

# Machine Learning.

Date: \_\_\_\_\_

Page No.: \_\_\_\_\_

## Day 1. Introduction to Machine Learning.

### - Supervised Learning

- Linear regression, logistic regression, ridge and lasso's regression, Decision Tree Making algorithm, Random Forest Algorithm
- K-nearest neighbour., Support Vector Machine
- Naive Bayes, boosting ....

### - Unsupervised learning

- kmeans, PCA, Clustering, DBScan.

### - Deep Learning -

- Time series
- Natural Language Processing (NLP)
- Neural Networks (ANN, CNN, etc)

## Day 2.

### \* Numpy Library. :- Numerical Python.

1) It is a Library consisting of multidimensional array objects and a collection of routines for processing those array.

### ① Identity Matrix.

→ diagonal matrix will be one and rest everything will be zero.

## Day 3. Numpy Part 2.

1) Indexing And slicing:

eg:-  $a = \text{np.array}([[1, 2, 3], [4, 5, 6]])$

a

$a[0,0]$  # 1<sup>st</sup> Row, 2<sup>nd</sup> Column

2) Slicing of 2-D Array:

eg:-  $a = \text{np.array}([[1, 2, 3], [4, 5, 6], [7, 8, 9]])$

a

$a[1:2, 1:2]$  # 1<sup>st</sup> Row, 2<sup>nd</sup> Column

3) Append.

eg:-  $a = \text{np.array}([[1, 2, 3], [3, 4]])$

a

$b = \text{np.append}(a, [[5, 6]], axis=0)$

b.

4) np.mean(a)

5) np.median(a)

6) np.std(a).

## Day 4 Pandas.

1) Used For data Manipulation and data Analysis.

2) Pandas is build on the top of the numpy package.

3) Panel Data ----> Pandas

4) In Pandas we have 2 data structure ---->

i) Series ii) DataFrame.

## # Series

it has a sequence of data values.  
eg:- pd.Series([1, 2, 3, 4, 5])

## # DataFrame

A DataFrame is nothing but a table.

eg:- pd.DataFrame({'Yes': [50, 21, 56], 'No': [101, 20, 11]})

eg:- pwd() # For File path

eg:- df = pd.read\_csv(x"Data.csv").

df.

df.columns

df.shape

df.head()

df.tail()

df.info()

df.isnull()

df.isnull().sum()

## DAY 5 Pandas Part 2.

## # Drop.

eg:- df.drop('c', axis=1, inplace=True).

df

## #

df['math'].max()

min()

mean()

median()

penkraft

## value\_count.

eg:- df.name.value\_counts().

# unique & nunique.

df['result'].nunique() #count.

df['result'].unique() #exact value.

# From numpy.random import randn as rn.

rn(5,3).

# Seed

np.random.seed(1)

r = rn(5,3)

r

ndf = pd.DataFrame(r)

ndf.

# loc() & iloc()

df.loc[[2,4],['name','result']] #Value.

→ select by label

df.iloc[[2,4],[0,3]] # index.

→ select by index or position.

ndf.iloc[0:3,0:2]

df [df > .5].

# Fillna()

df = df.fillna(5).

df.

df = df.fillna('good')

# Concat()

df1 = pd.DataFrame({ 'A': ['A0', 'A1', 'A2'], 'B': ['B0', 'B1', 'B2'], 'C': ['C0', 'C1', 'C2'] }).

df1

df2 = pd.DataFrame({ 'A': [ 'A3', 'A4', 'A5'], 'B': [ 'B3', 'B4', 'B5'], 'C': [ 'C3', 'C4', 'C5'] }).

Day 6. Pandas.

# Merging.

eg: left = pd.DataFrame({ 'key': ['k0', 'k1', 'k2'], 'A': ['A1', 'A2', 'A3'], 'B': ['B1', 'B2', 'B3'] }).

left

right = pd.DataFrame({ 'key': ['k0', 'k1', 'k2'], 'C': ['C1', 'C2', 'C3'], 'D': ['D1', 'D2', 'D3'] }).

right

eg: left

eg: right.

eg: inner = pd.merge(left, right, how='inner', on='key')

inner

eg: leftjoin = pd.merge(left, right, how='left', on='key')

leftjoin.

eg:- rightjoin = pd.merge(left, right, how='right', on='key')

rightjoin.

## # groupby

eg:- data = {'Company': ['GOOG', 'GOOG', 'MSFT', 'MSFT', 'FB', 'FB'],  
'Person': ['SAM', 'CHARLIE', 'AMY', 'VENSSA', 'CARL', 'SARAH'],  
'SALES': [200, 120, 340, 124, 243, 350]}

df = pd.DataFrame(data)

df

eg:- bycomp = df.groupby('Company')

bycomp

eg:- bycomp.mean()

bycomp.median()

bycomp.std()

## # Sorting

eg:- df = pd.DataFrame(ages, columns=['Country', 'Year', 'Population', 'Continent'])

df

eg df.sort\_values(by=['Country'])

eg df.sort\_values(by=['Population'], ascending=False).

eg: df.sort\_index(axis=1)

# Banking Credit Risk Analysis.

Date: \_\_\_\_\_  
Page No.: \_\_\_\_\_

## Day 7. Data cleaning.

eg:- import numpy as np.  
import pandas as pd.  
import matplotlib.pyplot as plt.  
import seaborn as sns

eg: cr = pd.read\_csv(r"CreditRisk.csv")

cr

eg: cr.head(2).

eg: cr.shape

eg: cr.info()

eg: cr.isnull().sum()

eg: cr.Education.value\_counts() / cr.shape[0] \* 100  
# For percentage

eg: cr.Education.unique()

eg: cr.Self\_Employed.value\_counts()

eg: cr.Self\_Employed = cr.Self\_Employed.fillna('no')

eg: cr.isnull().sum()

eg: cr.Self\_Employed = cr.Self\_Employed.replace({'no': 'No'})

eg: cr['New\_App\_Inc'] = 2 \* cr.ApplicantIncome

eg: cr.shape

eg: cr.head(2)

eg: cr.LoanAmount.value\_counts()

eg: cr.LoanAmount = cr.LoanAmount.fillna(cr.LoanAmount.mean())

# Replace by mean.

eg: cr.isnull().sum()

eg: cr.iloc[[1, 100, 500, 600, 800], [1, 6]]

penkraft

## Day 8. What is Sampling And Why we do it?

- i) Sampling :- Sampling is the process of selecting a subset of data from a large dataset to use in ML model.
- ii) Studies with inadequate sample suffer from overfitting of data and have a lower probability of producing true effects.
- iii) The increment in sample size increases the accuracy of prediction but may not cause a significant change after a certain sample size.

# There are main types of sampling techniques

- i) Random Sampling
- ii) Stratified Sampling.

i) Random Sampling

in which each item in the population has an equal chance of being selected in the sample.

In sampling we will divide ~~data~~<sup>it into</sup> in train data and test data.

\* Because on train data built the model on test data, test the model.

More data in train & less data in test.

e.g:- placement training should be more (more knowledge) then we will get the job.

\* The ratio of the train : test is ?

70% : 30%

# Why Random.

→ Random sampling to avoid the biasness.

# Why train & test ?

→ Testing should not be done on same data which you ~~have~~ building the model.

eg:- Teaching      Training.  
Exam                Testing.

If you do the testing on the same data on which you have builded the model, you will get the accuracy which very high but it will be a high false accuracy.

## Day 9. Steps involved in Model Building.

1. Understanding the problem statement or requirement.
2. Get the data / Load the data.
3. Data pre-processing / ~~Load the~~ Data cleaning
4. Sampling (Random sampling i.e. Dividing in train and test).
5. Building the model.
6. Test the model and check the performance
7. Re run the model from some previous step.

step1 Reading the Data.

step2 Data Cleaning / Data preprocessing

Note :- if there is ~~is~~ the column there is more than 60% null values we want to drop the table.

Step3 Sampling.

eg:- From sklearn.model\_selection import train\_test\_s

eg:- cr\_train, cr\_test = train\_test\_split(cr, test\_size = 0.2)

eg:- cr\_train.shape

eg 981 \* 0.8

eg cr\_test.shape

eg. 981 \* .2

## Day 10 Introduction to Linear Regression

# Machine Learning :-

1. Supervised Learning
2. Unsupervised Learning.

# Properties of Supervised Learning.-

1. Target variable is present.
2. Sampling is possible in Supervised Learning

# Properties of Unsupervised Learning

1. Target variable is not present.
2. Sampling is not possible in Unsupervised Learning

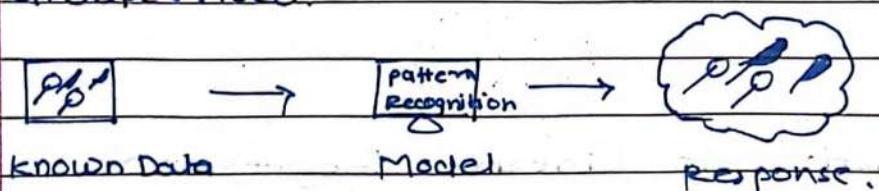
1) The machine that learns to that input new data.

\* Supervised Learning.

2) The machine will be train on know data or label data. it will predict data.

3) Supervised learning often we for classification problem, Regression problem and object detection problem.

\* Unsupervised.



The unsupervised <sup>unlabel</sup> Data while train in the model it will create the patterns.

Model will find out the pattern is known as clustering

The unsupervised learning is often use to perform the task as clustering, dimensionality Reduction,

1) Supervised Learning

It was classified into two categories of algorithm.

## 1. Regression -

A regression problem is when the output variable is a real and continuous value.

e.g.: weight & height (Integer)

## 2. Classification -

A classification problem is when the output variable is a category

e.g.: Gender, Colour (Object).

### 1. Regression:-

A. Linear Regression - "Straight Line" to do the prediction.  
( $\therefore$  Supervised).

\* Equation of Linear Regression :-

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + e$$

$\uparrow$   
(error terms)

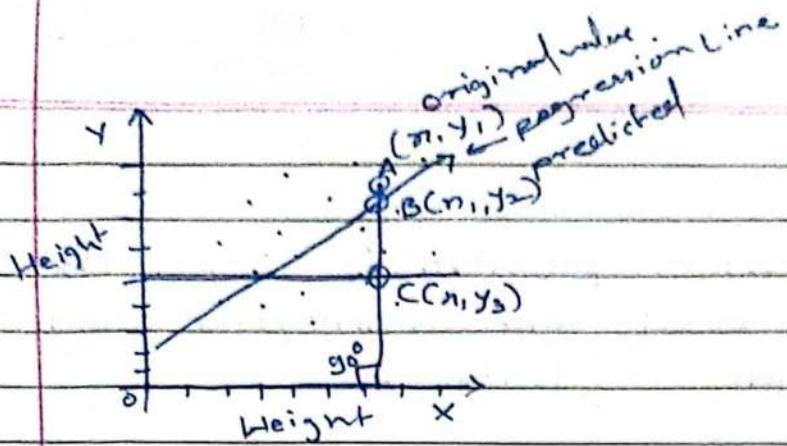
Consider I have a data.

X	Y
Weight	Mileage.
1000	15.
1200	13.
1100	13.5
1300	11.9.

Here; Y depends on X

Y is function of X.

Any value falls on regression line should be predicted value.



$$\text{Error} = \text{Actual} - \text{Predicted}$$

## Day 11. Linear Regression Part 1.

Actual	Predicted	Error	Sum of error (SSe)
40	43	-3	9
70	68	2	4
80	79	1	1
75	72	3	9
42	50	-3	9
Sum = 0			32.

