

II Machine Learning II

1. Introduction

Data Science :

- First of all Data science is all about data only.
- It is a process of extracting knowledge and insights from data by using few methods or scientific methods.
- Data science is domain of study that concern deal with vast no. of data only using some modern tools & technique to find some hidden patterns, to derive some meaningful information, to make business decisions. To make such decision on the base of data it uses complex machine learning algorithms to build the predictive model.

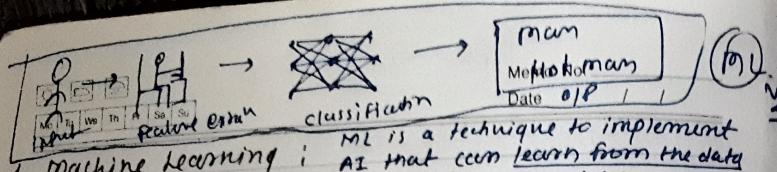
P.Q Fruits :

Orange >> orange, circular, sweet

Chikoo >> brown, oval, sweet

Orange >> orange, circular, li sour

Chikoo >> brown, oval, li sour



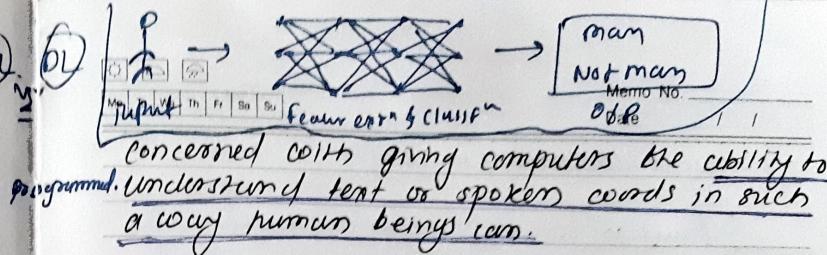
- It is a branch of AI and computer science which focuses on the use of data & algorithms to imitate the way that human learn, gradually improving its accuracy.
- It is a study of computer algorithms that improve automatically through experience and use of data.

2. Deep Learning : DL is a subfield of ML that uses Artificial Neural Networks to learn from the data.

- Deep learning is a subset of a machine learning, which is essentially a neural network with three or more layers. It uses for many AI applications and service that improve automation performing some analytical and physical tasks without human intervention.
- Deep learning technology lies behind so many everyday products and services such as, digital assistants, voice enabled TV remotes.

3. Natural language processing :

- It refers to the branch of AI again, which



4. Statistics

- This branch of applied mathematics that involves the collection, descriptions, analysis and inference of conclusion from the data.

5. Data Visualization

- Data visualization is a graphical representation of information and data. By using elements like charts, pie charts, graphs, maps, data visualization tools provide an accessible way to see and understand trends, outliers, patterns in data.

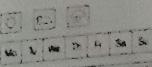
Data :

- Structured : CSV, EXCEL
- Semistructured : JSON, HTML
- Unstructured : text, images, audios, videos

#1. Machine Learning

Data ~~sets~~

- ↳ ① Training dataset (70%)
- ↳ ② Testing dataset (30%) unseen data.



e.g. Fruit understanding machine.

Fruit \Rightarrow colour shape taste

orange \Rightarrow	orange, circular, sweet
chikoo \Rightarrow	brown, oval, sweet
chikoo \Rightarrow	brown, oval, (i) sour
orange \Rightarrow	orange, circular, sweet
orange \Rightarrow	orange, circular, (i) sour
chikoo \Rightarrow	brown, oval, (i)sour
chikoo \Rightarrow	brown, oval, sweet
orange \Rightarrow	orange, circular, sweet

Application of ML

1. Chatbot (voice based, text based)
2. Recommendation system (flipkart, amazon, tutmix)
3. Medical field (cancer, covid)
4. Finance domain (bank)
5. Insurance Domain
6. Stock market
7. Self driving car
8. House price prediction
9. Object detection
10. Face recognition
11. Weather forecasting
12. Speech recognition

