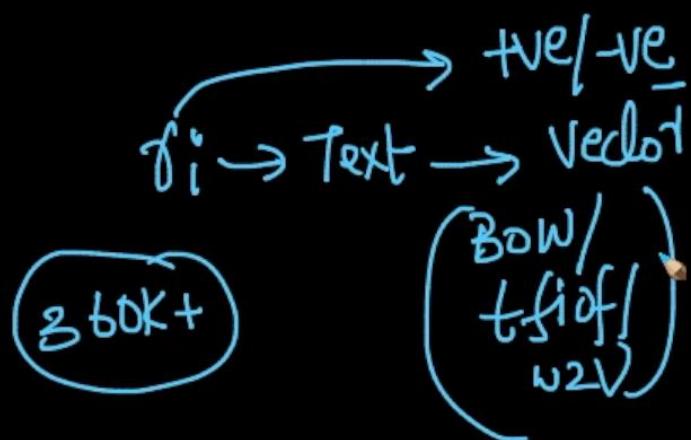


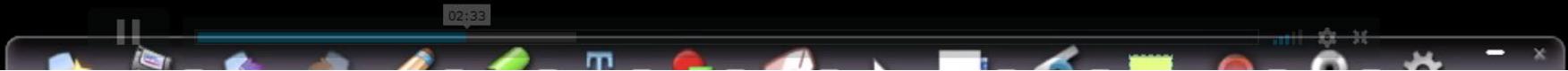
Classification & Regression : (K-NN)

How "classification" works



{ Amazon Food reviews }
↓
real-world eg.

{ MNIST (PCA/tSNE)
[0, 1, 2, ..., 9]



Classification :-

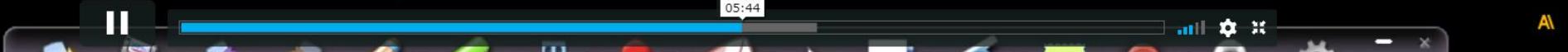
$364K \quad r_i \rightarrow (x_i, +ve/-ve)$

$y_i \rightarrow +ve/-ve$
classifying
 \uparrow \uparrow

{ given a new review-text, determine/predict
if the review is +ve or -ve

$$y = f(x) \rightarrow \begin{array}{l} \text{central concept} \\ \text{review-text} \end{array}$$

+ve/-ve

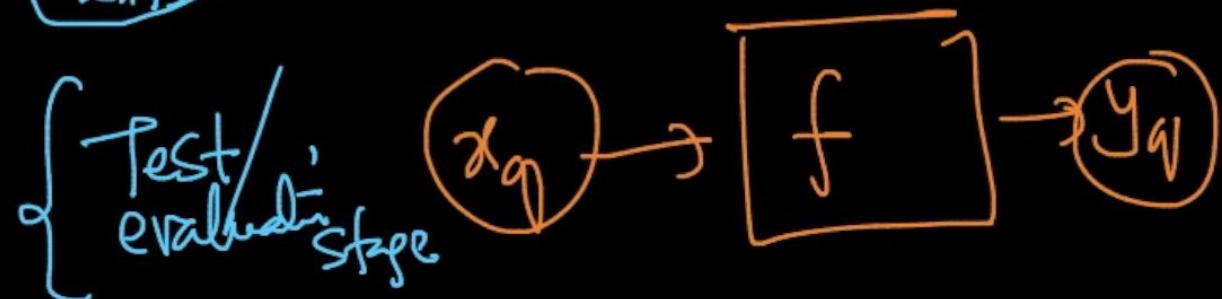
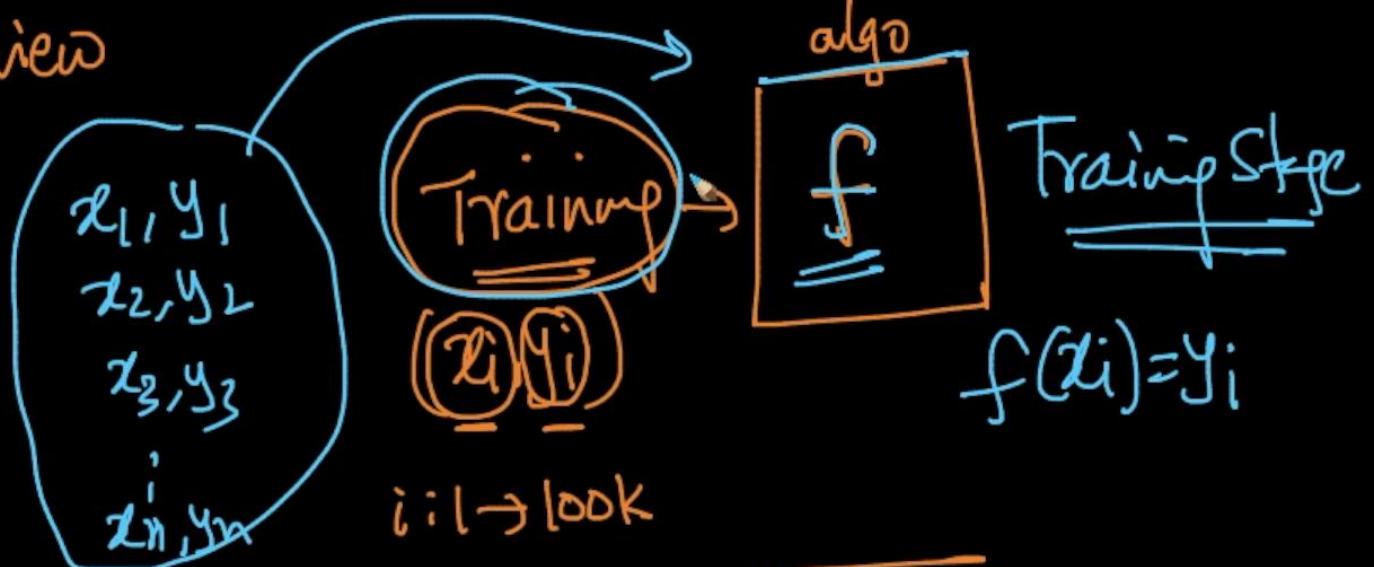


few classes

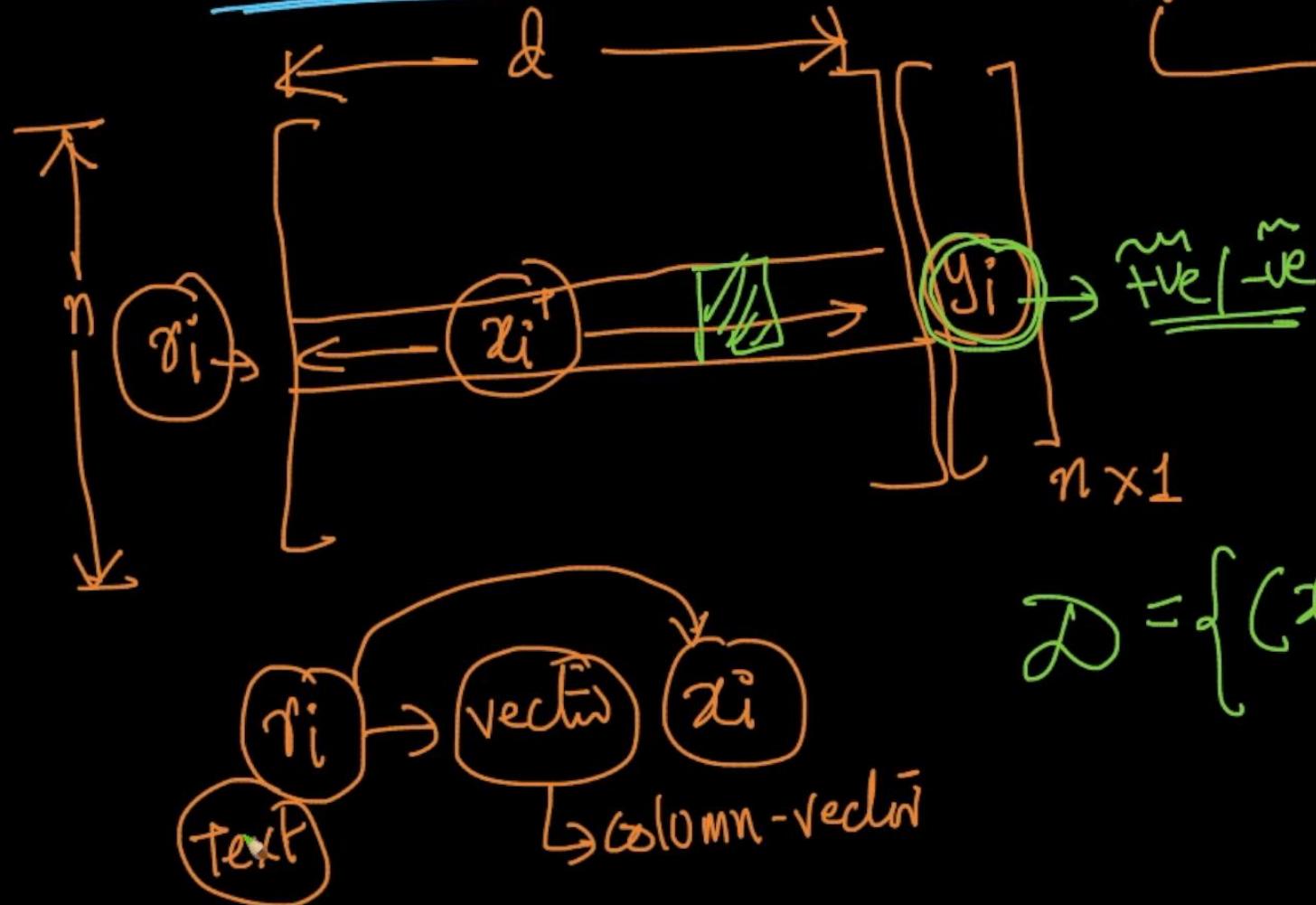
$$y_{qj} = f(x_{qj})$$

query-review

+ve/-ve



Notation:



$$\mathcal{D} = \left\{ (x_i, y_i) \right\}_{i=1}^n \quad \begin{array}{|l} x_i \in \mathbb{R}^d \\ y_i \in \{0, 1\} \end{array}$$

↑
+ve
↑
-ve

\mathcal{D} $\xrightarrow{\text{set}}$ $\left\{ (x_i, y_i) \right\}_{i=1}^n$ such that $x_i \in \mathbb{R}^d$, $y_i \in \{0, 1\}$

\downarrow

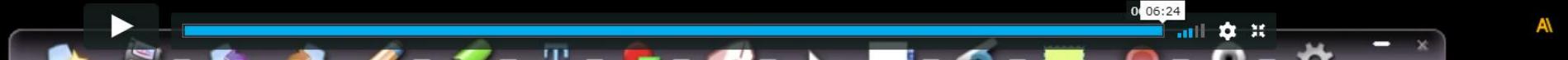
y_i $\xrightarrow{\text{vec}}$ x_i

$\hookrightarrow \{0, 1\} (y_i)$

$\xrightarrow{-ve} \xrightarrow{+ve}$

classification:

$$f(x_i) = y_i$$



Classification vs Regression:

$$\mathcal{D} = \left\{ (\underline{x}_i, y_i) \right\}_{i=1}^n \quad | \quad \underline{x}_i \in \mathbb{R}^d, y_i \in \{0, 1\} \}$$

$y_i \in \{0, 1\}$ ↓
↓
-ve +ve

→ 2 classes

Amazon Food reviews

→ 2 class - classification / binary

MNIST: $y_i \in \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\} \rightarrow$ 10-class / Multi-class
classification

K - Nearest Neighbors: (Geometric)

2D - toy dataset:



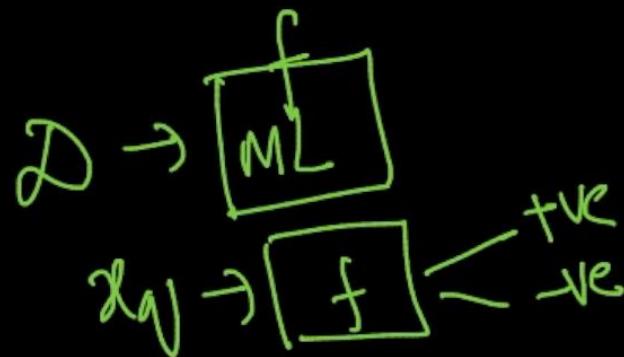
geom.: closest to x_q :
conclude $x_q \rightarrow$ blue (+ve)

$x_q \rightarrow y_q$ (binary classification)

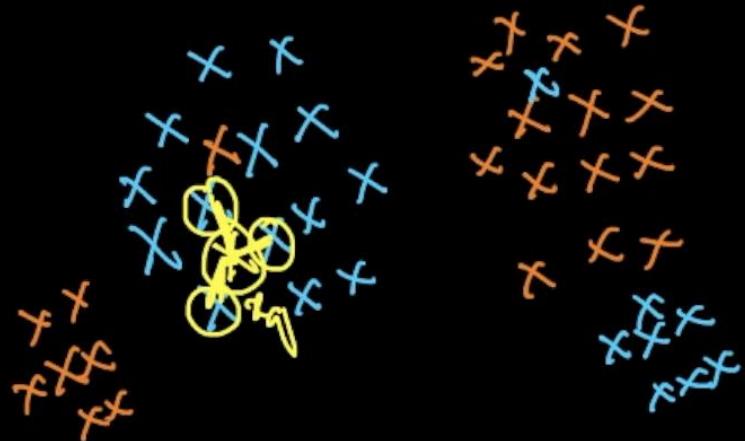
$\left\{ \begin{array}{l} x: +ve \text{ datapoint} \\ x: -ve \text{ class datapoint} \end{array} \right.$

$$\mathcal{D} = \{(x_i, y_i) \mid x_i \in \mathbb{R}^2, y_i \in \{0, 1\}\}$$

↓
-ve +ve



D
 $= \{(x_i, y_i)\}$



given $x_q \rightarrow$ find $y_q \in \{0, 1\}$
 ve +ve

① find K
 let $K=3$

"nearest" points to x_q in D
 $\{x_1, x_2, x_3\} \rightarrow$ 3 NN to x_q
 $\downarrow \quad \downarrow \quad \downarrow$
 $y_1 \quad y_2 \quad y_3$

② $\{y_1, y_2, y_3\} \rightarrow$ Majority vote

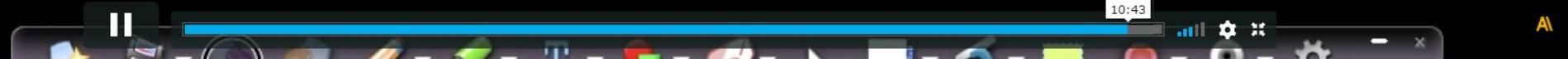
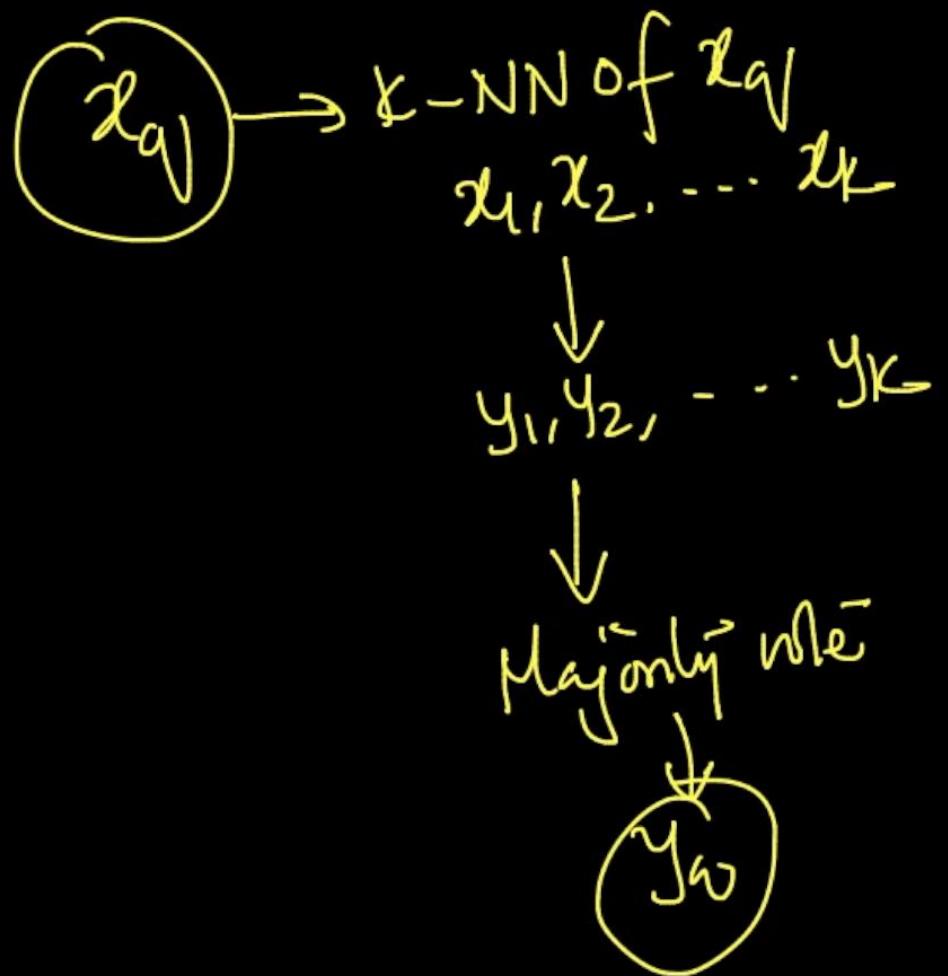


