

Notes on KNN :

- 1) Explain about Supervised Machine Learning?

Classification:

ex: Given a new review text, determine/predict if the review is +ve or -ve

Concept:  $y = f(x)$   $\rightarrow$  finding the function in a classification problem.

$$y_q = f(x_q) \quad (x_q \rightarrow \text{Query review})$$

1. Binary classification:-

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

such that  
 $f(x_i) = y_i$

↓  
-ve      ↓  
+ve

2. Multiclass classification:-

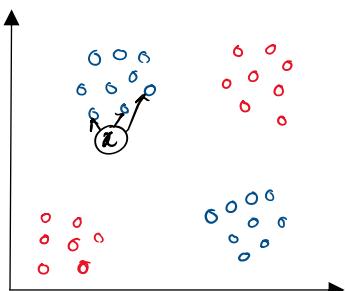
$$y_i \in \{0, 1, 2, 3, 4, 5, 6\}$$

3. Regression:  $y_i \in \mathbb{R}$  ---- Real number

- 2) Explain about K-Nearest Neighbours?

KNN : k-nearest neighbours

$$D = (x_i, y_i)_{i=1}^n ; \text{ given } x_q, \text{ find } y_q \in \{0, 1\}$$



① find k nearest point to  $x_q$  in D

let  $k=3$      $\{x_1, x_2, x_3\}$  3 nearest neighbours to  $x_q$   
 $\downarrow \quad \downarrow \quad \downarrow$   
 $y_1 \quad y_2 \quad y_3$

Notes on KNN :

② majority vote :-

{  $y_1, y_2, y_3$  } → majority or voting  
+ve    +ve    -ve

∴  $y_q$  for given  $x_q$ , will be +ve . . . . . (2/3 votes)

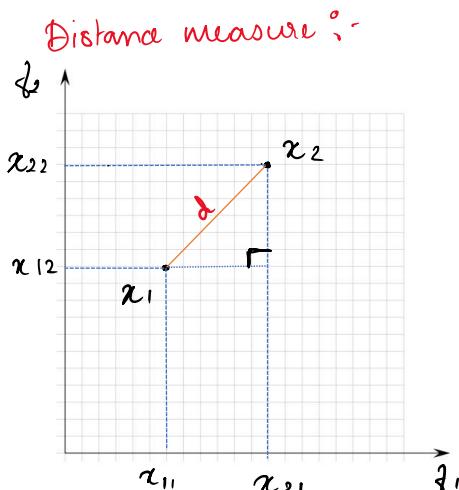
Note:  $k$ - should be odd number to get majority vote.

3) Failure cases of KNN?

Failure cases:

1. If there is randomness or noise in the data.
2. If  $(x_q)$  or -test value are far away from the cluster.

4) Define Distance measures: Euclidean(L2), Manhattan(L1), Minkowski, Hamming?



① Euclidean Distance :-

$d$  = length of shortest line from  $x_1$  to  $x_2$

$$\text{Euclidean Distance} = d = \sqrt{(x_{21} - x_{11})^2 + (x_{22} - x_{12})^2}$$
$$= \|x_1 - x_2\| \quad (c^2 = a^2 + b^2)$$

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

features       $f_1$        $f_2$       euclidean  $\rightarrow \|x_1 - x_2\|$

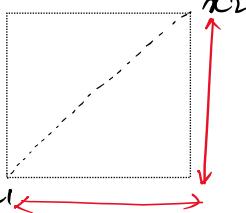
$$\begin{array}{lll} x_1 & x_{11} & x_{12} \\ x_2 & x_{21} & x_{22} \end{array}$$

$$= \left( \sum_{i=1}^d (x_{1(i)} - x_{2(i)})^2 \right)^{\frac{1}{2}}$$

Note :  $\rightarrow \|x_1 - x_2\|_2$  . . . . . L2 norm

Notes on KNN :

② Manhattan Distance :-


$$\sum_{i=1}^d |x_{1i} - x_{2i}| \quad \text{from origin}$$
$$\|x_1 - x_2\|_1 = \sum_{i=1}^d |x_{1i}| \quad \text{absolute}$$

③ Minkowski distance :- ( $L_p$ -norm)

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

If  $p=2$  minkowski distance = euclidean distance

$p=1$  minkowski distance = Manhattan distance

Note: \* distance are for 2 points

\* norms are for vector

i.e.: Minkowski  $\rightarrow L_p$ -norm of any vector  $(x_1 - x_2)$

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

④ Hamming distance:-

$x_1, x_2 \rightarrow$  boolean vector  $\rightarrow$  ex: Boolean BOW

Hamming distance  $(x_1, x_2) =$  # locations / dimensions where binary vector differs

ex:-

$$x_1 = [0, 1, 0, 1, 0, 1, 0, \dots]$$
$$x_2 = [1, 0, 0, 1, 0, 1, 1, \dots]$$

Hamming distance  $(x_1, x_2) = 3$

ex: Gene sequence : (AGTC)  $\rightarrow$  long sequences

$$x_1 = A \boxed{G} \boxed{A} T C T C \boxed{A} \boxed{G}$$
$$x_2 = A \boxed{A} \boxed{G} T C T C \boxed{G} \boxed{A}$$

Hamming distance = 4

Notes on KNN :

5) What is Cosine Distance & Cosine Similarity?

