

Machine Learning

(main)

Jupiter of ML

Unsupervised ML

Regression

- Linear

- Ridge & Lasso

- Polynomial Reg.

- Multi linear

- Non-linear

- OLR

- Time series forecasting

Classification

- Logistic

- KNN

- DT

- Naive Bayes

- RF

- XGBoost

- GBM

- SVM

- ANN

catboost, adaboost

Supervised ML →

We will start with Regression → linear

Linear Regression

Simple

linear Regression

(SLR)

Multiple

linear

Regression
(MLR)

polynomial

Regression

- Simple linear Regression (SLR) :- when for continuous output we have only one feature then it is called SLR

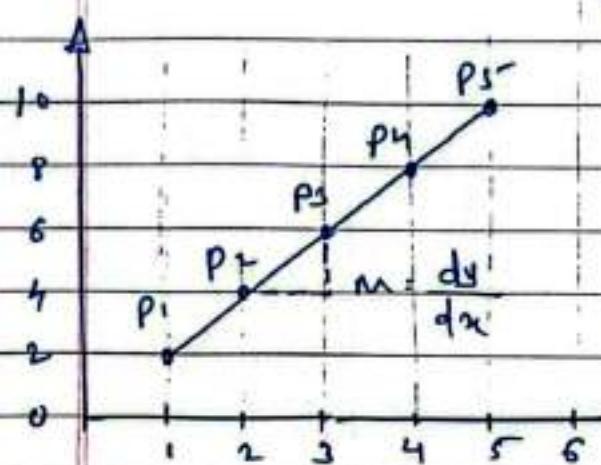
for Ex. CGPA | package

6.0 2 LPA

7.2 3 LPA

8.9 6 LPA.

Basics of linear Regression.



Equation of straight line

$$y = mx + c$$

find value of m and c

$$P_2 (2, 4) \rightarrow x_1, y_1$$

$$P_3 (3, 6) \rightarrow x_2, y_2$$

$$m = \text{slope} = \frac{y_2 - y_1}{x_2 - x_1} = \frac{6 - 4}{3 - 2} = 2$$

$$m = \frac{dy}{dx} = 2$$

slope means with every change in value of x how much changes in y

Intercept = c

$$P_4 (x_3, y_3) \rightarrow (4, 8)$$

$$y = mx + c$$

$$8 = 2(4) + c$$

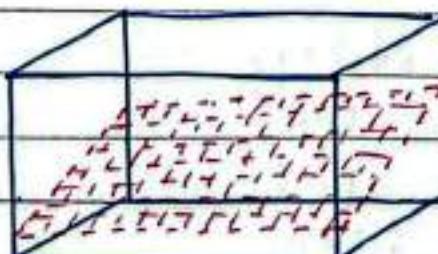
$$8 - 8 = c = 0$$

Inference : The above line eqn is function that relates x and y .

for given value of x we can find corresponding value of y .

what if we have two or more than two independent variable ?

→ then it is called multiple linear regression



(3D)

simple linear regression : $y = mx + c$

multiple linear regression

$$y = a_0x_0 + a_1x_1 + a_2x_2 + \dots + a_nx_n + b.$$

Advantages : 1) simple to implement

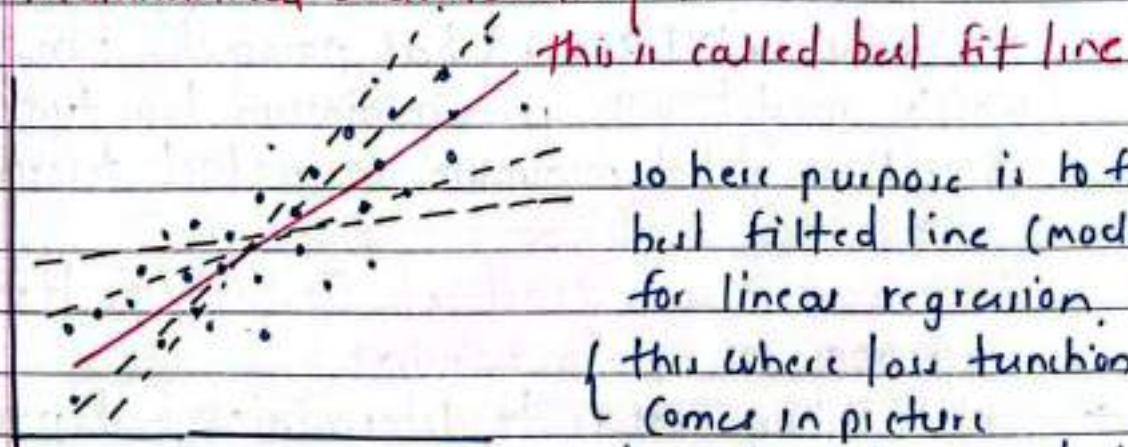
2) perform well on the data with linear relationship.

