

# CLASSIFICATION AND Regression MODELS

## K-Nearest Neighbors

18.1 How "classification" works?

Amazon fire food review

Review  $\rightarrow$  Text  $\rightarrow$  Vector (Bow/tfidf)  $w_{2v}$

Classification:

we have 364K review  $(x_i) \rightarrow (x_i^T, +ve/-ve)$   
Vector representation.

$\rightarrow$  given a new review predict if the review is +ve & -ve

$\rightarrow$  finding a function  $y = f(x)$   
+ve/-ve review text

Classification algo working:

Training  $\rightarrow$   $f$  algo  
 $(x_i, y_i)$   
 $i = 1 \text{ to } 100K$

Testing/Evaluation  $x_q \rightarrow f \rightarrow y_q$

18.2 Data matrix notation:

$\begin{array}{c} \uparrow n \\ \downarrow \\ x_i^T \end{array} \xrightarrow{\quad d \quad} \left[ \begin{array}{c} \xleftarrow{\quad n \quad} \\ \xrightarrow{\quad n \times d \quad} \end{array} \right] \xrightarrow{\quad y_i \quad} \left[ \begin{array}{c} \xleftarrow{\quad n \times 1 \quad} \end{array} \right]$

$x_i \rightarrow \text{Vector } \& x_i^T$

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

### 18.3 Classification vs Regression (Examples)

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

Amazon food reviews  
MOIST:

There are two classes (2-class classification)  
(binary classification)

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

what is  $y$   $y_i \in \mathbb{R}$

$\hookrightarrow y_i$  is no more part of a small finite set of classes

Ex:

$i = 1 \rightarrow 10k$

$x_i: <\text{weight, age, gender, race}>$

$y_i = \text{height}$  (real number)

$y_i = f(x)$

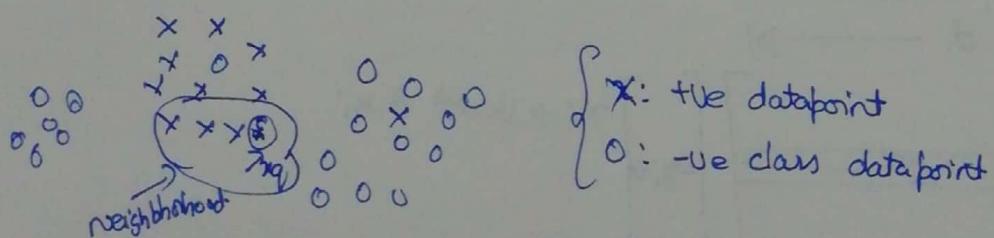
where  $y_i$  is a real number then it is called Regression algorithm.

### 18.4 $k$ nearest neighbours Geometric Intuition with a toy Example

$k$ -Nearest Neighbors (kNN)

2D toy dataset

Let's take a look of binary classification



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

geometric since  $x_q$  is close to  $\textcircled{X}$  we can conclude that  $x_q$  is positive

### Steps

① Find  $K$  nearest points to  $x_q$  is  $\emptyset$

Let  $K=3$  3 nearest points are

$(x_1, x_2, x_3)$   
 $y_1, y_2, y_3$

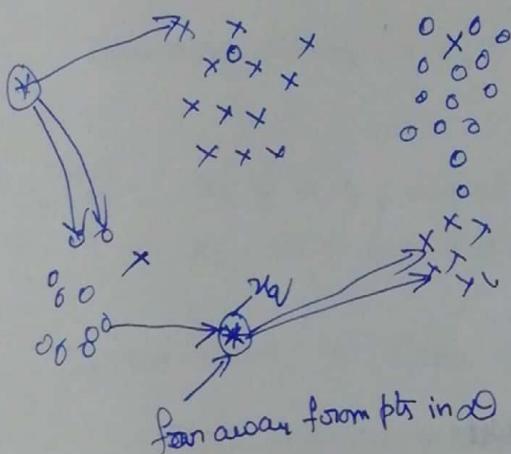
②

② Take all the class labels  $(y_1, y_2, y_3)$

Based on the majority we will declare the class

18.5

### Failure Cases of KNN



$x_0 x_0 x_0 x_0 x_0 x_0$

$x x_0 x_0 x_0 x_0$

$x_0 0 x x_0 x_0 x_0 x_0$

$x x_0 x_0 x_0 x_0 x_0$

↓

→ Jumble & +ve / -ve

→ Randomly spread

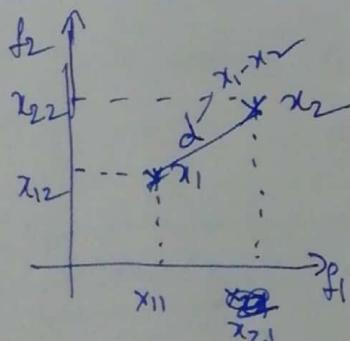
→ There is no useful information

→ As the  $x_q$  is far away, then it's

Very good say don't know

→ There is usefull information

18.6 Distance measure: Euclidean ( $L_2$ ), Manhattan ( $L_1$ ), Minkowski, Hamming.



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

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

$$\rightarrow d = \sqrt{(x_{21} - x_{11})^2 + (x_{22} - x_{12})^2} = \sqrt{\sum_{i=1}^2 (x_{2i} - x_{1i})^2}$$

by Pythagorean theorem. Euclidean distance.

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

$$\|x_1 - x_2\|_2 = \sqrt{\sum_{i=1}^d (x_{1i} - x_{2i})^2}$$

$$\|x_1 - x_2\|_2 = L_2 \text{ norm}$$

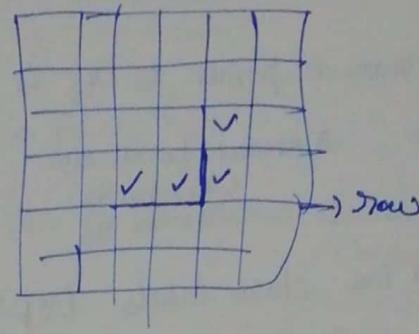
$$\|x_1\|_2 \rightarrow \text{dist of } x_1 \text{ from origin}$$
$$= \sqrt{\sum_{i=1}^d x_{1i}^2}$$

Manhattan dist:

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

$L_1$  norm of vector  $(x_1 - x_2)$

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



$L_1$  norm of  $x_1$  is

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

Generalization of  $L_1$  &  $L_2$  are  $L_p$  norms

$L_p$ -norm  $\rightarrow$  Minkowski dist

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

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

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

$$\|x_1\|_p = \left( \sum_{i=1}^d |x_{1i}|^p \right)^{1/p} \quad p \neq 0$$

$\rightarrow$  Distances are always computed b/w two points  
and norms are always for a vector.

$$\begin{aligned} \text{Eucl dist}(x_1, x_2) &= 22 \text{ norm of } (x_1 - x_2) \\ &= \|x_1 - x_2\|_2 \end{aligned}$$

$$\text{Manhattan}(x_1, x_2) = \|(x_1 - x_2)\|_1$$

$$\text{minkowski}(x_1, x_2) = \|(x_1 - x_2)\|_p$$

## Hamming Distance

→ This is used in text processing when we have a boolean vector

$x_1, x_2 \rightarrow$  boolean vector  $\rightarrow$  Binary Bow

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

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

Hamming distance( $x_1, x_2$ ) = # location/dimensions where binary vectors differ

