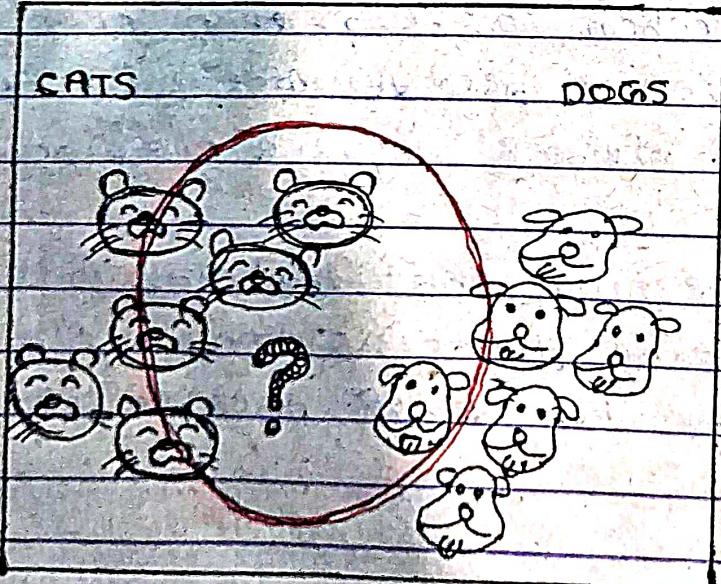


* K-nearest Neighbour

- The distance based machine learning algorithm.
- Non-parametric, supervised ml algorithm.
- Can be used to solve both classification & regression problem statements.
- If we often noticed that you share many characteristics with your peers, whether it be your thinking process, working etiquette, philosophies or other factors.
As a result we build friendship with people we deem similar to us.

The K-NN employs the same principle it aims to locate all of the closest neighbours around a new unknown data point in order to figure out what class it belongs to.

eg



In this value of k is 3. Since, there are 2 cat and just one dog

In the proximity of 3 closest neighbours, the algorithm would predict that it's a cat base on

proximity of the three closest neighbours in red circle boundary.

Here, ' k ' is the hyperparameter for KNN, for proper classification/prediction, the value of ' k ' must be fine tuned.

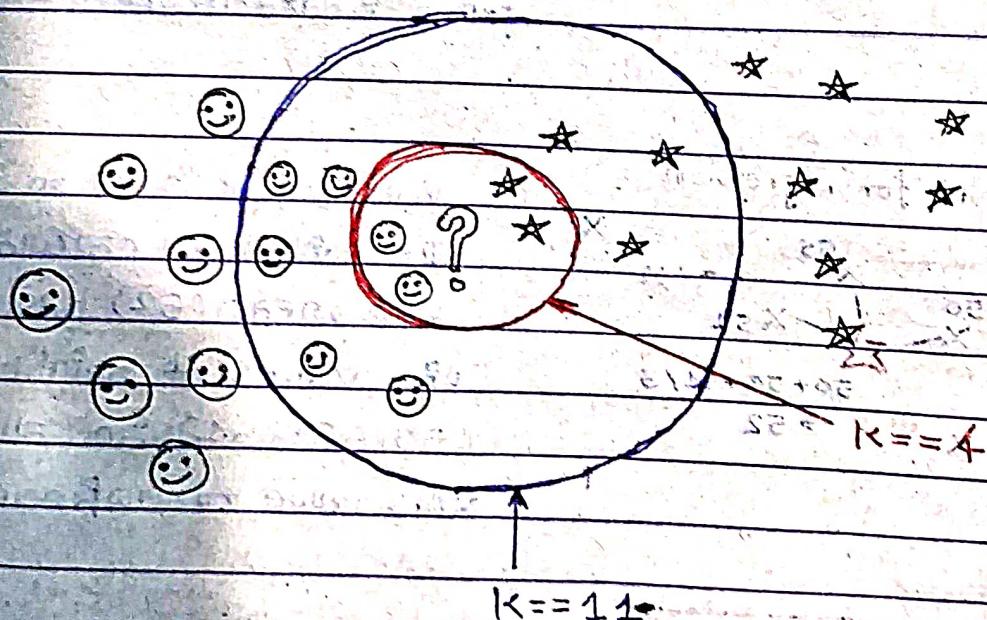
But, how do we select right value of ' k '?