$k = \text{odd}$

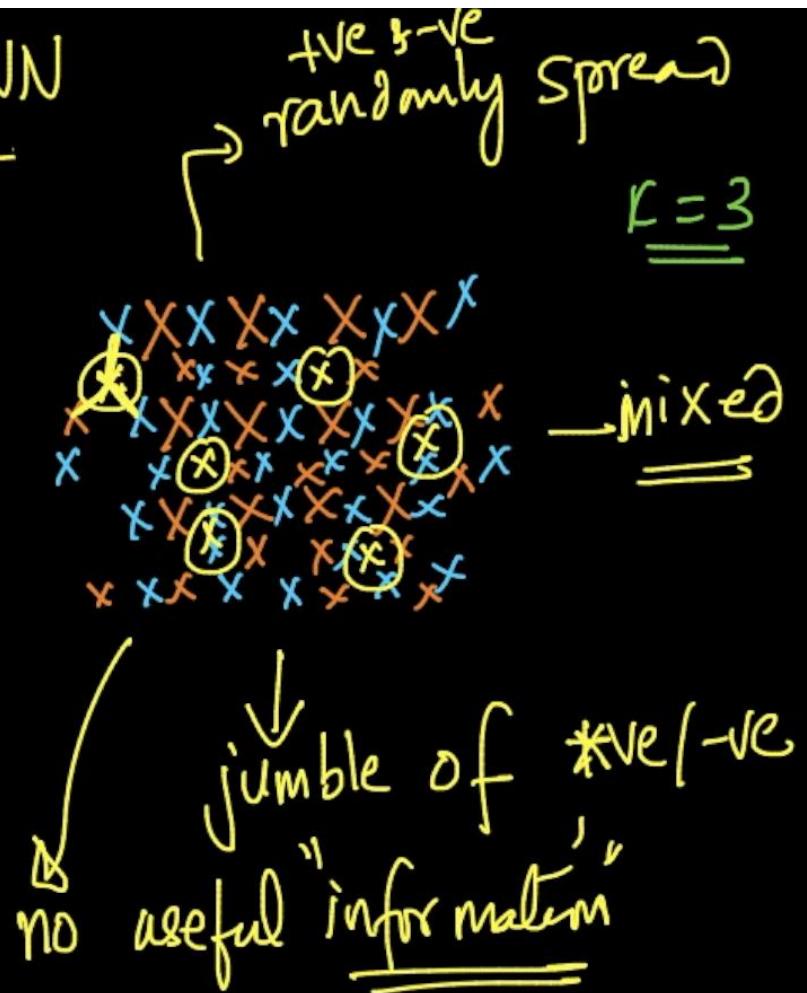
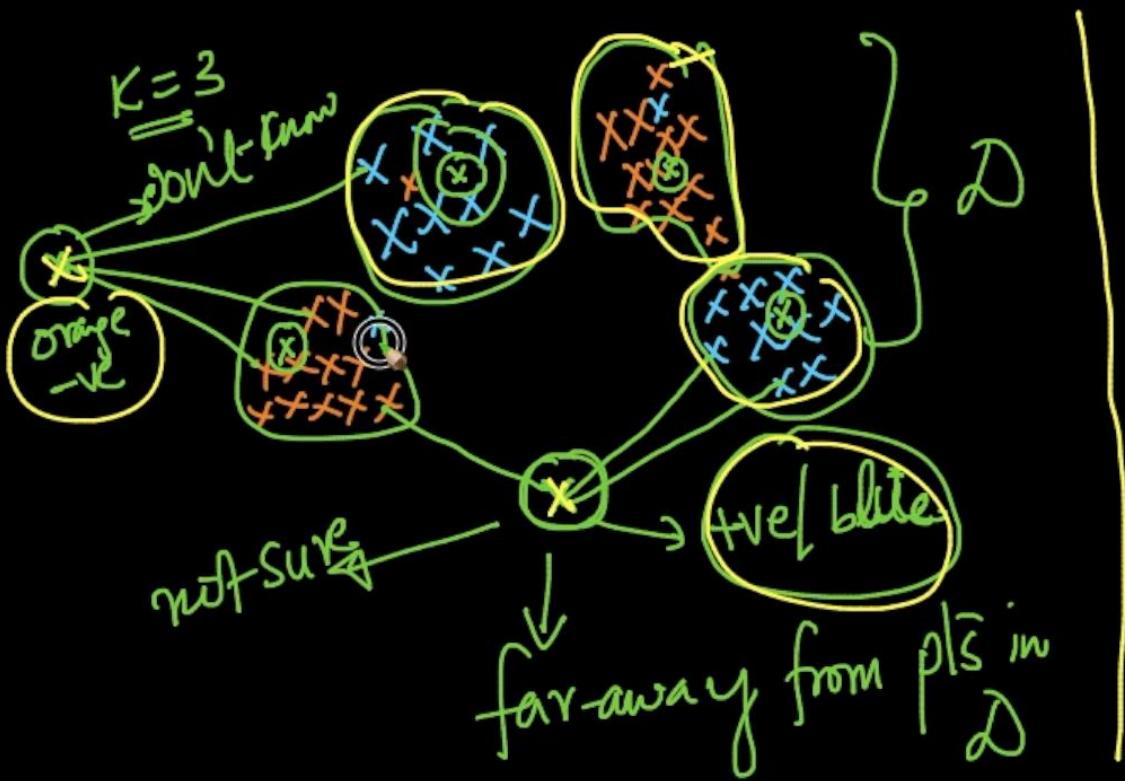
$$\begin{array}{c} \{y_1, y_2, y_3\} \\ +, +, + \\ +, +, - \end{array} \rightarrow \begin{array}{c} \text{Majority rule} \\ (+) \end{array} \rightarrow \begin{array}{l} y_{qj} = \text{+ve (or) blue} \\ y_{qj} = \text{+ve (or) blue} \end{array}$$

~~(X)~~ if $k = 4$
 $+, +, -, -$ $\rightarrow (+) \text{ or } (-)$





"Failure" Case of K-NN



$$f(x_0) = y_0$$

Train

noisy outliers

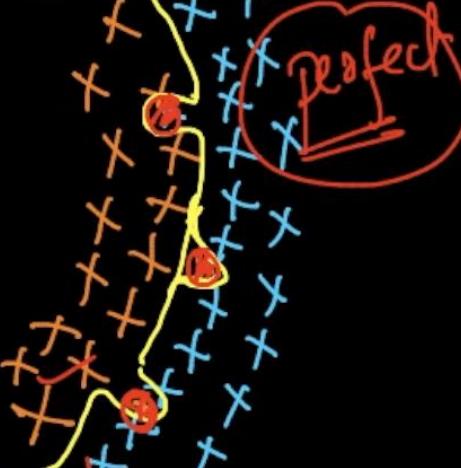
Over-fitting

over-work

$K=1 \rightarrow$
✓ no-mistakes

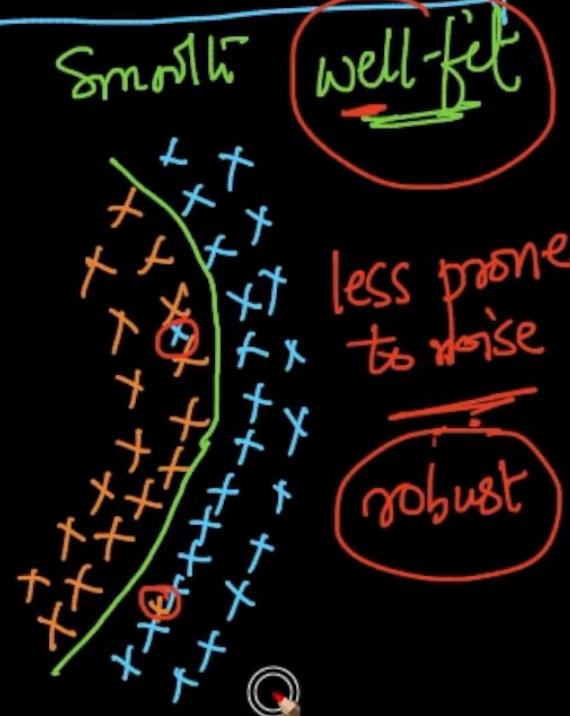
'Over' fitting

non-smooth



& 'under' fitting: (fundamental)

Smooth



well-fit

less prone
to noise

robust

\Rightarrow

(Under-fitting)
imperfect
every query point
belongs to the
majority class

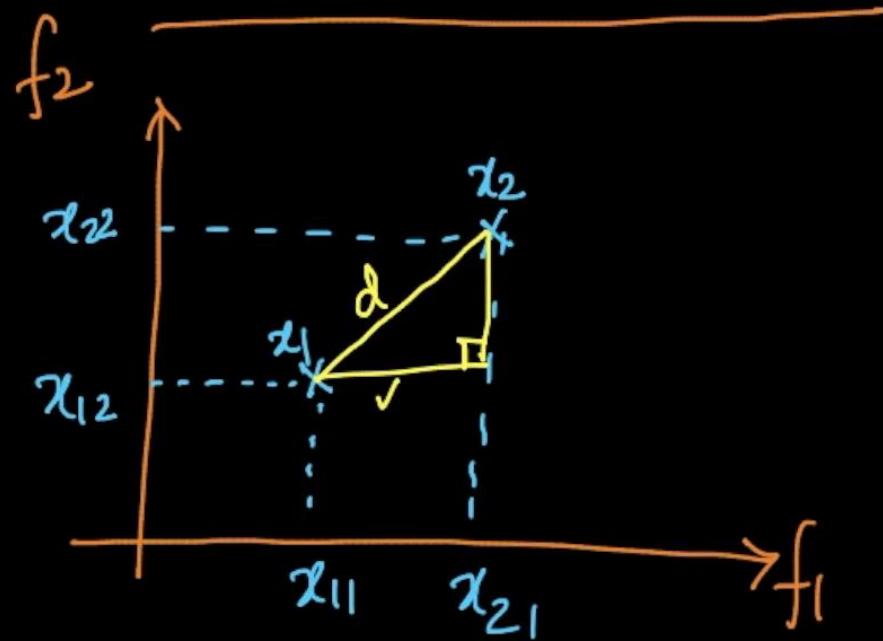
+ve
 $n_1 > n_2$

under-work

$K=n$ (lazy)



"Distance" Measures:



$$x_1 = (x_{11}, x_{12})$$
$$x_2 = (x_{21}, x_{22})$$

d = len of shortest line from x_1 to x_2

Euclidean dist $\leftarrow d = \sqrt{(x_{21} - x_{11})^2 + (x_{22} - x_{12})^2} = \|x_1 - x_2\|$



$x_1 \in \mathbb{R}^d, x_2 \in \mathbb{R}^d$

Euc-dist: $\|x_1 - x_2\|_2 = \left(\sum_{i=1}^d (x_{1i} - x_{2i})^2 \right)^{1/2}$

$\|(x_1 - x_2)\|_2 \rightarrow L_2 \text{ norm}$

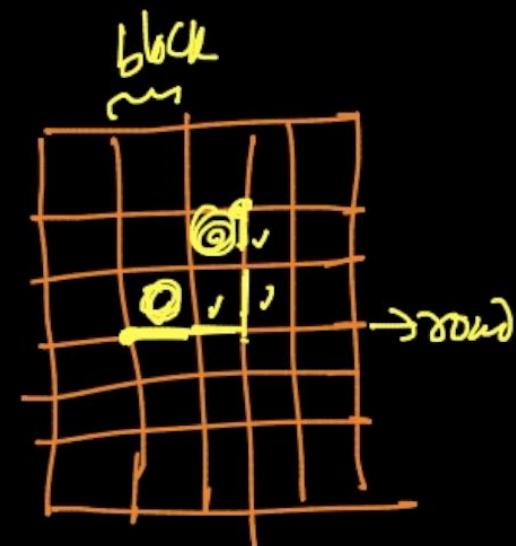
$\|x_1\|_2 = \text{dist of } x_1 \text{ from origin} = \left(\sum_{i=1}^d x_{1i}^2 \right)^{1/2}$

✓ "Manhattan" dist:

$$\sum_{i=1}^d \left| x_{1i} - x_{2i} \right|$$

L₁-norm of vector $(x_1 - x_2)$

$$\|x_1\|_1 = \sum_{i=1}^d \left| x_{1i} \right|$$

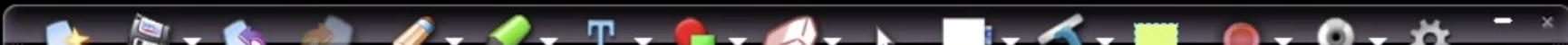


L_p -norms \nrightarrow Minkowski dist

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

= L_p -norm of $(x_1 - x_2)$

$p=2 \rightarrow$ Minkowski dist \rightarrow Eucl. dist
 $p=1 \rightarrow$ " \rightarrow Manhattan dist



$$\text{L}^p \text{ norm} \rightarrow \|x_i\|_p = \left(\sum_{i=1}^d |x_{ii}|^p \right)^{\frac{1}{p}} \quad p \neq 0$$

$\begin{cases} \text{dist}(x_1, x_2) \\ \text{norm} \rightarrow \text{vector} \end{cases}$

distances → ID's of types
of dist



✓ Hamming dist (boolean Vectr)

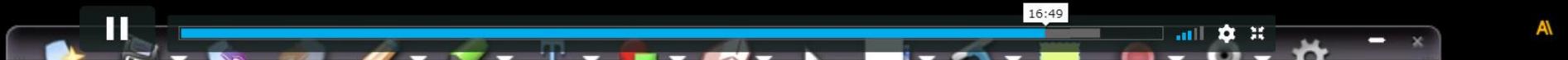
$x_1, x_2 \rightarrow$ boolean Vectr \rightarrow Binary BOW

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

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

Hamming-dist(x_1, x_2) = # locations where
binary vectors differ

↳ 3



Cosine-Similarity & Cosine-distance:

x_1, x_2

Similarity

distance

↑ inc
↓ dec

(opposite)