13. Cyber security.

14. Google maps.

15. Human activity recognition.

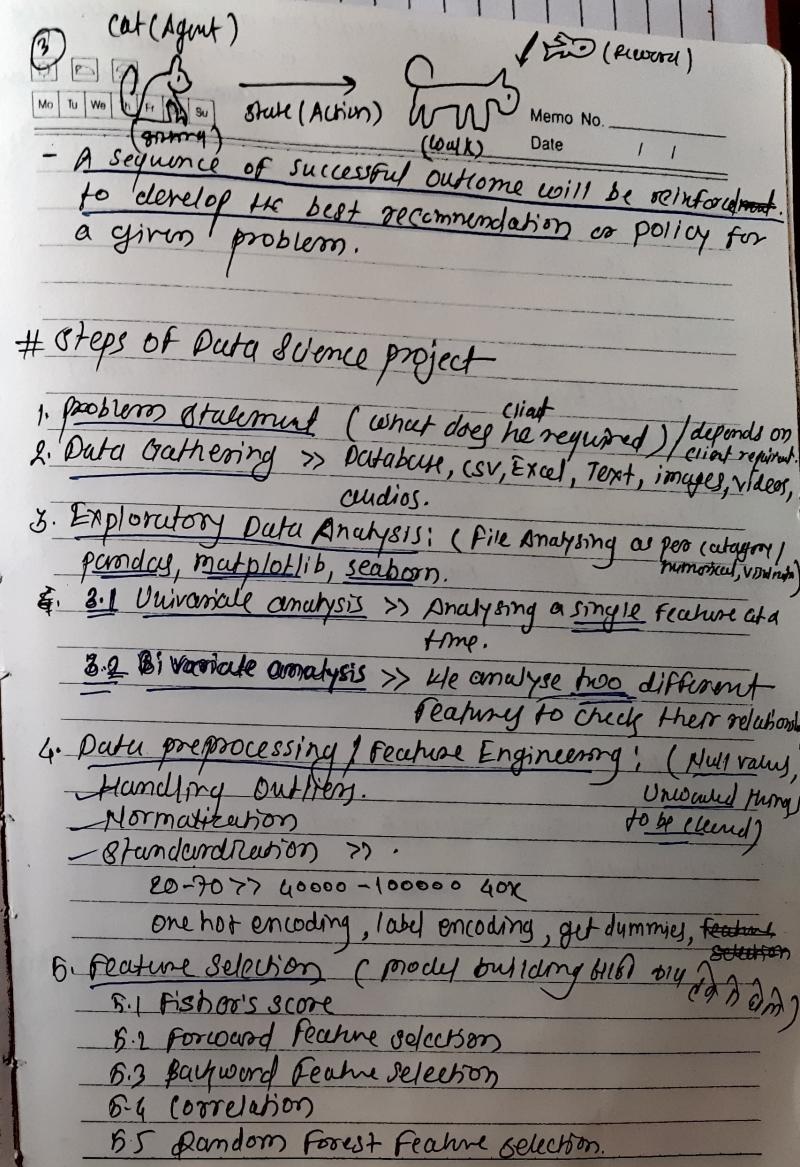
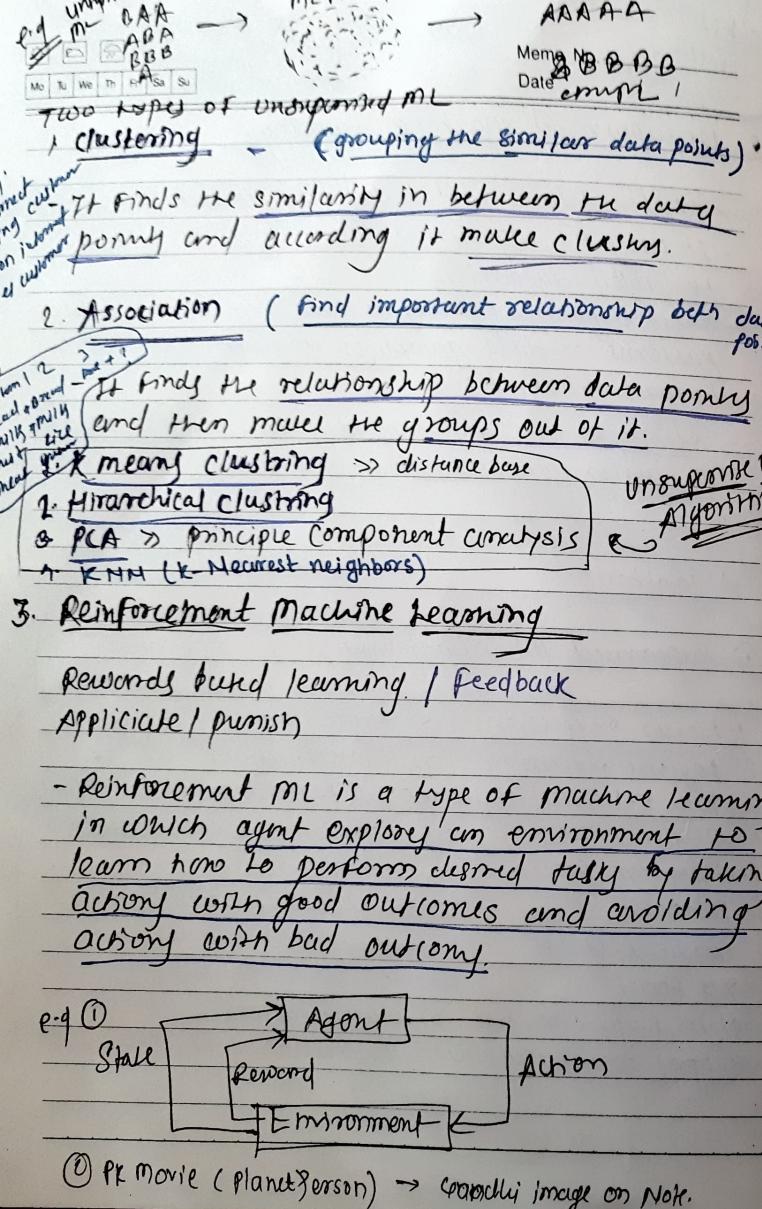
Types of Machine Learning

- (1) Supervised ML — The ML learning algorithm from labeled data
- (2) Unsupervised ML — Unlabelled data
- (3) Reinforcement ML — It is an area of ML concerned with how intelligent agents take actions in an environment to maximize its rewards.

(*) Independent Variable \Rightarrow which are given as input for prediction, predictors.

(*) Dependent Variable \Rightarrow Target variable, Outcome variable.

Training data \Rightarrow Independent + Dependent
Testing data \Rightarrow Independent



Classification :- It is about predicting a class or discrete values e.g. male or female; true or false.

Regression :- It is about predicting a quantity or continuous values e.g. salary, age, price.

6. Model training : (ML algorithm)

6.1 Linear regression

6.2 Logistic regression

6.3 KNN

6.4 DT

7. Model Evaluation : (Algorithm on test data)

(try to understand) Classification > Confusion matrix, classification report, accuracy score.

(try to understand) Regression > MSE, R2 Score

8. Deployment : (AWS, GCP, Azure)

Things to focus in each of every algorithm:

1. Concept.

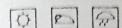
2. Assumption

3. Effect of missing values & and outliers.

4. Advantages

5. Disadvantages

6. Applications.



Mon Tu We Th Fr Sa Su

Day - 62

Always's
target
memo no
scribble is only
one
Date 26/15/2022

V7.0.0

1. Linear Regression Algorithm

- It is type of supervised ML algorithm.

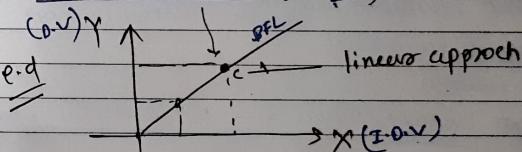
- This is a very basic algorithm.

- It is the algorithm which highly used initially and even highly studied as well.

