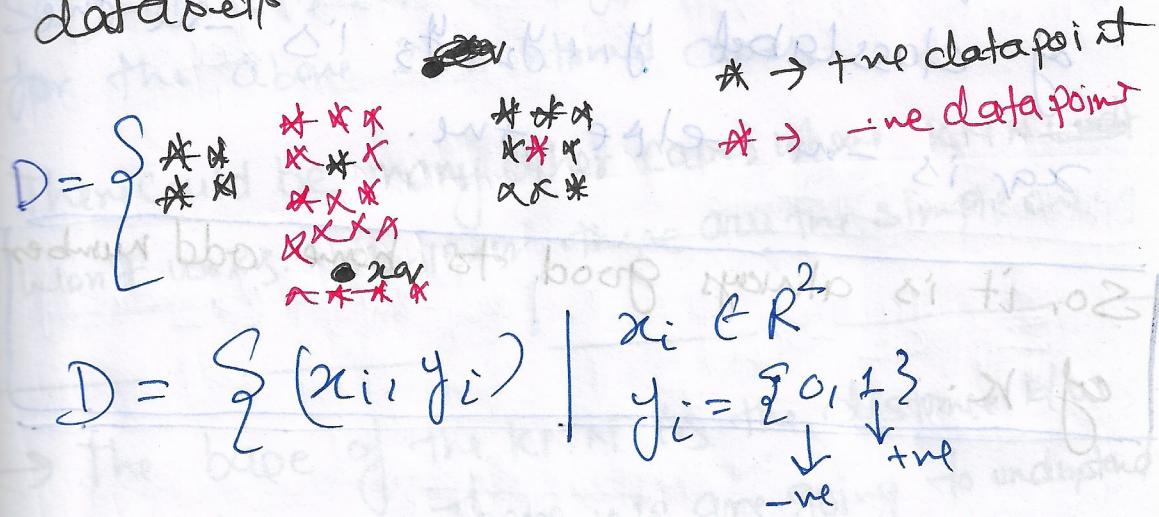


K-Nearest Neighbour (KNN)

Let's take an example to understand the concept. I'm having 2D toy datasets.



Let's I'm having a query point x_{qr} , now I want to know, whether the x_{qr} is +ve or -ve.
i.e. $x_{qr} \rightarrow y_{qr}$
→ If we see the neighbourhood point of x_{qr} , we will see mostly $+$ point (-ve). So, by just upping the proximity, I can say that x_{qr} is also -ve.

This is what the KNN is.

How KNN works

#1 Find k - "nearest" points to x_q in D.

#2 Let $k=3$ (How to select the value of k , we will see later) (x_1, x_2, x_3)

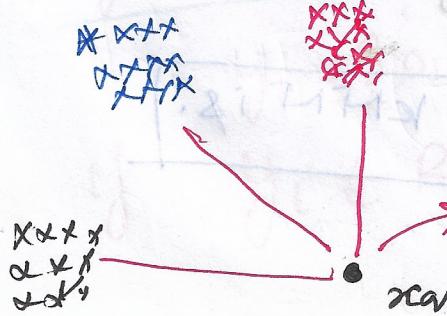
#3 If the majority of ~~point~~-best of vote of class label y_1, y_2, y_3 is $-ve$ the x_q is $-ve$ else $+ve$.

So, it is always good to have odd number of k .

Now, let's understand where ~~K-NNT~~ K-NNT

fails:

Cape I My query point is very far away from my data set, in that case it doesn't make any sense to calculate the behavior base on the neighbour point.



However KNN based on the k value predict the behavior, that is untrue.

In this case we can't say ~~what this is~~ anything about this point.

Cape II The query point is so much jumbled across each other.