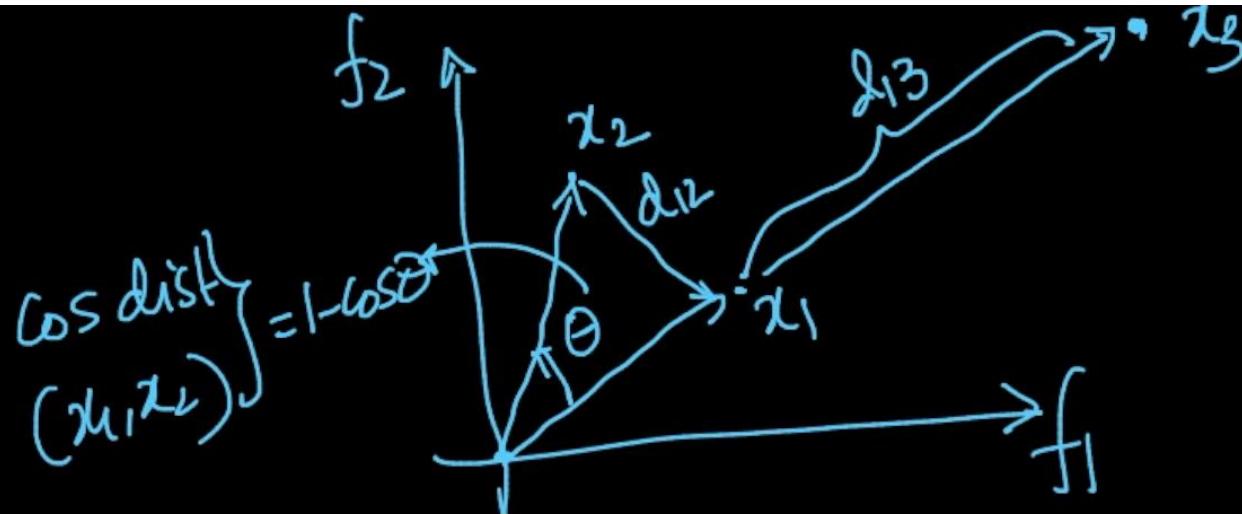
let

$$\text{v.-simila} \cos\text{-sim}(x_1, x_2) = (+1)$$

$$\text{v.-dissimila} \cos\text{-sim}(x_1, x_2) = (-1)$$

$$1 - \cos\text{-sim}(x_1, x_2) = \cos\text{-dist}(x_1, x_2)$$

let $[-1, 1]$



$$\cos\text{dist}(x_1, x_2) = 1 - \cos\theta$$

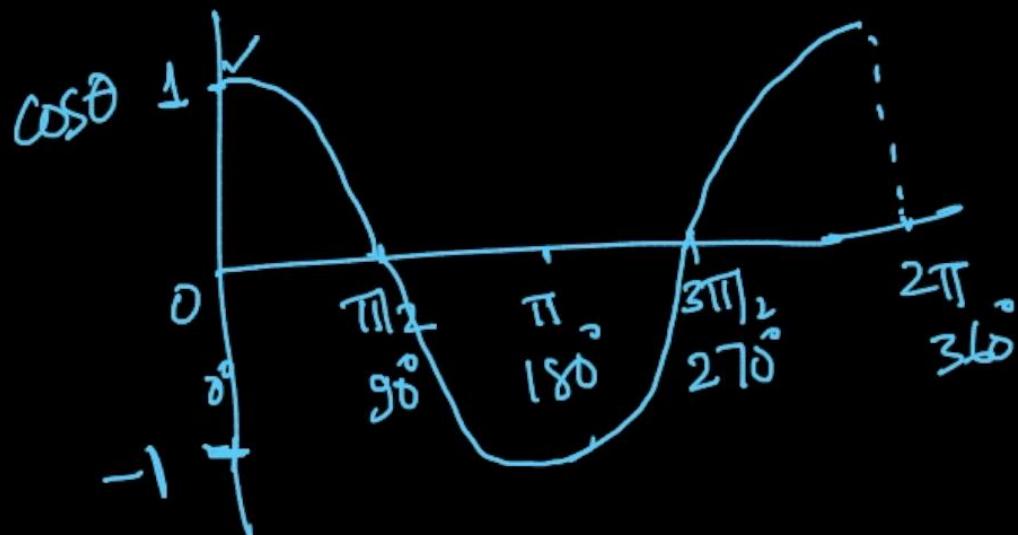
$$\cos\text{-dist}(x_1, x_3) = 1 - 1 = 0$$

$$\cos\text{-sim}(x_1, x_2) = \cos\theta$$

$$\cos\text{-sim}(x_1, x_3) = 1$$

$$\Theta_{x_1, x_3} = 0^\circ$$

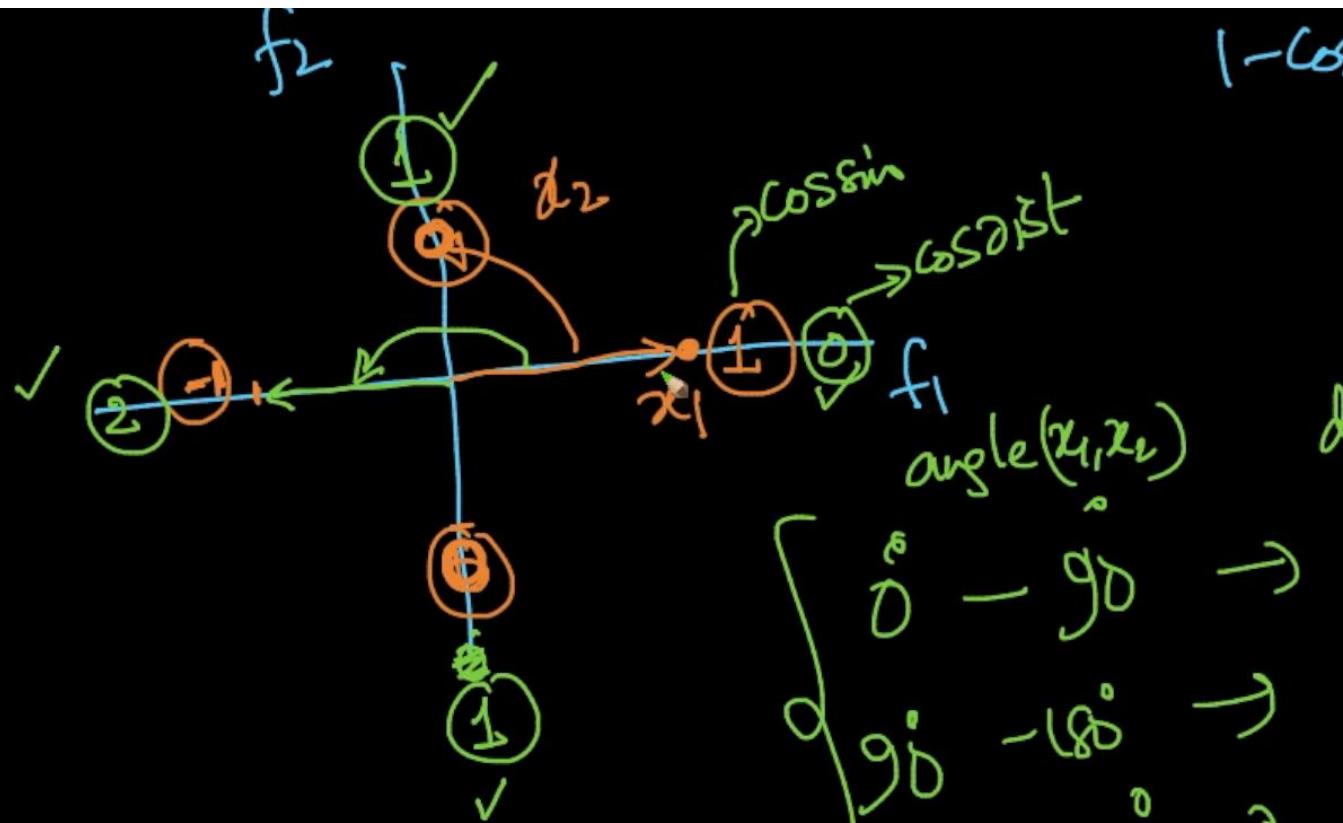
CCW-OK



$$\left\{ \begin{array}{l} d_{13} > d_{12} \\ \cos\text{-dist}_{13} < \cos\text{-dist}_{12} \end{array} \right.$$



$$1 - \text{cos sim}(x_1, x_2) = \text{cos dist}(x_1, x_2)$$



dist

- $0^\circ - 90^\circ \rightarrow 0 \rightarrow 1$
- $90^\circ - 180^\circ \rightarrow 1 \rightarrow 2$
- $180^\circ - 270^\circ \rightarrow 2 \rightarrow 1$
- $270^\circ - 360^\circ \rightarrow 1 \rightarrow 0$

x_1, x_2

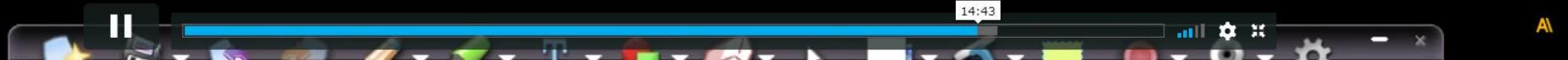
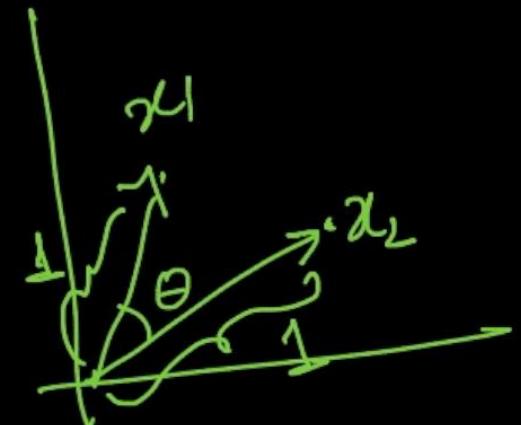
$$\cos(\theta) = \frac{x_1 \cdot x_2}{\sqrt{\|x_1\|_2 \|x_2\|_2}}$$

$\|x_1\|_2$ norm of x_1
 $\|x_2\|_2$ norm of x_2

① If x_1 & x_2 are unit vec

$$\|x_1\|_2 = \|x_2\|_2 = 1$$

$$\boxed{\cos \theta = x_1 \cdot x_2}$$



relationship b/w euc-dist & cos-sim
↑
(cos-dist)

θ = angle b/w

x_1 & x_2

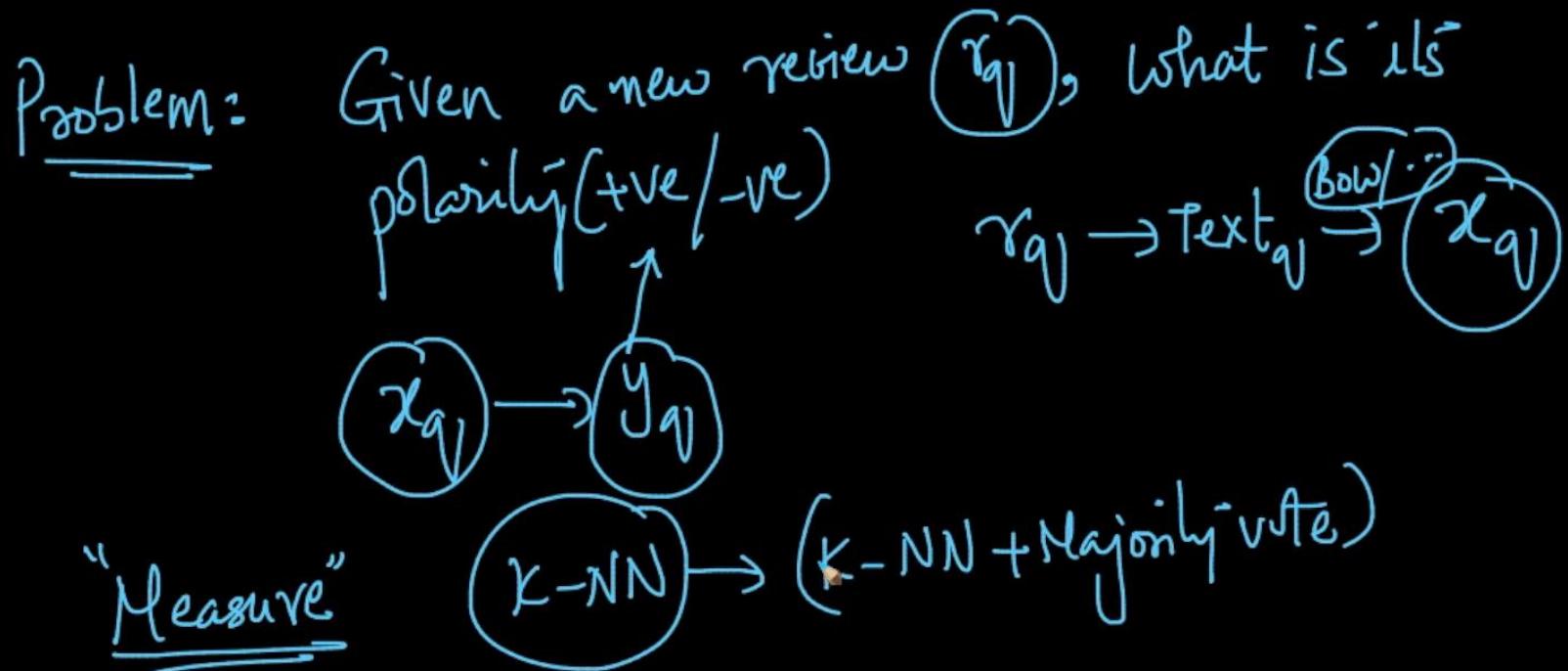
cos-sim

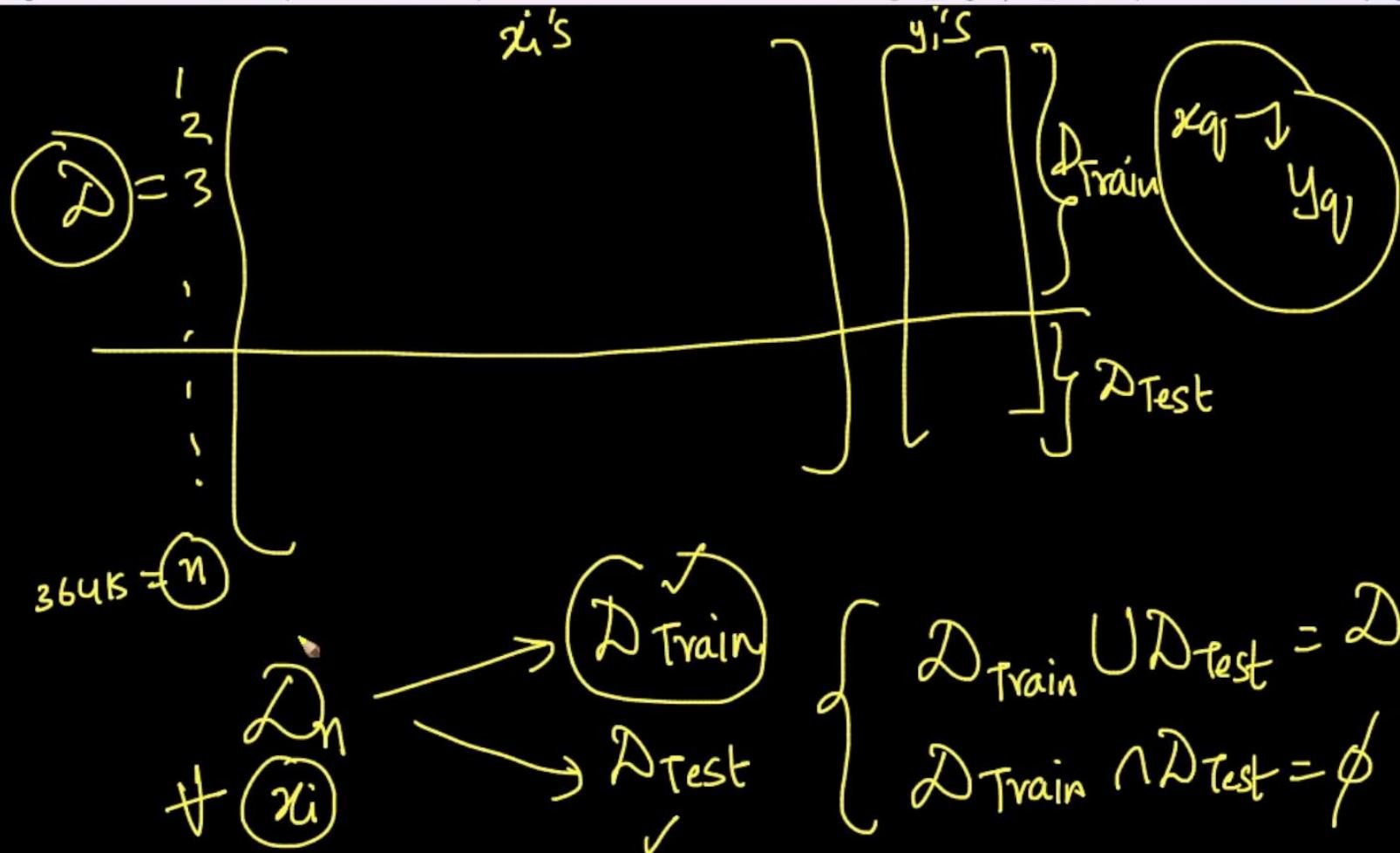
$\cos(\theta)$

cos-dist

- if x_1 & x_2 are unit vect

$$[\text{euc-dist}(x_1, x_2)]^2 = 2 \left(1 - \underbrace{\cos(\theta)}_{\text{cos-dist}} \right)$$
$$\hookrightarrow \boxed{2 \cos(\theta)}$$





