

# Machine learning

(1.08, 0.23, 0.23, 2.53) edit

it is a branch of AI that develops the algorithm by learning the hidden patterns of the data set used to make prediction on new data points

We can divide machine learning in 3 parts →

① **Supervised learning**

② **Unsupervised learning**

③ **Reinforcement learning**

## 1) **Supervised learning**

1) Supervised learning model is trained on **labeled data**.

label data → input data + Output data.

2) The differentiation is done on the basis of the output data.

3) It has Two main categories → A) **Regression** (continuous data  
22.5, 25.3, 28.5)

B) **Classification** (discrete data  
1, 2, 3, Yes, No)

A) Regression → 1) It deals with continuous target variable

like  $(22.5, 25.6, 28.9, 30.1)$

→ e.g. House price prediction based on its area, size

B) Classification → 1) It deals with categorical/discrete target variable

like  $(\text{yes}, \text{No}, 1, 2, 3)$  whole data

e.g. email spam or not

loan approved or not.

Note → 1) If targeted variable has continuous data then use regression algorithms.

2) If targeted variable contains discrete or categorical data then use classification Algorithm

### Regression Algorithm

- 1) linear Regression
- 2) Polynomial Regression
- 3) Decision Tree
- 4) Random forest.

### Classification Algorithm

- 1) logistic Regression
- 2) Support Vector Machine
- 3) decision tree
- 4) k-nearest neighbours
- 5) Naive Bayes

## 2) Unsupervised Machine learning

- 1) Unsupervised learning model is trained on unlabeled data set
- 2) Unlabeled data set → only input data (No target variable)
- 3) This differentiation is done on the basis of the output data.
- 4) It has two main categories →
  - 1) Clustering
  - 2) Association

A) Clustering → In these processes, there are different groups created.

It is a process of grouping data points into clusters based on similarity.

B) Association → It is a technique that is used for discovering relationships between items of dataset.

### Clustering Algorithm

- 1) k-mean clustering algorithm
- 2) PCA - (Principal Component analysis)

### application

- 1) Dimension Reduction
- 2) Recommendation System
- 3) image & video compression

### Association Algorithm

- 1) frequent pattern (FP) growth algorithm.

### 3) Reinforcement

It is based on Feedback and Punishment

Algorithm → 1) Q learning

2) Deep-Q-Network

3) Policy Gradient Methods

4) Actor-Critic Method

5) Proximal Policy Optimization

6) Trust Region Policy Optimization

7) Monte Carlo Tree Search

Application → 1) Gaming

2) Autonomous Vehicles

3) Robotics

4) Natural language processing.

## Supervised

## Unsupervised

- |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |
|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <ul style="list-style-type: none"><li>1) Labeled Data.</li><li>2) It is used for prediction (Predict Output).</li><li>3) Input data is provided to the model along with output data.</li><li>4) It needs supervision to train the model.</li><li>5) Supervised is categorized into<ul style="list-style-type: none"><li>A) Regression</li><li>B) Classification</li></ul></li><li>6) It takes direct feedback to check if it is predicting correct output or not.</li><li>7) It includes various algorithms like linear, logistic, SVM, decision tree.</li></ul> | <ul style="list-style-type: none"><li>1) Unlabeled Data.</li><li>2) It is used for finding Hidden pattern in data.</li><li>3) Only input data is provided to the model.</li><li>4) It does not need any supervision to train the model.</li><li>5) Unsupervised is classified into<ul style="list-style-type: none"><li>A) Clustering</li><li>B) Association</li></ul></li><li>6) It does not take any feedback.</li><li>7) It includes various algorithms like k-mean clustering, PCA.</li></ul> |
|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|

# Regression

Different types of Regression are as →

1) Simple linear Regression

2) Multiple linear Regression

3) Polynomial linear Regression

① Simple linear Regression

This is used to predict the value of a variable based on the value of another variable (same definition of Regression)

The variable we want to predict → Dependant

all other variables → Independant

This is based on a best fit line drawn on x,y plot.

Eq of line

$$y = mx + c$$

slope

$m$  = slope

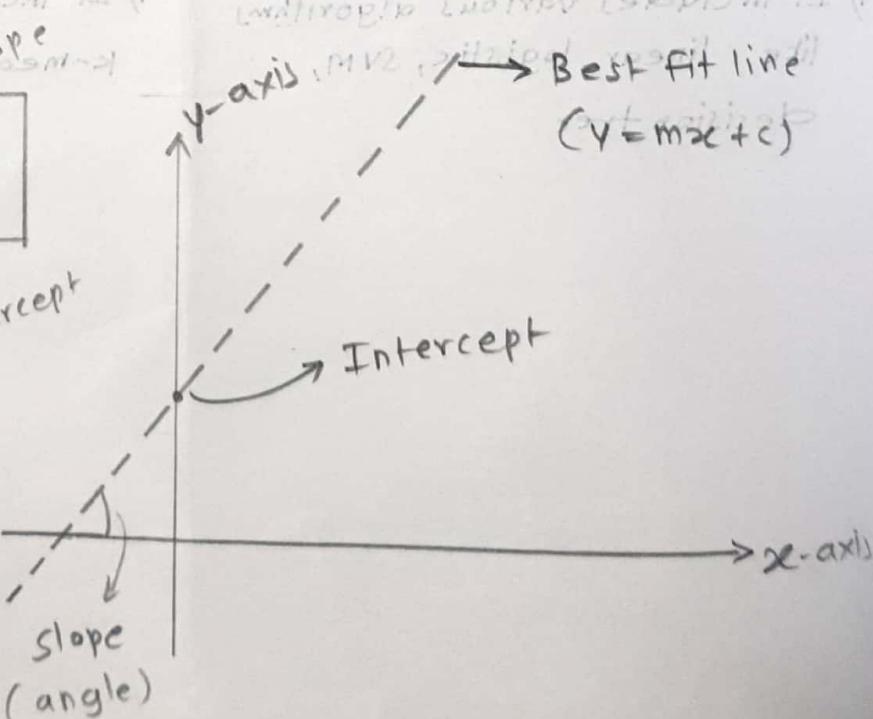
$c$  = intercept

$x$  = Independant

$y$  = dependant.

Intercept

Slope  
(angle)



Intercept → on y axis which point it touch that point is Intercept point.

Slope → The angle created by that line with x-axis is

$$\text{Slope} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\text{Intercept} = \bar{y} - m\bar{x}$$

Relation in linear Regression is as

Independent ↑ then dependent ↑

Independent ↓ then dependent ↓

Two types in  
linear Regression

simple linear Regression

Multiple linear Regression

1) simple linear → one Independent & dependent

2) Multiple linear → 1 Independent multiple dependent.

We need to calculate the Best fit line with line eqn

Best fit line → The line that has least error which means the error b/w predicted value & actual value should be minimum.

Residual (error) → The difference b/w observed value & predicted value is called Residual.

$$\text{Residual} = Y_{\text{pred}} - Y_i$$

$Y_{\text{pred}}$  = Predicted Value

$\bar{Y}$  = Observed

$Y_i$  = Actual Value

$Y_{\text{pred}}$  in simple linear Regression

$$Y_{\text{pred}} = \theta_0 + \theta_1 x$$

↑ Intercept      ↓ Independent Variable  
↓ Intercept       $\theta_0$  = Intercept  
                         $\theta_1$  = Slope  
                         $x$  = Independent

$Y_{\text{pred}}$  in Multiple linear Regression

$$Y_{\text{pred}} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

↑ Intercept       $\theta_0$  = Intercept  
                         $\theta_1, \theta_2, \theta_3, \dots, \theta_n$  = Slopes

$Y_{\text{pred}}$  in Polynomial Regression

$$Y_{\text{pred}} = \theta_0 + \theta_1 x_1 + \theta_2 x^2 + \theta_3 x^3$$

## Cost function

The function that we are trying to minimize it helps to find out the optimal values of  $\theta_1$  &  $\theta_0$ .

We are using three different cost function in linear Regression.

1) Mean Squared Error

2) R square

3) Root mean Squared Error

4) Residual standard Error.

1) Mean Squared Error (MSE)

$$MSE = \frac{1}{N} \sum (y_{pred} - y_i)^2$$

2) Root mean squared Error (RMSE)

$$RMSE = \sqrt{MSE} = \sqrt{\frac{RSS}{n}} = \sqrt{\frac{\sum (y_i - y_{pred})^2}{n}}$$

It specifies the absolute fit of the model to the data that is How close the observed value points are to the predicted values.