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

Ex: 2  $x_1 = a \underset{b}{\cancel{c}} \underset{d}{\cancel{a}} \underset{e}{\cancel{d}} \underset{f}{\cancel{e}} \underset{g}{\cancel{f}} \underset{h}{\cancel{g}} \underset{i}{\cancel{h}} \underset{j}{\cancel{i}} \underset{k}{\cancel{k}}$

$$x_2 = a \underset{c}{\cancel{b}} \underset{d}{\cancel{a}} \underset{e}{\cancel{d}} \underset{f}{\cancel{e}} \underset{g}{\cancel{f}} \underset{h}{\cancel{g}} \underset{i}{\cancel{h}} \underset{j}{\cancel{i}} \underset{k}{\cancel{k}}$$

If we have a Gene codes/seq of AGTC characters

$$x_1 = A A G T C T C A G \dots$$

$$x_2 = A G A T C T C C A \dots$$

## 18.7 Cosine distance & Cosine similarity

Similarity vs distance

→  $x_1, x_2$  distance increases, similarity decreases

→ " decreases " increases

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

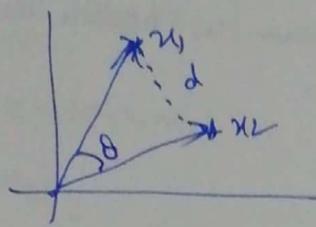
Cosine is b/w [-1, 1]

→ Cosine of  $(x_1, x_2) = +1$  when they are very similar

→ " = -1 " very dissimilar

## Cosine Similarity

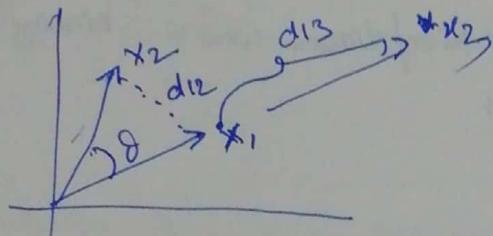
$d = \text{Euc. dist}$



$$\text{Cos-sim} = \cos \theta$$

where  $\theta$  is angle b/w  $x_1, x_2$

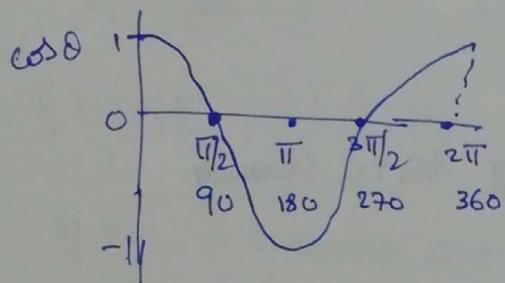
→ Difference b/w Euc. dist & Cos. similarity.



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

$$\text{Cos-sim}(u_1, x_3) = 1$$

$$\theta x_1, x_3 = 0^\circ$$



$$\begin{aligned} \text{Cosdist}(x_1, x_3) &= 1 - \text{Cos-sim}(x_1, x_3) \\ &= 1 - 1 \\ &= 0 \end{aligned}$$

Ex

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

① If  $x_1, x_2$  are unit vectors (length = 1)

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

## Relation b/w Euc. dist vs Cos-sim

If  $x_1, x_2$  are unit vectors

$$\begin{aligned} [\text{eucdist}(x_1, x_2)]^2 &= 2 \left( 1 - \underbrace{\cos(\theta)}_{\text{Cos dist}} \right) \\ &= 2 (\text{Cosdist}(x_1, x_2)) \end{aligned}$$

## 18.8 How to Measure The Effectiveness of knn

Amazon Fine food review

$$D = \frac{1}{n} \sum_{i=1}^n \|x_i - x_q\|^2$$

364K

Problem : Given a new food review ( $x_q$ ) what is its polarity (positive)

$$x_q \rightarrow \text{Text} q \rightarrow (x_q)$$

measure:

$$D_n \xrightarrow{\text{Randomly split}} D_{train} \cup D_{test} = D_n$$

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

70%  
30%

$(x_i)$

Randomly split

$$D_n \xrightarrow{\text{Randomly split}} D_{train} \rightarrow (k-n) \text{ v}$$

$$(x_i, y_i)_{i=1}^{n_1}$$

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

what is training in  $D_{train}$ ?

$$\begin{array}{c} x \ x \ y \ y \\ x \ y \ y \ y \\ x \ y \end{array}$$

$$\begin{array}{c} 0 \ 0 \ 0 \ 0 \\ 0 \ 0 \ 0 \end{array}$$

$$\begin{array}{c} 0 \ 0 \ 0 \\ 0 \ 0 \end{array}$$

$$\begin{array}{c} p \ p \ x \\ x \ x \\ x \end{array}$$

Each point  $x_i$  in  $D_{train} = (x_i, y_i)_{i=1}^{n_1}$

$$x_q = (x_q)$$

for each point in  $D_{test}$   $x_q = p$  use  $D_{train}$  & knn to determine.

$$y \ y_q == y_{pt}$$

Count  $\pm 1$

End

Count = # of pts for which  $\text{d}_{\text{train}} + \text{knw}$  gave a correct class label

$$\text{Accuracy} = \frac{\text{Count}}{n_2} \rightarrow \# \text{ pts. for which } \text{d}_{\text{train}} + \text{knw} \text{ gave a correct class label}$$

$\xrightarrow{\text{pts. in Start}}$

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

Acc = 0.91 = 91% of times it's predicting correctly

### 18.9 Test/Evaluation time & Space Complexity

We are given a query point  $x_q$

$$x_q \rightarrow y_q$$

Input =  $\mathcal{D}_{\text{train}}, k, x_q \in \mathbb{R}^d$ ; output =  $y_q$

$$\text{knwpts} = [ ]$$

for each  $x_i$  in  $\mathcal{D}_{\text{train}}$ :  $\rightarrow n_{\text{pts}}: d\text{-dim}$

- Compute  $d(x_i, x_q) \rightarrow d_i$   $\rightarrow \text{knwpts} [ ]$
- $\rightarrow$  Keep the  $k$  smallest distances  $\rightarrow (x_i, y_i, d_i) \in \text{tuple}$ .  
\*  $k$  is small (5 & 10)

$$\text{Count\_pt} = 0, \text{Count\_neg} = 0$$

for each  $x_i$  in knwpts

$$y \ y_i \text{ is true}$$

$$\text{Count\_pt} += 1$$

else

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

If  $\text{Count\_pt} > \text{Count\_neg}$

$$\text{return } y_q = 1$$

$$\text{else } y_q = 0$$

Time Complexity  $\hat{=} O(nd) + O(1) + O(1)$

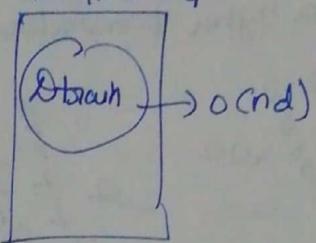
Amazon food reviews

$n \approx 364K$

$d \approx 100K$  (BOW)

300 (tfidf)

