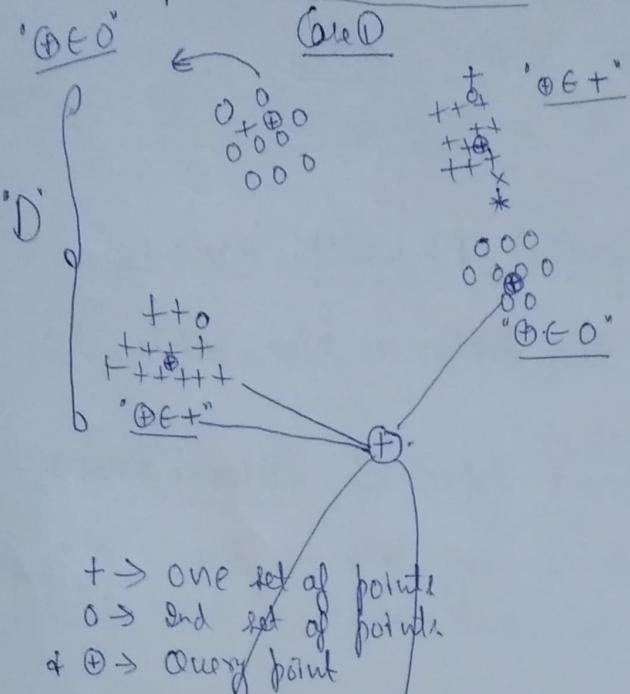


"Failure Cases of K-NN"

Let's now understand where KNN fails \Rightarrow Let's understand it with the help of examples.

Let's take two cases & let's assume $K=3$



Now let's assume our query point is away from the rest of the points.

This point (Q-point) is far away from the points in our dataset.

If we take 3 nearest neighbors for this point, we might conclude that this is a point belonging to "+ class".

Since this point is really far away from the rest of the points, does it really make sense to call the query point as either "+" or "0"?

It looks little absurd or untrue.

bcz we don't have any point really close to our query point \oplus . All these points are very far away from our query point.

Say, when we have our query point lying very far away from our dataset, we are not very sure of its class label.

This is one case where 'K-NN' may not work as per our expectations.

Case ② \Rightarrow "Jumbled data"

+ 0 + 0 + 0 + 0 +
0 + 0 + 0 + 0 + 0
+ + 0 + 0 + 0 +
+ 0 + 0 + 0 + 0 +
0 + 0 + 0 + 0 + 0

$\xrightarrow{\text{randomly spread in } d\text{-d space}}$

↓ Jumble of '+ve points' and '-ve points'

If we take a [query point] & apply 'K-NN' on it, we will get a class label. Not 100% sure, but most likely that it going to be a wrong prediction, bcz data is slightly jumbled up or intermixed.

So, if our data is completely random (randomly distributed), there is no useful information at all in it.

KNN miserably fails here. Actually most of the algorithms miserably fail here.

These are the two simplest cases where "K-NN" fails.

→ "Distance Measure"

Euclidean distance (L_2), Manhattan distance (L_1), "Minkowski", "Hamming"

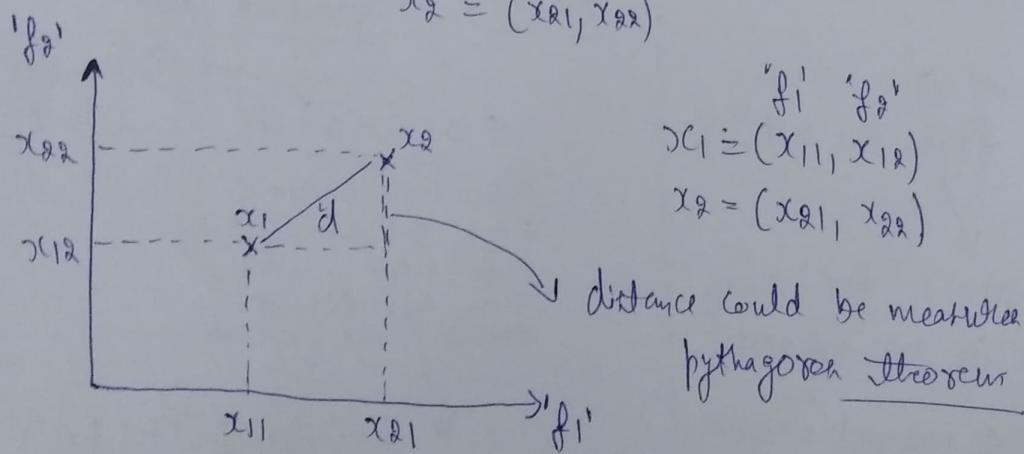
In earlier topic we have measured the distance b/w points to get the nearest neighbours.

Now let's look at some different types of distance measure.

Let's have two points ' x_1 ' & ' x_2 ' & let's say both of them are two-dimensional

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

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



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

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

distance could be measured very simply
pythagorean theorem

Simplest understanding of distance is basically the shortest (length of shortest) distance

$d = \text{length of the shortest line from } x_1 \text{ to } x_2$

That's what we typically say when we call distance.

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

length

"Euclidean distance"

↳ gives the distance of shortest distance

(3)

Now let's assume that $x_1 \in \mathbb{R}^d$, $x_2 \in \mathbb{R}^d$, $x_1 + x_2$ belongs to "d-dimensional" space.