5) Mean Absolute Error (MAE)

$$MAE = \frac{1}{N} \sum |y_{pred} - y_i|$$

### 3) R square ( $R^2$ )

It is the number that explain the amount of variation that is explained by the developed model.

- 2) always lies between 0 to 1
- 3) Higher the value of  $R^2$  that means that better the model fits for data.

$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

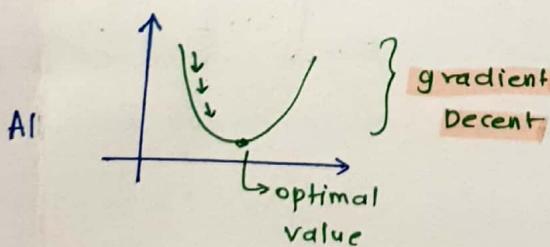
$RSS = \text{Residual Sum Square}$   
 $= \sum (y_i - \hat{y}_i)^2$

$TSS = \text{Total sum square}$   
 $= \sum (y_i - \bar{y})^2$

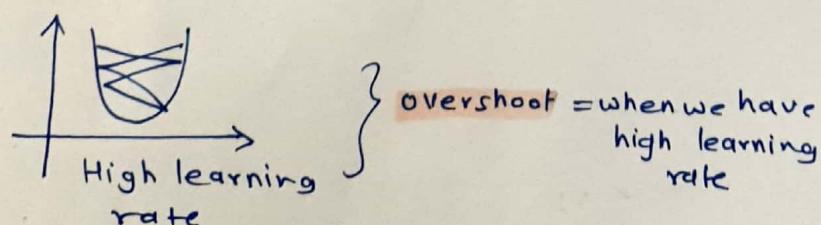
### 4) Residual Standard Error (RSE)

$$RSE = \sqrt{\frac{RSS}{n-2}} = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n-2}}$$

learning rate → How fast algorithm reach to final state  
 It decides how fast the algo converge to minimum value.



In gradient descent  
 No of steps can be considered as learning rate



Overshoot = when we have high learning rate

# EDA

Exploratory Data Analysis

EDA → Exploratory Data Analysis.

The data we collect from different sources is not in good form so we need to clean that data and set all the parameters in one scale so our machine can understand that data.

1) Check All the data which is as per data type.

e.g. Column holding Numbers but data type is object

so like this situation need to be solve `df[col].astype(DataType)`

2) Check the Null values `df.isnull().sum()`

3) Check Duplicate Values `df.duplicated().sum()`

4) Find outliers and remove them (use boxplot)

→ find the index of that outlier and then drop it from table

end fence all are outlier Identify them from box plot and drop them

Count Plot → Categorical Data

Dist Plot → Numerical Data

5) Convert the categorical Data into Numerical Data with encoding.

6) Then check the distribution of data points if there skewness is right or left then use skewness and convert Data into Normal Skew. (Pandas)

7) If Unit of data/columns are different then use scaling to convert them in one scale  
(Sklearn)

## \* Encoding \*

A43

Convert All Categorical Data in Numerical Data.

from sklearn.preprocessing } library of tools for data wrangling  
Pd.get\_dummies } methods at basic level  
in wrangling with the data frame

1) **label encoding**

2) **get dummies**

1) **label encoding** → 1) from sklearn.preprocessing import

LabelEncoder

2) le = LabelEncoder()

3) le.fit\_transform(col)

(Nominal)  
Data

- 1) Convert Categorical Data into numerical data
- 2) Starts from 0 & gives the number upto the total categories.
- 3) red, green, blue three categories so encoding will be like red=0, green=1, blue=2
- 4) No new columns will be created.

2) **get dummies** → 1) Pd.get\_dummies(col, dtype=int)

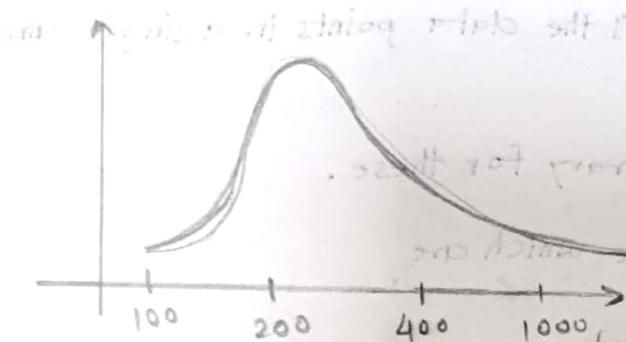
for every category there will be a new column which has only two values which are 0 & 1

Encoding is based on target variable

If target variable is categorical  $\rightarrow$  perform EDA before train test split.

If target variable is not categorical  $\rightarrow$  train test split.

### \* Skewness \*



In these data points have huge distance between them so the curve will be Not good or it could be left or Right skewed.

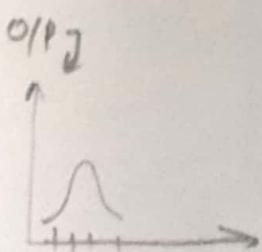
To solve this problem we perform skewness. So the distance between these points are minimized and a normal curved could be formed.

Two types  $\rightarrow$  1) `sqrt`  $\rightarrow$  `np.sqrt` 2) `log`  $\rightarrow$  `np.log` { present in Numpy library.

1) `sqrt`  $\rightarrow$  takes the square root of that value and place it on x-axis.

2) `log`  $\rightarrow$  takes the log of that value and place it on x-axis.

→ Done for every value present in the column.



## \* Scaling \*

- 1) Machine don't understand units
- 2) sometime  $1M = 100\text{ cm}$  this is not identified by machine and it assume both as different value
- 3) so to solve these problem we perform scaling
- 4) In these we will convert all the data points in a single same unit
- 5) we are using sklearn library for these.
- 6) There are two types in these which are

### 1) Normalization (min max)

### 2) Standardization

```
N → from sklearn.preprocessing import MinMaxScaler
```

```
S → from sklearn.preprocessing import StandardScaler
```

- 1) Normalization →
  - 1) adjust the values of data and fit it within a range
  - 2) range is most the time  $[0,1]$

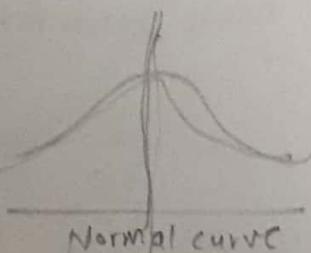
$$\text{Normalization} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

$x$  = original value

$x_{\min}$  = minimum value from dataset.

Median = Median = Mode

$x_{\max}$  = maximum value from dataset.



Normal curve

## 2) Standardization →

standard deviation = 1

probabilistic quantity related to mean and standard deviation

mean = 0      } standard Normal Distribution  
std = 1      } Happen

not been used like std diff to find standard deviation related ratio

Median = 0      } Median is a point which divides data into two equal parts  
diff      } Median = 0  
diff      } mode = 0

standard deviation is ratio of std diff to its mean for z score - 1

if in Z score sd below positive not possible else diff

else see negative diff to success probability becomes zero if diff

else diff is positive success probability becomes one if diff

else diff is negative success probability becomes zero if diff

else diff is positive success probability becomes one if diff

$\text{f1q2\_test\_diff} = \text{f1q1\_normal2\_lsbm\_areal2\_m01}$   
 $(\text{z0} = \text{std\_diff}, \text{v}, \text{x})$        $\text{f1q2\_test\_diff} = \text{f1q1\_normal2\_lsbm\_areal2\_m01}$   
 $(\text{z0} = \text{std\_diff}, \text{m01})$

## Train Test Split used in every model $\leftarrow \rightarrow 1500$ data point

- Now Data is clean or after getting clean data
- 1) we need to separate that Data into 2 parts which are independent & dependent.
  - 2) After these some part of these Data will be used for training on which we will build our model and other part will be used for testing.
  - 3) if these is not done and all the data is used for training then the accuracy for prediction could be 100% or it can face overfitting Because of that reason we will divide the data.

```
from sklearn.model_selection import train_test_split  
X-train, X-test, Y-train, Y-test = train_test_split(x, y, test_size=0.20,  
random_state=42)
```

## Cross Validation

500k data point

- 1) Some times data is not that huge or in small scale
- 2) So at that time we could not use train test split because some categories could be missed in this or some pattern will be undescoved
- 3) So at that time we will use k-fold validation

### process →

default  $k=3$

- 1) first take value of  $k$   
 $k=4$

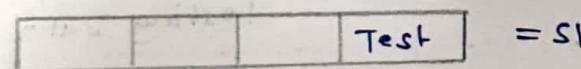
2) Now the hole dataset will be divided into 4 equal parts

- 3) Now from these 4 parts → 3 used for training  
1 used for testing

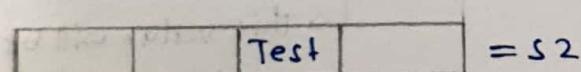
- 4) Step 3 will be repeated 4 time

Number of part = that much iteration.