Ram/memory

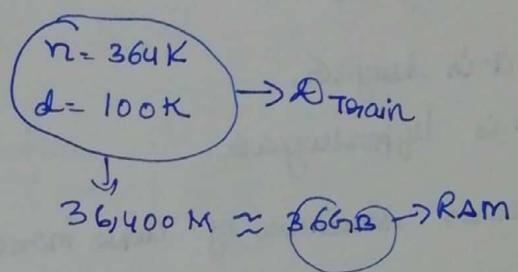


### 18.10 k-nn Limitations

#### Large space Complexity

Time Complexity  $\hat{=} O(nd)$

Space  $\hat{=} O(nd)$



#### Time Complexity

36 Billion Computation.

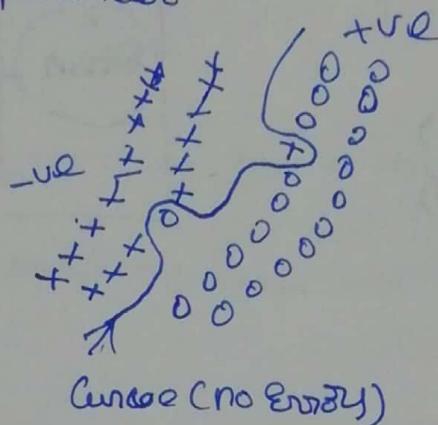
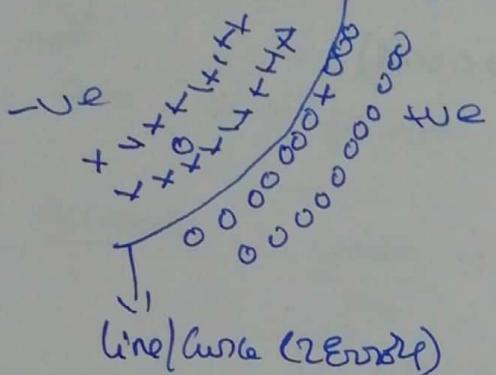
Review  $\xrightarrow{\text{1ms}}$  true/false. (Internet application)

Low-latency  $\hat{=} \text{Time it takes to predict } y_q \text{ given } x_q$  (It should be fast)  
Ex: Trading, Finance

- Simple implementation we have k-nn  $\rightarrow O(nd), O(nd)$
- Ppl don't use k-nn because of Time & Space Complexity

18.11 Decision Surface for k-nn as k changes

→ Here  $k$  is a hyper parameter



→ These Curve +ve form -ve & vice versa are called decision surfaces

→ In 3D it is surface

→ In nD it is hypersurface.

→ As  $k$  increases smoothness of curve increase.

Let's say  $k = n$

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

$n =$  Total number of points

$$n_1 = +ve (600)$$

$$n_2 = -ve (400)$$

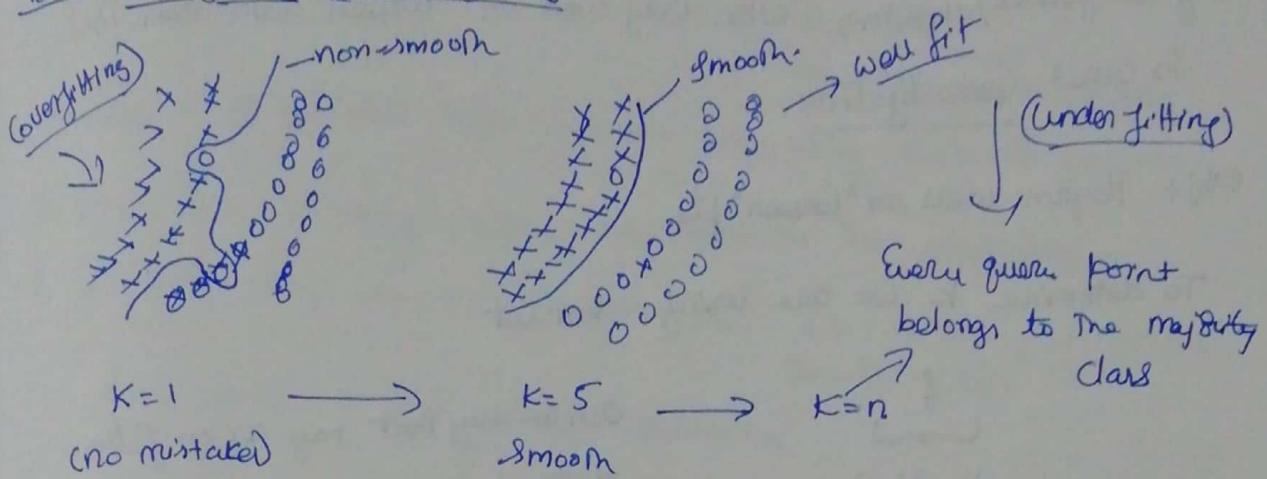
$$n_1 + n_2 = n$$

Let's say  $n_1 > n_2$

If my  $k = 1000$ , 1000 neighbour ~~one~~ should be considered which means we have to consider all the points.

Hence here the majority class is  $n_1$ , our model will predict  
Every input as positive only

## 18.12 Overfitting & Underfitting



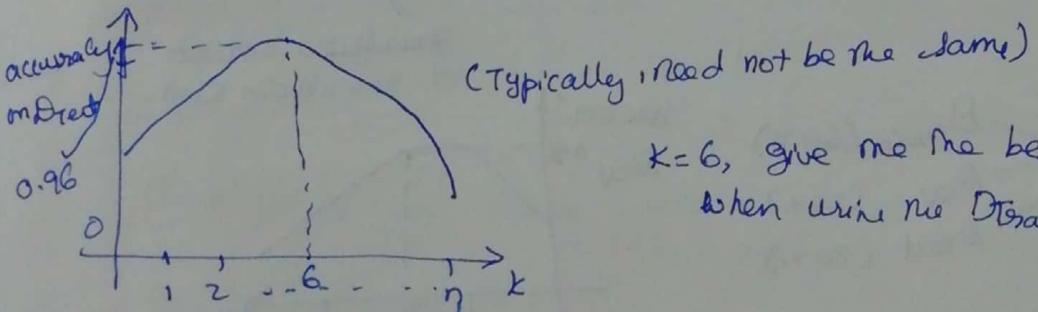
fitting because  $f(x_q) = y_q$ . The process of finding a function is called fitting

## 18.13 Need for Cross Validation

$D_n$   $\begin{cases} D_{\text{Train}} (70\%) \\ D_{\text{Test}} (30\%) \end{cases}$

one-idea:

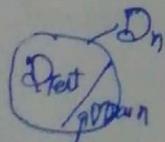
$K=1$	$T_{\text{Train}}$	accuracy on Test
	$D_{\text{Train}}$	0.78
$K=2$	"	0.82
$K=3$	"	0.85



$K=6$ , give me the best accuracy when using the  $D_{\text{Train}}$ .

$$K=6 \rightarrow (6 - \text{num}) \rightarrow 96\%.$$

Using  $D_{\text{Train}}$  & 6-fold on Amazon review dataset. It got an acc of 96%.

