

K-Nearest Neighbours

- 1 -

How classification works?

~~x_i : Text Vector~~
prediction of a class whether +ve or -ve

$$\boxed{y = f(x)}$$

Notations

Each vector by default is column vector



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

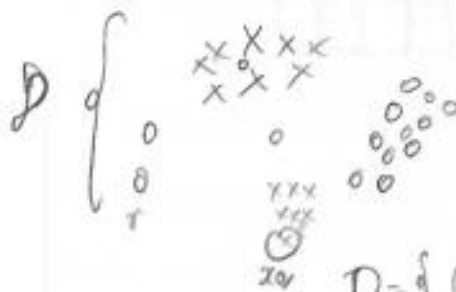
such that

-ve +ve

K-Nearest Neighbours:

2D - toy datasets:

Binary classification

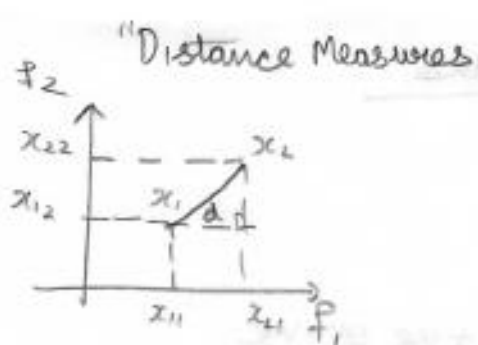


$+$: +ve data pt.
 o : -ve data pt.

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

-ve +ve

Take pt close to this pt $x_2 \rightarrow y_2$



$$f_1 \quad f_2$$

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

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

d = len of shortest line from x_1 to x_2

Euclidean distance

$$d = \sqrt{(x_{21} - x_{11})^2 + (x_{22} - x_{12})^2}$$

$$= \underbrace{\|x_1 - x_2\|}_{\text{length}}$$

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

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

$\|x_1 - x_2\|_2 \rightarrow L2 \text{ norm}$

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

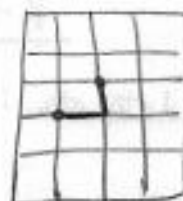
◇ Manhattan distance:

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

L_1 -norm vector $(x_1 - x_2)$

$$\|x_1 - x_2\|_1$$

$$\|x_1\|_1 = \sum_{i=1}^d \overbrace{|x_{1i}|}^{\text{absolute}}$$



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 distance \rightarrow Euclidean distance

$p = 1 \rightarrow$ Minkowski distance \rightarrow Manhattan distance

$$\|x_p\| = \sum_{i=1}^d (|x_i|^p)^{1/p} \quad \begin{array}{l} p \neq 0 \\ p > 0 \end{array}$$

$$\begin{aligned} \text{Euc. dist}(x_1, x_2) &= L_2 \text{ norm of } (x_1 - x_2) \\ &= \|x_1 - x_2\|_2 \end{aligned}$$

Distance between 2 pts \rightarrow Euclidean $(x_1, x_2) = \|x_1 - x_2\|$
 Norms is for a vector $\rightarrow L_2, L_1$ norms

* Hamming distance

Test passing and boolean Vectors \rightarrow Binary Base

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

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

Hamming dist $(x_1, x_2) = \# \text{ locations / dimensions where binary vectors differ.}$

$$\text{In } x_1 \& x_2 \text{ Hamming} = 3$$

$$x_1 = a b c a d e f g h i k \leftarrow \text{Gene code / seq}$$

$$x_2 = a c b a d e g f h i k \leftarrow$$

$$\text{Hamming dist}(x_1, x_2) = 4$$

→ Cosine distance & Cosine Similarity

Similarity $\xrightarrow{\downarrow \text{dec}} \xrightarrow{\uparrow \text{inc}}$ distance x_1, x_2
 (opposite)
 $\uparrow \text{inc} \xrightarrow{\quad} \downarrow \text{dec}$

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

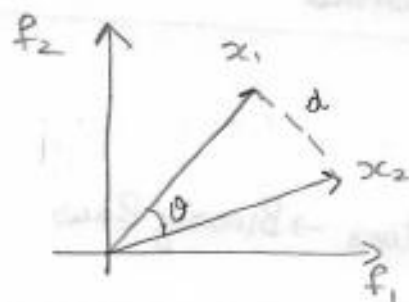
$[-1, 1]$

cos-sim^{distance}_n(x_1, x_2) very similar = +1

cos-sim^{distance}_n(x_1, x_2) very dissimilar = -1



→ what is cosine similarity

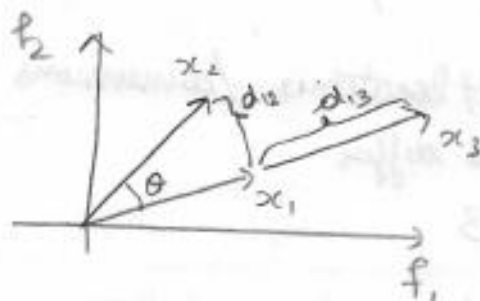


$d = \text{euclidean distance}$

cos-sim^{similarity}_n(x_1, x_2) = $\cos \theta$

θ : angle between x_1 & x_2

? Difference between euclidean distance and cosine similarity



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

cos-sim(x_1, x_3) = $\cos 0^\circ = 1$

but euclidean distance
 $d_{13} > d_{12}$

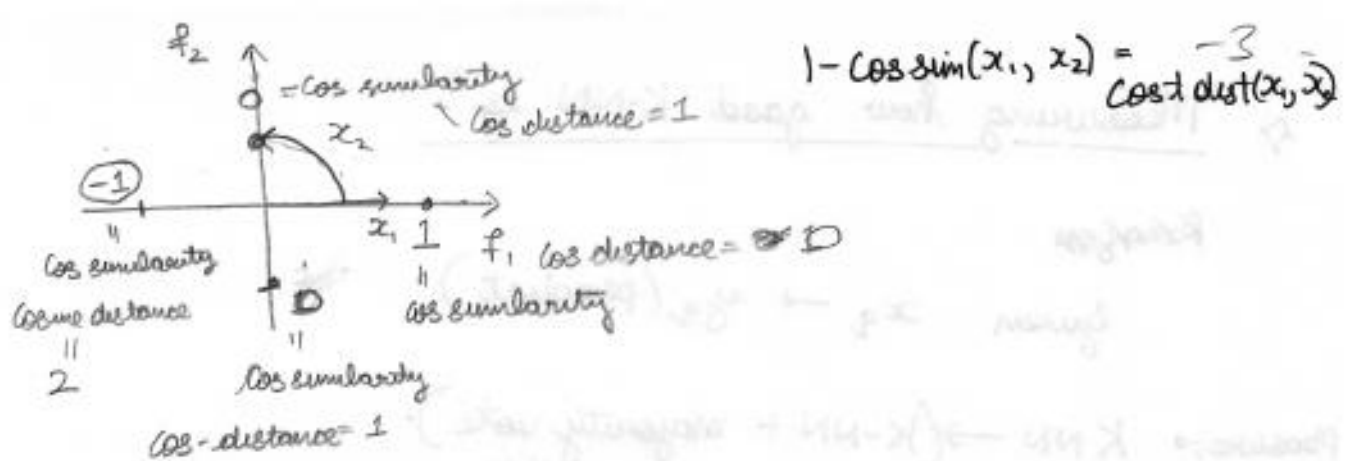
→ cos-distance(x_1, x_3) = $1 - 1 = 0$

cos-distance(x_1, x_2) = $1 - \cos \theta$

cos distan(x_1, x_3) < cosdis(x_1, x_2)

NOTE

$d_{13} > d_{12}$ on the other hand



| angle (x_1, x_2) | dist |
|-------------------------|-------------------------------|
| $0 - 90^\circ$ | $\rightarrow 0 \rightarrow 1$ |
| $90^\circ - 180^\circ$ | $\rightarrow 1 \rightarrow 2$ |
| $180^\circ - 270^\circ$ | $\rightarrow 2 \rightarrow 1$ |
| $270^\circ - 360^\circ$ | $\rightarrow 1 \rightarrow 0$ |

$$\begin{aligned}
 1 - \cos 90^\circ &= 1 \text{ (cos distance)} \\
 1 - \cos 180^\circ &= 2 \text{ (cos distance)} \\
 1 - \cos 270^\circ &= 1 \text{ (cos distance)} \\
 1 - \cos 360^\circ &= 0 \text{ (cos distance)}
 \end{aligned}$$

