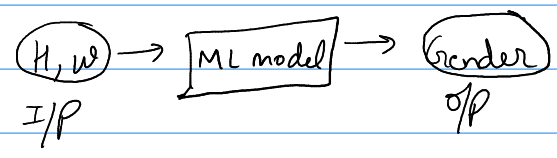# K - Nearest Neighbours (KNN)

↳ Type: Supervised Learning Algorithm → finding relationship b/w I/P & O/P

Task: Classification & Regression

eg. Problem Statement → Given a datapoint of Height & Weight, predict the Gender.

$(H, W)$ → [ML model] → (Gender)

I/P                O/P

| H | W | G |
|-----|-----|---|
| 150 | 45 | F |
| 180 | 80 | M |
| 145 | 50 | F |
| 190 | 75 | M |

Input → H, W

Output → G

Classification or Regression

① Distance based Approach → KNN

② Boundary   "   " → linear Reg, logistic Reg, SVM
                          (Reg).     (class)    (Both)

③ Rule   "   " → Decision Tree

④ Probabilistic "   " → Naive Bayes
                                   (class)

⑤ Ensemble "   "

⑥ Deep Learning "   "

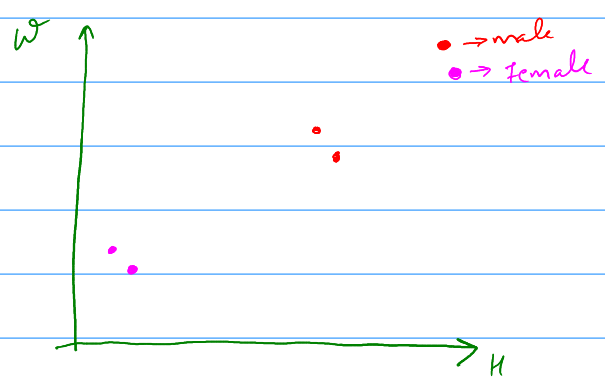Q. Is the data available? → yes

Q. Does the data have 'y'? → yes

Now Algo which can be used ⌐

Evaluation Metrics

↳ Accuracy, precision, F1 Score, Confusion matrix. etc.

| H | W | G |
|-----|-----|---|
| 150 | 45 | F |
| 180 | 80 | M |
| 145 | 50 | F |
| 190 | 75 | M |

Collection of rows / datapoints

(Dataframe) → matrix

Visualizatⁿ →

W ↑

• → male
• → Female

→ H

Now, given a query point, predict the gender
↓
(H, w)

| H | w | G |
|---|---|---|
| 150 | 45 | F |
| 180 | 80 | M |
| 145 | 50 | F |
| 190 | 75 | M |

→ Visualizat^ →



• → male
• → female

$q_1 \to (H_1, W_1) \to$ Female
$q_2 \to (H_2, W_2) \to$ Male

Basic KNN Steps

Given a query point
↳ find its nearest neighbour.
↳ Check the target values of the nearest neighbours.
↳ (Count the no. of male & females)
↳ Majority class wins and query point becomes that class

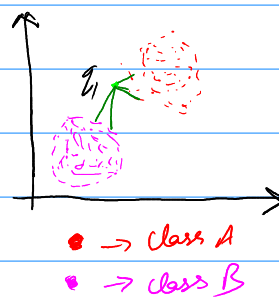Given a query point & the historical data

(pre a) Decide the value of K ← Done using hyperparameter tuning.

* Sklearn by default takes k=5

↳ K decides how many nearest neighbours to look at.

K=3 [odd]
↓
only for binary classification
doesn't work for multi class classificat^

(a) Find the nearest neighbours
↳ K nearest Neighbours.



• → Class A
• → Class B

(b) Voting (Counting)
↳ Count the nearest neighbours belonging to each class.

(c) Majority Vote wins → decide the class of query point based on majority voting

**#** Algorithm (K-NN) ↲

① Decide K
   ↑ no. of nearest values / nearest neighbours.

② distances = [ ]
   ↑ To store all the distances of query point from data set points.

③ KNN pts = [ ]

find's
K
neighbours ④ for each datapoint in distances :
   ↳ Compute the distance b/w query point & datapoint.
      ↳ Euclidean (by default, but can be changed as well)
   → dist = Distance (query point, datapoint).