Then distance b/w $x_1 + x_2$ is the length of the vector

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

↓ Euclidean distance is it is represented as $\|x_1 - x_2\|_2$,

in subscript we put "2".

This is also called as "L2 norm of a vector"

So, if we have a vector let's say ' x_i ' then L₂ norm of the vector

is $\|x_i\|_2$ = Euclidean distance of x_i from origin,

below ' x_i ' is a point & a point can be called a vector,

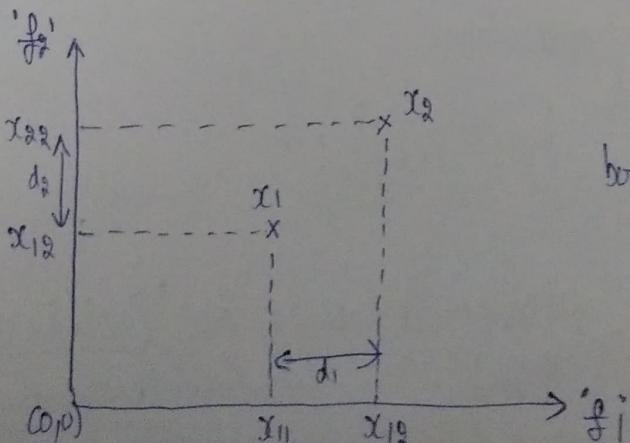
$\therefore \|x_i\|_2$ = Euclidean distance of x_i from origin = $\left(\sum_{i=1}^d x_{ii}^2 \right)^{\frac{1}{2}}$

(2) Distance Metric (Manhattan distance)

Manhattan distance b/w two points " x_1 & x_2 " in d-dimension

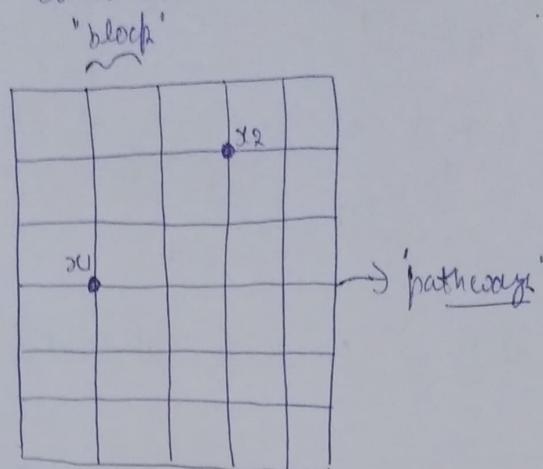
$$\left\{ \sum_{i=1}^d |x_{1i} - x_{2i}| \right\} \xrightarrow{\text{"Absolute value"}}$$

Geometrically, it means what is shown below \Rightarrow



both these $d_1 + d_2$ have to be added to get the manhattan distance b/w ' x_1 & x_2 '.

Sum of these two distances is considered as Manhattan distance.
"block"



Now if we have to go from " x_1 " to " x_2 "
 we have to cross 2 blocks in "horizontal direction"
 & 2 blocks in the vertical direction to reach " x_2 " from x_1

We can't go diagonally because there are buildings there.

Manhattan distance is, in each direction we take the absolute value.

$$\|x_1 - x_2\|_1 \rightarrow \text{"L1 norm" (also called)}$$

Now there is a generalization of ' L_1 -norm' & ' L_2 -norm' called ' L_p Norm' (p could be any no. here)

Now for Lp norm the corresponding distance is called [ℓ^p principal dist.]

Let's first write the formula, the "minnowphi dstance" between any two points is written as

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

If $P=2$ \rightarrow Min power distance is same as "Euclidean-distance"

If $P=1$ \rightarrow $\parallel \quad \parallel$ Manhattan-distance

(5)

$$\|x_1\|_p = \left(\sum_{i=1}^d |x_{1i}|^p \right)^{1/p} \rightarrow L_p \text{ norm of a vector } x_1$$

for L_p norm ' p cannot be zero' $\|x_1\|_0$ is not defined.

Notes: "Distances" are always b/w two points & "norm" are always b/w two vectors. (as a vector)

④ "Hamming distance" \Rightarrow It is mostly used in "text-processing".
Or when we have a binary vector or a boolean vector.

Imagine we have two points $[x_1]$ & $[x_2]$ & both of these points are Boolean vectors.

$$x_1, x_2 \rightarrow \text{"boolean vector"} \rightarrow \begin{matrix} \text{Used for converting text} \\ \text{into vector} \end{matrix}$$

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

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

Hamming distance b/w two vectors (x_1 & x_2) is the number of locations/dimensions where the binary vector differ.

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

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

Hamming distance b/w $[x_1]$ & $[x_2]$ is 3

\Downarrow being in three locations they are different

Hamming distance is also used when we have strings.

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

$$x_2 = acbadegfhik \quad \left[\begin{array}{l} x_1 = AAGTCTACG \\ x_2 = AGATCTCGA \end{array} \right]$$

In bioinformatics also we use Hamming distance.

→ "Cosine similarity & Cosine-distance" :-

There is one more distance measure called ["Cosine-distance"] & the corresponding similarity is ["Cosine-similarity"].

First let's understand the relationship between similarity & distance.

"similarity"

"distance"

" x_1 & x_2 "

As the distance b/w two points ' x_1 ' & ' x_2 ' increases, this means they are less similar.

$\downarrow (x_1, x_2)$

$\uparrow (x_1, x_2)$

+ As the similarity b/w two points increases, the distance between them decreases.