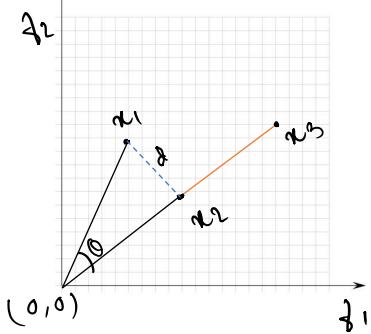
as Cosine Similarity ↑

Cosine distance ↓

as Cosine Similarity ↓

Cosine distance ↑

$$1 - \text{Cosine Similarity } (x_1, x_2) = \text{cosine distance } (x_1, x_2)$$



$d$  = Euclidean distance

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

where  $\theta$  = angle between  $x_1$  &  $x_2$

now lets take  $x_3$ :

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

$$\begin{aligned} \text{cos-sim } (x_2, x_3) &= 1 \quad \left\{ \text{because } \theta = 0 \right. \\ \cos \theta &= \cos 0 = 1 \end{aligned}$$

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

$$= 1 - 1$$

$$\therefore \text{cosine dist} = 0 \quad \text{--- for } x_2 \& x_3$$

Cosine Similarity and cosine distance uses angle to measure.

$$\cos(\theta) = \frac{x_1 \cdot x_2}{\|x_1\|_2 \|x_2\|_2} \quad \begin{array}{l} \text{If } x_1 \& x_2 \text{ are unit vector} \\ \text{then } \cos \theta = x_1 \cdot x_2 \end{array}$$

\* Relationship between Cosine Similarity and Euclidean Distance?

$$\therefore \text{cosine dist} = (1 - \cos \theta)$$

$$d^2 = 2 * (1 - \cos \theta)$$

Square of Euclidean distance =  $2 \times$  cosine distance

$$d_{(x_1, x_2)}^2 = 2 * \text{cos-dist } (x_1, x_2)$$

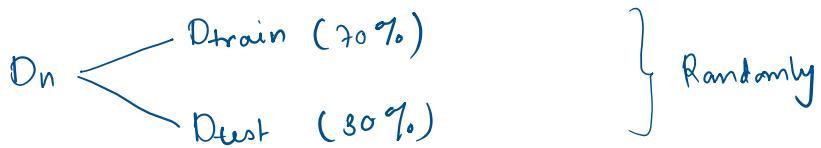
Notes on KNN :

6) How to measure the effectiveness of k-NN?

ex:

problem : Given a new review ( $x_q$ ), what is the polarity (+ve/-ve)

We have  $D_n \quad \therefore x_i \rightarrow y_i$



① For each point in  $D_{test}$

$x_q \rightarrow$  point

use  $D_{train} + kNN$  to detect  $y_p$

② If  $y_q = y_p$

count += 1

where, count = # pts which  $D_{train} + kNN$  gave correct class label

③ Accuracy

$$\text{Accuracy} = \frac{\text{count}}{n_{test}} \quad \therefore n_{test} = \# \text{ of pts in } D_{test}$$

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

# Test / Evaluation time & Space complexity

Input :-  $D_{train}, k, x_q \in \mathbb{R}^d$

Output :-  $y_q$

If  $k$  is small

$\therefore n = \text{no. of pt in Dataset}$

Time & Space complexity =  $O(nd)$

$d = \text{dimensions}$

Note it is roughly equal to  $O(nd)$ , if  $d$  is small - then  $O(n)$

Notes on KNN :

7) Limitations of KNN?

### 1. Space-time Complexity:

for a example where  $k$  so small

$$\begin{array}{l} \text{Time : } O(nd) \\ \text{Space : } O(nd) \end{array} \quad \left. \right\}$$

expectation :-  $x_q \xrightarrow{\text{fast}} y_q \quad \left. \right\} \text{ low latency system}$

$$\begin{array}{l} k\text{-NN} \rightarrow \text{Simple implementation} \rightarrow O(nd) ; O(nd) \\ \qquad \qquad \qquad \text{space} \qquad \qquad \text{time} \\ \text{ex: } n \approx 364k \qquad \qquad \qquad \downarrow \qquad \qquad \qquad \downarrow \\ \qquad \qquad \qquad d \approx 100k \qquad \qquad \qquad 36 \text{ GB ram} \qquad \qquad \qquad 36 \text{ Billion Computation} \end{array}$$

which is not acceptable

### ② Decision Surface for K-NN changes as $k$ changes

' $k$ '  $\rightarrow$  hyperparameter

\* Curves / lines that separates +ve & -ve points are called **decision surfaces**.