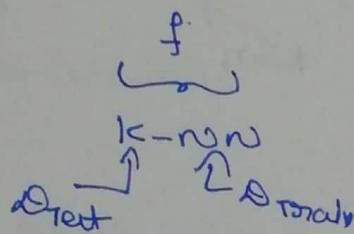


Obj + For a feature unseen point  $x_q$ , how accurate I can predict  $y_q$

→ If a function/algorithm work very well on unseen data then they is called Generalization

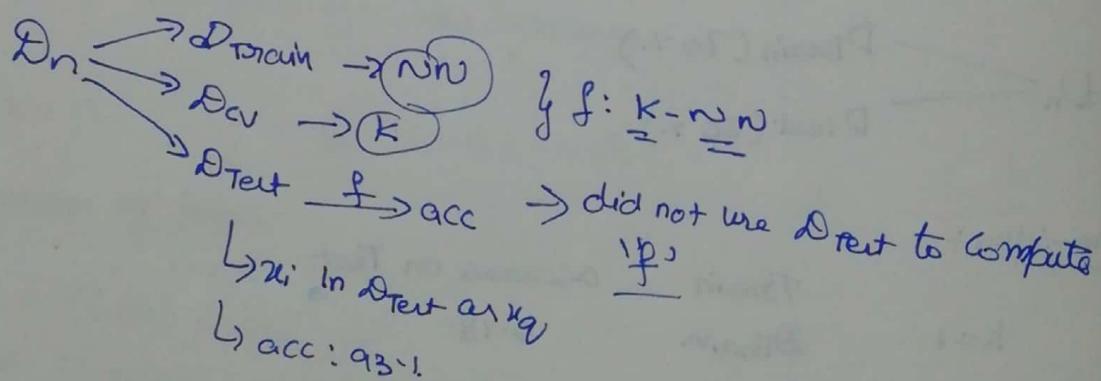
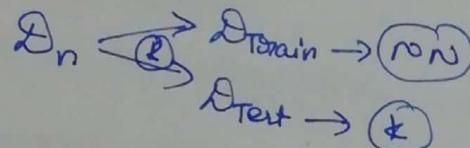
Obj: Perform well on "unseen" pts

To determine  $k$  we are using Test set



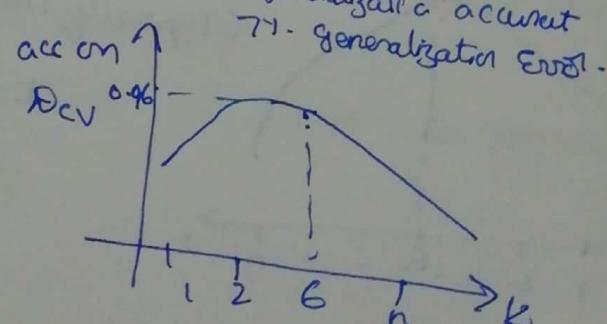
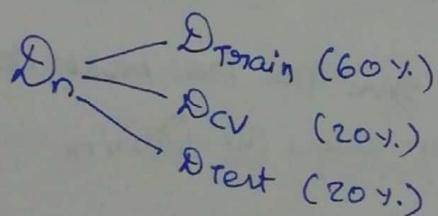
Can we say that my acc on future data is 96%.  
Ans: we cannot guarantee

To overcome this problem we use Cross Validation.



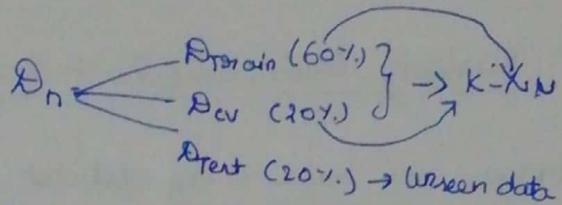
→ Can we say that 6-nn has an acc of 93% on unseen data?

\* Yes



19.14

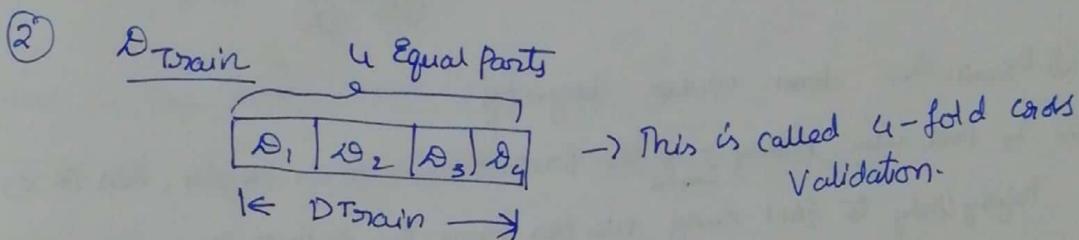
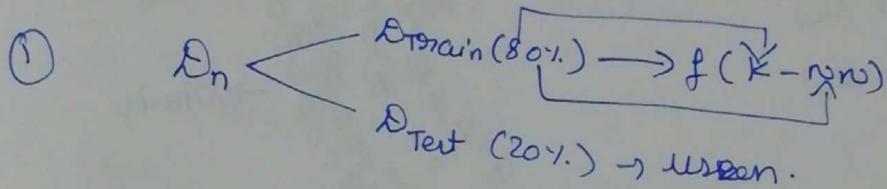
## K-Fold Cross validation



Problem : of the total data we are only using 60% to find error

→ More the Training data, the better is the algorithm.

→ Is there any way to combine the 60% & 20% , it can be  
 (Training) (CV)  
 done using K-fold Cross Validation.



	Train	CV	Acc. on CV	
4 times	$D_1   D_2   D_3$	$D_4$	$a_4$	$\left. \begin{array}{c} \text{avg}(a_1, a_2, a_3, a_4) \\ = a_{K=1} \end{array} \right\}$
	$D_1   D_2   D_4$	$D_3$	$a_3$	
	$D_2   D_3   D_4$	$D_1$	$a_1$	
	$D_2   D_3   D_4$	$D_1$	$a_1$	
4 times	$K=2$			
	$\vdots$			
4 times	$K=3$			

③  $K^1$  fold CV →  $D_{Train} \xrightarrow{Train} K_{in \text{ known}}$

→ what is the right  $K^1$

$$K^1 = 4; \quad K^1 = 10; \quad K^1 = 100$$

$\downarrow$

rule of thumb : 10 fold CV

We are repeating multiple times, to compute the best  $K_{in \text{ known}}$   
 Increases  $K^1$  time when we are doing  $K^1$  CV

18.15

## Visualizing Train, Validation & Test datasets

Imagine we have a big dataset  $D_n$

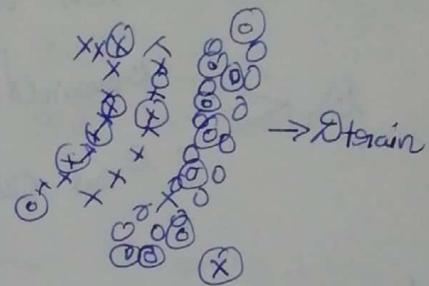
$D_{Train} (60\%)$   
 $D_{CV} (20\%)$   
 $D_{Test} (20\%)$

$\times$  : -ve class datapoint in  $D_{Train}$   
 $\circ$  : +ve class datapoint in  $D_{Train}$

$D_{Train}, CV, Test$  :  $(\vec{x}_i, y_i)$

Input data point  
 class label

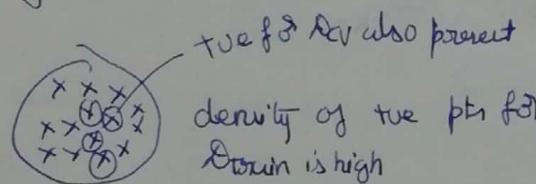
- ①  $\times$  -ve point  $D_{CV}$
- ②  $\circ$  +ve point  $D_{CV}$



### Tips

- ①  $D_{Train}, D_{CV}$  do not overlap perfectly
- ② If there are many +ve pts from  $D_{Test}$  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  $D_{Test}$  then it is very unlikely to find +ve/-ve from  $D_{CV}$  in that region.