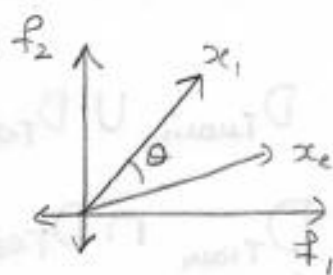
$$\cos \theta = \frac{x_1 \cdot x_2}{\|x_1\|_2 \|x_2\|_2}$$

L2 norm

① If x_1 & x_2 are unit vector

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

$$\cos \theta = x_1 \cdot x_2$$



? Relationship between Euclidean & cos-similarity/distance

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

If x_1 & x_2 are unit vectors

Using

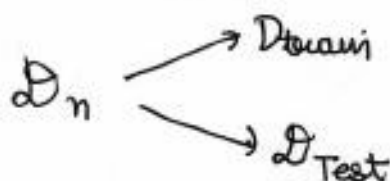
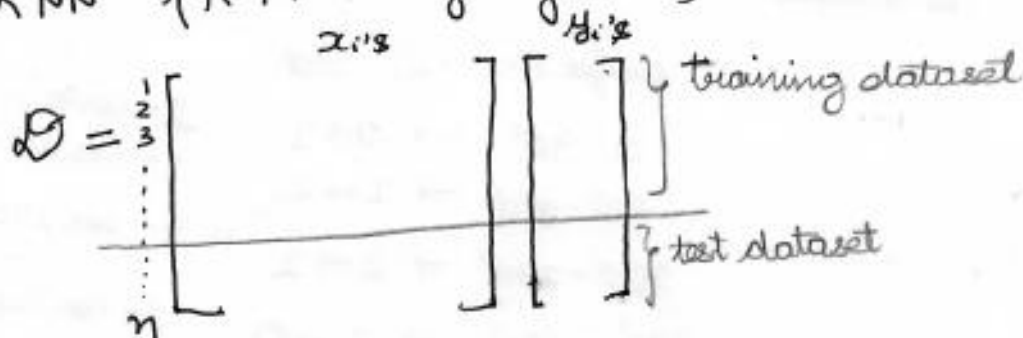
$$\begin{aligned}
 \text{From } \|A - B\|^2 &= (A - B)^T (A - B) = \|A\|^2 + \|B\|^2 - 2 A^T B \Rightarrow \|A\| = \|B\| = 1 \text{ (cosine similarity)} \\
 &= 2 - 2 A^T B = 2(1 - A^T B) = 2(1 - \cos(A, B))
 \end{aligned}$$

☆ 'Measuring how good K-NN is'

Assumption

Given $x_q \rightarrow y_q$ (predict)

Measure: \rightarrow KNN \rightarrow (K-NN + majority vote)



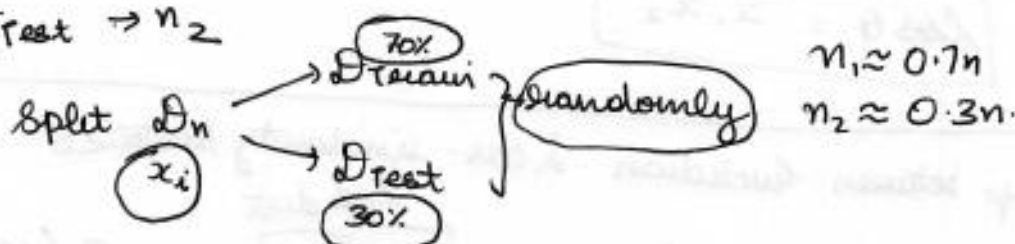
$$D_{Train} \cup D_{Test} = D_n$$

$$D_{Train} \cap D_{Test} = \phi$$

For example

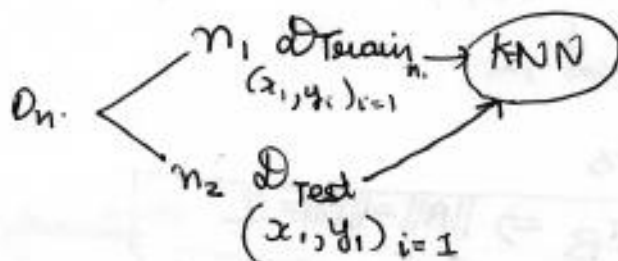
$$D_{Train} \rightarrow n_1, \quad n_1 + n_2 = n$$

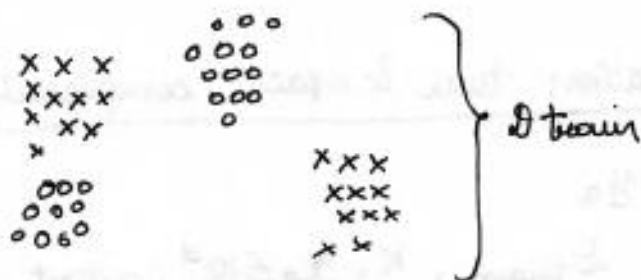
$$D_{Test} \rightarrow n_2$$



$$n_1 \approx 0.7n$$

$$n_2 \approx 0.3n$$



D_{train} 

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

$x_q = (x_i)$
 $y_q = \text{Blue (+ve)}$

for each point in D_{test}

→ $x_q = pt$

→ use D_{train} & KNN to predict

y_q

Algorithm:

count = 0;

for each pt in D_{test} :

$x_q = pt$

use D_{train} & K-NN to determine y_q

if $y_q == y_{pt}$

count += 1

number of points for which D_{train} + KNN gave correct class label.

$$\text{Accuracy} = \frac{\text{count}}{n_2}$$

↙
pts in D_{test}

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

Accuracy = 0.91 ⇒ 91% of times

$x_q \rightarrow y_q$

Test / Evaluation time & space complexity :

$$x_q \rightarrow y_q$$

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

KNN pts = [] $\rightarrow n$ pts (samples); d -dimensional (features)

for each x_i in D_{train} :

$O(d)$ - Compute $d(x_i, x_q) \rightarrow d_i$

$O(k)$ - Keep the smallest k -distances $\rightarrow (x_i, y_i, d_i)$ \rightarrow KNNpts []

Time comp-
lexity

K is small $\rightarrow 5$ or 10

Count-pos = 0; Count-neg = 0

for each x_i in KNNpts

if y_i is +ve

Count-pos += 1

else

Count-neg += 1

if Count-pos > Count-neg
return $y_q = 1 \rightarrow +ve$

else

$y_q = 0 - ve$

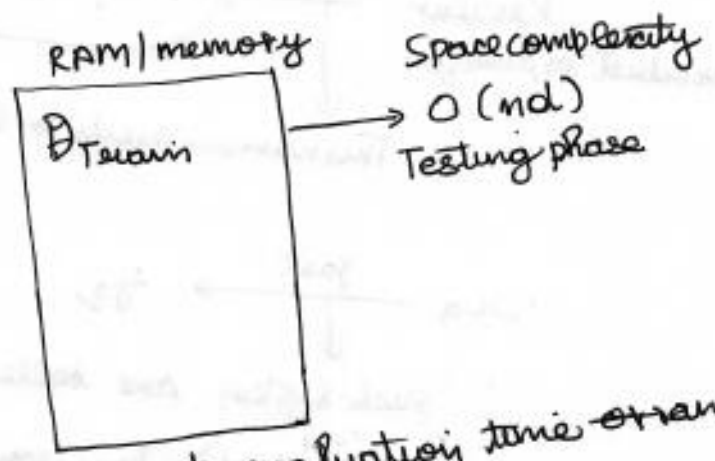
Time complexity :- $O(nd) + O(1) + O(1)$
 $O(nd)$

if d is small
if $d \ll n$

$O(n) \rightarrow$ Time

deploy:

$$x_q \rightarrow y_q$$



It will take lot of memory. at evaluation time and run time.