$D_n \rightarrow D_{Train} \rightarrow n_1$
 $D_n \rightarrow D_{Test} \rightarrow n_2$

$$n_1 + n_2 = n$$

$$n_1 \approx 0.7n$$

$$n_2 \approx 0.3n$$

Split D_n $\xrightarrow{R} D_{Train}$ $\xrightarrow{R} D_{Test}$ randomly

$n_1 \approx 0.7n$
 $n_2 \approx 0.3n$



$\text{cnt} = 0;$
for each pt in D_{test} :

$$x_i \rightarrow y_i$$

$$x_q = pt$$

use D_{Train} & KNN to determine y_q

- if $y_q = y_{pt}$
 $\text{cnt} += 1$

(end)
 $\text{cnt} = \# \text{ pts for which } D_{\text{Train}} + \text{KNN}$
gave a correct classification



Accuracy = $\frac{\text{cnt}}{n_2}$ → # pts for which
D train + KNN gave a correct class label
pts in D test

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

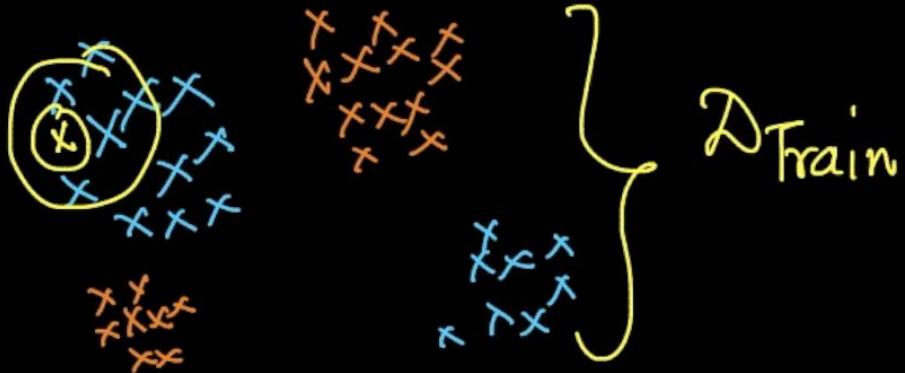
D test

ACC = 0.91 \Rightarrow 91% of times

$x_n \rightarrow y_g$

$K=3$

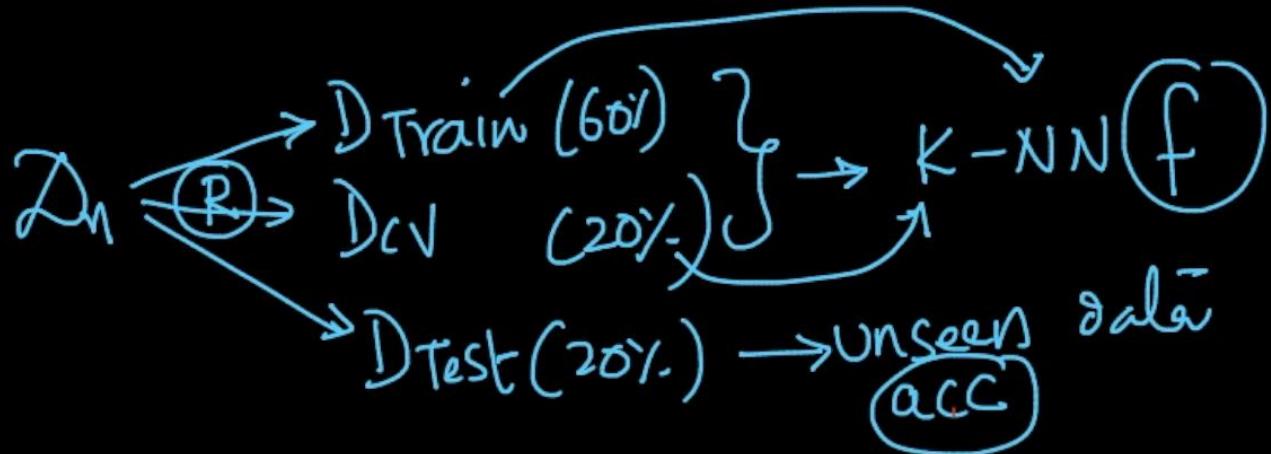
D_{Train}



each point x_i in $\underline{D_{Test}} = \{(x_i, y_i)\}_{i=1}^{n_2}\}$

$x_{qj} = x_1$ for each pt in D_{Test}
 $y_{qj} = \text{Blue } \leftarrow \text{+ve}$ $\rightarrow x_{qj} = \text{pt}$
 $\rightarrow \text{use } D_{Train} \& KNN \text{ to predict } y_{qj}$

K-fold cross-validation

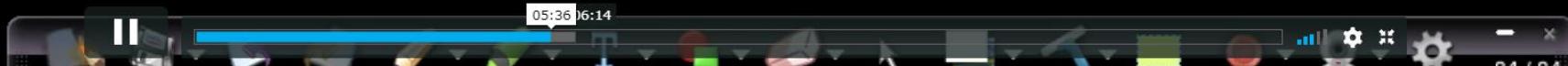
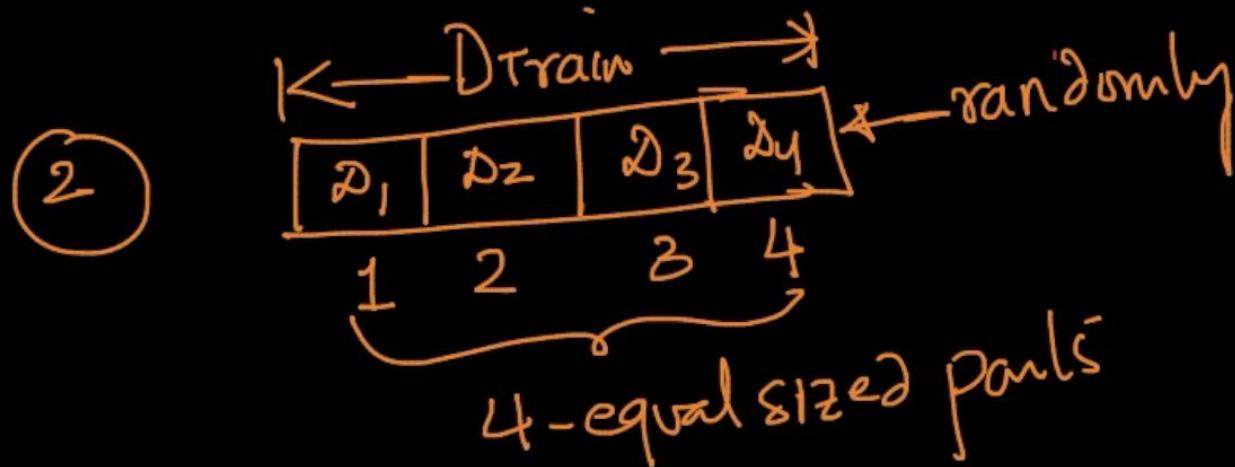
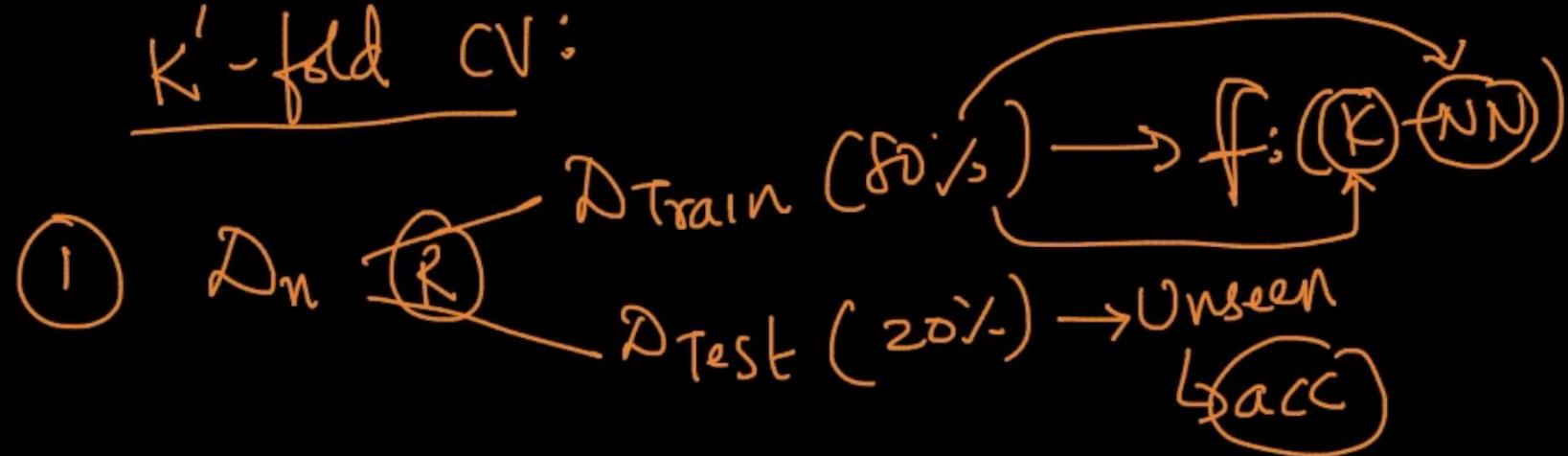


Problem: Only 60% of data (NN) } 80% of data to compute NN
20% CV → K }

20% Test → unseen

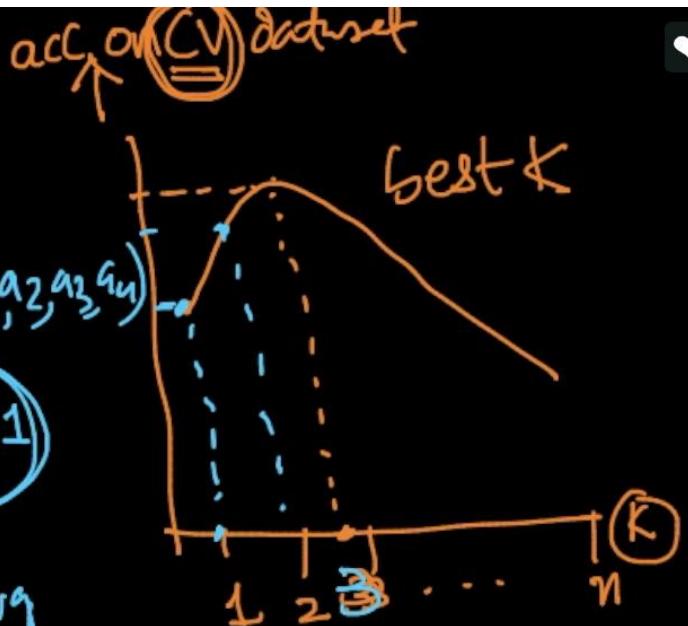


K' -fold CV:



③ $K = 4$ -fold CV
4-times

	Train	CV	acc. on CV dataset
			acc. on CV dataset
$K=1$	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
			$\text{avg}(a_1, a_2, a_3, a_4)$
			$a_{K=1}$
$K=2$	D_1, D_2, D_3	D_4	a_4'
$K=2$	D_1, D_2, D_4	D_3	a_3'
$K=2$	D_1, D_3, D_4	D_2	a_2'
$K=2$	D_2, D_3, D_4	D_1	a_1'
			avg
			$a_{K=2}$



{ acc. on K -NN
for $K=1$ on CV-dataset



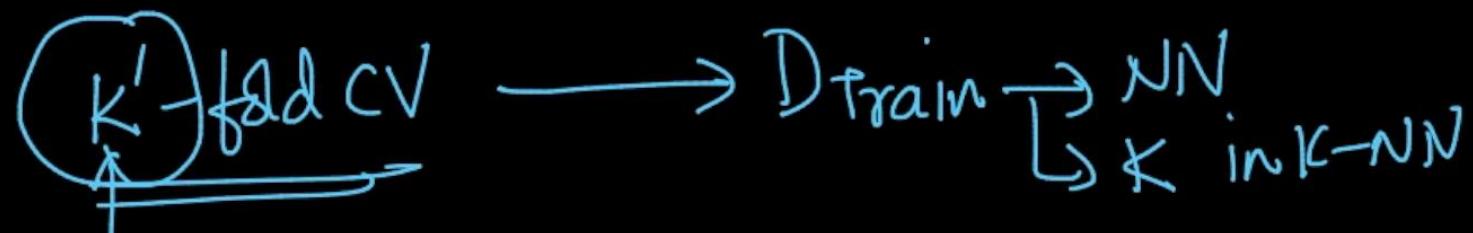
$f : \text{3-NN} \rightarrow D_{\text{train}}$

avg-acc on 4-fold CV

$x \rightarrow D_{\text{Train}} \xrightarrow{\text{NN}} \text{3 nearest neighbors}$

acc \checkmark D_{test} \rightarrow acc of 3-NN on D_{test}

\downarrow
93% on Unseen data



what is the right K'

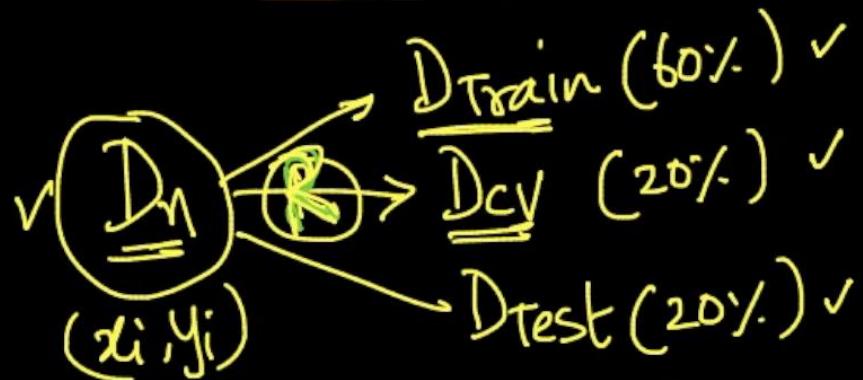
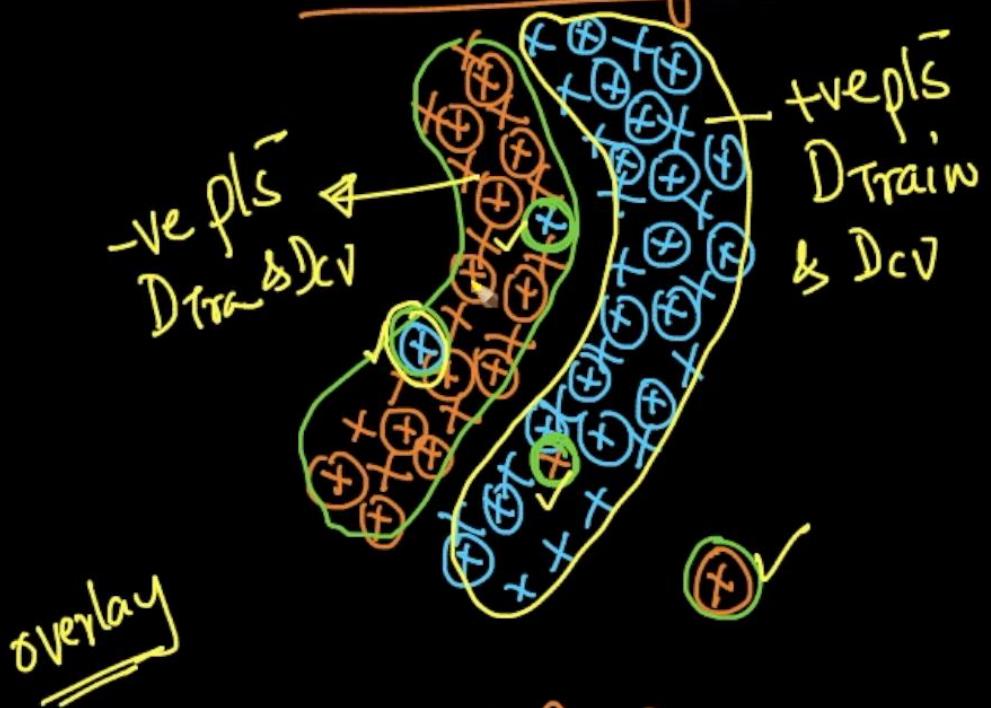
$$K' = 4$$

$$K = 10$$

$$K' = 10$$

rule of thumb : 10-fold CV

Visualizing Train, CV & test datasets:



$\begin{cases} x: -ve \text{ datapoint in } D_{\text{train}} \\ x: +ve \text{ datapoint in } D_{\text{train}} \end{cases}$

$\begin{cases} \textcircled{x}: -ve \text{ point from } D_{\text{cv}} \\ \textcircled{x}: +ve \text{ point from } D_{\text{cv}} \end{cases}$

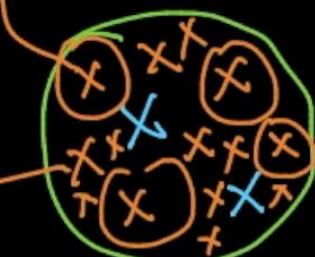
Train, CV, test: (x_i, y_i)
 ↑ class-label

- ✓ ① $D_{train} \cup D_{CV}^{test}$ don't overlap perfectly when randomly sampled
 - ② If there are many +ve pls for D_{train} in a region, then it is highly likely to find many +ve/pls from D_{CV} in that region
 - ③ If there are very few +ve/-ve pls in a region from D_{train} , then it is very unlikely to find +ve/-ve from D_{CV} in that region
- noise
error
outlier*

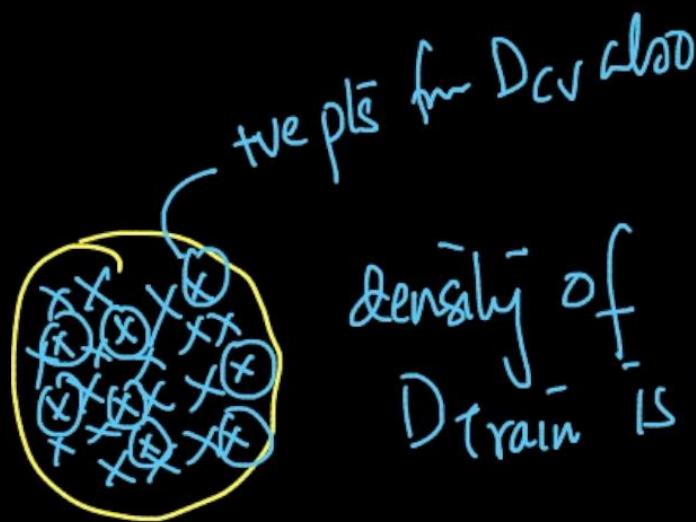


intuitively,

-ve pls from D_{train}
density A
 $D_{\text{train}} -$ ve pls is high



Very low density
may not find any pls from D_{CV} +ve class

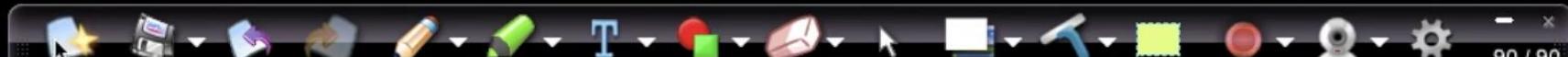


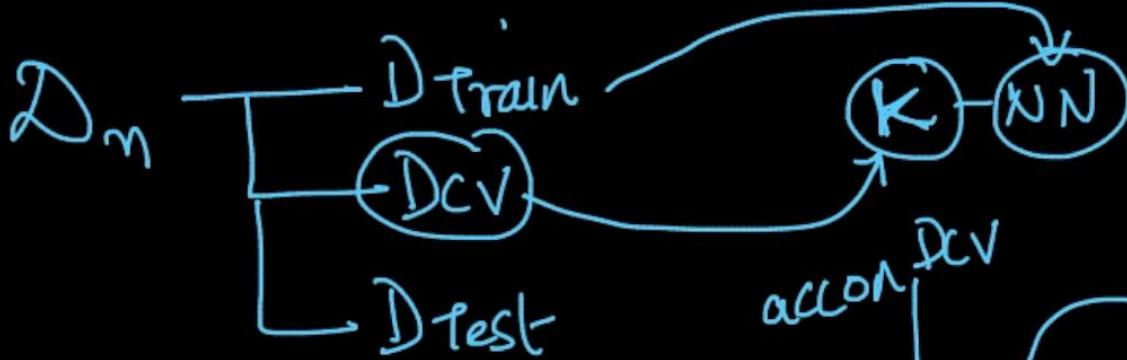
+ve pls for D_{CV} class
density of +ve pls for D_{train} is high



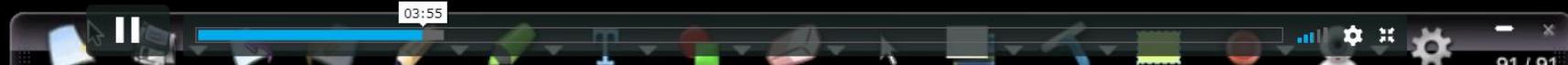
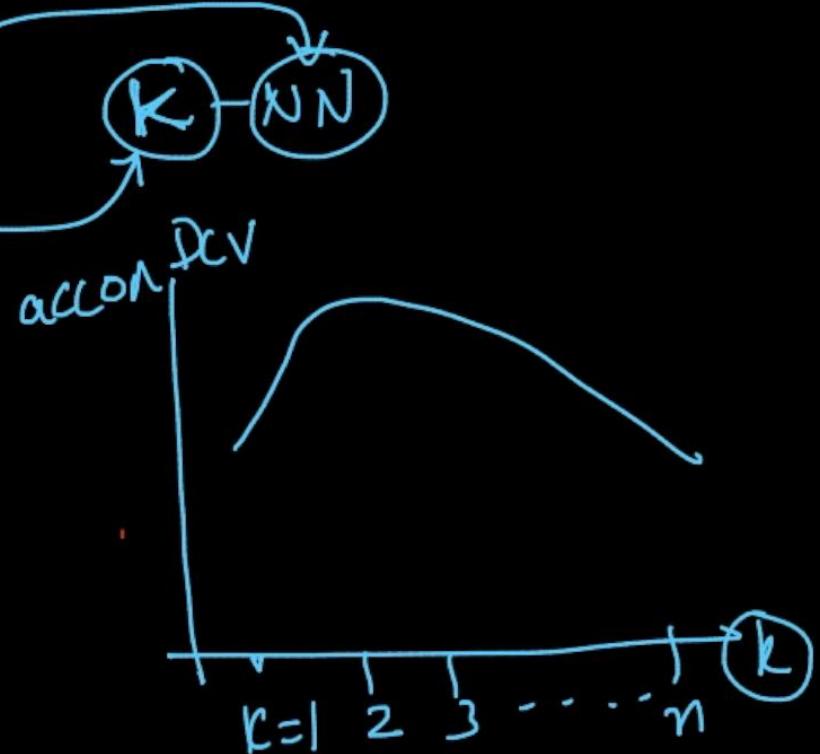
$$(\text{Max}) \uparrow \text{acc} = \frac{\# \text{ correctly classified pls}}{\text{Total \# pls}} = 0.93$$

$$(\text{Min}) \downarrow \text{(error)} = 1 - \text{acc} = 0.07 \rightarrow 7\%$$



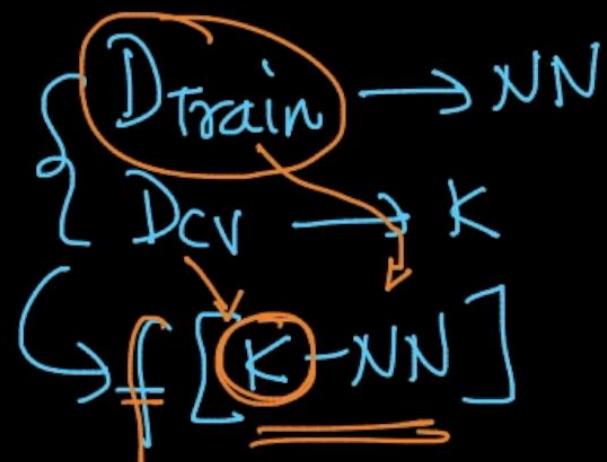


$\text{error on } D_{CV} = 1 - \text{acc on } D_{CV}$
 $\text{acc . on } D_{CV}$
 { for each pt in D_{CV}
 $x_i = \text{pt}$
 $y_q \rightarrow y_i$



Training error:

- { for each x_i in D_{Train}
- find 2 nearest neighbor to x_i from D_{Train}
 - majority rule to get the class label
 - if $y_1 = \text{class label}$
 accurate
 else error



what is the train error
if $2-NN$

$\{ D_{Train} \rightarrow \text{Training}$
 $D_{Train} \rightarrow \text{accuracy}$



Overfitting vs Underfitting

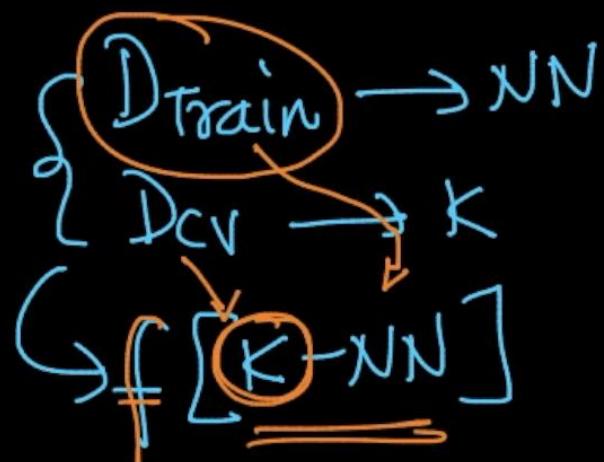
Press Esc to exit full screen

{ K-fold CV \textcircled{N} Dcv \rightarrow best K \rightarrow neither overfit nor underfit

✓(Q) how can we be sure that we are not underfitting
(or) overfitting? \rightarrow plot

Training error:

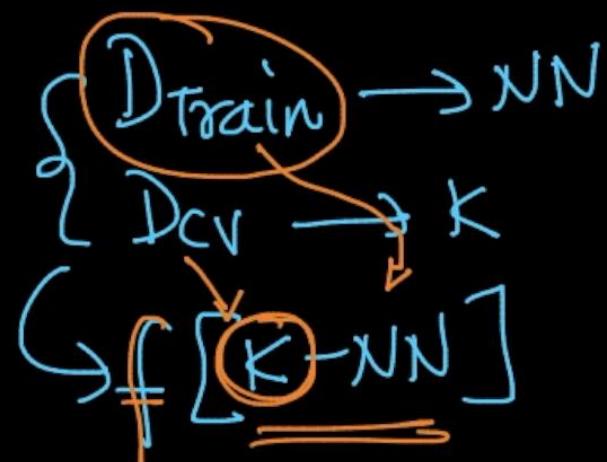
- { for each x_i in D_{Train}
- find 2 nearest neighbor to x_i from D_{Train}
 - majority rule to get the class label
 - if $y_1 = \text{class label}$ accurate
else error



what is the train-error
if $2-NN$

Training error:

- { for each x_i in D_{Train}
- find 2 nearest neighbor to x_i from D_{Train}
 - majority rule to get the class label
 - if $y_1 = \text{class label}$
 accurate
 else error



What is the train-error
if $\underline{2-NN}$

- { $D_{Train} \rightarrow \text{Training}$
 $D_{Train} \rightarrow \text{acc/error}$



error: train-error, Validation-error

error
-> acc



f
5-NN

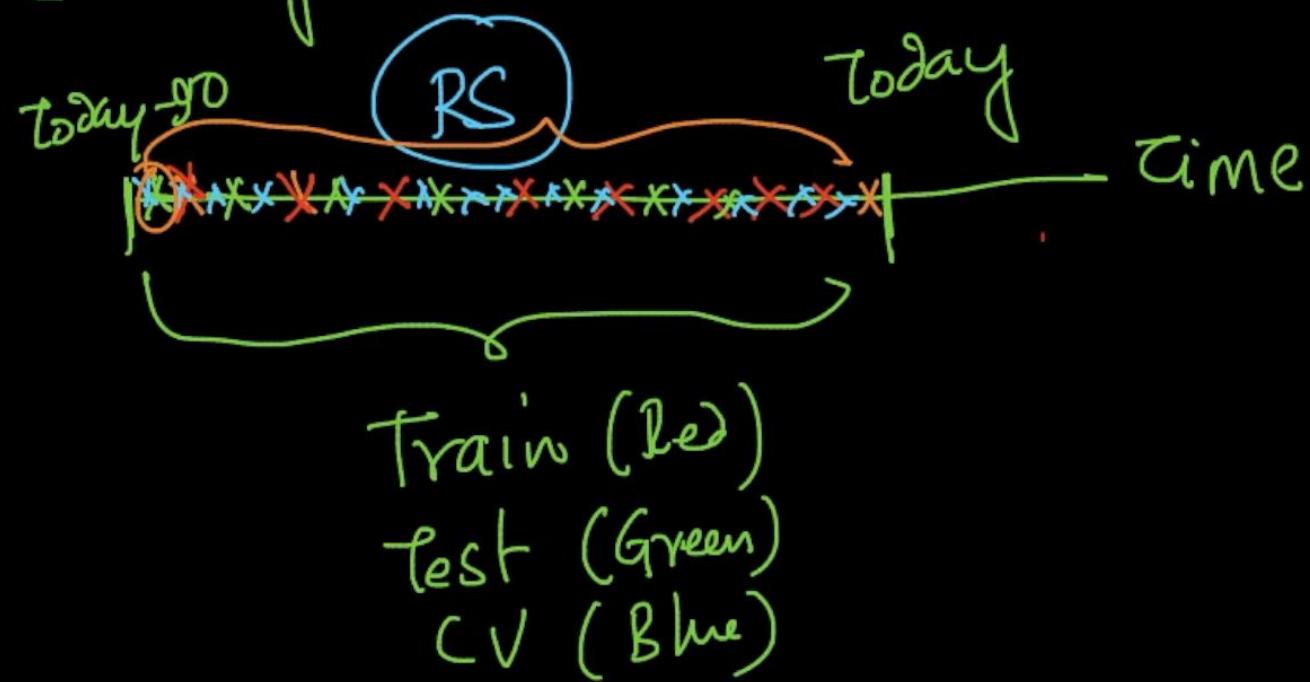


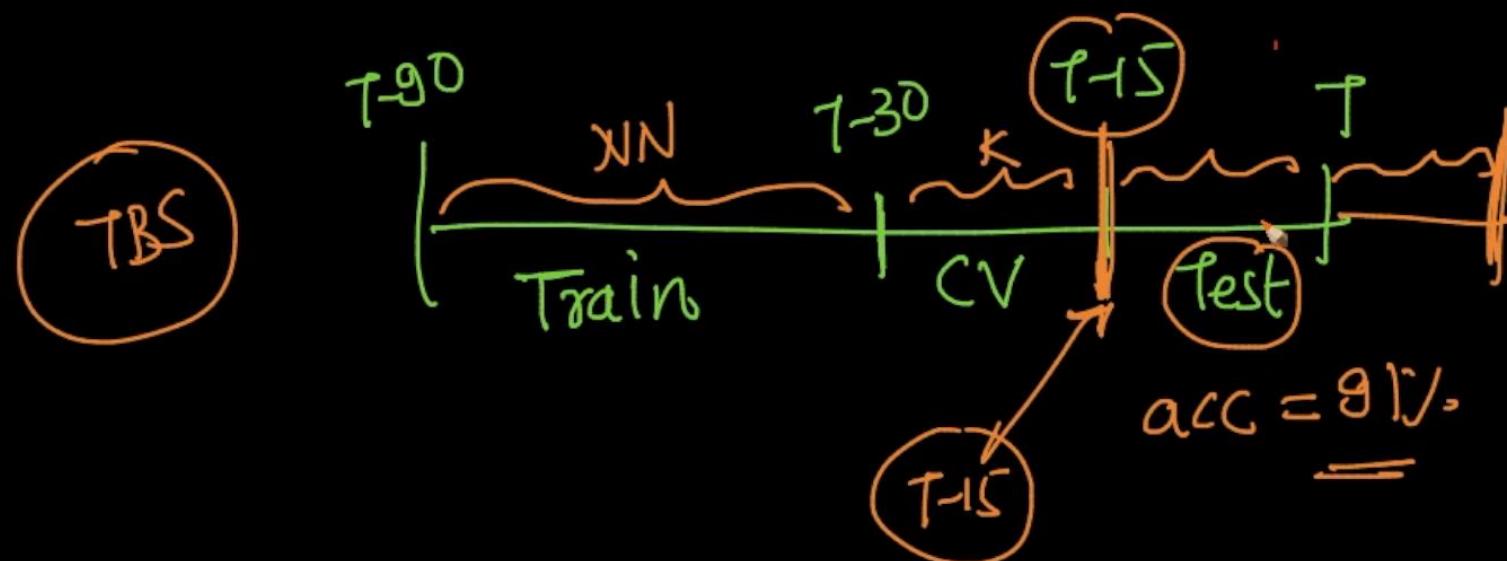
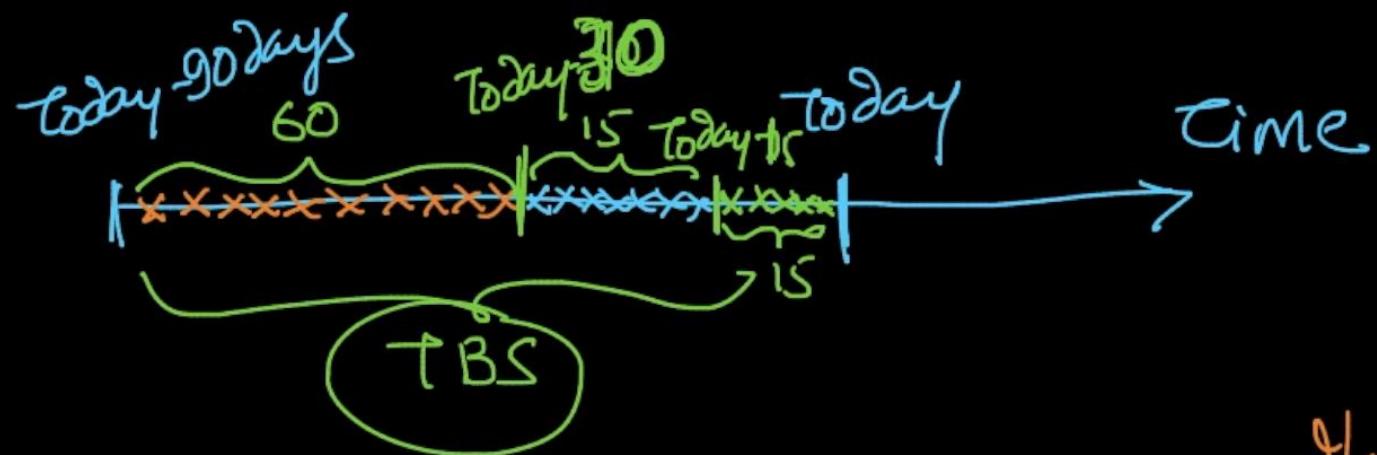