Disadvantages : 1) Not suitable for data having non-linear relationship.

2) Underfitting issue

3) sensitive to outliers

Mathematical Understanding :-



So here purpose is to find
best fitted line (model)
for linear regression.

{ this where loss function
comes in picture }

- loss function measures how far an estimated value from its true value.
- it is helpful to determine which model performs better and which parameters are better.
(m, c)

$$\text{Loss} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Let try to understand with the ex. $m = 3, c = 2$

$$\hat{y} = 3x + 2$$

x	4	4
2	10	8
3	14	11
4	18	14
5	22	17
6	26	20

Now find its loss function.

$$\begin{aligned} \text{Loss} &= \frac{[(10-8)^2 + (14-11)^2 + (18-14)^2 + (22-17)^2 + (26-20)^2]}{5} \\ &= \frac{4+9+16+25+36}{5} = \frac{90}{5} = 18 \end{aligned}$$

{ reason to take square: so that positive and negative value may not cancel each other }

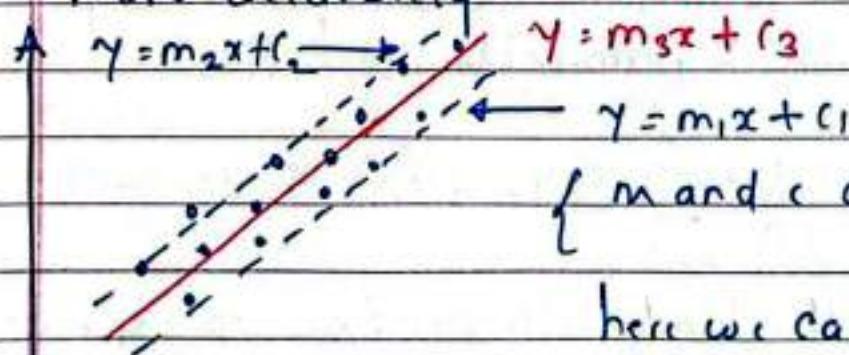
low loss value \rightarrow High Accuracy