- # why name is linear

Source) showing if not able to explain clearly - Linear regression is a linear approach to modelling the relationship between dependent and independent variable.

- Main target of linear regression is to find out the best fit line. & look for linear regression with O.I.V & I.O.V



1 Day - 62

Types of Linear Regression Algorithm

1.1 Simple Linear Regression. one O.V & one I.O.V

Dependent variable ~~Dependent Variable~~ and Independent Variable ~~Independent Variable~~

1.2. Multiple Linear Regression \rightarrow (1) single variable I.O.V
one dependent variable, more than one independent variable.
 \rightarrow D.I.V
 \rightarrow I.O.V

- Target variable should be continuous in nature.

Outcome - Target variable / outcome variable

(Y) Dependent Variable
Memo No. _____
Date / /

- Linear regression is an algorithm where we basically try to understand the linear relationship between dependent variable and one or more independent variables.

- It is based on line

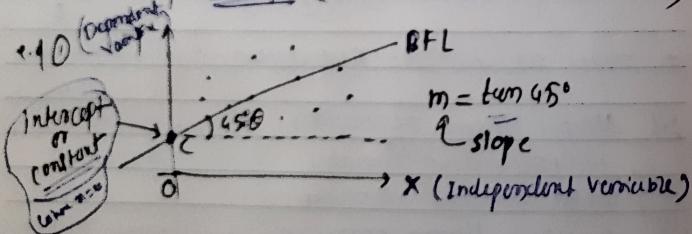
so when it comes to linear regression we basically have a straight line. and straight line is nothing but linear equation. That is why this algorithm is known as Linear regression.

- equation of line $\Rightarrow y = ax + b$ or $y = mx + c$

y = Dependent variable / outcome variable simple linear regression eqn.
 x = independent variable

m = Slope / position

c = slope intercept (where line of y when x is zero)

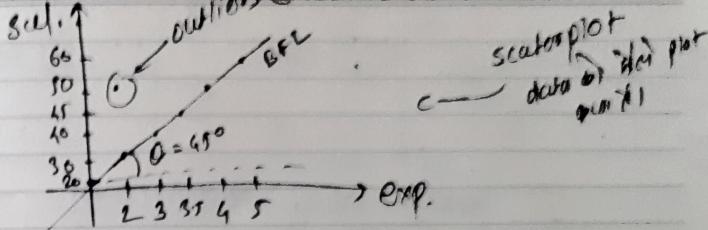


e.g.) Experience vs salary

1	20K	5	60K
3	40K		
3.5	55K		
4	50K		
2	50K	>> outliers	

x = independent variable.
 y = dependent variable

(\uparrow Intercept \Rightarrow where line of y when x is zero)
Memo No. _____
Date / /



If we want to find out $45 \Rightarrow ?$ scaling.

Simple linear regression

$$y = mx + c$$

$$mx + c = 1 \times 45 + 20 \\ = 20004.5$$

$$\therefore \tan 45 = 1$$

$$4.5 \text{ exp} \gg 20004.5 \quad (\text{So } k \approx 3.14159 \text{ in } \pi \text{ is used})$$

Multiple linear regression

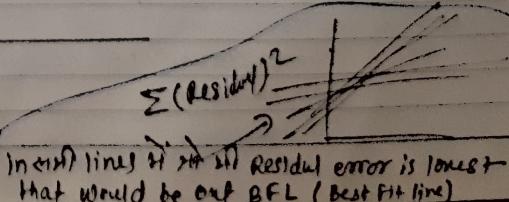
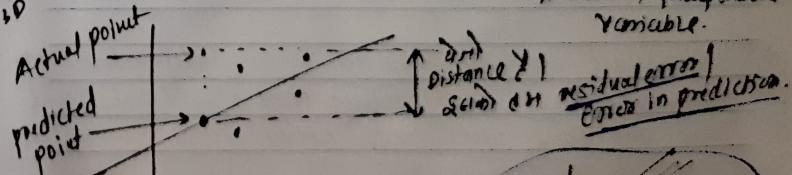
Suppose we have two independent variables.

multiple linear regression

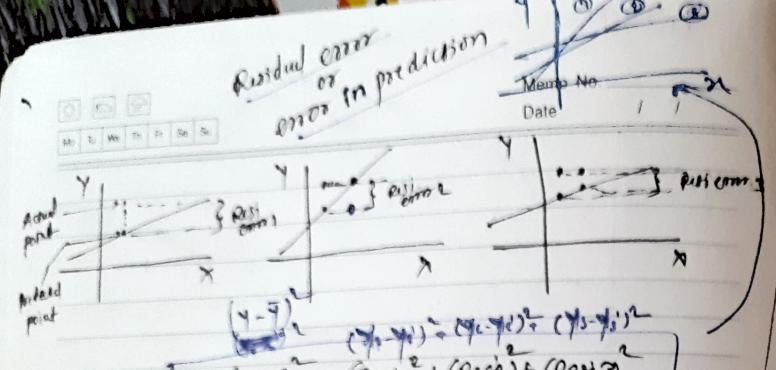
$$y = m_1 x_1 + m_2 x_2 + c$$

$$y = m_1 x_1 + m_2 x_2 + m_3 x_3 + \dots + m_n x_n + c$$

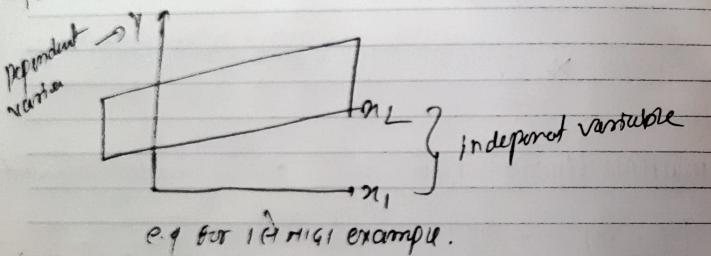
n = no. of independent variables.



In case lines fit \Rightarrow if residual error is lowest & that would be our BFL (Best Fit Line).