In K-NN ; The Smoothness of decision surface increases as  $k$  increases

Case 1: If  $k=1$

There is high chance that based on the nearest neighbour the pt will be missclassified

Case 2: If  $k=5$

Majority rule will play a role in the classification

Case 3: If  $k=n$  or  $k$  so very large

If  $k$  so large then algorithm will be biased towards majority class.

$$\begin{array}{llll} \text{ca} & n_1 = +ve & n_2 = -ve & \text{where } n_1 > n_2 \qquad \qquad n = 1000 = k \\ & 600 & 400 & \end{array}$$

Notes on KNN :

### 8) How to handle Overfitting and Underfitting in KNN?

| <u>Overfit</u>   | <u>wellfit</u>        | <u>underfit</u>                         |
|------------------|-----------------------|---|
| $k = 1$          | $k = 5$               | $k = n$                                 |
| ↳ no mistake     | ↳ some mistakes       | ↳ every point belongs to majority class |
| ↳ non smooth     | ↳ smooth              | ↳ Doesn't cares                         |
| ↳ prone to noise | ↳ less prone to noise |   |

**Overfitting:** function or Algorithm tries to overtrain to not make any mistake on training data (includes noise or outlier)

**underfitting:** Algorithm or function underworks on training data and will not care about mistakes.

### 9) Need for Cross validation?

**Objective :-** Future unseen point  $x_q \xrightarrow{\text{accuracy}} y_q$

**Generalization :-** If it performs better on unseen data.

**Earlier :-**

① we split  $D_n$  (total dataset) in two parts

(a)  $D_{test}$       (b)  $D_{train}$

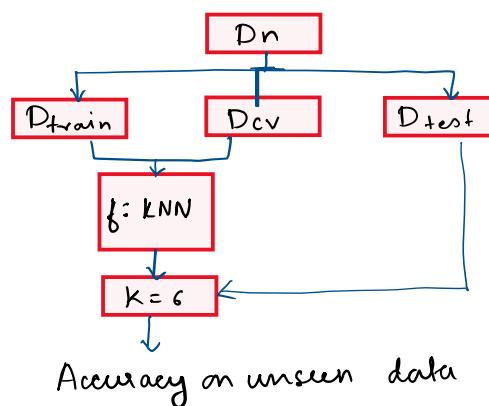
② we used  $D_{train}$  to train the data

③ we used  $D_{test}$  to test the data, based on the accuracy we decided the value of  $k=6$  as hyperparameter.

**But :-** Our model has seen both train and test data, so we are not sure that our function / Model will generalize better on future unseen data

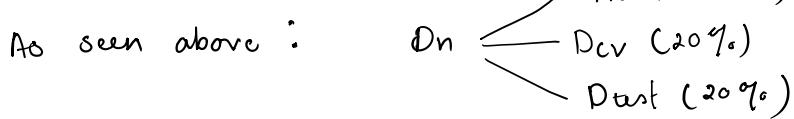
**Cross-validation :-**

$D_n$    
 ↗  $D_{train}$  (60%)  
 ↗  $D_{cv}$  (20%)  
 ↗  $D_{test}$  (20%)



Notes on KNN :

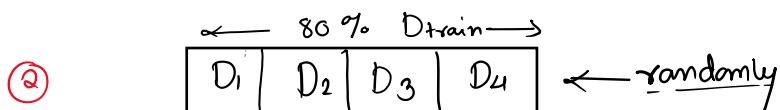
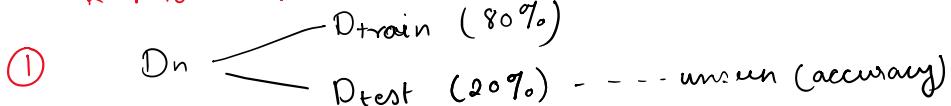
10) What is K-fold cross validation?



