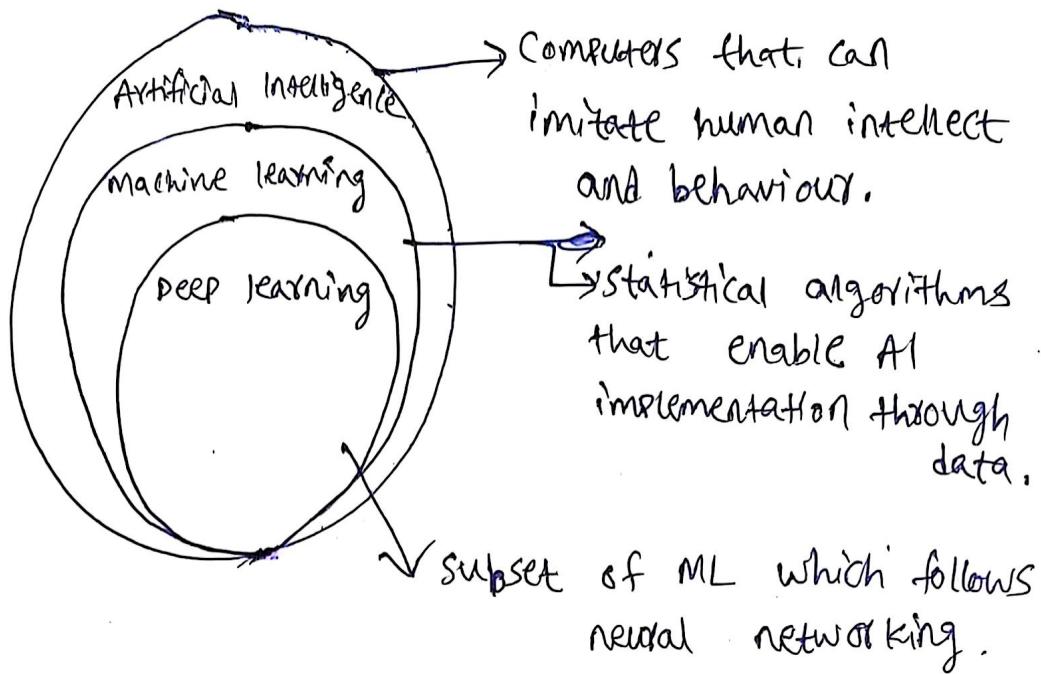


# AI VS ML VS DL



## AI (Artificial intelligence)

- AI is basically the mechanism to incorporate human intelligence into machines through a set of rules, (algorithms).
- AI is the broader family consisting of ML and DL as its components.
- The efficiency of AI is basically the efficiency provided by ML and DL respectively,
- Examples - Amazon's Alexa, Siri

## ML (Machine learning)

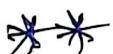
- ML is the study that uses statistical methods enabling machines to improve with experience.
- ML is subset of AI.
- Less efficient than DL as it can't work for longer dimensions or higher amount of data.
- Examples → Car price prediction, Diabetes prediction.

## DL (Deep Learning)

- DL makes use of neural networks (similar to human neurons) to imitate functionality just like a human brain.
- DL is the subset of ML.
- More powerful than ML as it can easily work for larger set of data.
- It attains the highest rank in terms of accuracy when it is trained with large amount of data.
- Examples - Image classification, Image de-noising etc.

### Note

\* \* DALL-E and DALL-E 2 are machine learning models developed by openAI to generate digital images from natural language descriptions.



Variables = features = dimensions = Attributes = Column

Records = Rows

## SUPERVISED Learning !

- It is the learning of the model where with input variable (say,  $x$ ) and an output variable (say,  $y$ ) and an algorithm to map the input to the output i.e.  $Y=f(x)$ .
- The basic aim is to approximate the mapping function (Mentioned above) that when there is new input data( $x$ ) then the corresponding o/p variable can be predicted.

- It is called supervised learning because the process of learning (from the training dataset) can be thought of as a supervisor who is supervising the entire learning process. Thus, the "learning algorithm" iteratively makes predictions on the training data and is corrected by the "supervisor", and the learning stops when the algorithm achieves an acceptable level of performance
- It is called supervised because data is labelled ( $y$  column is present)
- Examples of supervised learning algorithms - Regression and Classification

$$y = f(x)$$

$y$  = dependent variable

$x$  = independent variable

\* If there are two categories in target variable, then it is binary classification.

### Terminologies:

- 1) Training Data — The data onto which, we will train our ML model
- 2) Testing Data — The data onto which the ML model (already trained on training data) will be used to generate predictions.

### Prediction Pipeline:

- 1) Load the data
- 2) Feature Engineering — data cleansing and domain understanding  
(Removal/Filtering of null values)

b) Removal of duplicates (rows or columns)

c) Encoding categorical variables

d) Standardising or scaling variables.

Age (10-99)

Income - 100, 1000, 10000, ...

e) Creating custom features out of existing variables

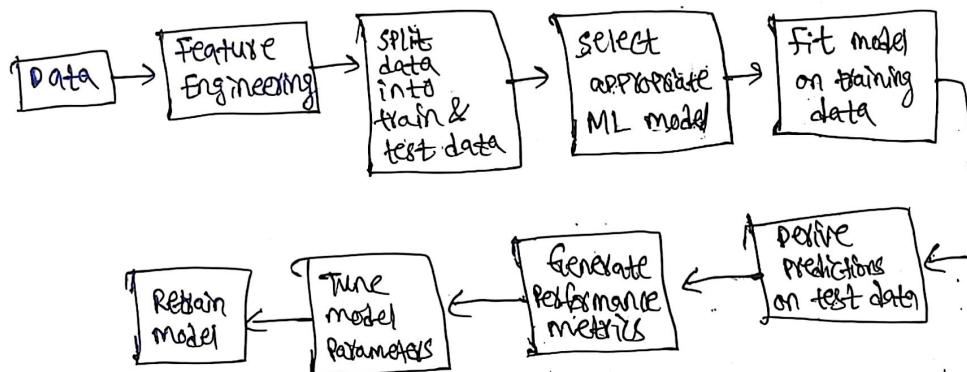
Date (2019-03-28)  $\Rightarrow$  from date, we can create custom features Year, Quarter, Month, Day

f) Binning

Salary - 12K, 15K, 20K, 18K, 25K, ...

Salary bracket - Low (10K-24K), Med (25K-49K), High (50K-70K)

## Prediction Pipeline



## Unsupervised Learning

$\rightarrow$  Unsupervised learning is where only the input data (say,  $X$ ) is present and no corresponding output variable is present.

$\rightarrow$  The main aim of Unsupervised learning is to model the distribution in the data in order to learn more about the data.

$\rightarrow$  It is called so, because there is no correct answer and there is no such supervisor (unlike supervised learning methods)

are left to their own devices to find the interesting structure in the data.

$\rightarrow$  It is unsupervised as data

$\rightarrow$  Examples - Clustering and Dimensionality Reduction

### Clustering Example



$(X_1, X_2)$

After clustering

~~$X_1$~~

$Y_1$

\* After performing becomes suitable

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt.
```

### DataFrame:

$\rightarrow$  Data represented in rows.

$\rightarrow$  DataFrame consists of index

$\rightarrow$  A row can be an index of a

### Series:

$\rightarrow$  Single column data is Series

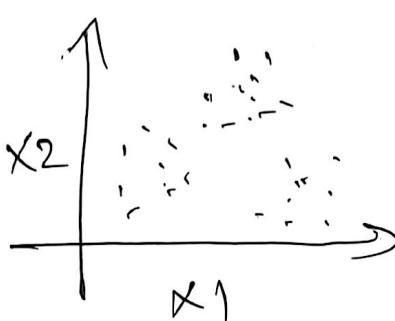
$\rightarrow$  Series also has index

are left to their own devices to discover and present the interesting structure in the data.

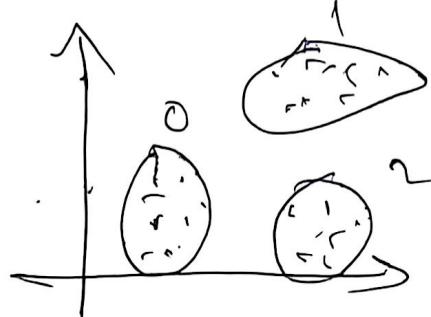
→ It is unsupervised as data is not labelled, ( $y$  column is present).

→ Examples - Clustering and Dimensionality Reduction.

### Clustering example:



$\xrightarrow{\text{K-means}}$



$X(x_1, x_2)$

After clustering

~~$X$~~   $X: [x_1, x_2]$

$Y \in \{0, 1, 2\}$

\* After performing clustering, the data becomes supervised for classification

import numpy as np

ii Pandas as pd

ii matplotlib.pyplot as plt.

### Dataframe:

→ Data represented in rows and columns.

→ Dataframe consists of index as well.

→ A II can be II of a single column.

### Series:

→ Single column data is series

→ Series also bear index as well.

P) `data = {'Name': ['Ankit', 'Shubham', 'Mohit', 'Shreya'], 'Age': [23, 26, 34, 20], 'City': ['Delhi', 'Bangalore', 'Delhi', 'Pune']}`

`df = pd.DataFrame(data)`

`df`

O/P:-

	Name	Age	City
0	Ankit	23	Delhi
1	Shubham	26	Bangalore
2	Mohit	34	Delhi
3	Shreya	20	Pune

P) `df.columns`

O/P `Index(['Name', 'Age', 'City'], dtype='object')`

P) `df.values`

O/P `([[{'Name': 'Ankit', 'Age': 23, 'City': 'Delhi'}, {'Name': 'Shubham', 'Age': 26, 'City': 'Bangalore'}, {'Name': 'Mohit', 'Age': 34, 'City': 'Delhi'}, {'Name': 'Shreya', 'Age': 20, 'City': 'Pune}]] , dtype='object')`

P) `df.shape`

O/P `(4, 3)`

P) `df.sample()` → if we pass any number here, that number of random rows will displayed,

O/P `It returns any random row`

Q) print(type(df))

O/P:-

`<class 'pandas.core.frame.`

`Dataframe'>`

Q) Accessing single column

df['Name']

O/P:-

0 ANKIT

1 shubham

2 Mohit

3 Shreya

Name: Name, dtype: object

Q) print type(df['name'])

O/P:- `<class 'pandas.core.series.Series'>`

# Selecting multiple columns

Q) df[['name', 'city']]

O/P:-

	Name	City
0	ANKIT	Delhi
1		
2		
3	Shreya	Pune

## Regression:

→ It is a predictive modeling technique which investigates the relationship b/w dependent & independent variable (one or more)

→ Dependent variable is continuous in nature, eg - sales, weight, profit, revenue, price etc --

⇒

→  $y$  = dependent variable / output  
 $-x$  = independent variable / input

## Simple Linear Regression :-

→ It is a regression model that estimates the relationship between one independent variable and one dependent variable using a straight line.

→  $y = ax + b$  or  $y = mx + c$

Where

$x$  = independent variable / input feature

'input attribut'  
" " column

$y$  = dependent variable / output feature

$a/m$  -----  $b$  after  
= slope or coefficient or weight  
" " column

$m$  = slope or coefficient or weight  
or how much we expect  $y$  to  
change as  $x$  changes.

$b/c = \text{intercept/constant/bias}$

$$m = a = \frac{(y_2 - y_1)}{(x_2 - x_1)} = \frac{dy}{dx}$$

Total variation is made up of two parts!

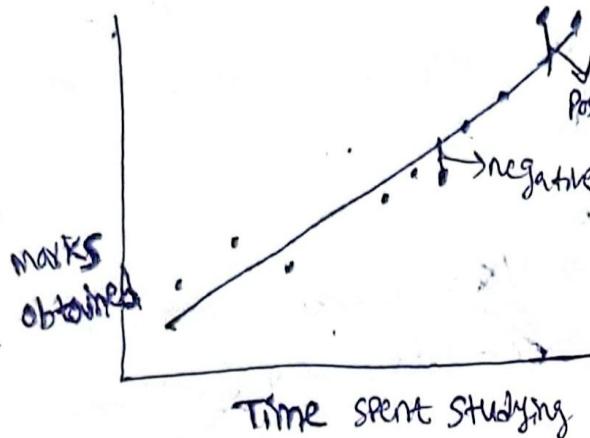
$$SST = SSE + SSR$$

$$SST = \sum (x_i - \bar{Y})^2$$

$$SSR = \sum (\hat{Y}_i - \bar{Y})^2$$

$$SSE = \sum (y_i - \hat{Y}_i)^2$$

Line of best fit:



Where

$\bar{Y}$  = Average value of dependent variable

$y_i$  = observed values of dependent variables.

$\hat{Y}_i$  = predicted value of  $y$  for the given  $x_i$  value.

R<sup>2</sup> score (Coefficient of determination or goodness of fit):

→ The coefficient of determination is the portion of the total variation in the dependent variable that is explained by variation in the independent variable.

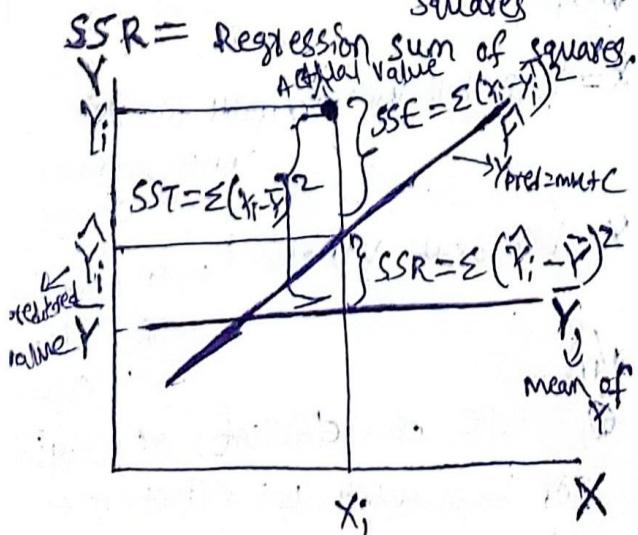
→ This is used to evaluate the performance of the model

→ The coefficient of determination is also called r-squared and is denoted as  $r^2$ .

$$r^2 = \frac{SSR}{SST} = \frac{\text{regression sum of squares}}{\text{total sum of squares}}$$

Note:  $0 \leq r^2 \leq 1$

$$\rightarrow SST = SSE + SSR$$



$$R^2 \text{ score} = \frac{SSR}{SST} = \frac{(SST - SSE)}{SST} = 1 - \frac{SSE}{SST}$$

when  $SSE=0$ ,  $R^2$  score = 1 (Best case scenario) → It has an eqn of the form

{when  $SSE=SST$ ,  $R^2$  score = 0 (Worst case scenario)}

$$Y = m_1x_1 + m_2x_2 + \dots + m_nx_n + b$$

where

$x_1, x_2, \dots, x_n$  = independent variables / input features.

## Regression Metrics:-

### Mean Absolute Error (MAE)

It is the mean of the absolute value of the errors.

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$y$  = dependent variable / output feature

$m_1, m_2, \dots, m_n$  = coefficients / slope corresponding to  $x_1, x_2, \dots, x_n$

$b$  = intercept / constant / bias

→ It is practically out of scope to represent multiple linear regression on a scatter chart with actual data points and a regression line as the scatter chart would have to span multiple dimensions corresponding to each independent variable and likewise the regression line would also span across multiple dimensions.

### Root Mean Squared Error (RMSE)

It is the square root of the squared errors.

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} = \sqrt{MSE} = \sqrt{\frac{SSE}{n}}$$

This paragraph means that if we have 5 dimensions, we can't plot them in 2-dimensional graph.

## MULTIPLE LINEAR REGRESSION

→ MLR is used to estimate relationship b/w two or more independent variables and one dependent variable.

### Pros:-

→ Linear Regression performs well when the dataset is linearly dependent.

→ Linear Regression is easy to implement, interpret and

efficient to train.

### Cons :-

- Main limitation of Linear Regression is the assumption of linearity between the dependent variable and the independent variables. In the real world, the data may not always be linearly dependent.

### For every regression Model

- Model only accepts int or float as input for X and Y
- Model will not accept any null value.
- X has to be either a Dataframe or a 2D numpy array or a list of list.
- Y has to be either a Series or a 1D numpy array or a list.

### Dataframe :-

- a) In this, data is represented in rows and columns with index. Generally Dataframe consists of multiple column data.
- b) Dataframe can be of single column with index, however

### Series

- Single column data with index.
- Series cannot contain multiple column, however it can contain multiple indexes.

### Simple Linear Regression Implementation