### Intuitively



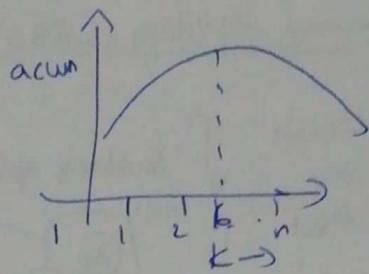
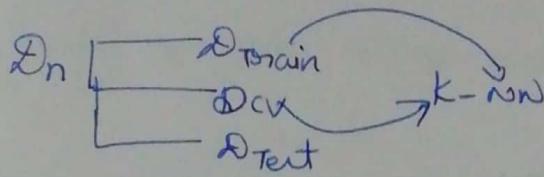
## 18.16 How to determine Underfitting & Overfitting

k-fold on  $D_{CV} \rightarrow$  best k  $\rightarrow$  neither overfit  
 underfit

- (Q) How can we be sure about O.F & U.F

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

$$\text{Error} = 1 - \text{accuracy} = 0.07$$



Training Error

what is the train error for 2-nn?

for each  $x_i$  in  $D_{train}$

- find 2 nearest neighbor to  $x_i$  from  $D_{train}$
- majority vote to get the class label
- if  $y_i$  == class label  
accurate

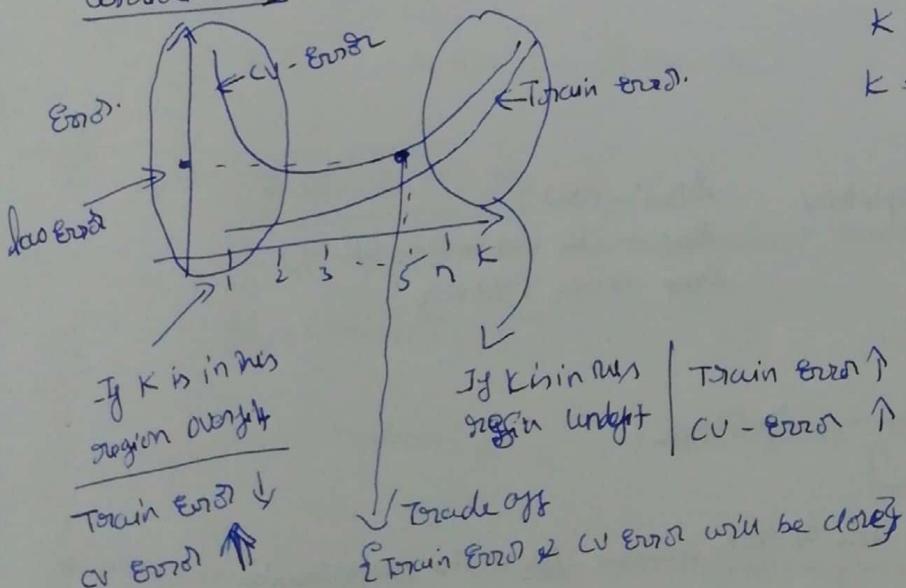
else

error

what is the  $\frac{D_{cv}}{D_{cv}}$  error for 2-nn?

→ same as above, calculate the error & accuracy

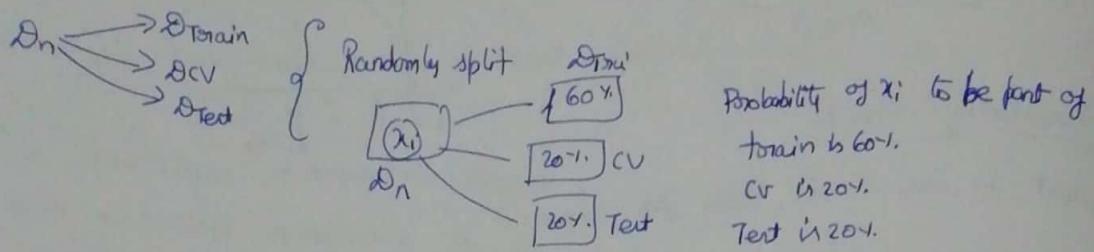
Consider  $k=3$



$k=1$  overfit

$k=n$  underfit

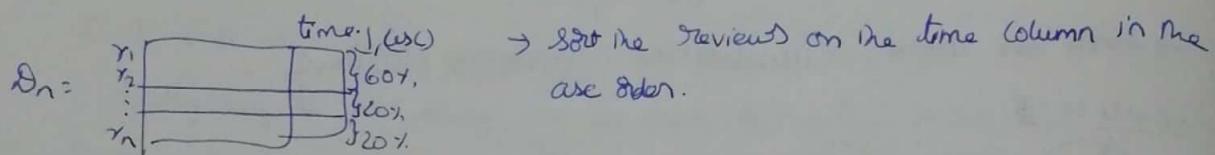
## 18.17 Time based splitting (TBS)



For Amazon food reviews, time based splitting is better than Random-splitting

$$D = \begin{bmatrix} f_1 & f_2 & \dots & f_d \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