$\uparrow (x_1, x_2)$

$\downarrow (x_1, x_2)$

There is the opposite relationship b/w "similarity" & "distance".

Let's first write down the relationship b/w ["Cosine similarity"] & ["Cosine-distance"].

$$1 - \text{Cos-Sim}(x_1, x_2) = \text{Cos-Distance}(x_1, x_2)$$

Let's say Cos-similarity lies b/w [-1, 1]

{ Let's say if two points are very similar, $\text{Cos-Sim} = +1$ }
{ If they are dissimilar, $\text{Cos-Sim} = -1$ }

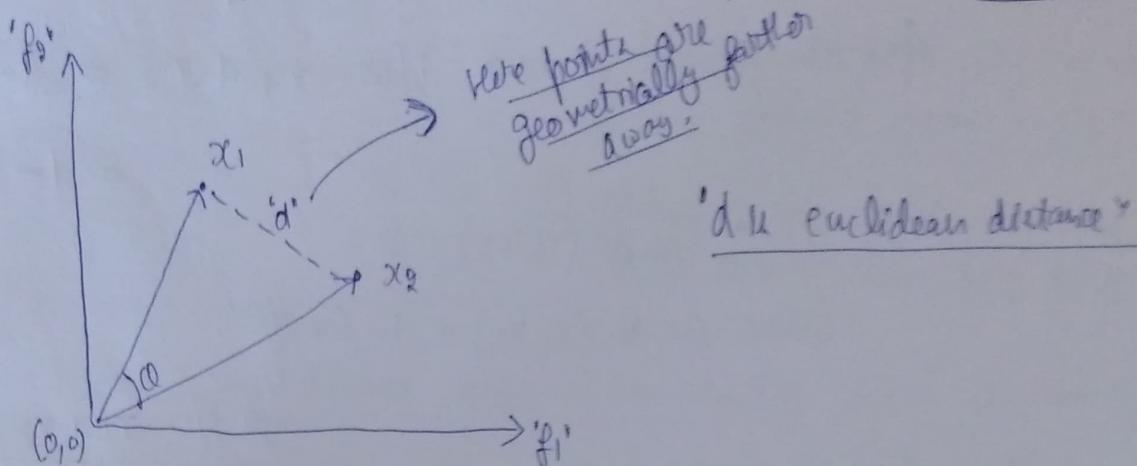
'+1' means they are perfectly similar.

8

'-1' means they are perfectly dissimilar.

(7)

With this background, now let's understand what is "Cosine Similarity".
Let's say we have two points $(x_1 + x_2)$ in 2-d space as shown below:



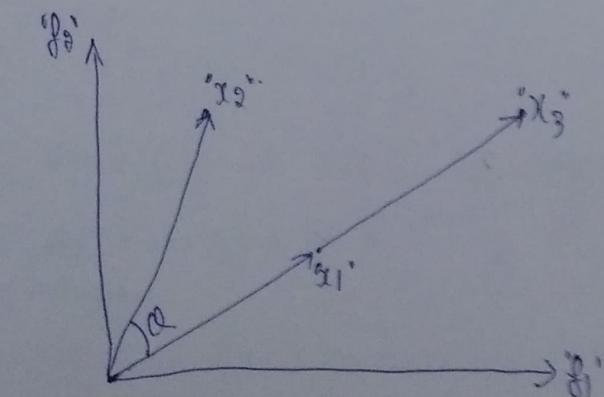
"Cosine Distance" says that instead of taking "geometric distance", what if we take angle " α " as shown above.

Cosine Similarity is ["Cos of angle ' α ']

Cos-Sim = Cos α
 (x_1, x_2) " α : angle b/w x_1 & x_2 "

Here points are angularly farther away.

Let's understand the diff b/w Euclidean distance & cosine similarity

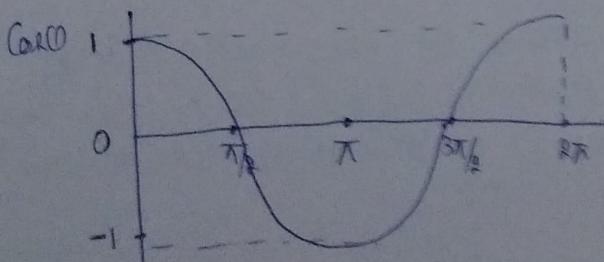


" x_1 & x_3 " are in the same direction

$$\text{Cos-Sim}(x_1, x_3) = \cos \alpha$$

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

$$\text{if } \alpha_{x_1 x_2} = 0^\circ$$



for "x₁ & x₃" since the angle b/w them is 0°

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

but the Cosine distance b/w x₁ & x₃ = 1 - Cos-dim(x₁, x₃)
 $= 1 - 1 = 0$

$$\therefore \text{Cos-dim}(x_1, x_3) = 0$$

while $\text{Cos-dim}(x_1, x_2) = 1 - \cos 0$

if we look at the "Euclidean distance" b/w "x₁ & x₃", the b/w (x₁ & x₃)

