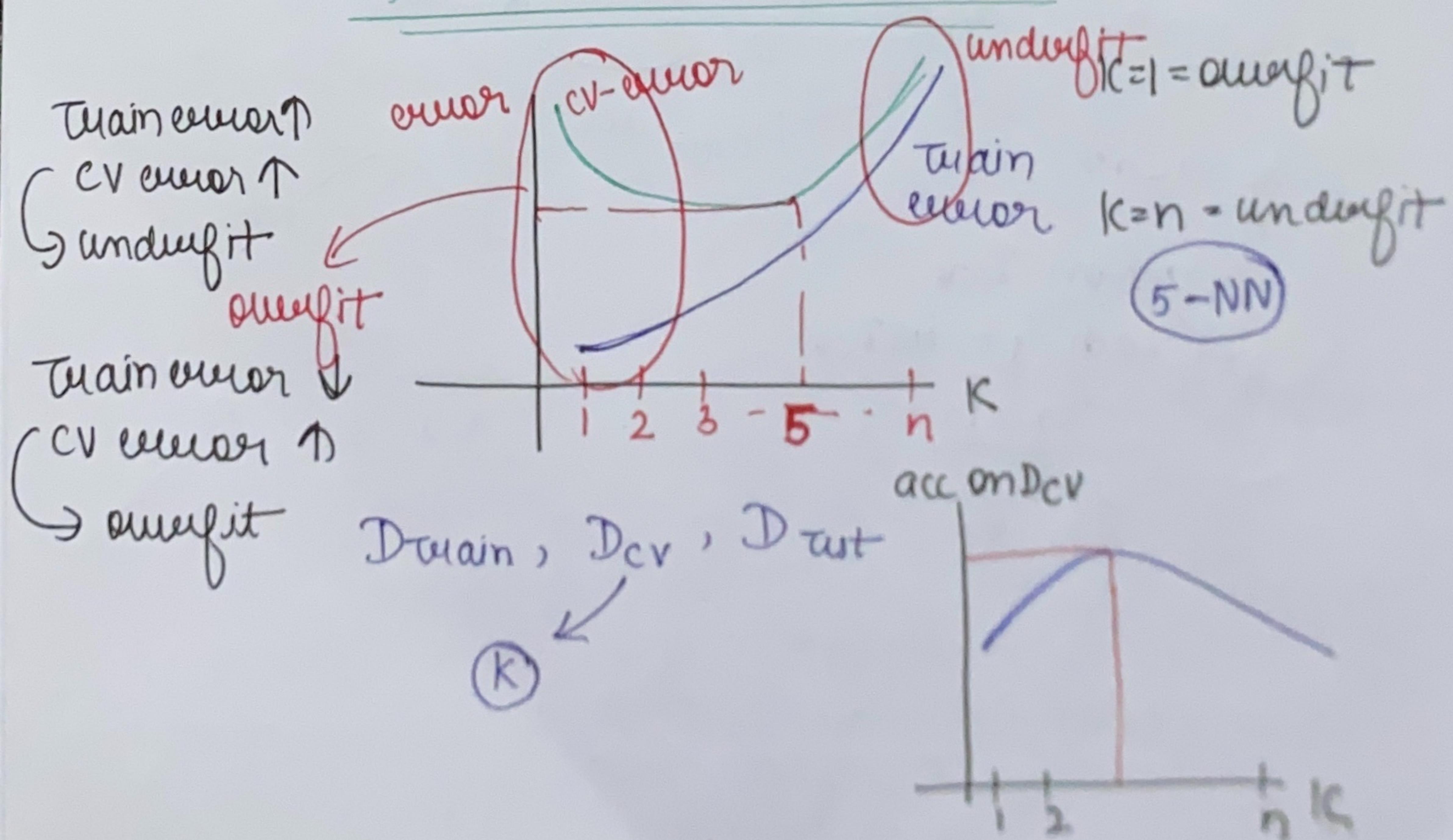
- majority vote to get class label
- if  $y_i = \text{class label}$   
accurate
- else error

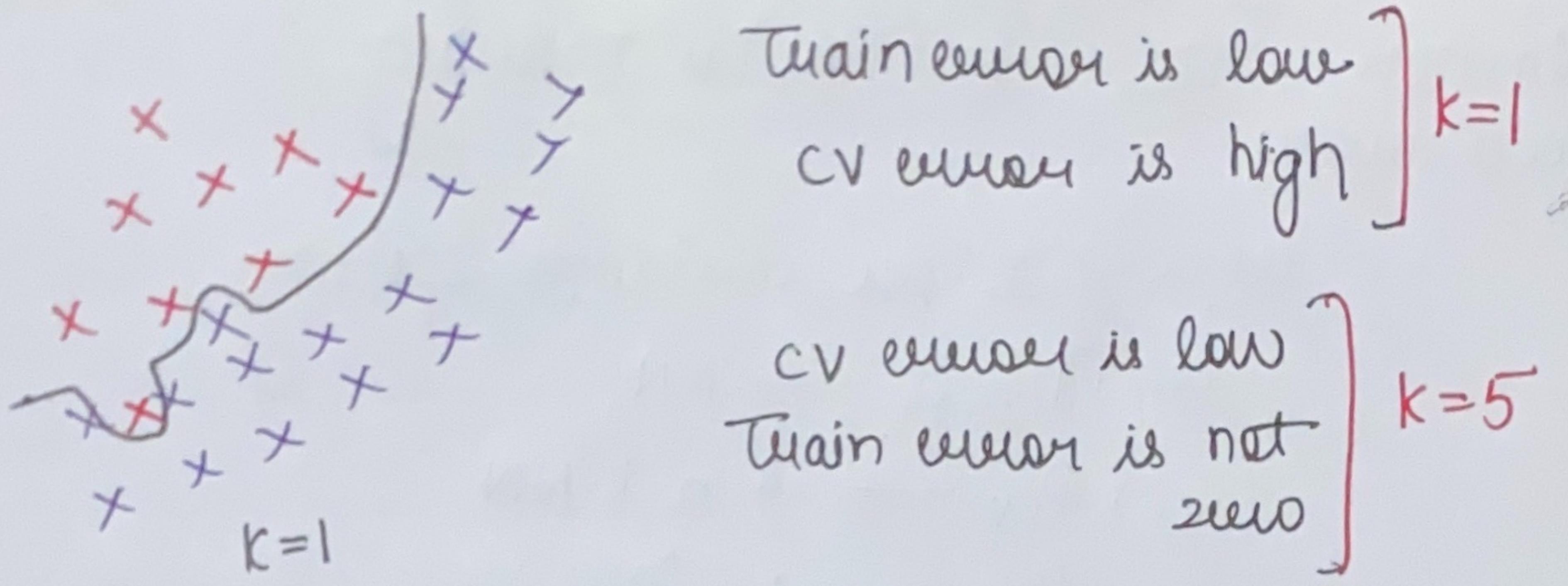
what is error of 2-NN on  $D_{\text{cv}}$

$D_{\text{train}} = D_{\text{Train}}$

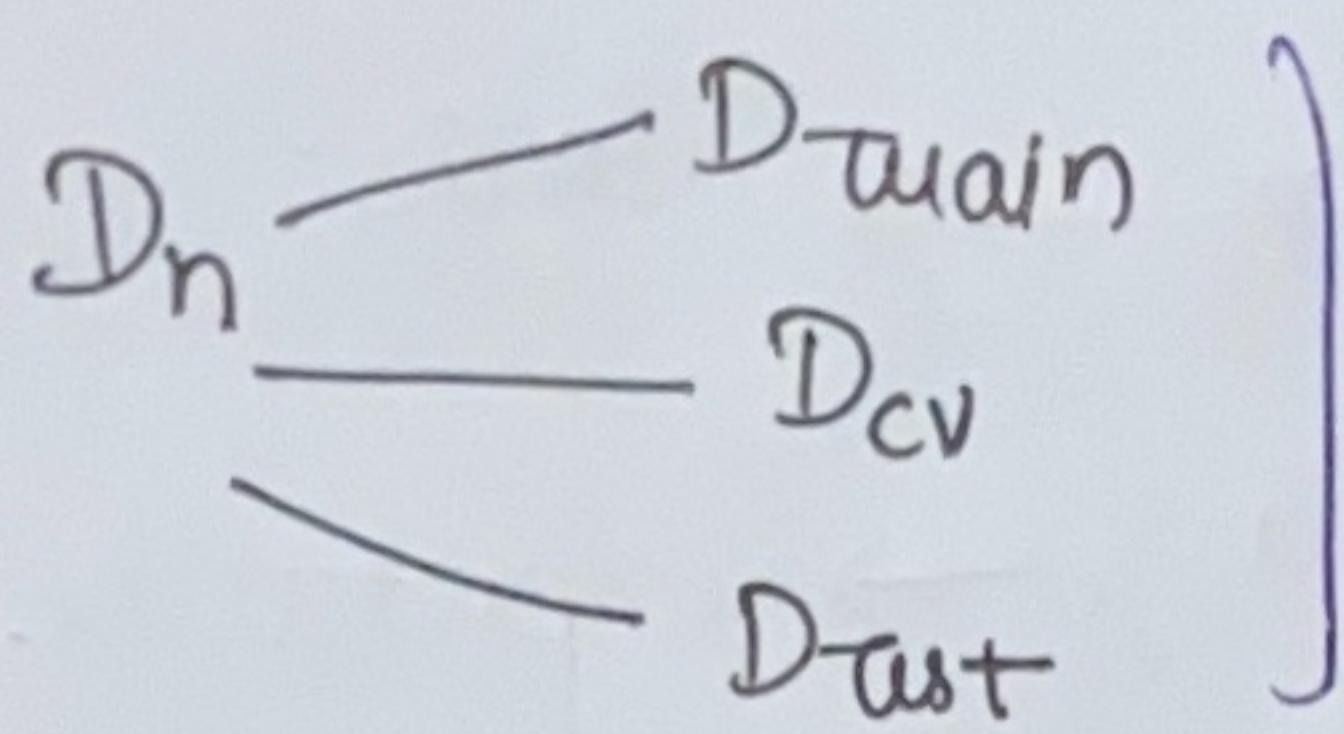
$D_{\text{cv}}$ : error on  $D_{\text{cv}}$