We don't have a particular method to determine value of ' k ', here we will try to test the models accuracy for different ' k ' values.

The value of ' k ' that delivers best Accuracy for both training & testing data is selected.

Note: It's recommended to always select an odd-value of ' k '.

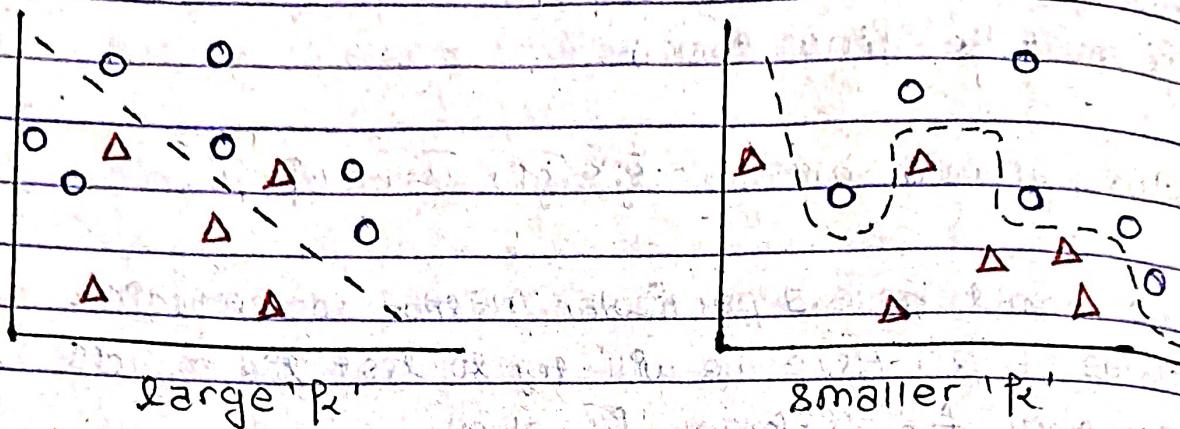
When the value of ' k ' is even a situation may arise for which elements from both groups are equal.



KNN for classification statements

KNN tends to use the concept of "majority voting". Within the given range of ' k ' values, the class with the most votes is chosen.

★ Impact of selecting smaller & larger k -value :



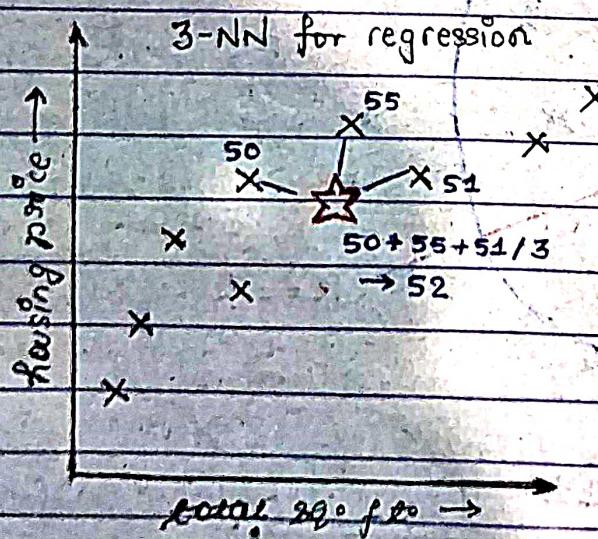
The case of "Underfitting" occurs when the value of k is increased.

The case of "Overfitting" occurs when the value of k is smaller.

In this case model would be unable to correctly learn on training data.

The model will capture all of the training data including noise.

The model will perform poorly for test data in this scenario.



Value of k is set to 3, It will now calculate the mean (52) based on the values of neighbors (50, 55 & 51) and allocate this value to unknown data.

[KNN works for regression statement]

• K-NN is non-parametric, meaning it doesn't make any underlying assumptions about the distribution of data (as opposed to other algorithms such as LMM which assume gaussian dist. of given data.)

• Intuition Behind K-NN Algorithm?