Mean (Actual)	Predicted	Error	$\frac{\text{Sum of (Actual (mean) - Pred)}}{\text{Sum of error}}$
62.4	43	19.4	376.36
62.4	68	-5.6	31.36
62.4	79	-16.6	275.56
62.4	72	-9.6	92.16
62.4	50	12.4	153.76
			185.94 (SSR)
Sum sq regression			

Sum square regression is nothing but the mean value of  $\text{sq. of } (\text{mean} - \text{Pred})$

## Predictive Modeling :-

- \* Evaluate the Model
- \* No model is 100% accurate.
- \* If it is accurate it has overfitting or done any mistake.
- \* Evaluating parameter & matrix

### Day 12 Linear Regression Part 2

- \* Linear Regression is Predictive Modelling
- \* Evaluating parameter  
For model evaluation,

1. We will use the matrix known as Rsquare or  $R^2$ .

$$\begin{aligned} R^2 &= \text{SSR} / \text{SST} \\ &= \text{SSR} / (\text{SSR} + \text{SSE}) \quad \dots \text{Highly imp.} \end{aligned}$$

$R^2$  range is between 0 to 1

$R^2$  will always be a positive number.

Higher the  $R^2$  will better the model

→ If it is near to 1 is better

If it is near to 0 is not better

Choose the model which has higher  $R^2$  value

Q What is the problem associated with  $R^2$  ?

→ We consider that higher the  $R^2$  better the model but at a same time you keep on adding insignificant X variables and  $R^2$  keep on increasing

penkraft

rossing

MPG<sub>(Y)</sub> Wt of Car Cyl Fuel traffic colour interior

Model 1 - Consider only wt of the car

$$R^2 = 0.6.$$

Model 2 - Consider wt of the car + cyl + fuel

$$R^2 = > 0.6.$$

Model 3 - Consider wt of car + cyl + fuel + traffic + color  
+ interior

$$R^2 > .$$

\* To solve the problem of  $R^2$

2. Adjusted  $R^2$ .

$$\text{Adj } R^2 = 1 - \frac{(1 - R^2) \times (N - 1)}{N - k - 1}$$

N = No of Record or Rows

k = No of column or x variable.

- it will just give correct answer

- Adjusted  $R^2$  can never be greater than  $R^2$

$$0 \leq \text{Adj } R^2 \leq 1$$

penkraft

- Range of adj  $R^2$  is -ve to 1
- If  $R^2$  is moving towards 0 then Adj  $R^2$  is moving towards -ve.
- then we have to recheck the model.
- Higher the adj  $R^2$  better is the model.

use of MSE Mean Square Error.

MSE use for the comparison

also called as

### 3. MSE (Mean Square Error) (Cost Function)

Advantage of taking MSE.

1. Negative & positive error are getting cancel out.
2. MSE is used to punish the model for large errors.
3. Lower the MSE better is your model.

MSE range will be  $\underline{0}$  to infinity.

### Day 13 Linear Regression Part 3.

- MSE depends upon the range and magnitude of target variable.

### 4. RMSE (Root mean square error) -

$$\text{RMSE} = \text{sq root MSE}.$$

\* Why we need RMSE ???



Actual Y <sub>actual</sub>	Pred Y <sub>pred</sub>	Error Y <sub>actual</sub> - Y <sub>pred</sub>	MSE Error <sup>2</sup>	RMSE $\sqrt{MSE}$
7	4	3	9	3

- The value will be shoot up.

### 5. MAPE (Mean Absolute Percentage Error)

- RMSE gives same answer in actual unit value are there while MSE gives answer in square unit

- sometime interpreting MSE would be difficult because no. are getting shootup that why we RMSE.

- Lower the RMSE better is the model.

### 5) MAPE (Mean absolute percentage Error)

here,

$$\% \text{ error} = \frac{\text{error}}{\text{Actual}} \times 100$$

- Lower the MAPE better is the model

7%

$$\text{Accuracy} = 100 - 7 = 93\%.$$

## # Assumptions in Linear Regression

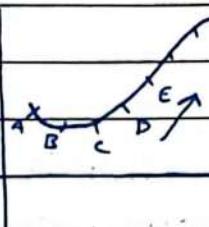
Assumption may be true or may not be true

penkraft

imp

1. There is a linear relationship between  $X$  &  $Y$ .
2. Error term are normally distributed and there is no pattern among them.
3. Minimum multicollinearity among  $X$  variable.
4. Homoscedasticity - variance around the regression line is same for all the predicted values.

Q Why we assume it to be linear only?



— Decision making should be difficult.

## Day 14. Assumption in Linear Regression.

### 3. Multicollinearity :-

There should be high correlation between x and y variable.

Production of Machine Raw Man electricity total  
Factory (Y) x<sub>1</sub> material power x<sub>2</sub> x<sub>3</sub> labour

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + E$$

2. Error terms are normally distributed and there is no pattern among them.

4. Homoscedasticity - variance around the regression line should be same for all the predicted values.

\* While building the model remove outliers to maintain the accuracy of the model.

\* First Model of Linear Regression.

### # Mandatory conditions for Linear Regression.

1. Target variable should always be numeric and continuous.
2. Number of rows (N) should always be greater than no. of columns (k)

eg:- lcnm = pd.read\_csv("LungCapData.csv")  
lcn.

# Requirement is to built the predictive model  
which can predict the Lung capacity.

# 3 step :- Data Cleaning.

eg:- lcn.shape.

eg:- lcn.info()

eg:- lcn.isnull().sum().

eg:- lcn.describe()

eg:- lcn.Gender.value\_counts()

eg:- lcn.Gender = lcn.Gender.replace({'male':1, 'Female':0})

eg:- lcn

eg:- lcn.Smoke.value\_counts()

eg:- lcn.Smoke = lcn.Smoke.replace({'yes':1, 'no':0})

eg:- lcn

eg:- lcn.Caesarean.value\_counts()

eg:- lcn.Caesarean = lcn.Caesarean.replace({'yes':1, 'no':0})

eg:- lcn

eg:- lcn.info()

Day 15 Linear regression model on lcn data.

# Step 4 :- Sampling.

eg:- From sklearn.model\_selection import train\_test\_split

eg:- lcn\_train, lcn\_test = train\_test\_split(lcn, test\_size=.2)

eg:- lcn\_train.shape # it will select Randomly data 580.

eg:- lcn\_test.shape # → ← →

eg:- lcn\_train\_x = lcn\_train.iloc[:, 1:]

eg:- lcn\_train\_y = lcn\_train.iloc[:, 0]

eg:- lcn\_train\_x

eg:- lcn\_train\_y

eg:- lcn\_test\_x = lcn\_test.iloc[:, 1:]

eg:- lcn\_test\_y = lcn\_test.iloc[:, 0]

eg:- lcn\_test\_x

eg:- lcn\_test\_y

## # Step5 :- Building the model

# - Import the necessary Function From  
sklearn

# - Create an object of that Function

# - Run the FIT Function and then model is created

# - Run the prediction (predict) the Function

eg:- From sklearn.linear\_model import LinearRegression

eg:- linereg = LinearRegression()

eg:- linereg.fit(lcn\_train\_x, lcn\_train\_y)

eg:- pred\_train = linereg.predict(lcn\_train\_x)

eg:- pred\_test = linereg.predict(lcn\_test\_x)

eg:- pred\_train.shape

eg:- pred\_test.shape

eg:- error\_test = lcn\_test\_y - pred\_test

eg:- error\_test.

## # Step 6 :- Evaluate the model

# R<sup>2</sup>.

eg:- Rsquare = linereg.score(lcn\_train\_x, lcn\_train\_y)

eg:- Rsquare  
penkraft

eg:- `N = len_train_x.shape[0]`

eg:- `k = len_train_x.shape[1]`.

# Adj-Rsquare.

eg:- `adj_Rsquare = 1 - (1 - Rsquare) * (N-1) / (N-k-1)`

eg:- `adj_Rsquare`.

# MSE.

eg:- `mse = np.mean(np.square(error_test))`

eg:- `mse`

# RMSE.

eg:- `rmse = np.sqrt(mse)`.

eg:- `rmse`

# MAPE

eg:- `mape = np.mean(np.abs(error_test * 100 / len_test_y))`

mape.

Day 16

eg:- `import pandas as pd`

`pd.set_option('display.max_columns', None)`

`pd.set_option('display.max_rows', None)`

`import numpy as np`

eg:- `ppt = pd.read_csv(r"Property_Price_Train.csv")`

`ppt`.

eg:- `ppt.head()`

## # Data Cleaning.

eg:- ppt.shape.

eg:- ppt.info()

eg:- ppt.isnull().sum()

eg:- len(ppt.isnull().sum()) [ppt.isnull().sum() &gt; 0]

eg:- ppt.isnull().sum() [ppt.isnull().sum() \* 100 / ppt.shape[0] &gt; 40]

eg:- ppt.Lot\_Extent.value\_counts()

eg:- ppt.LotExtent = ppt.LotExtent.fillna(ppt.Lot\_Extent.mean())  
ppt.Lot\_Extent

eg:- ppt.Brick\_Veneer\_Type.value\_counts()

ppt.Brick\_Veneer\_Type = ppt.Brick\_Veneer\_Type.fillna('None')

eg:- ppt = ppt.drop(['Pool\_Quality', ' ' ... ], axis=1)

eg:- ppt.shape

eg:- len(ppt.isnull().sum()) [ppt.isnull().sum() &gt; 0]

eg:- ppt = ppt.iloc[:, 1:]

eg:- ppt.head()

eg:- ppt.shape

IMP

## # Label Encoder.

eg:- `From sklearn.preprocessing import LabelEncoder`

eg:- `le = LabelEncoder()`

eg:- `pp1[pp1.select_dtypes(include=['object']).columns] = pp1[pp1.select_dtypes(include=['object']).columns].apply(lambda x: le.fit_transform(x))`

eg:- `pp1.head()`

eg:- `pp1.info()`

## # Sampling

eg:- `From sklearn.model_selection import train_test_split`

eg:- `pp1_train, pp1_test = train_test_split(pp1, test_size=0.25)`

eg:- `pp1_train_x = pp1_train.iloc[:, 0:-1]`

`pp1_train_y = pp1_train.iloc[:, -1]`

eg:- `pp1_train_x`

eg:- `pp1_test_x = pp1_test.iloc[:, 0:-1]`

`pp1_test_y = pp1_test.iloc[:, -1]`

eg:- `pp1-test_x`

## Day 17. Checking Assumptions on the PPT data.

# Building the model.

eg:- From sklearn.linear\_model import LinearRegression  
linereg = LinearRegression()

eg:- linereg.fit(ppt\_train\_x, ppt\_train\_y)

eg:- pred\_train = linereg.predict(ppt\_train\_x)  
pred<sub>test</sub> = linereg.predict(ppt\_test\_x)

eg:- error<sub>test</sub> = ppt<sub>test\_y</sub> - pred<sub>test</sub>.  
err<sub>test</sub>.

eg:- Rsquare = linereg.score(ppt\_train\_x, ppt\_train\_y)  
Rsquare

eg:- N = ppt\_train\_x.shape[0]  
k = ppt\_train\_x.shape[1]

eg:- adj\_Rsquare =  $1 - (1 - \text{Rsquare})^*(N-1) / (N-k-1)$   
adj\_Rsquare

eg:- linereg.coef #  $\beta$  values.

eg:- linereg.intercept\_ #  $\beta_0$  value.

eg:- mse = np.mean(np.square(error<sub>test</sub>))  
mse

eg:-  $\text{rmse.} = \text{np.sqrt}(\text{mse})$

$\text{rmse.}$

eg:-  $\text{mape} = \text{np.mean}(\text{np.abs}(\text{err-test} * 100 / \text{ppt-train}))$

$\text{mape.}$

# Lets check the assumptions

# Assumption on the train data.

eg.  $\text{err\_train} = \text{ppt-trainy} - \text{pred\_train}$

$\text{err\_train}$

# Error terms are normally distributed

eg:- `import matplotlib.pyplot as plt.`

eg:- `plt.plot(err_train, '*')`

eg:- `plt.hist(err_train, bins=70, edgecolor='r')`

# this is not a normally distributed graph.

eg. `pred actual df = pd.DataFrame()`

eg:- `pred actual df['Actual'] = ppt_train_y`

~~pred actual df['Predicted'] = ppt + pred\_train~~

eg:- `pred actual df`

eg `import seaborn as sns`

eg `sns.jointplot(x='Actual', y='Predicted', data=pred actual df, kind='reg').`

## Day 18. LCN data and Homework.

Continue From Day 15 ...

eg:-  $\text{err\_train} = \text{lcn\_train\_y} - \text{pred\_train}$   
 $\text{err\_train}$ .