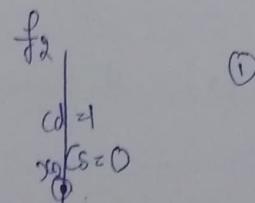
we will find out that

$$d_{13} > d_{12}$$

$$\text{but } \text{Cos-dim}_{13} < \text{Cos-dim}_{12}$$

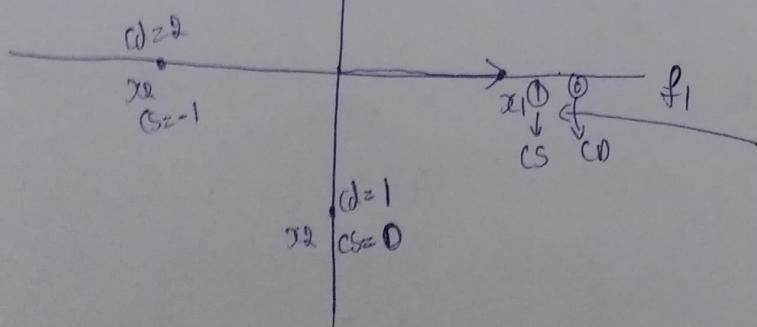
"Cos-distance" uses the angle b/w two points & not the distance as we imagine it geometrically.

let's understand the properties of "Cosine-similarity" :-



$$\text{At } [\text{Cos-dim}(x_1, x_2) = 1 - \text{Cos}(90^\circ)]$$

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



If x₂ is swapped with x₁ after cosine similarity is ①

Imagine a 2d space having two features as shown above, & imagine we have two feature "f₁" & "f₂". Let us assume that there is a vector "x₁" as shown above. Now imagine any other point (vector)x₂ in ① quadrant, as the angle between x₁ & x₂ increases the distance also increases.

(9)

$\text{angle}(x_1, x_2)$	$\text{distn}(x_1, x_2)$	
$0^\circ \rightarrow 90^\circ$	$0 \rightarrow 1$	{ Becoz angularly separating we are } { going further & <u>further away</u> }
$90^\circ \rightarrow 180^\circ$	$1 \rightarrow 2$	
$180^\circ \rightarrow 270^\circ$	$2 \rightarrow 1 \rightarrow$	distance is going to reduce becoz we are coming closer back to ' x_1 '
$270^\circ \rightarrow 360^\circ$	$1 \rightarrow 0$	

Then it's how the cosine distance works, or how exactly cosine similarity works.

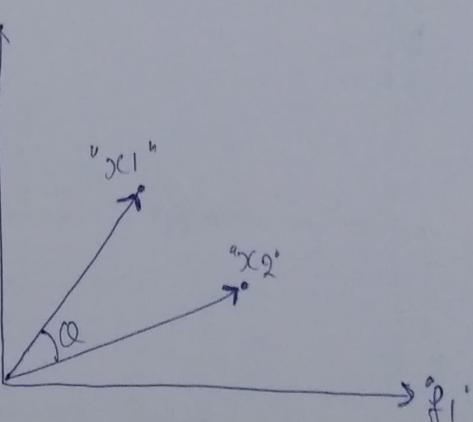
"Cosine similarity" is nothing but the 'cos of angle b/w two points'.

Let's assume ' x_1 ' & ' x_2 ' are two vectors.

& in linear algebra we have learnt that

$$\text{Cor}(\alpha) = \frac{x_1 \cdot x_2}{\|x_1\| \|x_2\|}$$

↓ length of $x_1 + x_2$
 ↓ L2 norm of x_2
 ↓ L2 norm of x_1



Now if ' x_1 ' & ' x_2 ' are unit vectors, i.e. if their length is 1 unit

then
$$\boxed{\text{Cor}(\alpha) = x_1 \cdot x_2}$$

"Cosine similarity" is nothing but the dot product b/w these two vectors ($x_1 + x_2$)

→ Relationship b/w "Euclidean distance" & "Cosine similarity":

Let ' A ' & ' B ' be two vectors. & Euclidean distance is denoted by

$$\|A - B\|^2 = (A - B)^T (A - B) = \|A\|^2 + \|B\|^2 - 2A^T B$$

When ' A ' & ' B ' are normalized to unit vectors $\|A\|^2 = \|B\|^2 = 1$

$$= 2(1 - \text{Cor}(A, B)) = 2 \times \underline{\text{Cos dist}(A, B)}$$

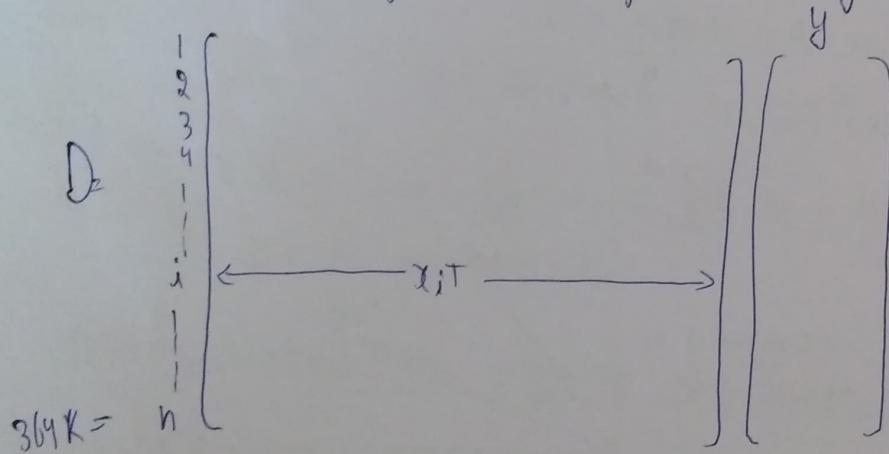
If x_1 & x_2 are unit vector)

$$[\text{euc-distance } (x_1, x_2)]^2 = 2 \underbrace{(1 - \text{cor}(x))}_{\text{cor-dist}}$$

$$\therefore [\text{Euclidean distance } (x_1, x_2)]^2 = \underbrace{2 \text{ cor-distance}(x_1, x_2)}_{=}$$

