# K-Nearest Neighbour

   ① KNN Classifier
   ② KNN Regressor.

## ✳ KNN Classifier
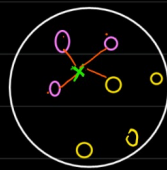
Seen-1

x (Assign a class)

Scn-1

* 4 nearest point

⇓

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

Seen-2

instead of
4 nearest dp,
take 7 nearest
dp.

* → Yellow.

Pink    Yellow

3       ⇓

       1

majority ⇒ Pink    * ⇒ Pink class

K → nearest dp.

K = 5

→ As k changes,
the class of new dp
might also change.

⇓

K is a hyperparameter.

4 nearest dp

| o | — | 1 |
|---|---|---|
| o | — | 2 |
| □ | — | 3 |
| o | — | 3 |
| o | — | 5 |
| o | — | 7 |
| o | — | 9 |

→ Prediction
will be

o

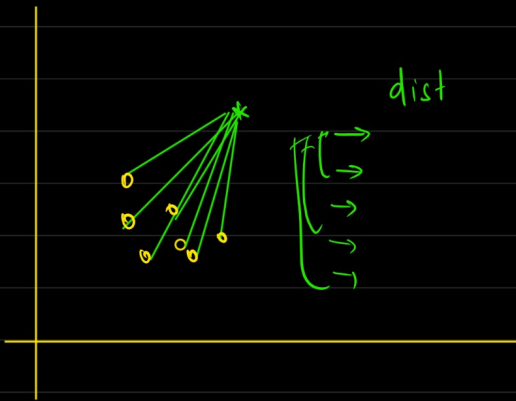binary class. $\longrightarrow$ K as odd value.

$\overset{l_s}{\underset{5, 7, 9, 11, 13}{}}$

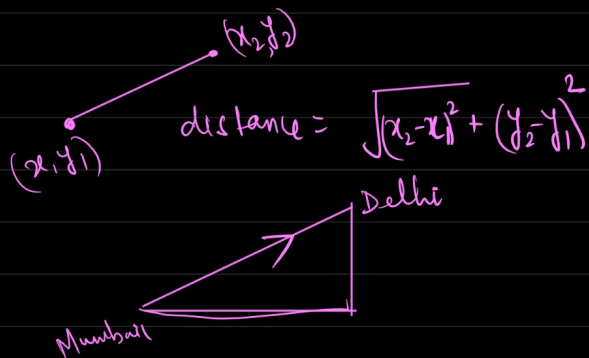2 class:

How K?? $\Longrightarrow$ K is hyperparameter.

## Algorithm

① Plot the dp in nd-space.

② Initialise the k-value.
(No of neighbours You want to consider)

$K \in 1 \longrightarrow \alpha$ (Generally k > 3)

$\rightarrow$ Calculate the distance of new dp w.r.t to all dp's.

$\rightarrow$ Sort the distance

$\rightarrow$ based on k find class of that k nearest dp's
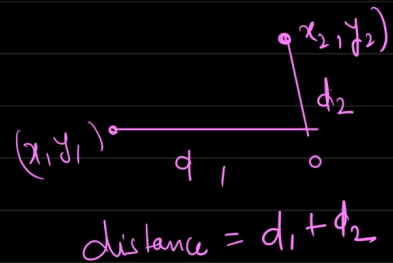
③ Find the mode of the class

④ Assign the class.

dist
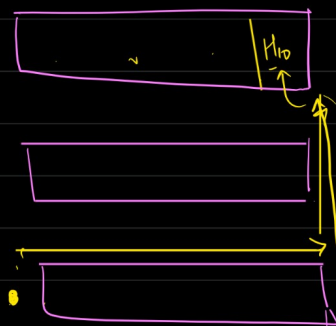


## distance

① Euclidean distance.

② Manhattan distance.



$(x_2, y_2)$

$(x_1, y_1)$

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

Delhi

Mumbai

$(x_2, y_2)$

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

$distance = d_1 + d_2$

$d(x, y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \cdots \cdots (x_n - y_n)^2}$

## How to decide K.



training accuracy (y-axis), K (x-axis with values 1 2 3 4 5 6 7 8 9 10)

---

\* **KNN Regressor**



price (y-axis), No. of rooms (x-axis)

→ Prediction for new dp will be avg (y (price) of nearest k dp.
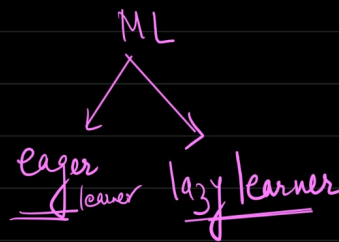
---

\* Advantage of KNN

→ Easy to Understand | very Intuitive.