# Error are normally distributed or not?

eg:- import matplotlib.pyplot as plt.

eg:- plt.plot(err\_train, 'o')

eg:- plt.hist(err\_train, bins=30, edgecolor='g');

eg:- new\_data = pd.DataFrame()

eg:- new\_data['Actual'] = lcn\_train\_y  
new\_data['Predicted'] = pred\_train

eg:- new\_data

eg:- import seaborn as sns

eg:- sns.jointplot(x='Actual', y='Predicted', data=new\_data,  
kind='reg').

## Day 19 Car price data and Removing Outliers.

eg:- import pandas as pd.  
import numpy as np.

eg:- car = pd.read\_csv(r"car-price.csv")  
car

eg:- car.shape.

eg:- car.head()

# We have special characters in our data

# We will convert this special characters to null values

# Once we convert into nulls then we already know how to handle the nulls.

eg: car = car.replace({'?': np.nan})

eg: car.head()

eg: car.isnull().sum()

eg: car.isnull().sum()[(car.isnull().sum() > 0)]

eg: car.normalized\_losses.value\_counts()

eg:- car.normalized\_losses = car.normalized\_losses.astype('float')  
car.bore = car.bore.astype('float')

car.stroke = car.stroke.astype('float')

car.horsepower =

car.peak\_rpm =

car.price =

eg:- car.normalized\_losses = car.normalized\_losses.fillna(car.normalized\_losses.mean())

eg:- car.bore = car.bore.fillna(car.bore.mean())

car.stroke =

car.horsepower =

car.peak\_rpm =

eg:- car.isnull.sum()

eg:- car.num\_of\_doors.value\_counts()

eg:- car.num\_of\_doors = car.num\_of\_doors.fillna('Four')

eg:- car = car.dropna() # it will drop the rows

eg car.shape

eg:- From sklearn.preprocessing import LabelEncoder  
le\_ = LabelEncoder()

eg:- car[car.select\_dtypes(include=['object']).columns] =  
car[car.select\_dtypes(include=['object'])].columns  
.apply(le\_.fit\_transform)

eg:- car.info()

eg:- car.head()

eg:- From sklearn.model\_selection import train\_test\_split

eg:- car = df1 # store in df1 for remove outliers

eg:- car\_train, car\_test = train\_test\_split(car, test\_size=0.2)

eg:- car\_train\_x = car\_train.iloc[:, :-1]

car\_train\_y = car\_train.iloc[:, -1]

eg car\_train\_x

car\_train\_y

eg:- car\_test\_x = car\_test.iloc[:, :-1]

car\_test\_y = car\_test.iloc[:, -1]

eg car\_test\_x

penkraft

eg:- `from sklearn.linear_model import LinearRegression  
lr = LinearRegression()  
lr.fit(car_train_x, car_train_y)`

eg:- `pred_train = lr.predict(car_train_x)  
pred_test = lr.predict(car_test_x)`

eg:- `err_test = car_test_y - pred_test  
err_train`

eg:- `Rsquare = lr.score(car_train_x, car_train_y)`  
Rsquare

eg:- `N = car_train_x.shape[0]  
k = car_train_x.shape[1]`

eg:- `adj_Rsquare = 1 - (1 - Rsquare) * (N - 1) / (N - k - 1)`  
adj-Rsquare

eg:- `mse = np.mean(np.square(err_test))`  
mse

eg:- `rmse = np.sqrt(mse)`  
rmse

eg:- `mape = np.mean(np.abs(err_test * 100 / car_test_y))`  
mape

eg:- `err_train = car_train_y - pred_train`  
err-train

eg:- `import matplotlib.pyplot as plt.`

eg:- `plt.plot(err_test, '.')`

eg:- `plt.hist(err_train, bins=30, edgecolor='r');`

eg:- preel-actual\_dF = pd.DataFrame()

eg:- preel-actual\_dF ['Actual'] = car\_train\_y  
 preel-actual\_dF ['Predicted'] = preel-train

eg:- preel-actual\_dF

eg:- import seaborn as sns

eg:- sns.jointplot(x='Actual', y='Predicted', data=preel-actual\_dF, kind='reg')

## # For removing Outliers

eg:- def remove\_outlier(dF, col, k):

mean = dF[col].mean()

global dF1.

sd = dF[col].std()

Final\_list = [x for x in dF[col] if (x > mean - k \* sd)]  
 #Upper Limit.

(#Lower Limit

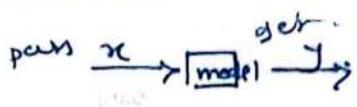
Final\_list = [x for x in dF[col] if (x < mean + k \* sd)]

dF1 = dF.loc[dF[col].isin(Final\_list)];

print(dF1.shape)

print('Number of outliers removed ---->', dF.shape[0] - dF1.shape[0]).

eg:- remove\_outlier(car, 'price', 2)



Date: \_\_\_\_\_  
Page No.: \_\_\_\_\_

## Day 20 Regularization Technique.

### # Regularization Technique -

- Ridge and Lasso's regression

Note:- There one problem in Linear regression is that there is high chance that model will get Overfit.

- Overfitting is ~~the~~ when the model will give very good results on train and gives bad results on test data.

In Overfitting the model will ~~not~~<sup>try to mimic</sup> fit the data. (∴ mean it will adjust the all data point cover all the datapoints)



Let's say there are two points in my training dataset.  
And I am having two features or variables  $x$  &  $y$ .

#### \* Overfitting :-

Train Accuracy - 90 - 98 %.

Test Accuracy - 66 - 65 %.

Low bias

High Variance

penkraft

Handmade

### \* Under-Fitting

Train Accuracy - 60%.

Test Accuracy - 62%.

High bias

High Variance.

### Balanced Fit / Generalized

Train Accuracy - 90%.

Test Accuracy - 88%.

Low bias

Low Variance.

### \* Reason For Overfitting.

1. Data used for training is not cleaned properly and it contains noise in it.
2. The model has a high variance there will be chance get your model Overfit.
3. The size of training cluster used is very small.
- 4) The model is too complex.

### \* Ways to tackle of.

1. Regularization Technique
2. Cross-validation Technique
3. Training your model with sufficient data.

### \* Theory of Regularization Technique

Q. How this ridge & lasso's avoid the overfitting

→ To prevent the issue of overfitting we are

penkraft

going to use ridge regression

- \* - Ridge regression can also be called as L2 regularization technique.

MSE

$$1/n * (x - y)^2 = 1/n * 0 = 0.$$

Ridge regression will add two values in this formula.  
 $= (x - y)^2 + \text{lambda} * (m)^2$   
↑ slope.

Lambda slope square is just penalizing the equation so that we should not get zero accuracy!

- \* - Lasso's Regression (L1 regularization).

$$= (x - y)^2 + \text{lambda} | \text{slope} |$$

↑ Constant absolute of slope.

- \* For the Feature selection we use the Lasso's regression.

- Feature selection

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$$

## Day 21 Ridge and Lasso's on Train Data.

eg:- import pandas as pd.  
pd.set\_option('display.max\_columns', None)  
pd.set\_option('display.max\_rows', None)  
import numpy as np

eg:- tt = pd.read\_csv(r"train.csv").  
eg:- tt.head()  
eg:- tt.drop(['ID'], axis=1)  
eg:- tt.isnull().sum().  
eg:- tt.shape.

eg:- From sklearn.preprocessing import LabelEncoder  
eg:- le = LabelEncoder()  
eg:- tt[tt.select\_dtypes(include=['object']).columns] =  
tt[tt.select\_dtypes(include=['object']).columns].  
.apply(le.fit\_transform).  
eg:- tt.head()  
eg:- tt.info()

eg:- From sklearn.model\_selection import train\_test\_split  
eg:- tt\_train, tt\_test = train\_test\_split(tt, test\_size=.2)  
eg:- tt\_train\_x = tt\_train.iloc[:, 1:]  
tt\_train\_y = tt\_train.iloc[:, 0]  
eg:- tt\_test\_x = tt\_test.iloc[:, 1:]  
tt\_test\_y = tt\_test.iloc[:, 0].

eg:- From sklearn.linear\_model import LinearRegression  
eg:- linereg = LinearRegression()  
eg:- linereg.fit(tt\_train\_x, tt\_train\_y)

eg:-  $R\text{square} = \text{linereg.score}(tt\_trainx, tt\_trainy)$   
 $R\text{square}$

eg:-  $N = tt\_trainx.\text{shape}[0]$   
 $k = tt\_trainx.\text{shape}[1]$

eg:-  $\text{adj\_Rsquare} = 1 - (1 - R\text{square}) * (N - 1) / (N - k - 1)$   
 $\text{adj\_Rsquare}$

eg:-  $\text{linereg.coef}$

eg:-  $\text{linereg.intercept}$

eg:-  $\text{pred\_train} = \text{linereg.predict}(tt\_train\_x)$   
 $\text{pred\_test} = \text{linereg.predict}(tt\_test\_x)$

eg:-  $\text{err\_test} = tt\_test\_y - \text{pred\_test}$   
 $\text{err\_train} = tt\_train\_y - \text{pred\_train}$

eg:-  $\text{mse\_test} = \text{np.mean}(\text{np.square}(\text{err\_test}))$   
 $\text{mse\_test}$

eg:-  $\text{mse\_train} = \text{np.mean}(\text{np.square}(\text{err\_train}))$   
 $\text{mse\_train}$

eg:-  $\text{rmse\_test} = \text{np.sqrt}(\text{mse\_test})$   
 $\text{rmse\_test}$

eg:-  $\text{rmse\_train} = \text{np.sqrt}(\text{mse\_train})$   
 $\text{rmse\_train}$

eg:-  $\text{mape\_test} = \text{np.mean}(\text{np.abs}(\text{err\_test}) * 100 / tt\_test\_y)$   
 $\text{mape\_test}$

eg:-  $\text{mape\_train} = \text{np.mean}(\text{np.abs}(\text{err\_train}) * 100 / tt\_train\_y)$   
 $\text{mape\_train}$

eg:-  $\text{pred\_test}.max()$

eg:-  $\text{pred\_test}.min()$

# We are facing the problem of overfitting so avoid this problem we ~~the~~ need to use the regularization technique.

## \* Ridge Regression.

eg:- From sklearn.linear\_model import Ridge

eg:- ridreg = Ridge()

eg:- ridreg.fit(tt\_train\_x, tt\_train\_y)

eg:- ridreg.score(tt\_train\_x, tt\_train\_y)

eg:- pre1\_rid\_train = ridreg.predict(tt\_train\_x)

pre1\_rid\_test = ridreg.predict(tt\_test\_x)

eg:- err\_ridge\_train = tt\_train\_y - pre1\_ridge\_train

err\_ridge\_test = tt\_test\_y - pre1\_ridge\_test

eg:- mse\_train\_rid = np.mean(np.square(err\_ridge\_train))

mse\_train\_rid = np.mean(np.square(err\_ridge\_train))

eg:- mae\_train\_rid = np.mean(np.abs(err\_ridge\_train))

eg:- mae\_train\_rid = np.mean(np.abs(err\_ridge\_train))

eg:- mape\_train\_rid = np.mean(np.abs((err\_ridge\_train \* 100) / tt\_train\_y))

mape\_train\_rid = np.mean(np.abs((err\_ridge\_train \* 100) / tt\_train\_y))

eg:- mape\_train\_rid = np.mean(np.abs((err\_ridge\_train \* 100) / tt\_train\_y))

eg:- mape\_train\_rid = np.mean(np.abs((err\_ridge\_train \* 100) / tt\_train\_y))

## \* Lasso's Regression.

eg:- From sklearn.linear\_model import Lasso

eg:- lasreg = Lasso()

eg:- lasreg.fit(tt\_train\_x, tt\_train\_y)

eg:- lasreg.score(tt\_train\_x, tt\_train\_y)

eg:- pre1\_train\_los = lasreg.predict(tt\_train\_x)

pre1\_test\_los = lasreg.predict(tt\_test\_x)

eg:- err\_train\_los = tt\_train\_y - pre1\_train\_los

err\_test\_los = tt\_test\_y - pre1\_test\_los

eg:- mse\_train\_los = np.mean(np.square(err\_train\_los))

mse\_train\_los = np.mean(np.square(err\_train\_los))

penkraft

## Day 22 Confusion Matrix.

1) - Confusion Matrix is a evaluation Matrix we use for evaluating the classification algorithms.