error; train error; validation error

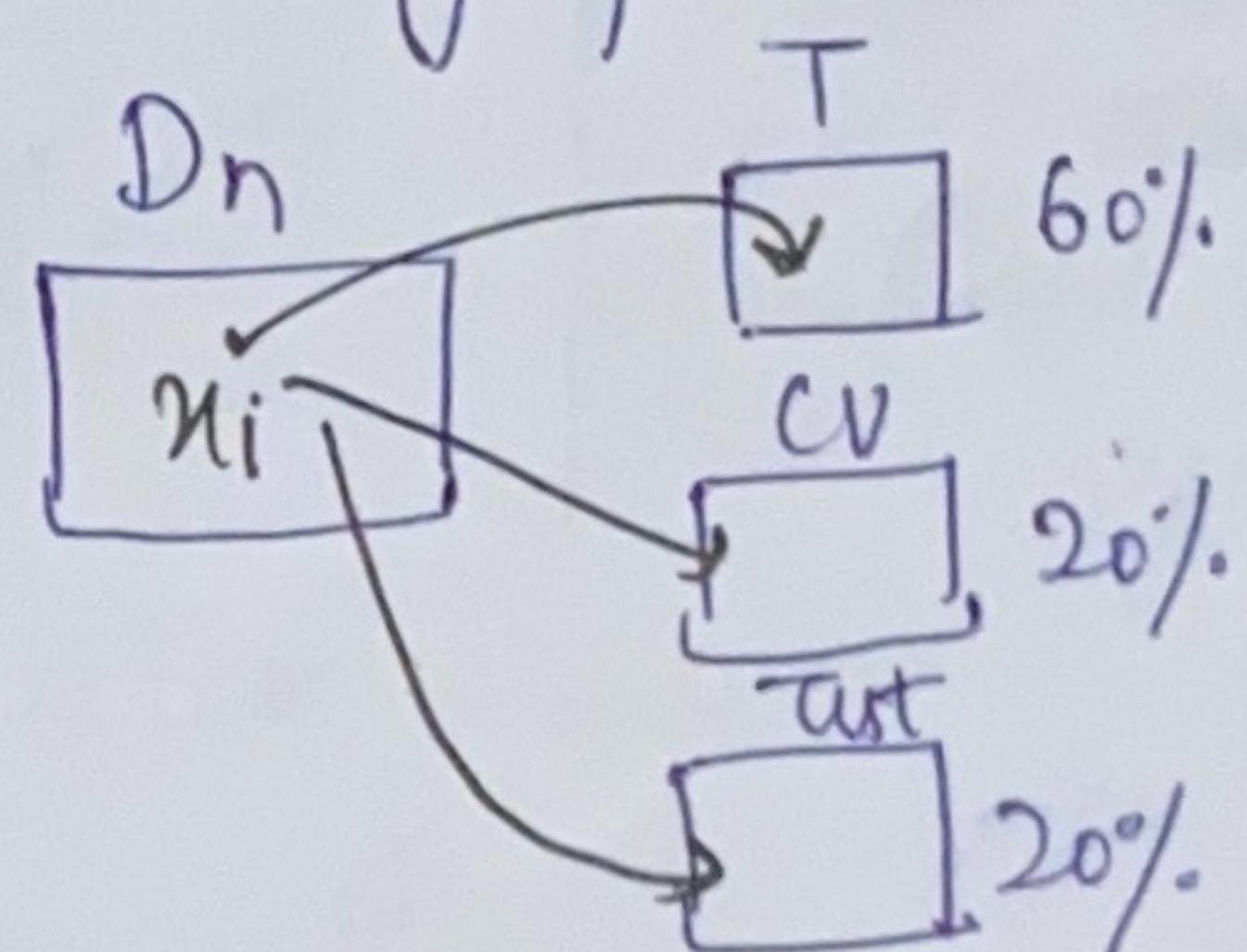




## Amazon Food Review



Randomly split

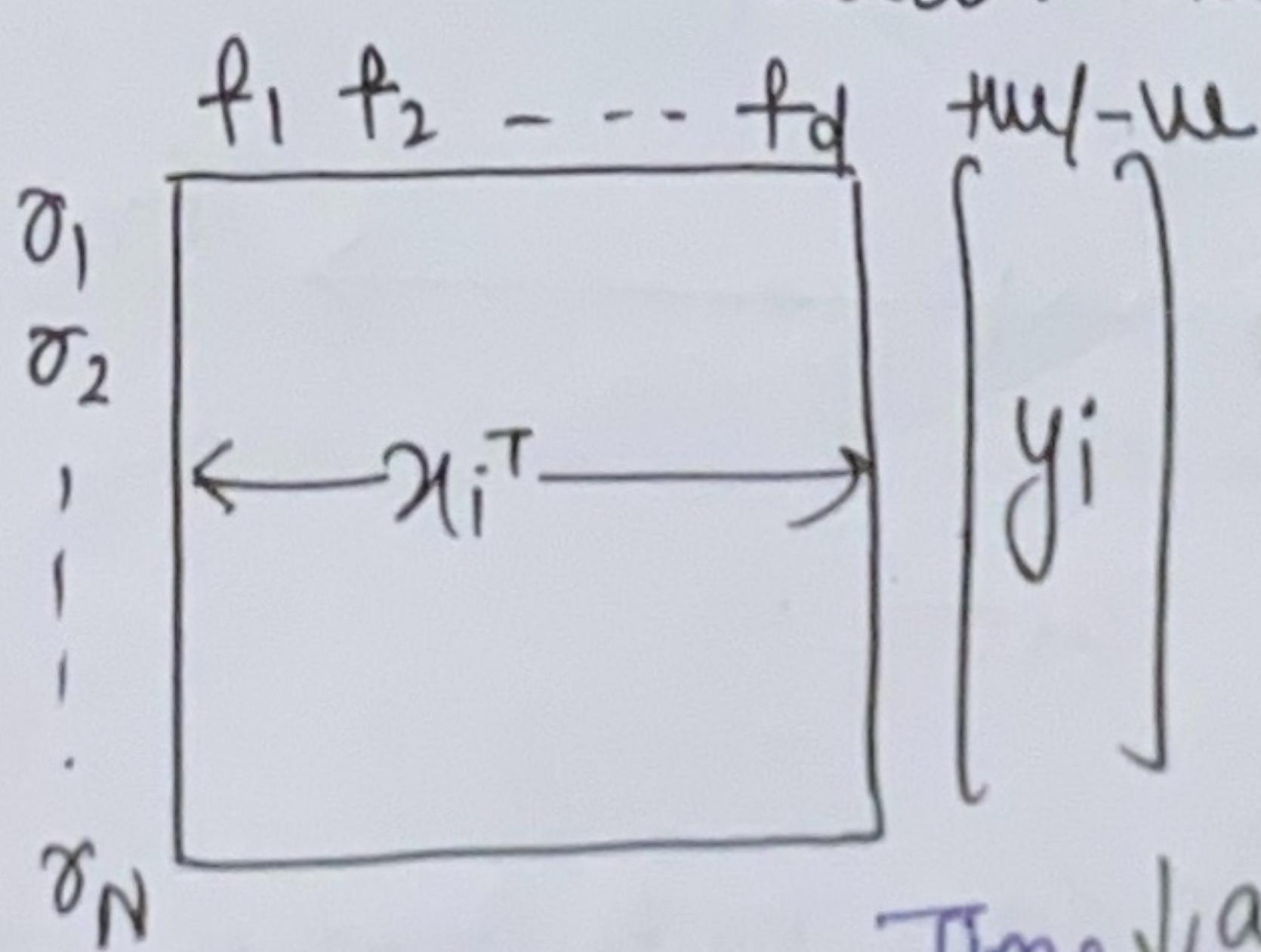


## Time Based Splitting

→ depends on case

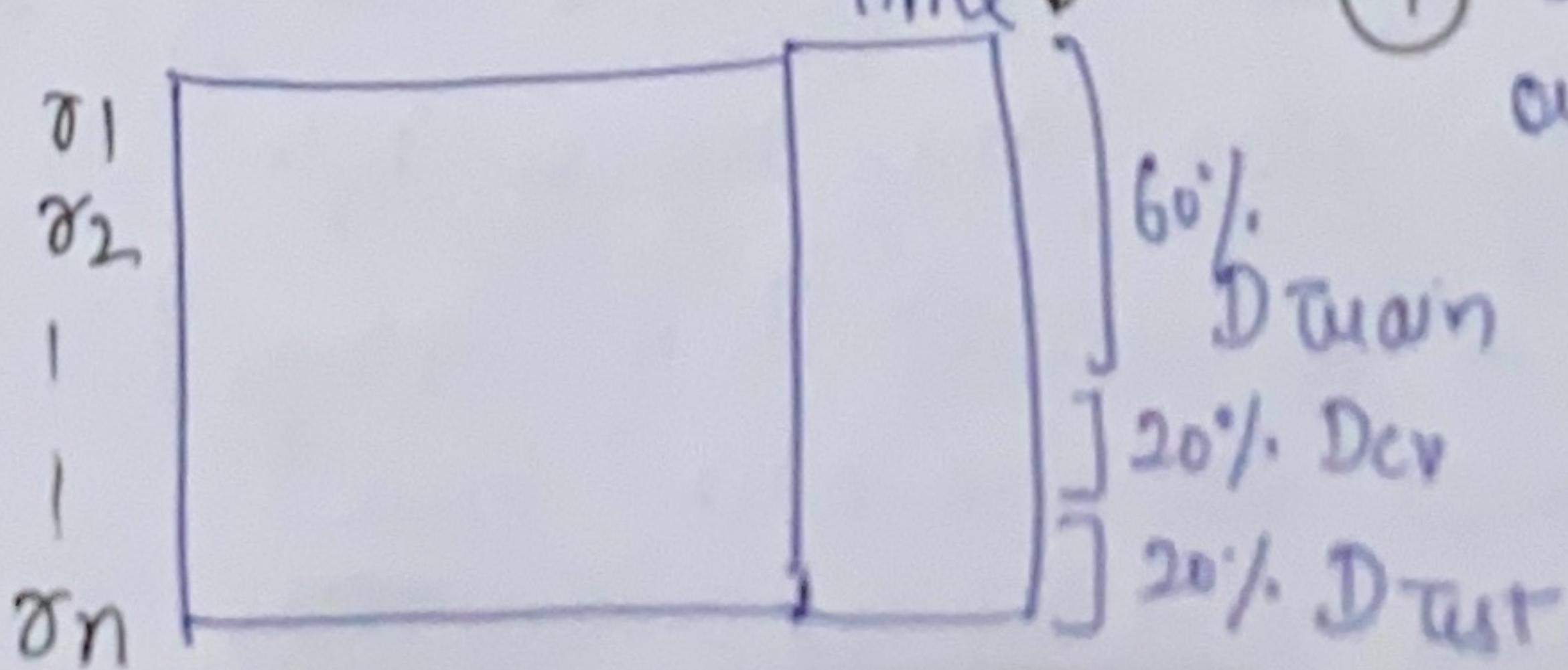
better than random-splitting

②



Time ↓ asc

TBS



① Sort D<sub>n</sub> in asc order of Time

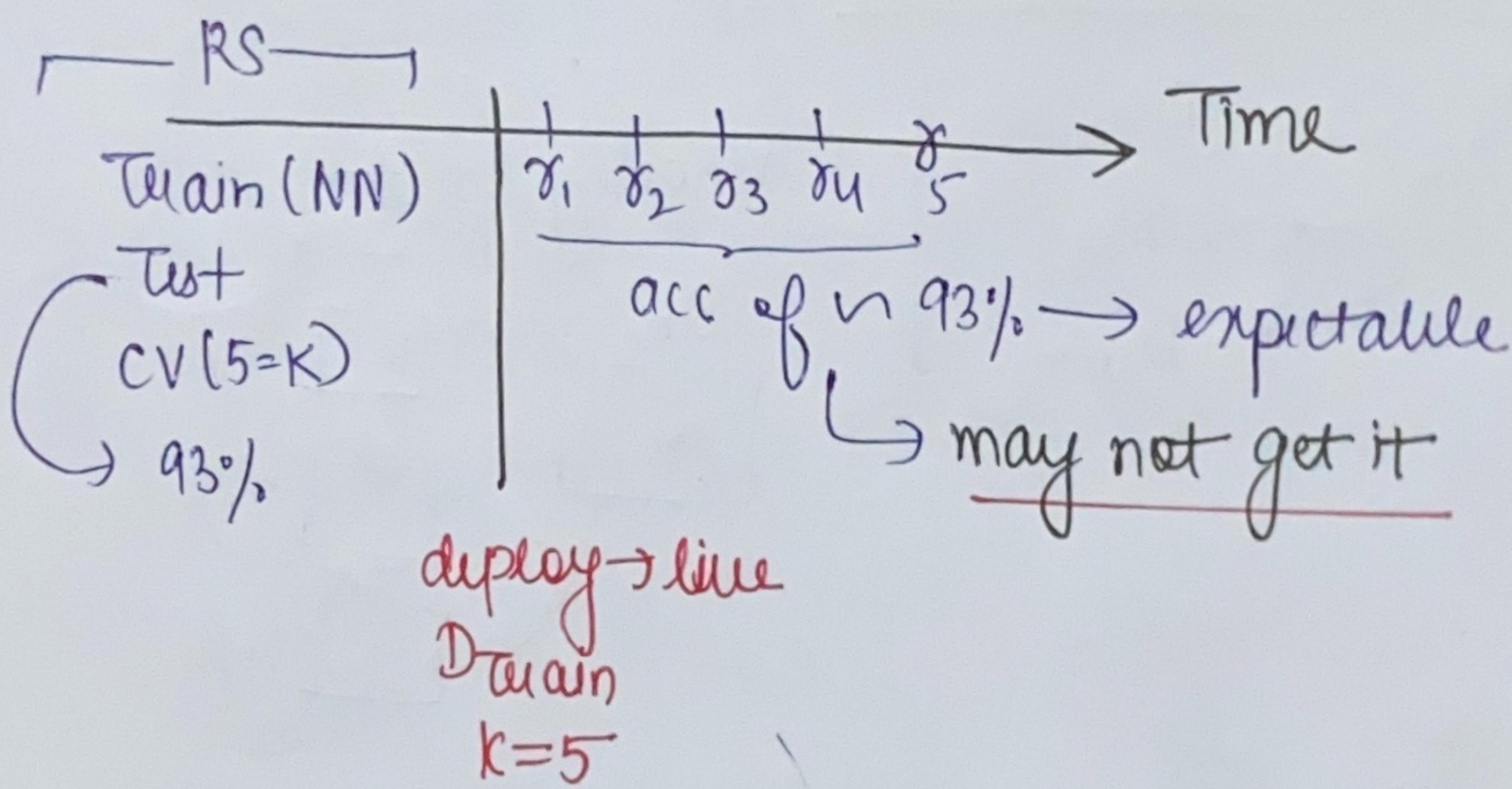
Amazon TBS  $\rightarrow$  (Time) is necessary  
food review

let say I use random splitting

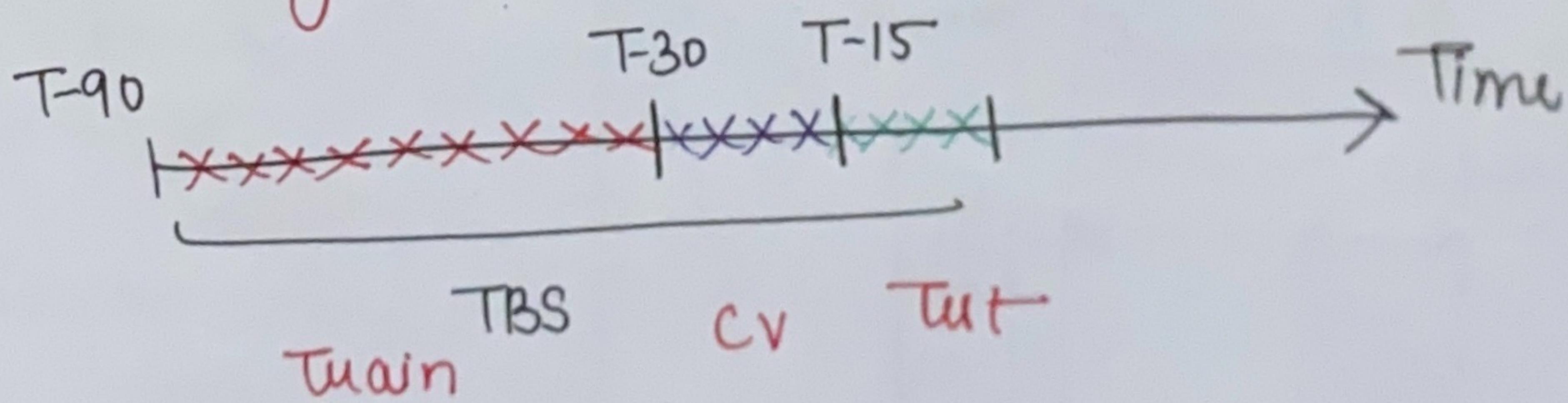
(60%)  $D_{\text{train}} \rightarrow \text{NN}$

(20%)  $D_{\text{cv}} \rightarrow k \text{ in } k\text{-NN}$

(20%)  $D_{\text{test}} \rightarrow \text{acc } (93\%)$



with time, my products and their reviews change



If I trained my model 15 days back and if I tested it on last 15 days data, my acc is 91%.

$\rightarrow$  whenever time is available & if things / behaviour of data changes over time then time based splitting is preferable.

## K-NN for regression

classification  
2(class)

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

$x_q \rightarrow y_q$

Regression

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

$x_q \rightarrow y_q \rightarrow \text{number}$

① given  $x_q$ , find K-nearest neighbour

$$(x_1, y_1), (x_2, y_2), \dots, (x_k, y_k)$$

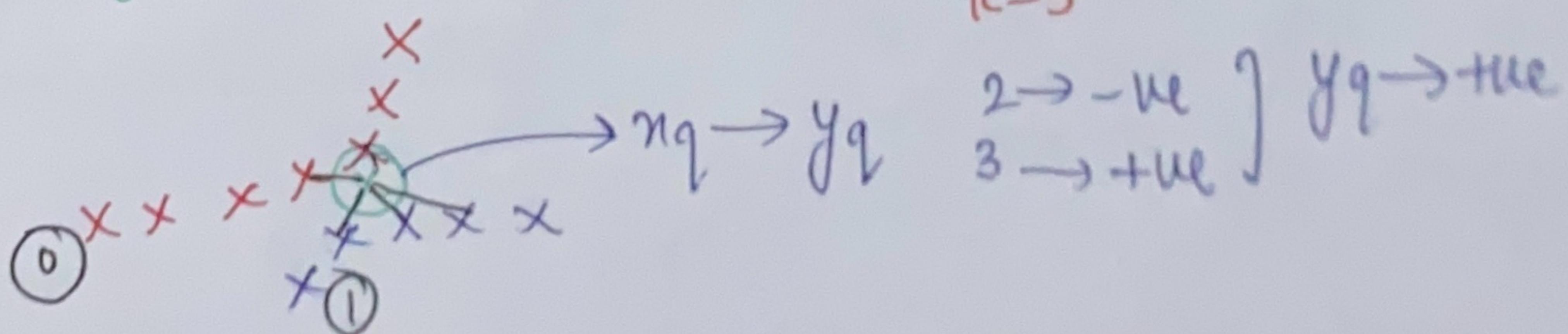
②  $\underbrace{y_1, y_2, y_3, \dots, y_k}_{\hookrightarrow \text{mean values}} \rightarrow y_q$

$$y_q = \text{mean } (y_i)_{i=1}^k$$

$$y_q = \text{median } (y_i)_{i=1}^k \rightarrow \text{less prone to outliers}$$

extend  
classif → regression

Weighted K-NN



weighted  
K-NN

$$n_q \rightarrow x_1, y_1, d_1 (0.1) -ve \quad \text{10} \quad \boxed{15}$$

$$x_2, y_2, d_2 (0.2) -ve \quad \boxed{5}$$

$$w_i = \frac{1}{d_i} \left[ \begin{array}{l} \text{Simplist weight function} \\ \text{for } n_3, y_3, d_3 (1.0) +ve \\ n_4, y_4, d_4 (2.0) +ve \\ n_5, y_5, d_5 (4.0) +ve \end{array} \right] \boxed{1.75}$$

$d_i \uparrow, w_i \downarrow$

$d_i \downarrow, w_i \uparrow$

$$y_q = -ve \text{ as } 15 > 1.75$$

## Voronoi Diagram (K-NN (K=1))

Voronoi diagram is partitioning of a plane into regions based on distance to points in a specific subset of the plane. That set of points (seeds, sites or generators) is specified beforehand and for each seed there is a corresponding region consisting of all points closer to that seed than to any region. These regions are called Voronoi cells.