**problem:** Only 60% of the data is used for training

- \* More training data = the better model

k-fold CV :-



③

| Fold | Train           | CV    | accuracy |
|------|-----------------|-------|----------|
| K1   | $D_1, D_2, D_3$ | $D_4$ | A4       |
| K2   | $D_1, D_2, D_4$ | $D_3$ | A3       |
| K3   | $D_1, D_3, D_4$ | $D_2$ | A2       |
| K4   | $D_2, D_3, D_4$ | $D_1$ | A1       |

}  $k = 4$  fold validation  
avg =  $A_k$

\* Accuracy (Avg accuracy using)  
 $\downarrow$   
4 fold CV  $\rightarrow D_{test} \rightarrow$  Accuracy of 3NN on  $D_{test}$

\*  $k$  in k-fold  $\neq k$  in KNN

\* Time it takes to compute optimal  $k$  in KNN increases by k-fold times if we use k-fold CV

Advantages :-

- ①  $D_{train}$  &  $D_{cv}$  don't overlap perfectly when randomly sampled.
- ② If there are many +ve points from  $D_{train}$  in a region, then there is high-likely to find many negative (-ve) points from  $D_{cv}$  in that region
- ③ If there are very few points (+ve) in a region from  $D_{train}$ , then it is very unlikely to find (+ve) points from  $D_{cv}$  in that region.

Notes on KNN :

Question: How can we be really sure that we are not underfitting or overfitting?

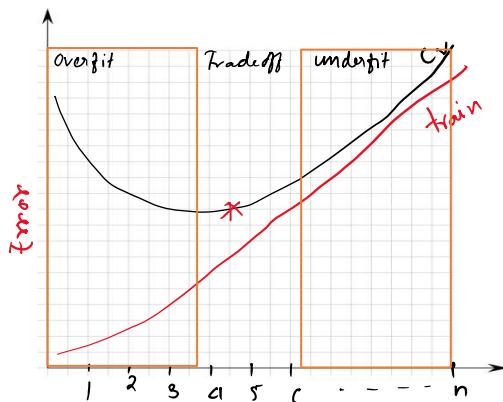
★ Accuracy =  $\frac{\text{No. of correctly classified pts}}{\text{Total no. of points}}$

★ Error =  $1 - \text{Accuracy}$

\* Now we can calculate Train error ; Validation error  
 $D_{\text{train}}$ ,  $D_{\text{cv}}$ ,  $D_{\text{test}}$

If Train Error ↑ { underfit  
CV Error ↑

If Train Error ↓ { overfit  
CV Error ↑



### 11) What is Time based splitting?

whenever time is available and if thing/behaviour and data changes overtime,  
Then we should use TBS (time based splitting)

TBS is alternate to Random Splitting:

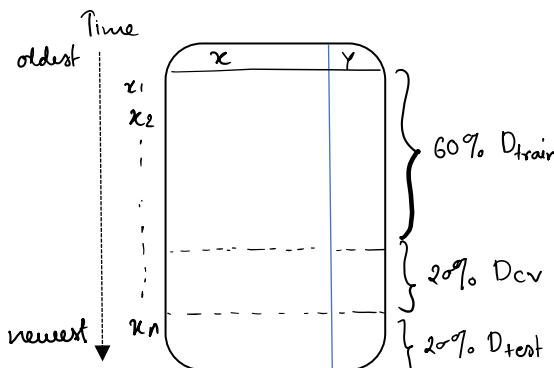
① Sort  $D_n$  in ascending order of time

② Split data in  $D_{\text{train}}$ ,  $D_{\text{cv}}$ ,  $D_{\text{test}}$   
based on time

$D_{\text{train}}$  → oldest

$D_{\text{cv}}$  → old

$D_{\text{test}}$  → newest



Example:



Notes on KNN :

### 12) Explain k-NN for regression?

$$\text{Regression} \rightarrow D = \left\{ (x_i, y_i)_{i=1}^n \mid x_i \in \mathbb{R}^d; y \in \mathbb{R} \right\}$$

for regression :-

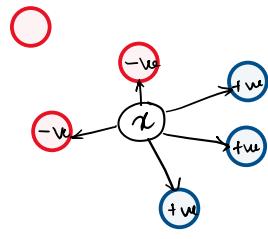
- 1.] given  $x_q$ , find k-nearest neighbour  
 $(x_1, y_1), (x_2, y_2), (x_3, y_3) \dots (x_k, y_k)$

2.]  $y_q \leftarrow y_1, y_2, y_3 \dots y_k$

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

$$x_q = \text{median } (y_i)_{i=1}^k \quad \leftarrow \text{less prone to outlier}$$

### 13) What is Weighted k-NN?



by majority vote

3+ve, 2-ve

$$\underline{y_q = +ve}$$

by weighted k-NN

$$w_1 (+ve) = 1.75$$

$$w_2 (-ve) = 15$$

$$\underline{y_q = -ve}$$

As seen in the example it is unfair to decide class just based upon majority vote