2) Original or actual values - Y or N & 0 or 1  
Predicted value. Y or N & 0 or 1

		Pred		Actual	
		0	1	0	1
Actual	0	22	14	9	4
	1	36	47	5	7

\* Accuracy = Total no of correct prediction / Total record.

$$= \frac{11}{20} = 0.55 \rightarrow 55\%$$

		0	1
0	5000	200	
1	240	50	

		0	1
0	4000	1200	
1	110	280	

Model A.

Model B.

$$\text{Accuracy Model A} = \frac{5050}{5590} = 0.90 = 90\%$$

$$\text{Accuracy Model B} = \frac{4280}{5590} = 0.76 = 76\%$$

- \* Accuracy cannot be only parameter to judge the model
  - \* We only do not focus on Accuracy we also want to focus on classes
- penkraft

Model B is good than Model A.

		Precl
		0      1
Actual	0	5000 (TN) 200 (FP)
	1	1340 (FN) 50 (TP)

- \* Class 0 will also be called as Class N (Negative)
- \* Class 1 can also be called as P (Positive)

1) TPR (True Positive Ratio) (Recall)

$$= \frac{TP}{TP + FN} = \frac{50}{50 + 340} = \frac{50}{390} = 0.12$$

\* TPR is always be greater

2) FPR (False Positive Ratio)

$$= \frac{FP}{FP + TN} = \frac{200}{200 + 5000} = \frac{200}{5200} = 0.0384$$

\* FPR to find the error in class 0

\* FPR is always be lesser.

## Day 23 Confusion Matrix part 2.

Q3 How to choose the best model?

- Choose a model which is balanced balanced :- performance on both the model should be good.
- Choose the model which meets your problem statement or requirement.

		Pred.		Pred.	
		Appro	Reject.	Actual	App
Actual	Appro	4200(TN)	150(FP)	Actual	3200(TN)
	Reject	190(FN)	50(TP)		750(FP)

Model A    Model B.

Acc 92

TPr 20.83

FPr 3

Acc 81.8

TPr 83.3

FPr 18

		0	1			0	1		
		0	5000(TN)	200(FP)			0	4600(TN)	300(FP)
0	1	340(FN)	50(TP)	1	110(FN)	280(TP)			
	0	4660(TN)	250(FP)		120(FN)	330(TP)			

If you are working for a telecomm major (jio or airtel) and there is problem that many customers are leaving and you want to find those customers which are ~~more~~ likely to leave so that we can take the preventive action class 0 not leaving and class 1 ppl leaving.

Acc 90.33

TPr 12.82

FPr 3

92

71.71

6

penkraft

5. Hit Ratio =  $\frac{TP}{TP + FP}$

\* Higher the precision the model is good.

## 5. AUROC CURVE.

Scenario 1 = You are getting 96% and want 99%.

Scenario 2 = You are getting 50% and want 75%.

Which scenario is easy to implement?

If TPR increase but FPR increase slowly after that FPR will not stop it will increase very fast.

Beyond a certain point when you try to increase the TPR TPR will increase very slowly AND FPR increase very fast on the certain points things will change ~~Drastically~~ Drastically.

## 6. F1 ratio

Harmonic mean bet pre and recall.

section A - 150

1/150

✓

section B - 400

1/400

x

$$F1 \text{ ratio} = 2 * \text{Precision} * \text{recall} / (\text{Precision} + \text{Recall})$$

High the F1 ratio better the model

range is like 0 to 1

penkraft

1. Accuracy - Higher
2. TPR - Higher
3. FPR - Lower
4. Precision - Higher
5. AUROC - Higher
6. F1 Ratio - Higher.

Day 24 Class Imbalanced:

	0	1
0	7523(TN)	340(FP)
1	110(FN)	30(TP)

Loan Approval System

1000 records

960 Rejected

40 Approved

Accuracy 94

TPR 21.82%

FPR 4.3

Precision 8.1

F1 ratio 11.7

Class imbalance is representation of one class is much more than the other class, in such scenario your model do very good for the class which is having more representation & will perform bad on the class which has less representation.

In class imbalance overall accuracy but as we know that accuracy is not only parameter to judge the model.

Solution for class imbalance,

There two solution -

penkraft

1. Over Sampling
2. Under Sampling

1. Over Sampling :- we increase record of the class which is under represented.
  2. Under Sampling we will decrease the record of the class which is over ~~represented~~ represented.
- Oversampling or undersampling should be done only over the sampling.
- Over & under sampling should be done only on train data because in oversampling we make duplicate of the record which are under represented & if you take duplicate in test data it will give false accuracy that why so we are avoiding test data.

## Day 25 Logistic Regression.

### \* Logistic Regression.

1. Logistic is supervised machine learning algorithm.
2. Target variable will be present
3. Sampling is possible.

penkraft

2. Logistic regression target variable is always binary.  
i.e. Yes or no, 0 and 1.
3. This Confusion matrix will be used to evaluate the logistic regression model.
4. X variable can be anything
5. In logistic reg... Y can take only 2 values i.e. 0 or 1 & its work on principle of probability.  
there will be 2 probability.  
Class 1 you will have one probability  
Class 0 you will have one probability.

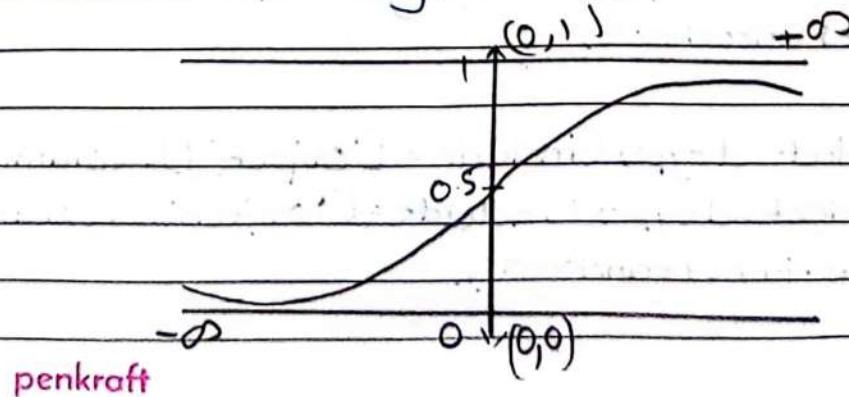
$$p + q = 1$$

$$0.4 + q = 1$$

$$0.6.$$

Whichever probability is more the Final

S curve or sigmoid curve :-



eg:- import pandas as pd.  
import numpy as np.  
pd.set\_option('display.max\_columns', None)  
pd.set\_option('display.max\_rows', None)

eg:- cr = pd.read\_csv(r"Credit\_Risk.csv")

eg:- cr.head

eg:- cr.shape

eg:- cr.isnull().sum

eg:- cr.Gender = cr.Gender.fillna('Male')

cr.Married = cr.Married.fillna('No')

cr.Dependents = cr.Dependents.fillna(0)

cr.Self\_Employed = cr.Self\_Employed.fillna('Yes')

cr.LoanAmount = cr.LoanAmount.fillna(cr.LoanAmount.mean())

cr.LoanAmount\_Term = cr.LoanAmount\_Term.fillna  
(cr.Loan\_Amount\_Term.mean())

cr.Credit\_History = cr.Credit\_History.fillna(0).

eg:- cr.isnull().sum()

eg:- cr = cr.drop(['Loan\_ID'], axis=1)

eg:- cr.shape

eg:- cr.Loan\_Status = cr.Loan\_Status.replace({'Y': 1, 'N': 0})

eg:- cr.info()

eg:- from sklearn.preprocessing import LabelEncoder.  
le = LabelEncoder()

eg:- cr[cr.select\_dtypes(include=['object']).columns] =  
cr[cr.select\_dtypes(include=['object']).columns].  
apply(le.fit\_transform).

eg:- cr.info()

eg:- From sklearn.model\_selection import train\_test\_split  
eg:- cr\_train, cr\_test = train\_test\_split(cr, test\_size = .2)  
eg:- cr\_train\_x = cr\_train.iloc[:, 0:-1]  
cr\_train\_y = cr\_train.iloc[:, -1]  
eg:- cr\_test\_x = cr\_test.iloc[:, 0:-1]  
cr\_test\_y = cr\_test.iloc[:, -1]

eg:- cr\_test\_x

eg:- cr\_test\_y

eg:- cr\_train\_x

eg:- cr\_train\_y

eg:- From sklearn.linear\_model import LogisticRegression

eg:- lr = LogisticRegression

eg:- lr.fit(cr\_train\_x, cr\_train\_y)

eg:- pred = lr.predict(cr\_train\_x)

pred = lr.predict(cr\_test\_x)

eg:- pred

eg:- From sklearn.metrics import confusion\_matrix

eg:- cr\_confuse = confusion\_matrix(cr\_test\_y, pred)

cr\_confuse

eg:- cr\_confuse.diagonal().sum() \* 100 / cr\_confuse.sum()  
# Accuracy

eg:- From sklearn.metrics import accuracy\_score

eg:- accuracy\_score(cr\_test\_y, pred) \* 100

## Day 26. Logistic regression model on CR data.

Continue ... (Day 25)

# Recall

# TP / TP + FN

eg:- From sklearn.metrics import recall\_score

eg:- recall\_score (cr-test-y, pred) \* 100

# Precision

# TP / TP + FP

eg:- From sklearn.metrics import precision\_score

eg:- precision\_score (cr-test-y, pred) \* 100

# F1 score

#  $\frac{2 * \text{precision} * \text{recall}}{\text{precision} + \text{recall}}$

eg:- from sklearn.metrics import f1\_score

eg:- f1\_score (cr-test-y, pred) \* 100

eg:- pred = lr.predict (cr-test-x)  
pred

eg:- pred\_proba = lr.predict\_proba (cr-test-x)

pred\_proba

eg:- from sklearn.metrics import roc\_auc\_score

From sklearn.metrics import roc\_curve

eg:- roc\_auc\_score (cr-test-y, pred) \* 100

eg:- pred\_proba[:, 1]

eg:- Fpr, tpr, ther = roc\_curve (cr-test-y, pred\_proba[:, 1])

eg:- import matplotlib.pyplot as plt

eg:- plt.plot (Fpr, tpr)

```

eg:- plt.plot(Fpr,tpr,color='r')
plt.xlabel('Fpr')
plt.ylabel('Tpr')
plt.title('Auroc on credit risk data')
plt.grid()
plt.text(x=.5,y=.5,s='Auroc is 73.59')

```

## Day 27 Logistic regression on Churn and Over sampling

```

eg:- import pandas as pd
import numpy as np
pd.set_option('display.max_columns', None)
pd.set_option('display.max_rows', None)

```

eg:- ch = pd.read\_csv(r"Churn.csv")

eg:- ch.head()

eg:- ch.shape

eg:- ch.isnull().sum()

eg:- ch.info()

eg:- from sklearn.preprocessing import LabelEncoder

eg:- le=LabelEncoder()

```

eg:- ch[ch.select_dtypes(include=['object']).columns]-
ch[ch.select_dtypes(include=['object']).columns].apply(le.fit_transform)