{ Train error is high } — Underfit
Val- error is high }

{ Train error is low } — overfit
Val- error is high }

{ With time, my products & their reviews
change





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



{ whenever time is available & if things / behavior
or data changes over time
then time based splitting is preferable



(K-NN) for regression:

2-class
classfn

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

$f(x_q) \rightarrow y_q \rightarrow$ class-label

✓ find
→ K-NN
↳ May only
use

Reg.

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

$x_q \rightarrow y_q \rightarrow$ number



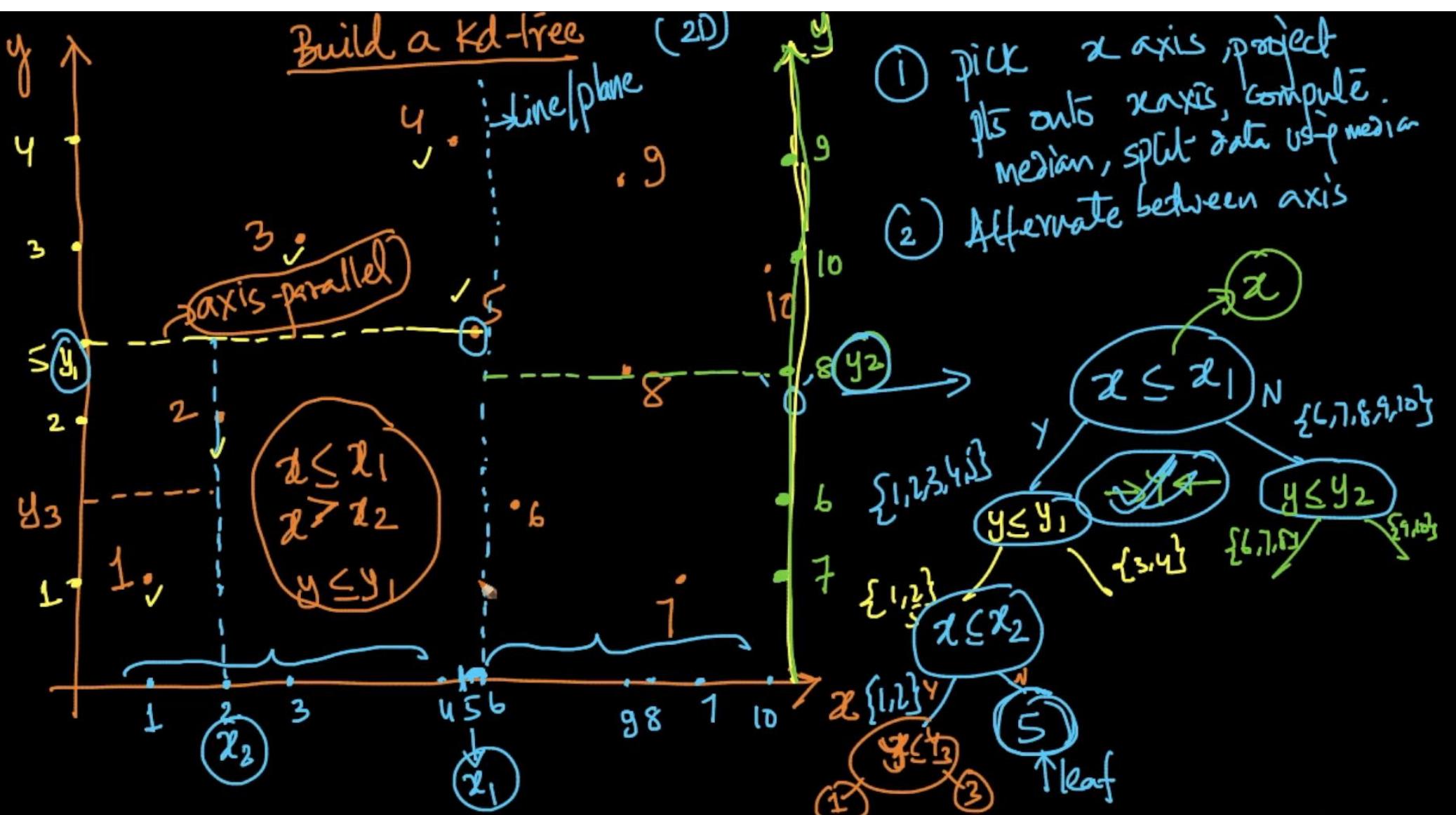
① given (x_q) , find k -nearest neighbors
 $(x_1, y_1), (x_2, y_2) \dots (x_k, y_k)$

~~Class → Majority vote~~
Repsn ② y_q \leftarrow $y_1, y_2, y_3 \dots y_k$ not $\rightarrow n^1$ IR

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

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

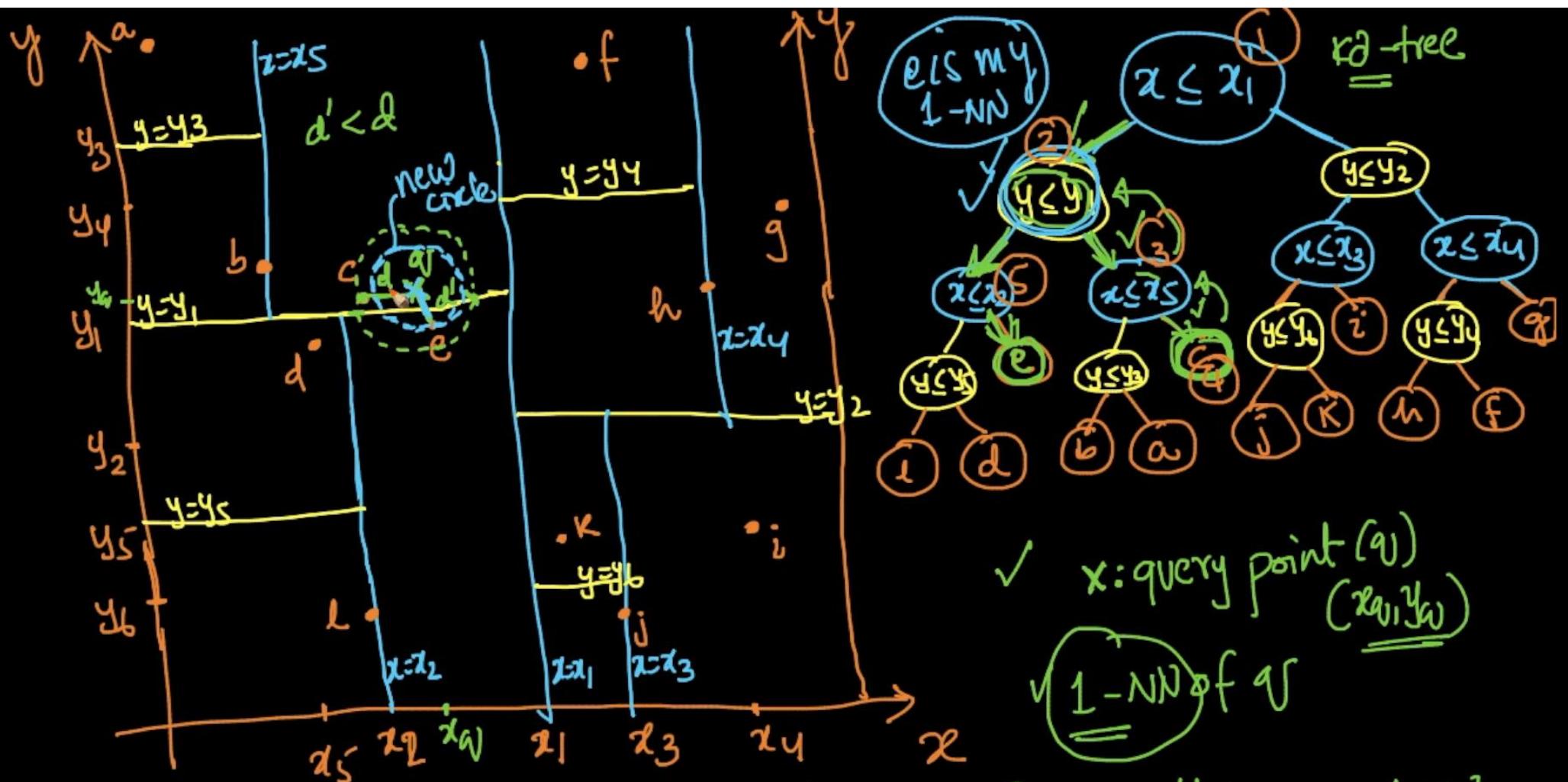




Kd-tree → go through each axis one after the other till leaf node
↳ breaking up my space using axis-parallel lines / planes into subcubes / hyper-cubes

{
2D → axis parallel lines → rectangles
3D → " " " planes → cuboids
nD → " " " hyperplanes → hypercuboids



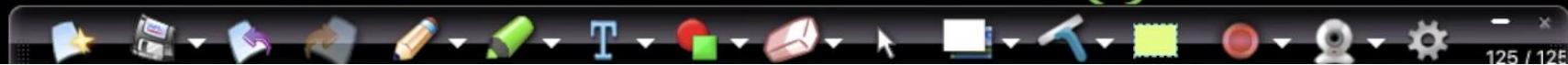


Circle/hyper- Δ phane: $\text{rad} - d$
 $C - \mathcal{O}$

\checkmark c : Could be my 1-NN
 \checkmark e : Could be my 1-NN

\checkmark x : query point (q)
 $(\underline{x_q}, \underline{y_q})$

\checkmark 1-NN of q



Time-complex

~~at~~ comparisons to find 1-NN

✓ best-case: $O(\lg(n))$

n : #pls

✓ worst-case: $O(n)$

$O(\log n)$

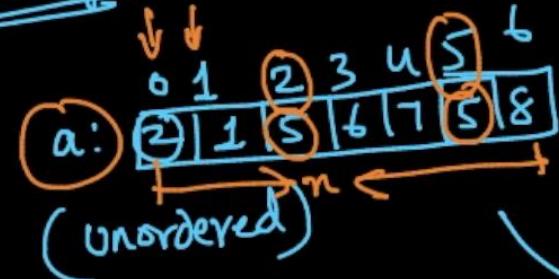
{ perfect analysis:

very
complex



Hashing & LSH ↳ locality sensitive hashing → d is large

DS / Algo



hash fn
 $h(x)$

Hashtable (Dict)
 K, V

K	V
2	1
5	2, 5
1	1
.	.

indexes in a

$h(5)$
 $O(1)$ time

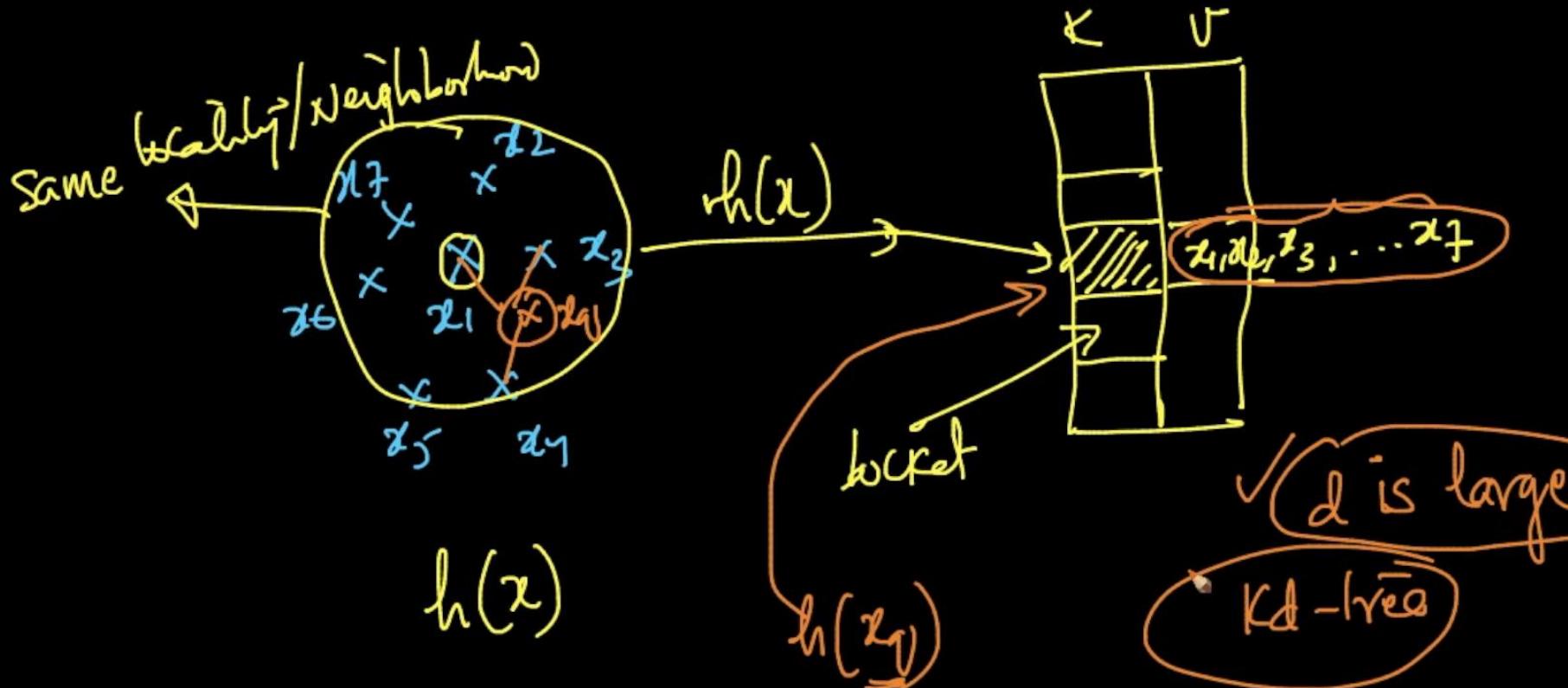
find if 5 exists in ~
Sequentially search
 $O(n)$