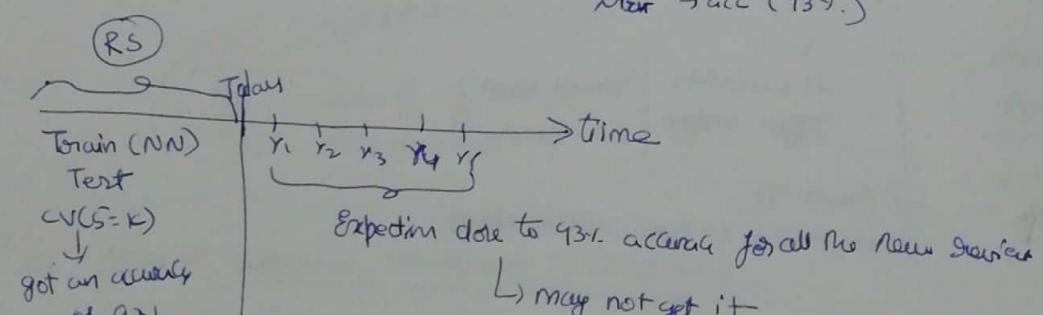
In the actual dataset



### Why TBS?

Let's say I used random splitting

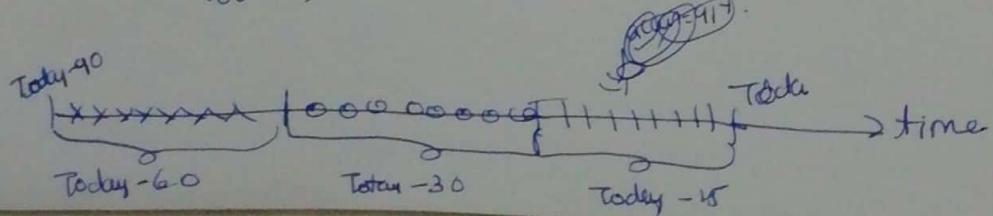
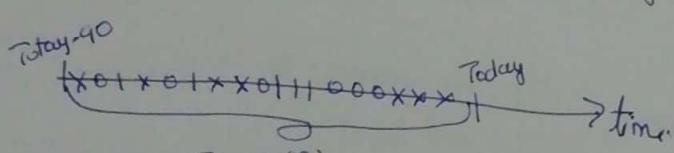
$D_{train} - N(1)$   
 $D_{cv} - K \text{ in known } (K=5)$   
 $D_{test} \rightarrow \text{acc } (93\%)$



Deployed-like ( $D_{train}$ ) ( $K=5$ )

\* the new reviews may be quite different from Train & CV

\* with time my products & review change



If I train my model 15 days back and I tested on last 15 days accuracy is 91%.

We can predict that when we predict the next 15 days data we can expect 91% of accuracy.

When even time is available y things/behaviour so data change over time then time based splitting is preferable than random splitting.

### 18.18 K-nn for regression

Classification  $D = \{(x_i, y_i)_{i=1}^n \mid x_i \in \mathbb{R}^d, y_i \in \{0, 1\}\}$   
 $x_q \rightarrow y_q \rightarrow \text{class label}$

Regression  $D = \{(x_i, y_i)_{i=1}^n \mid x_i \in \mathbb{R}^d, y_i \in \mathbb{R}\}$   
 $x_q \rightarrow y_q \rightarrow \text{number}$

① Given  $x_q$ , find K-nearest neighbors

$(x_1, y_1), (x_2, y_2), \dots, (x_K, y_K)$

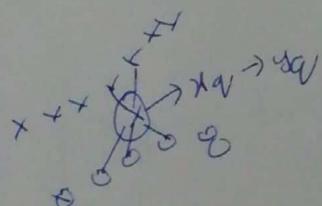
②  $y_q \leftarrow \underbrace{y_1, y_2, y_3, \dots, y_K}_{\text{all are real valued numbers.}}$

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

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

\* For some algorithm we can ~~not~~ extend/modify it to regression.

### 18.19 Weighted K-nn



5 nn  
 $K=5$

2  $\rightarrow$  we  $\{ \rightarrow y_q$   
3  $\rightarrow$  we  $\{ \rightarrow y_q$

$x_q \rightarrow x_1, y_1, d_1 \rightarrow 0.1 \text{ we}$   
 $x_2, y_2, d_2 \rightarrow 0.2 \text{ we}$   
 $x_3, y_3, d_3 \rightarrow 1.0 \text{ we}$   
 $x_4, y_4, d_4 \rightarrow 2.0 \text{ we}$   
 $x_5, y_5, d_5 \rightarrow 4.0 \text{ we}$

$x_i$	dist $\alpha$	$w_i$	
$x_1$	$\rightarrow x_1, -ve -0.1$	10	
$x_2$	$\rightarrow x_2, -ve -0.2$	5	
$x_3$	$\rightarrow x_3, -ve -1.0$	1	
$x_4$	$\rightarrow x_4, -ve -1.0$	0.5	
$x_5$	$\rightarrow x_5, -ve \rightarrow 4.0$	0.26	
		15	
		1.75	

$$w_i = 1/d_i$$

One way to find weights

$$\rightarrow w_i = \frac{1}{d_i} \quad \text{Simplest weight function}$$

### 18.20 Voronoi diagram

→ Related to k-nn ( $k=1$ )

→ It rays

### 18.21 Binary Search Tree

$k$ -nn +  $O(n)$  If  $d$  is small  $\rightarrow$  Time Complexity  
 $k$  is small

$\Theta(n)$  → Space Complexity

Time  $\vdash O(n) \rightarrow O(\lg(n)) \rightarrow$  Kd-trees

$$n = 1024 \rightarrow \lg(n) = 10$$

Applied in

Computer graphics, geometry, Applied math,

### Binary search tree (BST)

→ Data structure

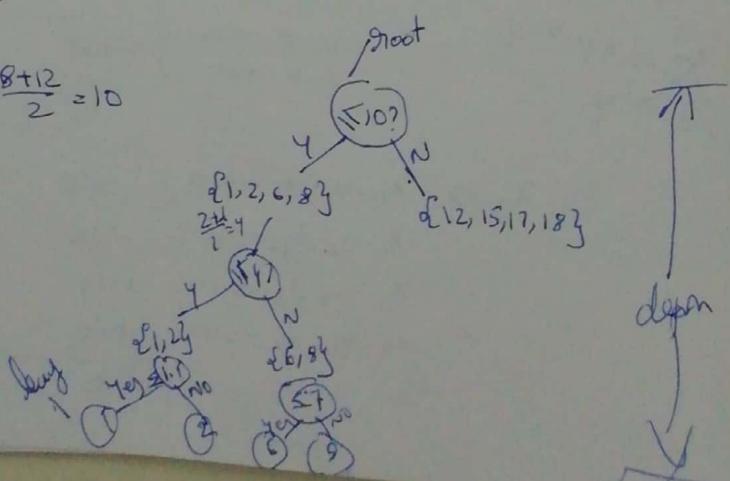
→ Binary Search

Prob: given a sorted array, find a number is present in array or not

$$A: [1, 2, 3, 4, 5, 6, 7, 8]$$

① Construct a BST

$$\text{Take median: } \frac{8+12}{2} = 10$$

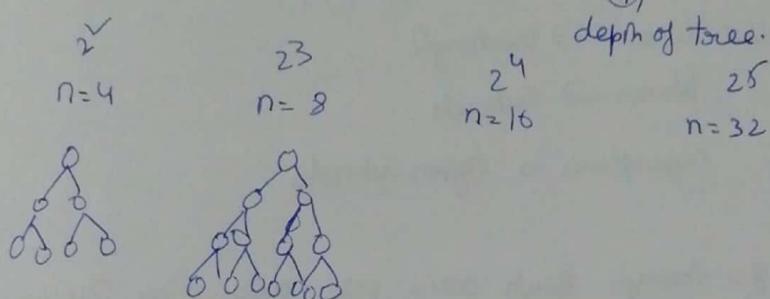


② find an element in an array

$$q = 15$$

To find 15, with comparison we can find 15 exists or not

BST: Time Complex to search is  $O(\log(n))$



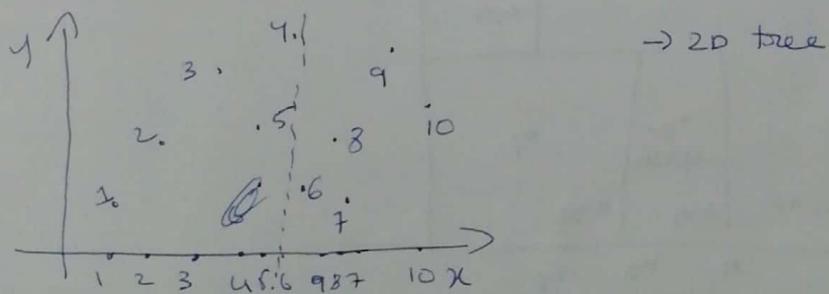
### Q2 How to build a kd-tree

BST: given a list of sorted numbers  $\rightarrow$  Given a tree.