```

eg:- ch.info()

eg:- from sklearn.model\_selection import train\_test\_split

eg:- ch\_train,ch\_test = train\_test\_split(ch, test\_size = .2)

eg:- ch\_train\_x = ch\_train.iloc[:, :-1]

ch\_train\_y = ch\_train.iloc[:, -1]

eg:- ch\_test\_x = ch\_test.iloc[:, :-1]

ch\_test\_y = ch\_test.iloc[:, -1]

penkraft

eg:- ch-train-x

eg:- ch-train-y

eg:- ch-test-x

eg:- ch-test-y

# Overloading.

eg:- dF1 = ch.train [ch.train-y == 1]

DF1

eg:- DF1.shape

eg:- ch-train = pd.concat([ch-train, DF1, DF1])

eg:- ch-train.Churn.value\_counts()

eg:- ch-train.shape

eg:- From ~~sklearn~~<sup>8K learn</sup>.linear\_model import LogisticRegression

eg:- lr = LogisticRegression()

eg:- lr.fit(ch-train-x, ch-train-y)

eg:- pred = lr.predict(ch-train-x)

eg:- pred = lr.predict(ch-test-x)

eg:- pred

eg:- From sklearn.metrics import confusion\_matrix

eg:- ch-confuse = confusion\_matrix(ch-test-y, pred)

ch-confuse

eg:- ch-confuse.diagonal().sum() \* 100 / ch-confuse.sum()

eg:- from sklearn.metrics import accuracy\_score

eg:- accuracy\_score(ch-test-y, pred) \* 100 # Accuracy

eg:- # Recall (TPR)

# TP / TP + FN

eg:- from sklearn.metrics import recall\_score

eg:- recall\_score(ch-test-y, pred) \* 100

# Precision

#  $\frac{TP}{TP + FP}$

eg:- From sklearn.metrics import precision\_score  
 $precision\_score(ch\_test\_y, pred) * 100$

# F1 score

#  $\frac{2 * precision * recall}{precision + recall}$ .

eg:- From sklearn.metrics import f1\_score  
 $f1score(ch\_test\_y, pred) * 100$ .

eg:- pred = lr.predict(ch\_test\_x)

pred

eg:- pred\_proba = lr.predict\_proba(ch\_test\_x)

pred\_proba

eg:- From sklearn.metrics import roc\_auc\_score

From sklearn.metrics import roc\_curve

eg:- roc\_auc\_score(ch\_test\_y, pred) \* 100

eg:- pred\_proba[:, 1]

eg:- fpr, tpr, thres = roc\_curve(ch\_test\_y, pred\_proba[:, 1])

eg:- import matplotlib.pyplot as plt.

eg:- plt.plot(FPR, TPR)

eg:- plt.plot(FPR, TPR, color='y')

plt.xlabel('FPR')

plt.ylabel('TPR')

plt.title('Auroc on churn dataset')

plt.grid()

plt.text(x=0.5, y=0.5, s='Auroc is 76.25')

## Day 28 Model building on titanic data.

```
eg:- import pandas as pd  
import numpy as np.  
pd.set_option('display.max_columns', None)  
pd.set_option('display.max_rows', None)
```

```
eg:- Tt = pd.read_excel(r"titanic3.xls")
```

```
eg:- Tt.head()
```

```
eg:- Tt.shape
```

```
eg:- Tt.isnull().sum()
```

```
eg:- Columns_to_drop = Tt.columns[Tt.isnull().sum() * 100  
/ Tt.shape[0] > 40]
```

```
print('\n---- DROPPED COLUMNS ----\n', Columns_to_drop)
```

```
eg:- Tt = Tt.drop(columns=['cabin', 'boat', 'body', 'home.dest',  
'name'])
```

```
eg:- Tt.shape
```

```
eg:- # Check unique values
```

```
# Select only object columns
```

```
object_columns = Tt.select_dtypes(include=['object'])
```

```
# Check for non-numeric values in each object
```

```
column & display unique symbols
```

```
for column in object_columns:
```

```
    non_numeric_values = Tt[column].str.extract(r'\d+',
```

```
, expand=False).dropna().unique()
```

```
    if len(non_numeric_values) > 0:
```

```
        print(f"unique non-numeric symbols found in  
column '{column}':")
```

```
        print(non_numeric_values).
```

eg:- `Tt.age = Tt.age.fillna(Tt.age.mean())`

`Tt.Fare = Tt.Fare.fillna(Tt.Fare.mean())`

`Tt.embarked = Tt.embarked.fillna(Tt.embarked.mode()[0])`

eg:- `Tt.embarked.mode()`

eg:- `Tt.embarked.value_counts()`

eg:- `Tt.isnull().sum()`

## # Preprocessing

eg:- `from sklearn.preprocessing import LabelEncoder`

`le = LabelEncoder()`

`object_columns = Tt.select_dtypes(include=['object']).columns`

`Tt[object_columns] = Tt[object_columns].astype(str)`

`Tt[Tt.select_dtypes(include=['object']).columns] =`

`Tt[Tt.select_dtypes(include=['object']).columns].apply(le.fit_transform)`.

eg:- `Tt.corr()`

## # Rearrange the columns.

`Tt = Tt[['survived', 'pclass', 'sex', 'age', 'sibsp', 'parch',  
'ticket', 'Fare', 'embarked']]`

## # Model Building.

eg:- `from sklearn.model_selection import train_test_split`

eg:- `Tt_train, Tt_test = train_test_split(Tt, test_size=0.3)`

eg:- `Tt_train.head()`

eg:- `Tt_test.head()`

eg:-  $\text{Tt-train-x} = \text{Tt-train}.iloc[:, 1::]$

$\text{Tt-train-y} = \text{Tt-train}.iloc[:, 0]$

$\text{Tt-test-x} = \text{Tt-test}.iloc[:, 1::]$

$\text{Tt-test-y} = \text{Tt-test}.iloc[:, 0]$

eg :-  $\text{Tt-train-x}$

eg :-  $\text{Tt-train-y}$

eg :-  $\text{Tt-test-x}$

eg :-  $\text{Tt-test-y}$

## # Building the model

eg:- `From sklearn.linear_model import LogisticRegression.`

`lr = LogisticRegression()`

eg:- `lr.fit(Tt-train-x, Tt-train-y)`

eg:- `pred = lr.predict(Tt-test-x)`

`pred`

## # Confusion Matrix.

eg:- `From sklearn.metrics import confusion_matrix`

eg:- `Tt-tab = confusion_matrix(Tt-testy, pred)`

`Tt-tab`

## # Accuracy

eg:- `From sklearn.metrics import accuracy_score`

`accuracy_score(Tt-test-y, pred) * 100`

penkraft

### # Recall

eg:- `From sklearn.metrics import recall_score  
recall_score(Tt-test-y, pred) * 100`

### # Precision

eg:- `From sklearn.metrics import precision_score  
precision_score(Tt-test-y, pred) * 100`

### # F1 ratio

eg:- `From sklearn.metrics import f1_score  
f1_score(Tt-test-y, pred) * 100.`

eg:- `pred_proba = lr.predict_proba(Tt-test-x)  
pred_proba`

### # Auroc

eg:- `From sklearn.metrics import roc_auc_score  
From sklearn.metrics import roc_curve  
roc_auc_score(Tt-test-y, pred) * 100`

### # Matplotlib.pyplot

eg:- `Fpr, tpr, _ = roc_curve(Tt-test-y, pred_proba  
[:, 1])  
import matplotlib.pyplot as plt  
plt.plot(Fpr, tpr)`

## Day 29 Decision Tree.

- \* Decision Tree can be built on binary class and Multi class also.
- \* Whenever we build the LR we can't build the decision tree but the vice versa is not possible.
- \* ~~Decision Tree has Target variable also~~
- \* If all the variables are not of same importance then we need to check the most significant variable first to start the splitting of the tree.
- \* In decision tree decision will be based on tree.
- \* Whenever parameter or condition is most important that should be checked first and then the second most ~~go~~ on.

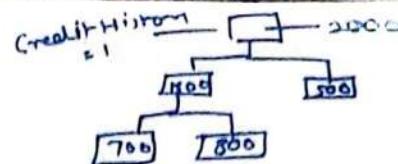
→ There are two Mathematical formula

→ How to decide parameters are most significant?

→ It will be decided based upon the mathematical formula (Entropy & Gini)

→ Either we use Entropy or gini to decision tree but we can't use both at same time.

Entropy } hyper parameter  
Gini  
penkraft



Date: \_\_\_\_\_  
Page No.: \_\_\_\_\_

## 1. Entropy.

- Entropy is nothing but measure of randomness in data
- Entropy tells how homogeneous & heterogeneous sample is.
- If sample data is complete homogeneous information gain is 'zero'
- If sample data is complete heterogeneous information gain is 'one'

Formula

$$\text{Entropy} = - \sum p_i \log_2(p_i)$$

$$IG(Y, X) = E(Y) - E(Y|X)$$

with respect to

Ex:- Let's say I have 1000 records, 600 belongs to class 1 and 400 belongs to class 0 ?

$$\text{prob for class 1} = 600/1000 \rightarrow 6/10$$

$$\text{prob for class 0} = 400/1000 \rightarrow 4/10$$

$$\text{Entropy of } Y, = - \frac{6}{10} \times \log_2\left(\frac{6}{10}\right) + \frac{4}{10} \times \log_2\left(\frac{4}{10}\right)$$

$$E(Y) = 0.97$$

penkraft

560 CH is good, 440 CH is bad

560 Good CH, the loan is approved for 530 ppl  
and 30 ppl rejected

440 bad CH, the approved for 70 ppl and 370 loan  
rejected.

0	0	1
0	370	70
1	30	530

~~30~~  $\frac{70}{440} = 15\%$

$$E(Y|X) =$$

1) loan status will P&P Good CH.

$$\rightarrow \frac{30}{560} \log \frac{30}{560} + -\frac{530}{560} \log \frac{530}{560}$$
$$= 0.305$$

(3)

2) loan status will P&P Bad CH

$$\rightarrow -\frac{370}{440} \log \frac{370}{440} + -\frac{70}{440} \log \frac{70}{440}$$
$$= 0.624$$

For entire population

$$= \frac{560}{1000} * 0.305 + \frac{440}{1000} * 0.624$$
$$= 0.1708 + 0.27456$$

$$E(Y|X) = 0.4453$$

$$TG(Y|X) = 0.97 - 0.4453$$

$$TG(Y|X) = 0.5247$$

## GINI Index or GINI Score.

Formula of GINI

$$\left[ 1 - \sum p_i^2 \right]$$

Day 30 Decision Tree Part - 1

Gini Index or Gini score

$$1 - \sum p_i^2$$

Loan status is our target variable,  
Calculating the Gini between Loan status  
and Credit History

Education	Credit History	Property Area	Loan status
Grad	1	Urban	N
Grad	1	Rural	N
Grad	0	Urban	Y
Grad	1	Urban	Y
Grad	1	Urban	Y
Non Grad	1	Urban	Y
Non Grad	1	Urban	Y
Grad	0	Semiurban	N
Non Grad	1	Urban	Y
Grad	0	Semiurban	Y
Grad	0	Semiurban	A.

$$CH = 1$$

True

7 record are satisfying  
Hochkond condition

False

4 record are not  
Satisfying the condition

Here 5 record loan is approved  
For 2 record loan is rejected

Here, 2 records loan is app  
and 1 record loan is rej

Here we are gonna use the Gini  
formula

$$= 1 - [(5/7)^2 + (2/7)^2]$$
$$= 1 - \left[ \frac{25}{49} + \frac{4}{49} \right]$$

$$= 1 - \frac{29}{49}$$

$$= 0.40$$

$$1 - [(2/4)^2 + (1/4)^2]$$

$$= 0.68.$$

For entire population

$$7/11 * 0.40 + 4/11 * 0.60$$

$$0.392$$

Day 31 DT on Churn data and Hyper parameter.

eg:- From sklearn.tree import DecisionTreeClassifier  
`dtree = DecisionTreeClassifier()`

eg:- `dtree.fit(BCCO_train_x, BCCO_train_y)`

eg:- `dec_pred = dtree.predict(BCCO_test_x)`  
`dec_pred`

eg:- `dtree.tab = confusion_matrix(BCCO_test_y, dec_pred)`  
`dtree.tab`

### \* Hyper Parameter or Hyper parameter tuning.

- Hyper Parameter are the parameters which are not part of data and user can control it or change it is called Hyperparameter.