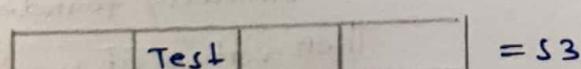
Sharing 1000 →  $(test-X)$  to 1000 =  $1000/4$



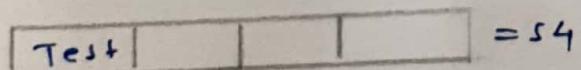
$$= S_1$$



$$= S_2$$



$$= S_3$$



$$= S_4$$



Take average of  
 $S_1, S_2, S_3, S_4$   
this

$(test, X-test)$  → error, loss

## Model Building

from sklearn.linear\_model import LinearRegression

Model = LinearRegression()

$$y = mx + c$$

All value will be found with this

→ 2 dimensional value needed.

Model.fit(x-train, y-train) → only takes the training data

To see { model.intercept\_ }

the values { model.coef\_ } These 2 value will be found with .fit

y-pred = Model.Predict(x-test) → only provide testing Data

$$y = mx + c$$

→ this value will be put and m & c are already founded

Then value of y will be calculated.

Model. ~~Score~~

Model.Score(x-test, y-test)

Model. Score (x-test, y-test)

→ The testing score

Model. Score (x\_train, y\_train)

→ The training score  
(Not that needed)

## find Error

MSE }

```
from sklearn.metrics import mean_squared_error
```

Mean-Square-Error ( $y_{\text{test}}, y_{\text{pred}}$ ) → The error is found

actual  $\downarrow$  predicted  $\downarrow$  actual value  $\downarrow$  predicted value

$$MSE = \text{mean\_square\_error}(y_{\text{test}}, y_{\text{pred}}) \quad \quad \quad RMSE$$

np.  $\sqrt{mSE}$

```
from sklearn.metrics import r2_score
```

r2-score (Y-test, Y-Pred)

$R^2$  score

## \*Rules Need to be follow while creating a ML Model \*

### (for linear Regression)

- 1) There must be linear Relationship bet<sup>n</sup> independent & dependant variable.
- 2) All variable are normally distributed.
- 3) There must be no outliers.
- 4) The mean of distribution of error is zero.
- 5) The variance of error is constant across all the level of independent variables. (Homocidity)
- 6) The distribution of error is normal & all errors are independent. (error-X, test-Y) good - 100%

### Corelation →

lies between 0 to 1

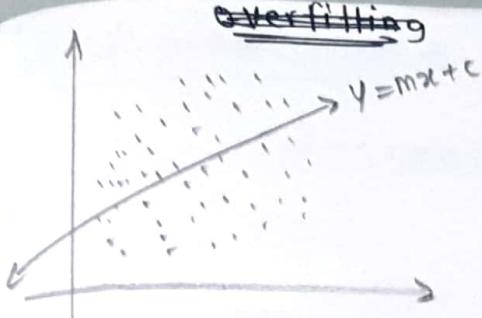
1 is high

feature selection → find corelation bet<sup>n</sup> independent if any independent variable has more correlation with other independent variable then remove

independent ← corelation → Independent

↳ then remove any one less imp.

dependent ← corelation → Independent  
↳ good.



In this situation Underfitting

MSE is very high

$R^2$  score is very low

Because of that reason Underfitting is occurred.

Main Reason  $\rightarrow$  low complexity

training score is low

Overfitting  
Underfitting

In this situation

MSE is very low  $\rightarrow$  fitted to every point

$R^2$  score is very high

Because of that Reason Overfitting is occurred

Main Reason  $\rightarrow$  High complexity

test score > training score

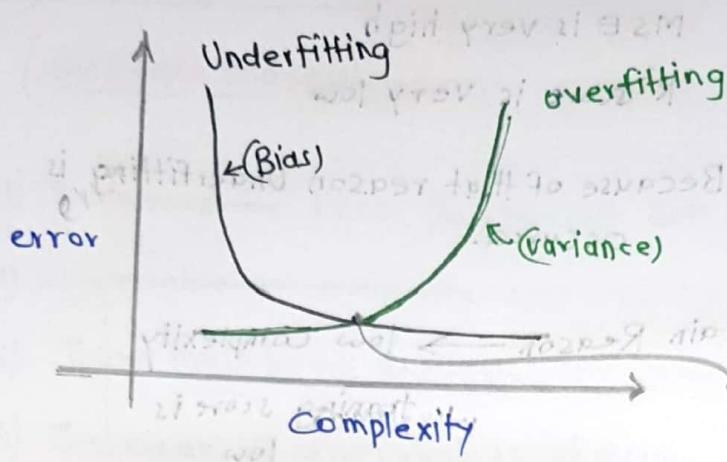
Complexity

$$\text{complexity} = \frac{\text{No of features}}{\text{No of independent col}}$$

Complexity  $\uparrow$  error  $\downarrow$  (Underfitting)

Complexity  $\downarrow$  error  $\uparrow$  (overfitting)

## \* Trade Off (Trade Variance) \*



Need to find the features which has low complexity

As we increase complexity the error is decreased.

- 1) If we increase the complexity a lot then the error will again start gradually increase.
- 2) The error in underfitting situation is called **Bias** as we increase the complexity the Bias will be decrease.
- 3) In overfitting situation the error start increase then that is known as **Variance**.
- 4) If we increase complexity then Variance is also increased.

## \* Identify Underfitting \*

- 1) Take data
  - 2) Split it into training & testing data
  - 3) train ~~data~~ on training data model
  - 4) find training score on training data
  - 5) if score is not good then underfitting
- Training score vs. Model complexity
- $\Rightarrow 80 \rightarrow \text{avg}$
- $< 70 \rightarrow \text{bad}$

## \* Sol<sup>n</sup> Underfitting \*

(22/01) notesirup99 13 (1)

Increase the Number of features.

$$10|3x + (b_{reg} - \bar{y})^3 = 0.201$$

## \* Identify Overfitting \*

Same still step 4

5) find training score for training data

6) find testing score for testing data.

7) If training score > testing score

8) If true then overfitting.

## \* Sol<sup>n</sup> Overfitting \*

Decrease Number of feature

Add some extra error

Two methods present for these.

1) L1 Regularization (Lasso)

2) L2 Regularization (Ridge)

We are adding external error in training.

1) MSE + error then line shift toward the testing data & give us the best fit line

2) As we are adding very small amount of external error in training MSE which helps us to reduce the testing error.

external error + training MSE = Regularization

## 1) L1 Regularization (Lasso)

$$\text{Lasso} = \frac{\sum (y - y_{\text{pred}})^2}{n} + \lambda \sum |\theta|$$

$\theta$  = coefficient

$\lambda$  = hyperparameter to direct off the unwanted

It will minimize the coefficient of unwanted feature to 0

## 2) L2 Regularization (Ridge)

$$\text{Ridge} = \frac{\sum (y - y_{\text{pred}})^2}{n} + \lambda \sum \theta^2$$

Ridge, Lasso → add extra error

remove unwanted features.

$\lambda \rightarrow$  decide by developer  
it can not take negative

best value  $\lambda \rightarrow$  Ridge  $\lambda (\theta^2)$

best value  $\lambda \rightarrow$  Lasso  $= \lambda (|\theta|)$

Ridge  $\rightarrow$  approach near to zero  $\lambda = \text{low}$

Lasso  $\rightarrow$  Make direct 0  $\lambda = \text{high}$

always take testing data for Ridge & Lasso

Delete Unwanted feature + Add some error = overfitting Handling

Mostly use the Lasso  $\rightarrow$  Because it directly make column value 0 so directly deleted.

Bad model  $\rightarrow$  accuracy below or near 80%. then it is ok but not that good.

accuracy below 70%. then bad model.