high loss value \rightarrow low Accuracy

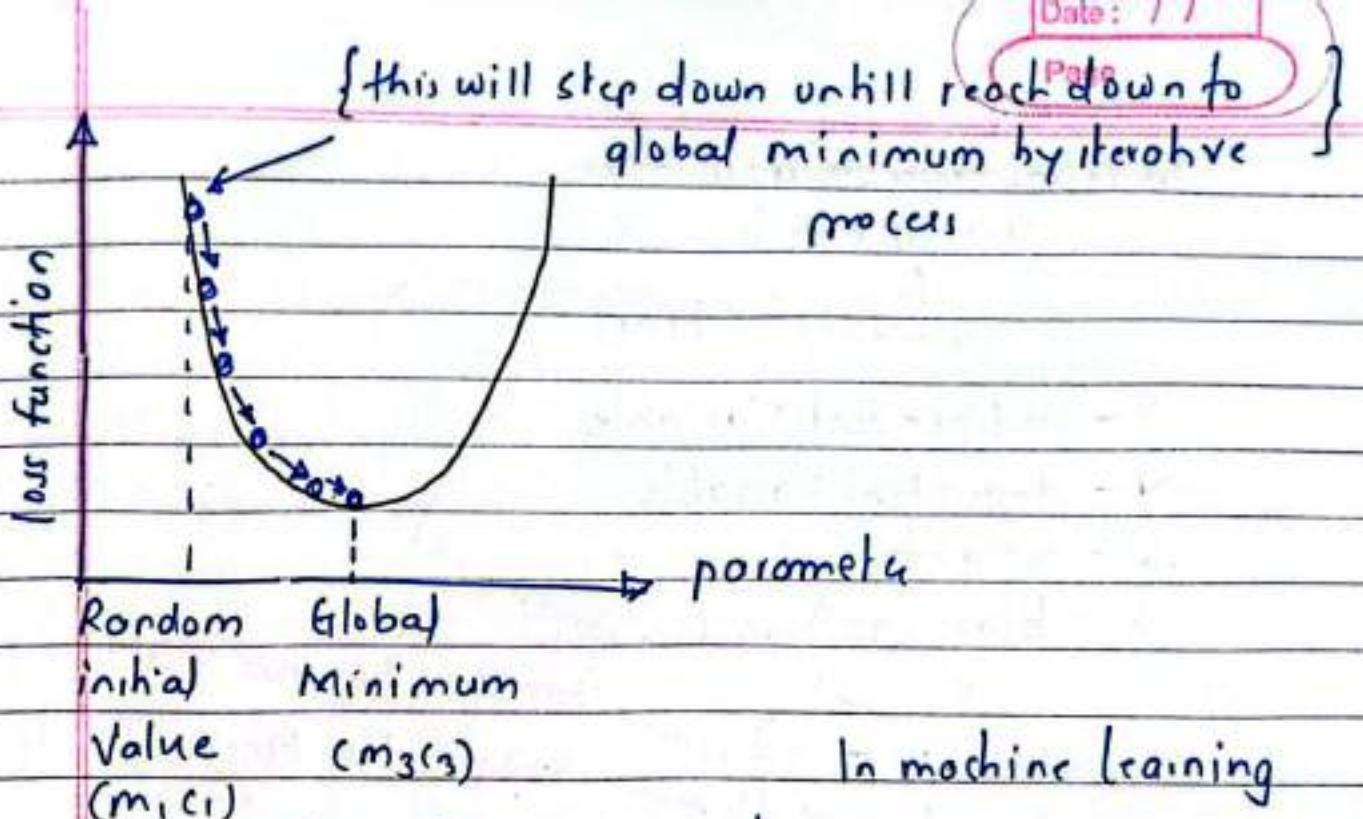
We can improve the model by some optimization technique called as "gradient descent" where repeat the process iteratively until we get best parameter (m, c) for which model will give minimum loss function. (called as "global minimum" in "gradient descent")

How we can use gradient Descent for linear regression for optimization.

→ optimization refers to determining best parameter for model such that loss function of the model decreases as result of which model can predict more accurately.



here we can find model 3 is best fit since the loss function is least and thus this model is optimum this where we use gradient descent to optimize



In machine learning

$$\begin{aligned}
 m &= m_i - L D_m & w &= w - L d_w \\
 c &= c_i - L D_c & b &= b - L d_b \\
 \text{m - slope} && w &= \text{weight} \\
 \text{c - Intercept} && b &= \text{bias.} \\
 \end{aligned}$$

L - learning rate : it is magnitude of change that you want in parameter during iteration.

D_m : partial derivative of loss function w.r.t m

D_c : partial derivative of loss function w.r.t c.

$$\begin{aligned}
 D_m &= \frac{\partial (\text{loss function})}{\partial m} = \frac{\partial}{\partial m} \left(\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \right) \\
 &= -2 \sum_{i=0}^n x_i (y_i - \hat{y}_i)
 \end{aligned}$$

$$\begin{aligned}
 \text{Ily } D_c &= \frac{\partial (\text{loss function})}{\partial c} = \frac{\partial}{\partial c} \left(\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \right) \\
 &= -2 \sum_{i=0}^n (y_i - \hat{y}_i)
 \end{aligned}$$

Intercept, weight and bias

$$y = w_0 + w_1 x_1 + b$$

w_0 (weight) b (bias)

x - independent Variable

y - dependent Variable

w - weight

b - bias

$$\begin{aligned} w &= w - \alpha \frac{\partial L}{\partial w} && \text{change in Loss function w.r.t } w \\ b &= b - \alpha \frac{\partial L}{\partial b} && \text{change in Loss function w.r.t } b. \end{aligned}$$

Learning rate :- it is tuning parameter in an optimization algorithm that determine step size at each iteration while moving toward minimum of loss function.

$$\frac{\partial L}{\partial w} = -\frac{2}{n} \sum_{i=0}^n x_i (y_i - \hat{y}_i)$$

$$\frac{\partial L}{\partial b} = -\frac{2}{n} \sum_{i=0}^n (y_i - \hat{y}_i).$$

for **Multiple linear Regression**. for one target we have two or more than two independent variables.

prediction Equation for MLR

$$\hat{y} = q_1 x_1 + q_2 x_2 + \dots + q_n x_n + b.$$

and,

Regression Equation

$$y = \underbrace{q_1 x_1 + q_2 x_2 + q_3 x_3 + \dots + q_n x_n + b}_{\text{predicted value}} + \epsilon$$

Actual

Value .