→ depend on Criteria like model on Entropy or gini. (it is not part of Data) it depends on we.

→ In DT its our choice whether we want to build ~~a~~ tree or model using gini or entropy it is not a part of data but we can control it that why they are the hyper parameters.

hyper parameters are in infinite number of hyper parameters.

penkraft

## Day 32 Decision on CTG. and Over sampling (wine data.)

eg:-

```
import pandas as pd
import numpy as np
pd.set_option('display.max_columns', None)
pd.set_option('display.max_rows', None)
```

eg:- W = pd.read\_csv(r'wine.csv')

eg:- W.head()

eg:- W.shape()

eg:- W.isnull().sum()

### # Sampling.

eg:-

```
from sklearn.model_selection import train_test_split
W_train, W_test = train_test_split(W, test_size=25)
```

eg:-

```
W_train_x = W_train.iloc[:, :-1]
W_train_y = W_train.iloc[:, -1] # Selected by position
```

W\_test\_x = W\_test.iloc[:, :-1]

W\_test\_y = W\_test.iloc[:, -1]

eg:-

```
W_train_x.head()
W_train_y.head()
```

W\_test\_x.head()

W\_test\_y.head()

### # Decision Tree

eg:- From sklearn.tree import DecisionTreeClassifier  
dt = DecisionTreeClassifier(criterion='gini', max\_depth=None)  
dt.Fit (W\_train\_x, W\_train\_y)

eg:- dec\_prccl = dt.predict(W\_test\_x)  
dec\_prccl

## # Confusion Matrix.

eg:- From sklearn.metrics import confusion\_matrix.

eg:- dt\_tab = confusion\_matrix(W\_test\_y, dec\_prccl)  
dt\_tab

eg:- From sklearn.metrics import accuracy\_score  
accuracy\_score(W\_test\_y, dec\_prccl) \* 100.

## Day 32 Feature selection using DT on Credit Risk

### \* Feature Selection

Higher the Feature importance value more imp or more significant is that variable or Feature

### - Feature selection

- Lasso's regression is also used for the Feature selection.

eg:-  
 import pandas as pd  
 import numpy as np.  
 pd.set\_option('display.max\_columns', None)  
 pd.set\_option('display.max\_rows', None)

eg:- cr = pd.read\_csv(r"CreditRisk.csv")

eg:- cr.head()

eg:- cr.shape

eg:- cr.isnull().sum()

eg:- cr.Gender = cr.Gender.fillna('Male')

cr.Married = cr.Married.fillna('No')

cr.Dependents = cr.Dependents.fillna(0)

cr.Self\_Employed = cr.Self\_Employed.fillna('Yes')

cr.LoanAmount = cr.LoanAmount.fillna(cr.LoanAmount.mean())

cr.Loan\_Amount\_Term = -h

cr.Credit\_History = cr.Credit\_History.fillna(0)

eg:- cr.isnull().sum()

eg:- cr = cr.drop(['Loan\_ID'], axis=1)

eg:- cr.shape  
penkraft

## # Data Cleaning.

eg:- `cr.loan_status = cr.loan_status.replace({'Y':1,'N':0})`  
eg:- `cr.info()`

eg:- `From sklearn.preprocessing import LabelEncoder  
le = LabelEncoder()`

eg:- `Cr[cr.select_dtypes(include=['object']).columns] =  
cr[cr.select_dtypes(include=['object']).columns]  
apply(le.fit_transform)`

eg:- `cr.info()`

## # Sampling.

eg:- `From sklearn.model_selection import train_test_split.  
cr_train, cr_test = train_test_split(cr, test_size=.2)`

eg:- `cr_train_x = cr_train.iloc[:, 0:-1]  
cr_train_y = cr_train.iloc[:, -1]`

`cr_test_x = cr_test.iloc[:, 0:-1]  
cr_test_y = cr_test.iloc[:, -1]`

## # Decision Tree.

eg:- `From sklearn.tree import DecisionTreeClassifier  
dtree = DecisionTreeClassifier()  
dtree.fit(cr_train_x, cr_train_y)`

eg:- `dec_pred = dtree.predict(cr_test_x)  
dec_pred`

penkraft

## # Confusion Matrix.

eg:- From sklearn.metrics import confusion\_matrix

eg:- dtree.tab = confusion\_matrix(cr-test-y, dcl-pred)  
dtree.tab.

## # Accuracy.

eg:- From sklearn.metrics import accuracy\_score  
accuracy\_score(cr-test-y, dcl-pred)

## # Feature Selection

### # Feature Importance

eg:- dtree.feature\_importances\_

eg:- dtree.feature\_importances\_.sum()

# Higher the value of feature importance more important or significant that variable or feature is

eg:- cr-train-x.columns

eg:- feature\_imp = pd.DataFrame({'Feature': cr-train-x.columns, 'Importance': dtree.feature\_importances\_})

eg:- feature\_imp

eg:- feature\_imp.sort\_values('Importance', ascending=False)

eg:- cr-train, cr-test = train-test split(cr-test-size=.2)

eg:- cr-train-x = cr-train.iloc[:, [9, 5, 7, 6, 10, 2, 8]]

cr-train-y = cr-train.iloc[:, -1]

penkraft

cr-test-x =   
cr-test-y = 

cr-test-x.head()

## # DecisionTree Using Entropy.

eg:- From sklearn.tree import DecisionTreeClassifier  
dtree = DecisionTreeClassifier(criterion='entropy',  
max\_depth=7, min\_samples\_split=40, class\_weight  
='balanced')

eg:- dtree.Fit(cr-train-x, cr-train-y)

eg:- dec-pred = dtree.predict(cr-test-x)  
dec-pred

eg:- confusion matrix(cr-test-y, dec-pred)

eg:- accuracy score(cr-test-y, dec-pred)

eg:- From six import StringIO  
From IPython.display import Image  
from sklearn.tree import export\_graphviz  
import pydotplus  
import pydot.

eg:- dot-data = StringIO()

eg:- `export_graphviz(dtree, outfile=dot_data, feature_names=crc_train_x.columns)`

`graph = pydotplus.graph_from_dot_data(dot_data.getvalue())`

`(graph,) = pydot.graph_from_dot_data(dot_data.getvalue())`

`Image(graph.create_png())`

## Day 33 Random Forest.

- RF target variable can be Binary or multiclass
  - RF is collection of Decision Trees.
- \* How random Forest work  
→ In RF the algorithm uses the three techniques
1. Ensemble
  2. Bootstrapping
  3. Bagging

### 1) Ensemble

It's nothing but N-numbers of decision tree i.e. group of N-number of samples

### 2) Bootstrapping

It is creating the samples with the replacement.

### 3) Bagging -

It is nothing but aggregate or taking the vote.

N is nothing but no of sample

## # Random Forest.

eg:- From sklearn.ensemble import RandomForestClassifier  
ran\_for = RandomForestClassifier()  
ran\_for.fit(ctg\_train\_x, ctg\_train\_y)

eg:- pred\_ran = ran\_for.predict(ctg\_test\_x)  
pred\_ran

eg:- tab\_ran = confusion\_matrix(ctg\_test\_y, pred\_ran)  
tab\_ran

eg:- accuracy\_score(ctg\_test\_y, pred\_ran)

## Day 34 Random Forest on GridCV (before model building)

eg:- From sklearn.tree import DecisionTreeClassifier  
dt\_cr = DecisionTreeClassifier()

eg:- From sklearn.model\_selection import GridSearchCV

eg:- search\_dict = {'Criterion': ['gini', 'entropy'],  
'max\_depth': (5, 6, 7, 8, 9, 10),  
'min\_samples\_split': (40, 50, 60, 75, 90, 120)}

eg:- grid = GridSearchCV(dt\_cr, param\_grid=search\_dict)  
grid.fit(ctg\_train\_x, ctg\_train\_y)

eg:- grid.best\_params\_

eg:- dt\_cr = DecisionTreeClassifier(criterion=  
penalty= max\_depth  
min\_samples\_leaf=

## Day 35 KNN Algorithm.

### \* KNN (K-nearest Neighbour) :-

- KNN is supervised Learning algorithm
- Distance is used to find the neighbours.
- If distance is less then chances of being neighbour is very high.
- In KNN if you deciding the neighbours then you are not deciding it by passing the distance but you are deciding by passing the value of K
- k in KNN is nothing but the no. of neighbours
- Euclidean Distance is used to find the distance between two records

	Income	Saving	Return	Networth
A	10000	5000	2000	100000
B	10200	4500	3100	101500

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

$$= \sqrt{(10000 - 10200)^2 + (5000 - 4500)^2 + (2000 - 3100)^2 + (100000 - 101500)^2}$$

- k is the hyperparameter. (integer & positive)

\* How KNN actually works?

Lets say i have 500 records

	$X_1$	$X_2$	$X_3$	$Y$	$\rightarrow 0 \leftarrow 1$
1	10	8	20	1	
2					
3					
:					
500					

As of now i will consider the value of  $k=5$

Now comes our test data A,

A 11 9 30 ?

distance b/w Record 1st & person A

$$= \sqrt{(10-11)^2 + (8-9)^2 + (20-30)^2}$$
$$= 102$$

Q - The value of k will not be very less.  
→ If its less then your decision would be  
bias

## Day 36 KNN model on Adult kNN data.

- Q How to find the exact value of k in KNN?
- We will use the trade off technique.
  - We will do the trade off b/w accuracy and stability.
  - Elbow.
  - Normally we should avoid KNN.

eg:- import pandas as pd  
import numpy as np  
pd.set\_option('display.max\_columns', None)  
pd.set\_option('display.max\_rows', None)

eg:- A = pd.read\_csv(r'adultkNN.csv')

eg:- A.head()

eg:- A.isnull().sum()

eg:- A.replace('?', pd.NA, inplace=True)

eg:- A.isnull().sum()

eg:- A.workclass.value\_counts()

eg:- A.workclass = A.workclass.fillna('Private')

eg:- A.occupation.value\_counts()

eg:- A.occupation = A.occupation.fillna('Prof-specialty')

eg:- A['Native-country'].value\_counts()

eg:- A['Native-country'] = A['Native-country'].fillna('United States')

eg:- A.isnull().sum()

eg:- A.info.