∴ weighted KNN

Weighted KNN gives more weight to the nearest point and less weight to farthest neighbour

$$w_i = \frac{1}{d_i}$$

where  $d_i$  is distance of the neighbour

$d_i \uparrow ; w_i \downarrow$

$d_i \downarrow ; w_i \uparrow$

| neighbour | class | distance | weight |          |
|-----------|-------|----------|--------|----------|
| $x_1$     | -ve   | 0.1      | 10     |          |
| $x_2$     | -ve   | 0.2      | 5      |          |
| $x_3$     | +ve   | 1.0      | 1      |          |
| $x_4$     | +ve   | 1.0      | 0.5    |          |
| $x_5$     | +ve   | 4.0      | 0.25   |          |
|           |       |          |        | { 15 }   |
|           |       |          |        | { 1.75 } |

Notes on KNN :

#### 14) How to build a kd-tree?

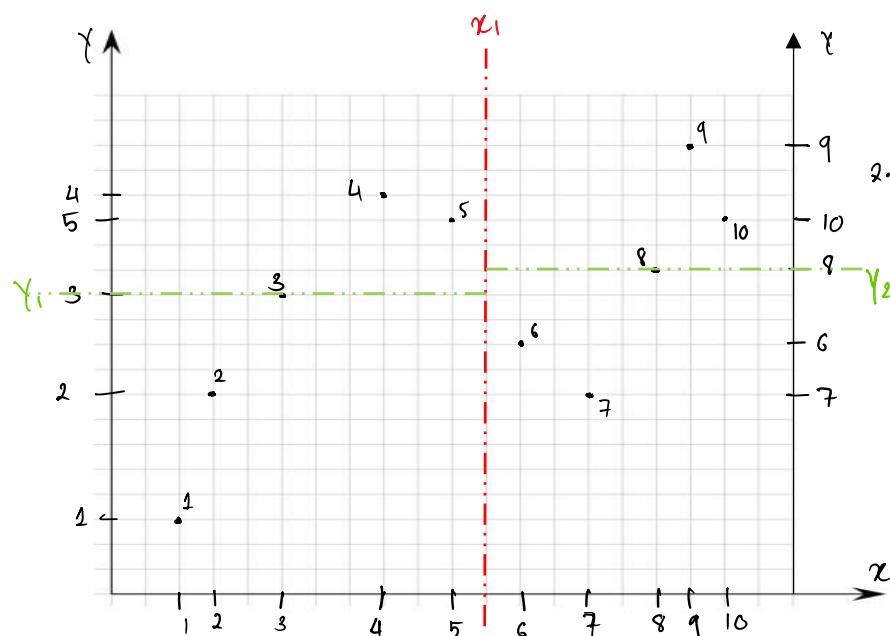
why?

As we see that time and space complexity of a simplified KNN is roughly  $O(nd)$

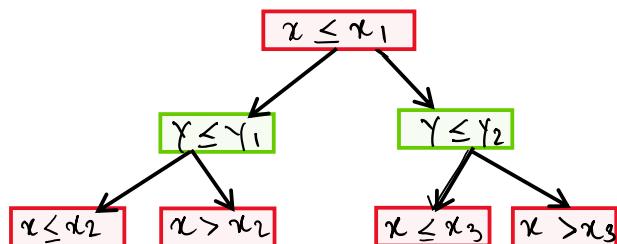
Can we somehow bring down time and space complexity from  $O(nd)$  to  $O(\lg n)$ ?  
using concept of Binary Search Tree (BST) algorithm

\* Such a model is called as a kd-tree

# Kd - tree :-



- 1.) pick  $x$ -axis, project points on  $x$ -axis, compute median and split using median
- 2.) pick  $y$ -axis, project points on  $y$ -axis, compute median and split using median
- 3.) Continue splitting alternate axis till we reach leaf node



Kd - tree :- breaking up spaces using axis-parallel lines/planes into rectangle, cube or cuboid.

Notes on KNN:

## Limitations of Kd-tree :-

- ① when point is near boundary we need to look at the adjoining region

$$2D = 4 \text{ regions}$$

$$3D = 8 \text{ regions}$$

$$d = 2^d$$

- ② when d is not small

$$d=10; 2^d = 1024$$

$$d=20; 2^d \approx \text{million}$$

$d \uparrow$ ; worst case # adjoining cell  $2^d \uparrow$

$$n = 1 \text{ million} \rightarrow 2^d \approx n$$

no longer  $O(\lg n)$

- ③ when d is not small {2,3,4,5}

Time complexity :  $O(\lg n)$  1-NN

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

## Disadvantages:-

- ① when d is 10, 11, 12 - - -

The time complexity increases dramatically

- ②  $O(\lg n)$  ← when d is small and data is uniformly distributed

But in real world data is not uniformly distributed

uniform =  $O(\lg n)$  ↗

simple implementation =  $O(n)$  ↗

## Extensions of Kd-tree :-

1. implicit

2. min/max kd-tree

3. Relaxed kd-tree

4. Ntree

5. Quad-tree

6. Ball-tree

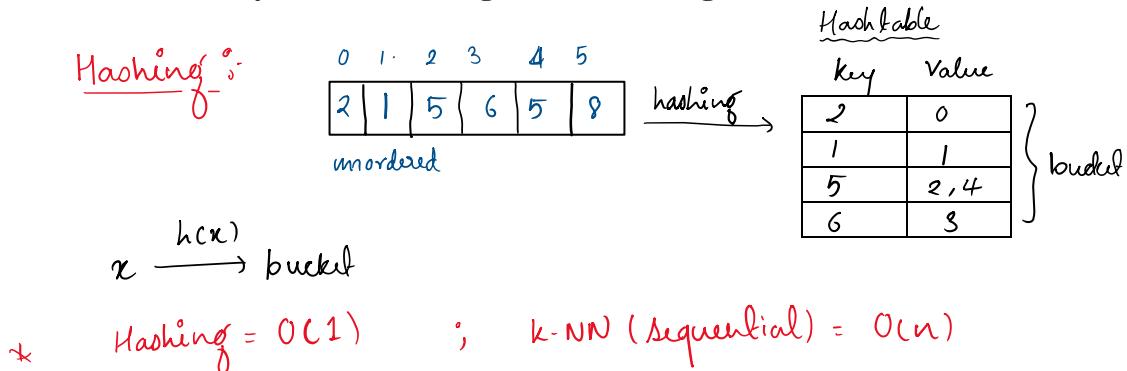
7. Oct-tree

8. R-tree

9. vantage point tree

Notes on KNN :

### 15) What is Locality sensitive Hashing (LSH)? (Hashing vs LSH?)



### Locality Sensitive Hashing :-

- ① It will create a hash function and will store all the nearest points in just one bucket.
- ② Given the new query point ( $x_q$ ) it will find nearest point and hash function will lead us to the bucket with all the nearest value.

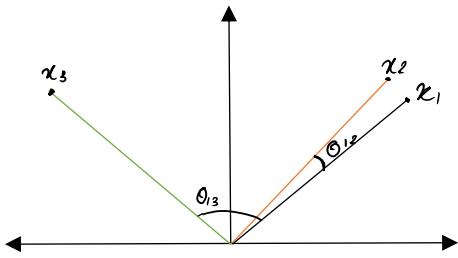
### 16) LSH for cosine similarity?

#### LSH for cosine similarity :-

Cosine Similarity = Angular Similarity

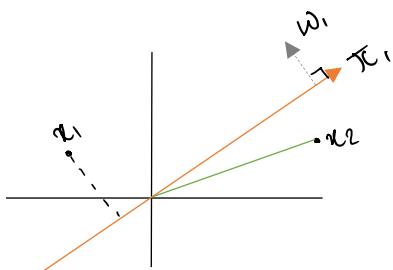
$\theta_{12} \approx 0 \Rightarrow x_1 \& x_2$  are very similar

$\theta_{13}$  is large  $\therefore x_1 \& x_3$  are very dissimilar



#### LSH (Randomized Algorithm)

It says → "It will not always give correct answer but it will always give correct answer with high probability."



$$w^T x_1 \geq 0 \quad (+ve)$$

$$w^T x_2 \leq 0 \quad (-ve)$$

$w_1 \rightarrow$  unit normal vector to plane  $\Pi_1$

#### 1. Random Hyperplane

$$w^T x = 0$$

$w$  : normal to the plane

$\therefore w^T x + b = 0$   
↳ as passing through origin

Notes on KNN :

Random numbers sampled from normal distribution ; where mean = 0 & variance = 1 ( $N(0, 1)$ )

$$w = \begin{bmatrix} 0 & 1 & \dots & 1 & \dots & d-1 \\ | & | & \dots & | & \dots & | \\ 0 & 1 & \dots & \text{---} & \dots & \text{---} \end{bmatrix}$$