→ evaluation / Test phase

Time complexity :- $O(nd)$

Space complexity :- ~~sample~~ space that is need to evaluate $O(nd)$

KNN limitations

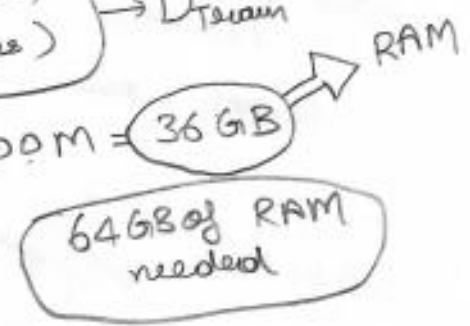
(Amazon) Fine Food reviews:

Production → live system

$$\begin{aligned} \text{Time} &:- O(nd) \\ \text{Space} &:- O(nd) \end{aligned}$$

$$\begin{aligned} n &\approx 364K \text{ (sample)} \\ d &\approx 100K \text{ (features)} \end{aligned} \rightarrow D_{train}$$

$$n \times d = 36,400m = 36GB$$



① Large space complexity

2.) Time Complexity : 36 Billion's Computations

Revenue of product in amazon $\xrightarrow{1 \text{ msec}}$ +ve / -ve
 This number needs to be small.

$x_q \xrightarrow{\text{fast}} y_q$

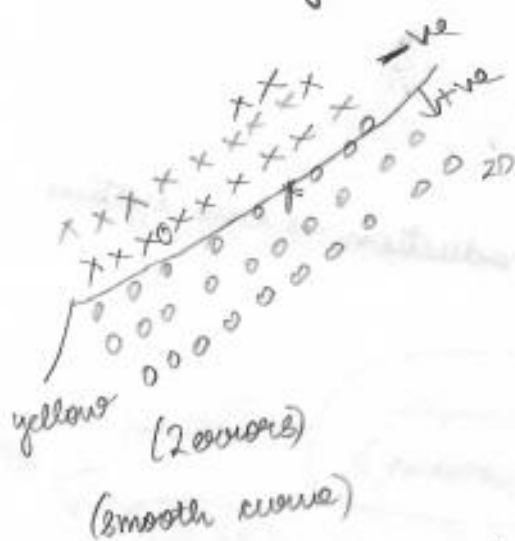
such system are called low latency system.

$O(nd)$ $O(nd)$ Time complexity is high.

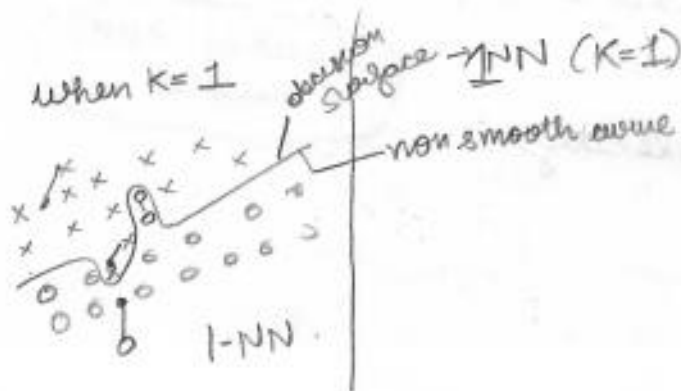
Decision surface for KNN as K changes

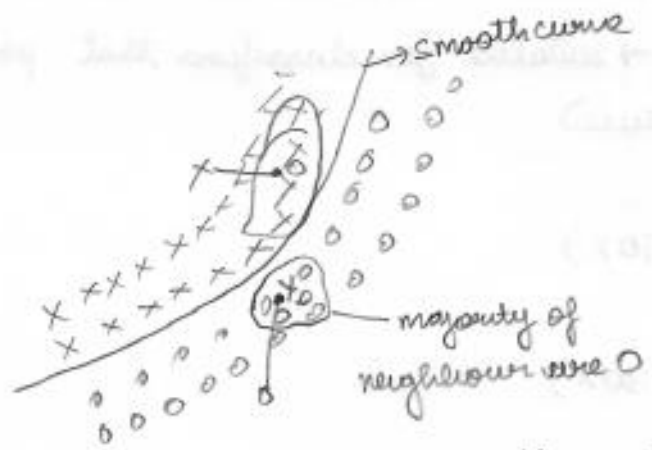
How to pick right K?

'K' is also referred as hyperparameter.



decision surfaces yellow and green.





Let $k=5$.

S-NN

Nearest neighbour
majority vote

NOTE: k increases smoothness of curve increases

In k nearest neighbour. the smoothness of the decision surfaces increases as k increases

lets say $k=n$

$k=1, 2, 3, 4, \dots, n$

max value of $k=n$.

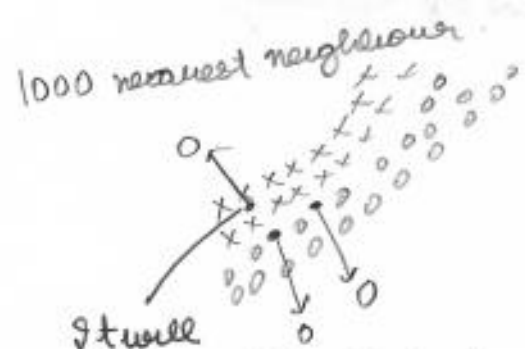
n is total number of pts

n = Total number of points

1000 \leftarrow
 $600 \leftarrow n_1 = +ve$
 $400 \leftarrow n_2 = -ve$

$$n_1 + n_2 = n$$

lets say $n_1 > n_2$
 $600 > 400$



$k=n=1000$
 $n_1 = +ve = 600$
 $n_2 = -ve = 400$
 $n_1 > n_2$
 $600 > 400$

NOTE: when k becomes n than the everything becomes majority class.

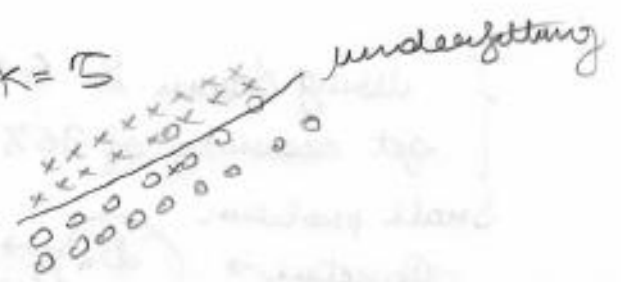
if $n_1 > n_2$ so it will be considered $+$ instead of x .

Overfitting & underfitting

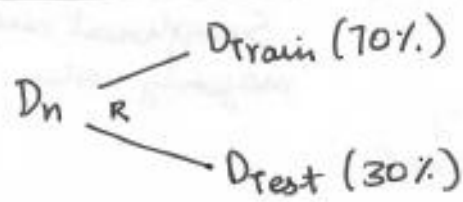
$k=1$



$k=5$



Need For Cross Validation → suitable for classifiers that predict labels (positive or negative)
How to determine k

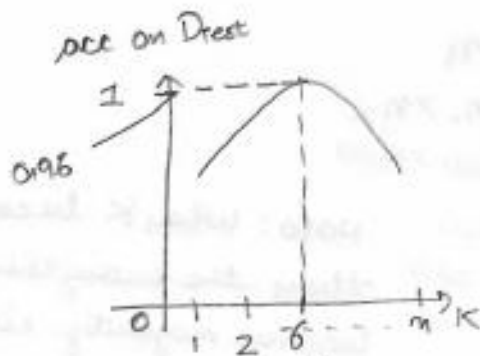


One idea

| | Train | accuracy Dtest |
|-----|--------|----------------|
| K=1 | Dtrain | 0.78 |
| K=2 | | 0.82 |
| K=3 | | 0.85 |

num of correctly classified files

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



(Typically)

Calculating value of k using test data

best accuracy on my Dtest when using Dtrain as training data

K=6 (6-NN)