→ Performance of model in terms evaluatic metric is good

---

\* disadvantage

→ Lazy learner.
   ‖
Computation time is more.

ML
 ↙        ↘
eager        lazy learner
 learner
                    → you are
• fit → Coeff's      calculating
• predict.           Parameter
                     of model
$\beta_0 + \beta_1 x_1 + \beta_2 x_2$   in real

\* All the model           time.
Parameter are calculated      |
while training & used      KNN.
for prediction.

Some more lazy learners

1. Locally weighted learning (LWL)
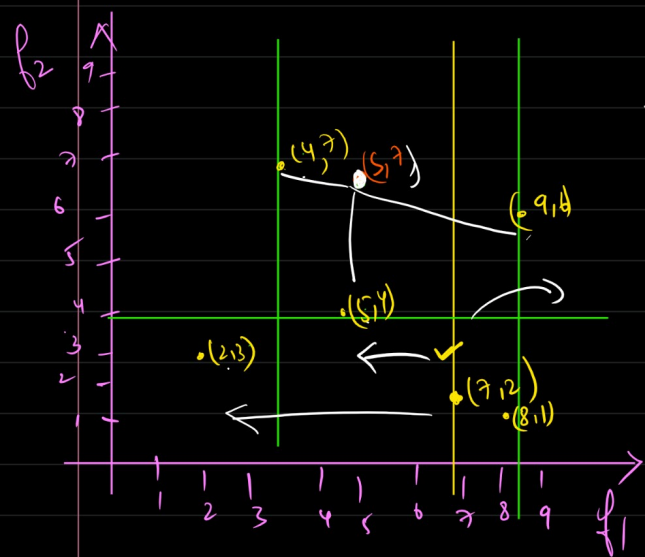2. Case-Based Reasoning
3. Lazy Bayesian rules. ✓

✦ Variants of KNN.

1. K-D tree
2. Ball tree.

why KNN is lazy learner ?

⟶ In real time the distance of test data from each of dp is calculated ⟹ Brute Force

* Variants of KNN

| $f_1$ | $f_2$ |
|-------|-------|
| 7 | 2 |
| 5 | 4 |
| 9 | 6 |
| 2 | 3 |
| 4 | 7 |
| 8 | 1 |

→ $f_1$ - 2, 4, 5, 7, 8, 9.
→ median → 7 as median

*
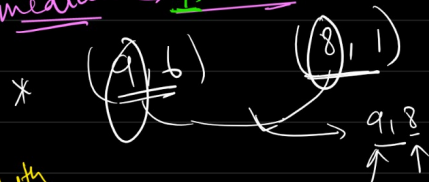
* K-D-tree.
  ↳ Partition the data in a binary tree.

