

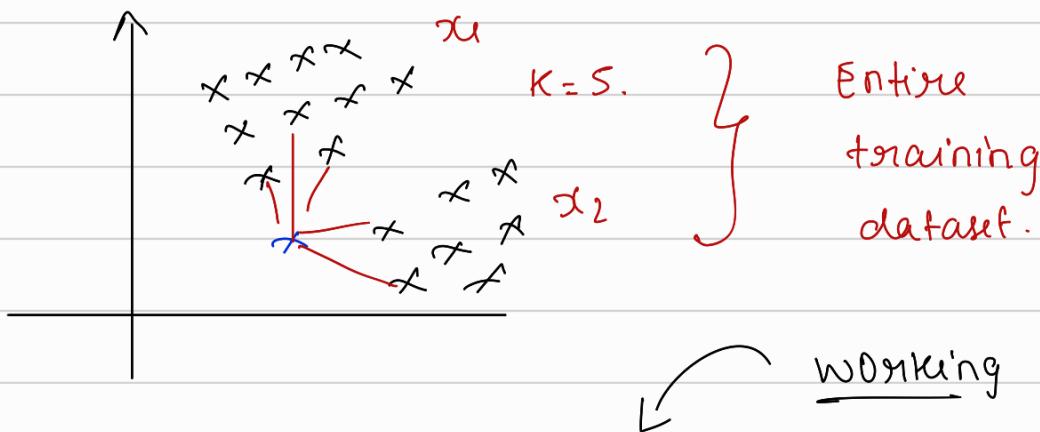
K- Nearest Neighbor

K-Nearest Neighbor is a simple, non-parametric machine learning algorithm used for classification and regression tasks. It works by finding the K closest training data points to a new unseen data points and making predictions based on their values.

It is a supervised machine learning method. KNN does not require an explicit training phase where a model is learned. The algorithm simply stores the entire training dataset.

Classification,

Let x_1, x_2 be the features and y be Binary output.



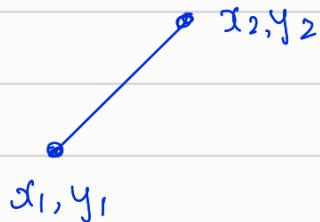
- ① we have to initialize the 'k' values
 $k > 0, 1 \dots \infty$ (its hyperparameter)
- ② Find the K-Nearest neighbour for test data.
- ③ from those $k=5$, how many neighbours belong to 0 category and 1 category

To find the distance between test data and k-nearest neighbor.

1. Euclidean distance.

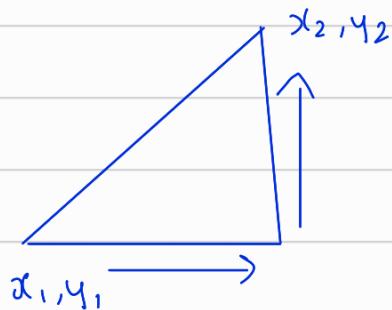
2. Manhattan distance.

Euclidean :
distance



$$\text{distance} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Manhattan :
distance



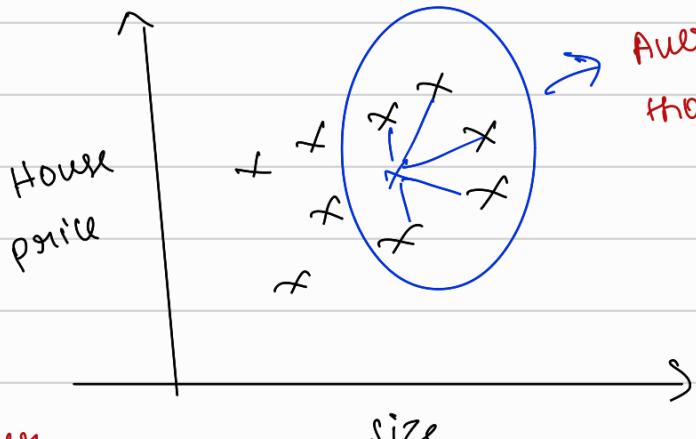
$$\text{distance} = |x_1 - x_2| + |y_1 - y_2|$$

Regression,

The algorithm predicts a continuous value based on the k-nearest neighbor of a given input point.