What are assumptions of linear Regression?

there are five main assumption in linear regression.

1. Linear Relation between input and output.

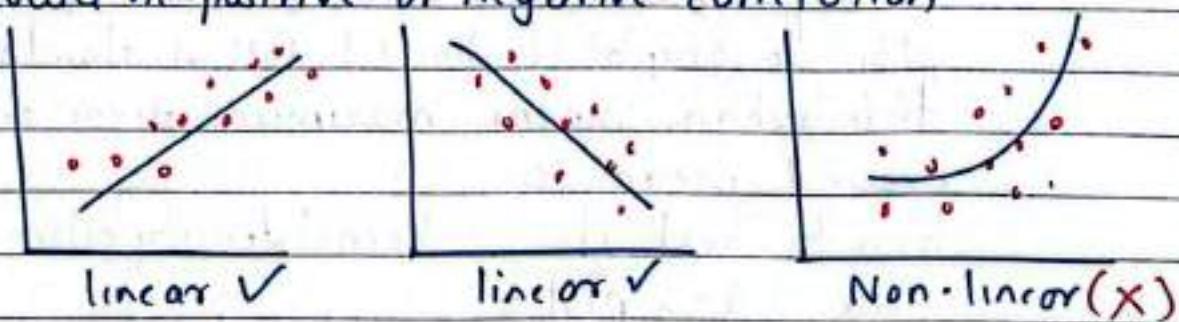
2. No multicollinearity.

3. Normality of Residual.

4. Homoscedasticity.

5. No autocorrelation of error

1. Assumption 1: there should be linear relation between individual feature and target (output) it could be positive or negative correlation



applicable to multiple linear problem as well. Not applicable

how to check this : scatter plot : (feature 1 v/s target) (feature 2 v/s target).

2. Assumption 2: Multicollinearity: it mean there all the feature should be independent or should not have any correlation among themselves.

why, what the problem?

In multi linear Regression Model for 3D we draw a hyperplane.

$$y = q_1x_1 + q_2x_2 + q_3x_3 + b$$

where q_i represent what will be the change in y with respect to x_i assuming x_2 and x_3 constant. but if it violate this assumption then model will not perform good. (Ex of two phys sci)

how to check multicollinearity :-

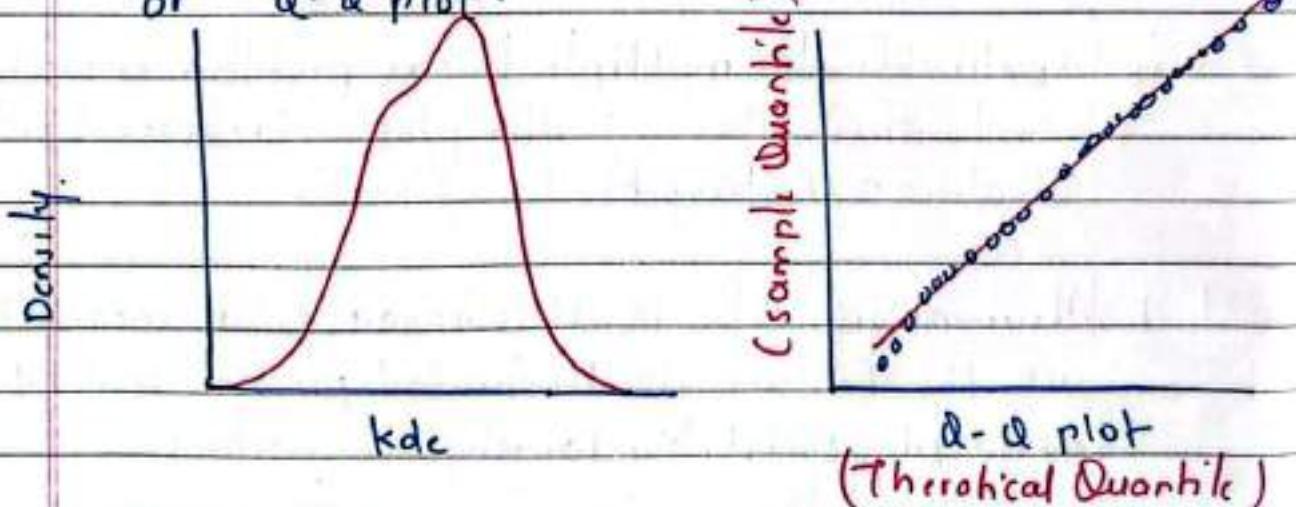
- 1) VIF (Variance Inflation factor) :-
if it is around 1 then features don't have the issue
and if it is 5 or more than that then that
particular feature has multicollinearity issue and
need to remove it.
- 2) another method is to find out correlation between all
the features (Heatmap)