## Kd-tree

Simple  
Implm

K-NN:  $O(n)$  if d is small  $\rightarrow$  time  
k is small  
 $O(n) \rightarrow$  space

Time:  $O(n) \rightarrow O(\lg n) \rightarrow$  Kd-tree v1975  
 $n=1024 \quad \lg(n)=10$   
Computer graphics  
geometry

## Binary Search Tree (BST)

$\rightarrow$  binary search (Python)

prob: given a sorted array, find a number is present in array or not.  $q=15$

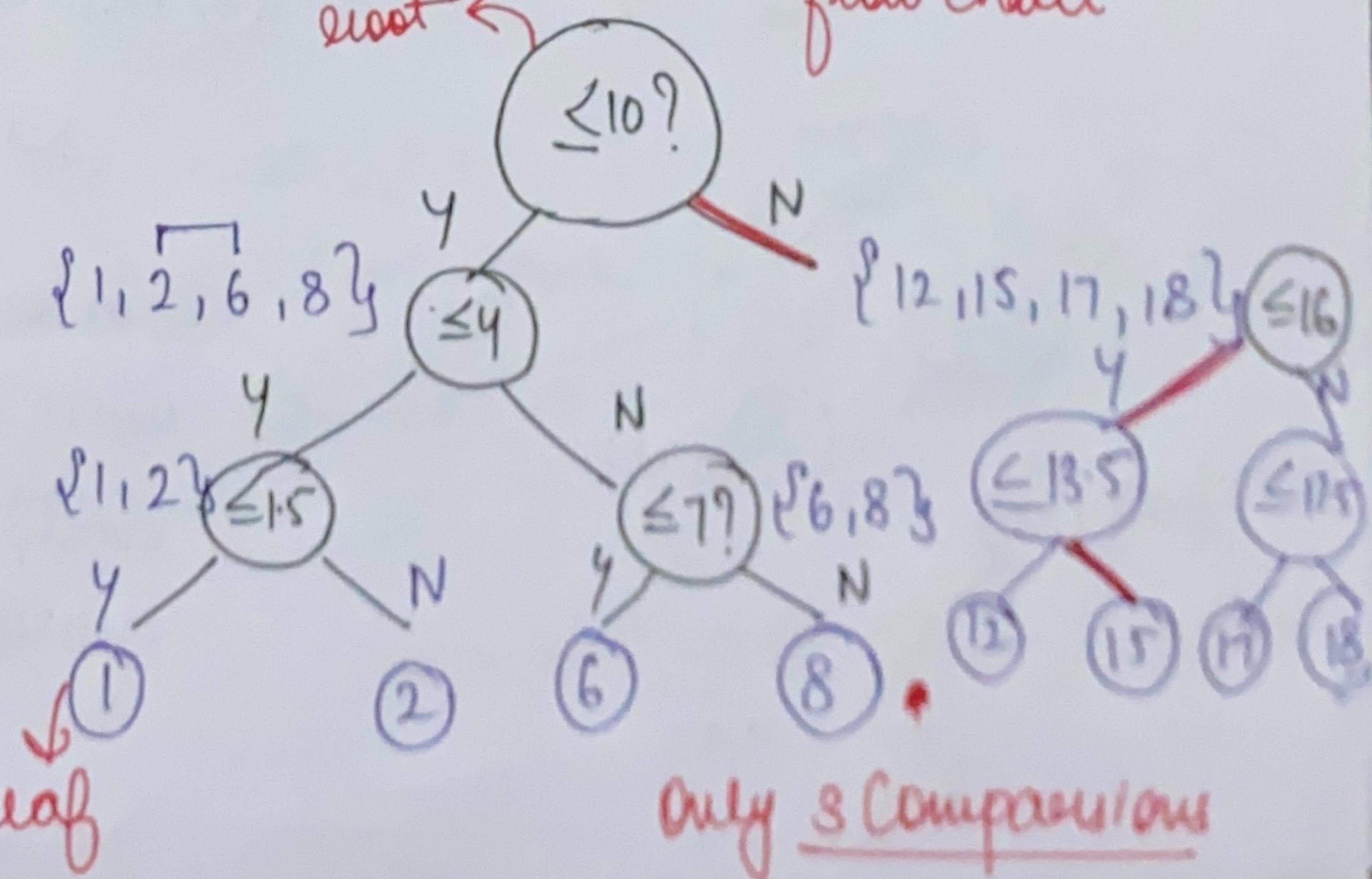
A:  $\boxed{1|2|6|8|12|15|17|18} \leftarrow$  Sorted

① Construct a BST

Time:

root  $\leftarrow$

flow chart



Linear Search: no of comparisons of  $O(n)$

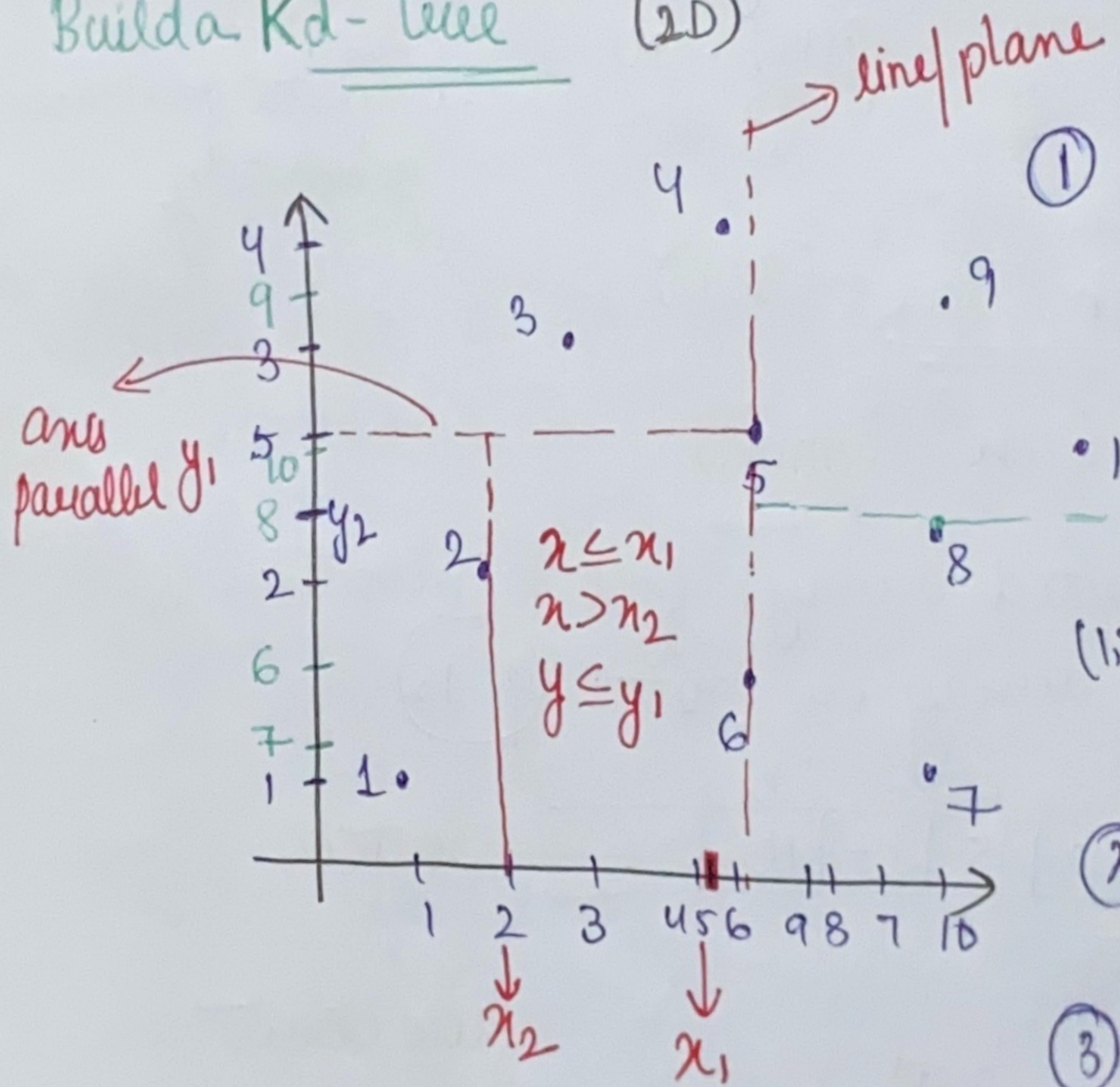
Binary Search :  $O(\log n)$

Time complexity to search is  $O(\log n)$

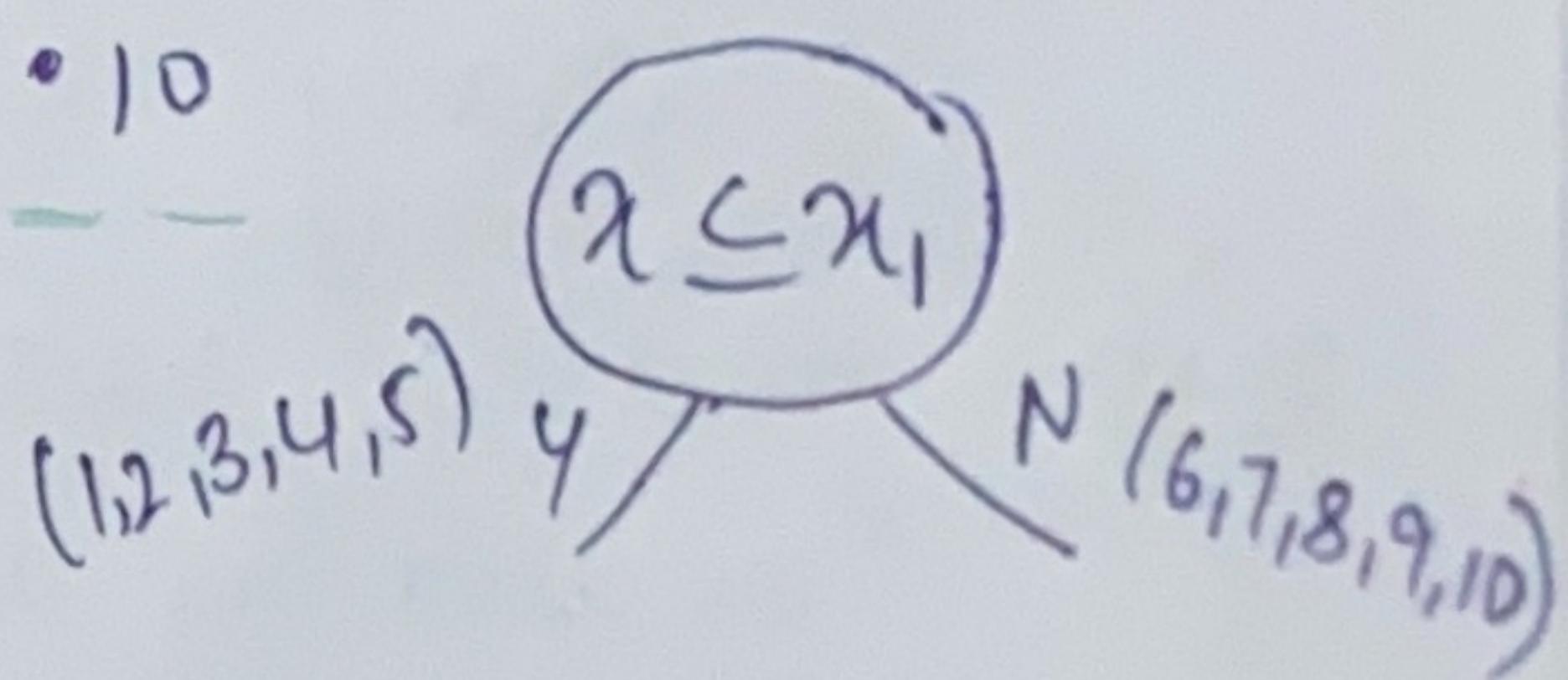
↓  
depth of the tree

Build a Kd-Tree

(2D)