```
import pandas as pd  
u numpy as np  
u matplotlib.pyplot as plt.
```

#### 1) Load the data

```
d = {'Area(sf)': [2600, 3000, 3200,  
                  3600, 4000],  
     'Price': [550000, 565000, 610000,  
               680000, 725000]}
```

```
print(d)
```

```
print(type(d))
```

o/p :-

```
{'Area(sf)': [- - -], 'Price': [- - -]  
<class 'dict'>}
```

#### 2) df = pd.DataFrame(d)

```
df
```

o/p :-

	Area(sf)	Price
0	2600	-
1	-	-
2	-	-
3	-	-
4	-	-

Print statement - Based on Area(sf) → print(x.shape)  
predict the price when print(y.shape)

a) Area = 1400 sqft

b) Area = 3450 sqft.

OP:

(5,1)

(5, )

P) ~~df['price']~~

OP:  
0 550000  
1  
2  
3  
4

Name: price, dtype: int64

P) single column data frame.

df[['price']]

type(df[['price']])

OP:  
0 price  
0 550000  
1  
2  
3  
4

P) ~~df.shape~~

OP: (5,2)

P)  $x = df[['Area(sf)']]$

$y = df['Price']$

print(type(x))

print(type(y))

OP:

<class 'pandas.core.frame.DataFrame'>

<class 'pandas.core.series.Series'>

P) plt.scatter(df['Area(sf)'],

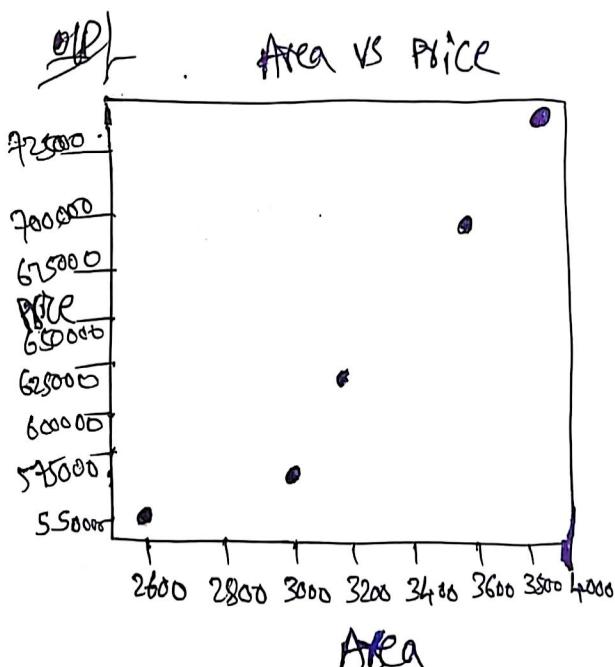
df['Price'], color='magenta')

plt.xlabel('Area')

plt.ylabel('Price')

plt.title('Area vs Price')

plt.show()



lets create linear regression model.

P) from sklearn.linear\_model import LinearRegression

m1 = LinearRegression()

m1.fit(x,y)

#model.fit() - It trains the ML model on the training data.

# model score() gives R<sup>2</sup> value P)

for Regression.

plt.scatter(df['Area(sqft)'], df['Price'])

color = 'maroon', label = 'Actual data points'

label = 'Predicted line')

P) print (model score, m1.score(x,y))

O/P:-

Model score 0.958430118199486

plt.xlabel('Area')

plt.ylabel('Price')

plt.title('Area vs Price')

plt.legend()

plt.show()

P) ypred\_m1 = m1.predict(x)

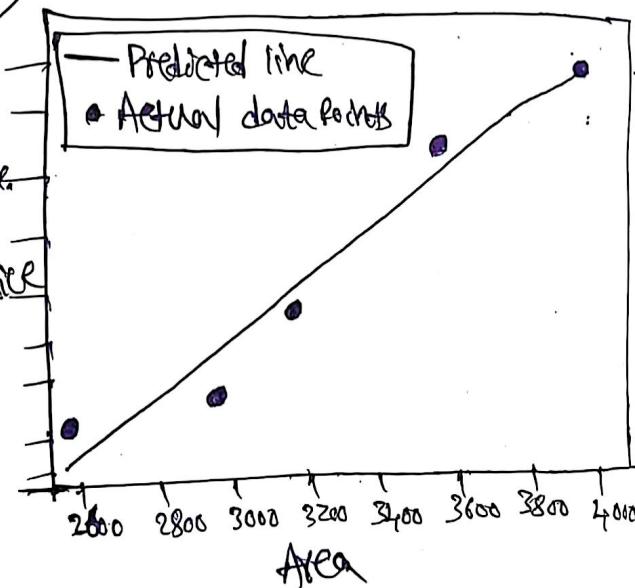
ypred\_m1

O/P:-

array([533664.38356164, 587979.452227, ..., ..., ...])

O/P:-

Area vs Price.



O/P:-

Area(sqft) Price Pred\_Price

	Area(sqft)	Price	Pred_Price
0	2600	580000	533664.383562
1	2700	580000	533664.383562
2	2800	580000	533664.383562
3	2900	580000	533664.383562
4	3000	580000	533664.383562

P) # YPred = mx + c.

m = m1.coef\_

c = m1.intercept\_

print ('slope', m)

print ('Intercept', c)

O/P:-

slope [35.78767028]

Intercept 180616.13835616432

p)  $\text{df}[\text{pred\_price\_eqn}] = \text{df}[\text{Area}(\text{sqft})] * m + c$ .

`df.head()`

~~df~~

	Area (sqft)	PRICE	PRED_PRICE	PRED_PRICE_EQN
0	2600	550000	533664.383562	533664.383562
1	3000	:	587979.452055	587979.452055
2	3200	:	:	:
3	:	:	:	:
4	:	:	723767.123288	723767.123288

p) from sklearn.metrics import mean\_squared\_error, mean\_absolute\_error, r2\_score

def eval\_metrics(ytest, ypred):

mae = mean\_absolute\_error(ytest, ypred)

mse = mean\_squared\_error(ytest, ypred)

rmse = np.sqrt(mean\_squared\_error(ytest, ypred))

r2s = r2\_score(ytest, ypred)

print('MAE', mae)

11 ('MSE', mse)

11 ('RMSE', rmse)

11 ('R2-score', r2s)

p) eval\_metric(y, ypred\_m1)

MAE 11246.57534246575

MSE 186815068.493509

RMSE 13668.030893042014

R2-Score 0.9584361138199486.

p) # a) Area = 1400 sqft.

# b) Area = 3450 sqft.

$y_{pred} - 1400 - a1 = m1 \cdot \text{Predict}$  (1400)

$y_{pred} - 1400 - a2 = 1400 \cdot m + c$

Print(ypred - 1400 - a1)

Print(ypred - 1400 - a2)

~~df~~

[778682.19178082]

[770087.10170007]

## Multiple Linear Regression

P) df1 = pd.read\_csv('houseprice2.csv')  
df1.head()

O/P:-

	area	bedrooms	age	price
0	2600	3.0	20	580000
1	3000	4.0	15	650000
2	3200	NaN	18	
3	3600	3.0	30	
4	4000	5.0	8	

To check if null values is present or not.

P) df1.isnull().sum()

O/P:-

area	0
bedrooms	1
age	0
price	0
dtype?	int64

P) print(df1['bedrooms'].mean()),  
print(df1['bedrooms'].median())

O/P:-

4.2
4.0

Filling the null values.

P) df1['bedrooms'].fillna(df1['bedrooms'].median(), inplace=True)  
df1.isnull().sum()

O/P:-

area	0
bedrooms	0
age	0
price	0

## Confusion Matrix:-

1) It is a performance metric for classification statements.

2) It is a square matrix made up of 4 terms - TP, TN, FP, FN,

3) If target variable has n categories then shape of the confusion matrix will be (n,n)

## Terminologies

### 1) TP (True Positive)

Actual value is positive, ML model also predicted a positive value.

### 2) FN (False Negative)

Actual value is positive, ML model predicted a negative value

### 3) FP (False Positive)

Actual value is negative, ML model predicted a positive value

### 4) TN (True Negative)

Actual value is negative, ML model also predicted a negative value.

## Note

- sum of all the actual positive values/cases =  $TP + FN$
  - sum of all the actual negative values/cases =  $FP + TN$
  - sum of all the positively predicted values/cases =  $TP + FP$
  - sum of all the negatively predicted values/cases =  $FN + TN$
- #  $x = [\text{Age}, \text{Gender}, \text{BMI}, \text{Body\_mass}, \text{Blood-Glucose\_Level}]$
- #  $y = \text{Diabetic (0) or Non-Diabetic (1)}$

## Classification Metrics:

1) Precision -  $\frac{TP}{(TP+FP)}, \frac{TN}{(TN+FN)}$

→ from all the positively predicted cases, how many are actually positive.

→ From all the negatively predicted cases, how many are actually negative.

2) Recall -  $\frac{TP}{(TP+FN)}, \frac{TN}{(TN+FP)}$

→ from all the actual positive cases, how many has the ML model predicted positive.

→ from all the actual negative cases, how many has the ML model predicted negative.

# Binary Classification Example

# 0 - Positive (True), 1 - Negative (False).

$y_{\text{true}} = [1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 1, 1, 0, 1, 0]$

$\begin{bmatrix} 0, 1 \\ \text{# Actual values} \end{bmatrix}$

$y_{\text{pred}} = [0, 1, 0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 1, 0]$

$\begin{bmatrix} 1, 0, 0 \\ \text{# Predicted values} \end{bmatrix}$

$\text{print}(\text{len}(y_{\text{true}}), \text{len}(y_{\text{pred}}))$

3) F1-Score -  $\frac{2 * \text{precision} * \text{recall}}{(\text{precision} + \text{recall})}$

It is the harmonic mean between precision and recall.

4) Accuracy -  $\frac{(TP+TN)}{(TP+FN+FP+TN)}$

→ Out of all the values, how many the model has predicted correctly.

O/P

18 18

A)  $\text{print}(y_{\text{true}})$   
 $\text{print}(y_{\text{pred}})$ .

O/P

$[1, 0, 1, 1, 0, 0, \dots, 0, 1]$

$[0, 1, 0, 1, 0, 0, \dots, 0, 0]$

P) from sklearn.metrics import confusion\_matrix, classification\_report  
 $\text{cm} = \text{confusion\_matrix}(y_{\text{true}}, y_{\text{pred}})$

$$\# AV = 0, PV = 0 \Rightarrow TP = 4$$

$$\# AV = 0, PV = 1 \Rightarrow FN = 5$$

$$\# AV = 1, PV = 0 \Rightarrow FP = 6$$

$$\# AV = 1, PV = 1 \Rightarrow TN = 3$$

Print(cm)

# [TP FN]

# [FP TN]

$$\begin{array}{|c|c|} \hline \text{O} & \begin{bmatrix} 4 & 5 \\ 6 & 3 \end{bmatrix} \\ \hline \end{array}$$

P) Print (classification\_report(y\_true, y\_pred))

	Precision	Recall	f1-score	Support
(corabetic) 0	0.40	0.44	0.42	9
(NonDiacritic) 1	0.38	0.33	0.35	9
accuracy			0.39	18
Macro avg	0.39	0.39	0.39	18
Weighted avg	0.39	0.39	0.39	18

$$\# \text{Precision} = \frac{\overbrace{TP}^1}{\overbrace{(TP+FP)}^2} / \frac{\overbrace{TN}^1}{\overbrace{(TN+FN)}^2}$$

$$\# \text{Recall} = \frac{\overbrace{TP}^1}{\overbrace{(TP+FN)}^2} / \frac{\overbrace{TN}^1}{\overbrace{(TN+FP)}^2}$$

$$\# \text{f1-score} = 2 \cdot \frac{\overbrace{\text{Precision} \times \text{Recall}}^1}{\overbrace{(\text{Precision} + \text{Recall})}^2}$$

$$\# \text{Accuracy} = \frac{\overbrace{(TP+TN)}^1}{\overbrace{(TP+TN+FP+FN)}^2}$$

## Validation

Print(cm)

~~OP~~  $\begin{bmatrix} [4\ 5] \\ [6\ 3] \end{bmatrix}$

# [TP FN]  
# [FP TN]

$$\text{Pre}0 = \frac{4}{(4+6)}$$

$$\text{Pre}1 = \frac{3}{(3+5)}$$

$$\text{Rec}0 = \frac{4}{(4+5)}$$

$$\text{Rec}1 = \frac{3}{(3+5)}$$

Print values

~~OP~~:

$$\text{Pre}0 \ 0.4$$

$$\text{Pre}1 \ 0.375$$

$$\text{Rec}0 \ 0.44444$$

$$\text{Rec}1 \ 0.33333$$

D)  $\text{F1}0 = 2 * \text{Pre}0 * \text{Rec}0 / (\text{Pre}0 + \text{Rec}0)$

$$\text{F1}1 = 2 * \text{Pre}1 * \text{Rec}1 / (\text{Pre}1 + \text{Rec}1)$$

Ans

~~OP~~

$$\text{F1-score}0 \ 0.421052631$$

$$\text{F1-score}1 \ 0.35294$$

E)  $\text{acc} = (4+3) / (4+3+5+6)$

Print(acc)

~~OP~~

$$0.3888$$

## Logistic Regression

→ It is a supervised learning algorithm used for classification.

→ LR is a classification algorithm

used to assign observations to a discrete set of classes.

→ It is used when the dependent variable (target) is categorical.

→ LR transforms its OP using the logistic or sigmoid function to return a probability value which can then be mapped to two or more discrete classes.

→ It is approximate algorithm when dependent variable is categorical and consists of two categories (binary).

→ It is a special case of linear regression where the target variable is categorical in nature. It uses a log of odds as the dependent variable.