96%

{ using Dtrain & 6-NN on amazon food review dataset
 get accuracy of 96%

Small problem

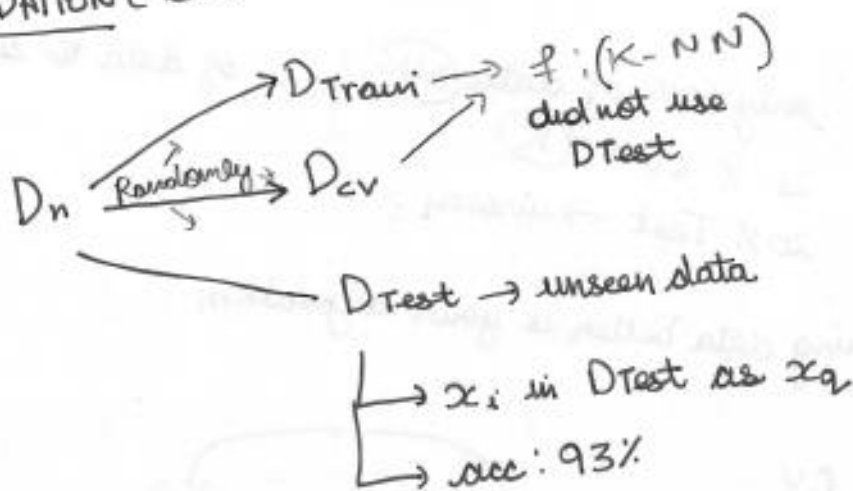
Objective → $D_n \rightarrow f(\text{Train}) (6\text{-NN})$
 $\rightarrow \text{accuracy } (D_{\text{testing}})$

future, unseen point $x_q \rightarrow y_q$

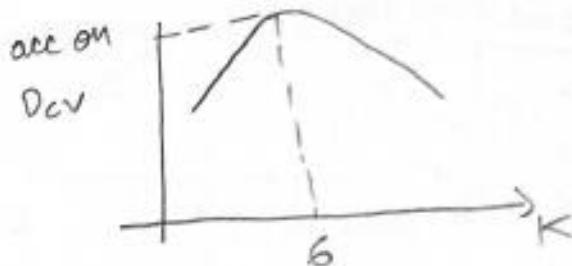
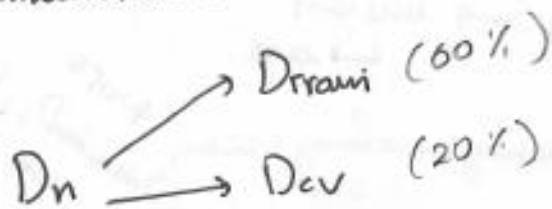
Objective \rightarrow accurate on future, unseen point ~~than it~~

Generalization \rightarrow when algorithm does well on future unseen point not only for D-test.

CROSS-VALIDATION (CV)



Yes { Can I now say that 6-NN has an accuracy of 93% on unseen data?



NOTE first we were finding value of K using test set but now we use cross validation to find best value of K

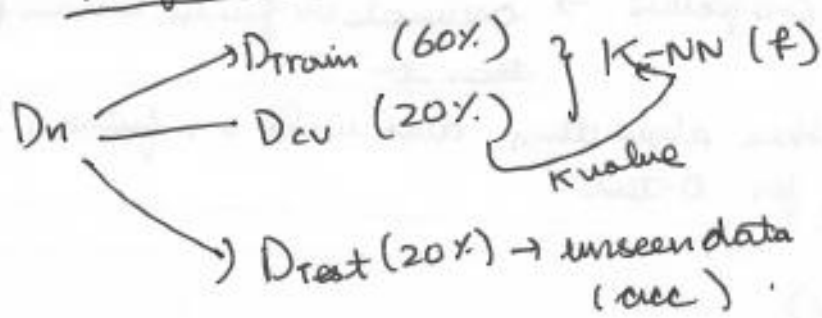
$\rightarrow K=6$

D_{test}

93% accuracy on test data on future on unseen data

What is right K ?

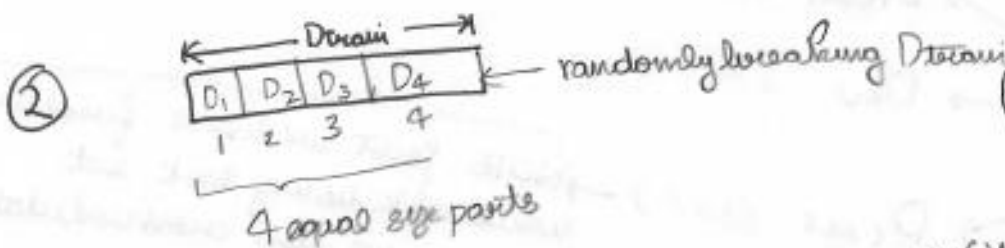
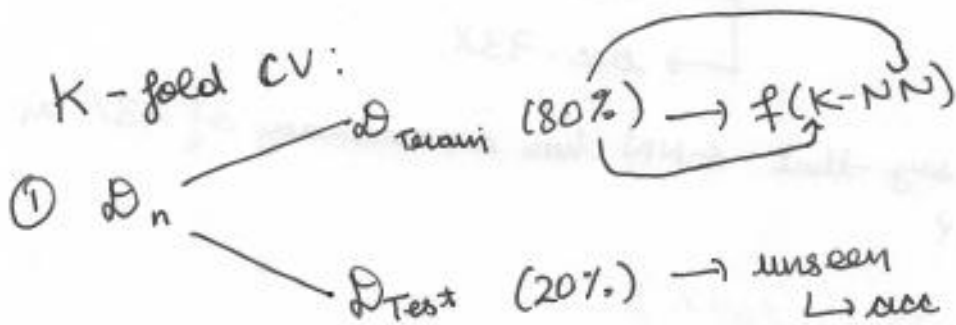
K-fold cross validation



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

Rule: More training data better is your algorithm.

K-fold CV:

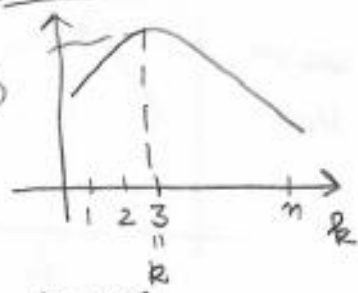


③

| | 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 |

avg (a_1, a_2, a_3, a_4)

acc on CV dataset



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

D_1, D_2, D_3, D_4
 each of 20% data

4 times

| | train | CV | acc on CV |
|-----|---------------|-------|-----------|
| 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 of (a_1, a_2, a_3, a_4)
 $a_{K=2}$

4 times

| | train | CV | acc on CV |
|-----|---------------|-------|-----------|
| K=3 | $D_1 D_2 D_3$ | D_4 | a_4 |
| K=3 | $D_1 D_2 D_4$ | D_3 | a_3 |
| K=3 | $D_1 D_3 D_4$ | D_2 | a_2 |
| K=3 | $D_2 D_3 D_4$ | D_1 | a_1 |

avg of (a_1, a_2, a_3, a_4)
 $a_{K=3}$

acc on CV data set

→ This is called 4 fold - cross validation

Suppose $K=3$ you got the best accuracy.

$f: 3\text{-NN}$

↑
 D_{train}

→ avg-acc on 4 factor.

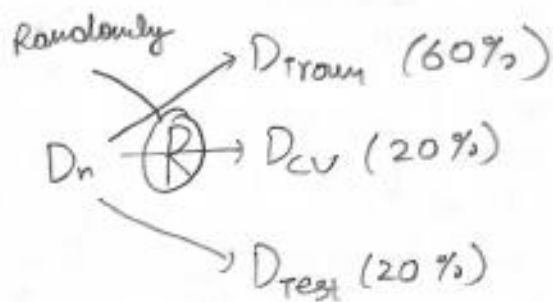
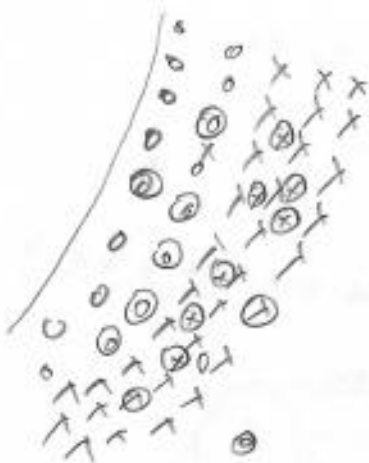
K' fold CV → D_{train} → NN
 → K in $K\text{-NN}$

What is right K'

$K'=4; K=10; K'=100$

↓
 rule of thumb 10 fold CV.