Distance metric used in KNN algorithm.

Euclidean Distance

$$d(p, q) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$

cartesian dist. b/w two points which are in plane / hyperplane.

Manhattan Distance

$$d(x, y) = \sum_{i=1}^n |x_i - y_i|$$

summing the absolute difference b/w the coord. of the points in n-dimension

① Choose value of 'K'

If your data has more outlier/noise \downarrow
Higher value of K would be better.

use cross-validation methods-

Selecting the optimal value of K, K represent nos of nearest neighbours that needs to be considered while making prediction.

② calculating distance

③ finding nearest neighbour

④ voting for classification or taking average for regression

Exercise

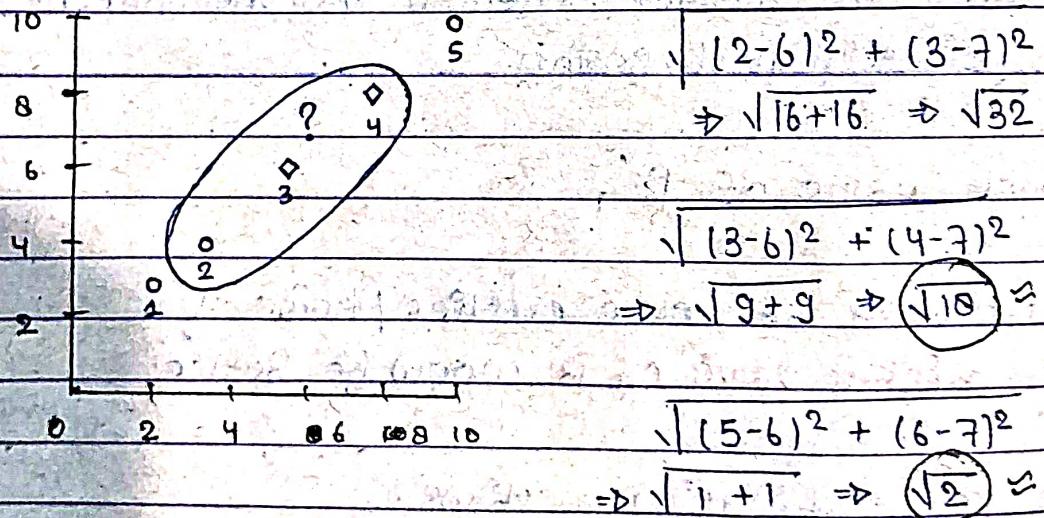
1.

| Data Point | x | y | Label |
|------------|----|----|-------|
| 1 | 2 | 3 | A |
| 2 | 3 | 4 | A |
| 3 | 5 | 6 | B |
| 4 | 7 | 8 | B |
| 5 | 10 | 10 | A |

Now, consider a new data point

with $x_1 = 6, y_1 = 7$

Using KNN with $k=3$ predict label for this new data point



$$\sqrt{(7-6)^2 + (8-7)^2} \Rightarrow \sqrt{1+1} \Rightarrow \sqrt{2} \approx$$
$$\sqrt{(10-6)^2 + (10-7)^2} \Rightarrow \sqrt{(4)^2 + (3)^2} \Rightarrow \sqrt{16+9} = \sqrt{25}$$

choose $(k=3)$ min dist

predicted through majority voting
class $\in B$

Q) Consider a set of five training examples given as $((x_i, y_i), c_i)$ values, where x_i, y_i are two discrete values (possible integers) and c_i is the binary class label.

Classify a test eg at coordn $(3, 6)$ using k-NN classifier with $k=3$, and manhattan dist. defined by $d((x_1, y_1), (x_2, y_2)) = |x_1 - x_2| + |y_1 - y_2|$.

$$\textcircled{1} ((1, 2), -1) \rightarrow d(1-3+1-6) \rightarrow 1-2+1-5 \Rightarrow 7$$

$$\textcircled{2} ((1, 7), +1) \quad d(1-3+1-6) \rightarrow 1-2+1-1 \Rightarrow 3$$

$$\textcircled{3} ((3, 3), +1) \quad d(3-3+1-6) \rightarrow 1+1+1-5 \Rightarrow 3$$

$$\textcircled{4} ((5, 4), -1) \quad d(3-5+1-6) \rightarrow 1-2+1-2 \Rightarrow 4$$

$$\textcircled{5} ((2, 5), -1) \quad d(2-3+1-6) \rightarrow 1-1+1-1 \Rightarrow 2$$