## Commands

→ pros of your approach ← Ridge

Convert  
Categorical  
into  
Numerical

1) Label  
Encoding

edit file  
edit my

from sklearn.preprocessing import LabelEncoder  
le = LabelEncoder()  
le.fit\_transform (col)

→ model building block  
+ map

2) get

dummies

import pandas as pd  
pd.get\_dummies (\_\_\_\_\_, dtype=int)

Skewness

1) sqrt

b = np.sqrt (Numpy array)

b = pd.DataFrame (b) → convert it into df  
after that to put in  
next operation.

2) log

b = np.log (Numpy array)

b = pd.DataFrame (b) → convert

Scaling

1) Normalization  
(MinMax)

from sklearn.preprocessing import MinMaxScaler  
min\_max = MinMaxScaler ()  
min\_max.fit\_transform (col )

2) Standardization

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
scaler.fit_transform(X)
```

Data points distribution is std-Normal  
ie  $\text{mean} = 0$  &  $\text{std} = 1$

### Train Test Split

```
from sklearn.model_selection import train_test_split
X_train, X_test, Y_train, Y_test = train_test_split(X, y,
                                                    test_size=0.2, random_state=42)
```

data > 500

### Cross Validation

data < 500

```
from sklearn.model_selection import cross_val_score
```

$CV = \text{cross_val_score}(lr, X, y, cv=5)$

for which model

we are using obj of that model

↳ No of folds.

### Model Building and Prediction

#### 1) import Model

```
from sklearn.linear_model import LinearRegression
lr = LinearRegression()
```

→ for simple and Multiple regression

from sklearn.preprocessing import PolynomialFeatures

$Poly = \text{PolynomialFeatures(degree=2)}$

↳ for polynomial regression import these.

2). fit	$\text{lr.fit}(X\text{-train}, Y\text{-train})$	↳ independent data
slope (m)		
intercept(c)	$\text{poly\_x} = \text{poly.fit\_transform}(X)$	↳ what to put here.

3). predict	$y\text{-pred} = \text{lr.predict}(X\text{-test})$	
put all in formula & find Y $y = mx + c$	$y\text{-pred} = \text{lr.predict}(Poly\_x)$	

4). score	$\text{lr.score}(X\text{-test}, Y\text{-test}) \rightarrow \text{Testing score}$	↳ simple, multiple
only check for regression for classification use classification report.	$\text{lr.score}(Poly\_x, Y)$	↳ polynomial

Error	1) MSE	$\text{from sklearn.metrics import mean_squared_error}$
		$\text{mean_squared_error}(Y, Y\text{-pred})$

actual data point      predicted data point

2) MAE

```
from sklearn.metrics import mean_absolute_error
```

```
mean_absolute_error(y, y-pred)
```

modeling no better than a flat line

3) R<sup>2</sup>

```
from sklearn.metrics import r2_score
```

```
r2_score(y, y-pred)
```

modeling no better than a flat line

Regularization

(Ridge L2)  
(Lasso L1)

Ridge

(L2)

Lasso

(L1)

```
from sklearn.linear_model import Ridge
```

```
L2 = Ridge(alpha=2)
```

```
L2.fit(X-train, Y-train)
```

modeling no better than a flat line

```
from sklearn.linear_model import Lasso
```

```
L1 = Lasso(alpha=150)
```

```
L1.fit(X-train, Y-train)
```

modeling no better than a flat line

Classification Report

(only for classification)

Accuracy

Precision

Recall

F1 score