# Cleaning (ELT)

eg:- From sklearn.preprocessing import LabelEncoder  
le = LabelEncoder()

penkraft

eg:- `A[A.select_dtypes(include=['object']).columns] = [ ]`. Apply (i.e. Fit-transform)

eg:- `A.info()`

eg:- `A.head()`

## # Sampling

eg:- `From sklearn.model_selection import train_test_split  
A_train, A_test = train_test_split(A, test_size=.2)`

eg:- `A_train_x = A_train.iloc[:, 0:-1]`

`A_train_y = A_train.iloc[:, -1]`

`A_test_x = [ ]`

`A_test_y = [ ]`

## # K - Nearest Neighbors

eg:- `From sklearn.neighbors import KNeighborsClassifier`

eg:- `Aknn = KNeighborsClassifier(n_neighbors=20, weight="distance")`

eg:- `Aknn.fit(A_train_x, A_train_y)`

eg:- `pred_knn = Aknn.predict(A_test_x)`

`pred_knn`

## # Confusion Matrix & Accuracy.

eg:- `from sklearn.metrics import confusion_matrix  
tab_knn = confusion_matrix(A_test_y, pred_knn)  
tab_knn`

eg:- `from sklearn.metrics import accuracy_score  
accuracy score(A_test_y, pred_knn)*100`

# How to Find the value of k

# For this we will ~~get~~ plot a graph betw Acc & Stability

IMP

eg:- `acc = []`

For k in range(1, 51):

`Aknn = KNeighborsClassifier(n_neighbors=k)`

`Aknn.fit(A_train_x, A_train_y)`

`pred_knn = Aknn.predict(A_test_x)`

`tab_knn = confusion_matrix(A_test_y, pred_knn)`

`acc.append(tab_knn.diagonal().sum() * 100 / tab_knn.sum())`

eg:- `k_values = list(range(1, 51))`

eg:- `import matplotlib.pyplot as plt.`

eg:- `plt.plot(k_values, acc).`

## Day 37 Naive Bayes Classifier.

- This algorithm works on the conditional probability.
  - Bayesian Theorem. on basis of this theorem NBC works.
  - We mostly use NBC. It does not have hyperparameters.
- Q What is Conditional Probability?
- $P(A/B) \rightarrow$  Prob of occurrence of event 'A' when event 'B' has already occurred.

$$P(A/B) = \frac{P(A \cap B)}{P(B)} : \text{Conditional Probability Formula.}$$

eg:-  
import pandas as pd  
import numpy as np  
pd.set\_option('display.max\_columns', None)  
pd.set\_option('display.max\_rows', None)

eg:- s = pd.read\_csv(r'svmFile.csv')  
eg:- s.head()  
eg:- s.shape  
eg:- s.isnull().sum()  
eg:- s.info()

eg:- From sklearn.preprocessing import LabelEncoder  
le = LabelEncoder()

eg:- s.marital.value\_counts()

eg:- `s.job = le.fit_transform(s.job)`  
`s.marital = le.fit_transform(s.marital)`  
`s['cons.conf.idx'] = le.fit_transform(s['cons.conf.idx'])`  
`s['emp.var.rate'] = le.fit_transform(s['emp.var.rate'])`

eg:- `s.education.replace({'illiterate': 0, 'unknown': 0, 'basic.6y': 1, 'basic.4y': 2, 'basic.9y': 3, ...}, inplace=True)`

eg:- `s.default.replace({...}, -1)`  
`s.bowing.replace({...}, -1)`  
`s.loan.replace({...}, -1)`

eg:- `s.contact = le.fit_transform(s.contact)`

eg:- `s.month.replace({...}, -1)`  
`s.outcome.replace({...}, -1)`  
`s.day_of_week.replace({...}, -1)`  
`s.y.replace({...}, -1)`

eg:- `s.info()`

# Sampling.

From `sklearn.model_selection import train_test_split`  
`s_train, s_test = train_test_split(s, test_size=.2)`

# Over Sampling

`df1 = s_train[s_train.y == 1]`

`df1.shape`

`s_train = pd.concat([s_train, df1, df1, df1])`

eg:- ~~s~~ train\_x = s\_train.iloc[:, :-1]  
~~s~~ train\_y = s\_train.iloc[:, -1]

~~s~~ test\_x = →  
~~s~~ test\_y = →

## # Naive Bayes Classifier.

eg:- from sklearn.naive\_bayes import MultinomialNB  
ct\_s = MultinomialNB()

eg:- ct\_s.fit(s\_train\_x, s\_train\_y)

eg:- predict\_ct = ct\_s.predict(s\_test\_x)  
predict

eg:- From sklearn.metrics import confusion\_matrix  
confusion\_matrix(s\_testy, predict\_ct)

eg:- From sklearn.metrics import accuracy\_score  
accuracy\_score(s\_test\_y, predict\_ct)

... Continue to SVM

## Day 38 SVM Algorithm.

— SVM (Support Vector Machine)

\* Supervised Learning main two property.

— Sampling is possible

— Target variable are present.

\*

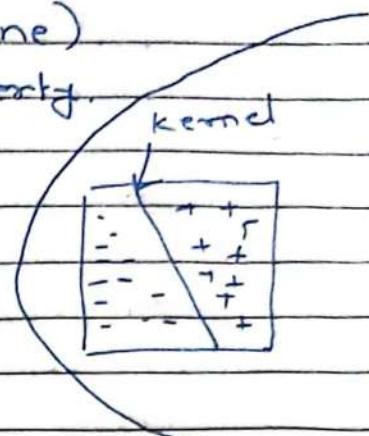
eg:- From sklearn.svm import SVC  
svc\_s = SVC()

eg:- svc\_s.fit(s\_train\_x, s\_train\_y)

eg:- pred\_svc = svc\_s.predict(s\_test\_x)  
pred\_svc

eg:- confusion\_matrix(s\_test\_y, pred\_svc)

eg:- accuracy\_score(s\_test\_y, pred\_svc)



## Day 39 Solving the ML paper - 1

e.g:-

```
import pandas as pd  
import numpy as np
```

```
from sklearn.preprocessing import LabelEncoder
```

```
from sklearn.ensemble import RandomForestClassifier
```

```
from sklearn.feature_selection import RFE
```

```
import seaborn as sns
```

```
from sklearn.metrics import roc_curve
```

```
import matplotlib.patches as mpatches
```

```
import matplotlib.pyplot as plt
```

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

— h —

precision\_score

— h —

recall\_score

— h —

F1\_score

— h —

accuracy\_score

— h —

roc\_auc\_score

```
from sklearn.model_selection import train_test_split
```

```
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
```

```
import xgboost as XGBClassifier
```

```
from sklearn.ensemble import AdaBoostClassifier
```

```
from sklearn.naive_bayes import GaussianNB
```

```
from sklearn.naive_bayes import MultinomialNB
```

penkraft

import warnings  
warnings.filterwarnings('ignore')

e.g:-  
data = pd.read\_csv(r"bank.csv", sep =';')  
data.shape  
data.head(10)  
data.info()

Q17 What does the primary analysis of several categorical features reveal.

Day 40 Solving ML paper part - 2

Day 41 Discuss Question Paper

Day 42

## Unsupervised Learning.

1. kmeans
2. PCA
3. Hierarchical Clustering
4. VIF.

### - kmeans

Properties of unsupervised:-

1. Target variable is NOT present
2. Sampling is not possible.

Here in kmeans k is the no. of clusters (Groups).  
These groups are formed based upon similarity.

Q How to measure the similarity?

Distance will help us to find similarity.



Euclidean Distance.

Important -

1. kmeans unsupervised algorithm where k stands for no. of clusters.
  2. k in kmeans is a hyperparameter
  3. kmeans is an iterative process (Repetitive)
- kmeans will randomly select any 2 records and mark them as centroid.

25 63 37 45 52 41 29 33 31 56 44 42 53 59  
27 34 38 50 28

penkraft

- At the starting centroid were assigned Randomly.
- According to centroid we will find the distance between centroid and the record whichever record has the least distance with centroid it will get added in that cluster.
- On 2<sup>nd</sup> iteration centroid were calculated by taking the "MEAN"

#### Stopping Criteria:-

When mean gets kind of Fixed which means records will stop to move from the other cluster which means our model get converged. The model will stopped the

... In Final model will always be similar for all.

## Day 43 Kmeans on SNS Data.

eg:- import pandas as pd  
import numpy as np  
pd.set\_option('display.max\_columns', None)  
pd.set\_option('display.max\_row', None)

eg:- sns = pd.read\_csv(r'sns.csv')

eg:- sns.head()

eg:- sns.shape

# Problem Statement

# The many millions of teenage consumers using social media or networking sites have attracted the attention of marketers

# So this companies, or clients can avoid targetting ads to teens with no interest in the product being sold

# Build kmeans clustering model to classify the interest of teenagers by using various attributes

eg:- sns.info()

eg:- sns.isnull().sum()

eg:- sns.gender.value\_counts()

sns.gender = sns.gender.fillna('F')

eg:- sns.age.value\_counts()

sns.age = sns.age.fillna(sns.age.mean())

eg:- sns.isnull().sum()

penkraft

eg:- `sns.gender.replace({'F': 1, 'M': 0}) inplace=True`

`sns.info()`

# Unsupervised model does not require sampling

eg:- `From sklearn.cluster import KMeans`

eg:- `kmeans_sns = KMeans(n_clusters=4)`

eg:- `kmeans_sns.fit(sns)`

# If 4 clusters have been created so each record should have gone into one of the clusters

eg:- `kmeans_sns.labels_`

eg:- `len(kmeans_sns.labels_)`

# For each record u have 1 label bcz it's one on one mapping

eg:- `kmeans_sns.cluster_centers_`

# Here we got 40 centroid point bcz our data has 40 columns i.e. it is 40 dimensions

# So each column has separate coordinates so we get 40 values.

eg:- `sns_new = pd.DataFrame(kmeans_sns.cluster_centers_)`

# For finding the best value or optimum value of k we need to plot the elbow plot

## SSD (sum of squared dimension).

Date:

Page No.:

## We will build the plot  
## run a loop from k=1 to k=10

e.g.: - kmeans.sns.score (sns)

## While building the plot we need to ignore  
the minus sign i.e. we are going to take  
the mod.

e.g.: - sn = range (1, 11): print (sn)

kmeans sns = [KMeans(n\_clusters=i) for i in sn]

kmeans sns.

score = [kmeans sns[i].fit(sns).score(sns) \*\* -1 for  
i in range (len(kmeans sns))]

print(score)

e.g.: - import matplotlib.pyplot as plt.

e.g.: - plt.plot (sn, score, marker='\*')

plt.grid()

e.g.: - len(cns\_new)

e.g.: - sns['Label'] = kmeans sns.label

e.g.: - sns.head (15)

penkraft

## Day 44 PCA Algorithm.

- PCA (Principal Component Analysis.)
- PCA is Unsupervised Machine Learning Algorithm
- PCA is the Dimensionality reduction algorithm
- Dimensionality means large data, etc.
- PCA will always work on numeric data.
- PCA uses orthogonal transformation for dimensionality reduction

\* Let say if 'A' is matrix and 'At' its transpose,

$$A * At = \text{Identity matrix}$$

here,  $A = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}$        $At = \begin{bmatrix} a & d & g \\ b & e & h \\ c & f & i \end{bmatrix}$

but normally  $A * A^{-1} = \text{Identity matrix}$  and  
 $A * At = \text{Identity matrix}$

So basically transpose ( $At$ ) is also acting as inverse ( $A^{-1}$ )

OT.  $A^* X = \lambda \text{ambda} * X$   
 $N * N * N * 1 = \lambda \text{ambda} * N * 1$

-  $\lambda$  is nothing but here eigen value  
 $X$  is eigen vector ( $N * 1$ )  
 $A = (N * N)$  matrix.

- PCA can also be used when you have the problem of multicollinearity.

- PCA can be used for dimensionality reduction  
but PCA cannot be used for Feature selection
- In PCA the first principal component is always very important followed by the 2nd one and 3rd one

## Day 45 PCA on PPT Data.

```
eg:- import pandas as pd  
import numpy as np  
pd.set_option('display.max_columns', None)  
pd.set_option('display.max_rows', None)
```

eg:- ppt = pd.read\_csv("Property\_Price\_Train.csv")

eg:- ppt.head()

### # Data Cleaning

eg:- ppt.shape

eg:- ppt.info()

eg:- ppt.isnull().sum()

eg:- len(ppt.isnull().sum()[ppt.isnull().sum() > 0])

eg:- ppt.isnull().sum()[ppt.isnull().sum() \* 100 / ppt.shape[0] > 40]

eg:- ppt.Lot\_Extent.value\_counts()

eg:- ppt.Lot\_Extent = ppt.Lot\_Extent.fillna(ppt.Lot\_Extent.mode())

eg:- ppt.Brick\_Vecter\_Type = ppt.Brick\_Vecter\_Type.fillna('None')

eg:- ppt.Brick\_Vecter\_Area = ppt.Brick\_Veneer\_Area.fillna(ppt.Brick\_Veneer\_Area.mean())

— 4 —

eg:- ppt = ppt.drop(['Pool\_Quality', 'Fence\_Quality', 'Miscellaneous\_Feature', 'Fireplace\_Quality', 'Lane\_Type'], axis=1)

eg:- ppt.shape

eg:- len(ppt.isnull().sum()[ppt.isnull().sum() > 0])

eg:- ppt = ppt.iloc[:, 1:]

penkraft

eg:- ppt.head()

## # Label Encoder

eg:- From sklearn.preprocessing import LabelEncoder

g:- le = LabelEncoder()

g:- ppt[ppt.select\_dtypes(include=['object']).columns] = ppt[ppt.select\_dtypes(include=['object']).columns].apply(le.fit\_transform)

eg:- ppt.head()

g:- ppt.info()