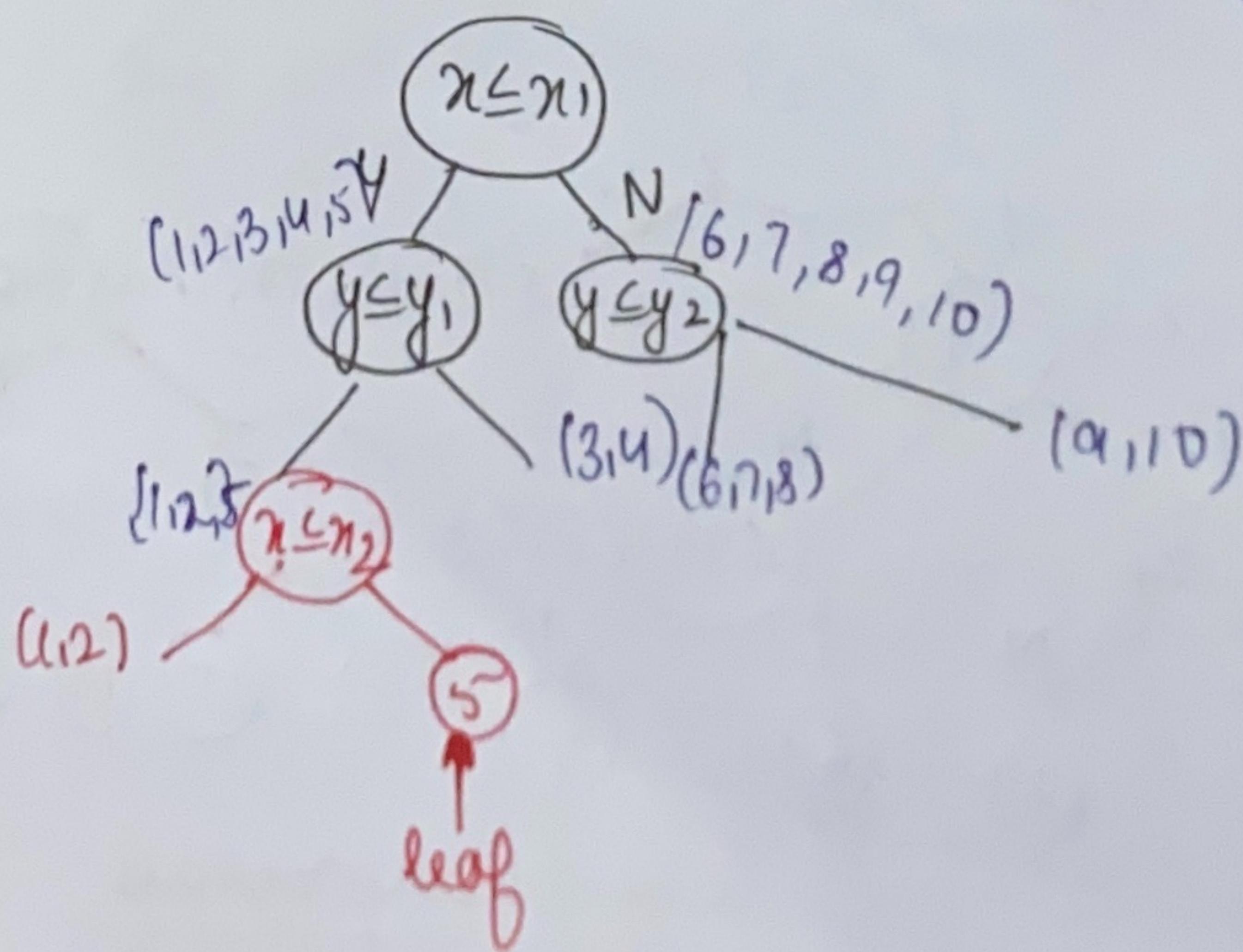


① pick  $n$  axis, project pts on  $n$  axis and compute median



② split data using median

③ Alternate b/w axis



Kd tree

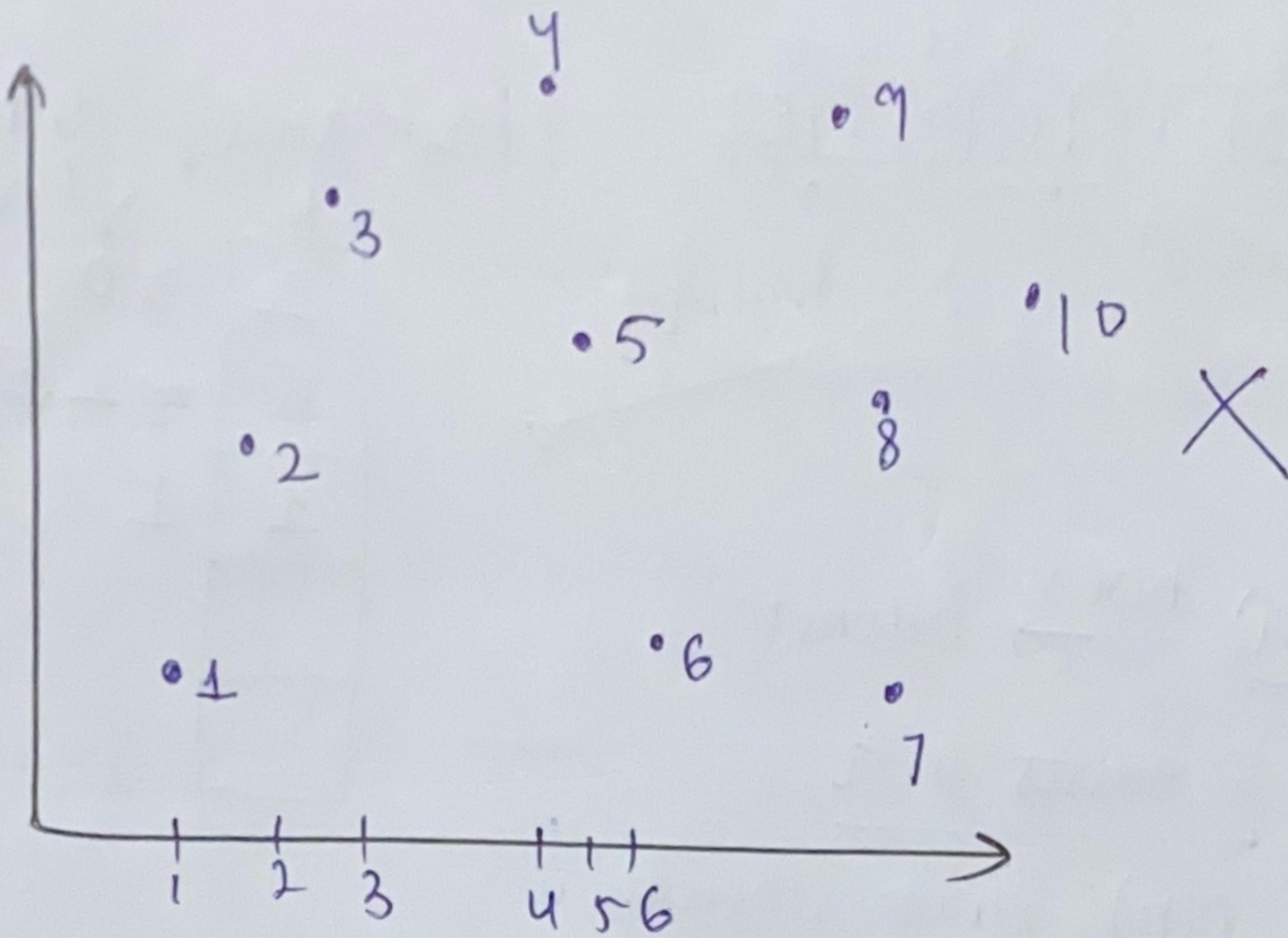
↓

breaking my space using axis parallel line/ plane into cuboids/ hypercuboids

2D  $\rightarrow$  axis parallel lines  $\rightarrow$  rectangles

3D  $\rightarrow$  " " planes  $\rightarrow$  cubeoids

nD  $\rightarrow$  " " hyperplanes  $\rightarrow$  hypercubeoids



### Limitation of kd tree

① When d is not small

$$d=10 \quad 2^d = 1024$$

$$d=20 \quad 2^d \approx 1M$$

$d \uparrow$ ; (worst-case) # adj cells  $2^d \uparrow$

When d is not small {2, 3, 4, 5}

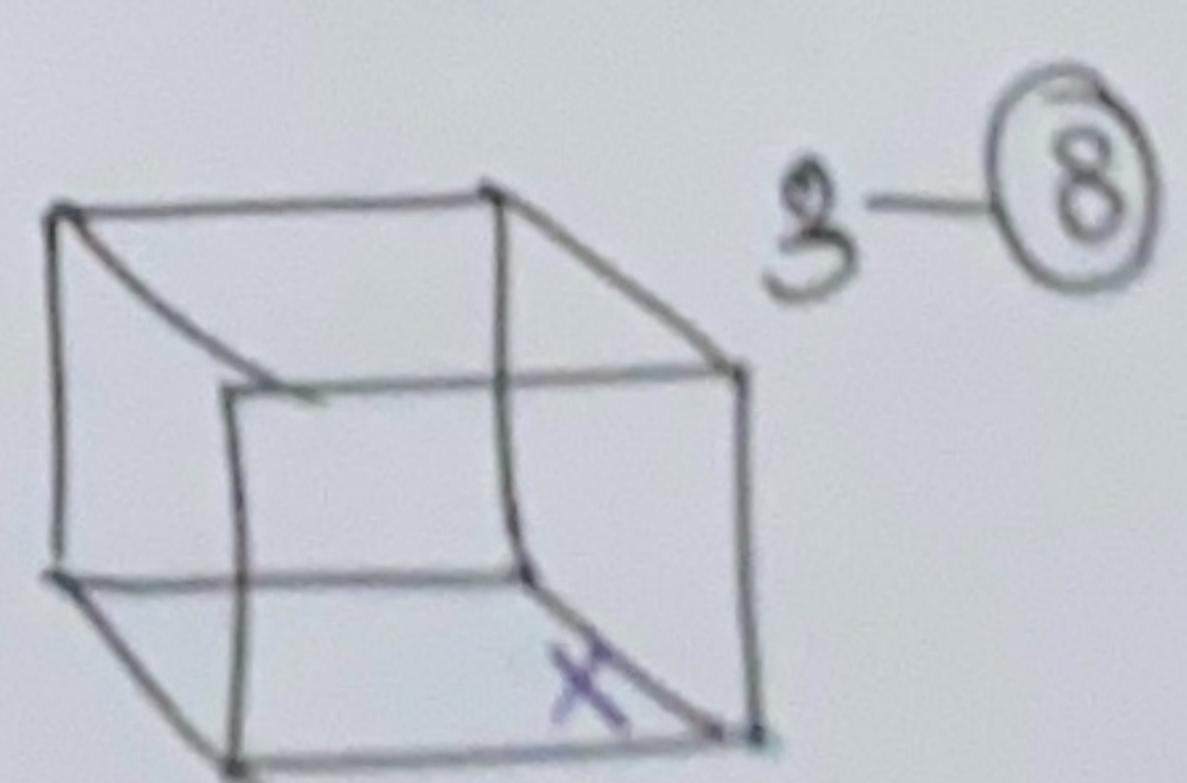
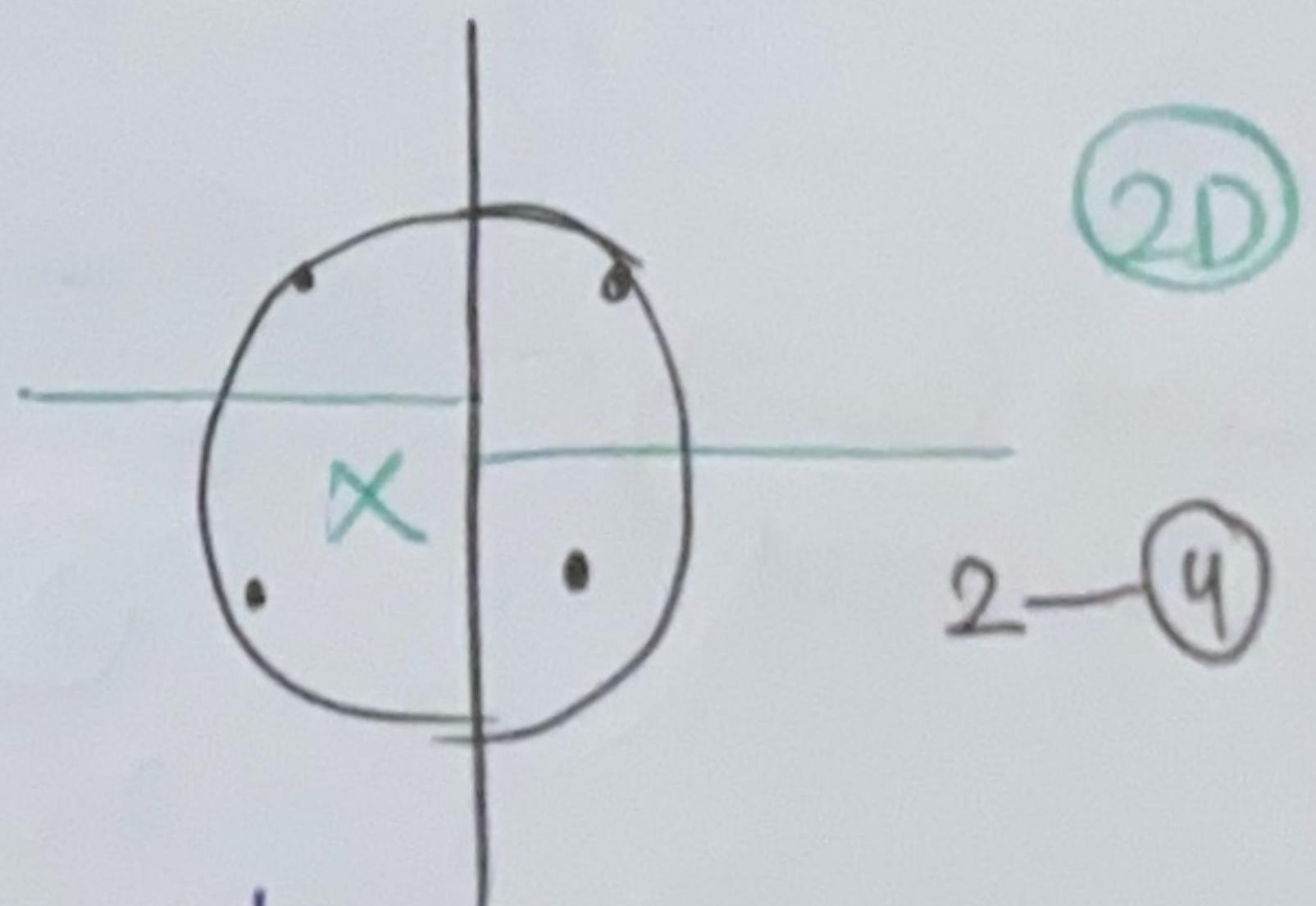
Time complex:  $O(\lg n) \leftarrow 1\text{-NN}$

$O(2^d \lg n) \leftarrow O(n \lg n)$

② When d is small

$O(\lg n)$

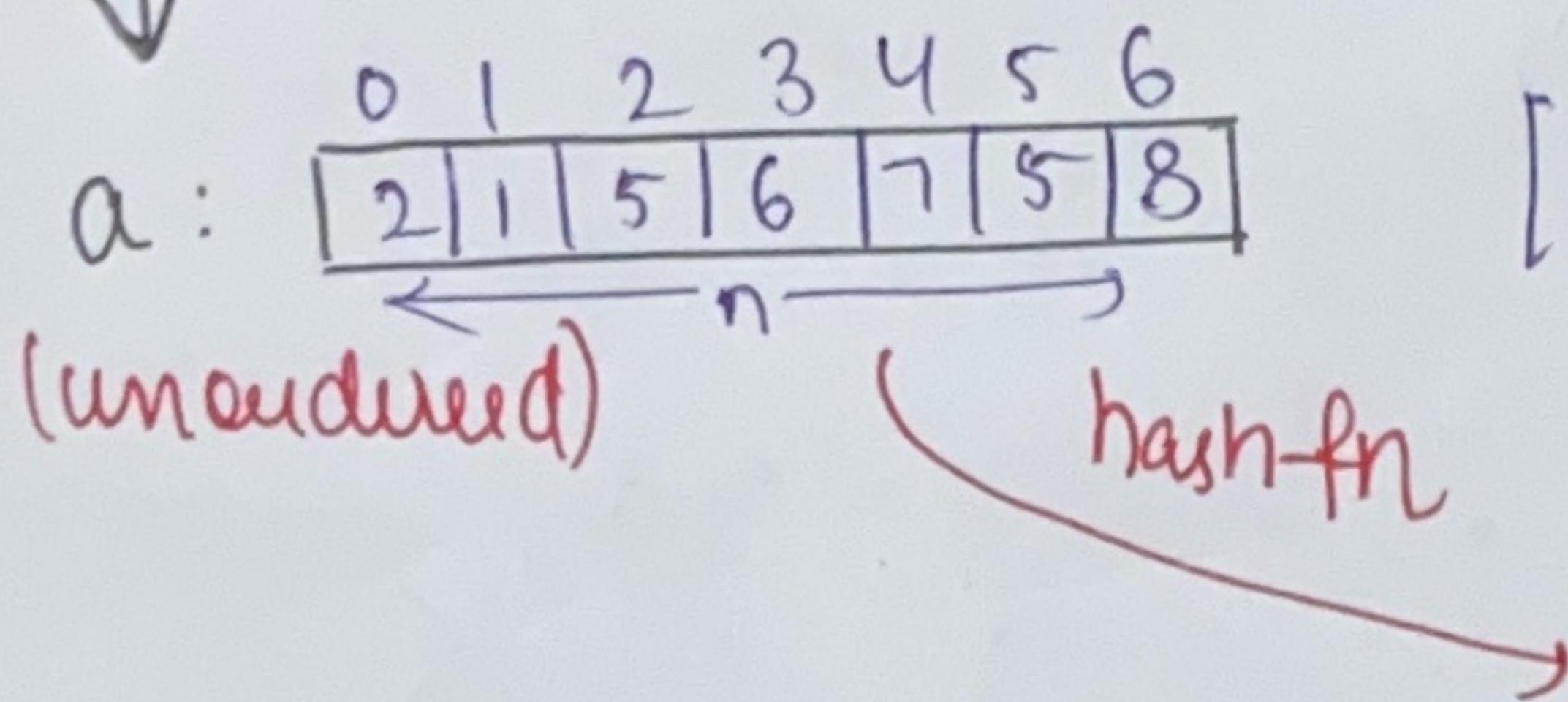
$\hookrightarrow$  only if data is uniformly distributed!