Step-1 sort the feature $f_1$ & $f_2$
Step-2 Calculate the median of $f_1$ and $f_2$.
Step-3 Partition the data based on median
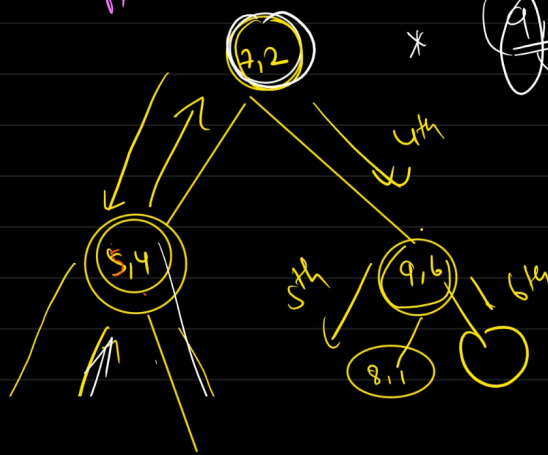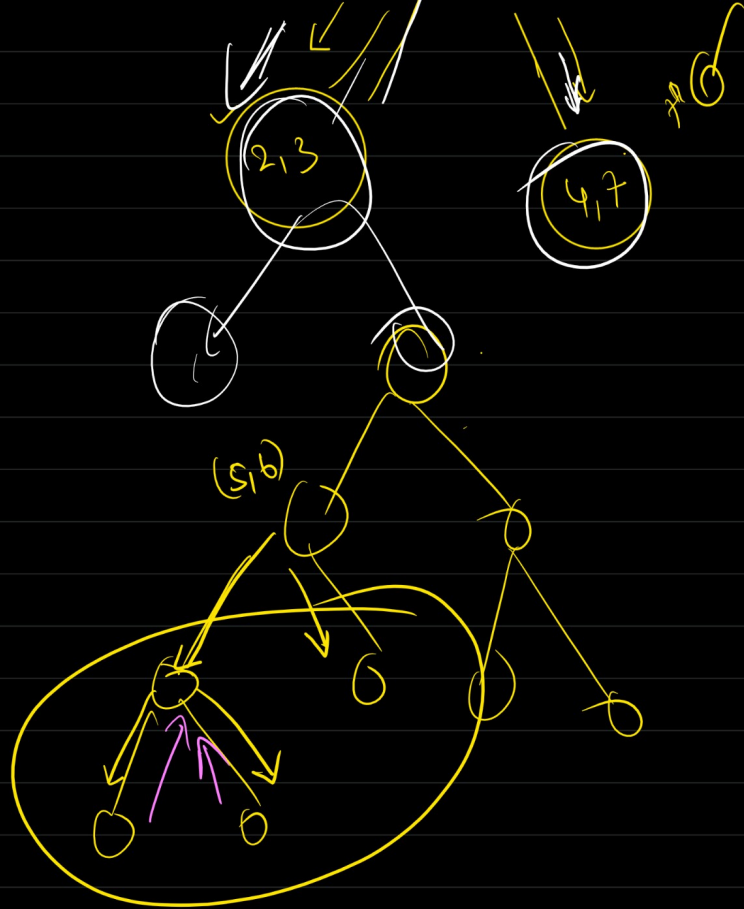Step-4 Partition all the data recursively in each of Subpartion.

Until In one partion only one dp is there

(2,3)

(4,7)

(5,6)

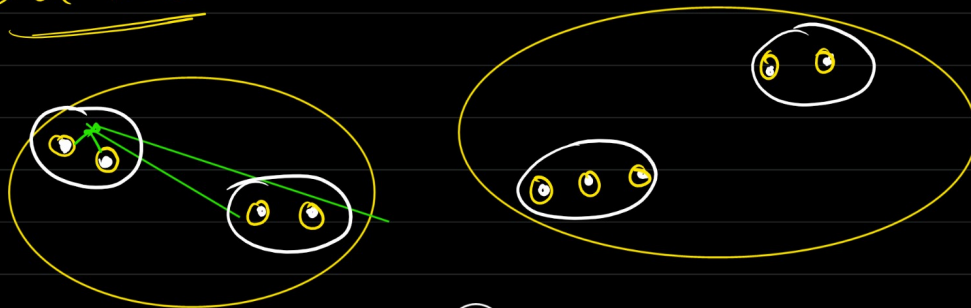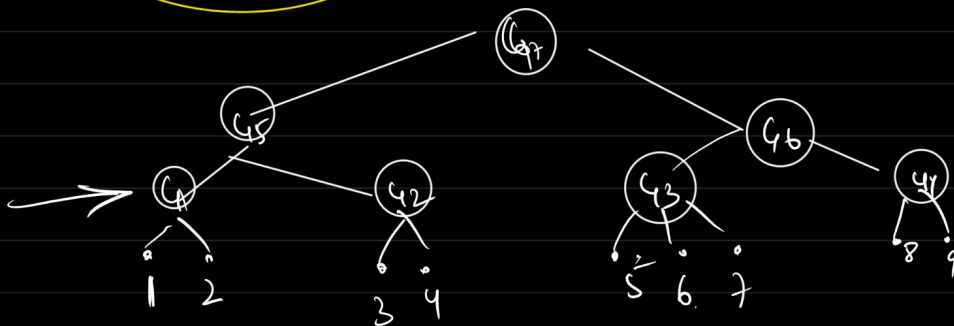BST

Binary $\Rightarrow O(\log n)$

worst $= O(m) \Rightarrow$ Brute forse

* Conclusion $\rightarrow$ reduced the search space.

Advatage

$\rightarrow$ Only need to calculate distance of the nearse group element.

② Ball tree

C7

C5                    C6

C1        C2      C3        C4

1  2      3  4    5  6  7    8  9

if Small dataset $\rightarrow$ Use kNN

→ read all csv in one go

☆ .csv

Concatenate