→ Logistic regression predicts the probability of occurrence of a binary event utilizing a logit or sigmoid function.

$$\rightarrow \text{Odds Ratio} = \frac{P(\text{event happening})}{1-P(\text{event happening})}$$

$$\rightarrow \log \text{ of odds} = \log \left( \frac{p}{1-p} \right) = Y = \text{ans}$$

## Derivation of sigmoid function

$$\log\left(\frac{P}{1-P}\right) = ax+bx$$

$$\frac{P}{1-P} = e^{ax+bx}$$

$$P = e^{ax+bx} - P e^{ax+bx}$$

$$P(1+e^{ax+bx}) = e^{ax+bx}$$

$$P = \frac{e^{ax+bx}}{1+e^{ax+bx}}$$

$$P = \frac{1}{1+e^{-ax-bx}}$$

↓

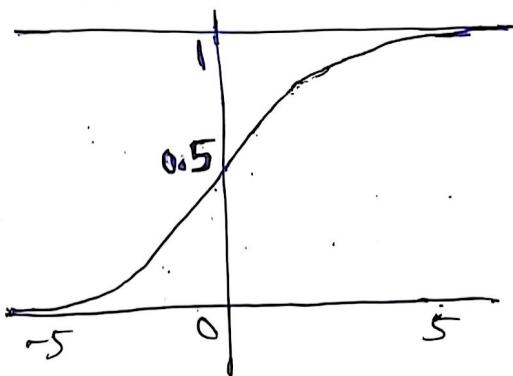
This eqn is called logistic or sigmoid or logit function

positive infinity,  $y$  predicted will become 1, and if the curve goes to negative infinity,  $y$  predicted will become 0.

→ If the O/P of sigmoid function is more than 0.5, we can classify the outcome as 1 or Yes and if it is less than 0.5, we can classify it as 0 or No.

$$S(z) = \frac{1}{1+e^{-z}}$$

where  $z = ax+bx$



## Sigmoid function,

→ In order to map predicted values to probabilities, we use the sigmoid function. The function maps any real value into another value between 0 and 1. In machine learning, we use sigmoid to map predictions to probabilities.

→ The sigmoid function, also called logistic function given

an 'S' shaped curve that can take any real-valued number and map it into a value b/w 0 & 1.

→ If the curve graph to

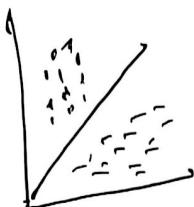
## Pros/Advantages

- Logistic regression is easier to implement, interpret and very efficient to train.
- It provides good accuracy for many simple data sets and it performs well when the dataset is linearly separable.

→ Due to its simple probabilistic interpretation, the training time of logistic regression algorithm comes out to be far less than most complex algorithms.

## Cons/Disadvantages

- only important and relevant features should be used to build the model otherwise the probabilistic predictions made by the model may be incorrect and the model's predictive value may degrade.
- doesn't perform well with non-linearly separable data



Linearly  
separable



Non-linearly  
separable

import libraries

p)

df = pd.read\_csv('insurance.csv')  
df.head()

out:

	age	bought_insurance
0	22	0
1	25	0
2	31	1
3	38	1
4	46	1

p) df.shape

out:

(27, 2)

p) df.isnull().sum()

out:

	age	bought_insurance
count	0	0
dtype	int64	int64

for every classification model:

- Model only accepts int or float as input for x.
- Model can accept input as strings/object (in certain models) as input for y
- Model will not accept any null value.
- x has to be either a Dataframe or a 2D numpy array or a list of list.
- y has to be either a series or a 1D numpy array or a list.

p) df.dtypes

	age	bought_insurance
dtype	int64	int64

dtype: object

p) x=df[['age']]

y=df['bought\_insurance']

print(x.shape)

print(y.shape)

print(type(x))

: (( ))

out:

(27, 1)

(27, 1)

<class 'pandas.core.frame.DataFrame'>

<class 'pandas.core.series.Series'>

B) from sklearn.model\_selection

import train-test-split.

x-train, x-test, y-train, y-test = train-test-split(x, y, test\_size  
= 0.25)

print(x-train.shape)

11 (x-test.shape)

11 (y-train.shape)

11 (y-test.shape)

OP)

(20, 1)

(7, 1)

(20, 1)

(7, 1)

P) x-train.head()

age

20 21

9 61

5 81

11 28

8 62

→ every output changes for  
above program because  
training is based on different  
rows

P) from sklearn.linear-model import Logistic Regression,

m1 = LogisticRegression()

m1.fit(x-train, y-train)

# accuracy names

print('Training score', m1.score(x-train, y-train))

print('Testing score', m1.score(x-test, y-test))

OP) Training score 0.9  
Testing score 0.78

P)  $y_{pred\_m1} = m1.predict(x\_test)$

$y_{pred\_m1}$

~~OR~~  
array([1, 1, 0, 1, 1, 1, 0], dtype=int64)

1) from sklearn.metrics import confusion\_matrix, classification\_report.

cm = confusion\_matrix(y\_test, y\_pred\_m1)

print(cm)

print(classification\_report(y\_test, y\_pred\_m1))

~~OR~~

$\begin{bmatrix} 2 & 1 \\ 0 & 4 \end{bmatrix}$

	precision	recall	f1-score	support
0	1.00	0.67	0.80	3
1	0.80	1.00	0.89	4
accuracy			0.86	7
macro avg	0.90	0.83	0.84	7
weighted avg	0.89	0.86	0.85	7

P) m = m1.coef\_

c = m1.intercept\_

print(m)

print(c)

~~OR~~

$\begin{bmatrix} 0.1575114 \end{bmatrix}$

$\begin{bmatrix} -5.81788525 \end{bmatrix}$

P) # sigmoid =  $1/(1+e^{-(m*x+c)})$

def sigmoid(x, m, c):

$\text{logistic} = 1/(1 + \text{np.exp}(-(m*x + c)))$

print(logistic)

## Problem Statement

- Generate predictions whether or not the person bought the insurance for given age.

(a) Age = 59

(b) Age = 23

b)  $y_{pred} - 59 = m1, predict([59])$

print( $y_{pred}-59$ )

sigmoid(59, m, c)

OP

[1]

[0.96997638]

b)  $y_{pred} - 23 = m1, predict([23])$

print( $y_{pred}-23$ )

sigmoid(23, m, c)

OP

[0]

[0.1001893]

## Multiclass classification

$y\_true = [1, 0, 1, 2, 0, 2, 1, 0, 2, 1, 0, 1, 1, 2, 1, 2, 2, 1]$

$y\_pred = [0, 1, 2, 1, 0, 2, 1, 1, 0, 1, 1, 0, 2, 1, 0, 1, 2, 0]$

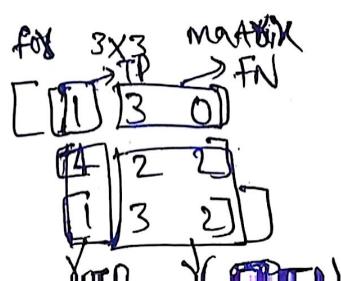
print( $ten(y\_true)$ ,  $ten(y\_pred)$ )

OP 18 18

b)  $cm2 = \text{confusion\_matrix}(y\_true, y\_pred)$

print(cm2)

OP [ [ 1 3 0 ]  
[ 4 2 2 ]  
[ 0 3 2 ] ]



`print(classification_report(y_true, y_pred))`

	Precision	Recall	f1-score	Support
0	0.17	0.25	0.20	4
1	0.25	0.25	0.25	8
2	0.50	0.33	0.40	6
accuracy			0.28	18
macro avg	0.31	0.28	0.28	18
weighted avg	0.31	0.28	0.29	18

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad \text{Recall} = \frac{\text{TN}}{\text{TN} + \text{FN}}$$

$$\text{Prc0} = \frac{1}{(1+5)}$$

$$\text{Prc1} = \frac{2}{(2+6)}$$

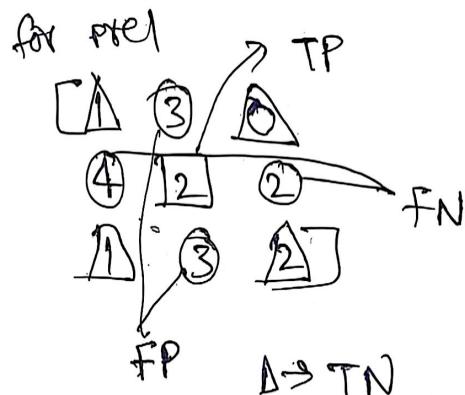
$$\text{Prc2} = \frac{2}{(2+2)}$$

$$\text{Prc0} = \frac{1}{(1+5)} \\ \text{Prc1} = \frac{2}{(2+6)} \\ \text{Prc2} = \frac{2}{(2+2)}$$

~~0.16666~~

$$0.25$$

$$0.5$$



## Decision Trees

- It falls under the category of supervised ML.
- It is used for both classification and regression.
- A decision tree is a flowchart-like structure in which each internal node represents a "test" on an attribute (e.g. whether a coin flip comes up heads or tails), each branch represents the outcome of the test, and each leaf node represents a class label (decision taken after computing all attributes).

→ The decision tree from the name itself signifies that it is used for making decisions from the given dataset. The concept behind the decision tree is that it helps to select appropriate features for splitting the tree into ~~different parts~~, subparts.

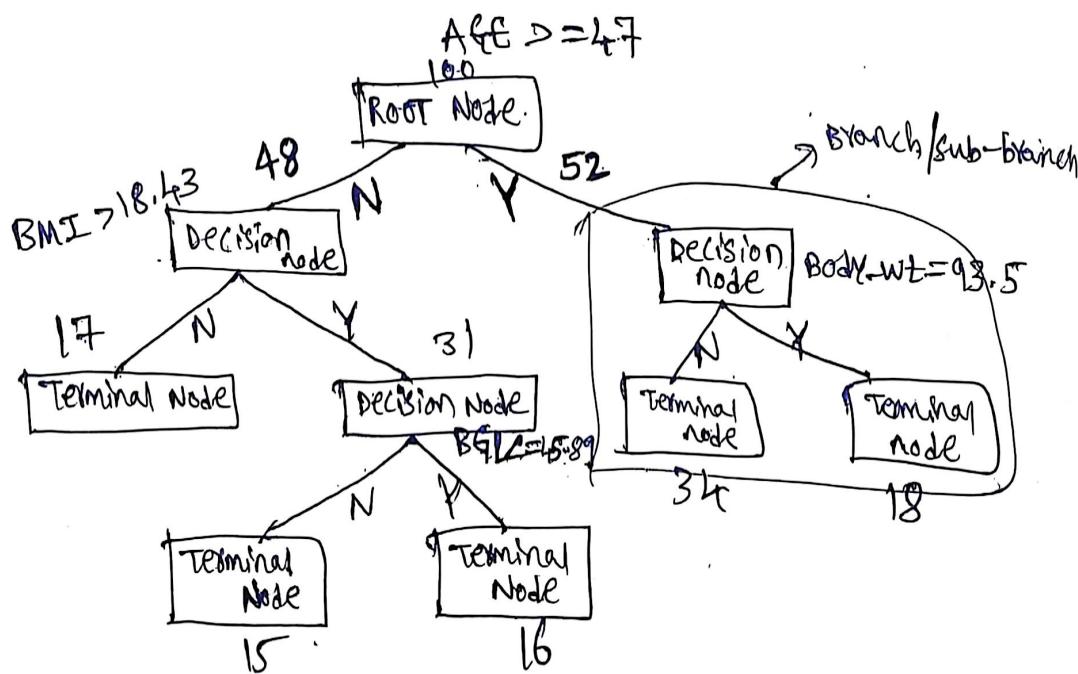
- It breaks down a dataset into smaller and smaller subsets while at the same time an associated decision tree is incrementally developed. The final result is a tree with decision nodes and leaf nodes.
- Decision tree uses the tree representation to solve the problem in which each leaf node corresponds to a class label (target variable) and attributes are represented on the internal node of the tree.

### Example:

We have 100 rows of data with

$X = \text{'Age'}, \text{'Gender'}, \text{'BMI'}, \text{'Body\_wt'}, \text{'Blood\_Glucose\_Level'}$ .

$Y = 0 \text{ or } 1$  ( $0 = \text{Non diabetic}$ ,  $1 = \text{Diabetic}$ )



## Decision Tree Algorithm:

$$E(\text{Target Variable}) = \sum p_i \log p_i$$

→ We use statistical methods (gini or entropy) for ordering attributes (information gain = entropy(target-variable) - entropy(target-variable, feature)) as root of the internal node.

### 1) Entropy:

→ It is the measure of uncertainty or disorder or impurity. It characterizes the impurity of an arbitrary collection of examples.

The higher the entropy more the information content.

→ If the sample is completely homogenous the entropy is zero and if the sample is equally divided it has entropy of one.

$$E = - \sum_{i=1}^n p_i \log_2(p_i)$$

Where

$p_i$  = probability of label i

n = number of categories in the label

### Information Gain:

→ It is based on decrease in entropy after dataset is split on an attribute (feature or column).

→ The split with highest information gain is chosen as the root node.

### 2) Gini

→ Gini index is a metric to measure how often a randomly chosen element would be incorrectly identified.

→ It means an attribute with lower gini index should be preferred.

$$\rightarrow G(\text{ini}) = 1 - \sum_{i=1}^n (p_i)^2$$

for 1 = 1 to  
no. of classes

→ Works well with numerical & categorical features.

→ Decision Tree can handle non-linear data.

→ Performs well on large datasets.

### Cons

→ It suffers from overfitting. Random Forest mitigates this issue.

→ Training time is high.

Overfitting, scattering b/w actual value  $y_i$  and predicted value  $\hat{y}_i$  now far is the actual value  $y_i$  from predicted  $\hat{y}_i$ .

- a) High variance and low bias
- b) High training score but low test score.