$$d \rightarrow 2^d$$

## Hashing & LSH

locality Sensitive Hashing → (d lauge)



[ hashtable ] (dict)

K	$\sqrt{K, N}$
2	0
5	2, 5
1	1

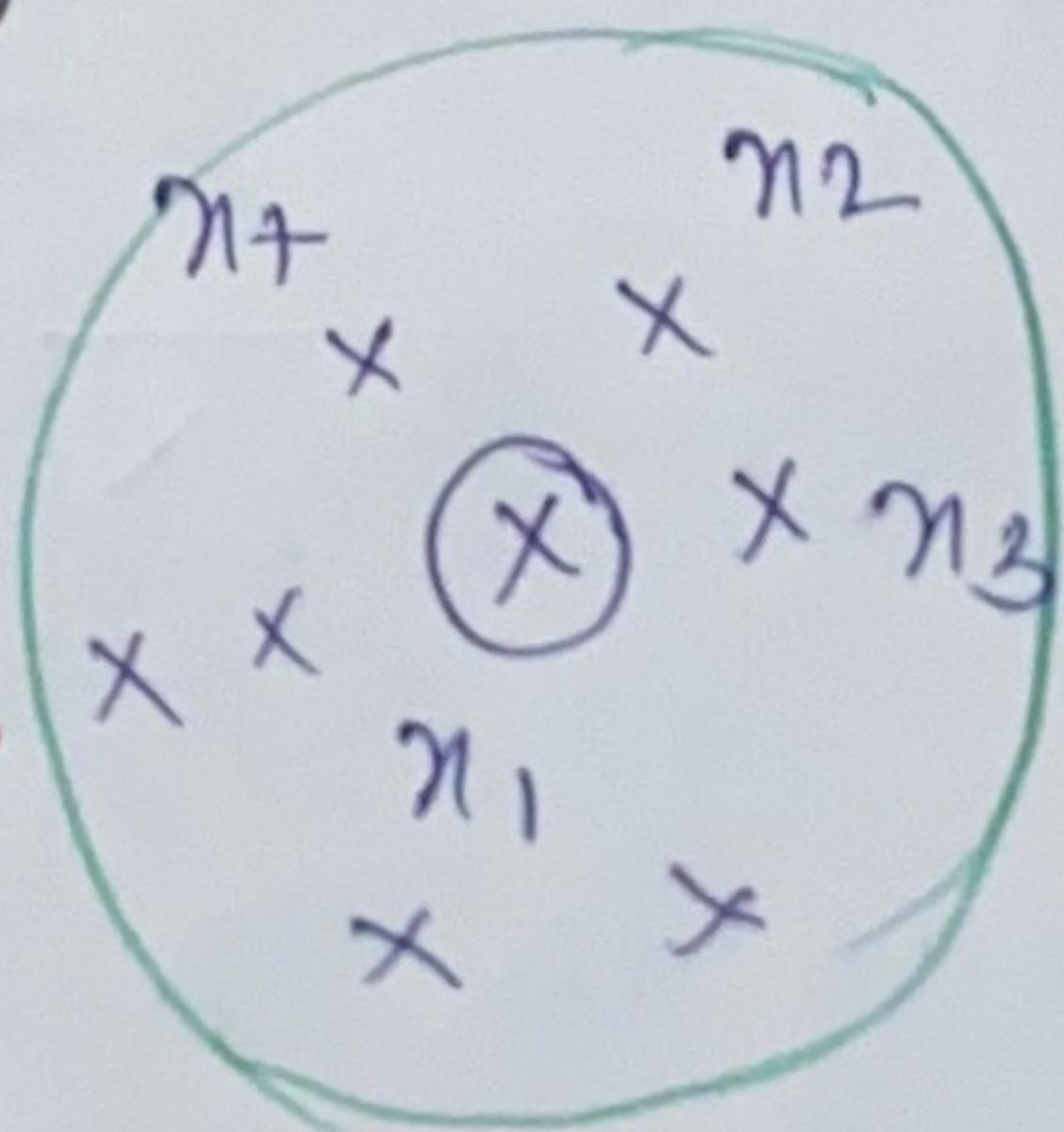
②  $\xrightarrow{h(x)}$  bucket

find if 5 exists in a

$O(n)$ : linear search

## Locality Sensitive Hashing

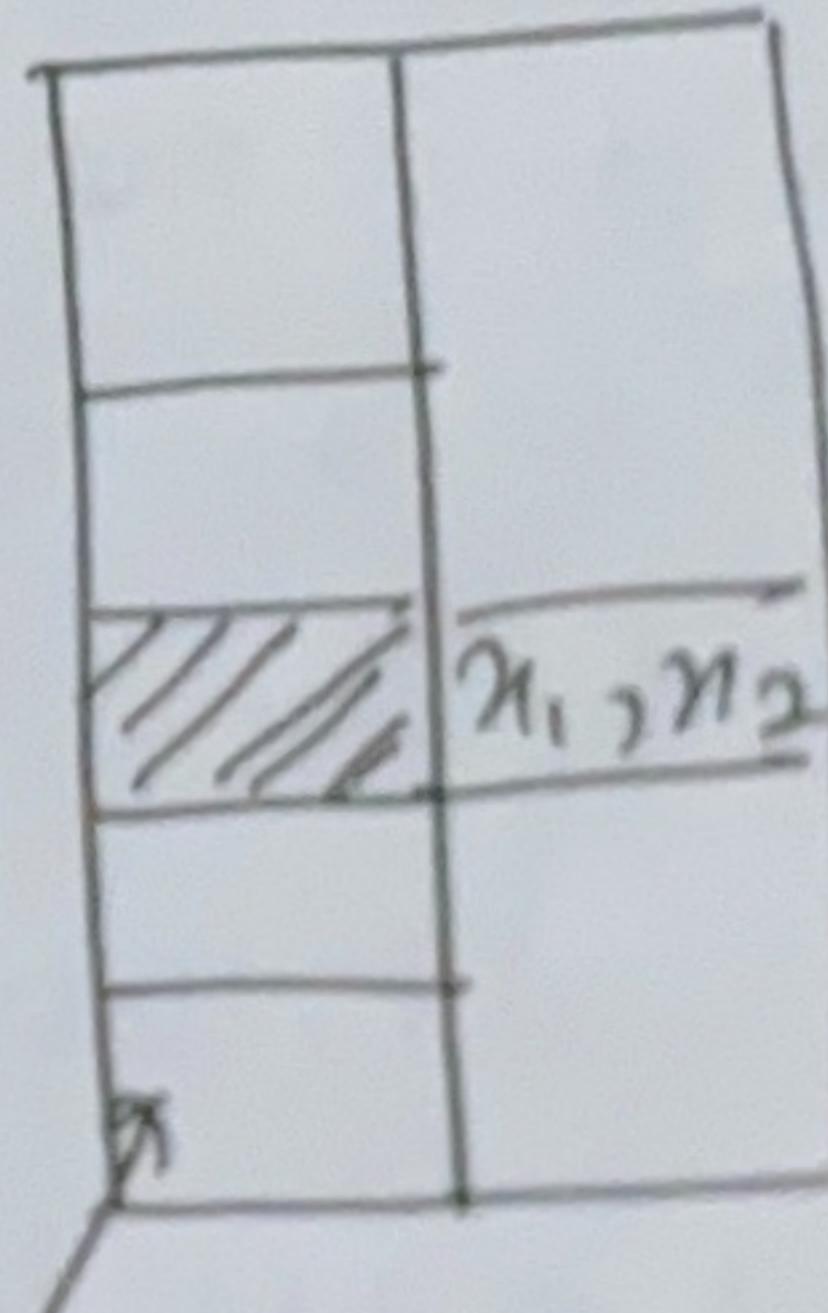
same locality/  
Neighborhood



$h(x)$

$h_m$

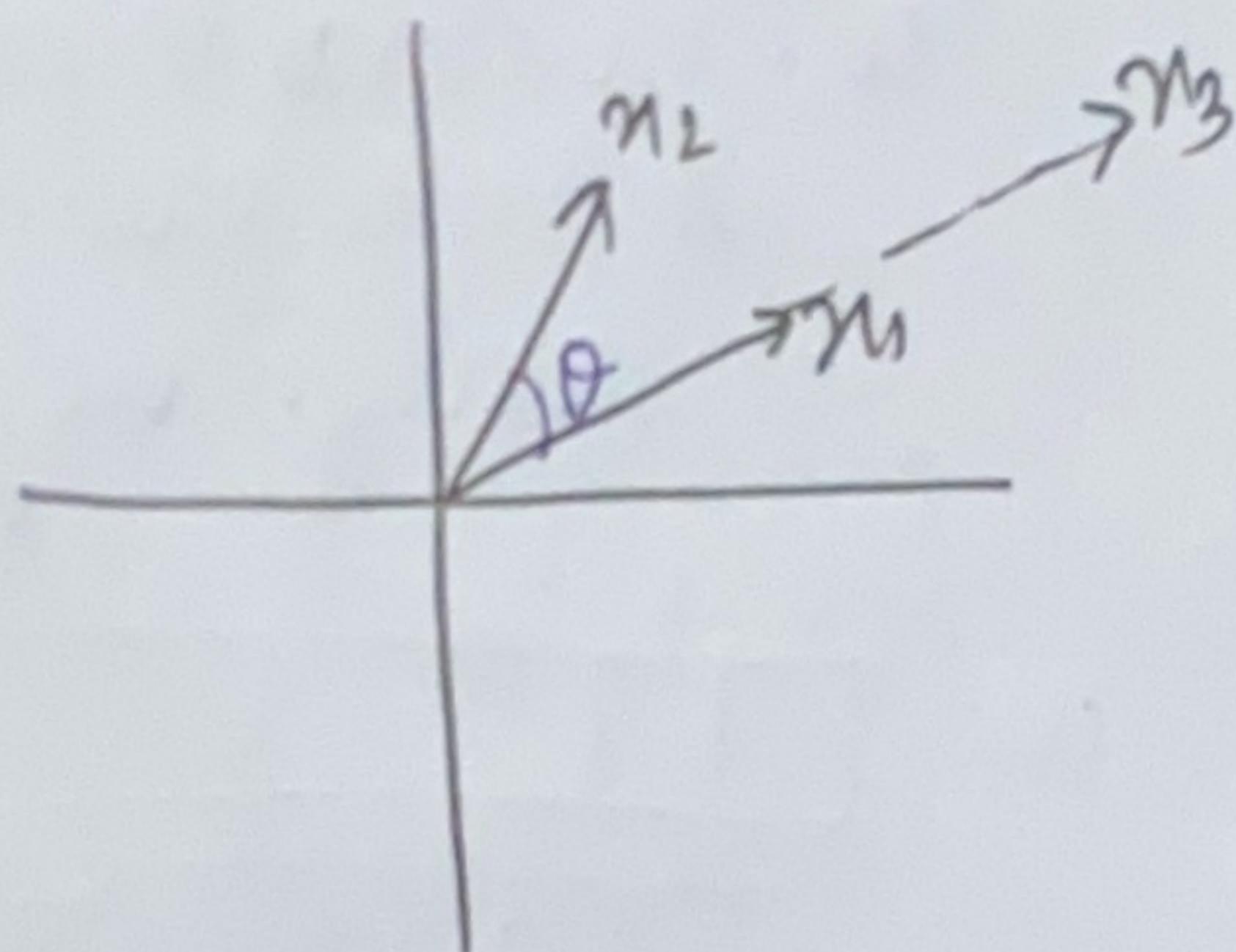
$h(x_q)$



bucket

o)