The predicted value for the test point is the average of the k-nearest neighbors target values.

$$y' = \frac{1}{k} \sum_{i=1}^k y_i$$



find out median
if there is a outlier.

Major issue is time complexity, finding the distance between test data and every point in the dataset and it takes $O(n)$

Now, we use

- ① KD tree
 - ② Ball tree
- } to optimize the KNN

These trees are created in form of binary tree.

Variants of KNN

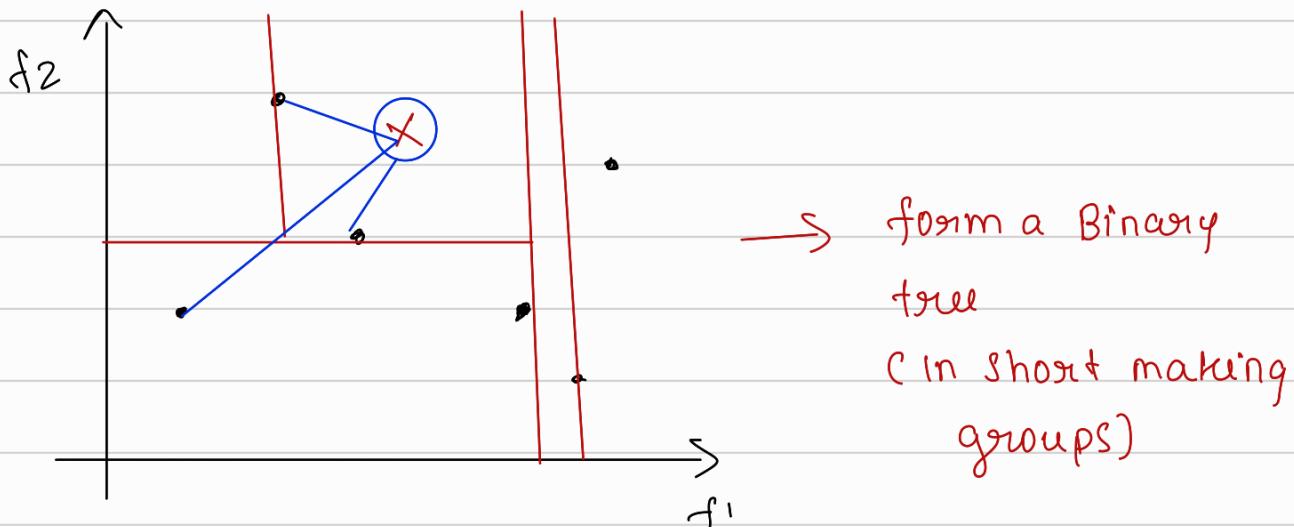
- ① KD-tree, is a space partitioning data structure used to organize points in k-dimensional space, it is particularly useful for efficient nearest neighbors searches and making it a popular choice for implementing KNN
 - * Binary tree structure, Each node in the tree represent a point in the dataset and splits the space into two half spaces
 - * The dataset is recursively divided along one dimension at each level, alternating to built the tree.
 - * Efficient search, Allows for faster range searches and nearest neighbor queries compared to brute force approach.

Working

- * Choose a dimension to split the data, (x, y, z for 3D)
- * Select a median point along the chosen dimension and ensure a balanced tree.
- * Recursively partition the data into two subsets.
 - (i) Left subtree, contains points less than the median on the chosen dimension.
 - (ii) Right subtree, contains points greater than median.

Example, (2D)

1. Split by x -coordinate
2. Split by y -coordinate.
3. Split by x -coordinate again and so, on



Even backtracking is used to find the next nearest neighbor.

(easy to find the nearest point)

Ball tree,

unlike KD tree, Ball trees partition the space into hyperspherical regions (balls) and making them more suitable for certain type of data distribution & distance metrics.

- * Hierarchical partitioning, the space is divided into nested hyperspheres, each node representing a ball (region containing data points)
- * Adaptability, works well with a variety of distance metrics.
- * Efficient in Searching

Working,

- * Select a subset of points to form a node.
- * compute a centroid for the points in the current node.
- * calculate the radius from the centre to any point
- * split the points into two subsets based on proximity to two chosen pivot points.
- * Repeat the process recursively.