Visualizing Train, CV & test datasets

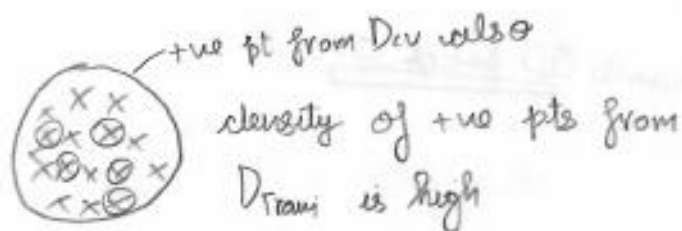


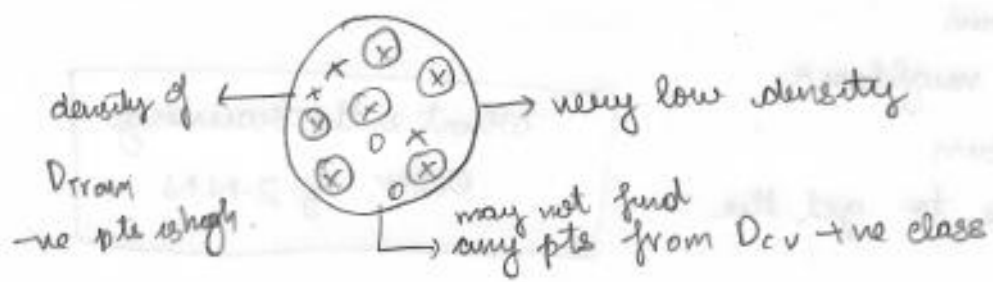
x : +ve datapoint in D_{train}
 \odot : -ve datapoints in D_{train}

\otimes → +ve pt from D_{cv}
 \odot → -ve pt from D_{cv}

NOTE

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





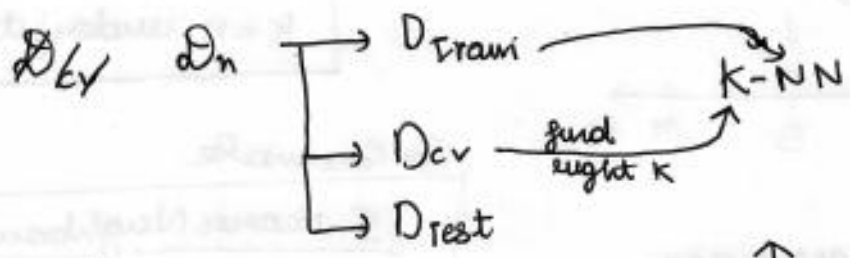
How to determine overfitting vs underfitting:

$\{ k \text{ fold } CV \text{ or } D_{cv} \rightarrow \text{best } k \rightarrow \begin{cases} \text{neither overfit} \\ \text{nor underfitting} \end{cases}$

Q How can we be sure that we are not underfitting nor overfitting

$$\text{accuracy} = \frac{\# \text{ correctly class pts}}{\text{Total \# pts}} = 0.93 \text{ (Suppose)}$$

$$\text{error} = 1 - \text{accuracy} = 1 - 0.93 = 0.07 = 7\%$$



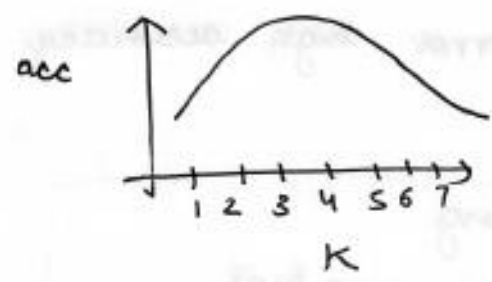
error on $D_{cv} = 1 - \text{accuracy}$

acc. on D_{cv}

for each pt in D_{cv}

$x_q = \text{pt}$

y_q y_i



Training error:

$D_{train} \rightarrow NN$

$D_{cv} \rightarrow k$

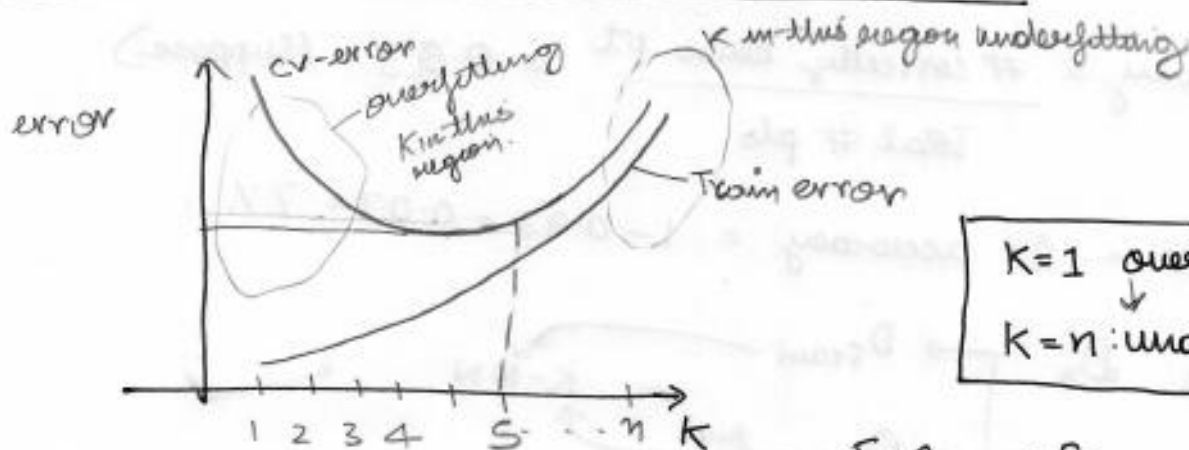
$f(K-NN)$

for each x_i in D_{train}

- find 2 nearest neighbours to x_i from D_{train}
- majority vote to get the class label
- if $y_i == \text{class label}$
 accurate
- else
 error.

what is the training error of 2-NN

error; train error, validation - error



$K=1$ overfit
 \downarrow
 $K=n$ underfit

low error high accuracy

Underfitting

Train error \rightarrow high
Cross validation error \rightarrow high

Overfitting

Train error \rightarrow low
Cross validation error \rightarrow high

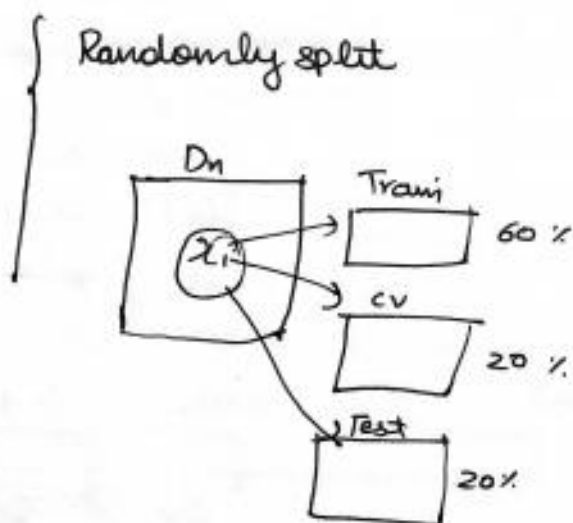
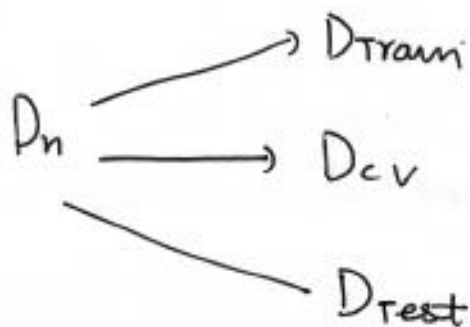
For example

5-Nearest Neighbour

lowest error from graph

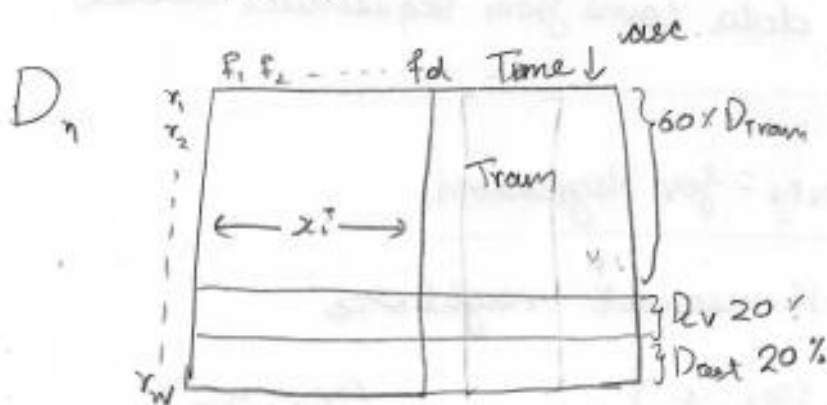
K is selected using cross validation with highest accuracy.