3 Assumption 3. - Normality of Residual.

it says that when error (actual - predicted)
plot or graph, it should follow standard normal
distribution, means maximum error should
around mean = 0.

how to check it : kernel density estimation (kde)

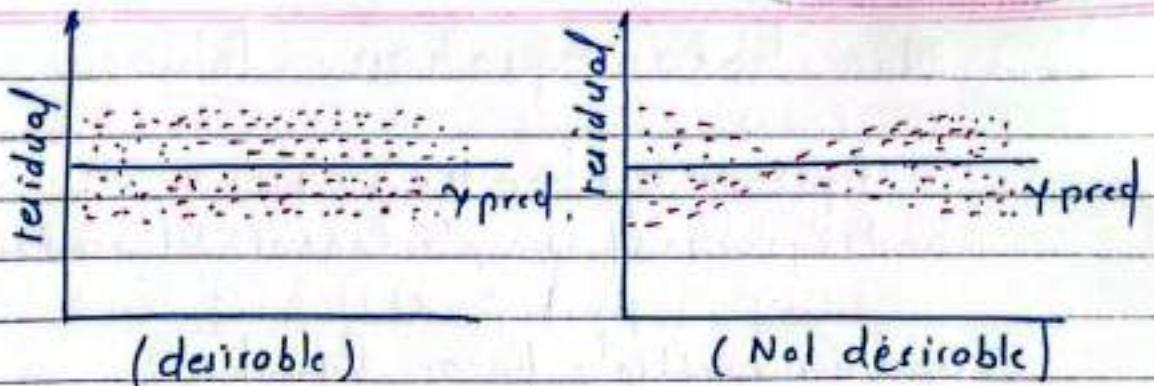
or Q-Q plot



4. Assumption 4: Homoscedasticity.

some scattered spread.

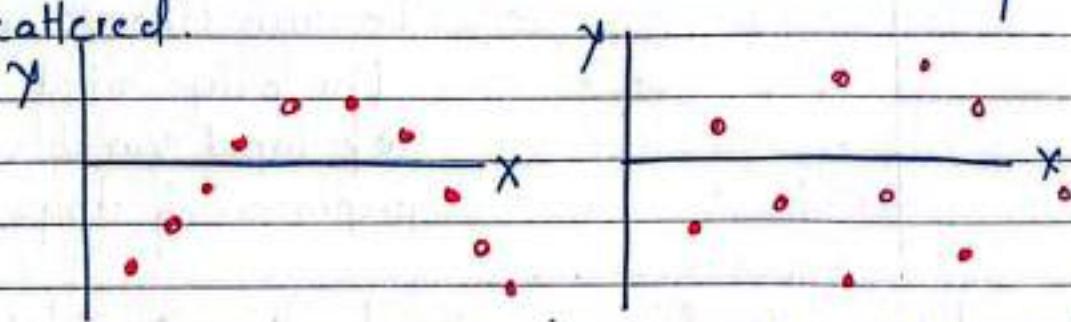
: spread of residual should be equally or uniformly
scattered. If it is not, it called as heteroscedastic
shity which is not desirable.



how to check:- scatter plot (y^{pred} , residual).

Assumption 5: No autocorrelation of error.

If we plot the error there should not be any specific pattern instead it should randomly scattered.



positive Autocorrelation
(Not desirable)

Negative Autocorrelation
(desirable)

how to check : p11. plot (residual)

Summary chart

Assumption	Severity	Prediction	Inference
------------	----------	------------	-----------

1) linear Relation	High	✓	✓
--------------------	------	---	---

2) Multicollinearity	Medium	✗	✓
----------------------	--------	---	---

3) Normality	low	✗	✓
--------------	-----	---	---

4) Homoscedasticity	High	✓	✓
---------------------	------	---	---

5) Autocorrelation of	-	-	-
-----------------------	---	---	---

Non-linear Equation - Polynomial Regression

we know

$$\text{Equation of line } y = mx + c$$

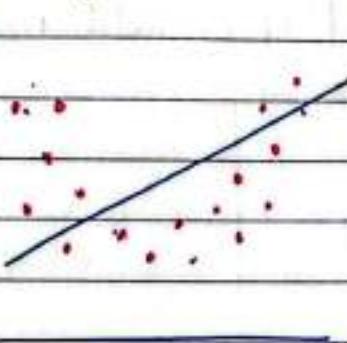
and Equation of simple Linear Regression

$$y = \beta_0 + \beta_1 x$$

and for multiple linear Equation.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n$$

this is applicable, only when data is linear but what if data is not linear?