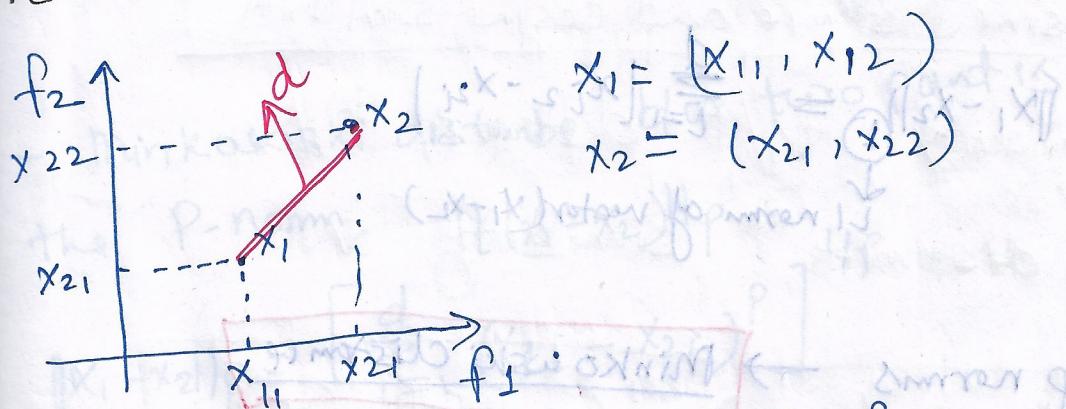
Whatever the KNN predict most probably wrong or ^{no} useful information.



So KNN is not the algorithm we should select for the above 2 mentioned cases.

There could be many other cases where KNN ~~won't work~~ won't work. However these are the simple ones.

→ The base of the KNN is the distance b/w the points. Now, here we are going to understand how to measure the distance b/w the points.



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

Pythagoras Theorem

→ Euclidean Distance \Rightarrow distance b/w shortest line b/w 2 data point.

see later

motorman $\leftarrow i=9$

envelope $\leftarrow s=9$

Similarly $x_i \in \mathbb{R}^d$,

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

\hookrightarrow Euclidean Distance. $\|x_1 - x_2\|_2$

\hookrightarrow L_2 norm of a vector

Manhattan Distance

$$\sum_{i=1}^d |x_{1i} - x_{2i}| \quad i.e. |(x_{21} - x_{11}) + (x_{22} - x_{12})|$$

algebra of prior knowledge

$$d = d_1 + d_2$$

$d = d_1 + d_2$

$$\|x_1 - x_2\|_1 = \sum_{i=1}^d |x_{1i} - x_{2i}|$$

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

L_p norms \rightarrow Minkowski Distance

P could be any norm

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

$P=1 \rightarrow$ Manhattan

$P=2 \rightarrow$ Euclidean distance.

L_p norm of a vector

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

So, distance b/w between 2 points, however
the Norms ~~are~~ is for a vector.

- So, we have learnt that Euclidean distance

b/w 2 points = L_2 norms of $(x_1 - x_2)$

$$\text{norms (choose)} = \|x_2 - x_1\|_2$$

- Similarly Manhattan distance b/w 2 points

$$L_1 \text{ norm } \|(x_1 - x_2)\|_1 = \sum_{i=1}^d |x_{1i} - x_{2i}|$$

- Minkowski distance b/w two points

the p -norm $\|(x_1 - x_2)\|_p$

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

→ There are many types of distances available,

here we discuss these we are going to
use a lot. When to use what we will

see later.

Hamming Distance \rightarrow Mainly going to use in text processing (boolean vector)

$x_1, x_2 \rightarrow$ boolean vector

Hamming-dist(x_1, x_2) = # locations/dimensions where the binary vectors differ.

Hamming distance also used when x_1 and x_2 are strings = # locations/features/dimensions they are different.