   → distances . append ( (datapoint, dist))

   ⇒ distances = $\left[ (dP1, 100), (dP2, 5), - - - - \right]$

⑤ Sort distances:
   distances = $\left[ (dP2, 5), (dp5, 6) - - - \right]$

⑥ KNN-pts = distances [ : K]

⑦ Voting :    count_male = 0
              count_female = 0
              for pt in KNN-pts :
                 ↳ if pt is male:
                       count male += 1
                 else
                    count-female += 1              Voting

⑧   If   count-male > count-female > 0
                print('male')

Else
                print('female')

Q   What is the range of K?
                K → 1 to n

The reason for KNN having such a fast training time is because it just memorizes the data during training time and there is no learning. All the work happens when a query point comes thats when it calculates the dist., sorts them and finds K nearest neighbours, as a result it has a large prediction time

Q   Why not even K?
         ↳ In case of binary classification, there can be a situation of equal voting as a result we fail to classify the query point
         ↳ In case of multi-class classification, the K value doesn't matter as
                        in either case we can get equal vote situation.

A→0 ⎫ n ⎫
B→2 ⎬   ⎬ → 0
C→2 ⎭

A→1 ⎫ n?
B→1 ⎬
C→1 ⎭

Q   What happens if 'k' is too small or too large?
    Range of K = 1 to n          ⌐K=1              K=n ⌐
                        Can cause overfitting      Can cause underfitting

Q   How to decide the value of K?
       ↳ Hyperparameter tuning.           [ It is a
                                          (hyperparameter
                                          of KNN)
                                                   ↓
⁎ KNN is a non-parametric ML algorithm    There
                                          is dist.
⁎ KNN is aka Lazy Learner.                metric
                                          also
⁎ KNN has a very high time & space complexity.
         ↳ KNN is not used in production.

X-train ⎫ → KNN → model.
y-train ⎭

                              Overfit  │  Underfit

                                  Best fit

⁎ Best Alternative for KNN
                        ⎧ ↳ ANN (Approximate Nearest Neighbours) ⎫
                        ⎨                                          ⎬
                        ⎩ ↳ K-d tree Data Structure               ⎭
                                ↳ Improve the efficiency of KNN

⁎ for KNN → Training Phase → very less time, high space complexity
aka lazy learner ← no training at all
                                                              → for each query point
                                                                all dist. in df are
                                                                calculated
                                                                        +
              → Testing Phase → almost all steps happen here ⟹ high prediction time  Sorting
                                                                                     is very time
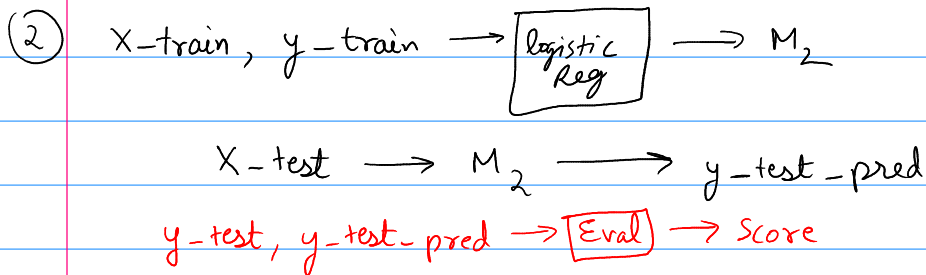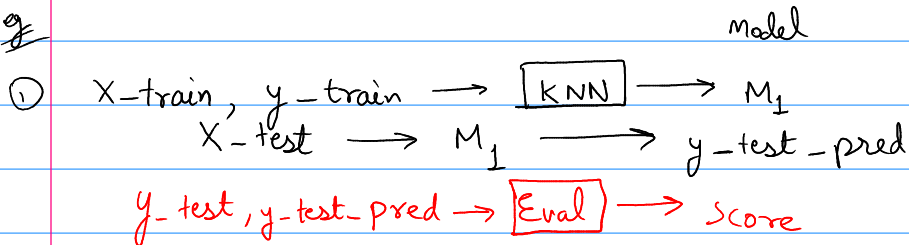                                                                                     consuming

Selection Criteria for Production ML model
  ↳ Model train time
      ↳ model prediction time
          ↳ Model Size

Model

① X-train, y-train ⟶ [KNN] ⟶ $M_1$
     X-test ⟶ $M_1$ ⟶ y-test-pred

   y-test, y-test-pred ⟶ [Eval] ⟶ Score

② X-train, y-train ⟶ [logistic Reg] ⟶ $M_2$

     X-test ⟶ $M_2$ ⟶ y-test-pred

   y-test, y-test-pred ⟶ [Eval] ⟶ Score

Now for the same data both the models give a score of 98%. then how to choose?

|  | $M_1$ | $M_2$ |
|---|---|---|
| Train time | 0.000001 sec | 5 hr |
| Test time | 1 hr. | 0.000001 sec |
| Size | 100 GB | 1 GB |
| Cost per pred. | Rs 1 | Rs 1000 |

Now to choose a model you need to understand the targeted user base behaviour:
  ↳ A daily internet user would expect faster output and the model would have
     to be light weight. ⟹ M2 will be chosen
  ↳ If the user base being targeted don't care about the execution time but
     want less cost per prediction
        ⟹ M1 will be chosen.