$h(x)$ ↳ bucket

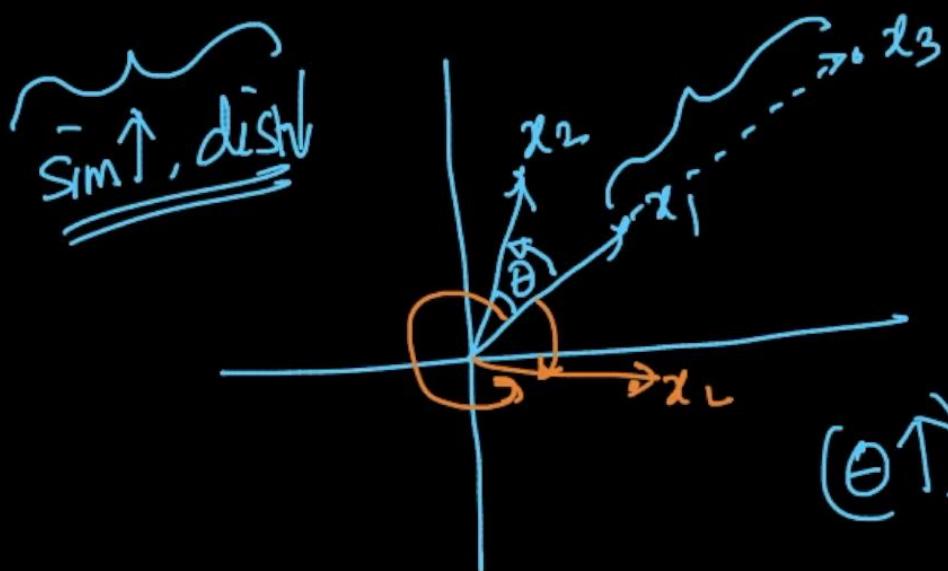


Locality sensitive hashing

$$\mathcal{D} = \{\underline{x_1, x_2, \dots, x_m}\}$$



Locality Sensitive hashing for Cosine Similarity



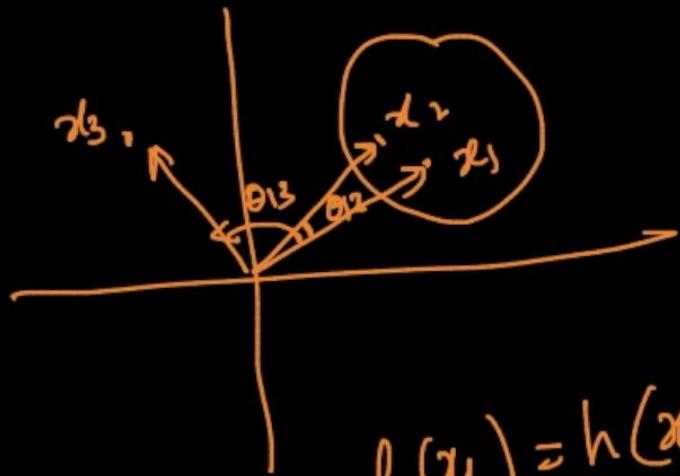
$$\text{cos-Sim}(x_1, x_2) = \cos(\theta)$$

$$\text{cos-Sim}(x_1, x_3) = \cos(0) = 1$$

$(\theta \uparrow)$; (x_1 & x_2 are "angular" \uparrow separated) \rightarrow dist \uparrow
 $\downarrow (\text{cosine} \downarrow)(\text{sim} \downarrow)$

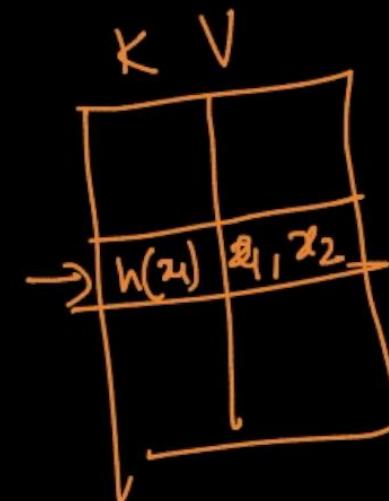
Cosine-SIM \approx angular similarity

LST for cos-Sim

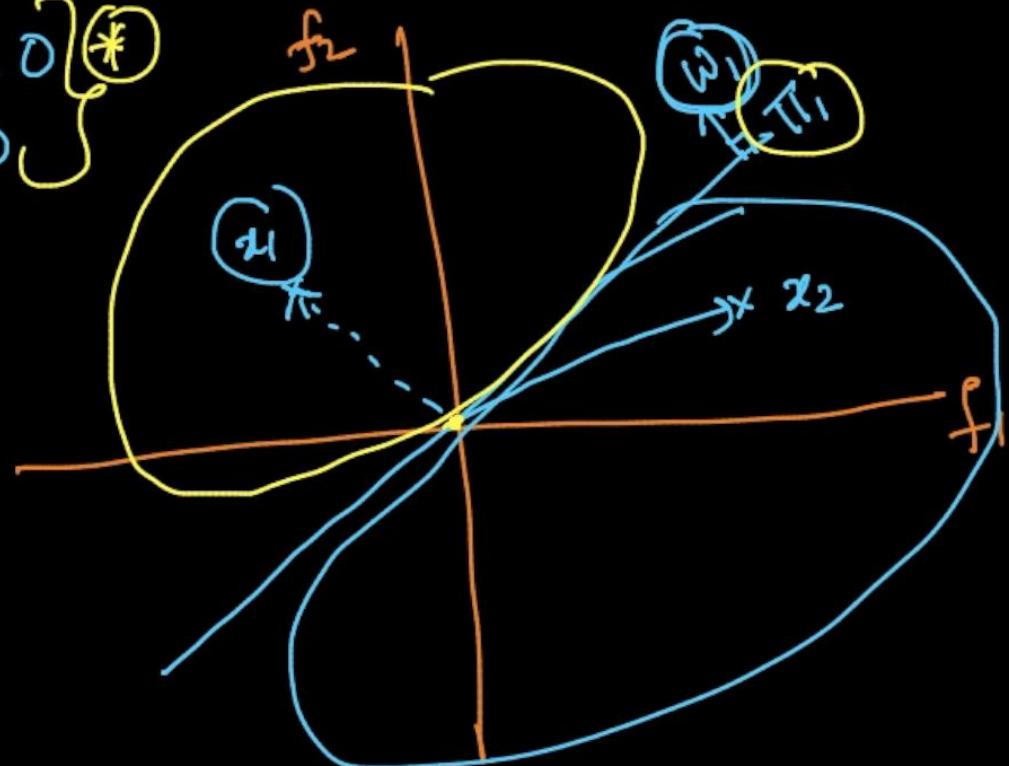


$$h(x_1) = h(x_2) = \text{key}$$

$\theta_{12} \sim 0$ (small) $\Rightarrow x_1 \text{ & } x_2 \text{ are similar/ close}$