→ 'How to measure effectiveness of K-NN' :

Let's say we have "Amazon food Review Dataset" comprising of "n" reviews & there are roughly 364K after removing data duplication



Now for each review $[i]$ we have a vector representation called " x_i^T " & for each review, we also have a corresponding " y_i ", where $y_i \in \{0, 1\}$

$\begin{cases} 0 \rightarrow \text{-ve review} \\ 1 \rightarrow \text{+ve review} \end{cases}$

Now we have a very fundamental question, as how good is K-NN in solving our problem.

Problem: Given a new review "ug" what is the polarity?

"polarity" means whether the review is ("+ve") or ("-ve")

(11)

This basically means, given $[x_q]$ we want to determine the value of $[y_q]$.
 $x_q \rightarrow y_q ?$

\hookrightarrow " y_q " is polarity.

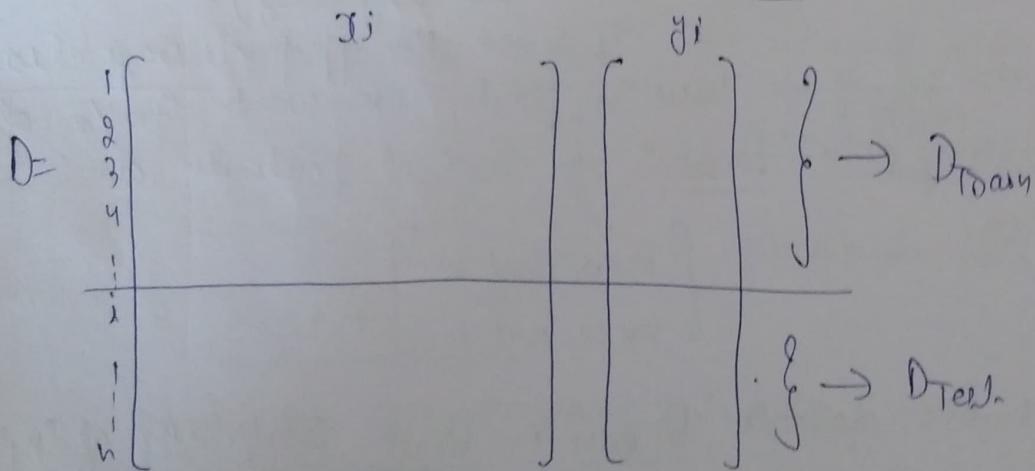
This is our problem, & for every problem, we want to measure the performance of the model we are building, here in our case it is 'K-NN'

In KNN, we find the ' K ' nearest neighbours & we do a majority vote.

Now let's see how to answer, (How good is KNN?) \rightarrow

let's devise a strategy to answer this question.

We have our huge dataset ' D ' of ' x_i ' & y_i '



So, what we do is, we break up our data into two parts.

So, given our big dataset ' D ' consisting of ' n ' points, we are going to break it into two datasets.

$$D_n \xrightarrow{\text{"}} D_{\text{Train}} \quad \text{such that} \quad D_{\text{Train}} \cup D_{\text{Test}} = D_n$$
$$D_n \xrightarrow{\text{"}} D_{\text{Test}} \quad \text{and} \quad D_{\text{Train}} \cap D_{\text{Test}} = \emptyset$$

This means given any point ' x_i ' in ' D_n ' it is either there in ' D_{Train} ' or ' D_{Test} ' but not both.

so for every point " x_i " in " D_n ", get either goes to " D_{Train} or D_{Test} " but not to both.

Now the question is how to split the "entire dataset" into two parts

" D_{Train} " & " D_{Test} "

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

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

$$\underline{n_1 + n_2 = n}$$

One simple way to do that, is to do it randomly.

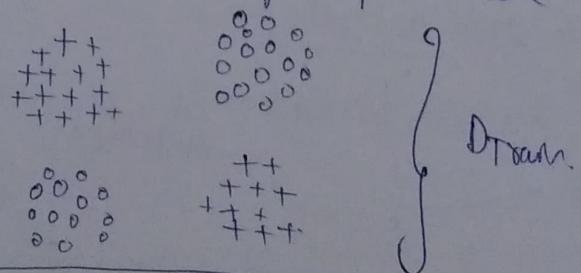
(It is one of the simpler ways, but not the only way).

for simplicity, let's say if we have " n " points, we want 70% of points in our "Training data" " D_{Train} " & 30% of points in our "Testing Data" " D_{Test} ".

$$\begin{bmatrix} n_1 \approx 70\% n \\ n_2 \approx 30\% n \end{bmatrix}$$

Now if we have two datasets " D_{Train} " & " D_{Test} " comprising of " n_1 " & " n_2 " points.

Now I can use " D_{Train} " to train the K-NN & to measure, how well it performs, we are going to use its " D_{Test} " data.