```
from sklearn.metrics import classification_report
```

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

# Classification

If targeted variable has discrete or continuous data then it is a classification problem.

e.g. patient has cancer or not

fraud detection

email spam

signature is real or fake.

Different Types of classification are →

1) k-nearest neighbour

2) Logistic Regression

3) Decision tree

4) k-Ne. Random forest

5) Support Vector Machine.

Matrix in classification

1) Accuracy

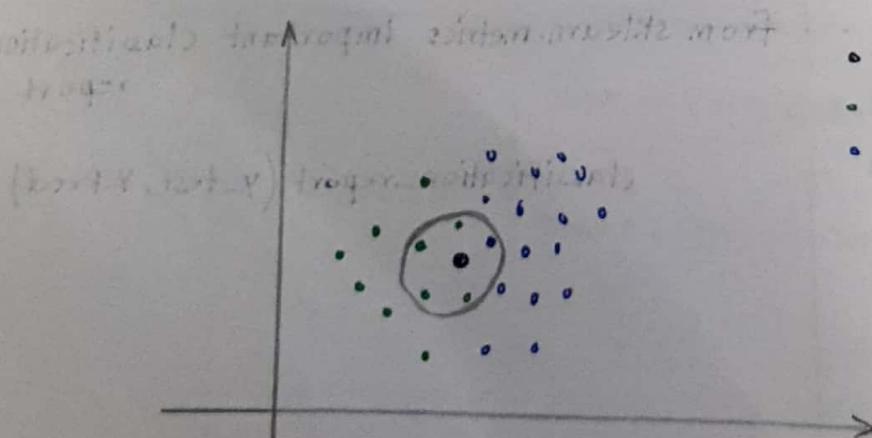
2) Recall

3) Precision

4) f1-score

1) K-nearest neighbour

It is based on distance.



- → we are predicting for
- → Yes
- → No

$k=5 \rightarrow$  other than the predicting one it choose 5 data points which are most near to it

Now Here 4 are yes and 1 is No  
so prediction will be Yes

default value of k = 5

\*Ways to find minimum distance \*

for the nearest 5 data points we can find distance through  
2 ways

1) Euclidean distance

2) Manhattan distance

1) Euclidean distance

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



It is simple  
distance

2) Manhattan Distance

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

$$\frac{9T}{4T + 9T} = 0.6666\overline{6}$$

## Confusion Matrix

		Actual	
		P	T
Predicted	1	TP	FP
	0	FN	TN

What does it mean? A small note

True and False notifications are

Predict Actual

P T

\* Significant confusion matrix of 2x2

agreement will be  $\frac{TP+TN}{P+T}$  and misclassification error is  $\frac{FP+FN}{P+T}$

so this is false

so that row is false negative

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

$(P - sY) + (sP - s^2)$

Sometime FN could be critical

At that time we can not use the accuracy

Because result of accuracy is always positive

So in these state we will use recall

man has cancer or  
Not

↳ so in this FN is  
critical

Because man has

cancer Harsh machine  
predicts it does not  
has cancer.

So for FN use  
recall

$$\text{Recall} = \frac{TP}{TP + FN}$$

$|P - sY| + |sP - s^2|$

IF recall is very high then model is good.

Sometime FP could be critical  
so at these time we will use precision

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

email is spam or not

↳ Here FP is critical

Because email was not spam but machine predict that it is spam.

So use precision

Sometime FP & FN both are critical  
so at these time we will use the f1 score.

f1 = Harmonic mean of Precision & recall

$$f1 = \frac{P \cdot R}{P + R}$$

Share market will crash or not

↳ it will not but machine predict it will

it will be crash but machine predict it will not

Both are harmful for us

- FN Critical → Recall
- FP Critical → Precision
- Both critical → f1 score

## KNN codes

perform all action of EDA and make data clean.

Use train test split

1) from sklearn.neighbors import KNeighborsClassifier  
knn = KNeighborsClassifier(n\_neighbors=5)

2) knn.fit(x\_train, y\_train)

3) y\_pred = knn.predict(x\_test)

4) from sklearn.metrics import classification\_report

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

tn correct ← recall

fp incorrect ← precision

fn correct ← 1 - recall

## 2) Logistic Regression

Probability =  $\frac{\text{No of Occurrence}}{\text{Total No of Value}}$

Odd ratio =  $\frac{\text{Number of Success}}{\text{Number of failure}}$

probability value is always positive

it always lies in 0 to 1

$$0 < \text{Probability} < 1$$

$$P = 0 \rightarrow \frac{1}{2} \rightarrow 1$$

$$\text{OR} = 0 \rightarrow 1 \rightarrow \infty$$

$$\log \text{odd} = -\infty \rightarrow 0 \rightarrow \infty$$

probability

odd ration

log odd

lies between

$$\text{prob success} = \frac{\text{No of success}}{N}$$

$$\text{Prob failure} = \frac{\text{No of failure}}{N}$$

from above formula we can find Number of success & failure

$$\text{No of Success} = \text{Prob of success} \times N$$

$$\text{No of failure} = \text{Prob of failure} \times N$$

let

$$\text{Prob success} = P$$

$$\text{Prob failure} = q$$

let odd ratio of them are as

$$\text{Odd ratio} = \frac{P \times P}{q \times q}$$

Total No  
of them

$$\therefore 1 = \frac{P}{q}$$

$$P + q = 1$$

$$q = 1 - P$$

$$P = q$$

$$\text{odd ratio} = \frac{P}{q} = \frac{P}{1-P}$$

$$\text{odd ratio} = \frac{P}{1-P}$$

$$\text{odd Ratio. } 0 \leq \frac{P}{1-P} \leq 1$$

$$\log(\text{odd Ratio}) \Rightarrow (-\infty, \infty)$$

$$\text{line Range} \Rightarrow (-\infty, \infty)$$

$$\therefore y = mx + c = \log\left(\frac{P}{1-P}\right)$$

$$\begin{aligned} Y &= mx + c \\ Y &= \alpha_0 + \alpha_1 x \end{aligned} \quad \left\{ \begin{array}{l} \text{Range is } -\infty \text{ to } \infty \\ \frac{Y_2 - Y_1}{Y_2 + 1} = q \end{array} \right.$$

odd ratio  $\rightarrow$  Range is  $-\infty$  to  $\infty$

$\therefore$  So we can say

line of eq<sup>n</sup> = odd ratio

$$Y = \log \left( \frac{P}{1-P} \right)$$

$$e^Y = e^{\log_e \left( \frac{P}{1-P} \right)} \quad \frac{1}{1+P} = q$$

$$\frac{e^Y}{1} = \frac{P}{1-P}$$

$$\left( \frac{1}{e^Y} - 1 \right) = \frac{1-P}{P} \quad \text{not diff}$$

$$\frac{1}{e^Y} = \frac{1}{P} - \frac{P}{1-P} \quad \text{to solve for } P$$

(a)

$$\frac{1}{e^Y} = \frac{1}{P} - 1 \quad \text{not diff}$$

$$\frac{1}{e^Y} + 1 = \frac{1}{P}$$

$$\frac{1}{P} = \frac{1+e^Y}{e^Y} \quad \text{not diff}$$

$$P = \frac{e^y}{1+e^y}$$

$$P = \frac{e^y/e^y}{1+e^y} = \frac{1}{1+e^{-y}}$$

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

$$\left(\frac{e^y}{e^y + 1}\right) e^y = P$$

$$P = \boxed{\frac{1}{1+e^{-y}}}$$

→ logic function  
→ sigmoid function

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

Cost function =  $y * \log(y_{\text{Pred}}) - (1-y) * \log(1-y_{\text{Pred}})$

So value of  $y$  can be 0/1

$y$  replaced in cost fun with 0/1

$$0 \rightarrow -\log(1-y_{\text{Pred}}) \rightarrow y \rightarrow 0$$

Substituting in the formula.

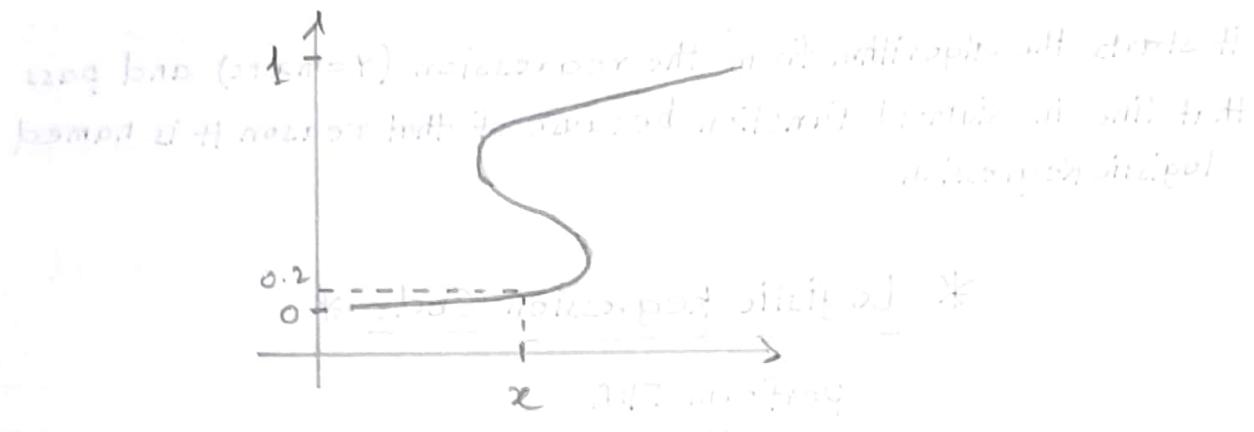
Cost

$$0 \rightarrow 0 \rightarrow 0$$

$$1 \rightarrow -y \log(y_{\text{Pred}}) \rightarrow y \rightarrow 1$$