Time Based splitting



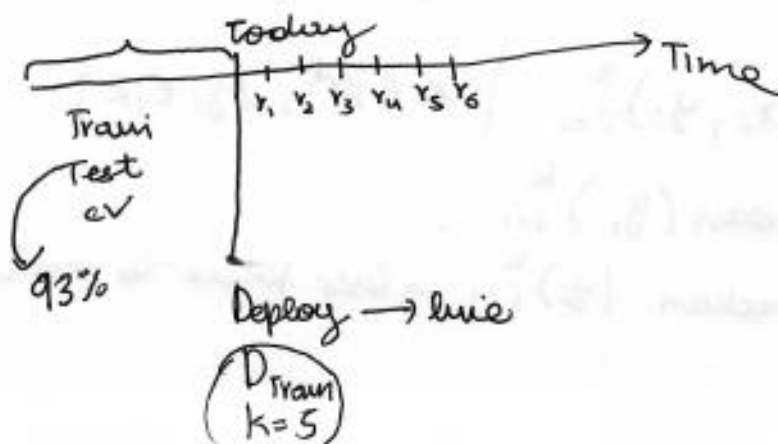
☆ Time based Splitting

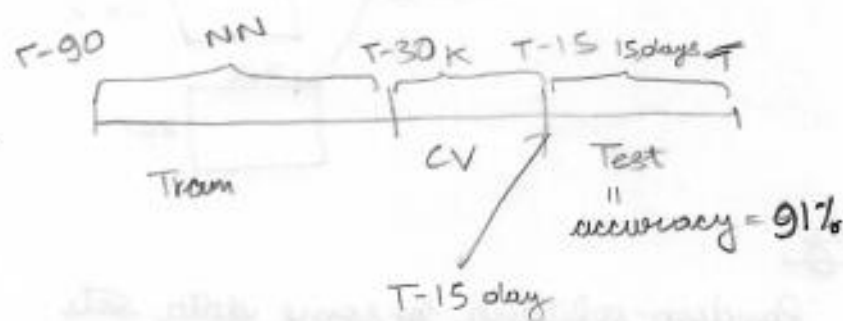
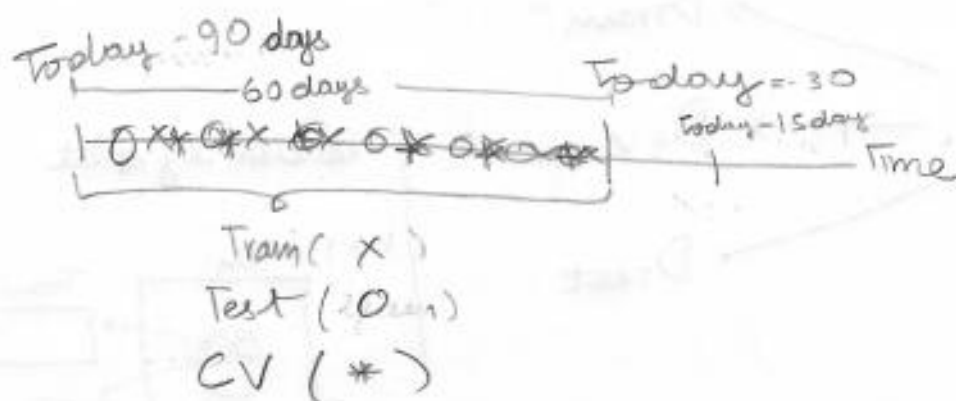
better than Random-splitting for some data sets



① sort D_n as asc order of time

(TBS) → can only be done if Time stamp is given





TBS: In this Test data come from sequential model

KNN - for Regression

given x_q , find K-nearest neighbors

① $(x_1, y_1), (x_2, y_2), \dots, (x_k, y_k)$

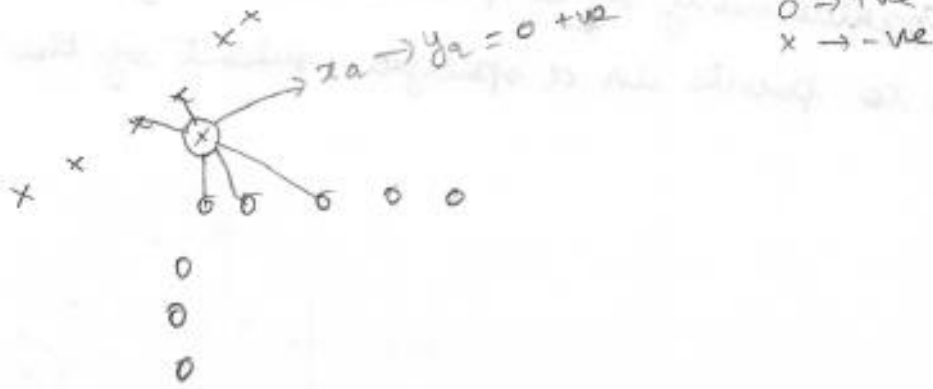
② $y_q \leftarrow y_1, y_2, y_3, \dots, y_k$

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

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

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

Weighted K-NN



SNN $K=5$

$x_a \rightarrow x_1, y_1, d_1, \begin{matrix} \text{Let distance} \\ 0.1 \\ 0.2 \end{matrix} \begin{matrix} -ve \\ -ve \end{matrix} \rightarrow \text{distance more near}$

$\begin{cases} x_3, y_3, d_3, 1.0 +ve \\ x_4, y_4, d_4, 2.0 +ve \\ x_5, y_5, d_5, 1.0 +ve \end{cases}$

One way

$$w_i = \frac{1}{d_i} \text{ one way}$$

$x_a \rightarrow$

| x_i | distance (d_i) | w_i | $w_i \rightarrow$ |
|------------|--------------------|-------|-------------------|
| $x_1, -ve$ | 0.1 | 10 | |
| $x_2, -ve$ | 0.2 | 5 | |
| $x_3, +ve$ | 1.0 | 1 | |
| $x_4, +ve$ | 2.0 | 0.5 | |
| $x_5, +ve$ | 1.0 | 1 | |

$\begin{cases} d_i \uparrow, w_i \downarrow \\ d_i \downarrow, w_i \uparrow \end{cases}$

majority rule

$y_a = +ve$ which is wrong

but by using weighted K-NN we can predict right prediction



Voronoi diagram

Used partitioning of a plane into regions based on distance to points in a specific subset of the plane.



Voronoi cells

cell : region.

Binary Search Tree (BST)

K. NN : $\begin{cases} O(n) & \text{if } d \text{ is small} \rightarrow \text{Time} \\ O(n) & K \text{ is small} \\ \downarrow \\ \text{space} \end{cases}$

Time : $\begin{matrix} O(n) \rightarrow O(\lg n) \\ n = 1024 \rightarrow \lg(n) = 10 \end{matrix}$ \rightarrow Kd-tree can be used to reduce it

prob: given a sorted array, find a number is present in the array or not.

| | | | | | | | | |
|---|---|---|---|----|----|----|----|--------|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | |
| 1 | 2 | 6 | 8 | 12 | 15 | 17 | 18 | Sorted |

