# K-Nearest Neighbour (KNN)
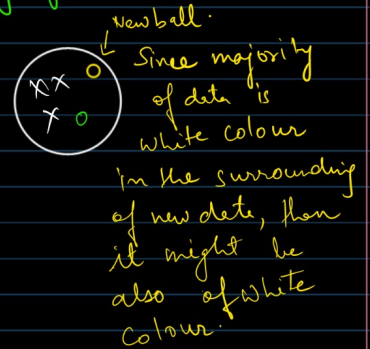
① KNN Classifier
② KNN regressor

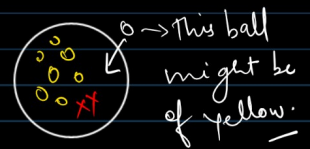You will be like people surrounding you.

Scenario 1

New ball.
Since majority of data is white colour in the surrounding of new data, then it might be also of white colour.

In both of the cases you have considered the majority colour ball in the bag as deciding factor.

Scenario 2.

0 → This ball might be of yellow.
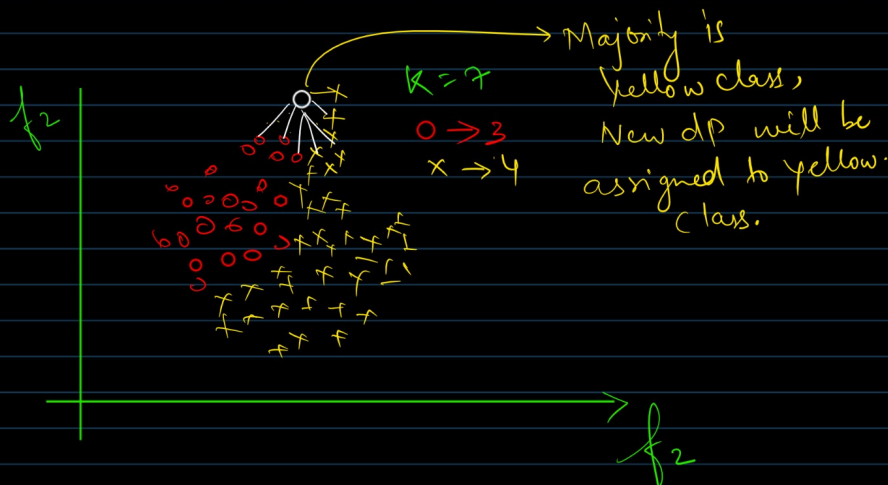
---

① KNN Classifier

K = 5 (How many nearest $dp$ you want to see)

→ This new $dp$ will be assigned to $0^{th}$ class as it is in majority.

x = 1
0 = 4

| $f_1$ | $f_2$ | $y$ |
|-------|-------|-----|
| — | — | 0 |
| — | — | 1 |
| — | — | 0 |
| — | — | 0 |
| — | — | 1 |

$f_2$

$f_2$

---

Majority is Yellow class, New $dp$ will be assigned to yellow. class.

K = 7

0 → 3
x → 4

$f_2$

$f_2$

* As k changes the class of new $dp$ might also change

⇓

K - hyper parameter.

# KNN Classifier Algorithm

① Plot the datapoint in n-d space.

② Initialise the k-value (No of Neighbours you want to Consider)

for each k, $k \in 1 \to \infty$ (Generally is taken greater than 3 and it should always be greater than 0)

→ Calculate the distance of new dp w·r·t to all datapoints.

→ Sort the distance

⇒ based on k find the class of that k nearest datapoint.

③ find mode of the class

④ Assign the class.



$f_2$

→ distance of new dp from all the dp's
→ Sort the distance.
→ Select k nearest distance dp.
→ take mode and assign the class.

training        $f_1$

newdp = (100, 50)

| $f_1$ | $f_2$ | y | distance of new dp with rest dp | sort | k=3 |
|---|---|---|---|---|---|
| | | | 40 | 38 (0) | Majority is 0 |
| | | 1 | 58 | 40 (1) | 1 for k=3, ∴ |
| | | 1 | 60 | 58 (1) | Prediction |
| | | 0 | 38 | 60 (1) | is class 1. |

## To Calculate distance

① Euclidean distance

$(x_2, y_2)$

$(x_1, y_1)$

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

A •——• B    by Aeroplane → Euclidean distance.

② Manhatta distance.

$(x_2, y_2)$

$(x_1, y_1)$   $d_1 \to$   $d_2$

$$distance = d_1 + d_2$$

## K - hyperparameter

→ for different k, you keep the track of train / test accuracy.

row of house

row of house

row of house

$P_1$ $P_2$

train accuracy

1 2 3 4 5 6 7 8 9 10 k

→ for the k which gives highest train accuracy, that will be optimal k.

## ② KNN Regression

price

No of rooms

**Summary of Classification**

→ calc dist
→ sort distance
→ select nearest k distance
→ Assign the mode to class

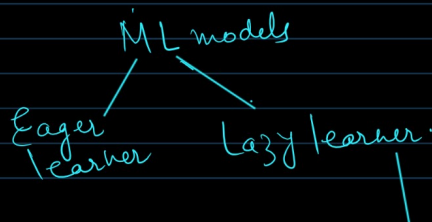→ Prediction of new dp will be average (y) price of (mean) house.

→ In case of outlier instead of Avg, you take median of k nearest dp.

**Advantage**

→ Easy to Understand / Very intuitive
→ Performance of model in terms of evaluation metric is good. (Rsq, acc)

**disadvantage**

ML models

Eager learner    Lazy learner

→ Lazy learner.

At real time
the distance of
test data from
each of dp is
Calculated $\Rightarrow$ Brute force.

$\rightarrow$ All the
model parameter
are calculated
While training
and used for prediction

$\rightarrow$ you are Calculating parameters
of model in real time.

training | computation time is more.

$\rightarrow$ n dp $\rightarrow$ n-1 distance you will be Calculating.
$\rightarrow$ Computationally expensive $\left( O(N) \right)$

* Variants of kNN

① kD tree $\Big\}$ $\rightarrow$ Optimizers
② Ball tree _____