$\underbrace{\hspace{1cm}}_{d \text{ dim}}$

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

$m = \# \text{ hyperplanes}$

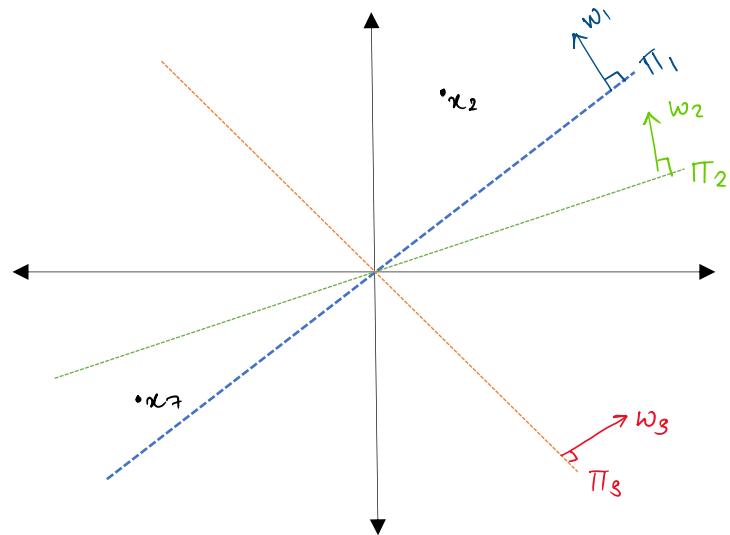
$\nabla$  slice  $\rightarrow m$ -dim vector

$$h(x) = \begin{bmatrix} 1 & 2 & 3 \\ +1 & +1 & +1 \end{bmatrix}$$

$\downarrow \text{Sign of } w_1^T x \quad \downarrow \text{Sign of } w_2^T x \quad \downarrow \text{Sign of } w_3^T x$

$$\text{ex: } h(x_2) = +1, +1, +1$$

$$h(x_7) = +1, -1, -1$$



Hash table :-

$$x \rightarrow h(x) \quad \begin{bmatrix} 1 & 2 & 3 \end{bmatrix}$$

Vector of size 'm' (here 3)

| key        | values               |
|------------|----------------------|
| +1, +1, +1 | $x_1, x_2, x_3, x_4$ |
| +1, -1, -1 | $x_6, x_7, x_8$      |
| .          | .                    |

Time and Space LSH :-

space :  $O(dn)$   $\rightarrow$  every point has to be stored

Time :  $O(md)$   $\rightarrow$  time to compute hash table

Time Complexity of Querying

$$x_q : h(x_q) \quad \begin{bmatrix} | & | & | & \dots & | \end{bmatrix}$$

$O(md)$

$x_q \rightarrow n'$  element in the bucket/slice  $O(md + n'd)$

If  $n'$  is small then  $n' \approx n$   $O(md)$

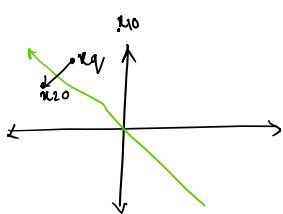
Typically  $m \approx \lg(n)$

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

Notes on KNN :

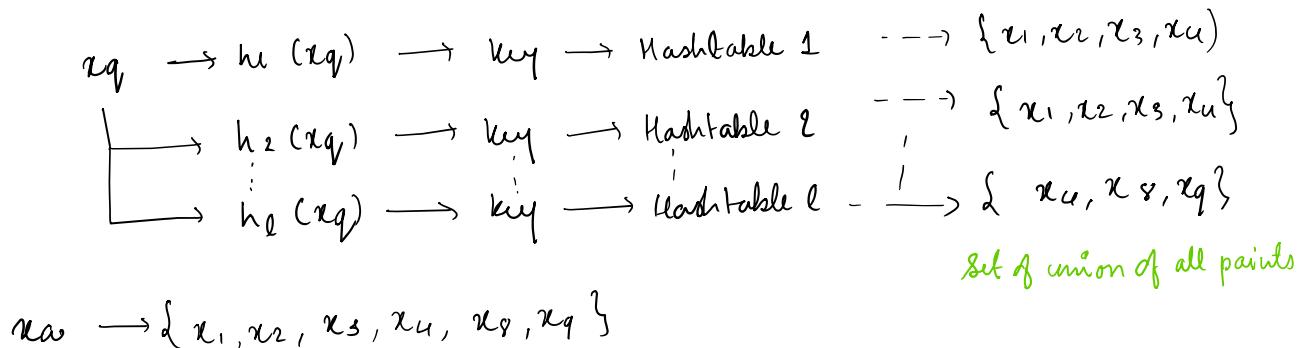
problems with LSH :-

You could miss the nearest neighbour in a cosine-similarity as your measure



$x_1$  and  $x_2$  will go in different bucket of hashtable and we will miss.

Fix :- Create multiple sets of hyperplane and repeat multiple times



$$x_q \rightarrow \{x_1, x_2, x_3, x_4, x_5, x_6\}$$

Now we can compute cosine similarity for nearest neighbours)

Time Complexity :-  $O(m d L)$        $m = O(\lg n)$

As no. hyperplane increases ; slices  $\uparrow$  ; # of points per slice  $\downarrow$   
 $\therefore$  we should choose 'm' wisely.