$\underbrace{\quad\quad\quad}_{\text{scalar}}$   
1D

Extend BST to kd-tree

$\hookrightarrow$  It will operate in 2D, 3D ... n

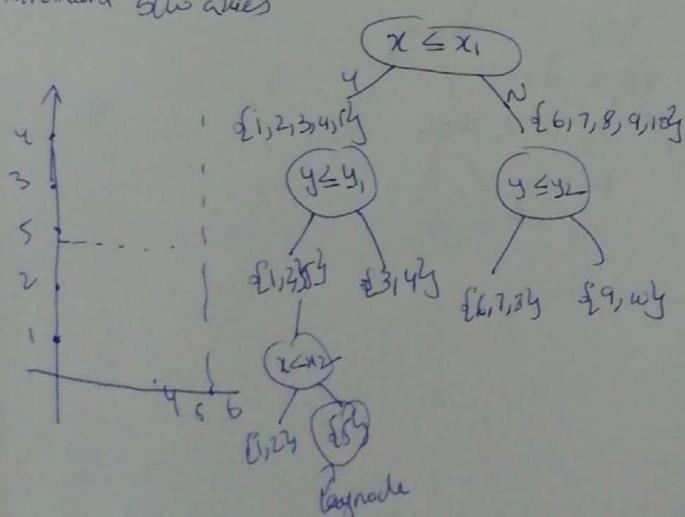


① Pick x-axis, project points on to x-axis.

② Compute the median (will be b/w 5 & 6)

$\hookrightarrow$  Half the points to the left of median. and the half to the right

③ Alternate axes



## 18.30 Code Sample : Decision Boundary

sklearn.neighbors.KNeighborsClassifier

class sklearn.neighbors.KNeighborsClassifier (n\_neighbors=5, weights='uniform',  
algorithm='auto', leaf\_size=30, p=2, metric='minkowski', metric\_params=None  
n\_jobs=1, \*\*kwargs)

n\_neighbors: int optional (default=5)

weights: str or callable, optional (default='uniform')

- 'uniform': uniform weights. All points in the neighborhood are weighted equally
- 'distance': weight points by the inverse of their distance. In this case, closer neighbors of a query point will have a greater influence than neighbors which are further away
- [callable]: a user-defined function which accepts an array of distances, and returns an array of the same shape containing the weights

algorithm: {'auto', 'ball-tree', 'kd-tree', 'brute'}, optional

Algorithm used to compute the nearest neighbors:

- \* 'ball-tree' will use BallTree
- \* 'kd-tree' will use KDTree
- \* 'brute' will use the brute-force search
- \* 'auto' will attempt to decide the most appropriate algorithm based on the values passed to fit method

Note: fitting on sparse input will override the setting of this parameter, use brute force

- \* Dimensionality is low it will pick KDTree
- \* If dataset is small it will pick Brute

leaf\_size: int, optional (default=30)

leaf\_size passed to BallTree & KDTree. This can affect the speed of the construction and query, as well as the memory required to store the tree.

The optimal value depends on the nature of the problem

p: integer, optional (default=2)

Power parameter for the Minkowski metric. When  $p=1$ , this is equivalent to using manhattan\_distance (L1) & euclidean\_distance (L2) for  $p=2$ . For arbitrary  $p$ , Minkowski\_distance ( $l_p$ ) is used

metric: string or callable, default 'minkowski'

The distance metric to use for the tree. The default metric is minkowski and with  $p=2$  is equivalent to the standard Euclidean metric.

See the documentation of the distance-metric class for a list of available metrics

metric\_params: dict, optional (default=None)

Additional keyword arguments for the metric function

n\_jobs: int, optional (default=1)

The number of parallel jobs to run for neighbour search.

If -1, then the number of jobs is set to no. of CPU cores.

Doesn't affect fit method.

## Example In sklearn documentation

```
x = [[0], [1], [2], [3]]
```

```
y = [0, 0, 1, 1]
```

from sklearn.neighbors import KNeighborsClassifier

```
neigh = KNeighborsClassifier(n_neighbors=3)
```

```
neigh.fit(x, y)
```

```
print(neigh.predict([[1.1]]))
```

O/p : [0]

```
print(neigh.predict_proba([[0.9]]))
```

```
[[0.666667, 0.333333]]
```

## 18.3) Code Sample: Cross Validation

```
import sklearn. cross cross-validation import cross_val_score.
```

```
names = ['x', 'y', 'class'] # define class names.
```

```
df = pd.read_csv('Concentriccircles.csv', header=None, names=names)
```

```
df.head()
```

```
# Create design matrix X and target vector y
```

```
X = np.array(df.iloc[:, 0:4]) # End index is exclusive
```

```
y = np.array(df['class']) # showing you two ways of indexing a pandas df
```

### Simple Cross Validation

```
# Split the dataset into train & test.
```

```
X_train, X_test, y_train, y_test = cross-validation.train_test_split
```

```
(X, y, test_size=0.3, random_state=42)
```

# Split the train data set into Cross Validation train & Cross Validation test

$x_{-tr}, x_{-cv}, y_{-tr}, y_{-cv} = \text{cross_validation. train\_test\_split}(x_{-1}, y_{-1}, \text{test\_size: } 0.2)$

for i in range (1, 30, 2):

knn = KNeighborsClassifier (n\_neighbors = i)

knn.fit (x\_tr, y\_tr)

pred = knn.predict (x\_cv) # prediction of cv train

# Evaluate cv accuracy

acc = accuracy\_score (y\_cv, pred, normalize = True) \* float (100)

print ('cv accuracy for k = %d is %.2f' % (i, acc))

knn = KNeighborsClassifier (1)

knn.fit (x\_tr, y\_tr)

pred = knn.predict (x\_test)

acc = accuracy\_score (y\_test, pred, normalize = True) \* float (100)

print ('Test accuracy')

O/P

K=1 100 %

3 100 %

5 100 %

# Creating odd list of k for knn.

myList: list (range (0, 50))

neighbors: list (filter (lambda x: x % 2 != 0, myList))

# Empty list that will hold cv scores.  
cv\_scores = []

for k in neighbors:

knn = KNeighborsClassifier (n\_neighbors = k)

scores = cross\_val\_score (knn, X\_train, y\_train, cv=10, scoring = 'accuracy')

cv\_scores.append (scores.mean())

# changing the misclassification error.

MSE = [1 -> for x in cv\_scores]

# determining best k

optimal\_k = neighbors (MSE. index (min (MSE)))

print ('The optimal number of neighbor is x.d.', x. optimal\_k)

plt. plot (neighbors, mse)

for xy in zip (neighbors, np.round (mse, 3)):

plt. annotate ('(x.s, y.s)', xy, xytext = xy, textcoords = 'data')

plt. xlabel ('number of neighbor k')

plt. ylabel ('misclassification error')

plt. show()