# Locality Sensitive Hashing for Cosine Similarity



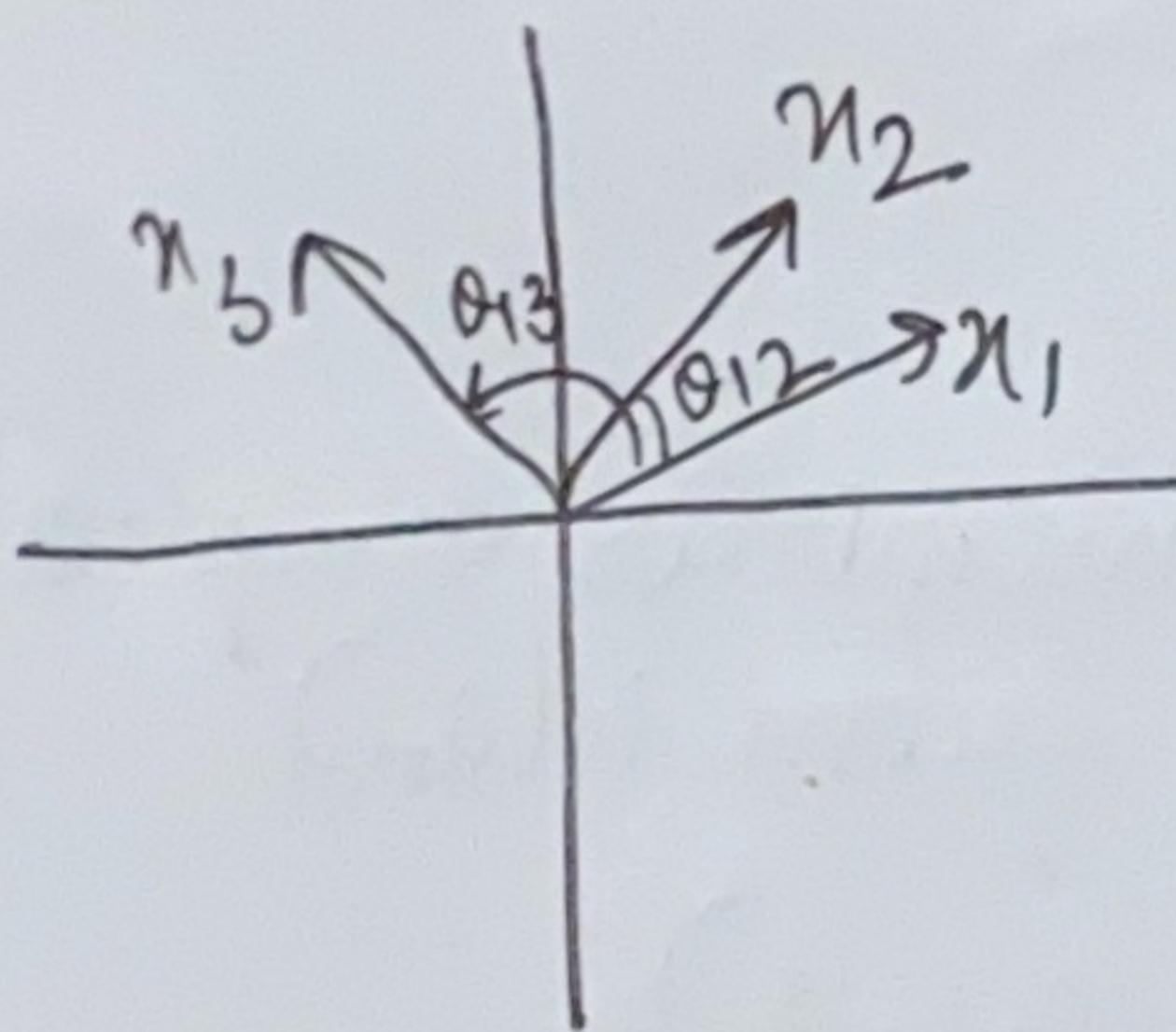
$$\text{cos sim}(n_1, n_2) = \cos \theta$$

$$\text{cos sim}(n_1, n_3) = \cos(\theta) = 1$$

$\theta \uparrow$ ;  $n_1$  &  $n_2$  are angularly separated [dist  $\uparrow$  cosine  $\downarrow$  sim]

Cosine-Sim  $\Leftrightarrow$  angular similarity

LSH for Cos-Sim



$\theta_{12} \vee \theta_{13} \Rightarrow n_1, n_2$  are similar  
↳ small

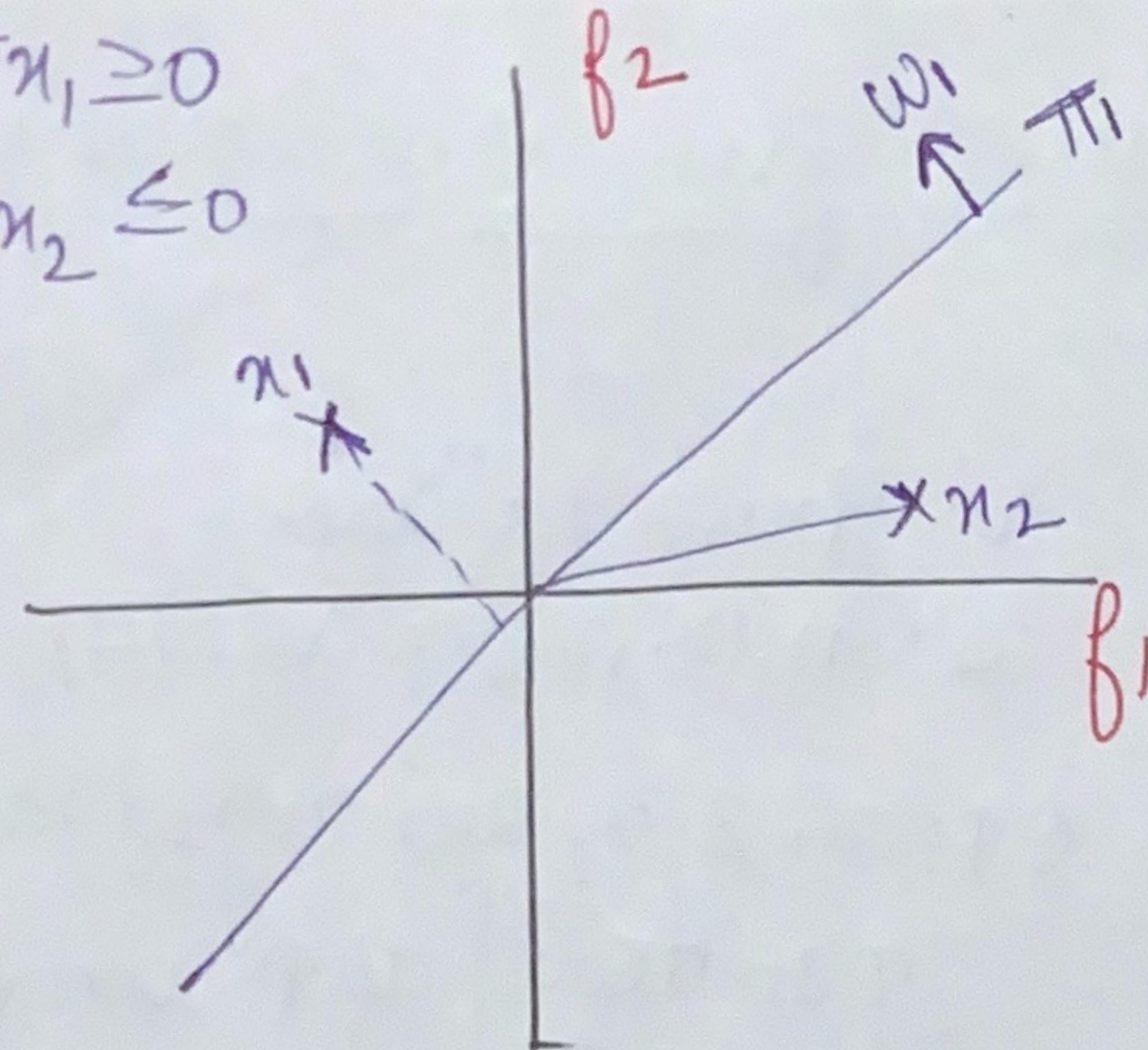
K	V
$h(n_1)$	$n_1, n_2$

$$h(n_1) = h(n_2) = \text{Key}$$

LSH: randomized algorithm

- I will not always give correct ans
- prob. I will give correct ans with high probability

$$\begin{aligned} w_1^T \pi_1 &\geq 0 \\ w_1^T \pi_2 &\leq 0 \end{aligned}$$



① Random-hyperplane

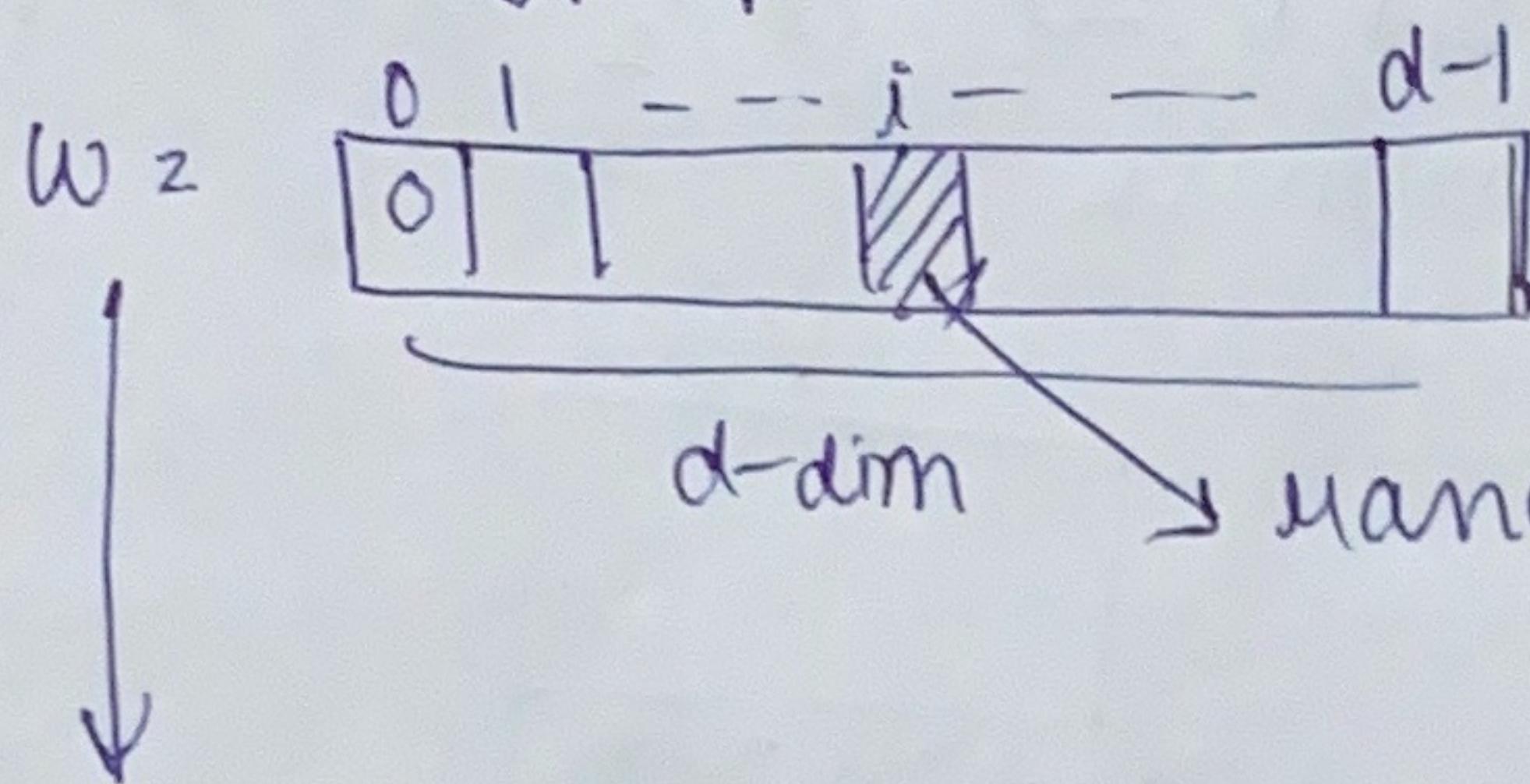
$$w^T x = 0$$

w: normal to plane  
δ plane passing  
through origin

$$w: \begin{array}{|c|c|c|c|c|c|} \hline 0 & 1 & 2 & \cdots & d-1 \\ \hline \end{array} \quad \uparrow \text{d-dim vector random}$$