Whenever ML models are deployed, there will always be tradeoffs but we need
to select the model based on tradeoffs that prove to be the best fit/ beneficial to
us.

# $ KNN Regression $

$$D_n = \{(x_i, y_i) \mid x_i \in R^{d-1}, y_i \in R\}$$

| Classification | Regression |
|---|---|
| **Problem Statement:** Given the height & weight of an individual predict the gender of the individual | **Problem Statement:** Given the height of an individual, predict the weight of an individual. |
| **Target (y)** → Gender | **Target (y)** → Weight |
| **Input** → Height & Weight | **Input** → Height |

**Data →**

| Height | Weight | Gender |
|---|---|---|

$n \times d$

X — y

$$D_n = \{(x_i, y_i) \mid x_i \in R^{d-1}, y_i \in \{G, M\}\}$$

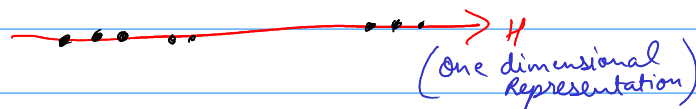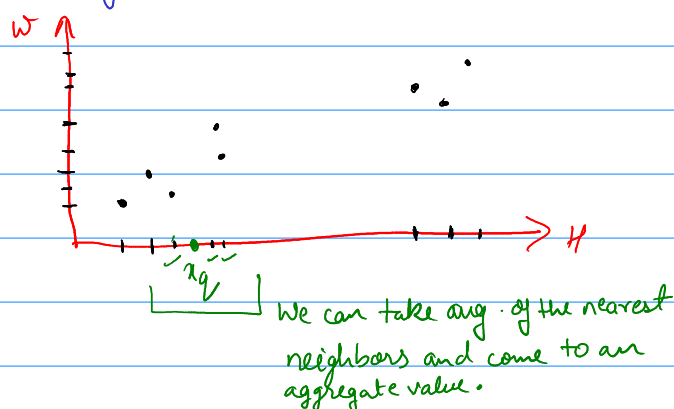Target → Discrete
↓
Classification

**Data →**

| Height | Weight |
|---|---|

$n \times d$

X — y

$$D_n = \{(x_i, y_i) \mid x_i \in R^{d-1}, y_i \in R\}$$
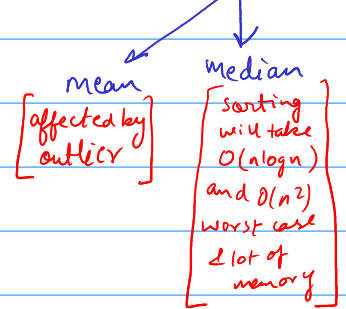
Target → Real Value /
Continuous value
↓
Regression

→ Female
→ Male

(One dimensional Representation)

Now, we try to visualize it in 2 dimensions

**Steps:**

(a) decide the value of K.

(b) Compute the distance of $x_q$ from each datapoint in the dataset

(c) Get K-nearest neighbours → Sort & slice

$x_q$

We can take avg. of the nearest neighbors and come to an aggregate value.

ⓓ Voting → count the number of datapoints for each class. [Take the mode of the classes]

ⓔ Class with majority wins.

→ Basically computing the mode

In the above example, we find the 3 nearest neighbor and take the aggregate value (Avg. Value).

Here we can use two ways to compute the target.

mean → [affected by outlier]

median → [sorting will take $O(n\log n)$ and $O(n^2)$ worst case & lot of memory]

Steps:

ⓐ Decide the value of K.

ⓑ Compute the distance of $x_q$ from each datapoint in the dataset.

ⓒ Get the K nearest neighbors → Sort & Slice.

ⓓ Compute the avg. from K-nearest neighbors.
 → Basically computing the mean / median.

♯ When implementing KNN in sklearn we need to specify → K
 → Distance to be used [Euclidean, manhattan, minkowski etc.]

Parameters Same for Regressor also.

Value of K

♯

class sklearn.neighbors.**KNeighborsClassifier**(n_neighbors=5, *, weights='uniform', algorithm='auto', leaf_size=30, p=2, metric='minkowski', metric_params=None, n_jobs=None)          [source]

for $p=2$, it by default computes Euclidean distances

assign importance to points by default everyone is equal.

Kd tree & other optimizing algo.

- - - - - - - - - - - - - - - - - - - - - - - - - - - -

♯ Hyperparameters for KNN

Their best value can be found using hyperparameter tuning.

{ ① value of K
② Distance metric by selecting the value of 'p'.

♯ Alternate to KNN
→ ANN (Approximate Nearest Neighbor)
→ K-d trees.