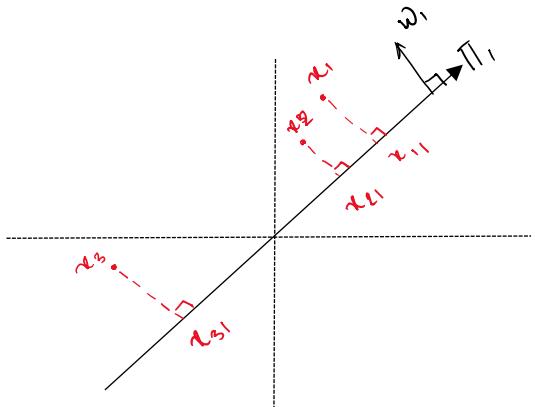
17) LSH for Euclidean distance?

$x_{11} \rightarrow$  projection of  $x_1$  on  $T_{11}$

$x_{21} \rightarrow$  projection of  $x_2$  on  $T_{11}$

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

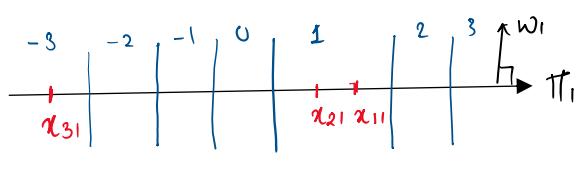
LSH → similar or closer points should go to the same bucket



Breaking plane into pieces/regions

# hyperplanes = m

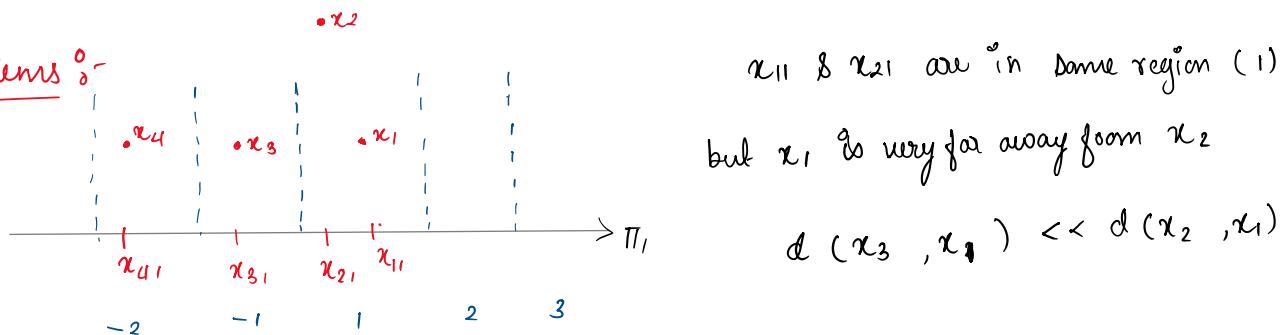
$$h(x_{11}) = \begin{array}{|c|c|c|c|c|c|c|} \hline & \text{region} & & & & & \\ \hline & 2 & | & | & - & - & - \\ \hline & T_{11} & T_{12} & & & & \\ \hline & m & & & & & \\ \hline \end{array}$$



Notes on KNN :

$$h(x_{31}) = \begin{array}{c|c|c|c|c|c} & 1 & 2 & 3 & \cdots & m \\ \hline -3 & | & | & | & \cdots & | \\ \hline x_1 & x_2 & x_3 & & & x_m \end{array}$$

Problems :-



Solution : Create multiple and random hyperplane

Note : LSH is not a perfect but probabilistic / randomized algorithm.

# KNN : Probabilistic Class Labels

example: 2 class classification

$$y^i \in \{0, 1\}$$

$x_q$ : 4 -ve pts ; 3 +ve pts ( $y_q = \text{-ve}$ )  $\leftarrow$  majority vote

$x_{q'}$ : 7 (-ve pts) ; ( $y_{q'} = \text{-ve}$ )

Quantify this uncertainty

$$P(y_q = \text{-ve}) = 4/7 = \frac{\# \text{-ve pts}}{\# \text{total pts}}$$

$$P(y_{q'} = \text{-ve}) = 7/7 = 100\%$$

similarity

$$P(y_q = \text{+ve}) = 3/7$$

$$P(y_{q'} = \text{+ve}) = 0/7$$

Probabilistic class label

$$x_q \rightarrow y_q \rightarrow \begin{cases} P(y_q = \text{+ve}) \\ P(y_q = \text{-ve}) \end{cases} \leftarrow \text{certainty}$$