## Parameters ↗

- 1) Criterion - entropy or gini
- 2) max\_depth - max depth of decision tree.
- 3) min\_samples\_split - Minimum number of samples beyond which an internal node gets converted to a leaf node.

and merges them into a strong estimator,

- 5) RF, like its name implies, consists of a large number of individual decision trees that operate as an ensemble. Each individual tree in the random forest splits out a class prediction and the class with the most votes becomes our model's prediction,

## Random Forest Classification

RF is an ensemble machine learning technique capable of performing both regression and classification tasks using multiple decision trees and a statistical technique called Bagging (Bootstrap Aggregation).

A RF besides just averaging the prediction of trees it uses two key concepts that give it the name random:

(a) Random sampling of training observations (rows) when building trees.

(b) Random subsets of features (columns) for splitting nodes.

3) Random forest builds multiple decision trees and merge their predictions together to get a more accurate and stable prediction rather than relying on individual decision trees.

4) It combines a number of weak

6) Random forest Intuition - A large no. of relatively uncorrelated models (trees) operating as a committee will outperform any of the individual constituent models,

7) It uses bootstrapping (selection with replacement) to select random subsets of records and/or columns for creating individual decision trees. It means some samples might be used multiple times in a single tree.

### Example :-

- 1) There are 100 rows of data.
- 2) The target variable (Y) has 2 categories - A and B
- 3) In random forest, we are constructing 7 DTs.
- 4) Suppose out of 7 DTs, row 37 is a part of DT1, DT3, DT4, DT5, DT6.
- 5) Predictions generated by decision trees for row no 37 are:

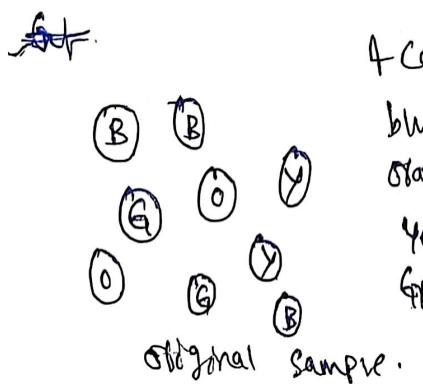
DT1 → A  
 DT3 → B  
 DT4 → A  
 DT5 → A  
 DT6 → B

6) Prediction out of RF for row no 37 is made of these predictions (3A and 2B)  
 $\rightarrow \text{mode}(3, 2) = 3.$

7) Prediction for Row no 37 is Category A according to RF;

### Bootstrapping:-

→ It is randomly selecting samples from training data with replacement. The samples so generated are called bootstrap samples.



### Working of Random Forest:-

- 1) Select random samples from a given dataset.
- 2) Construct a decision tree for each sample and get a prediction result from each decision tree.
- 3) Perform a vote for each predicted result.
- 4) Select the prediction result with the most votes as the final prediction.

(a) The random forest combines hundreds or thousands of decision trees, trains each one on a slightly different set of observations, splitting nodes in each tree considering a limited number of features.

(b) The final predictions of the random forest are made by averaging the predictions of each individual tree.

(c) The random forest is a classification algorithm consisting of many decision trees. It uses bagging and feature randomness when building each individual tree to try to create an uncorrelated forest of trees whose prediction by committee is more accurate than that of any individual tree.

## Working of Random Forest:

- Select random samples from a given dataset.
- Construct a decision tree for each sample and get a prediction result from each decision tree.
- Perform vote

on the individual tree.

b) sqrt: This option will take square root of total no. of features in individual tree. For instance, if the total no. of variables are 100, we can only take 10 of them in individual tree, "log2" is another similar type of option for max\_features,

c) o.x: This option allows the random forest to take  $x\%$  of variables in individual tree. We can assign any value in a format "o.x" where we want  $x\%$  of features to be considered.

### 5) bootstrap: true/false,

If it is set to true, samples are drawn with replacement and samples are drawn without replacement if bootstrap is set to false.

### PLOS:-

→ RF is considered as a highly accurate and robust method bcs of no. of decision trees participating in the process.

→ It overcomes the problem of overfitting faced by decision tree. The main reason is that it takes the average of all the predictions.

### 3) criterion: gini or entropy.

4) max\_features: These are the maximum no. of features RF is allowed to try in individual tree. The values it can have are:-

(a) Auto/None: This will simply take all the features which make sense in every tree. Here we simply do not put any restrictions

## CONS:

→ RF is slow in generating predictions bcz. it has multiple decision trees. whenever it makes a prediction, all the trees in the forest have to make a prediction for the same given input and then perform voting on it. This whole process is time-consuming.

) from sklearn.model\_selection import train\_test\_split  
 $x\_train, x\_test, y\_train, y\_test =$   
 $\text{train\_test\_split}(X, Y, \text{test\_size} = 0.25)$

print(x\_train.shape)

X-test - 31

y-train. 11

y-test 11

OR  
 $(111, 4)$

$(38, 4)$

$(111, 1)$

$(38, 1)$

## Implementation:-

) df = pd.read\_csv('iris.csv')  
df.head()

OR

	sepal.length	sepal.width	petal.length	petal.width	label
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa

) X = df.iloc[:, :-1]  
Y = df['label']  
X.shape  
Y.shape

all columns  
except last column

) from sklearn.tree import DecisionTreeClassifier

from sklearn.ensemble import RandomForestClassifier

) m1 = DecisionTreeClassifier(criterion=  
'gini',  
max\_depth=3, min\_samples\_split=15)

m1.fit(x\_train, y\_train)

) print('Training score:', m1.score(x\_train,  
y\_train))  
print('Testing score:', m1.score(x\_test, y\_test))

OR

$(150, 4) \rightarrow \text{dataframe}$

$(150, 1) \rightarrow \text{series.}$

OR

Training score 0.9819819

Testing 11 0.973684

P)

$$y_{\text{pred}} - m_1 = m_1 \cdot \text{predict}(x_{\text{test}})$$

$$\text{print}(y_{\text{pred}} - m_1)$$

~~Q1D~~)

[iris-varicolor, iris-setosa]

iris-virginical

- - - - - ]

P) from sklearn.metrics import

confusion\_matrix, classification\_report

P) def eval\_metrics(y\_test, y\_pred):

m = confusion\_matrix(y\_test, y\_pred)

print(m)

print(classification\_report(y\_test, y\_pred))

P) eval\_metrics(y\_test, y\_pred\_m1)

~~Q1D~~)

[[16 0 0  
0 17 1  
0 0 4]]

precision recall f1-score support

Iris-setosa 1.00 1.00 1.00 16

Iris-versicolor 1.00 0.94 0.97 98

Iris-virginica 0.80 1.00 0.89 4

accuracy 0.97 38

## Plotting decision Tree

P) from sklearn.tree import PlotTree

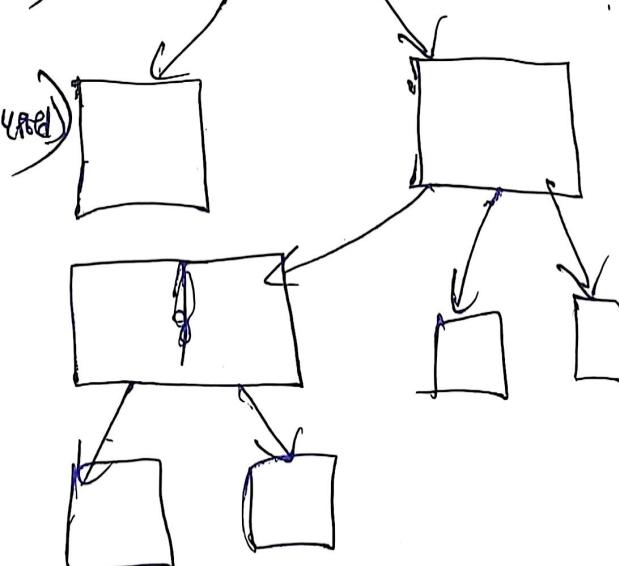
fn = x\_train.columns

cn = ['setosa', 'versicolor', 'virginica']

plot\_tree(m1, feature\_names=fn,  
class\_names=cn, filled=True)

plt.show()

Petal width <= 0.8  
gini = 0.657  
samples = 111  
value = [33, 32, 46]  
Class = virginica



## SVM (Support Vector Machines): find the optimal one? intuitively

→ A ~~SVM~~ is a supervised ML

→ A SVM takes in the data points (inputs) and outputs the hyperplane (which in two dimensions is simply a line) that best separates the data points.

→ The objective of the SVM algorithm is to find a hyperplane in an N-dimensional space (N - the no. of independent features) that distinctly classifies the data points.

→ It follows a technique called the kernel trick to transform the data and based on these transformations, it finds an optimal boundary between the possible outputs.

→ SVM objective - The main idea is to identify the optimal separating hyperplane which maximizes the margin of the training data.

→ There can be multiple hyperplanes, but which one of them is the best separating hyperplane?

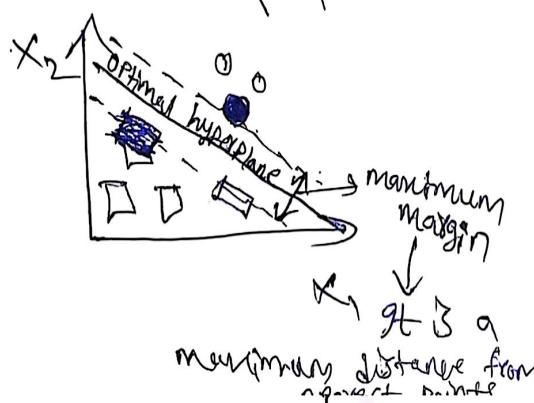
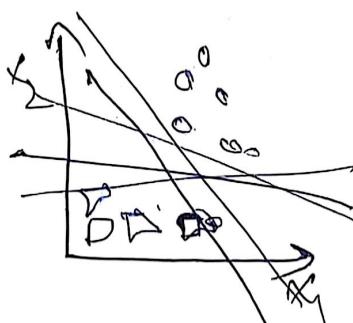
→ There can be multiple separating hyperplanes as well, that do we

if we select a hyperplane which is close to the data points of one class, then it might not generalize well (may not perform well on testing data).

→ So, the aim is to choose the hyperplane which is as far as possible from the data points of each category, therefore maximizing the distance b/w the nearest points of each class and the hyperplane would result in an optimal separating hyperplane,

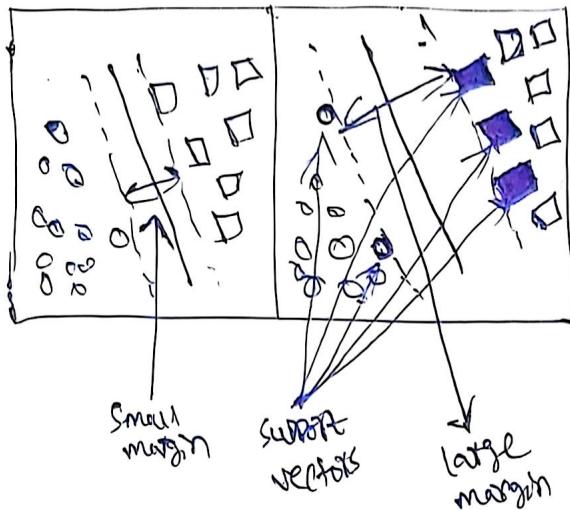
→ This distance is called margin. The goal of SVMs is to find the optimal hyperplane because it not only classifies the existing dataset but also helps predict the class of unseen data.

→ And the optimal hyperplane is the one which has the biggest margin.



→ SVM's objective is to find a plane that has the maximum margin, i.e. the maximum distance b/w data points of both classes.

Using these support vectors, the model maximizes the margin of the classifier. These are the points that help us build our SVM.



