

# K Nearest Neighbour (KNN)

→ K-nearest neighbour is one of the simplest machine learning algorithms

→ K-NN algorithm assumes the similarity between the new and available cases and put the new case into the category that is most similar to available categories.

→ K-NN algorithm can be used for **Regression** as well as **classification**

→ It is also called as lazy learner algorithm, because KNN algorithm at the training phase just stores the dataset and when it gets new data, then it classifies that data into a category that is much

similar to the new data.

→ Suppose we have new data point and we need to put it in the required category.

→ we have to choose the number of neighbours in this case I am considering  $K=5$ .

→ Now calculate the distance between new data point and nearest neighbours.

→ In order to calculate the distance we have two methods for distance metrics.

① Euclidean distance

② Manhattan distance

Fig: Before KNN

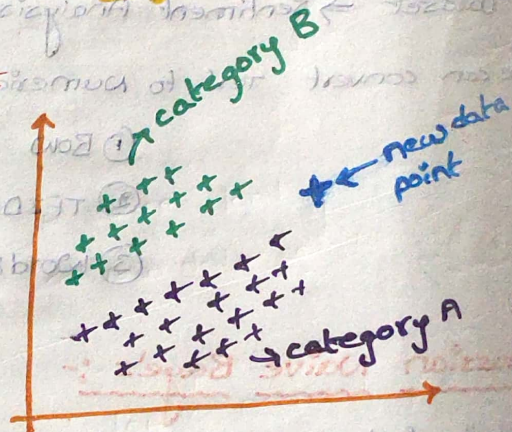
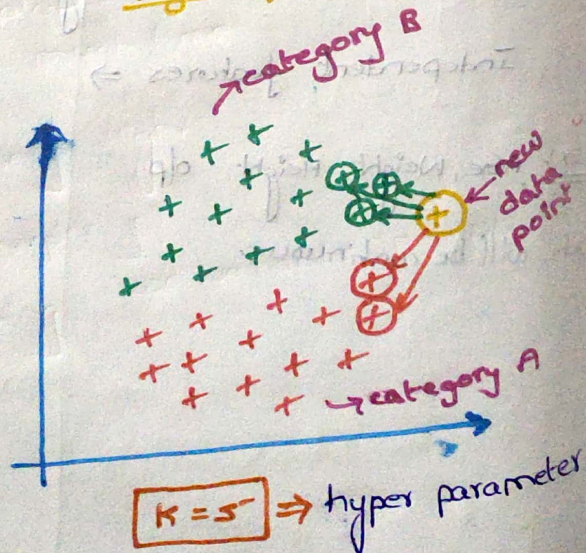
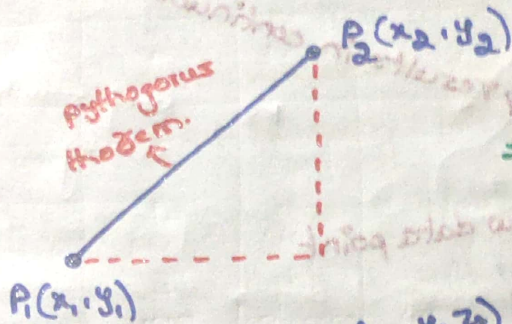


Fig: After KNN

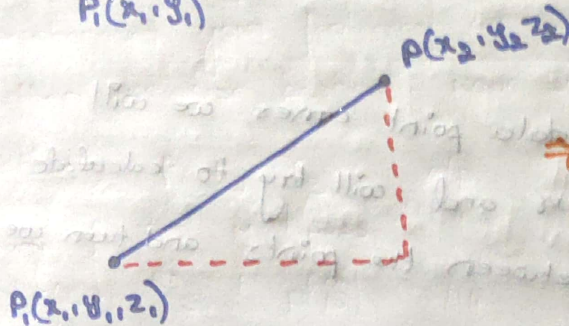




## ① Euclidean Distance:

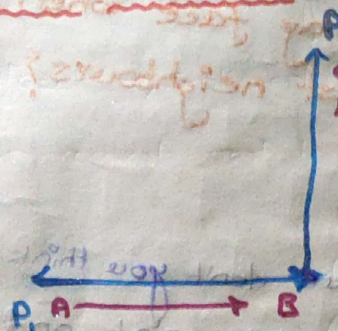


$$\Rightarrow \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$



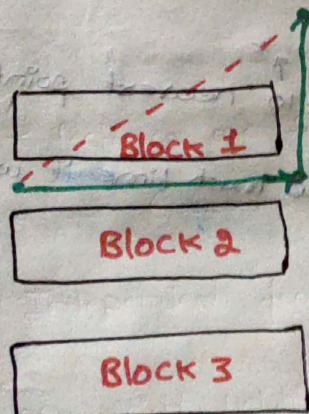
$$\Rightarrow \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2}$$