② 2D  $w_i$ : unit normal vector  
to  $\pi_i$

Random-hyperplane



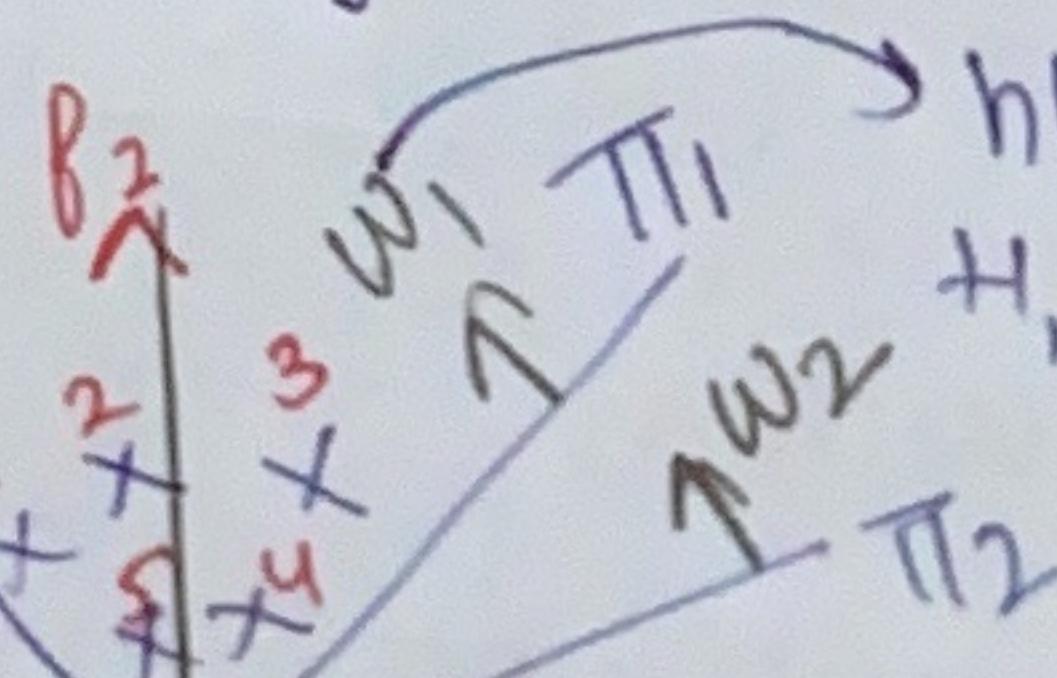
random numbers sampled  
from  $N(0,1)$

$w = \text{numpy.random.normal}(0, 1, d)$

$$\pi_1^T w_1 \geq 0 \rightarrow +1$$

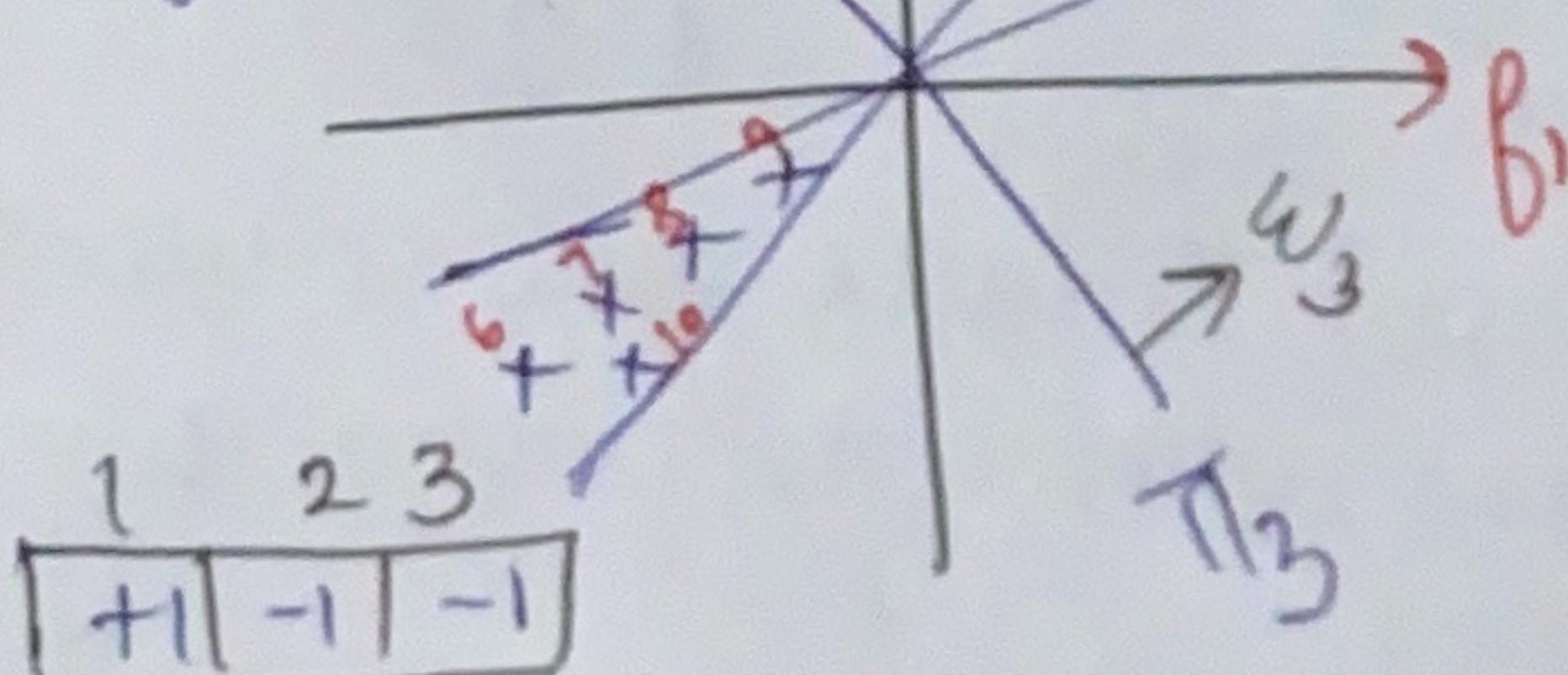
$$\pi_1^T w_2 \geq 0 \rightarrow +1$$

$$\pi_1^T w_3 \geq 0 \rightarrow +1$$



$$h(\pi_2) = +1, +1, +1 \quad n = \# \text{ hyperplanes}$$

Point  $\rightarrow$  slide



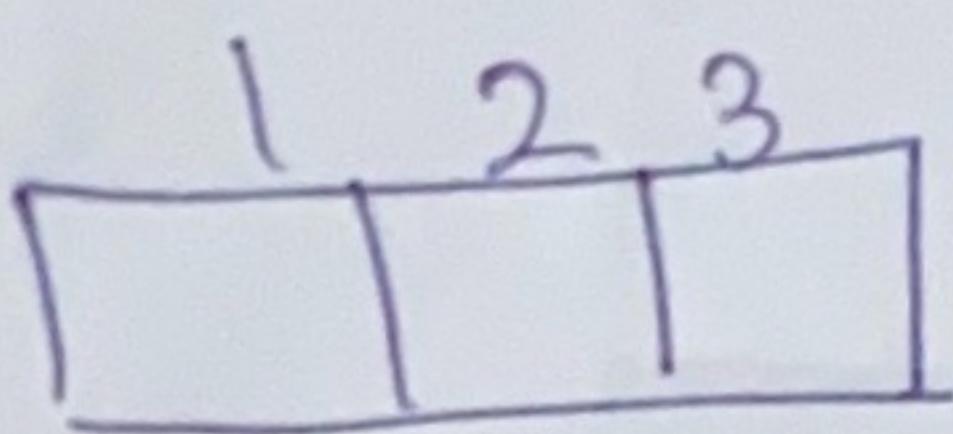
$$h(\pi_7) = (+1, -1, -1)$$

$$h(x) = \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline | & | & | \\ \hline \text{Sign}(w_1 \cdot x) & \text{Sign}(w_2 \cdot x) & \text{Sign}(w_3 \cdot x) \\ \hline \end{array}$$

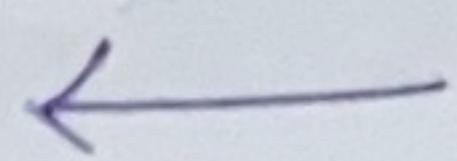
## hash-table

$x \rightarrow h(x)$

K	v
(+1,+1,+1)	$x_1, x_2, x_3, x_4, x_5$



vector of size m



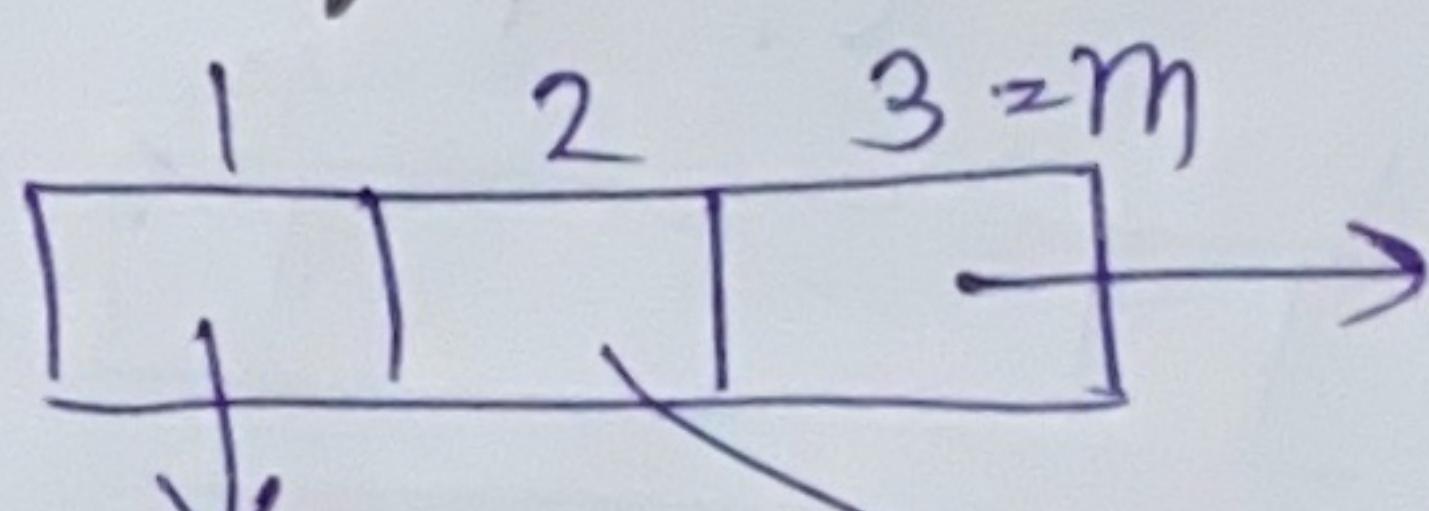
Time:  $O(md)$

$\approx O(mdn)$

Space:  $O(nd)$

Query point ( $x_q$ )

①  $h(x) =$



$\text{Sign}(w^T x_q)$

$\text{Sign}(w_2^T x_q)$

$h(x_q)$  as key in hash-table

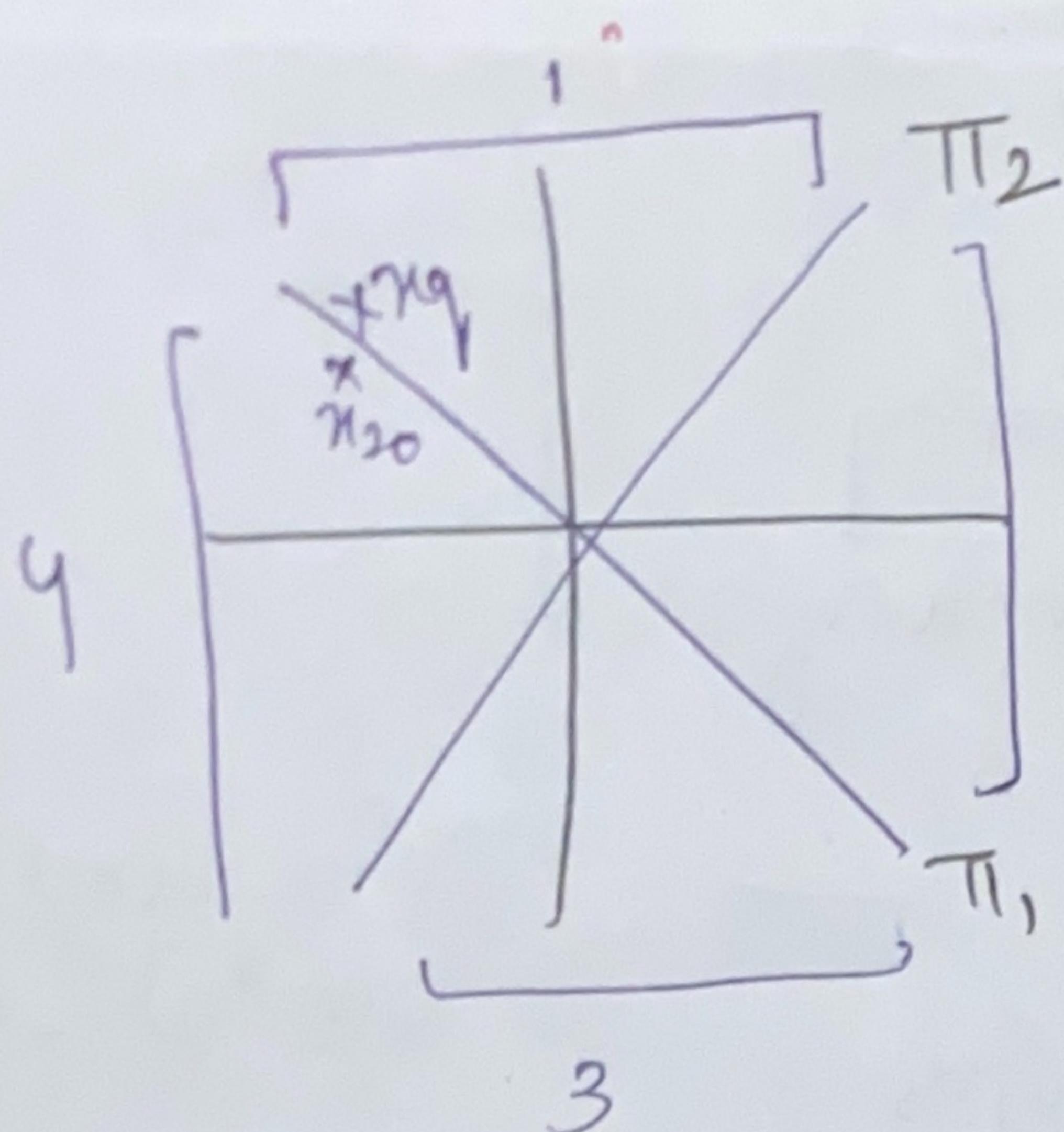
d.get( $h(x_q)$ )  $\rightarrow \{x_1, x_2, x_3, x_4, x_5\}$

Typically

# hyperplanes

$\downarrow$   
 $m \ll \log(n)$

Time:  $O(d * \lg n)$



Cos Sim  
{ $n_{20}$  very close to  $n_q$ )}

\* Could miss a NN  
in Cos Sim as your  
measure

$$\textcircled{1} \quad h_1(n) = \begin{array}{|c|c|c|c|c|} \hline & 1 & 2 & \cdots & m \\ \hline \end{array} \rightarrow \text{hash table}$$

$m$ -hypuplanus

Not same

$$\textcircled{2} \quad h_2(n) = \begin{array}{|c|c|c|c|c|} \hline & 1 & 2 & \cdots & m \\ \hline \end{array} \rightarrow \text{hash table}$$

$m$  hypuplanus

2 hash tables

$n_q \rightarrow h_1(n_q) \rightarrow \text{key} \rightarrow \text{hash table, } \rightarrow \{n_1, n_2, n_3, n_4\}$

$\downarrow$   
 $n_q \rightarrow h_2(n_q) \rightarrow \text{key} \rightarrow H_2 \rightarrow \{n_1, n_2, n_3, n_{20}\}$

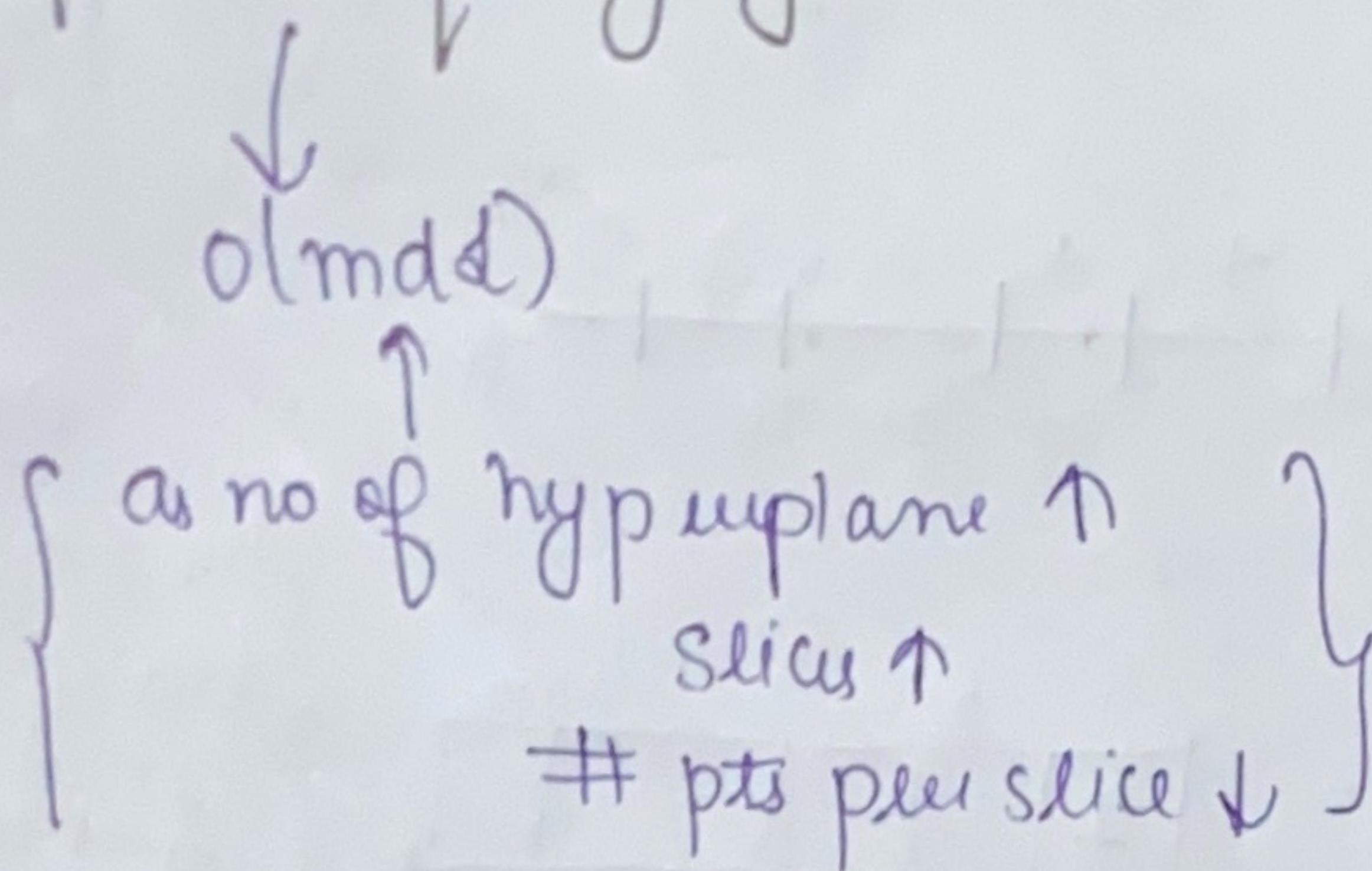
$\downarrow$   
 $h_2(n_2) \rightarrow \dots \{ \quad \}$