- sum of the value on X of the X child DFL.
 p. ⑧ \Rightarrow CH HIGH A MIGI ED to 3D mugale au GAN X



- The motto of linear regression is we need linear data.

Imp. Assumption
For 1-2. $X_1, X_2 \gg$ There should not be any relation b/w both of them.

2 independent variable supposed to be independent with each other. It simply means they should not have any influence upon each other.

Eg. Crt. \gg Height, HP, Mile/gallon etc.

Mon	Tue	Wed	Thu	Fri	Sat	Sun
-----	-----	-----	-----	-----	-----	-----

Memo No. _____ Date. / /

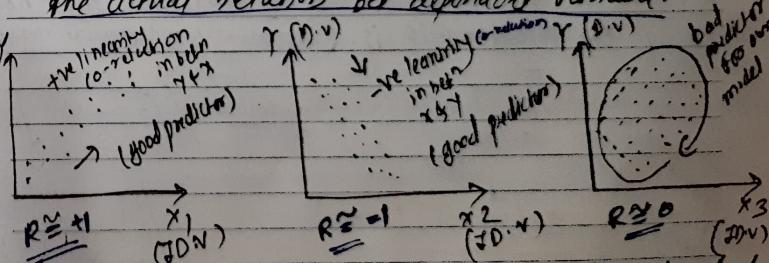
In real world more the wt, more the HP but our model always considers that there should be no relation b/w between this two variable. This problem is known as co-linearity. (contd)

$$Y = m_1 X_1 + m_2 X_2 + C$$

$m_1 \Rightarrow m_1$ is how much y changes whenever x_1 changes by one unit keeping x_2 constant.

$m_2 \Rightarrow m_2$ is how much y changes whenever x_2 changes by one unit keeping x_1 constant.

x_1 and x_2 (independent variable) \gg if this two variables do not have some relation or they are not independent of each other, then it becomes difficult understand the actual relation bet dependent variable.



- There should be high co-relation in betw dependent variable & independent variable.

- There should not be co-relation in betw independent variable.

K10 about machine learning of Iris
classification with Python

Day - B3

Memo No.

Date 27/05/2022

In Jupyter notebook
= file correlation.

import numpy as np
import pandas as pd

df = pd.read_csv('Iris.csv')
df.head() # note the column drop

df.drop(['ID', 'Species'], axis=1, inplace=True)
df.head()

Now we are checking co-relation in it,

import seaborn as sns

import matplotlib.pyplot as plt

df.head() # required column idnm

sns.pairplot(df)

df.corr() # coefficient of correlation.

R approx -1 >> good

R approx +1 >> good

R approx 0 >> bad

	A	B	C	D
A	1.00	-0.32	-0.32	-0.32
B	-0.32	1.00	-0.32	-0.32
C	-0.32	-0.32	1.00	-0.32
D	-0.32	-0.32	-0.32	1.00

sns.heatmap(df.corr(), annot=True) # graph visualization

○ ○ ○

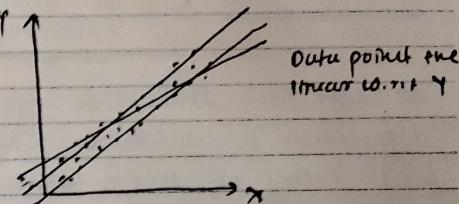
Mo Tu We Th Fr Sa Su

Memo No.

Date 1/1

Good predictor

(gradient descent algorithm)



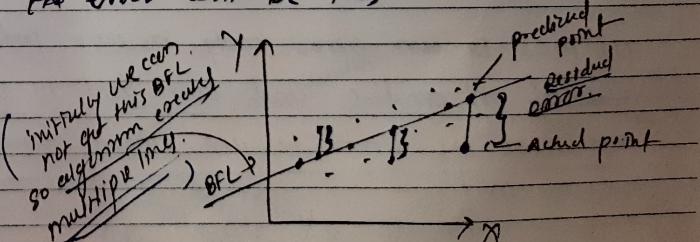
$$y = mx + c$$

(internally algorithm uses variation X/Y)

- It starts with the random value of m and c. and at first value of m and c will change, relationship b/w n & Y will also change.

- Different m & c value it basically evaluate different possible lines before finding the best fit line.

- BFL could be that which goes through the maximum no. of data points. and at some time it minimizes the distance b/w other point of the line. following the error will be less.



- so our algorithm initially creates multiple lines on the basis of different m and c values.



Memo No. _____
Date 1/1

- from all all multiplied but to find out this best fit line has to have best value of $m & c$ where we will be having least error.

- so there is a process or we can say an algorithm by which we can find out the best values of $m & c$ which is called as "gradient decent algorithm".

11 (line) $m_1 c_1 \Rightarrow$ suppose data point away from the line $\Rightarrow E_1$ (error) will be high

12 $m_2 c_2 \Rightarrow$ suppose data point bit away from the line $\Rightarrow E_2$ will be lower.

LBF MLE CBF \Rightarrow distance b/w mat of data point one on line, data points above & below the line are at minimal distance

\Rightarrow EBF least

error Best Fit.

Calculations of error (loss function / cost function)

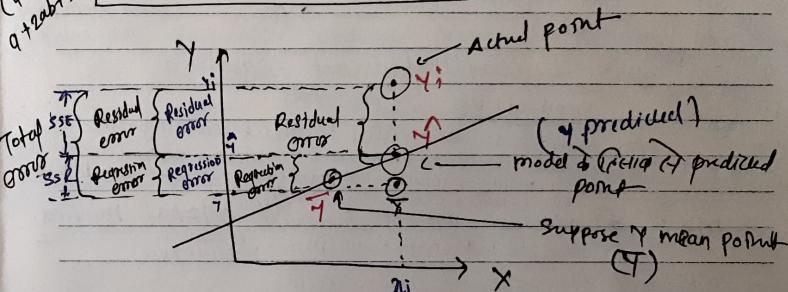
- The loss is an error only in our predicted value of m and c and our goal is to minimize the error to get most accurate value of $m & c$.