## ② Manhattan Distance:



$$\Rightarrow |(x_2 - x_1) + (y_2 - y_1)|$$

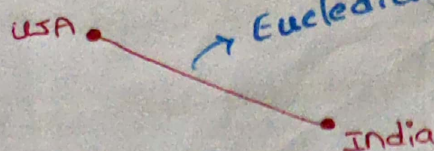
Ex 1:



Manhattan distance

Ex: 2

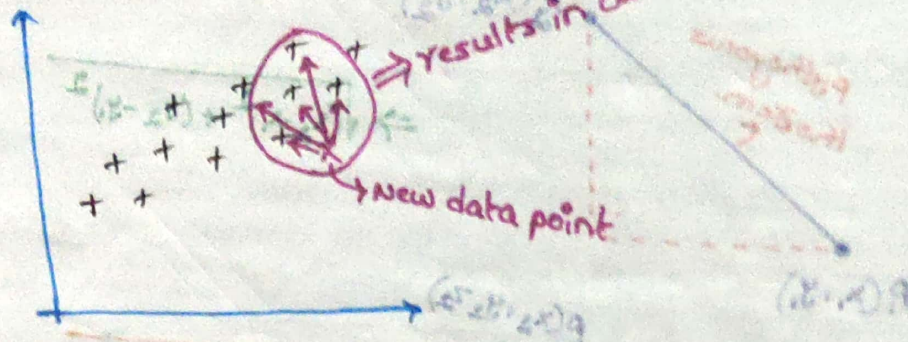
How do we compute the distance between USA and India?



Euclidean distance



## Regression:-



→ Let's say when a new data point comes we will find its nearest neighbours and will try to calculate the average distance between the points and then we will get the continuous values.

→ Let's say in my dataset I have 1 million datapoints. What is the problem that I may face when I am trying to find out the nearest neighbours?

→ Let's say  $K$  value = 10

→ Now in order to find out that  $K$ -value don't you think we need find the distance between a new data point and all the points.

→ Then sort those values

→ And then try to find out the top 10 nearest points.

→ In case of millions of datapoints how much time it will take obviously it will take much time.

→ Here time complexity is very high.

→ How do we fix time complexity?

1) we will fix this by using two variants in KNN.

## variants of KNN

①  $K$ , Dimension tree

② Ball tree



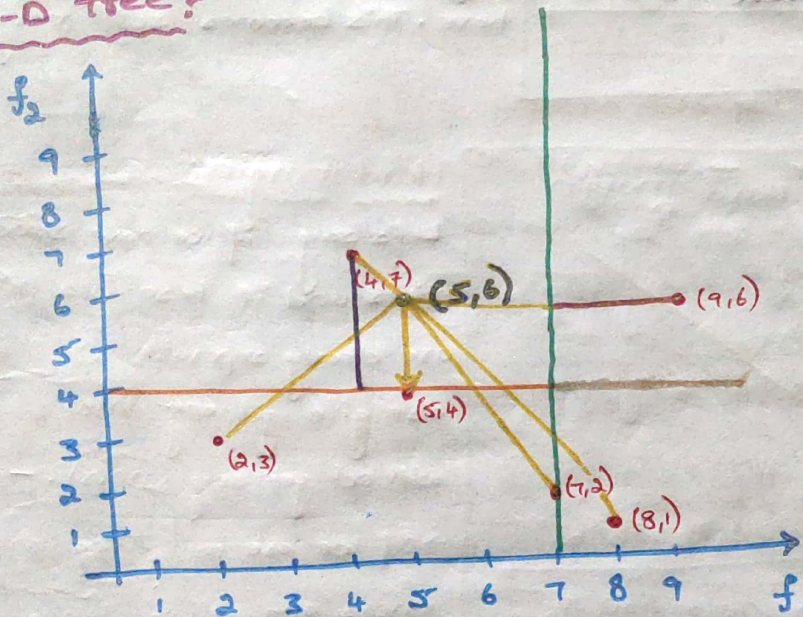
→ In both these techniques first we ~~create~~ construct a decision tree.

→ what will happen if we construct a decision tree?  
 → in this case we don't need to calculate the distance from one point to every point if we have decision tree.

→ what is the importance of binary tree or Decision tree?

- It takes  $\log n$  times it will take
- It will not take  $n$  time or  $n^2$  time.
- It reduces the time to reduce an element.