Now each point " x_i " in " D_{Test} " = $f(x_i, y_i)$ $\left\{ \begin{array}{l} n_2 \\ j=1 \\ \vdots \\ j=n_2 \end{array} \right\}$
 make it an x_q & find corresponding y_q $\left| \begin{array}{l} \rightarrow x_q = p_t \\ \text{Use } D_{Train} \text{ & K-NN to predict } y_q \end{array} \right.$ \rightarrow y_q \rightarrow p_t

Count = 0;
 for each p_t in D_{Test} :
 $y_q = p_t$
 Use D_{Train} & KNN to
 determine y_q .
 Now we compare predicted y_q' 's
 If $y_q = y_q'$
 Count = Count + 1

$\xrightarrow{\text{At end of this loop}}$
 Count = # of points for which D_{Train} & KNN gave a correct class label.

Now we can say a measure called Accuracy which is equal to

$$\text{Accuracy} = \frac{\text{Count}}{N_d}$$

No. of points for which
 D_{Train} & K-NN gave a correct
 class label

No. of points in D_{Test}

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

if let's say our accuracy is 0.91, it means that 91%
 of times KNN, with D_{Train} is able to correctly predict y_q .
 from test data set.

Now using all this data, we can conclude that.. K-NN on
 Amazon food reviews using " D_{Train} " gives an accuracy of
91%. (Toy example).

This is how we measure the effectiveness of "K-NN" on a given
 problem.

\rightarrow "Text/Evaluation Time & Space Complexity" \Rightarrow

Given a query point $[x_q]$, how much Time & space do we need to compute corresponding $[y_q]$.

let's write a simple "pseudo code".

Input: D_{Train} , K , $x_q \in \mathbb{R}^d$; Output: y_q .

$\text{kNNptr} = []$

for each x_i in D_{Train} :

$O(d)$ - Compute $d(x_i, x_q) \rightarrow d_i$ → Put in $\text{kNNptr}[]$
 $O(1)$ - Keep the smallest K -distances. $\rightarrow (x_i, y_i, d_i)$

At the end of this loop, we get a list of all the nearest points.

Now let's see how much time does it take.

for each dimension we have
to do subtraction &
ignore

Now since ' K is small' so we can optimize it to be a $O(1)$ time

∴ Total time complexity of this is $O(nd)$

$O(nd)$ is the time complexity to find K -NN $(K$ nearest neighbor)

\hookrightarrow $\left\{ \begin{array}{l} \text{Cnt_tot} = 0; \text{Cnt_neg} = 0 \\ \text{for each } x_i \text{ in } \text{kNNptr} \\ \quad \left\{ \begin{array}{l} \text{If } y_i \neq t \\ \quad \text{Cnt_tot} += 1; \\ \text{else} \\ \quad \text{Cnt_neg} += 1 \end{array} \right. \end{array} \right.$

$O(K)$
Now since ' K ' is small
 $O(K)$ becomes $O(1)$

PTO

0(1) {
 if $\text{cnt_pos} > \text{cnt_neg}$
 return $y_q = 1 \rightarrow \text{true}$
 else
 $y_q = 0 \rightarrow \text{false}$.

So most time taking part in this Algo is, finding the K nearest neighbors.

$$\begin{aligned}\text{Time Complexity} &:= O(nd) + O(1) + O(1) \\ &= \underline{O(nd)}\end{aligned}$$

if d is small the $d \ll n$

$$\text{then Time Complexity} = \underline{O(n)}$$

So, given a single query point, it takes roughly " $O(n)$ time".

This is typically alarmingly.

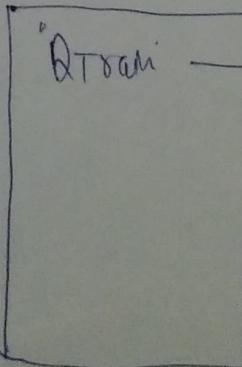
Let's take our "Amazon food reviews dataset".

Our $n \approx 36K$	$\left. \begin{array}{l} \text{lot of memory} \\ \text{Now imagine, we have to deploy} \\ \text{this system.} \end{array} \right\}$
$d \approx 100K$ (Bgs of words)	
≈ 300 (tfidf)	

So, given an $[x_q]$ if we have to compute corresponding $[y_q]$, what all things we need to store in memory.

$x_q \rightarrow y_q$

RAM/Memory



In ML, we will say what is the amount of memory consumed at runtime or evaluation time.

The more, it takes large out of memory.

Time complexity : $O(nd)$ straightforward

Space Complexity \Rightarrow what extra space is needed to evaluate.

to get yq from xq
 $O(nd)$

→ "Decision surface for K-NN or K-Nearest"