### 3) Margin :-

It is the distance b/w support vectors and the hyperplane.

→ A hyperplane in 2D is "line".

→ " " " " " 3D is "Plane".

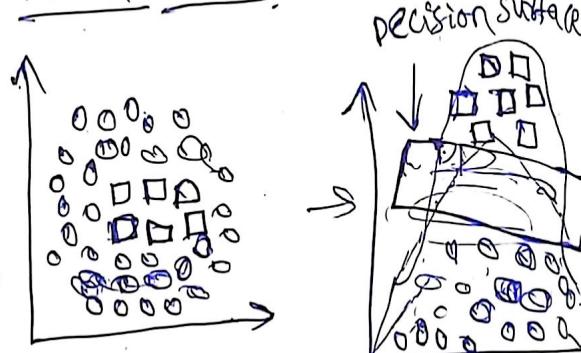
### Terminologies

1) Hyperplane - Hyperplanes are decision boundaries that help classify the data points. Data points falling on either side of the hyperplane can be attributed to different classes. Also, the dimension of the hyperplane depends upon the number of features. If the no. of input features is 2, then the hyperplane is just a line. If the no. of input features is 3, then the hyperplane becomes a two-dimensional plane. It becomes difficult when the no. of features exceeds 3.

### 2) Support Vectors:-

These are data points that are closest to the hyperplane and influence the position and orientation of the hyperplane.

### Kernel Trick :-



### SVM Parameters

#### 1) Gamma :-

→ The gamma parameter defines how far the influence of a single training example reaches, with low values meaning 'far' and high values meaning 'close'.

If gamma is low  $\Rightarrow$  data points far away from hyperplane, data points affect the hyperplane.

If gamma is high  $\Rightarrow$  data points close to hyperplane, data points affect the hyperplane.

Commonly used gamma

Kernel Types :

Values are  $[0.0001, 0.001, 0.01, 0.1, 1, 10, 100]$

$x_i, x_j$  represents data you are trying to classify.

## 2) C:

→ The C parameter tells the SVM optimization how much you want to avoid misclassifying each training example. Small values of C attribute to more misclassification and large values of C attribute to less misclassification.

Commonly used C values are

$[0.0001, 0.001, 0.01, 0.1, 1, 10, 100]$

2) Poly (Polynomial)  $K(x_i, x_j) = (x_i \cdot x_j)^d$

Here ". ." shows dot product of both the values and d denotes the degree of the polynomial.  $K(x_i, x_j)$  representing the decision boundary to separate the given classes.

## 3) Rbf (Radial Basis function):

$$K(x_i, x_j) = \exp(-\gamma * \|x_i - x_j\|^2)$$

The value of  $\gamma > 0$ .

$$\|x\| = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$$

SVM pros and cons

Pros :-

→ SVM is an algorithm which is suitable for both linearly and nonlinearly separable data (using Kernel trick)

→ SVM is very good when we have no idea about the data.

→ Kernel trick is what makes SVM unique

## 4) degree:

Defines degree of polynomial for 'poly' kernel

Linear kernel, Radial basis function kernel, Polynomial kernel (rbf) (Poly)

### Cons:-

- They are not suitable for larger datasets because the training time with SVMs can be high and much more computationally intensive.
- Tuning parameters can be time consuming

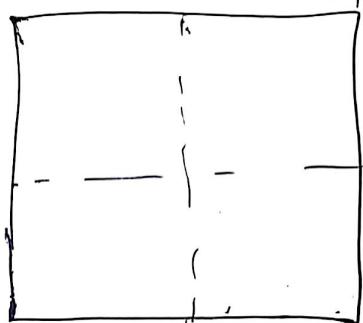
\* When we didn't know the data is linearly separable or not, just blindly apply SVM or decision tree or Random forest, then there is a possibility to get good results.

### SVM Simulator :-

Link  
<https://jgriemann.github.io/>

After opening link

Interactive demo of SVM



Kernel:

$$\gamma = \boxed{\quad} . G =$$

~~scribble~~

### K-Nearest Neighbors (KNN) : Alg.

→ It is supervised learning used for classification.

→ KNN Intuition - The KNN algorithm assumes that similar things exist in close proximity. In other words, similar things are near to each other.

→ KNN works by finding the distance between a query (test data) and all the examples in data (train data), selecting the specified no. of examples (k) closest to the query, then votes for the most frequent label (in the case of classification).

### KNN Classification Algorithm:-

- 1) Load the data.
- 2) Initialize k to your chosen no. of neighbors.
- 3) For getting the predicted class, iterate from 1 to total no. of training data points.

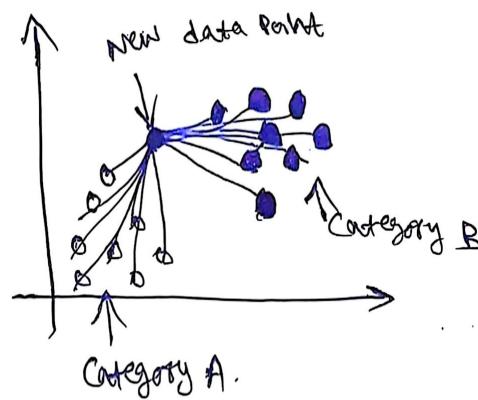
(a) Calculate the distance b/w test data and each row of training data. Use Euclidean distance as our distance metric since it's the most popular method. The other metrics that can be used are Manhattan, etc.

b) Sort the calculated distances in ascending order based on distance values.

c) Get top K rows from the sorted array.

d) Get the most frequent class (mode) of these rows.

e) Return the predicted class.



→ first compute distance from all data points

→ And then arrange in ascending order.

→ Get top K rows from sorted array.

→ Get the most frequent class (mode) of these rows as shown below.

→ Then return predicted class,

### Different Distance metrics in K-NN

#### 1) Euclidean Distance:

$$A = (x_1, y_1) \quad B = (x_2, y_2)$$

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

#### 2) Manhattan Distance:

$$AB = |x_2 - x_1| + |y_2 - y_1|$$

#### 3) Minkowski's Distance:

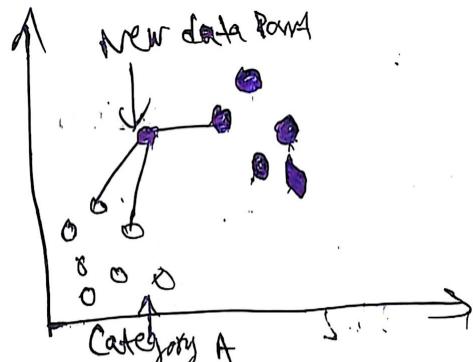
$$AB = \left( |(x_2 - x_1)|^p + |y_2 - y_1|^p \right)^{1/p}$$

Note

for  $p=2$ , Minkowski's distance equal to Euclidean distance

for  $p=1$ , " " " "

to Manhattan distance.



→ Then new data point belongs to category A.

#### Pros:

→ Low computation time.

→ Simple and easy to implement;

#### Cons:

→ It assumes that similar things occur in close proximity which can't be generalized for every single scenario.

## Titanic Survival Prediction :-

P) import libraries.

P) df = pd.read\_csv('titanic\_train.csv')  
df.head()

O/P:-

PassengerId survived name --- Embarked

0

1

2

3

4

P) df.isnull().sum()

O/P:-

PassengerId	0
survived	0
Pclass	0
Name	0
Sex	0
Age	177
SibSp	0
Parch	0
Ticket	0
Fare	0
Cabin	687
Embarked	2

P) df.shape

O/P:-

(891, 12)

P) df['survived'].value\_counts()

O/P:-

0 549

1 342

Name: survived, dtype: int64

P) df.dtypes.

O/P:-

PassengerId	int64
survived	int64
Pclass	int64
Name	object
Sex	object
Age	float64
SibSp	int64
Parch	int64
Ticket	object
Fare	float64
Cabin	object
Embarked	object
dtype: object	

\* In general, whenever we proceed with model building, we exclude the columns like PassengerId, SibSp, Registration number, Name, ~~Age~~, Cabin, Embarked, ~~dtype: object~~

## Handling Null Values.

P) `df['Age'].value_counts()`

O/P:-

24.00	30
22.00	27
18.00	26
- - -	- - -
55.50	1
74.00	1

To compute % of null values per column.

$$P) \frac{\text{len}[\text{"Null\_percentage"}]}{\text{len}[\text{Count}]} = \frac{\text{len}[\text{Count}] * 100}{\text{len}[\text{Count}], \text{sum}()}$$

O/P:-

	Count	Null_Percentage
PassengerId	0	0.000
Survived	0	0.000
Age	177	20.438799
Cabin	687	79.330254
Embarked	2	0.230947

P) `y = pd.DataFrame([df.isnull().sum()])`  
`y.columns = ['Count']`  
`y`

	Count
PassengerId	0
Survived	0
Pclass	0
Name	0
Sex	0
Age	177
SibSp	0
Parch	0
Ticket	0
Fare	0
Cabin	687
Embarked	2

To get rid of null values

P) `df['Age'].fillna(df['Age'].mean(), inplace=True)`

`df.isnull().sum.`

PassengerId	0
Age	0
Cabin	687
Embarked	2
<code>dtype: int64</code>	

Q) Out of total rows how many values are null.

$$P) \frac{\sum[\text{Age\_is\_null}]}{\text{len(df)}} = \frac{\sum[\text{Age} == \text{null}]}{\text{len(df)}} * 100$$

### Dropping non-significant columns

P) df.drop(['PassengerId', 'Cabin', 'Name'])

axis=1, inplace=True  
↓  
along columns.

O/P:-

	Count	NaN	Age\_of\_null
PassengerId	0	0.00	0.00
Survived	0	1	1
Age	177	20.438	19.86532
Cabin	687	79.338254	77.10457
Embarked	2	0.2309	0.124467

df.isnull().sum()

	survived	Pclass	Sex	Age	Embarked
survived	0	0	0	1	0

P) df['Embarked'].value\_counts()

O/P:-

S	644
C	168
Q	77

Name: Embarked, dtype: int64

P) df['Embarked'].fillna('S', inplace=True)

df.isnull().sum()

O/P:-

PassengerId	0
Age	0
Cabin	687
Embarked	0

### Label Encoding

P) df.dtypes

Survived	int64
sex	object
Age	float64
Ticket	object
Fare	float64
Embarked	object

P) df['sex'].value\_counts()

O/P:-

male	577
female	314

Name: Sex, dtype: int64

Q) df['Ticket'].unique()

O/P:-

681

It is irrelevant feature, drop it.

P) df.drop(['Ticket'], axis=1, inplace=True)

df.dtypes.

O/P:-

Survived	int64
Pclass	int64
Sex	object
Age	float64
Embarked	?

P) df.head()

O/P:-

	Survived	Pclass	Sex	Age	Embarked
0	1	1.	male	25	S
1	.	.	female	C	
2	.	.	female	S	
3	.	.	female	S	
4	1	.	male	S	

Convert Sex, Embarked into integer

O/P:-

Encoding Categorical Columns

sex	Sex-1b	Sex-M	Sex-F
m	0	1	0
f	1	0	1
M	0	1	0
M	0	1	0
f	1	0	1
f	1	0	1
m	0	1	0

label encoding

This approach is One hot encoding

P) from sklearn.preprocessing import LabelEncoder

P) lb = LabelEncoder()

df['Sex'] = lb.fit\_transform(df['Sex'])

df['Embarked'] = lb.fit\_transform(df['Embarked'])

P) df['Embarked'].value\_counts()

O/P:-

2	646	#S 646
0	168	#C 168
1	77	#Q 77

Name: Embarked, dtype: int64,

0,1,2 assig alphabetically.

P) df['Sex'].value\_counts()

O/P:-

1	577	#male 577
0	314	#female 314

then).

P) df.head()

Output

	Survived	Sex	Fare	Embarked
0	1	1	2	0
1	0	0	2	2
2	0	0	2	2
3	0	1	2	2
4	1			

Output

## Building the models

P) df.columns

Output