$n_q$

$\downarrow$   
 $\{n_1, n_2, n_3, n_4, n_{20}\}$

{ set union }  
{ all points }

Time complex to query given  $\Delta$  hash tables

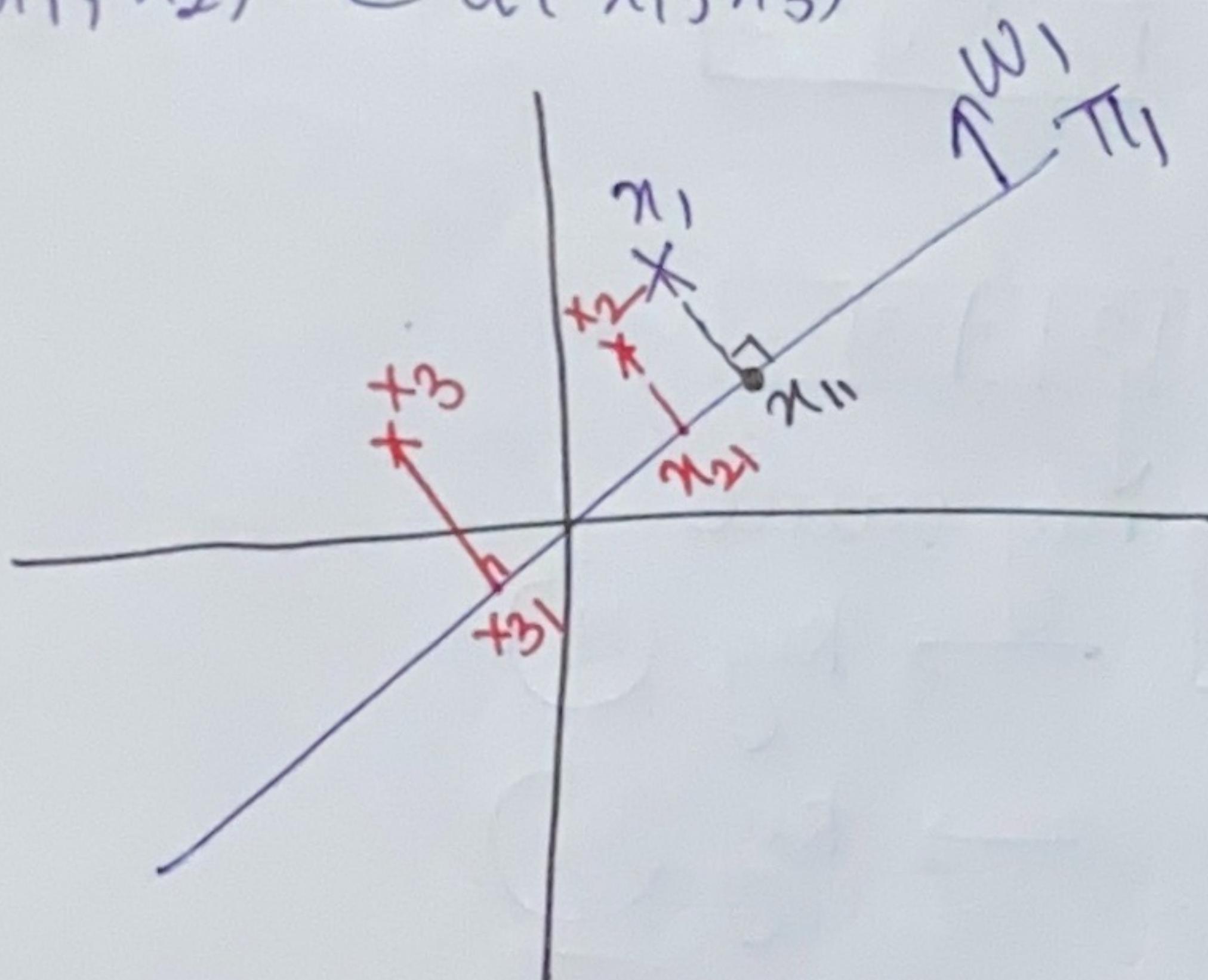


### LSH for Euclidean distance

→ LSH: Cos. Sim

→ eucl dist: Simple extension

$$d(x_1, x_2) \ll d(x_1, x_3)$$



$$h(n) =$$

1	2	3	---	m
---	---	---	-----	---

↓

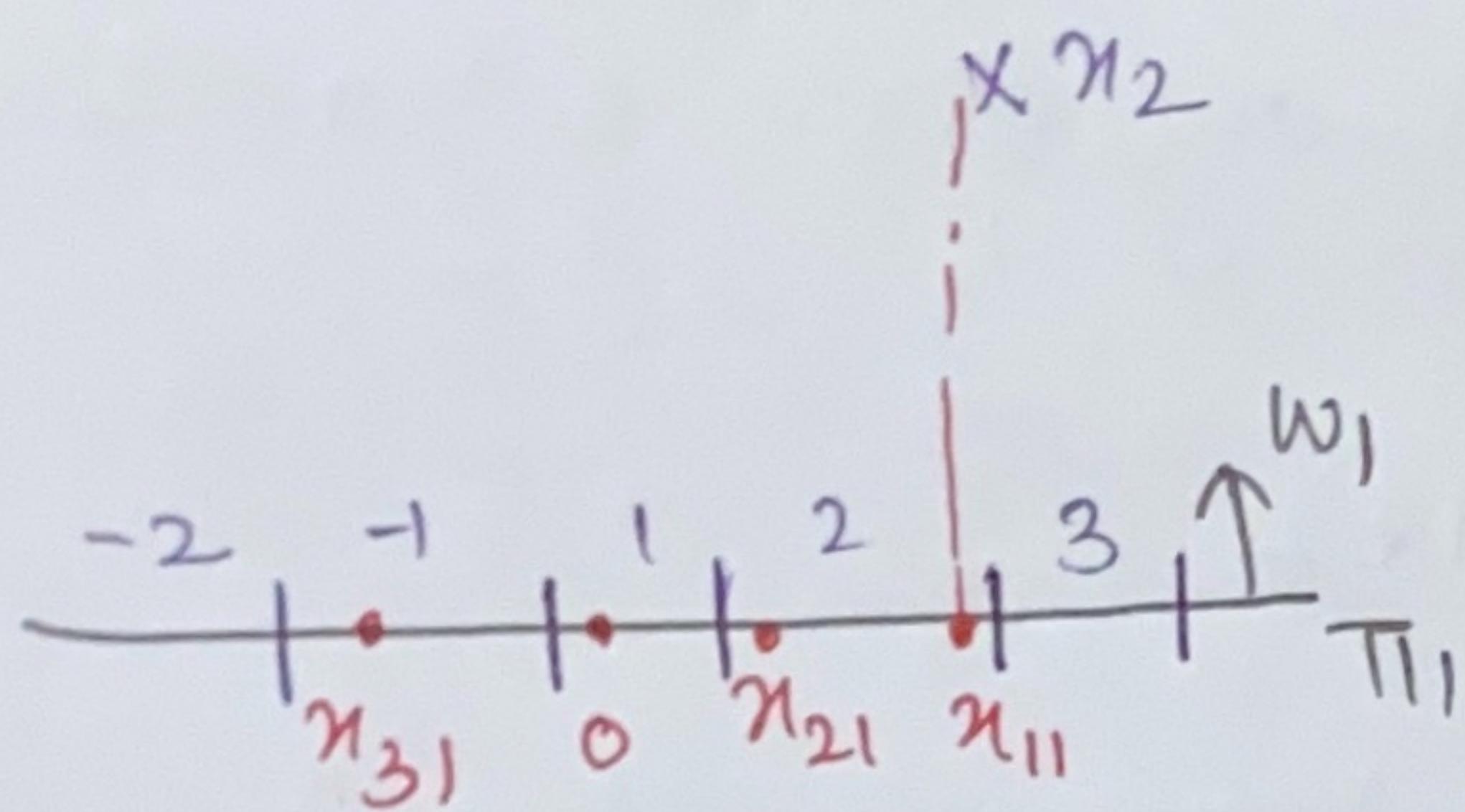
$+1/-1$   
(CosSim)

eucl. dist

$$x_{11} = \text{proj of } x_1 \text{ on } T_1,$$

LSH: Similar / closer points

Should go to same pocket



Breaking plane into pieces / regions  
 ↑ 'a' regions  
 $a = 8$

$$h(n) = \boxed{\begin{array}{cccc|c} 1 & 2 & \cdots & m \\ 2 & | & | & | & | \end{array}}$$

$\pi_1$

$\pi_{11} \& \pi_{21} \rightarrow$  same region  
 but  $\pi_1$  is very far from  $\pi_2$

$$h(\pi_{21}) = \boxed{\begin{array}{cccc|c} 1 & 2 & \cdots & m \\ 2 & | & | & | & | \end{array}}$$

$$h(\pi_{31}) = \boxed{\begin{array}{cccc|c} 1 & 2 & \cdots & m \\ -1 & | & | & | & | \end{array}}$$

Cos-Sim LSH:  $\boxed{+1 \mid -1 \mid \dots}$

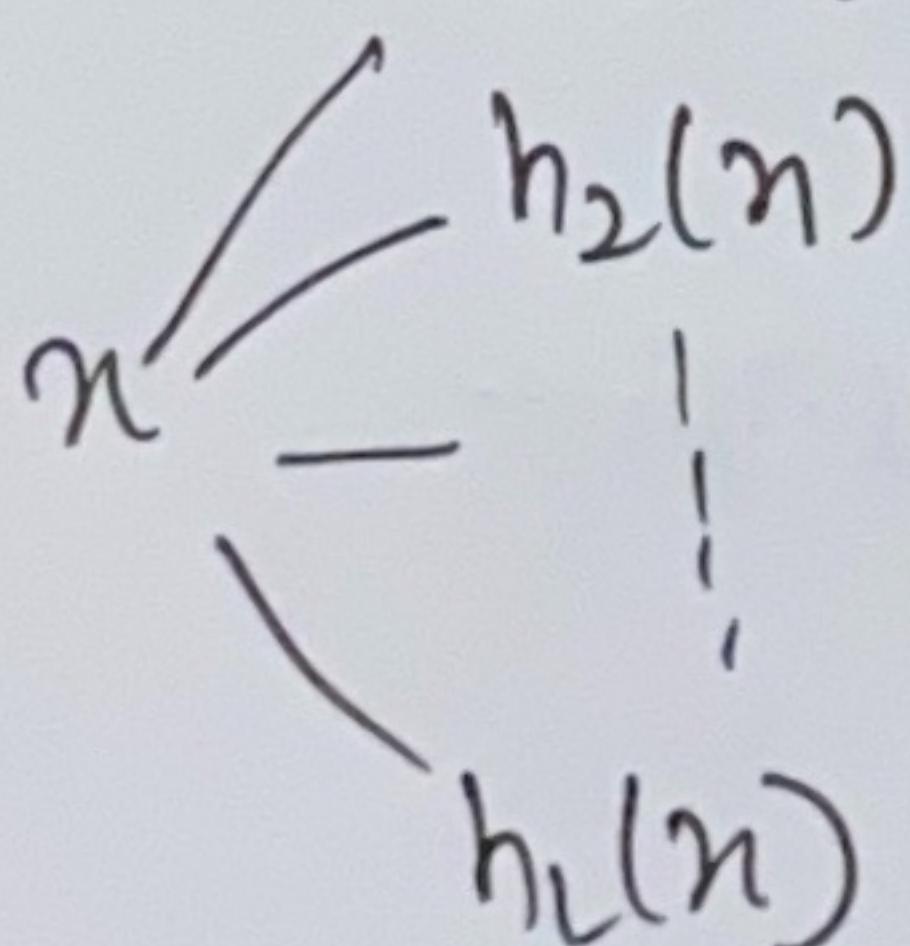
Euc dist.  $\rightarrow$

$$\boxed{1 \ 2 \ \mid -1 \mid -3 \mid \dots}$$

$$h_1(n) : \boxed{\begin{array}{cccc|c} 1 & 2 & 3 & \cdots & m \\ 3 & | & -1 & | & -2 & | & \dots \end{array}}$$

many value  $\rightarrow$  key  $(H_1)$

$\rightarrow$  key  $(H_2)$



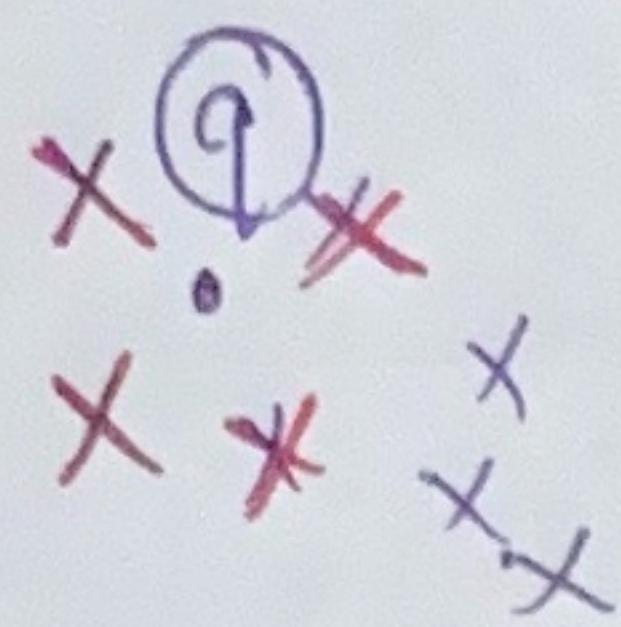
LSH is not perfect

↳ probabilistic/randomized algo

## K-NN: Probabilistic Class Label

2-class classfn

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



7-NN  
( $k=7$ )

Majority Vote

$$y_q = (-w) \cdot (x)$$

$$n_q = \{x, x, x, x, x, x, x\}$$

$n_q$ : 4 (-ve pts); 3 (+ve pts)

$$P(n_q = -ve) = \frac{4}{7} = \frac{\# -ve pts}{\text{Total } \# \text{ pts}}$$

✓ Probabilistic Output (class label)

↳ how certain we are