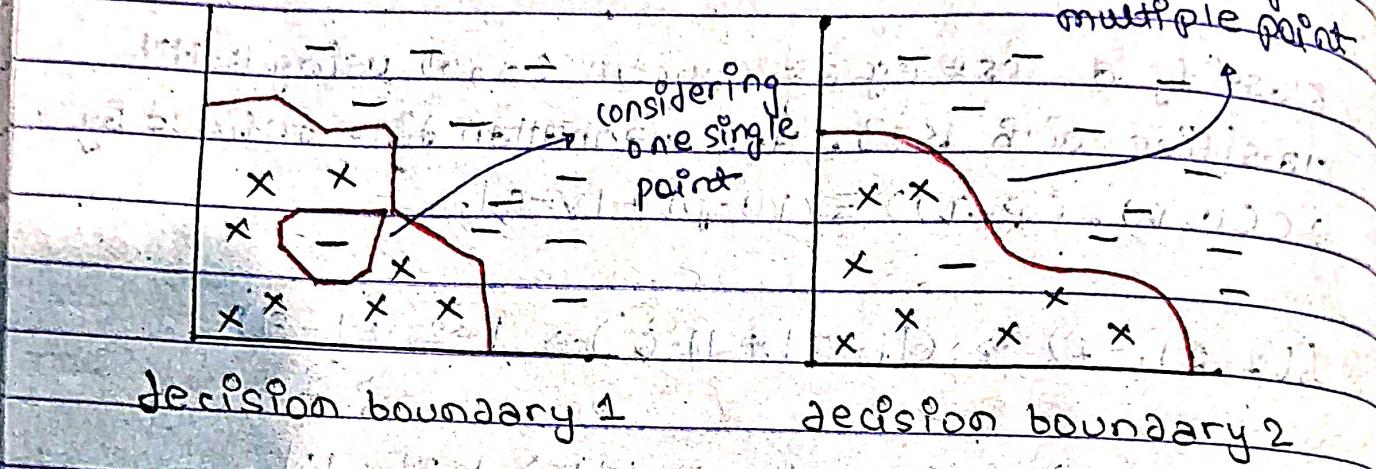
major voting \Rightarrow b/c 1, 2

+1 +1 -1 -

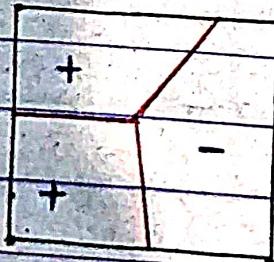
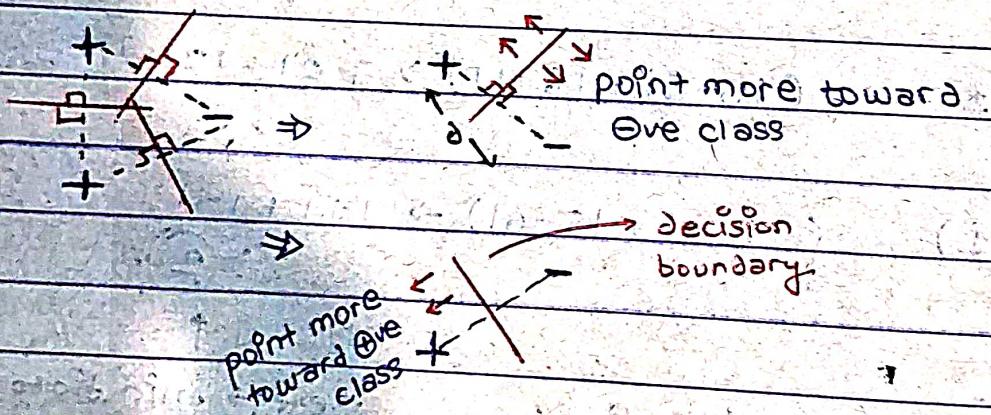
Select min disto