The cost function of linear regression
 \Rightarrow mean square error (MSE)

Residual error = in b/w actual point & predicted point
Regression error = in b/w mean value & predicted point
Mean square error (mean square error)
Memo No. _____
Date 1/1

quadratic equation (LSD)
 $E = \frac{1}{n} \sum (y_i - \bar{y})^2$
 $E = \frac{1}{n} \sum (y_i - (mx_i + c))^2$

This is cost function of linear regression



① $y_{actual} - y_{predicted} = \text{Residual error (SSE)}$

② $y_{predicted} - y_{mean} = \text{Regression error (SSR)}$

③ Total error = $y_{actual} - y_{mean} = \text{Residual error} + \text{Regression error}$
 $= SSE + SSR$

$SSE \Rightarrow$ sum of square error $\Rightarrow y_i - \bar{y} \Rightarrow \sum (y_i - \bar{y})^2$
 $SSR \Rightarrow$ sum of square regression error $\Rightarrow \bar{y} - \bar{y} \Rightarrow \sum (\bar{y} - \bar{y})^2$
 $SST \Rightarrow$ Total error / MSE $\Rightarrow y_i - \bar{y} = SSE + SSR$

$MSE \Rightarrow \frac{SSE}{n}$

Cost function of linear regression

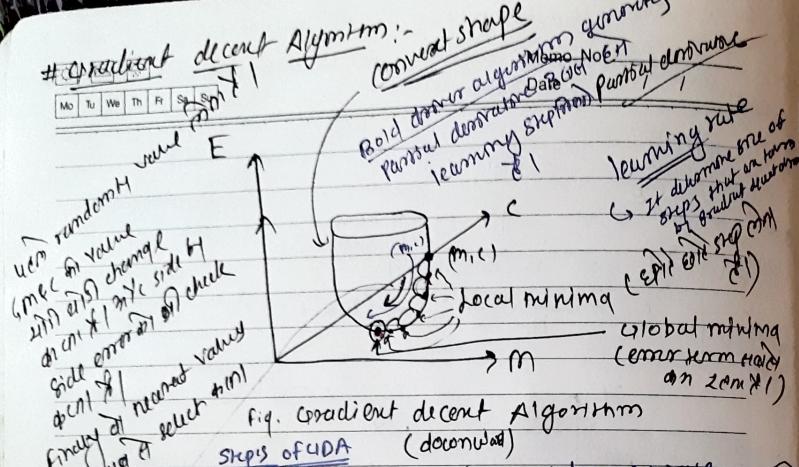
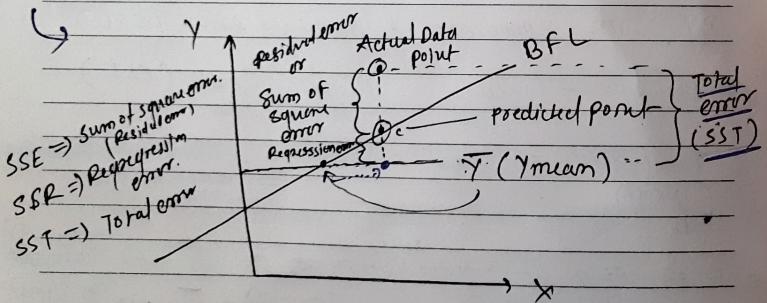


Fig. Gradient descent Algorithm

Steps of GDA (deconvoluted)

- ① Initialize random value to m, c + learning rate
 - ② Create new value of m, c \rightarrow (Using partial derivative of $J(m, c)$)
 $m_{new} = m_{old} - \text{learning rate} \times \text{partial sum w.r.t } m$
 $c_{new} = c_{old} - \text{learning rate} \times \text{partial sum w.r.t } c$
 - m, c \Rightarrow value goes to global minimum.
 - ③ Replace $\Sigma (y - \text{pred})^2$ cost funcn (MSE) near to zero.
- To findout the global minima



- Total error is distance betn actual data point or y mean / expected value.

Day - 5 G						
Mo	Tu	We	Th	Fr	Sa	Su
Sum of Regression square \Rightarrow SSR (sum of squares error)	Sum of square error \Rightarrow SSE					
Sum of square total error \Rightarrow SST						
Memo No. _____	Date 30/07/2022					

Total error e.g

gpttka people weight

one people weight

The difference b/w them called total error.

- error measures empty variance.

- The gradient descent working in finding the best fit line, where the BFL is that line where this error is minimum.

① So then model is most fit, it means all data points should lie on BFL, where $SSE = 0$

② SSR should be equal to SST i.e. $SSR/SST = 1$

③ BFL is a poor fit \Rightarrow mean SSE is large and SSR/SST will be close to zero.

- $\frac{SSR}{SST}$ it is nothing but R^2 or we can say R squared value.

- It is also called by Coefficient of determination.

$$\text{Coefficient of determination } R^2 = \frac{SSR}{SST} = \frac{1 - SSE}{SST}$$

Ideal value of $R^2 = 1$

R^2 \downarrow \rightarrow 0 to 1

R^2 \downarrow \rightarrow -1 to +1

- SSE is always less than total error (SST) i.e. $SSE < SST$
- $R^2 = 0.01$
 $R = -1.01$
- SSE is always less than total error (SST)

Memo No.

Date

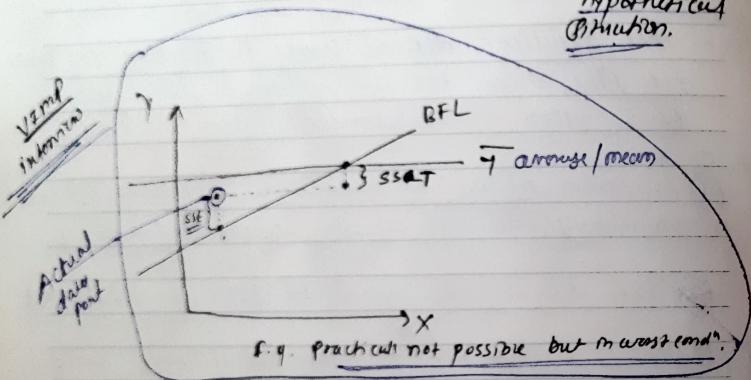
P

Allure's

(when $SSE = 0, R^2 = 1 \Rightarrow$ that is an ideal case.)