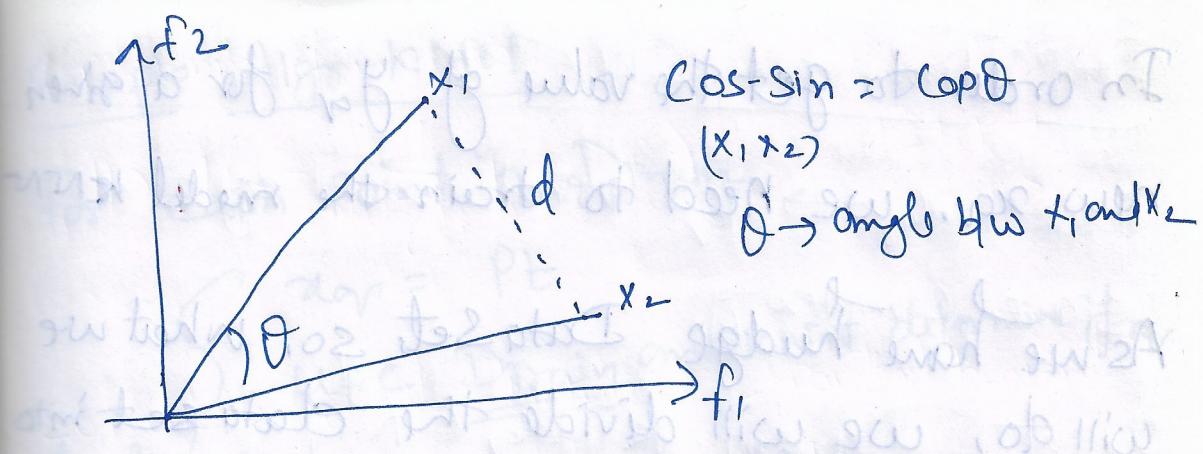
Cosine Distance and cosine Similarity

Similarity \downarrow / \uparrow distance \uparrow / \downarrow
As distances \uparrow similarity \downarrow .

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

\rightarrow If the 2 points are very similar $\Rightarrow \text{Cos}\theta = 1$,

If the 2 points are very dissimilar $\Rightarrow \text{Cos}\theta = 0$



→ How to measure distance in KNN

If we talk about the Amazon Dataset

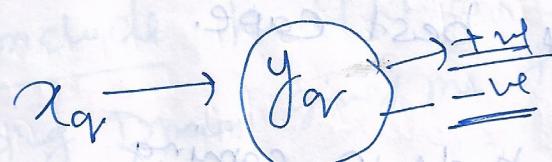
$$D = \left[\begin{array}{c} x^T \\ \vdots \end{array} \right] \left[\begin{array}{c} y_i \\ \vdots \end{array} \right] \quad y_i \in \{0, 1\}$$

↓
-ve +ve

→ Text Bow | tf | IDF | w2v

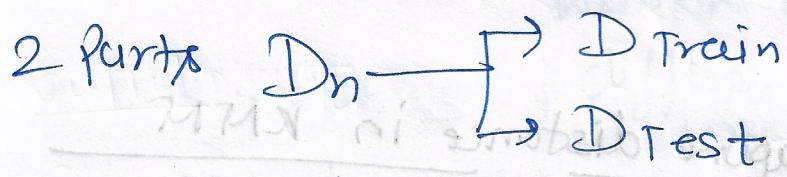
Problem This is the dataset we are having

What is the problem? Given a new fact
review, what is its polarity i.e. +ve/-ve.



In order to get the value of y_{or} for a given new x_{or} , we need to train the model KNN.

As we have huge Data Set, so, what we will do, we will divide the data set into



$$D_{\text{train}} + D_{\text{test}} = D_n$$

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

Any given data point can either belong to

D_{train} or D_{test} but not both.

How to split D_n to D_{train} and D_{test}

One simple answer is randomly. It could be but not always be the best case.

We will see more methods in coming
Section

Simple Algo kNN

for each point in D_{train}

$x_{ar} = pt$
 use $\underline{D_{train}}$ and kNN determine y_{ar}
 loop
 $y_{ar} == y_{pt}$
 count $+ \pm 1$

Count = # pt for which $D_{train} + kNN$ gave a correct class label

Accuracy = $\frac{\text{Count}}{n^2}$ # Pts which $D_{train} + kNN$ given a correct class label
 ↓
 # points in D_{test}

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

Say ACC = 0.91 \Rightarrow 91% of times we predict y_{ar} for a given x_{ar} .

Conclude k-NN on Amazon food review

using D_{train} gives me the accuracy of 91%.

(91% of people who like bacon give it 5 stars)

That is why people like bacon.