```
Index(['survived', 'Pclass', 'Sex',  
       ..., 'Embarked'], dtype='object')
```

P) x = df.drop(['survived', axis=1]. # df.iloc[:, 1:]  
 # Allows columns starting from 1

y = df['survived']

print(x.columns)

Output

```
Index(['Pclass', 'Sex', 'Age', ..., 'Embarked'], dtype='object')
```

P) print(type(x))

print(type(y))

Output

<class 'pandas.core.frame.DataFrame'>

<class 'pandas.core.series.Series'>

P) from sklearn.model\_selection import train\_test\_split.

x\_train, x\_test, y\_train, y\_test = train\_test\_split(x, y, test\_size=0.25)

print(x\_train.shape)

print(x\_test.shape)

```
print(y_train.shape)  
print(y_test.shape)
```

Output:-

```
(668, 7)  
(223, 7)  
(668, )  
(223, )
```

- p) from sklearn.linear\_model import LogisticRegression.  
from sklearn.tree import DecisionTreeClassifier  
from sklearn.ensemble import RandomForestClassifier  
from sklearn.neighbors import KNeighborsClassifier  
from sklearn.svm import SVC
- p) from sklearn.metrics import confusion\_matrix, classification\_report  
p) def eval\_metrics(ytest, ypred):  
 cm = confusion\_matrix(ytest, ypred)  
 print(cm)  
 print(classification\_report(ytest, ypred))

Logistic Regression:-

- p) m1 = LogisticRegression(max\_iter=1000)  
m1.fit(x\_train, y\_train)
- p) def printscore(model):  
 print('Training score:', model.score(x\_train, y\_train))  
 print('Testing Score:', model.score(x\_test, y\_test))

p) mscore (m)

o/p Training score 0.80389  
Testing score 0.7937  $\rightarrow$  accuracy.

p)  $y_{pred\_m} = m_1 \cdot \text{predict}(x_{-test})$ .

p) eval metrics ( $y_{-test}$ ,  $y_{pred\_m}$ )

#0  $\begin{bmatrix} 0 \\ TP & FN \end{bmatrix}$

#1  $\begin{bmatrix} FP & TN \end{bmatrix}$

O/P)

$$\begin{bmatrix} [12 & 18] \\ [28 & 65] \end{bmatrix}$$

	precision	recall	f1-score	support
0	0.80	0.86	0.83	130
1	0.78	0.70	0.74	93

accuracy

0.79

223

macro avg 0.79

0.78

0.78

223

weighted avg 0.79

0.79

0.79

223

#0 - no. of non-survivors, 1 - no. of survivors.

# TP - Actual value is 0, ML model predicted 0 = 122

# FN - 0, 11, 11, 0, 11, 11, 11, 11, 11, 1 = 18

# FP - 11, 11, 11, 1, 11, 11, 11, 11, 0 = 28

# TN - 11, 11, 11, 1, 11, 11, 11, 11, 11, 1 = 65

$$TP = 112, FN = 18$$

$$FP = 28, TN = 65$$

$$\# \text{P}r\text{eo} = \frac{TP}{(TP+FP)} = \frac{112}{(112+28)}$$

$$\# \text{P}re\text{l} = \frac{TN}{(TN+FN)} = \frac{65}{(65+18)}$$

$$\# \text{reco} = \frac{TP}{(TP+FN)} = \frac{112}{(112+18)}$$

$$\# \text{rec1} = \frac{TN}{(TN+FP)} = \frac{65}{(65+28)}$$

1) print("preo", preo)

!:

~~def~~

preo 0.8

prel 0.723

reco 0.861

rec1 0.6989

p)

eval-metrics(y-test, ypred-m2)

~~def~~-

$\begin{bmatrix} 112 & 18 \\ 30 & 63 \end{bmatrix}$

	precision	recall	f1-score	suppo
0	0.80	0.91	0.85	130
1	0.84	0.68	0.75	93
accuracy			0.81	223

macro avg 0.82 0.79 0.80 223  
weighted avg 0.82 0.81 0.81 223

## Random Forest

~~def~~

### Decision Tree

m2 = DecisionTreeClassifier(criterion='gini', max\_depth=6, min\_samples\_split=15)

criterion='entropy', min\_samples\_split=20,

=90,

max\_depth=8)

m2.fit(x-train, y-train)

m2.fit(x-train, y-train)

1) mscore(m2)

~~def~~- Training score 0.8802

Testing score 0.81165

p) mscore(m3)

~~def~~- Training score 0.8787

Testing score 0.8206

p) ypred-m2 = m2.predict(x-test)

p) ypred-m3 = m3.predict(x-test)

ypred-m2

p) eval-metrics(y-test, ypred-m3)

~~def~~- array([1, 0, 0, ..., ...])

## SVM

~~OP~~  
 $\begin{bmatrix} [123 & 7] \\ [33 & 60] \end{bmatrix}$

P)  $m5 = SVC(\text{kernel}='\text{linear}', C=1)$

$m5.\text{fit}(X\text{-train}, Y\text{-train})$

	Precision	Recall	F1-score	Support	P) mscore(m5)
0	0.79	0.95	0.86	130	<del>OP</del>
1	0.90	0.65	0.75	93	<del>OP</del>
accuracy			0.82	223	Training score 0.7889
					Testing score 0.7802
macro avg	0.84	0.80	0.81	223	P) $y_{\text{pred}} - m5 = m5.\text{predict}(X\text{-test})$
weighted avg	0.83	0.82	0.81	223	print( $y_{\text{pred}} - m5$ )

## KNN classifier

P)  $m4 = KNeighborsClassifier(n\_neighbors = 5)$

$m4.\text{fit}(X\text{-train}, Y\text{-train})$

P)  $m4.\text{mscore}(m4)$

number  
should be  
odd

~~OP~~  
Training score 0.71556

Testing 11 0.6816

P)  $y_{\text{pred}} - m4 = m4.\text{predict}(X\text{-test})$

P) eval-metrics ( $Y\text{-test}, y_{\text{pred}} - m4$ ) accuracy

~~OP~~

$\begin{bmatrix} 1 & 0 & 0 & - & - & - \end{bmatrix}$

$\vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots$

P) eval-metrics ( $Y\text{-test}, y_{\text{pred}} - m4$ )

~~OP~~

$\begin{bmatrix} [10 & 20] \end{bmatrix}$

$\begin{bmatrix} [29 & 64] \end{bmatrix}$

Precision recall f1-score support

0	0.79	0.85	0.82	130
1	0.76	0.69	0.72	93

P)  $y_{\text{pred}} - m4 = m4.\text{predict}(X\text{-test})$

P) eval-metrics ( $Y\text{-test}, y_{\text{pred}} - m4$ ) accuracy

0.78 223

~~OP~~  
 $\begin{bmatrix} [12 & 18] \\ [53 & 46] \end{bmatrix}$

macro avg 0.78 0.71 0.77 223

weighted avg 0.78 0.78 0.78 223

Precision recall f1-score support

0	0.68	0.86	0.76	130
1	0.69	0.43	0.53	93

accuracy 0.68 223

macro avg 0.68 0.65 0.64 223

weighted avg 0.68 0.65 0.64 223

~~OP~~, Best performing model

is Random Forest Classifier,

Gender wise passenger survival

Count