10 \rightarrow median

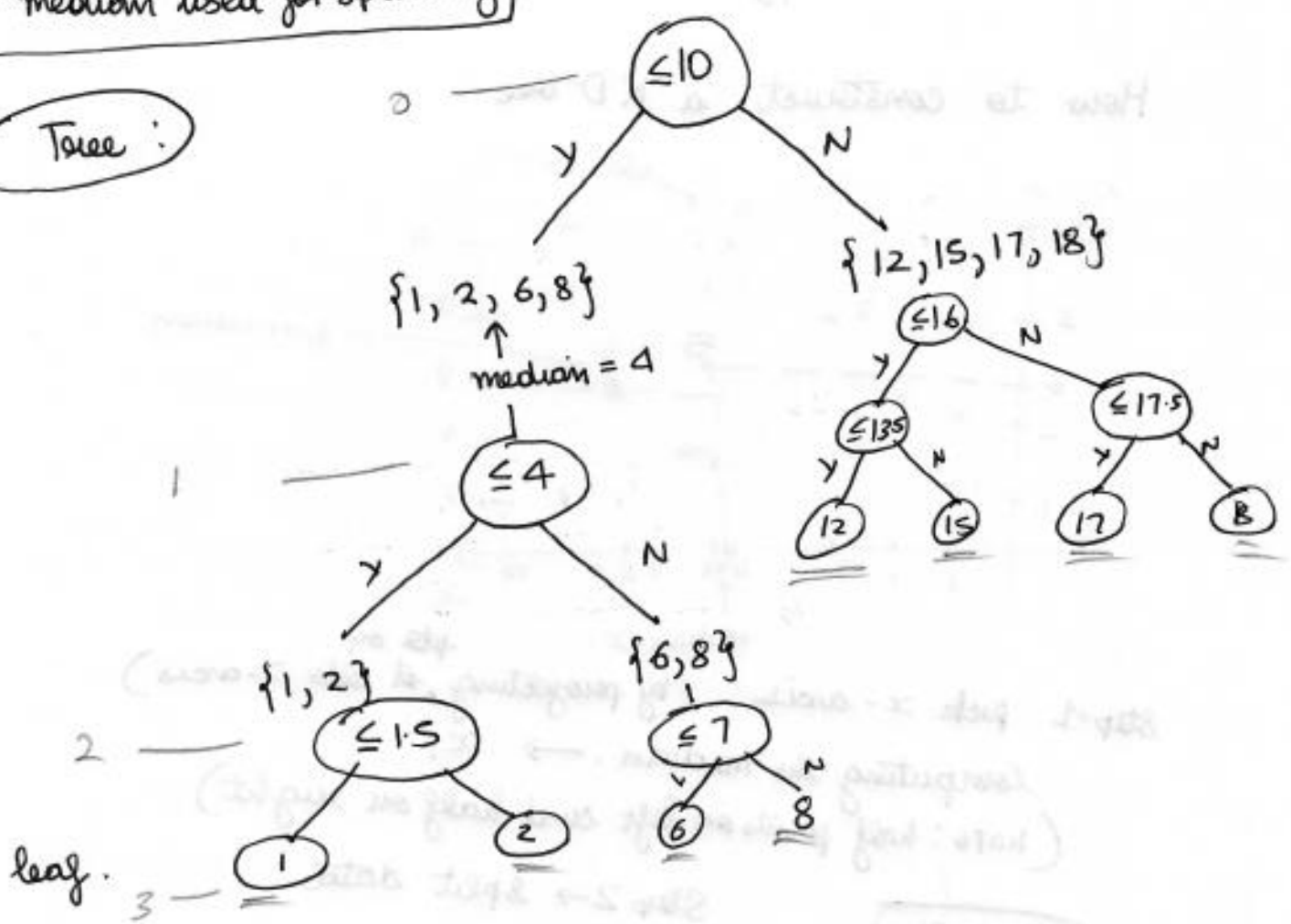
1.) Construct a BST

NOTE median used for splitting

median = 10

Flow chart

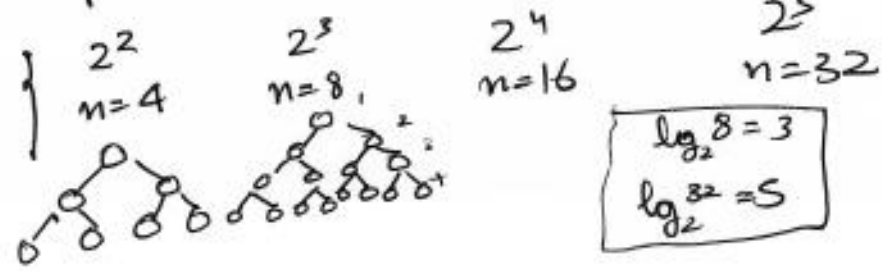
Tree :



Depth of tree = 3

B+ST ÷ Time complexity to search is $O(\lg n)$
↓
depth of binary search tree

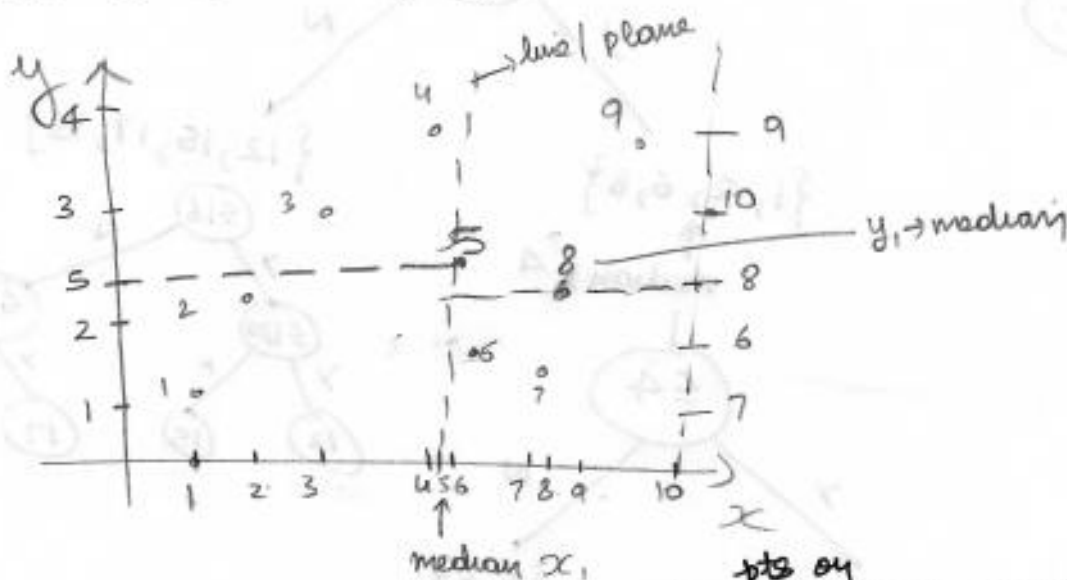
$\begin{cases} n \text{ elements} \\ \text{depth BST} = O(\lg n) \end{cases}$



Kd-tree :- (2D, 3D, ... nD)

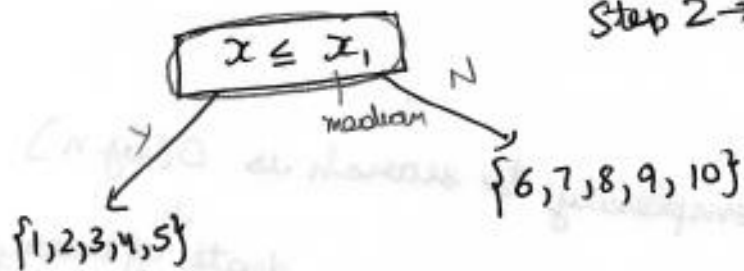
BST :- \rightarrow sort numbers \rightarrow Tree
scalars
'1D'