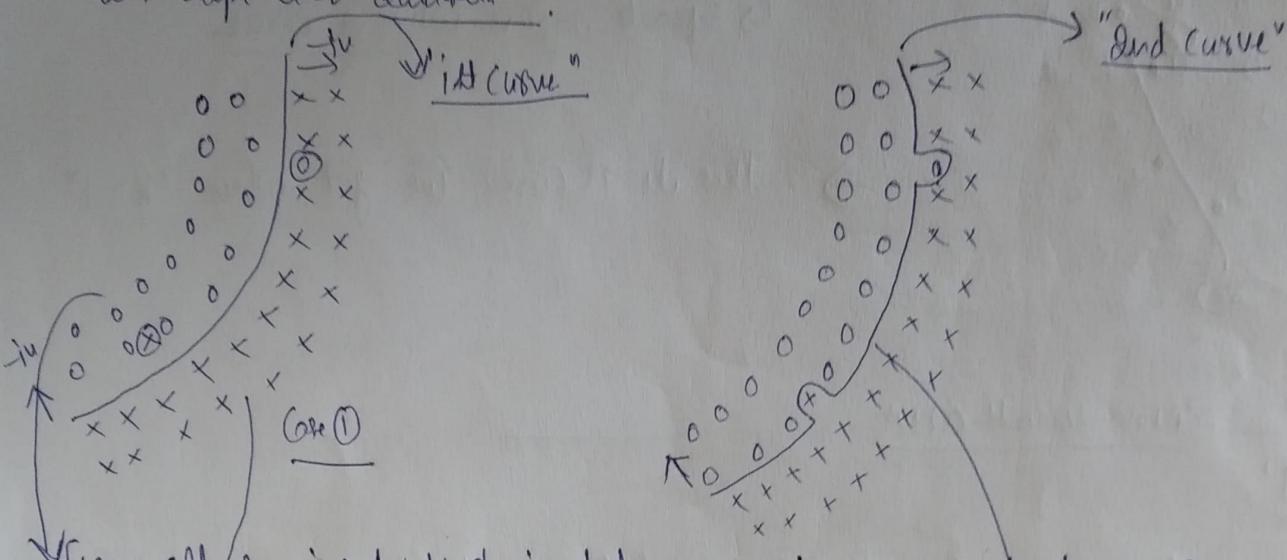
Now let's find out how do we pick the right " K ".

" K " in KNN is often defined as a "hyperparameter"

Now let's see how to determine the right " K " & let's find out how things change as " K " changes.

Let's understand it geometrically.

Let's take two datasets.



[Here all are fine points & in between we have a fine point.]

[Here all are fine points & in between somewhere, we have a fine point.]

If we ask a kid to draw a line to separate these points, he will draw the line/curve as shown above. He will say everything that lies on one side of the curve belongs to one 'set of points' & everything that lies on the opposite side of the curve belongs to other 'set of points'.

There is one other way to draw this curve, because if we consider the curve drawn earlier, the two points are going to be classified incorrectly.

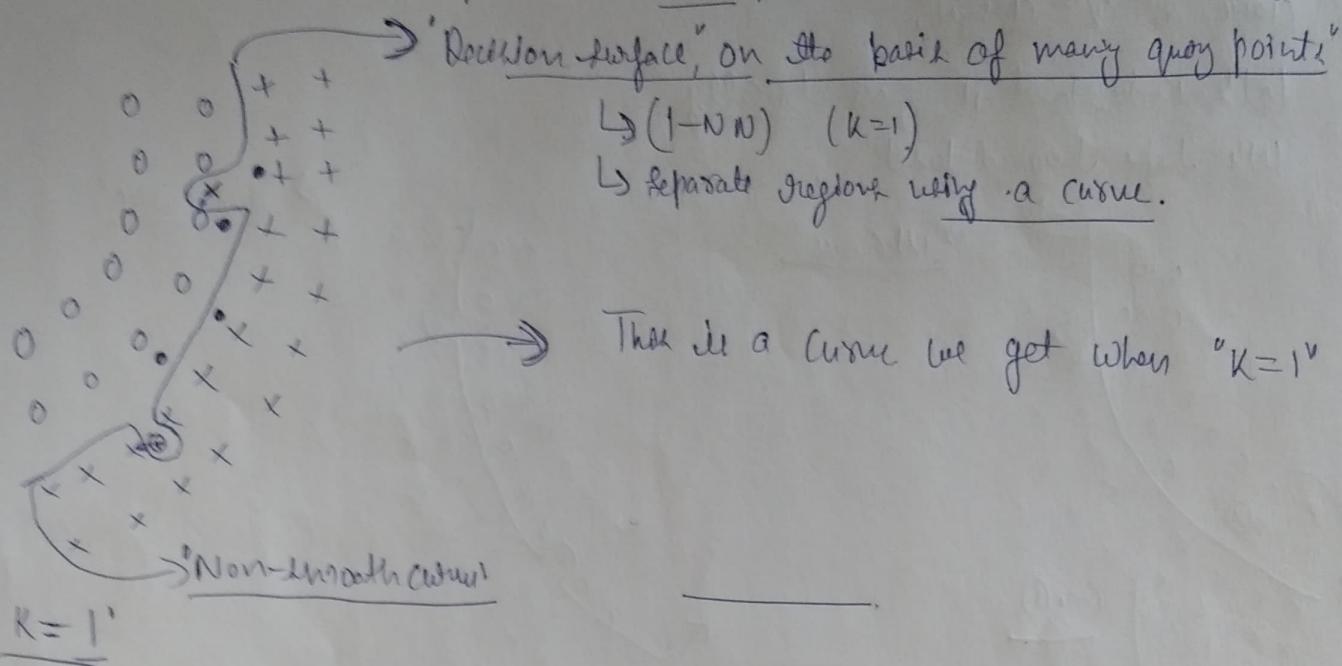
So, somebody who gives more importance to errors, he will draw the curve in such a way that no point is left behind, & he will say that everything which lies on one side of the curve belongs to one 'set of points' & everything which lies on other side of the curve belongs to the other 'set of points'.

Get curve make two errors & and curve makes no errors at all.

first curve is a smooth curve & not jagged, whereas 2nd curve is having all kinds of "twists" & "turns".

Note → The curve that separates "true points" from "false points" is called a decision surface.

The decision surface helps us understand, on what happens when "K" changes.