## KFold

class sklearn.model\_selection.KFold (n\_splits, shuffle=False, random\_state=None)

→ K-Folds Cross Validation

→ Provides train/test indices to split data in train/test sets. Split dataset into K consecutive folds (without shuffling by default)

→ Each fold is then used once as a validation while  $K-1$  remaining folds form the training set.

n\_splits : int, default=3

Number of folds. Must be at least 2.

shuffle : boolean, optional

Whether to shuffle data before splitting into batches.

random\_state : int optional.

from sklearn.model\_selection import KFold

X = np.array ([[1, 2], [3, 4], [1, 2], [3, 4]])

y = np.array ([1, 2, 3, 4])

Kf = KFold (n\_splits=2)

Kf.get\_n\_splits (X)

//2

print (Kf)

// KFold (n\_splits=2, random\_state=None, shuffle=False)

for train\_index, test\_index in kf.split (X):

X\_train, X\_test = X [train\_index], X [test\_index]

y\_train, y\_test = y [..], y [..]

## Stratified K Fold

Stratified K Fold (n\_splits=3, shuffle=False, random\_state=None)

- Stratified K-Folds cross Validation.
- provides train/test indices to split data in train/test sets
- This cross validation object is a variation of KFold that generates stratified folds.
- The folds are made by preserving the percentage of samples for each class.

### Notes

All the folds have size  $\text{train}(\text{n-samples})/\text{n-splits}$ , the last one has the complementary.

from sklearn.model\_selection import StratifiedKFold

$X = \text{np.array}([ [1, 2], [3, 4], [1, 2], [3, 4], [5, 6], [7, 8] ])$

$y = \text{np.array}([ 0, 0, 1, 1, ])$

$\text{skf} = \text{StratifiedKFold}(\text{n_splits}=2)$

$\text{skf.get_n_splits}(X, y)$

for train-index, test-index in skf.split(X, y):

$X_{\text{train}}, X_{\text{test}} = X[\text{train-index}], X[\text{test-index}]$

$y_{\text{train}}, y_{\text{test}} = y[ \dots ], y[ \dots ]$

## Repeated KFold

Repeated KFold (n\_splits=5, n\_repeats=10, random\_state=None)

- Repeats K-Fold n times with different randomization in each repetition
- $\text{rkf}$ : Repeated KFold (n\_splits=2, n\_repeats=2, random\_state=2652124)  
for train\_index, test\_index in  $\text{rkf}.\text{split}(X)$ :  
 $X_{\text{train}}, X_{\text{test}} = X[\text{train_index}], X[\text{test_index}]$   
 $y_{\text{train}}, y_{\text{test}} = y[\text{train_index}], y[\text{test_index}]$

## Leave One Out Cross Validation (LOOCV)

- In this approach, we reserve only one data point from the available dataset and train the model on the rest of the data.
- This process iterates for each data point.
- This also has its own advantages & disadvantages.
  - \* we make use of all data points, hence bias will be low
  - \* We repeat the Grid Validation process  $n$  times (where  $n$  is the number of data points) which results in high execution time.
  - \* This approach leads to higher variation in testing model effectiveness because we test against one data point. So, our estimations gets highly influenced by the data point. If the data point turns out to be an outlier it can lead to a higher variation.

==

```
from sklearn.model_selection import LeaveOneOut
```

```
X = np.array ([[1,2], [3,4]])
```

```
Y = np.array ([1,2])
```

```
loo = LeaveOneOut()
```

```
loo.get_n_splits(X)
```

```
for train_index, test_index in loo.split(X):
```

```
X_train, X_test = X[train_index], X[test_index]
```

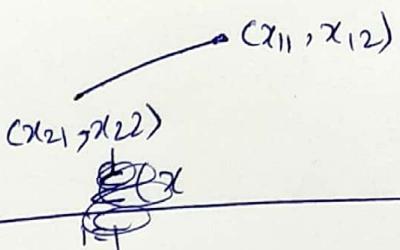
```
Y_train, Y_test = Y [ : , ], Y [ : , ]
```

$$\text{Cosdistance}(x_1, x_2) = 1 - \text{CosSim}(x_1, x_2)$$

$$\text{CosSim}(x_1, x_2) = \cos(\theta_{x_1, x_2})$$

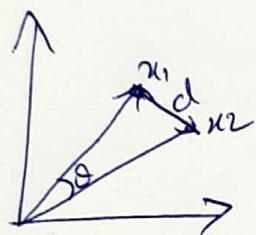
Euclidean distance

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

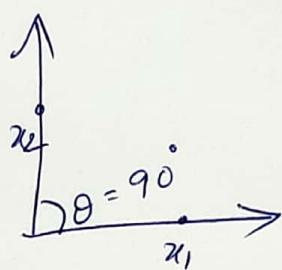


Example

①



②

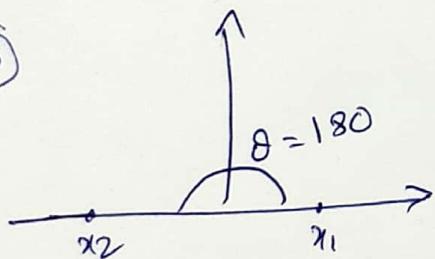


$$\cos(90^\circ) = 0$$

$$\text{CosSim}(x_1, x_2) = 0$$

$$\begin{aligned} \text{Cosdist}(x_1, x_2) &= 1 - 0 \\ &= 1 \end{aligned}$$

③

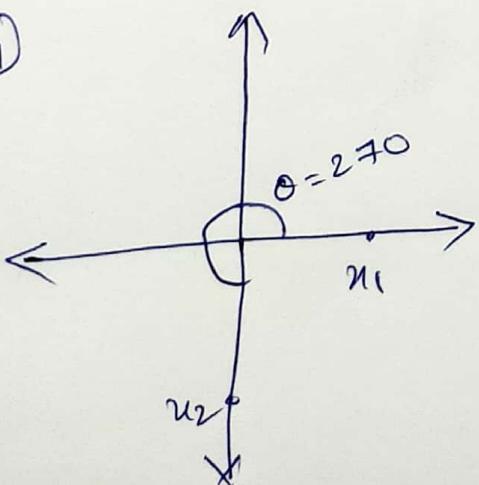


$$\cos(180^\circ) = -1$$

$$\text{CosSim}(x_1, x_2) = -1$$

$$\text{Cosdist}(x_1, x_2) = 1 - (-1) = 1 + 1 = 2$$

④

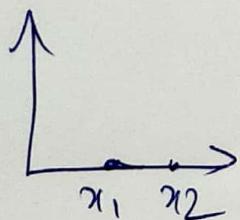


$$\cos(270^\circ) = 0$$

$$\text{CosSim}(x_1, x_2) = 0$$

$$\text{Cosdistance}(x_1, x_2) = 1 - 0 = 1$$

⑤



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

$$\text{CosSim}(x_1, x_2) = 1$$

$$\text{Cosdist}(x_1, x_2) = 1 - 1 = 0$$

Manhattan distance

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