How to construct a KD-tree

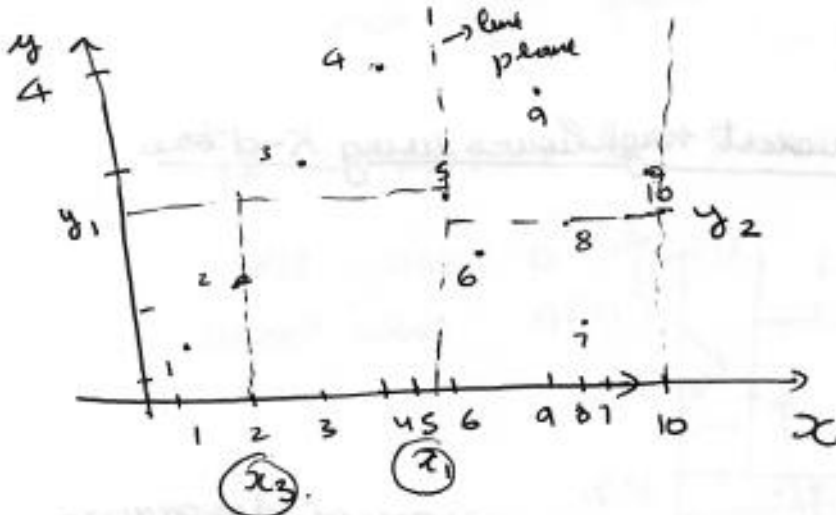
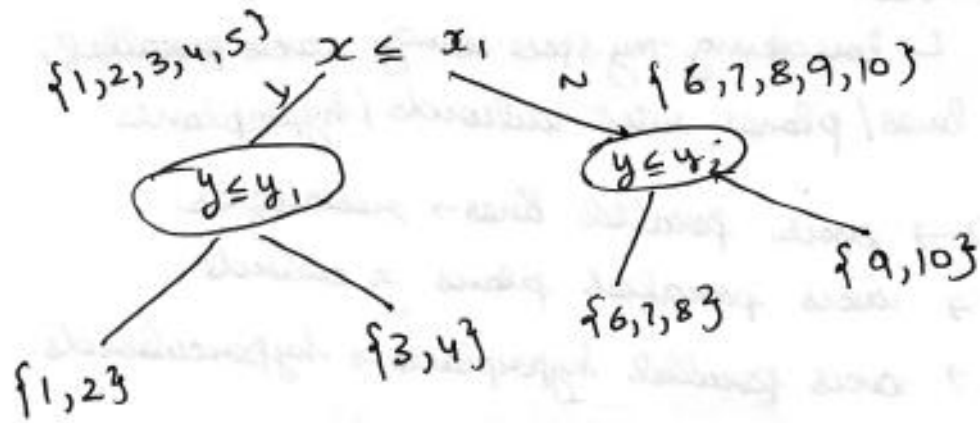


Step 1 pick x-axis. (by projecting ^{pts on} ~~at~~ into x-axis)
Computing the median. $\rightarrow x_1$
(NOTE: half points on left and half on right)

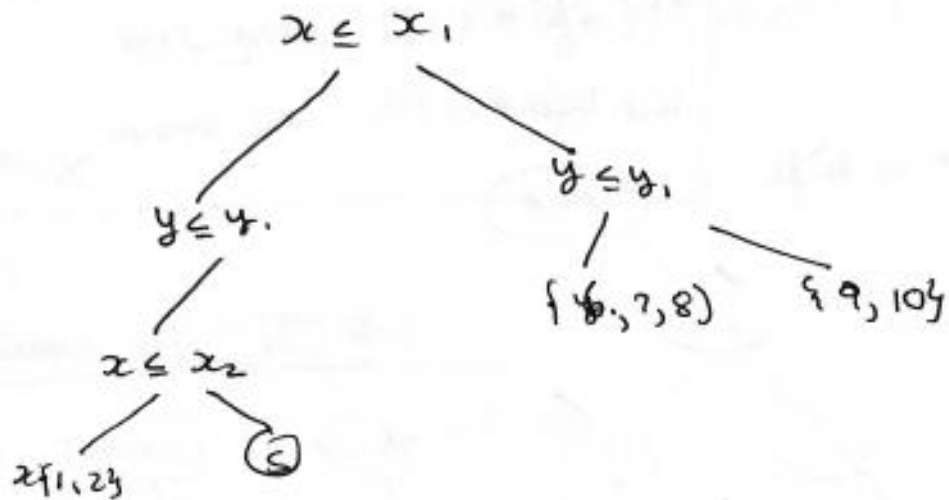
Step 2 \rightarrow split data



Step 3 \rightarrow Use Y axis (project pt on Y axis). Find median of Y axis on left and right side



These line are not parallel.



It is not an ideal kd tree will not work well.

Kd tree

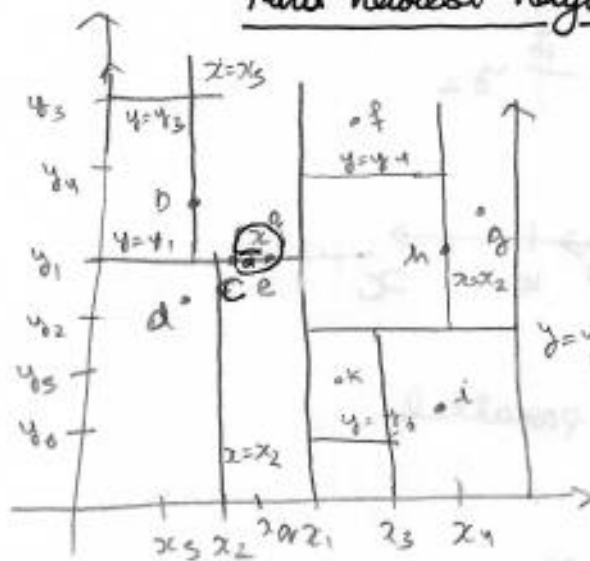
↳ breaking my space using axis parallel lines/planes into culoids/hyperplanes

2D → axis parallel lines → rectangles

3D → axis parallel planes → culoids

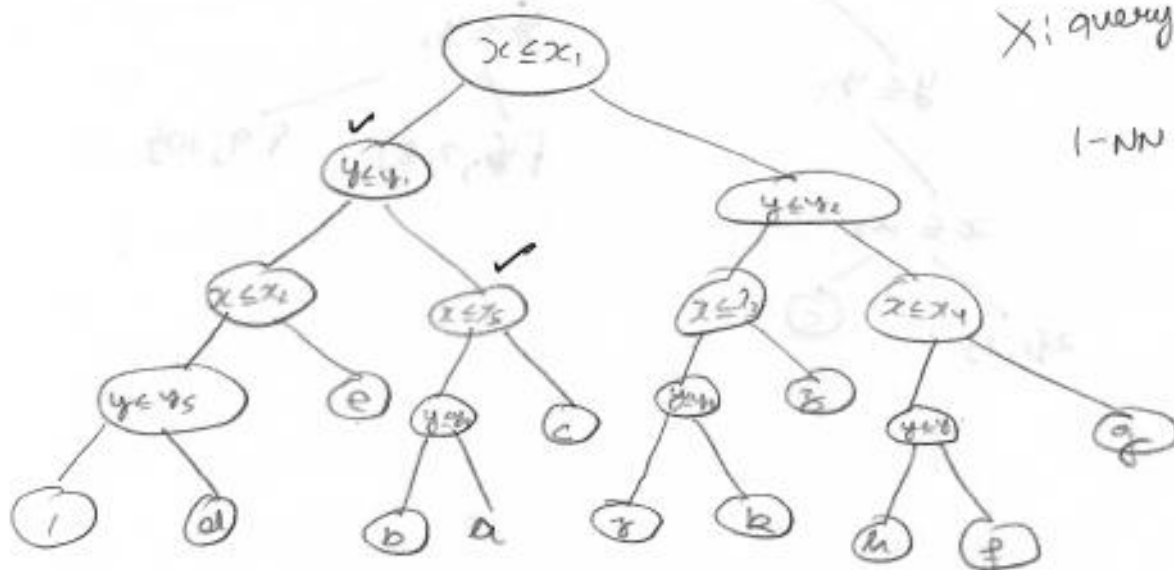
nd → axis parallel hyperplanes → hyperculoids.

Find nearest neighbours using K-d tree



- nearest diagram

Kd tree



x : query point (a)
(x_q, y_q)

1-NN of q

Circle / hyper-sphere : mod-d
C-q

comparisons to find 1-NN

best case : $O(\lg(n))$

worst case : $O(n)$



check in every box.

| | |
|------------|----------------|
| K-NN | |
| best case | $O(K * \lg n)$ |
| worst case | $O(K * n)$ |

if d is small

☆ Limitations of Kd tree

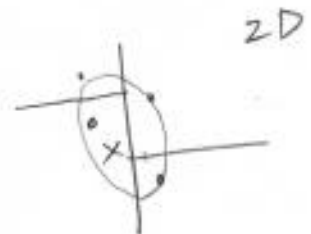
what happens if d is not small.

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

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

d ↑ worst case # adj cells 2^d

$$O(\lg n) \quad \text{if } n=1 \text{ Million} \quad d=20 \quad 2^d \approx n.$$



if d is not small Kd tree will not work well.