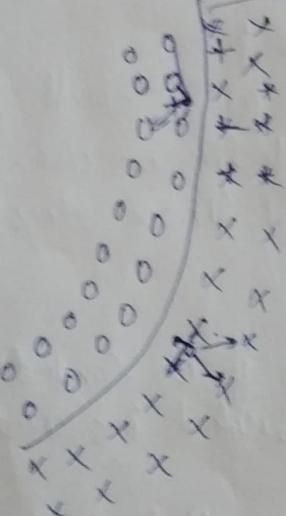


Let's draw the curve for the same dataset, but with " $K=5$ ".

Let $K=5$; Majority Vote

smooth curve

So, one thing we notice is as [K] increased,
the smoothness of the curve increased.



∴ In K-NN, the smoothness of the decision surface increased. As K increases.

→ worst case

We say ' $K=n$ ' n is the total no. of points

$n \rightarrow$ Total no of points in our dataset.

n_1 points are positive among these n points
 n_2 points are negative among these n points } a netting = n
 let's assume that $n_1 > n_2$

o o x

○ ○ × ×

0 x
 0 x x

○ X X

— 1 —

0 X X

O X

~~0~~

x^0 x

10

X

All its po

1

0 w of

Now $K = n$

$$n_1 = +1 \text{ u.e}$$

$$nq = -\ln p$$

CHIANG

'• It may profit'

Now let $n = 1000$, $n_1 = 600$ & $n_2 = 400$

Now we have to take "1000" nearest points. So all the points will come.

Out of all those points, 600 are true

"Yoo are -Iw"

so that point will be called a fine point.

Now let's take a query point deep inside the "five-class" region. Now if we take 4000 nearest neighbors of it, "500" will be the & the query point will be mapped to the

So wherever our query point is, our "1000. NN" will say every thing is negative positive. There is no more a decision surface, everything is +ve, becz we have more +ve than -ve.

So here we understood geometrically, at what happens when $k=1, 2, 3, \dots, n$ "

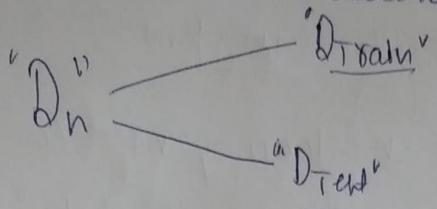
As 'K' increases our curve gets smoother & when "K" becomes "n", every step will become one cloth

→ "Need for cross validation"

We have already seen what happens when 'K' changes, i.e. when $[K=1]$ we get an overfit function, when $[K=5]$ we get a fit function & when $[K=n]$ we get an underfit function.

Now the big question here is how to determine the value of 'K'.

We have now two datasets, or our "big dataset" " D_n " is broken down into two datasets



& we choose it randomly

$\begin{bmatrix} 70\% \text{ of data in } "D_{Train}" \\ 30\% \text{ of data in } "D_{Test}" \end{bmatrix}$

One idea to determine 'K' could be

Let's $K=1$, we keep training data for training the model. Now using " D_{Test} " we will try to find out the accuracy.

→ $\left[\frac{\text{Mean No of points that are correctly classified}}{\text{Total no of points}} \right]$

	Train	Accuracy on D_{Test}
$K=1$	D_{Train}	

for $K=1$, we take each point in D_{Test} & we try to predict its label.

Remember our D_{Test} also contains pairs like

$$D_{Test} = \{(x_i, y_i) \mid x_i \in \mathbb{R}^d, y_i \in \{0, 1, 2\}\}_{i=1}^{n_{Test}}$$

We could map $x_0 = x_1, x_0 = x_2, \dots, x_0 = x_{n_{Test}}$ for each of these we are getting a corresponding y value.

(5)

Now for each of the query point, we get a "y_q" & now these "y_q" are compared with the "y" values that we already have in our D_{test} dataset.

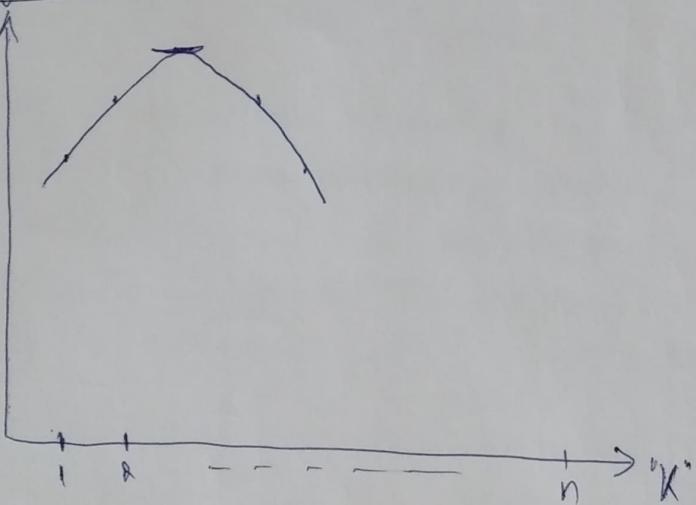
{ + if both agree, we say it is correctly classified else misclassified, we have seen this concept while computing accuracy. }

let's say.

	Team	Acc. on D _{test}
K=1	D <small>isain</small>	0.78
K=2		0.82
K=3	n	0.85

Now let's draw the value, on "x-axis" we draw the "K-value" & on "y-axis" we have our accuracy on D_{test}.

Accuracy on D_{test}



Typically what happens is, as K increases the performance will increase upto some point & then it will start decreasing.

Higher the accuracy the better it is.