* when $SSE < SST, R^2 = 0 to 1 \Rightarrow$ that is an practical case.

when $SSE > SST, R^2 = \text{negative value} \Rightarrow$ that is very hypothetical situation.



- when BFL is worst than average/mean line, in that case SSE will be greater than SST.
- R^2 is used to check the goodness of BFL, for feature which good predictors, R^2 will be high. for feature which bad predictors, R^2 will be very low/close to zero.

Disadvantage of R^2 is not a good metric for evaluation of models, because it easily used to get impacted

Adjusted R^2

by Elusive relationships.

Sometimes Value of R^2 get increases for a bad predictor as well, in that case we go for another concept called adjusted R^2 square value.

	R^2	Adjusted R^2
Initially	0.85	0.84
Good predictor	0.87	0.86
Bad predictor	0.875	0.80

some time value of R^2 increases for a bad predictor as well, in that case we go for another concept called adjusted R^2 value, which gets increase only when the good predicting has been to the model.

minimum the disadvantage of R^2

$$R^2_{adj} = 1 - \frac{(1-R^2)(n-1)}{n-p-1}$$

$R^2 = R$ - squared = coefficient of determination
 $n = \text{no. of samples/rows in the dataset}$
 $p = \text{no. of predictors / feature}$

Adjusted R^2 gets increases the only when the good predictor

$R = \text{coefficient of correlation}$



Mo Tu We Th Fr Sa Su

Assumptions of Linear Regression

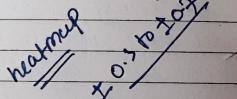
part 8/other part

1. Assumption of linearity: There should be a linear relationship between Dependent & independent variable.

~~(D) If R² is approx. zero \Rightarrow Non-linear. i.e. there is no linear relationship, it does not mean that there is no relationship at all.~~

- It is a property of mathematical relationships, which graphically represented as a straight line.

2. Small/no multicollinearity

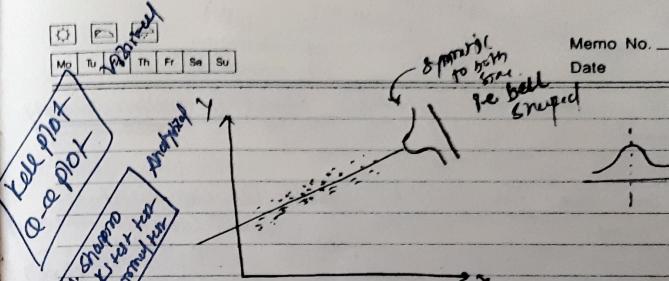
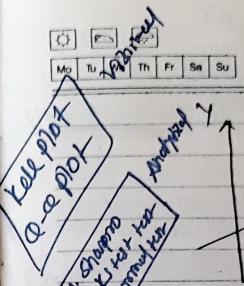


~~Small or no multicollinearity~~ It is a concept in which there is very little correlation between two independent variables. Due to multicollinearity it is difficult to get the relationship between target variable & predictor variable, or independent variable.

$x_1, x_2 \gg$ independent variables (linearity not)
 $y \gg$ dependent variable

3. Normal distribution in the error terms

~~Normal distribution of errors~~ - There should be normal distribution in error terms.
bell shaped curve
 $\mu = 0$
 $\sigma^2 = 1$



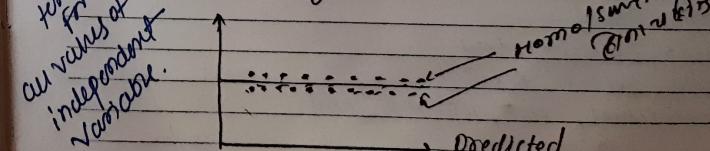
If we check variance on x-axis, we should get bell-shaped value (ideally). Means we should get as much good bell-shaped curve as possible!
mean = 0, std = 1

4. Homoscedasticity: (Same Variance)

Homo \Rightarrow same

Scadasticity \Rightarrow Variance (error)

~~same error term for all values of independent variable.~~



- error term should be almost same for all values of independent variable.

Date: 1/1
 Memo No.: 101
 R = coefficient of correlation
 $R^2 = 0.1$
 How to understand whether there is multicollinearity or not?
 There are few concepts by which we can get this thing:
 1. Correlation matrix (we have seen by using corr)
 $R = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m (x_{ij} - \bar{x}_i)(y_{ij} - \bar{y}_j)$
 2. VIF (Range of VIF is 1 to ∞) any possible value if VIF ≥ 5 > independent to each other / moderately independent if VIF ≥ 10 > greater than others is an issue of multicollinearity (or) we can say that independent variables are highly dependent.
 $VIF = \frac{1}{R^2}$
 1 to 5 -- moderately independent
 > 5 -- Highly independent. Correlate 1 -- Not correlated
 $VIF = 1 \rightarrow 5 \rightarrow$ Completely independent

Ordinal values
 i) high, medium, low
 ii) good, better, best
 iii) High school, HSE, BE, ME, PhD
 Nominal values
 pune, Delhi, Mumbai etc.