$$1 \rightarrow 0 \rightarrow \infty$$

## \* Threshold \*



So here we are predicting  $x_{\text{category}}$   
 ~~$x$  is 0.2~~ and it lies between 0.2

But our model only knows 0 & 1

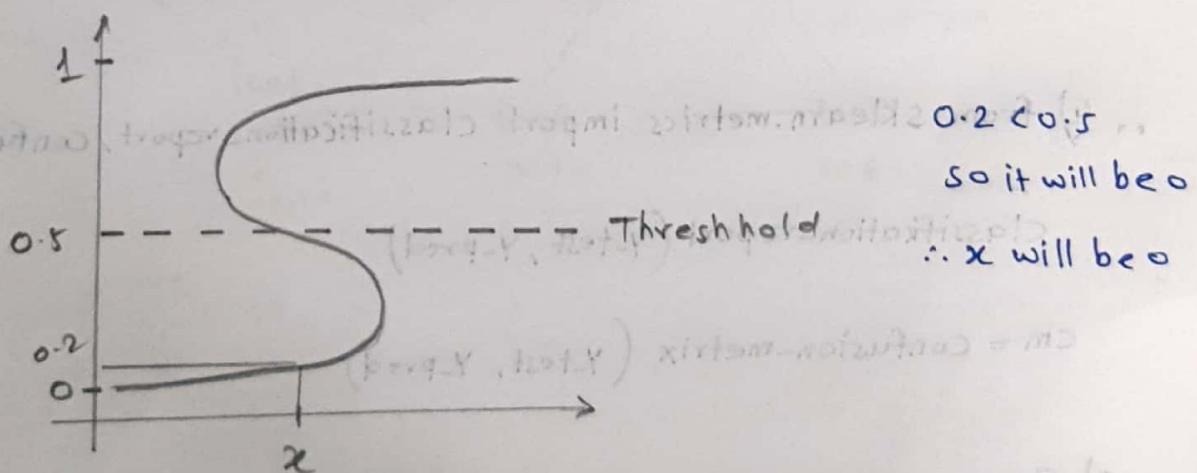
So for that we will create a threshold  
 if it lies at 0.5  $\rightarrow$  default

all the values below 0.5 threshold are 0

all the values above 0.5 threshold are 1

$$\boxed{\text{default threshold} = 0.5}$$

$$f(x) = \frac{1}{1 + e^{-x}}, \text{ if } f(x) < 0.5 \text{ then } y = 0, \text{ else } y = 1$$



$$0.2 < 0.5$$

so it will be 0

$\therefore x$  will be 0

## Why Regression Name is Present in logistic Regression

it starts the algorithm from the regression ( $y = mx + c$ ) and pass that line in sigmoid function because of that reason it is named logistic Regression.

### \* Logistic Regression Code \*

Perform EDA

Scaling

Skewness

Encoding

train test split

1) from sklearn.linear\_model import LogisticRegression

model = LogisticRegression()

2) model.fit(x-train, Y-train)

3) Y-pred = model.predict(x-test)

model.predict\_proba(x-test)[:, 1] → To check all probability score

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

classification\_report(Y-test, Y-pred)

cm = confusion\_matrix(Y-test, Y-pred)

5) from sklearn.metrics import recall\_score, precision\_score  
recall\_score(Y-test, Y-pred)

precision\_score(Y-test, Y-pred)

### 3) Decision Tree

open farming → open fields with red soil  
open fields

It is a supervised machine learning algorithm used for both classification & regression Problem

But Most of time we use for classification.

In decision tree there are mainly two nodes.

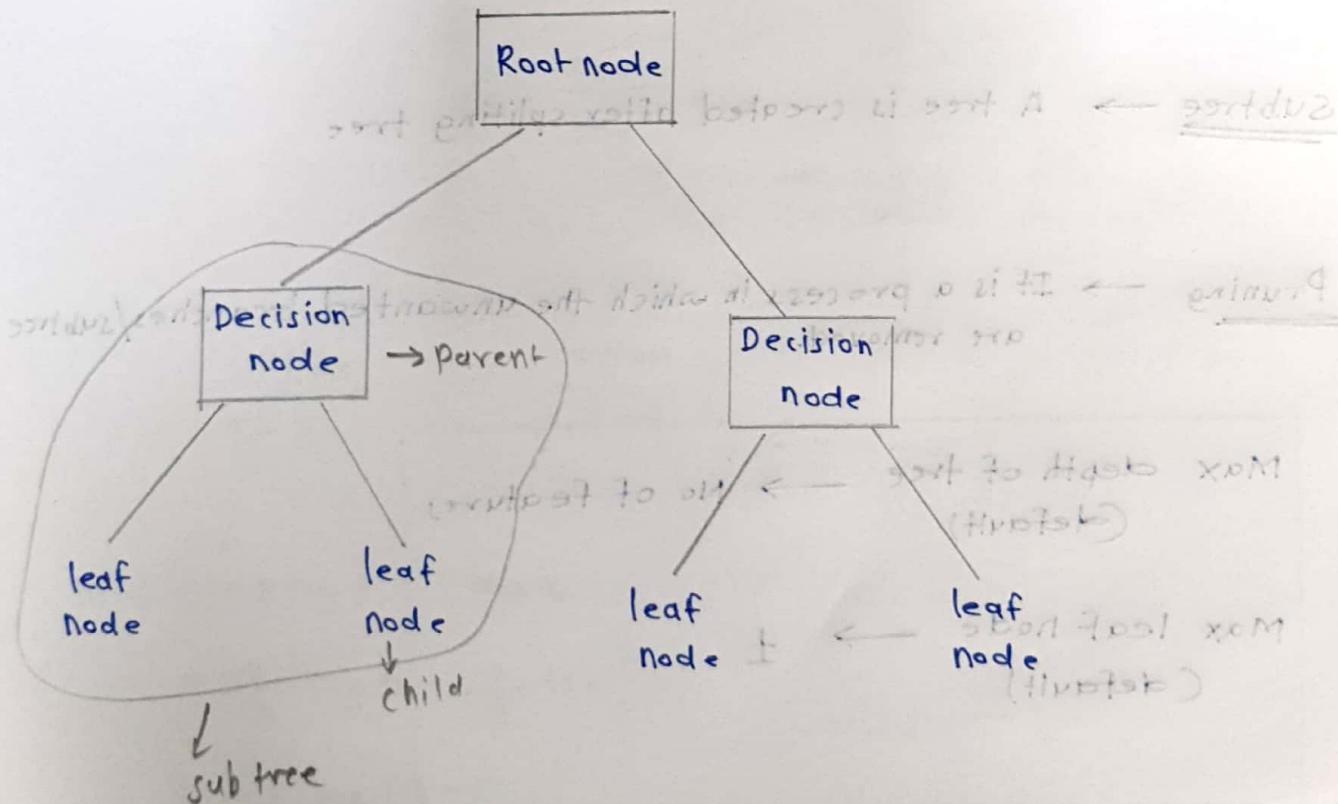
1) leaf node

2) decision node

decision node are used to make any decision and it has many branches.

leaf node is output of those decision & does not has further branches.

Decision node → Parent



There are two main node → parent node,  
child node.

- 1) Root node is called as, parent node
- 2) Decision tree start from that root node
- 3) Root node represent the entire data set.
- 4) which further get divide in two or more set.
- 5) other than root node all other are child node
- 6) leaf node are final output of the node

Splitting → It is a process to divide the root node or decision node into sub nodes according to the given condition.

Subtree → A tree is created after splitting tree

Pruning → It is a process in which the unwanted branches/subtree are removed.

Max depth of tree → No of features  
(default)

Max leaf node → 1  
(default)

foot  
shoe

foot  
shoe

foot  
shoe

foot  
shoe

Gini impurity and entropy Both are used to check impurity of Decision tree.

### \* GINI Impurity \*

$$\text{Gini} = 1 - \sum(p_i^2)$$

$$= \left( \frac{1}{2} \right) e^{0.5} \left( \frac{1}{2} \right) + \left( \frac{1}{2} \right) e^{0.5} \left( \frac{1}{2} \right) = \frac{1}{2} \left( 1 + e^{-0.5} \right)$$

$$= 1 - \left[ \left( \frac{1}{2} \right)^2 + \left( \frac{1}{2} \right)^2 \right] = \frac{1}{2}$$

Root for a 2-class tree example  $\leftarrow t = 100$

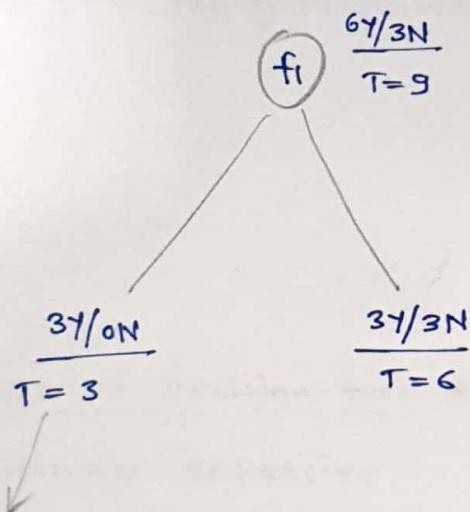
### \* Entropy \*

$$-p \log p - q \log q$$

$p = \text{positive probability}$

$q = \text{Neg}$

$$\left[ \left( \frac{1}{2} \right) + \left( \frac{1}{2} \right) \right] - 1 = 0$$



$$H(c_1) = -p \log p - q \log q$$

$$= -\frac{3}{10} \log \frac{3}{10} - \frac{7}{10} \log \frac{7}{10}$$

$$= -1 \log (1) - 0$$

$H(c_1) = 0 \rightarrow$  It is pure split  $\leftarrow 2 \cdot 0 = 0$   
 $\therefore$  it is a leaf node

$$H(c_2) = \frac{-3}{6} \log\left(\frac{3}{6}\right) - \left(\frac{3}{6}\right) \log\left(\frac{3}{6}\right)$$

$$= -\frac{1}{2} \log\left(\frac{1}{2}\right) - \left(\frac{1}{2}\right) \log\left(\frac{1}{2}\right)$$

$$= -\log\left(\frac{1}{2}\right)$$

$$= \log_2(2)$$

$H(c_2) = 1 \rightarrow$  Impure split so it is not leaf.



$$GI(c_1) = 1 - \sum(p^2)$$

$$= 1 - \left[ \left(\frac{3}{3}\right)^2 + \left(\frac{0}{3}\right)^2 \right]$$

$$= 1 - [1+0]$$

$$= 1 - 1$$

$(c_1) = 0 \rightarrow$  pure split

$$GI(c_2) = 1 - \sum(p^2)$$

$$= 1 - \left[ \left(\frac{3}{6}\right)^2 + \left(\frac{3}{6}\right)^2 \right]$$

$$= 1 - \left[ \frac{1}{4} + \frac{1}{4} \right]$$

$$= 1 - \frac{1}{2}$$

$$= \frac{1}{2}$$

$(c_2) = 0.5 \rightarrow$  Impure split if  $q < 0.5 = 0.5$

Gini impurity lies between 0 to 0.5  
Entropy lies between 0 to 1

All the bad features of tree are present at the end  
(overfitting)

To remove that features we use pruning.

Decision tree has more chances of overfitting.

### \* Advantage of Decision Tree

- 1) able to classify non linear separable data.
- 2) automatically select good feature.
- 3) less pre processing required (only need to do )

Table Encoding + Missing Values

### \* Disadvantage of Decision tree \*

- 1) Computationally expensive
- 2) High Variance (High chance of overfitting)

### \* Imbalanced Data \*

targeted Data has diff values and if these attribute data has imbalanced values then we are using the

- 1) Under
- 2) Over → random over sampling.

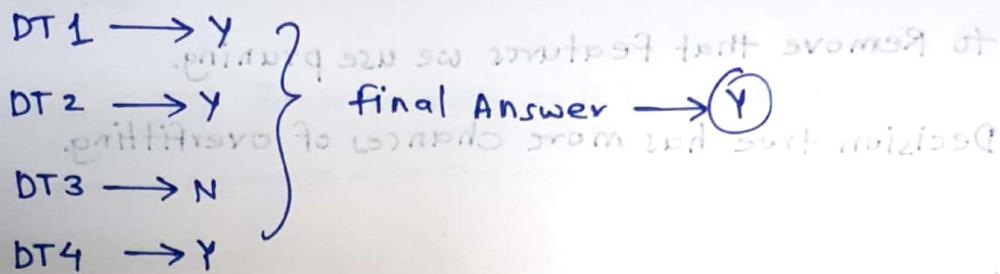
## \* Random forest \*

2.0 of a 4th soil type  
soil

Plot of a scattered soil types

Multiple Decision Trees are trained and after training output of them are found

After that Max Number occurrence output is taken.



Decision Tree & Random forest are not good for

Multi Imbalanced Data

(imbalance ratio) brings following error (e.g.)

minority + error  
majority error

\* Impacted Data \*

Plot showing with time taken this and also better not with grain size as well as by location and

error (

grain size metric → error (

## Boosting

\* follow sequence of values

There are three types in these

$\Rightarrow$  Ada Boosting = follow with weight

1) Ada Boosting

(start with different weight - 1) add = weight

2) XG Boosting

3) Gradient Boost.

1) Ada Boosting

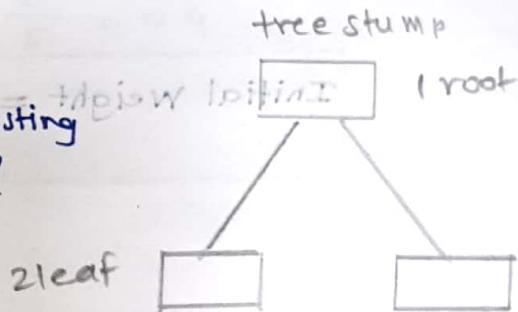
adaboosting grows to one adaboosting

1) full form  $\rightarrow$  Adaptive

2) we use tree stump in Ada Boosting