```
R> df.groupby(['Survived', 'Sex']).size
      .count()
```

Out:

Survived	Sex	Count
0	0	81
0	1	468
1	0	233
1	1	109

Name: Sex, type: int64.

\* If training score is high when compared to testing score, then the model is overfitted  
for example:-

Training Score 0.899

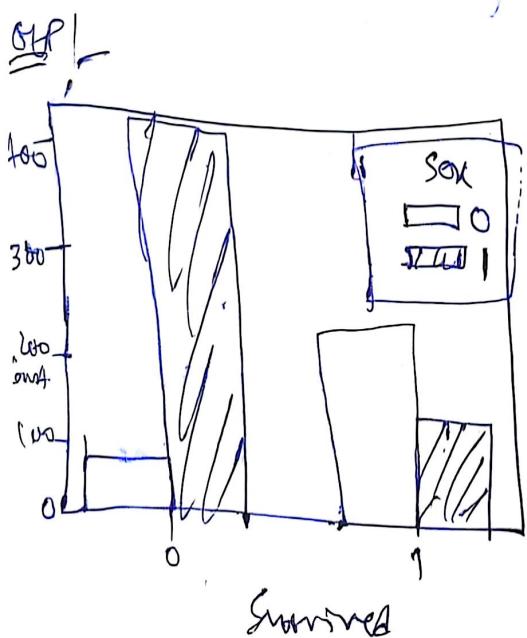
Testing Score 0.744

↙  
then this model is overfit model

We can use functions to write all the common things,

R> import seaborn as sns

: sns.countplot(x=df['Survived'], hue=df['Sex'])  
plt.show()



R> def model\_gen(model, x\_train, x\_test, y\_train, y\_test):

model.fit(x\_train, y\_train)

print('Training score', model.score(x\_train, y\_train))

print('Testing score', model.score(x\_test, y\_test))

y\_pred = model.predict(x\_test)

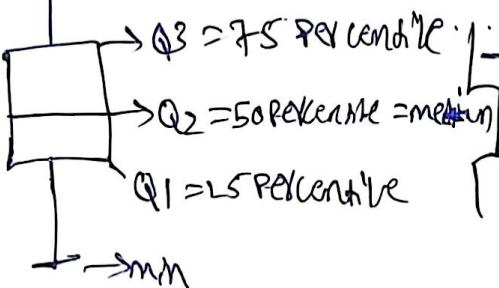
print(y\_pred)

cm = confusion\_matrix(y\_test, y\_pred)

print(Confusion\_Matrix\_In, cm)

print('Classification Report', classification\_report(y\_test, y\_pred))

About box plot)  
T → Max



without previously knowing what the output will be.

→ There is no need to split the data in training and testing dataset.

IQR = Inter Quartile Range

$$= Q_3 - Q_1$$

$$\text{Mean} = 9.3 + 1.5 * \text{IQ}$$

$$\text{Min} = q1 - 1.5 * \text{IQR}$$

→ Any values that are greater than mark and less than are called outliers.

## KMeans

→ K-means algorithm is an iterative algorithm that tries to partition the dataset into  $K$  pre-defined distinct non-overlapping subgroups (clusters) where each data point belongs to only one group.

# Clustering

→ clustering identifies similarities between objects, which it groups according to those characteristics in common and which differentiate them from other groups of objects. These groups are known as "clusters". Similar as possible while also keeping the clusters as different (far) as possible. It assigns a data points to a cluster such that the sum of the squared distance between the data points and the cluster

→ Clustering is framed in centroid (arithmetic mean of all unsupervised learning; that is, the data points that belong to for this type of algorithm we that cluster) is at the only have one set of input minimum. The less variation

more homogeneous (similar) the data points are within the same cluster.

Since clustering algorithms including KMeans which use distance based measurements to determine the similarity between data points, it's recommended to standardize or scale the data since almost always the features in any dataset would have different units of measurements for instance as age vs. income,

### K-Means Algorithm

- 1) Specify no. of clusters
- 2) Initialize centroids by first shuffling the dataset and then randomly selecting K data points for the centroids without replacement.
- 3) Keep iterating until there is no change to the centroids. i.e. assignment of data points to cluster isn't changing.
- 4) Compute the sum of the squared distance b/w data points and all centroids.

5) Assign each data point to the closest cluster (centroid)

6) Compute the centroids for the clusters by taking the average of the all data points that belong to each cluster.

Elbow Method to find the optimal number of clusters:-

→ It gives us an idea on what a good K no. of clusters would be based on the WCSS (within cluster sum of square) between data points and their assigned clusters centroid. We pick 'K' at the spot where WCSS starts to flatten out and forms an elbow.

WCSS = sum of squared distance b/w individual data points and its closest centroid for all the data points.

#### Pros

- Simple to implement & understand.
- It gives best results, when data sets are distinct.

#### Cons

- The user has to specify K (no. of clusters) in the beginning.

- Data is required to be scaled in most cases.
- Centroids can be dragged by the outliers.
- As the no. of dimensions (no. of columns/features) increases, it is recommended us to use PCA to reduce dimensions.

P) df.head()

	CustomerID	Gender	Age	AnnualIncome	SpendingScore
0	1	Male	19	15	39
1	2	Female	21	15	81
2	3	Male	26	16	6
3	4	Female	26	16	77
4	5	Male	27	17	50

P) df.iloc[:, -2:] = df  
df.head()

P) df = pd.read\_csv('customers.csv')

df.head()

	CustomerID	Gender	Age	AnnualIncome	SpendingScore
0	1	Male	19	15	39
1	2	Female	21	15	81
2	3	Male	26	16	6
3	4	Female	26	16	77
4	5	Male	27	17	50

: Renaming the columns.

P) df.columns = ['CustomerID', 'Gender', 'Age',  
'AnnualIncome', 'SpendingScore']

df.columns

Index(['CustomerID', 'Gender', 'Age',  
'AnnualIncome', 'SpendingScore'],  
dtype='object')

	CustomerID	Gender	Age	AnnualIncome	SpendingScore
0	1	Male	19	15	39
1	2	Female	21	15	81
2	3	Male	26	16	6
3	4	Female	26	16	77
4	5	Male	27	17	50

P) X = df.values

X[:5]

P) array([ [15, 39],  
[15, 81],  
[16, 6],  
[16, 77],  
[17, 50]], dtype=int64)

Clustering Annual Income and  
Spending Score!

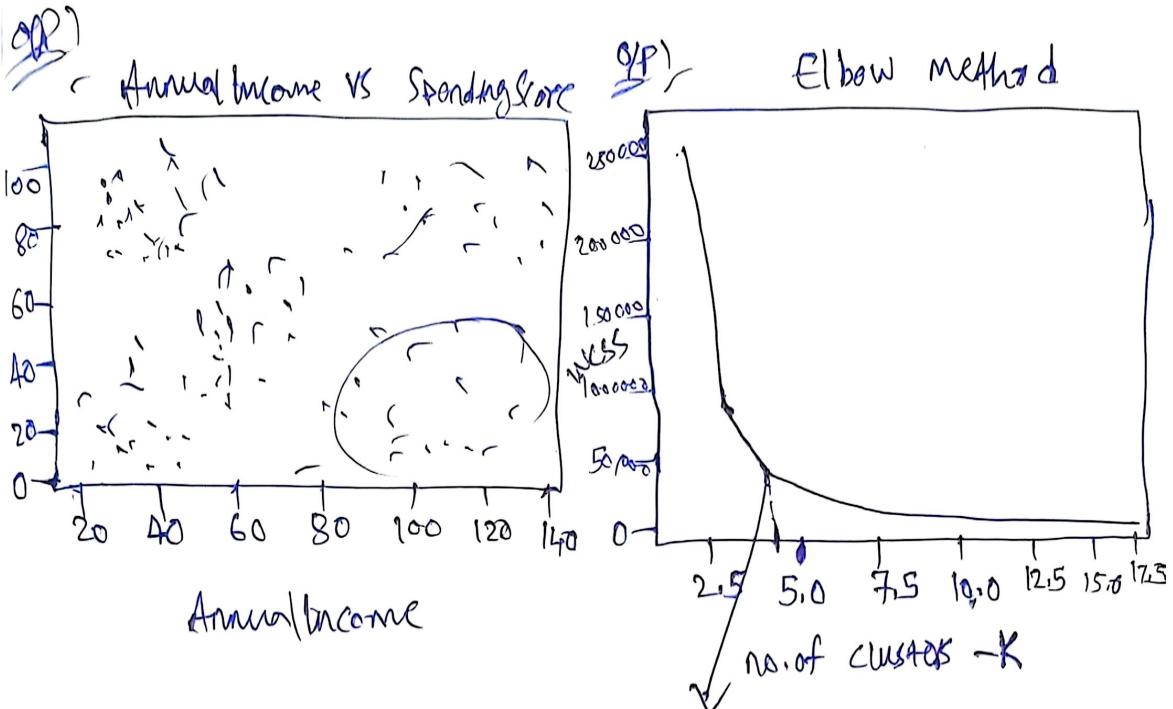
P) plt.scatter(X[:, 0], X[:, 1])

plt.title('AnnualIncome vs Spending Score')

plt.xlabel('AnnualIncome')

plt.ylabel('SpendingScore')

plt.show()



Elbow method

P) from sklearn.cluster import KMeans

$\text{WCSS} = []$

for i in range(1, 19):

m = KMeans(n\_clusters=i)

m.fit(x)

WCSS.append(m.inertia\_)

print('WCSS\n', WCSS)

plt.plot(list(range(1, 19)), WCSS)

plt.xlabel('Number of clusters - K')

plt.ylabel('WCSS')

plt.title('Elbow Method')

plt.show()

Q.P)

WCSS

[269981.28, 181363.5959,

106348.37306, ...]

Applying: KMeans at k=5

P) K5 = KMeans(n\_clusters=5)

K5.fit(x)

P) YPred5 = K5.predict(x)

YPred5

Q.P)

array ([4, 1, 4, 1, 4, 1, ...

----- 2, -----  
-----, 0, -----, 3])

P)

Cent5 = K5.cluster\_centers\_

Cent5

Q.P)

array ([ [88.2, 17.1142857],

[25.7273, 79.3636],

[..., ..., ..., ...])

P) ~~PLT-SCATTER~~ ( $x[6,1], x[1,1], c = y_{pred}[5]$ )

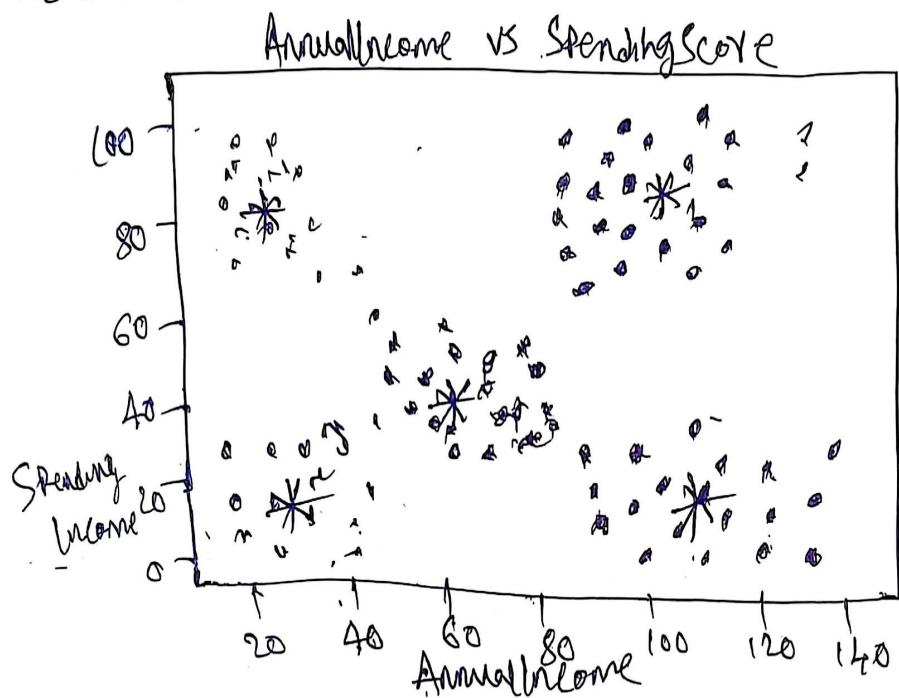
```
plt.scatter(cent5[:,0], cent5[:,1], color='red', s=130, marker=1)
```

plt.title('AnnualIncome vs SpendingScore')