Day-55

Encoding (we are discussing about categorical variable from sklearn.preprocessing import LabelEncoder)

Label encoding : Temperature (original value) 35/05/2022
 Temperature (from sklearn.preprocessing import LabelEncoder) 0
 Temperature (after applying label encoder) 1
 i.e. Categorical variable int/float values
 original value (After applying label encoder) continuous feature
 0 (Before applying label encoder) discrete

① Label encoding : - It is applied on target variable / Y-axis
 e.g. high (Original value) 2 (After applying label encoder)
 medium 1
 low 0

② One hot encoding
 e.g. City_pune. City_Delhi. City_Mumbai
 pune 1 0
 delhi 0 1
 mumbai 0 0

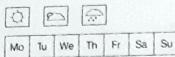
③ Nominal encoding
 We do not need to consider about association between categories because categories will have their own rank

e.g. label encoding
 $df = pd.DataFrame(\{'city': ['pune', 'mumbai', 'delhi', 'Delhi', 'pune', 'Mumbai']\})$
 df

0	City
0	pune
1	mumbai
2	delhi
3	Delhi
4	pune
5	Mumbai

e.g. one-hot encoding
 $df_city = pd.get_dummies(df['city'])$
 df_city

0	City_Delhi	City_mumbai	City_pune
0	0	0	1
1	0	1	0
2	1	0	0
3	0	0	0
4	0	1	0
5	0	0	0



06v-59

Memo No.

Date

6/6/22

This algorithm is the best suite for our model building.

8. Linear regression is prone to over-fitting, but with the help of some dimensionality reduction technique, regularization methods (L1 and L2) it can avoid.

Disadvantages of Linear Regression Algorithms

1. It gets effect because of outliers.
2. It has many assumptions.
3. Sometimes a lot of feature engineering is required.
4. Often get overfit.
5. It is sensitive to missing values.

Applications

1. House price prediction
2. Temperature prediction
3. Salary prediction
4. Stock price prediction
5. Weather forecasting
6. Cricket score prediction.

Cost

Q P Used for only classification problem
Mo Tu We Th Fr Sa Su not for regression.
Day-60

Classification problems
measurable dependent variable
in categorical form.
Memo No.

Date

2) Logistic Regression

- It is supervised type of machine learning algorithm for regression problem. Still the name given is logistic regression.
- Because of at the back end it also uses linear model (Or) we can also say that classification method uses the same concept as linear regression, that is the reason it is called as logistic regression.
- The word logistic represent here, the logit function. / Cost func (they shows error)

- Target variable should be in the categorical form. It can be two or more categories.

It does have two categories
→ spam/not spam
→ Yes/no
It does have more than two categories
→ low/high/medium >> ordered
→ mango/apple/banana >> multiclass
classification

① Binary Class Classification

e.g Suppose we do have two classes in target variable

Class 1

Class 2

>> Appraise

Decline

prob (1/0)

Model o/p

True sample (1)

False sample (0)

Boundary

Prob

Model o/p

True sample (0)

False sample (1)

cross function (or) log function aim to minimize error

shows the error
linear regression \Rightarrow cross sum (MSE) $\frac{1}{n} \sum (y_i - \hat{y}_i)^2$
logistic regression \Rightarrow log's sum (logloss) $\frac{1}{n} \sum [y_i \ln(\hat{y}_i) + (1-y_i) \ln(1-\hat{y}_i)]$

model.predict \Rightarrow class

It can also be used to find out probability which an observation belong to either class 1 or class 2

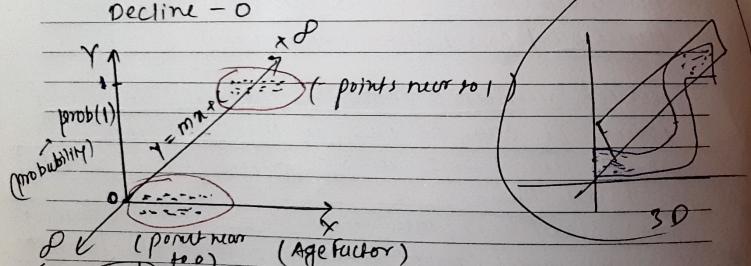
model.predict_proba \Rightarrow it predict the probability not class

Decision tree, logistic regression, Naive Bayes algorithm.

e.g. Loan Application Status

Approve - 1

Decline - 0



Linear model \Rightarrow probability always towards $+0$ & -0
- But here BFL line not came due to decision making tree. so we have solution that is sigmoid function.

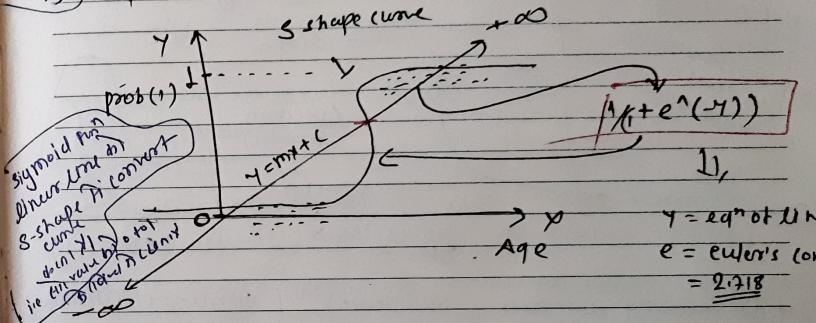
i.e. in linear model line can cross the boundary of class but BFL not come due to decision

Memo No. _____
Date 1/1

$$\frac{1}{1+e^{(mx+c)}} \quad s(y) = \frac{1}{1+e^{-y}}$$

Sigmoid function equation

$$\frac{1}{1+e^{-y}} \Rightarrow 1/(1+e^{(-y)})$$



$$y = \text{eq of line}$$

$$e = \text{euler's constant}$$

$$= 2.718$$

- Sigmoid function simply tries to convert independent variable into an expression of probabilities that ranges in between 0 & 1 with respect to dependent variable.

- In linear model probability always tends towards $+\infty$ & $-\infty$.

- Sigmoid function also an activation function in machine learning which is used to add non-linearity in a machine learning model. In simple words it decide which value to pass by 0 & what not to pass

Day - 61

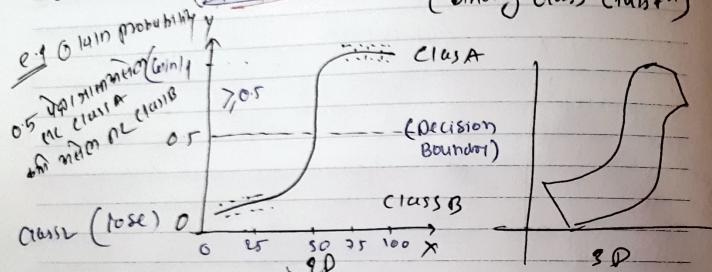
Memo No.