=

eg:- ppt1 = ppt

eg:- ppt = ppt.drop(['Sale Price'], axis=1)

g:- ppt.shape

eg:- ppt = ppt.sample(frac=1, axis=1)

↑

# This code is to shuffle the columns to change the position of the columns.

## # PCA

eg:- from sklearn import decomposition

from sklearn.decomposition import PCA

from sklearn.preprocessing import StandardScaler

# Before doing or applying the PCA, we will make our data Unitless

penkraft

eg:- `scaler = StandardScaler()`

eg:- `scaled_ppt = scaler.fit_transform(ppt)`

eg:- `scaled_ppt`

eg:- `scaled_ppt.shape`

eg:- `pca_ppt = PCA()`

eg:- `data_after_pca_transformation = pca_ppt.fit_transform(scaled_ppt)`

`data_after_pca_transformation`

eg:- `pca_ppt.explained_variance_ratio_`

eg:- `list(pca_ppt.explained_variance_ratio_)`

eg:- `pca_ppt.explained_variance_ratio_.sum()`

eg:- `len(pca_ppt.explained_variance_ratio_)`

eg:- `L1 = list(pca_ppt.explained_variance_ratio_)`

eg:- `np.sum(L1[0:53])`

# In any model there always be an error. (5 to 10%)

# We can select those number of PCA which are able to explain around 90% of the variance.

# We can ignore the rest of them.

# PCA can be use for dimensionality reduction  
not for Feature selection

# We will build our model using our first 55 PCA's

eg:- `type(data_after_pca_transformation)`

eg:- df1 = pd.DataFrame (data - after - pca - transformation)  
df1.head()

eg:- From sklearn.linear\_model import LinearRegression  
eg:- lr = LinearRegression()  
eg:- lr.fit(df1.iloc[:, 0:53], ppt1.sale\_price)  
eg:- lr.score(df1.iloc[:, 0:53], ppt1.sale\_price)

# PCA can be used to solve the problem of  
Multicollinearity.

eg:- df1.corr()

## Day 46 Hierarchical Clustering

- As we increase the distance, records which are far away will start to come in cluster
- As we decrease the distance then no of cluster will increase & records in cluster will decrease
- In hierarchical clustering everything depends on the distance.
- We use linkage function

```
eg:- import pandas as pd
      import numpy as np
      pd.set_option('display.max_columns',None)
      pd.set_option('display.max_rows',None)
```

```
eg:- mk = pd.read_csv(r'mall_kmeans.csv')
```

```
eg:- mk
```

```
eg:- mk.shape
```

```
eg:- mk.isnull().sum()
```

```
eg:- del mk['Customer ID']
```

```
eg:- mk.info()
```

```
eg:- mk.Gender.replace({'Female':0, 'Male':1}, inplace=True)
```

```
eg:- mk.head()
```

```
eg:- import scipy
```

```
eg:- from scipy.cluster.hierarchy import dendrogram, linkage
```

eg:- `z1 = linkage(mk, 'ward')`

eg:- `z1.shape`

eg:- `import matplotlib.pyplot as plt`

eg:- `plt.Figure(figsize=(10,10))`

~~denrogram(z1, p=16, truncate\_mode='lastp')~~

`plt.title('Hierarchical Clustering')`

`plt.xlabel('Cluster size')`

`plt.ylabel('distance')`

## # Agglomerative Clustering

eg:- `from sklearn.cluster import AgglomerativeClustering`

`clusters = AgglomerativeClustering(n_clusters=5)`

↑# Here default value is 2

eg:- `clusters.fit mk predict(mk)`

# Let's plot them and see how clusters are distributed.

eg:- `plt.scatter(mk['Age'], mk['Spending Score (1-100)'])`

`, c=clusters.labels_, cmap='rainbow')`

eg:- `Fig, ax = plt.subplots()`

`sc = ax.scatter(mk['Annual Income (k$)'], mk`

`[Spending Score (1-100)], c=clusters.labels_,`

`cmap='rainbow')`

`ax.legend(*sc.legend_elements(), title='clusters')`

## Day 47 Time Series Part 1.

- Time series is basically a Univariate Time series. Univariate means only one variable is present.
- When we compared time series to other ML algorithms like LR, Logistic, DT, RF & KNN we were having the separate X & Y variables.
- But in time series there is no separate X & Y variable. A single series will have both X & Y variable.
- For forecasting we use time series (e.g.: Weather Forecasting, stock market, etc.)

For eg:- We want to do the monthly Forecast of the sales, here I will be having the data from Jan to June,

	X variable					Y variable	
	Jan	Feb	Mar	Apr	May	Jun	July
	1000	1100	9500	8700	12000	14000	22?

Q) Why we need Forecasting ?

1. Patterns in Time Series :-

- i) Trend
- ii) Seasonality
- iii) Cyclic
- iv) Random.

penkraft

### 1) Trend :-

When the values will either increase or decrease with respect to time.

A) Up trend :- If values are continuously increasing then we called it as up trend

B) Down trend :- If values are continuously decreasing then we called it as down trend

### 2) Seasonality :-

The values will rise & fall with fixed time period.

### 3) Cyclic :-

The value will increase or decrease but not with fixed time period and generate time period is more than 2 years.

### 4) Random :-

If time series is not in trend nor seasonal and not even cyclic then we called random.

e.g. Stock Market.

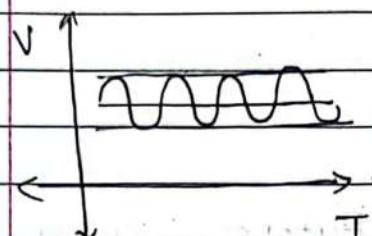
\* When values are random no time series model should be created.

## Day 48 Stationarity in Time Series.

### \* Stationarity in Time Series :-

- Stationarity is the main assumption in Time Series.

- mean or the variant are time variant they do not change with the time



Variance constant / Not changing  
mean constant / Not changing

- Before building the model

- If time series is not stationary then we will make it stationary.

- To make the series stationary we are going to use the Differencing Technique.

### \* Time series models :-

- 6 time series model

1) Auto regressive

2) Moving average

3) ARMA

4) ARIMA

5) SARIMA

6) SARIMAX.

## Day 49 Time Series model Part 1.

### 1) Auto Regressive (AR) :-

In AR model Future values are forecasted using the past values.

Eqn,

$$y_t = C + Y_1 y_{t-1} + Y_2 y_{t-2} + \dots + Y_n y_{t-n}$$

↓  
Previous  
Time  
period

— we consider the latest value

### 2) Moving Average:-

In MA model Future values are forecasted using the past error terms.

— Compare Actual & Predicted Values any error.

Eqn,

$$Y_t = Y_1 E_{t-1} + Y_2 E_{t-2} + \dots + Y_n E_{t-n}$$

### 3) ARMA :-(AR + MA)

— We are going to do the forecast using the previous actual values + previous error term also ARMA (p, q)

— Here AR can be mapped to 'p' & MA can be mapped to 'q'

penkraft

## Day 50 ARIMA and Differences Technique.

### 4) ARIMA Model

(Auto Regressive integrated moving average)

ARIMA can be denoted as  $(p(\text{AR}), d(\text{I}), q(\text{MA}))$

- \* Lets try to understand what is integrated means:
  - If series is not stationary we will make it stationary
  - Differencing Technique

Let say i have a series;

200, 210, 217, 240, 280, 360, 400, 205, 213, 220, 265, 300, 370. and this is not stationary.

We will be using Differencing Technique to make it stationary

Differencing means subtracting 1 term for the other terms.

After differencing,

10, 7, 23, 40, 80, 40, -195, 8, 7, 25, 55, 70.

... series after differencing

You will do 2<sup>nd</sup> level of differencing if series is not stationary even after 1<sup>st</sup> level of differencing

-3, 16, 4, 40, -40, -155 and

d= it will be decided by the level of differencing you are doing

- How to Find the value of  $p \& q$

By using ACF & PACF  
1) ACF (Auto correlation Function)  
2) PACF (Partial Auto Correlation Function)

e.g:- jan Feb mar April may  
540 510 530 540 580

The correlation will come from the lag of the series.

# Time Series on Air Passengers.

Date:

Page No.:

Day 51, 52, 53.

g:- import pandas as pd  
import numpy as np

eg:- airpas = pd.read\_csv("AirPassengers.csv")  
eg:- airpal

# PS is to Forecast the passengers From next 12 months (For year 1961)

# Random sampling in crime

# Breaking the series is also crime.

g:- airpas.info()

eg:- airpas.Month = pd.to\_datetime(airpas.Month, infer\_datetime\_format=True)

eg:- airpas.info()

eg:- airpal = airpal.set\_index(['Month'])  
airpal

eg:-

# We will check if our series is stationary or not.

eg:- import matplotlib.pyplot as plt.

eg:- plt.plot(airpal)

penkraft

# the series which we got is not stationary  
and we have to make it stationary

g:- plt.plot(airpas.diff().diff().diff())

# Even after multiple level of differencing  
series does not become stationary

# We are going to take the log of the  
series.

g:- airpas\_log = np.log(airpas)

g:- airpas\_log

g:- plt.plot(air\_log)

g:- plt.plot(air\_log.diff().diff().diff())

# After taking log and applying differencing we  
get the series stationary

# If you are using the log your model variance  
& mean will be constant.

# for building a model we need 3 thing one  
 $p, d, q$  and we know ( $p=?$ ,  $d=1$ ,  $q=?$ )

# To find the  $p & q$  we need to use ACF & PACF.

g:- from statsmodels.graphics.tsplots import plot\_acf,  
plot\_pacf.

g:- plot\_acf(airpas\_log)

g:- plot\_pacf(airpas\_log)  
penkraft

# When your ACF declines slowly or gradually and PACF declines sharply then we call it as AR process

Here are 2 lags we are getting correlation = 0 so this is AR(2) process

AR is Mapped to p

The value we got = (p=2, d=1, q=0)

## # ARIMA.

e.g.: From ~~statsmodels.tsa.arima.model~~ import ARIMA.

# The problem which associated with ARIMA  
ARIMA can't be handled the seasonality

## # SARIMA

Seasonal Auto Regressive Integrated moving average.

## # SARIMAX

Seasonal Auto Regressive Integrated moving average with Exogenous Variables

## # How to use SARIMA & SARIMAX

# We will use auto\_arima and whatever values we will get from auto\_arima we are going to pass it into our model.

e.g.: From ~~statsmodels.tsa.statespace.sarimax import SARIMAX~~

From ~~statsmodels.tsa.seasonal import seasonal\_decompose~~  
~~import pmndarima~~

from pmndarima import auto\_arima  
penkraft

eg:- result = seasonal\_decompose(airpus\_log)

eg:- result.plot();

# this plot is to just get the idea about data

eg:- auto\_arima(Airpus\_log, seasonal=True, m=12)

# m is the time period

# Whatever values will get from the auto\_arima we will pass them in SARIMAX.

eg:- model\_sarima = SARIMAX(Airpus\_log, order=(2,0,0), seasonal\_order=(0,1,1,12))

eg:- airpus\_model = model\_sarima.fit()

eg:- pred\_values = airpus\_model.predict(start=144, end=144) # predicted value.

eg:- pred\_values = np.exp(pred\_values)  
pred\_values

eg:- plt.plot(Airpus\_log)

plt.plot(pred\_values)

eg:- df = pd.DataFrame()

df['Predicted'] = pred\_values

date\_for\_pred = ['1961-01-01', ..., '1961-12-01']

df['Month'] = date\_for\_pred

df['Month'] = pd.to\_datetime(df['Month'], infer\_datetime\_format=True)

df = df.set\_index(['Month'])

penkraft

Nov 2021

g:- df.

g:- airpay['Passenger'].plot(legend=True, label='Original',  
color='green')  
df['Predicted'].plot(legend=True, label='Predicted',  
color='red')

Day 54 Time series on Alcohol Sales

or it is ...