Predicted class label $\Rightarrow +1$

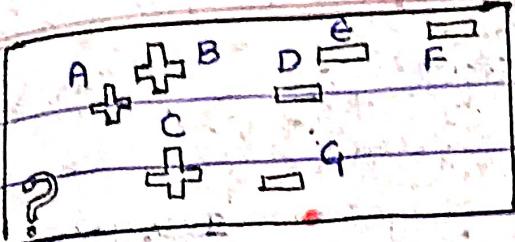
③ fig. 9 illustrates decision boundaries for two-on classifiers. Determine which one of the boundary belongs to 1-NN classifier & which one belongs to 3-NN classifier?



i.e.,
1-NN decision boundary
(using Euclidean distance metric)



④ Given following labelled dataset
For what (minimal) value of k will the query point "?" be negative?



$1NN \rightarrow NN \rightarrow A \rightarrow \text{Give}^+$

$2NN \rightarrow NN \rightarrow A, C \rightarrow \text{Give}^+$

$3NN \rightarrow NN \rightarrow A, C, B \rightarrow \text{Give}^+$

(max vote)

$4NN \rightarrow NN \rightarrow A, C, B, S \rightarrow \text{Give}^+ \text{ (max vote)}$

$5NN \rightarrow NN \rightarrow A, C, B, S, D \rightarrow \text{Give}^+$

$6NN \rightarrow NN \rightarrow A, C, B, D, G, E \rightarrow \text{Ignore}$

$7NN \rightarrow NN \rightarrow A, C, B, D, G, I, F \rightarrow \text{Give}^+$

(max vote)

Pros -

- easily implemented

- few hyperparameters

($C_F \rightarrow$ which is selected by validation)

Cons -

- Lazy algorithm

(takes lot of computing power as well as data storage)

- Curse of dimensionality

(face a hard time classifying the data-points properly when dimensionality is too high)

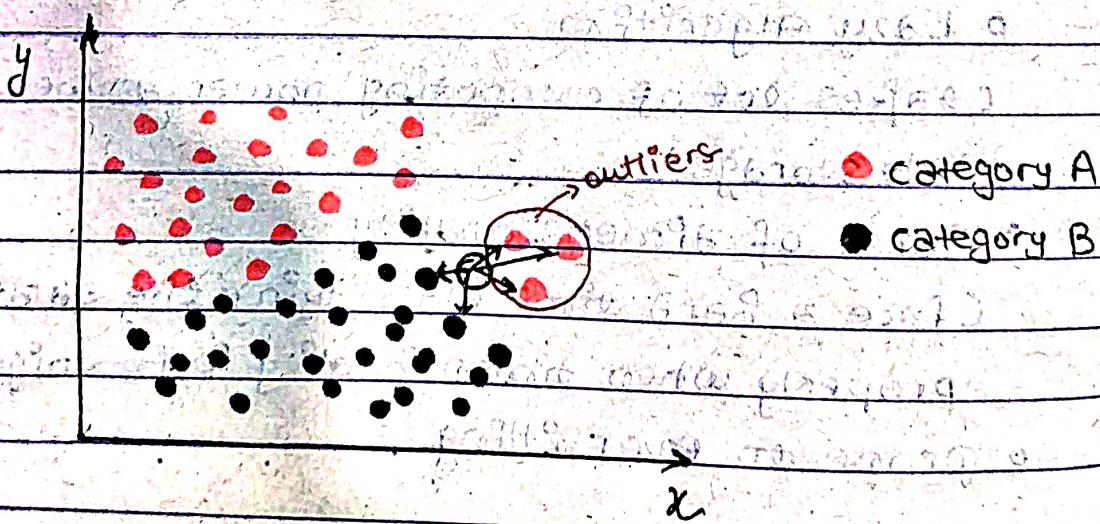
- prone to overfitting

o Impact of imbalanced dataset



- the model will be biased
- As a consequence, the bulk of the closest neighbours to this new point will be from "dominant class".
- we have to balance our dataset using "upsampling" or "downsampling" strategy

o Impact of outliers on dataset



- new datapoint should belong to category B, but due to outliers it will have category A