3) tree stump has one root node & one leaf node

4) Max depth of tree stump is 1



1) The next model try to fix the mistake done by the previous model this is called sequential training.

2) we just increase the weight of wrong prediction of the first model.

parallel training  $\rightarrow$  random forest

sequential training  $\rightarrow$  Ada Boosting

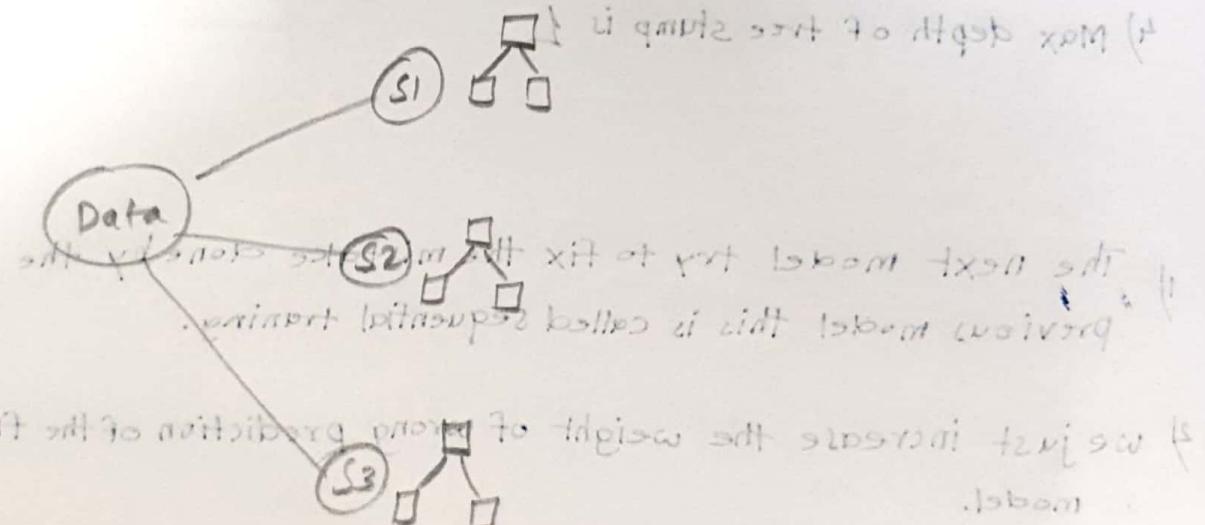
## \* Way to increase weight \*

$$\text{New Weight} = \text{old weight} \times e^{\text{Stage}}$$

$$\text{Stage} = \ln \left( \frac{1 - \text{Missclassification rate}}{\text{Missclassification rate}} \right)$$

$$\text{Misclassification rate} = \frac{\text{No of wrong prediction}}{\text{Total prediction}}$$

$$\text{Initial weight} = \frac{1}{\text{Total data points}}$$



format mobno <-- extract followed  
extract mobno <-- extract followed

	$X_1$	$X_2$	$X_3$	$Y$	$w$	pred	New $w$	final
1					0.2	✓	0.29	0.16
2					0.2	✓	0.29	0.16
3					0.2	✓	0.29	0.16
4					0.2	✗	0.29	0.24
5					0.2	✗	0.29	0.24
					1.0		1.81	0.96
								$\rightarrow 4 \text{ pt missing}$

$$\text{Initial weight} = \frac{1.0}{5} = 0.2$$

$$ps.0 = 1.0/5 = 0.2$$

$$21.0 = 1.0 \leftarrow 5.0$$

$$\text{Initial weight} = 0.2$$

$$\text{Misclassification rate} = \frac{\text{wrong}}{\text{total}} = \frac{2}{5} = 0.4$$

initially 4x5N  $\leftarrow$  4x5B  $\rightarrow$  2x5T  $\leftarrow$  2x5B  $\rightarrow$  1x5T  $\leftarrow$  1x5B

new 1x5B  $\leftarrow$  1x5T  $\leftarrow$  1x5B  $\leftarrow$  1x5B

1x5B  $\leftarrow$  1x5B  $\leftarrow$  1x5B

$$M - Y = 0.4$$

initially  $\leftarrow$  0.4

initially

initially  $\leftarrow$  0.4

$$\text{stage} = \ln\left(\frac{1 - M_Y}{M_Y}\right)$$

$$= \ln\left(\frac{1 - 0.4}{0.4}\right) = \ln\left(\frac{0.6}{0.4}\right)$$

initially  $\leftarrow$  0.4

initially  $\leftarrow$  0.4

$$\text{stage} = 0.40$$

initially  $\leftarrow$  0.4

initially  $\leftarrow$  0.4

initially  $\leftarrow$  0.4

initially  $\leftarrow$  0.4

$$\text{New weight} = \text{old weight} \times e^{\text{stage}}$$

$$= 0.2 \times e^{0.40}$$

initially  $\leftarrow$  0.4  $\leftarrow$  0.2  $\leftarrow$  0.2 = 0.25  $\leftarrow$  0.25

$$= 0.29 \quad \text{--- only place in wrong prediction}$$

$$\begin{aligned}
 \text{final} &= \frac{\text{New weight}}{\text{Total new weight}} \\
 &= \frac{0.29}{1.18} = \frac{0.24}{1.18} \\
 0.29 \rightarrow \text{final} &= 0.24 \\
 0.2 \rightarrow \text{final} &= 0.16
 \end{aligned}$$

$0.2 = \text{Adaboost IHT}$

$0.16 = \text{Adaboost IIHT}$

Gradient Boosting $\rightarrow$ full tree Ada Boosting $\rightarrow$ tree stump	Gradient Boosting $\rightarrow$ Next model train on residual given by previous model. Ada Boosting $\rightarrow$ Next model try to fix the mistake done by previous model.
------------------------------------------------------------------------------------	-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

Random forest $\rightarrow$ parallel training Ada Boosting $\rightarrow$ sequential training	Random forest $\rightarrow$ tree Ada Boosting $\rightarrow$ tree stump
-------------------------------------------------------------------------------------------------	---------------------------------------------------------------------------

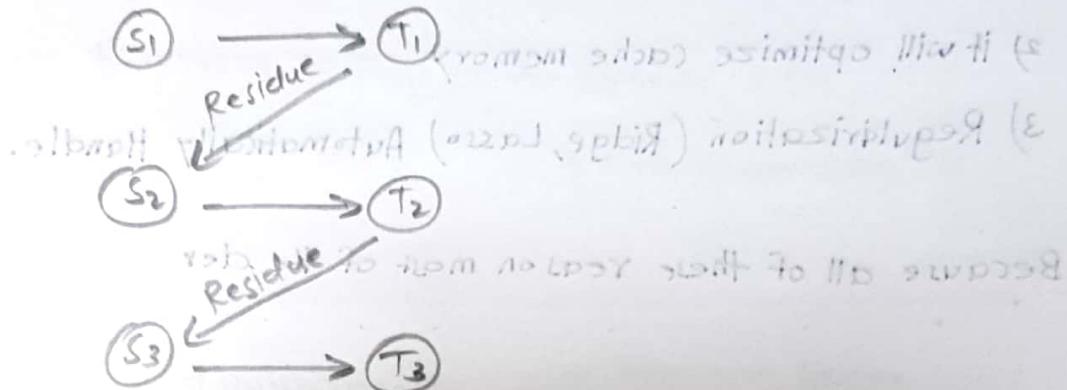
$n\text{-estimator} = 50 \rightarrow$  No of stump you are using in Ada Boosting.

Yes

## 2) Gradient Boosting

Next model train on residue of the first model.

At the end it tries to convert it into but not direct zero  
learning rate will not be able to do that



$$\hat{y}_2 = y + \text{learning rate} \times \text{residue}$$

$y_t$	$\hat{y}_1$	residue 1	$\hat{y}_2$	residue 2
170	166	4	170.4	0.4
171	166	5	171.5	0.5
162	166	-4	161.6	-0.4
160	166	-6	154.4	-0.6

$$170 + 0.1 \times 4$$

default tree = 100  
default learning rate = 0.1

### 3) XG Boost

## Extreme Gradient Boosting

This Boosting provides

- 1) Parallelization
  - 2) it will optimize cache memory
  - 3) Regularization (Ridge, Lasso) Automatically Handle.

Because all of these reasons most of the day

00) = ent theſe

1-0 = 0.000001 Hz/s

## SVM

$$y = w^T x + b \rightarrow \text{pos class}$$

- 1) It is a supervised learning Algorithm and it is used for classification
- 2) classification of two group.

There are two types in these → 1) linear SVM  
2) Non linear SVM

We found boundaries to separate two class

- line → linear
- circle → Non linear  
(other)

### linear SVM

If data is in 2 dimension and we can separate that Data using single straight line  
(The data points can be classified using a single straight line)

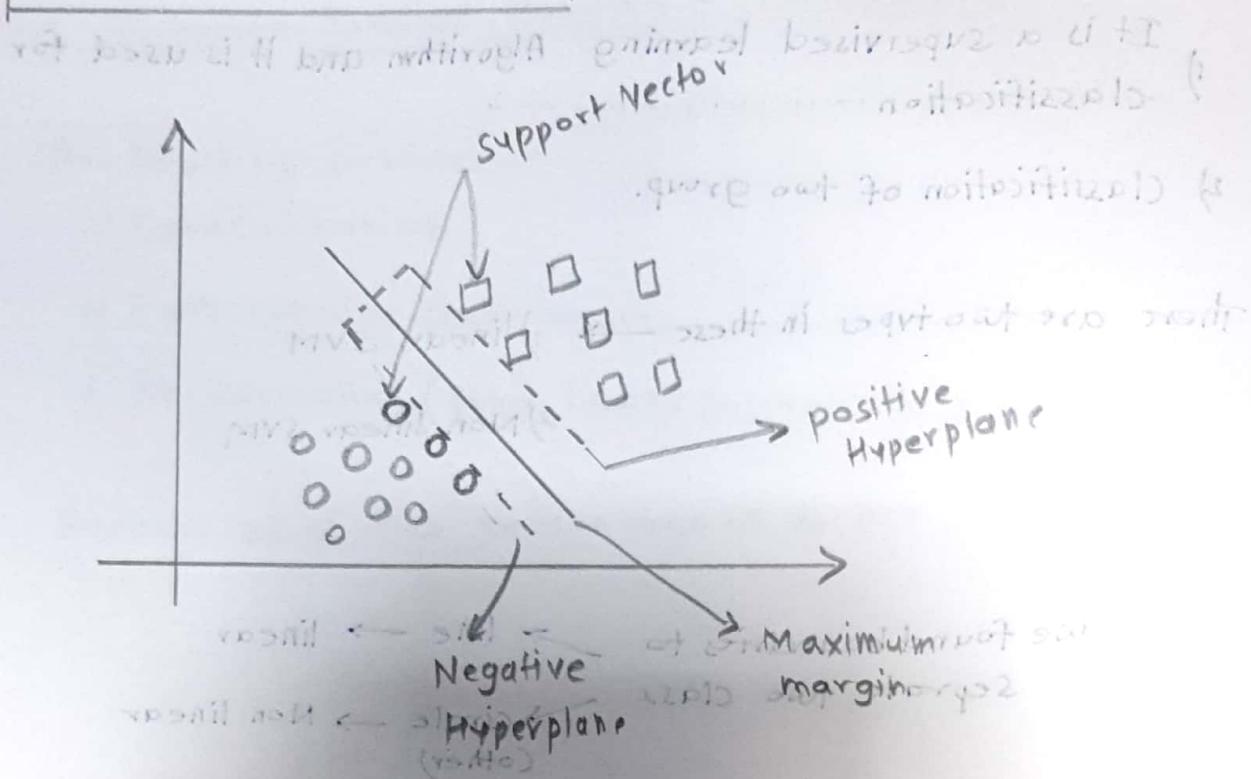
### Non linear SVM

If the Data is in 2 dimension and we can not separate that Data point with single straight line

for that we use some advance techniques like kernel tricks to classify them

$$\text{circle eq}^n = x^2 + y^2 = c^2$$

MV2



Margin  $\rightarrow$  distance between Hyperplane & support Vector

Support Vector  $\rightarrow$  The points that are closest to the hyperplane  
 There could be any number of support vector  
 It depends upon the data.

If graph 2D  $\rightarrow$  decision boundary is

straight line

If graph has more dimension  $\rightarrow$  decision boundary is

Hyperplane

Best Hyperplane  $\rightarrow$  Maximum distance from both the class.

$$\text{eqn of plane} = ax + by + cz = d$$

Kernal  $\rightarrow$  Kernal is used for transforming higher dimension

from Kernal we can transform complex data into higher dimensional simpler form

Advantages of SVM

- \* Advantage of SVM \*
- 1) works better when data is linear
- 2) More effective in High dimension
- 3) with the help of kernal tricks we can solve any complex problem.
- 4) SVM is not sensitive to outliers.
- 5) Help us with image classification.

\* Disadvantage of SVM \*

- 1) choosing best kernal is not very simple/easy
- 2) doesn't show good result on big dataset.
- 3) SVM hyperparameter are cost-c & gamma. It is not that easy to fine tune these hyper parameter.

## KNN