Date 09/06/22

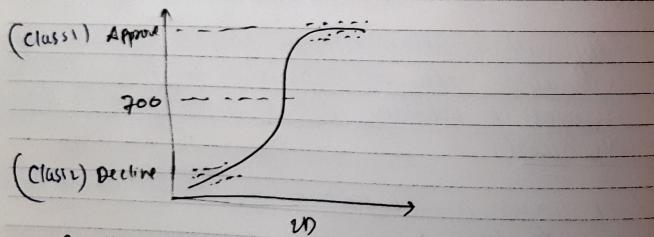
S(1) (M) & F(1)

- with the help of sigmoid function, it converts one linear line to S-shaped curve which is called as Sigmoid curve.

(binary class classifying)



e.g. Banking Domain (Cibil score) (binary classification)



If Cibil score - $> 700 \Rightarrow$ Approve
If Cibil score - $< 700 \Rightarrow$ Decline
even though Cibil score $> 700 \Rightarrow$ Decline (salary/income per month)

Binary class classification either ~~object~~ will belong to Class A or Class B

② Multiclass classification

Multiclass classification how it works

C1

C2

C3

It uses a concept called OVR \Rightarrow one versus rest.
(one versus all)

C1 or Not in C1
(C1 Vs C2 and C3)

either C1 or not in C1
or not in C1
ie C1 vs C2 & C3

- Logistic regression is about two steps

① The propensity to belong to class 1 is $P(Y|X)$

② It basically uses a cut-off or threshold to decide the class

③ It finds a BFL/Plane, then it goes through sigmoid transformation and convert it into S-shaped curve.

- That the S-shaped which has been formed can be shifted to right-left and upward & downward

It happens on the basis of m and c value.

up and down \Rightarrow on the basis of c value
right & left \Rightarrow m value.

- So basically here with values of m & c we get to know that which Sigmoid surface/curve/plane is best for our model.

$$-\frac{1}{M} \sum_{i=1}^M y_i \cdot \log(p(y_i)) + (1-y_i) \cdot \log(1-p(y_i))$$

Memo No. _____

Date _____

and to actually find out which is the best sigmoid surface, it uses one function which is called logloss function.

Logloss function

- to find out which is the best sigmoid surface, it uses one function which is called logloss function.

Logloss is indicative of how close the predicted probability is to the corresponding actual/pure values.

$$\begin{aligned} \text{logloss} &= -\frac{1}{M} \sum_{i=1}^M (y_i (\log p_i) + (1-y_i) \log(1-p_i)) \\ &= -\frac{1}{M} \sum_{i=1}^M y_i \log p_i + (1-y_i) \log(1-p_i) \end{aligned}$$

Y - Target variable, which ranges in between 0 to 1
P - Is our predicted probability.

$$y_i (\log p_i) + (1-y_i) \log(1-p_i) \Rightarrow \begin{cases} y_i (\log p(y_i)) + (1-y_i) \log(1-p(y_i)) \\ \text{probability } 1 \\ \text{probability } 0 \end{cases}$$

Case 1: Y=1, P=1 (high)
↳ correct classification

$p(y_i)$ is the probability of 1
 $1-p(y_i)$ is the probability of 0

Memo No. _____
Date _____

* $(Y=1, p \text{ is high} \Rightarrow \text{second part our eqn will be } 0)$
and as here p is high then $\log p_i = 0$
The whole logloss function will be 0
This is exactly basically a correct classification.

import numpy as np

np.log(1)

↳ 0.0

np.log(1)

↳ -inf

* Case 2: Y=0, P=0 (low) \Rightarrow correct class

$y=0$, first part will be 0, p is low

$(1-y_i) \log(1-p_i)$

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

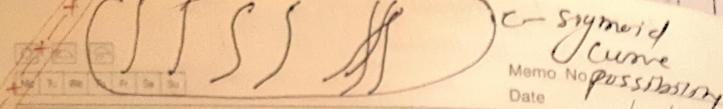
This is basically a correct classification

Case 3: Y=1, P=0 \Rightarrow incorrect class

$y=1, p \text{ is low}$

$y_i (\log p_i)$

logloss function will not turn to 0, we will get some value here.
This is incorrect classification



Case 4: $y=0, p=1 \Rightarrow$ incorrect classification.

$y=0, p$ is high

$(1-y_i) \log(1-p_i) \gg 1 \cdot \log(1-1)$
logistic function will not sum to 0 here.
This is incorrect classification

Day-62

Assumptions of logistic regression

- Dependent variable should be categorical in nature.
- No multicollinearity. ("in bet" independent variable)
- It construct linear boundaries.

(It means)



Day-62

Memo No. Date 10/6/22

peaks and distribution curve of two variables
should be separated from each other. If they are overlapped with each other, logistic regression will not do good.)

- Data size should be large (No. of observations should be large)

~~VRML in 10.1 classification problem~~

Evaluation Metrics

↳ Various evaluation metrics

1. Accuracy

VRML 2. Confusion matrix

3. Log-loss $\Rightarrow -\frac{1}{N} \sum y_i \cdot \log p(y_i) + (1-y_i) \cdot \log(1-p(y_i))$

4. precision and recall

5. F1-Score $\Rightarrow F_1(\text{beta}) \Rightarrow \frac{\text{beta}}{1+\text{beta}}$

6. AUC-ROC curve

VRML performance metric

Evaluation metrics quantify the performance of a predictive model or measure the quality of ML model.

1. Accuracy \Rightarrow Ratio of no. of correct prediction & total no. of prediction

- Accuracy is ~~ratio~~ measure of how often the classifier makes the correct prediction.

- It is a ratio between the no. of correct predictions and the total no. of predictions.

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