### ① K-D Tree:-



→ Let's consider few data points  
 (2,3) (4,7) (5,4) (7,2) (8,1) (9,6)

Aim:- whenever we get a new data point I should not calculate the distance ~~but~~ with every point.

→ now will try to find out the median of all the x-axis values

2, 4, 5, 7, 9

⇒ picking the value 7

$$\text{Median} = \frac{5+7}{2} = 6$$

→ Now try to find out the median of all the y-axis values

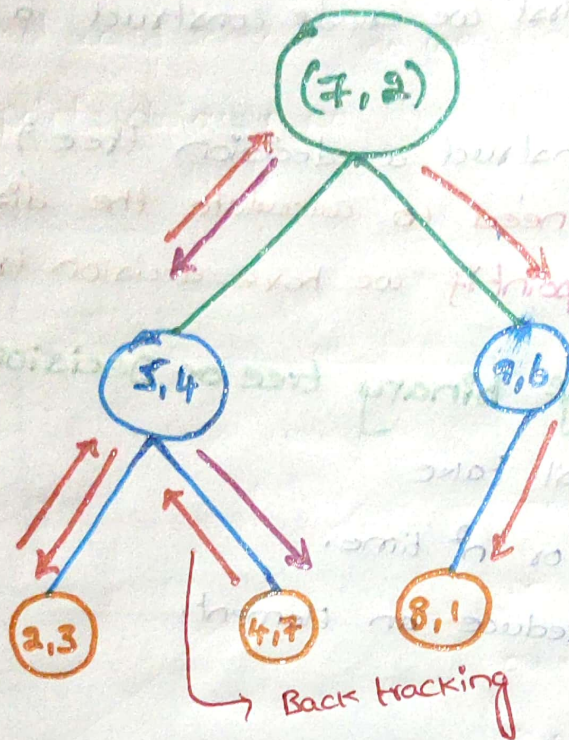
3, 7, 4, 2, 1, 6

① sort the values

1, 2, 3, 4, 6, 7

$$\text{Median} = \frac{3+4}{2} = 3.5$$



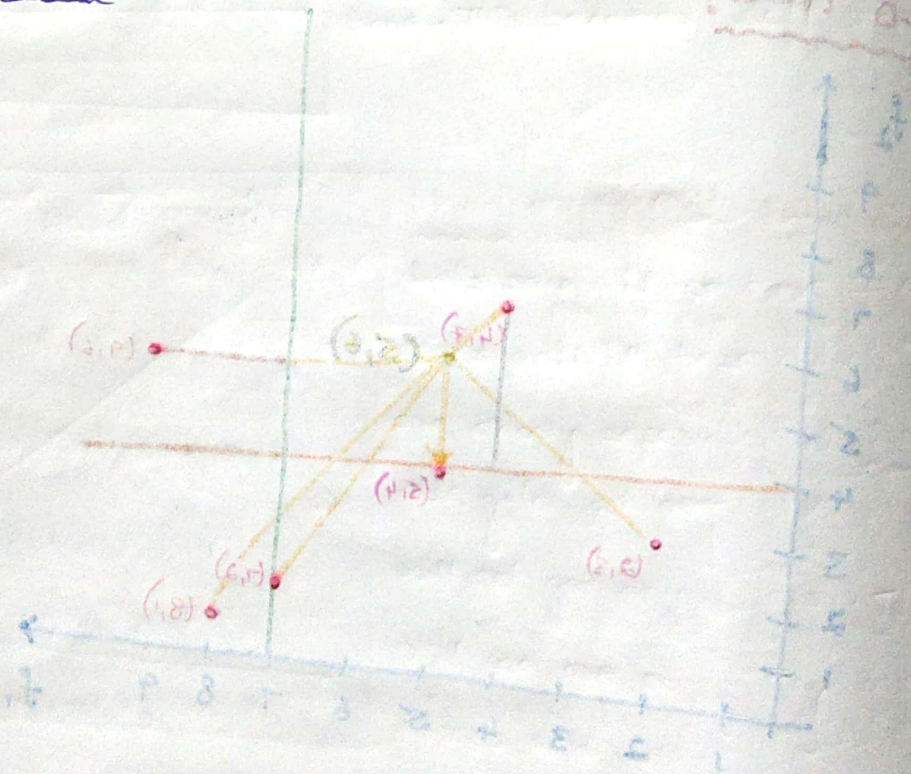


$k=6$

$k \downarrow$

no. of searches  
reducing.

→ From 7 calculate draw a vertical line. Root node is (7, 2)



2nd data point (5, 4) (9, 6) (1, 8) (6, 5) (3, 2) (7, 2)

Whenever we get a new data point, we calculate the distance from that point to all the points already in the list and find the point which is closest to it.