```
plt.xlabel('AnnualIncome')
```

plt.ylabel('Spending(\$k\$)').

plt.show()



## Applying k-means at K=3

P) K3 = KMeans (n<sub>clusters</sub> = 3)

K3. fit ( $x$ )

D)  $\text{y}_{\text{pred3}} = \text{k3}.\text{predict}(x)$

4 Feb 3

8)  $\text{Cent3} = x_3, \text{cluster\_center}_5$

Print(cent3)

~~(1P)~~ - [ 44.1544 49.8292 ]

[87. - 18.6315]

[86,5384 82.12820]

plt.scatter( $x[i, 0]$ ,  $x[i, 1]$ , c='red')

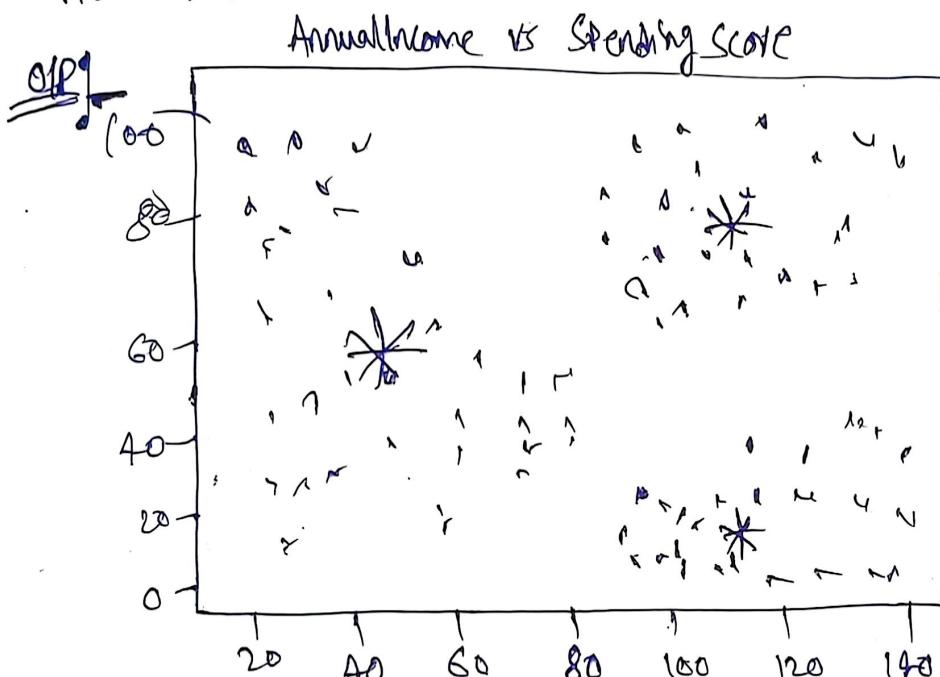
plt.scatter(cent3[:, 0], cent3[:, 1], color='red', s=150,

plt.title('AnnualIncome vs SpendingScore')  
marker='\*' )

plt.xlabel('AnnualIncome')

plt.ylabel('SpendingScore')

plt.show()



### Mean-Shift Clustering Algorithm

→ We don't have to specify the no. of clusters at the time of model creation.

→ It is based on a bandwidth parameter,

### KMeans

→ We have to specify the no. of clusters at the time of model creation,

→ It assigns data points to a cluster such that the sum of the squared distance b/w the data points and the cluster's centroid (arithmetic mean of all the data points that belong to that cluster) is at the minimum.

## Hyperparameters in different ML models:

- Random Forest (n\_estimators, max\_depth, criterion, min\_samples\_split)
- DecisionTreeClassifier (max\_depth, criterion, min\_samples\_split)
- SVM (kernel, degree, C, gamma)
- KNN (n\_neighbors)

P) Import libraries

P) df = pd.read\_csv('sonar.csv', header=None)  
df.head()

id	0	1	2	3	-	-	-	-	60
0	0.0200								R
1		1							R
2			1						R
3				1					M
4					1				M

P) df.shape

Out: (208, 6)

P) nv = df.isnull().sum()

nv

nv

Out: - Series ([], dtype: int64)

P) df[60].value\_counts()

#M - mines

#R - rocks

Out:

M 111

R 97

P) df.dtypes.value\_counts()

Out	float64	60
	object	1

P) x = df.iloc[:, :-1]

y = df.iloc[:, -1]

print(x.shape)

print(y.shape)

Out: (208, 6)

(208, 1)

```
from sklearn.model_selection import DecisionTree
```

```
train_test_split
```

```
x_train, x_test, y_train, y_test =
```

```
train_test_split(x, y, test_size=0.25)
```

```
print(x_train.shape)
```

```
print(x_test.shape)
```

```
print(y_train.shape)
```

```
print(y_test.shape)
```

```
OR {
```

```
(156, 60)
```

```
(52, 60)
```

```
(156, )
```

```
(52, )
```

```
from sklearn.tree import DecisionTreeClassifier
```

```
DecisionTreeClassifier
```

```
m1 = DecisionTreeClassifier(criterion='gini', max_depth=9, min_samples_split=15)
```

```
m1.fit(x_train, y_train)
```

```
m1.score(m1)
```

```
OR {
```

```
Training Score 0.9230
```

```
Testing Score 0.6730
```

```
from sklearn.metrics import confusion_matrix
```

```
classification_report.
```

```
y_pred_m1 = m1.predict(x_test)
```

```
print(y_pred_m1)
```

```
OR {
```

```
[R' M' M' R' - -  
- - - - - ]
```

```
def mscore(model):
```

```
print('Training Score', model.score(x_train, y_train))
```

```
print('Testing Score', model.score(x-test, y-test))
```

```
def gen_metrics(y-test, y-pred):
```

```
Cm = confusion_matrix(y-test, y-pred)
```

```
print(Cm)
```

```
print(classification_report(y-test, y-pred))
```

```
gen_metrics(y-test, y-pred)
```

```
OR {
```

```
[E 5 5]  
[8 14]
```

```
[ 24 5  
 02 10 ]
```

	Precision	Recall	F1-Score	Support
M	0.77	0.83	0.74	29
R	0.59	0.48	0.56	23

	accuracy	
	0.67	52

	macro avg	
	0.68	0.65

	weighted avg	
	0.68	0.67

## Hyperparameter Tuning

→ Tuning model hyperparameters such that the model generates the best performance.

These are 2 methods of Hyperparameter Tuning.

### 1) GridSearchCV

→ It takes into consideration all the permutations & combinations of the hyperparameters & returns the hyperparameters that would generate the best result.

→ It consumes a lot of time.

### 2) RandomizedSearchCV

→ It takes into consideration some random subset of hyperparameters and return the hyperparameters that would generate the best result out of those randomly sampled hyperparameters.

→ It consumes a comparatively less time.

P) from sklearn.model\_selection import

RandomizedSearchCV, GridSearchCV

### GridSearchCV on DecisionTree

P)

params\_df = {  
'criterion': ['gini', 'entropy']

'max\_depth': list(range(7, 12)),

'min\_samples\_split': [10, 12, 15, 17, 18, 20]

params\_df

OP)

{  
'criterion': ['gini', 'entropy'],

'max\_depth': [7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20]}

min\_samples\_split: [10, 12, 15, 17, 18, 20]

P) dt = DecisionTreeClassifier()

gs1 = GridSearchCV(estimator=dt),

param\_grid = params\_df, scoring =  
(accuracy)

gs1.fit(x\_train, y\_train)

OP)

OP) ➤ GridSearchCV

► estimator: DecisionTreeClassifier

① `print(gs1.best_params_)`  
`print(gs1.best_estimator_)`

Q1)

{'criterion': 'gini', 'max\_depth': 8,  
 'min\_samples\_split': 12}

DecisionTreeClassifier(max\_depth=8,

min\_samples\_split=12)

P) `y_pred_gs1 = gs1.predict(x-test)`

`print(y_pred_gs1)`

Q2)

[M] [M] 'R' ...

- — — ]

P) gen-metrics(y-test, y\_pred\_gs1)

Q3)

[23 4]

[2 15]

	Precision	Recall	F1-score	Support
M	0.68	0.86	0.76	29
R	0.73	0.48	0.58	23
accuracy			0.69	52
macro avg	0.70	0.67	0.67	52
weighted avg	0.70	0.69	0.68	52

## RandomizedSearchCV on Random Forest

\* Random forest & SVM takes more time when GridSearchCV is used. That's why RandomizedSearchCV is used for both RF & SVM.

P) from sklearn.ensemble import

RandomForestClassifier

P) `params_rf = {'n_estimators': [10, 80, 100, 120],`

'criterion': ['gini', 'entropy'],

'max\_depth': list(range(1, 12)),

'min\_samples\_split': [10, 12, 15, 17, 18, 20]}

P) `rf1 = RandomForestClassifier()`

`gs1 = RandomizedSearchCV(estimator=rf1,`

param\_distributions=params\_rf, scoring='accuracy', cv=8)

`gs1.fit(x-train, y-train)`

Q4)

► RandomizedSearchCV

► estimator: RandomForestClassifier

④ print(rs1.best\_params\_)

## Applying GridSearchCV on KNN

```
P) from sklearn.neighbors import  
    KNeighborsClassifier
```

1

$\{ \text{n-estimators} : 100, \text{min\_samples\_split} : 10, \text{max_depth} : 9, \text{criterion} : \text{gini} \}$  | P)   
 params\_knn = {n\_neighbors: list(range(5, 20))};

RandomForestClassifier(max\_depth=9,  
min\_samples\_split=10)

$\text{y\_pred\_ys1} = \text{ys1}.\text{predict}(x\_test)$

10

[R' M' R' R' R, -]  
- - - ]

```
P)  
Params-knn = { 'n_neighbors': list(range  
                           (5, 20)) }  
print(Params-knn)  
OP:
```

{'n-neighbors': [5, 6, 7, 8, ..., - - 27, 28]}

①) knnq = KNeighborsClassifier()

$gs2 = GridSearchCV(knn1, param_grid,$   
 $= param_knn, scoring='accuracy',$   
 $CV=10)$

p) gen-metrics (y-test, ypred-ys1) g52.fit (x-train, y-train)

5

[ [ 25 4 ] ]

[4 19]

	Precision	Recall	F1-score	Support
M	0.86	0.86	0.86	29
R	0.83	0.83	0.83	23
accuracy			0.85	52
macroavg	0.84	0.84	0.84	52
weightedavg	0.85	0.85	0.85	52

01P1

## → GridSearchCV

## Estimator: K Neighbors Classifier

P) print(gs2, best\_pairs -)

```
print(gs2.best_estimator_)
```

10

'n-neighbors' : 5 }

KNeighborsClassifier()

$$f) y_{pred} - g_{S2} = g_{S2} \cdot p_{S2} \text{dice2}(x - k_{S2})$$

$$\approx \text{mt}(y_{pred} - g_{S2})$$

~~OR~~)

$$[R' M' M' M' \dots]$$

$$[ \quad \quad \quad \quad \quad \quad ]$$

$$g) \text{gen-metrics} - (y_{test}, y_{pred} - g_{S2})$$

~~OR~~)

$$[[26 \ 3]]$$

$$[8 \ 15]]$$

	Precision	Recall	f1-score	Support
--	-----------	--------	----------	---------

M	0.76	0.90	0.83	29
---	------	------	------	----

R	0.83	0.65	0.73	23
---	------	------	------	----

accuracy		0.79	52
----------	--	------	----

macro avg	0.80	0.77	0.78	52
-----------	------	------	------	----

weighted avg	0.80	0.79	0.78	52
--------------	------	------	------	----