```
from sklearn.neighbors import KNeighborsClassifier  
knn = KNeighborsClassifier(n_neighbors=5)
```

b = 3x + 2y + 3z = apply to ps  
n\_neighbors=5 default

## logistic Regression

```
from sklearn.linear_model import LogisticRegression  
model = LogisticRegression()
```

## Decision Tree

```
from sklearn.tree import DecisionTreeClassifier  
model = DecisionTreeClassifier(max_depth=5,  
max_leaf_nodes=5)
```

## Random forest

```
from sklearn.ensemble import RandomForestClassifier  
RF = RandomForestClassifier(n_estimators=100,  
max_depth=4,  
max_leaf_node=5)
```

## Ada Boosting

```
AdaBoostClassifier(n_estimators=_,  
random_state=_,  
learning_rate=1)
```

## Gradient Boosting

```
from sklearn.ensemble import GradientBoostingClassifier  
GB = GradientBoostingClassifier(n_estimators=_____,  
                                 max_depth=_____,  
                                 learning_rate=0.1)  
  
for Regression problem ← Regressor  
for all the Boosting tech ↓
```

## XG IMP

### Boost

Extreme gradient

pip install xgboost.

```
from xgboost import XGBClassifier
```

```
XG = XGBClassifier(random_state=1, max_depth=_____,  
                    learning_rate=_____, n_estimators=_____)
```

## SVM

```
from sklearn.svm import SVC
```

```
SVM = SVC(kernel='rbf', gamma='scale')
```

## Random Over Sampling

Main focus is on Y variable always  
Y sampled then automatically other will be also sampled.

```
from imblearn.over_sampling import RandomOverSampler
```

```
rs = RandomOverSampler()
```

```
rs.fit_resample( )
```

## Random Under Sampling

```
from imblearn.under_sampling import RandomUnderSampler
```

```
rs = RandomUnderSampler(random_state=42)
```