$$\begin{cases} \omega_1^T x_1 \geq 0 \\ \omega_1^T x_2 \leq 0 \end{cases}$$

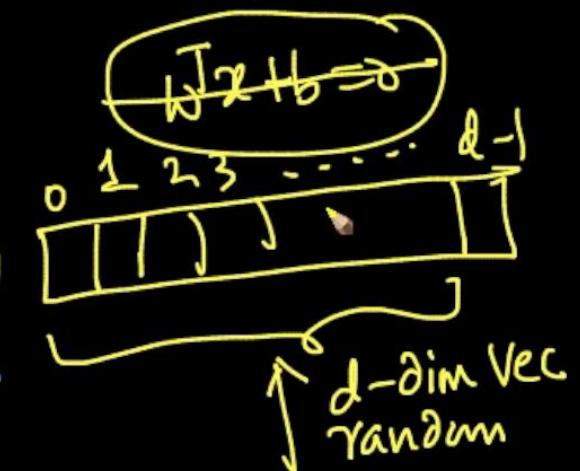


2D

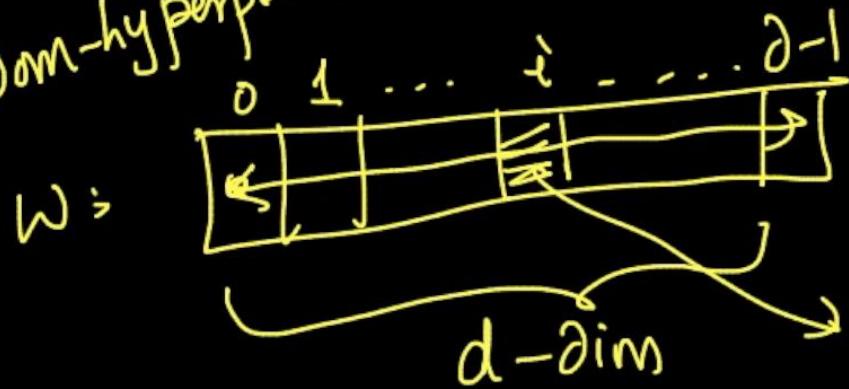
w_1 : unit normal vector
to T_{f_1}

① Random-hyperplane

$$\begin{cases} (\bar{w})^T x = 0 \\ w: \text{normal to the plane} \end{cases}$$



Random-hyperplane

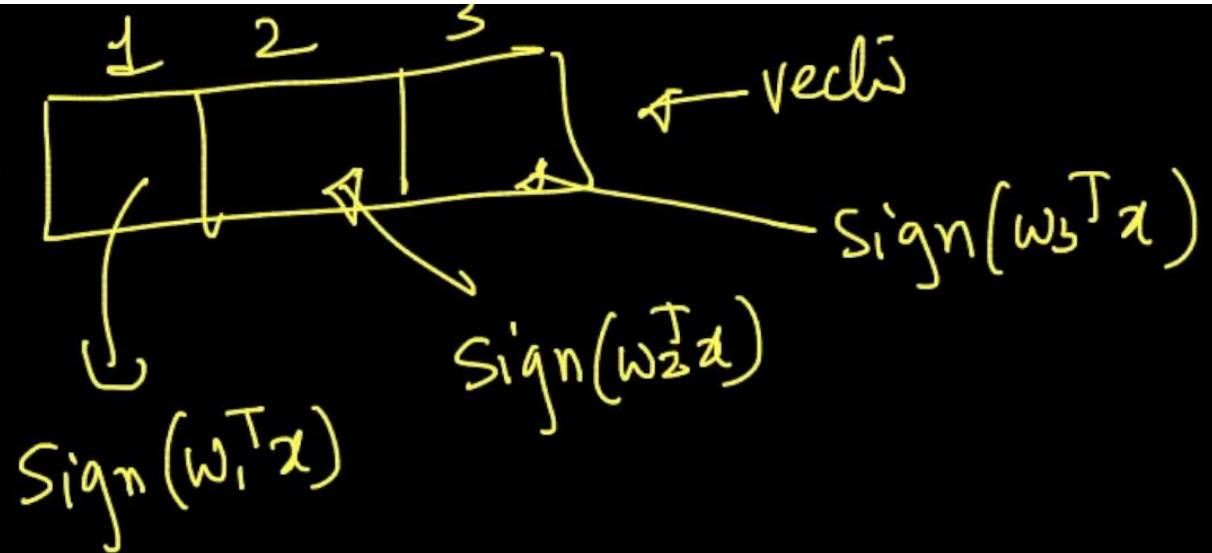


random numbers sampled from
 $\underline{\underline{N(0,1)}}$

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



$h(x) =$

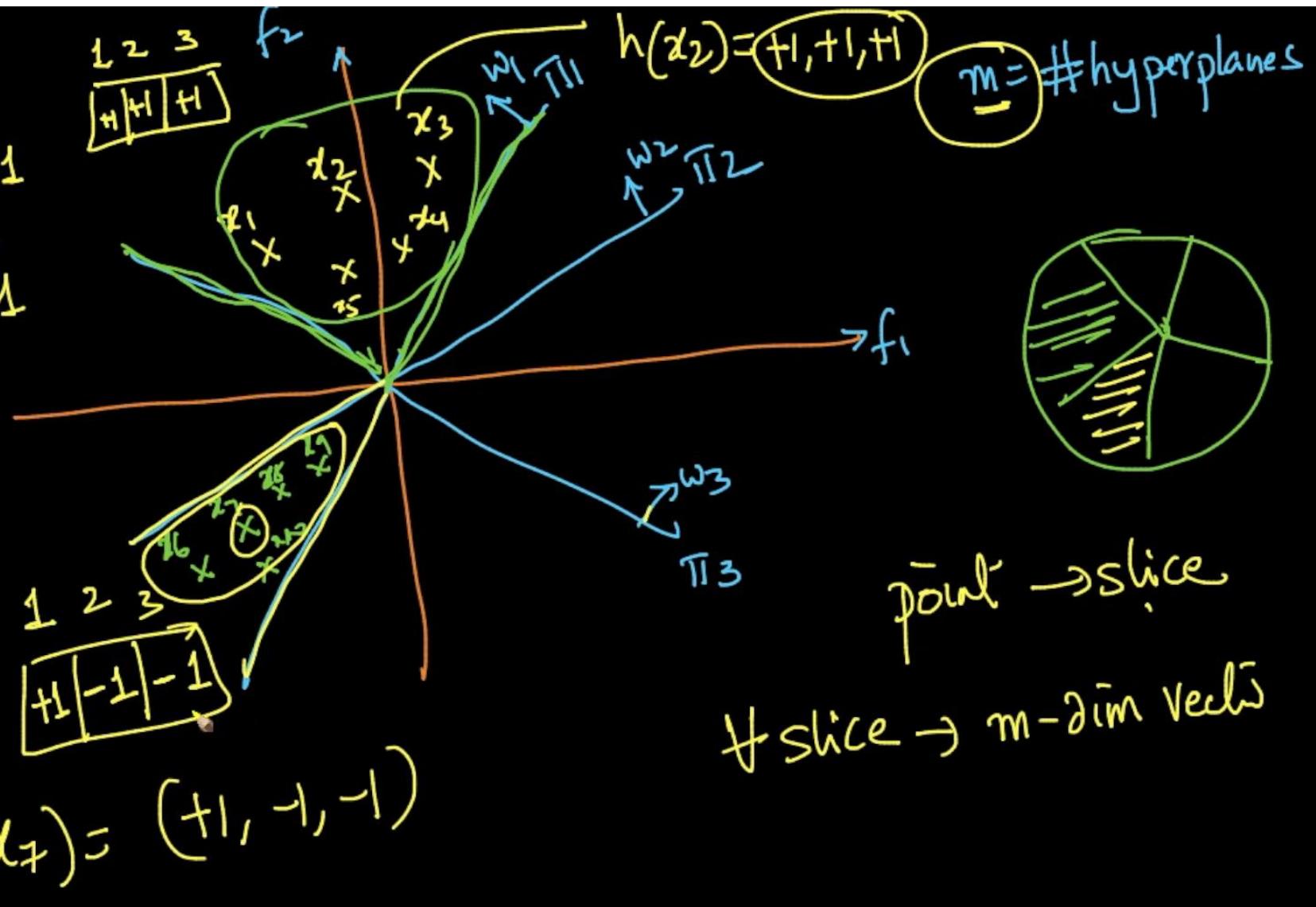


LSTH for $\cos \theta^m$

$$x_1^T w_1 > 0 \rightarrow +1$$

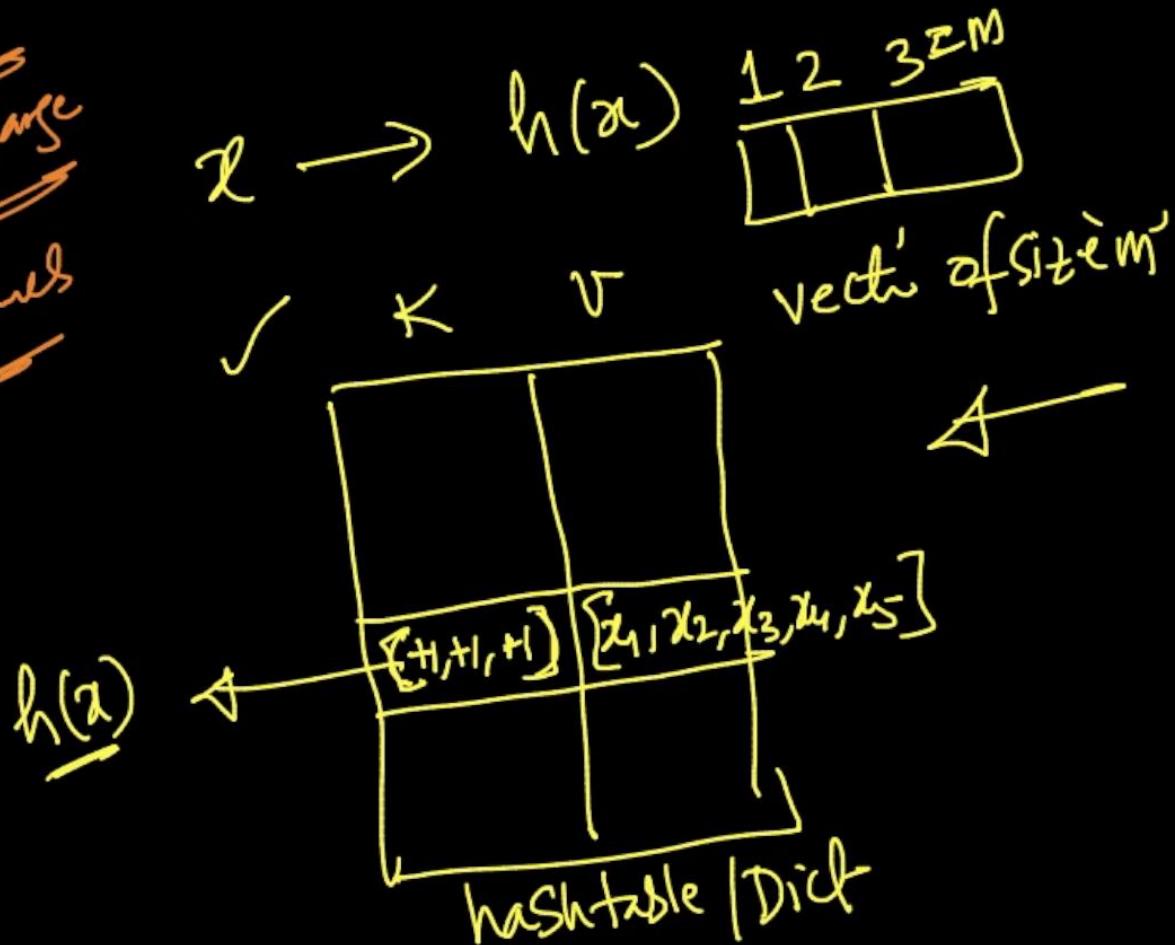
$$x_1^T w_2 > 0 \rightarrow +1$$

$$x_1^T w_3 > 0 \rightarrow +1$$



n points
 d -dim
 $\circ d$ could be large
 m -hyperplanes
 $w^T x$

hashtable



construct LSH
hash-table

Time for a hashtable
 $\sim O(mdn)$

Space:
 $O(nd)$

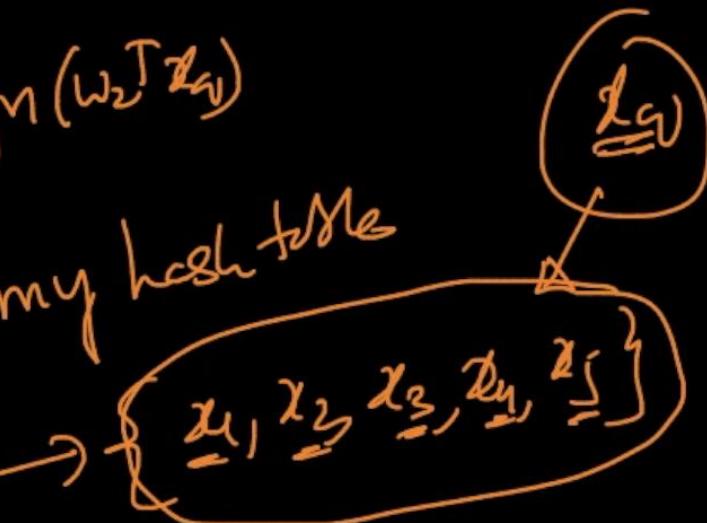
(Query point) (x_w)

① $h(x_w) =$



$h(x_w)$ as the key in my hash table

• do get ($h(x_w)$) →

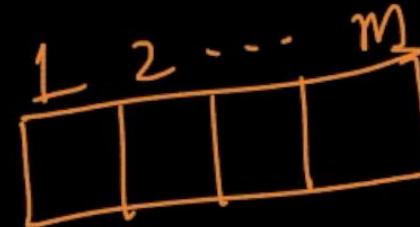


x_w

Given a hash-table:-

Time comp for querying:

$$x_{qj} : \underline{h(x_{qj})}$$



$$\mathcal{O}(md)$$

$$\mathcal{O}(md + \frac{n'd}{m})$$

Small

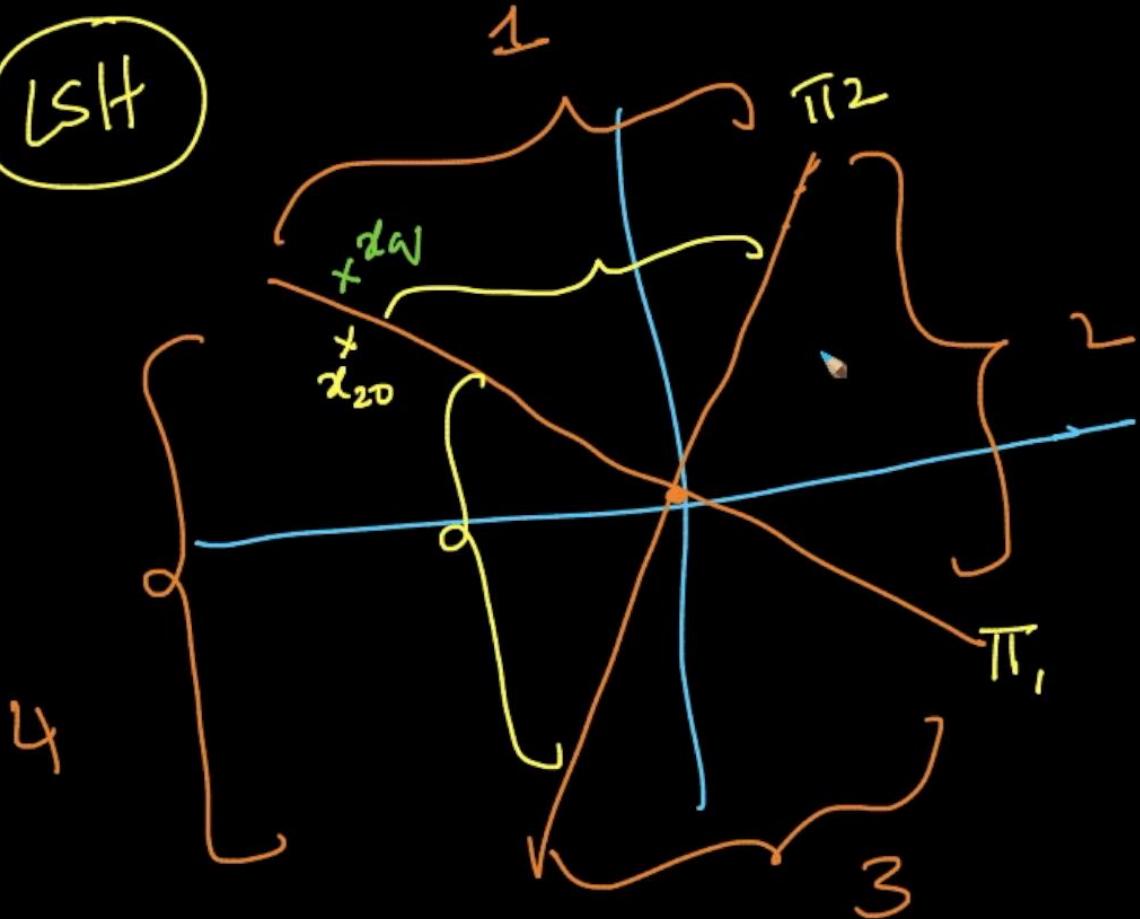
$$\mathcal{O}(md)$$

$x_{qj} \rightarrow n'$ elements in the bucket/slice
if n' is small

$$n' \approx m$$



LSH



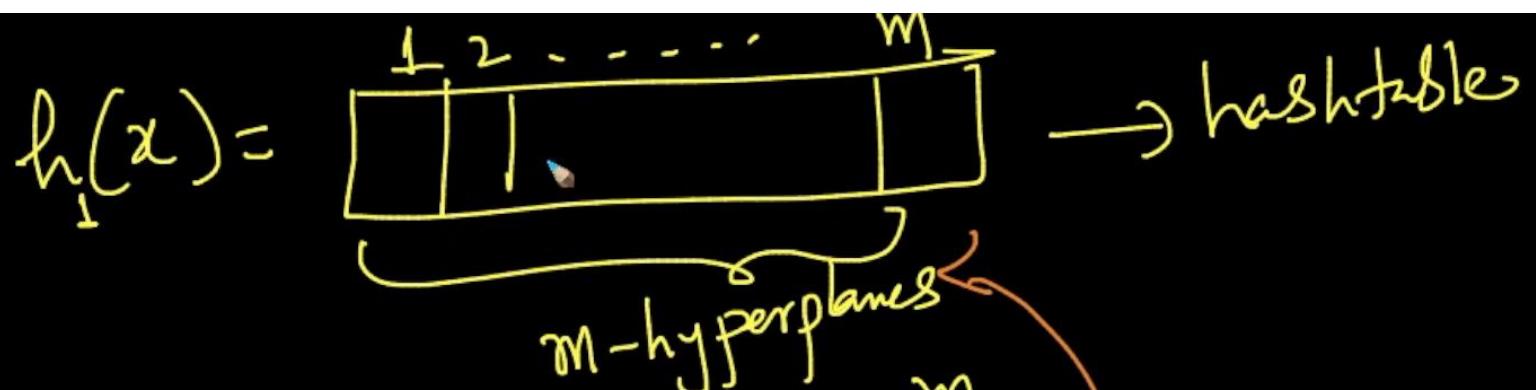
GSSIM

{ x_{20} is very close to x_q }



* could miss a NN
in GSSIM as your
measure

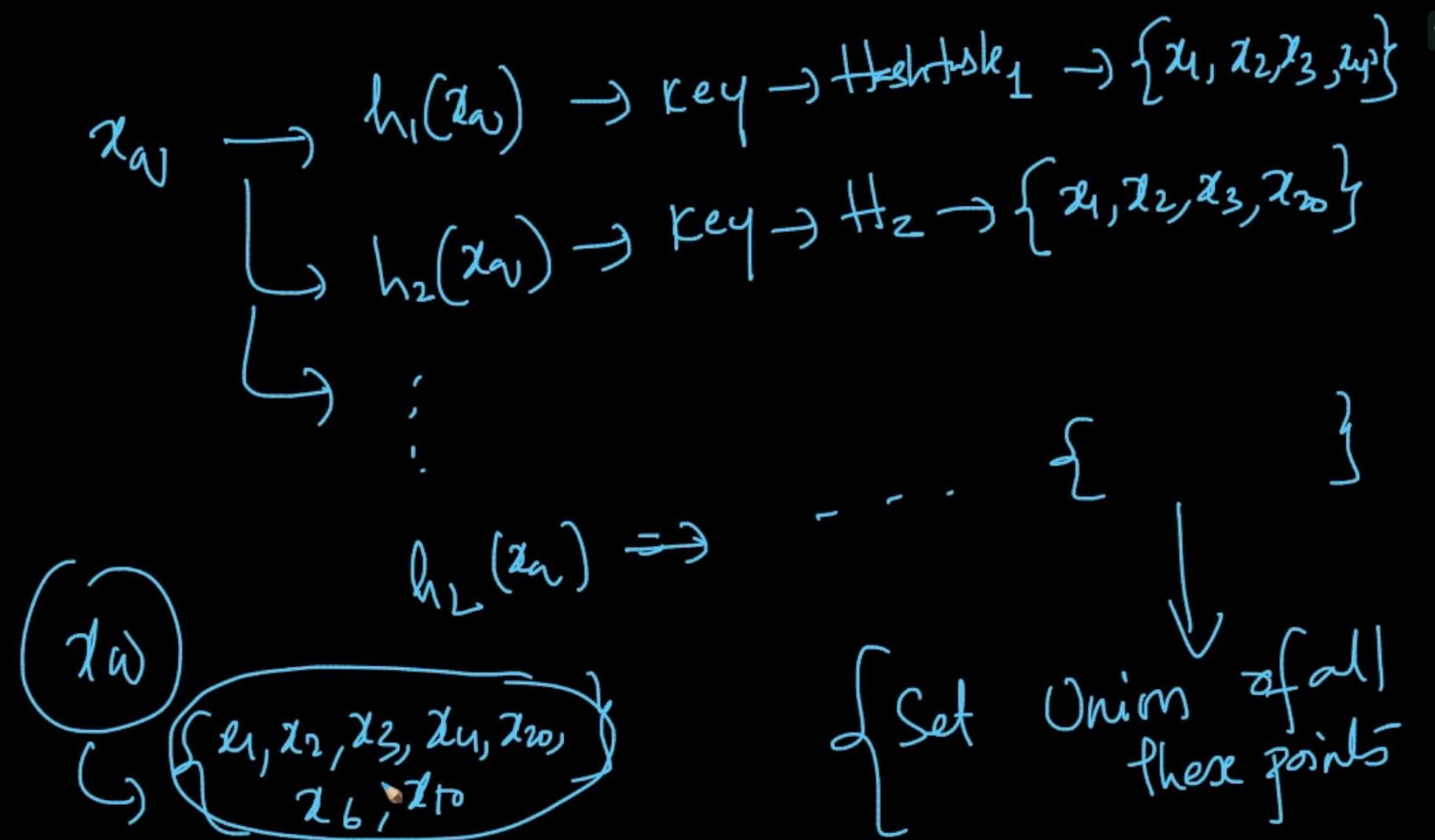
①



②



L-hashtables



Time Complicated to query given L hash-tables

$$m = O(\lg n)$$

$$O(m \lg L)$$

L is small

{ as number of hyperplanes ↑ ↓
slices ↑ ↓
 $\#$ pls per slice ↓ ↑ }

https://www.appliedaicourse.com/lecture/11/applied-machine-learning-online-course/2950/lsh-for-cosine-similarity/3/module-3-foundations-of-natural-lan...

I summarised the steps in short. Pls review it

Steps for locality sensitive hashing:

1. First make 'm' hyperplanes to split into regions and create slices such that cluster of points lie in a particular slice and be called their neighbourhood. typically $m = \log(n)$

2. Next for each point create a vector(also called hash function) by $W_1(\text{transpose}).point$. if it is greater than 0 , it lies on the same side of that hyperplane else other side. Based on that create a vector of m size. For eg the vector can be [1,-1,-1] denoting point x lies on same side of normal to hyperplane 1, opposite side to normal of hyperplane 2 and 3. Now this vector serves as key to the hash table and all the points with the same key or vector representation will go in the same bucket as they have similar vector representation denoting they lie in the neighbourhood of each other.

3. Now it may happen that two points which are very close fall on different slice due to placing of hyperplane and hence not considered as nearest neighbour. To resolve this problem, create l hash tables. (l is typically small). In other words repeat step 2 l times thus creating l hash tables and m random planes l times. So when a query point comes compute the hash function each of the 'l' times and get its neighbours from each of bucket. Union them and find the nearest neighbours from the list. So basically in each 'l' iterations create m hyperplanes and hence region splitting will be different thus vector representation or hash function of the same query point will be different in each of the representations. Thus the hash table will be different as points which lied on the same region in previous iteration might lie in a different region in this iteration and vice versa due to different placement of hyperplanes.

Time complexity is $O(m^*d^*l)$ for each query point. And for creating the hash table $O(m^*d^*l^*n)$ which is one time only.

Space complexity is $O(n)$

Reply    

Jan 27, 2019 16:25 PM

2.15 Visualizing train, validation and test datasets ✓ 13 min

2.16 How to determine overfitting and underfitting? ✓ 19 min

2.17 Time based splitting ✓ 19 min

2.18 k-NN for regression ✓ 5 min

2.19 Weighted k-NN ✓ 8 min

2.20 Voronoi diagram ✓ 4 min

2.21 Binary search tree ✓ 16 min

2.22 How to build a kd-tree ✓ 17 min

2.23 Find nearest neighbours using kd-tree ✓ 13 min

2.24 Limitations of Kd tree ✓ 9 min

2.25 Extensions ✓ 3 min

2.26 Hashing vs LSH ✓ 10 min

2.27 LSH for cosine similarity 40 min

2.28 LSH for euclidean distance 13 min

2.29 Probabilistic class label 8 min

2.30 Code Sample:Decision boundary . 23 min

Type here to search           23:38 07-07-2019 ENG 4