In such scenario we extract the polynomial feature of input variable in preprocessing stage

let say for ex $x | y \rightarrow 5 | 10$

- For x we want to make polynomial of degree 2 then we will convert $x \rightarrow x^0, x^1, x^2$, y
so, it will be 1, 5, 25
- this way we create a new data for training the extra polynomial feature try to extract their non-linear relationship.
- its formula become for simple polynomial regression.

$$y = \beta_0 + \beta_1 x + \beta_2 x^2$$

for degree 3

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$$

* Now how we will know the perfect value for the degree — since this is hyperparameter.

if we keep it low then may be it cause underfitting means, it may be not able to learn the all attributes. and if we select very high then there is chance of overfitting or overlearned that's why our job is to find out optimum value

- In case if we have two features x_1, x_2, Y then for degree 2 our simple polynomial Equation would become

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_2 + \beta_4 x_2^2$$

Interview Ques: why polynomial Eqⁿ is basically called as Linear Regression

Ans: when we talk about linear Regression we talk about relation between y and coefficients of features and degree of coefficient is still one and thus relation between y and coefficient is still linear

Ordinary least square :- (OLS Algorithm)

- it is method for estimating the parameters of linear regression model.
- its aim to find the values of the linear regression model's parameters (i.e coefficient) that minimize the sum of squared residual.
- the residuals are differences between observed values of dependent variable and predicted values of dependent variable given w.r.t independent variable
- OLS Algorithm assume that the errors are normally distributed with zero mean and constant variance and that there is no multicollinearity (high correlation) among the independent variables.
- other method like generalized least square or weighted least square, should be used in case where these assumption are not met.

Let understand with problem

x	1	2	3	4	5	6	7
y	1.5	3.8	6.7	9.0	11.2	13.6	16.

We will calculate the equation for the best fit line where all the point will be as close as possible by least square method.

x	y	xy	x^2	
1	1.5	1.5	1	$\sum x = 28 \quad \sum y = 61.8$
2	3.8	7.6	4	$\sum xy = 314.8 \quad \sum x^2 = 140$
3	6.7	20.1	9	
4	9.0	36	16	$n = 7$
5	11.2	56	25	(number of data points)
6	13.6	81.6	36	
7	16	112	49	
Σ	28	61.8	314.8	$y = mx + b$
			140	

$$m = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2} : \frac{7(314.8) - (28)(61.8)}{7(140) - (28)^2}$$

$$m = \frac{473.2}{196} : \approx 2.4142857.$$

$$b = \frac{\sum y - m \sum x}{n} = \frac{61.8 - 2.4142857(28)}{7}$$

$$b = -0.828571.$$

To get linear equation we should plug value. In

$$y = mx + b.$$

~~$y = 2.41x + b$~~ $y = 2.41x - 0.83$

I.e. take it for x \hat{y} y_{act}

$y = 2.41(2) - 0.83 = 3.99 \quad 3.8$

$y = 2.41(5) - 0.83 = 11.22 \quad 11.2$

$y = 2.41(7) - 0.83 = 16.09 \quad 16.$

Syntax: statsmodel.api.OLS(y, x)

y - dependent variable x - independent variable

- Import statsmodel.api as sm
import pandas as pd.

- # reading data from csv
 $dt = pd.read_csv('train.csv')$

- # defining the variables

$x = dt['x'].to_list()$

$y = dt['y'].to_list()$

adding the constant term in
x = sm.add_constant(x)

performing regression and fitting model.
result = sm.OLS(y, x).fit()

print summary table.
print(result.summary())

What is Ridge Regression (L2 Regularization).

Ridge regression is model tuning method, that is used to analyse any data that suffers from multicollinearity. This method performs L2 regularization when issue of multicollinearity occurs, least square are unbiased, and variance are large this result in predicted values being far away from actual values.

Regularization:- it is technique used to calibrate machine learning model to minimize adjusted loss function and avoid overfitting and underfitting.

there are three type of regularization techniques:

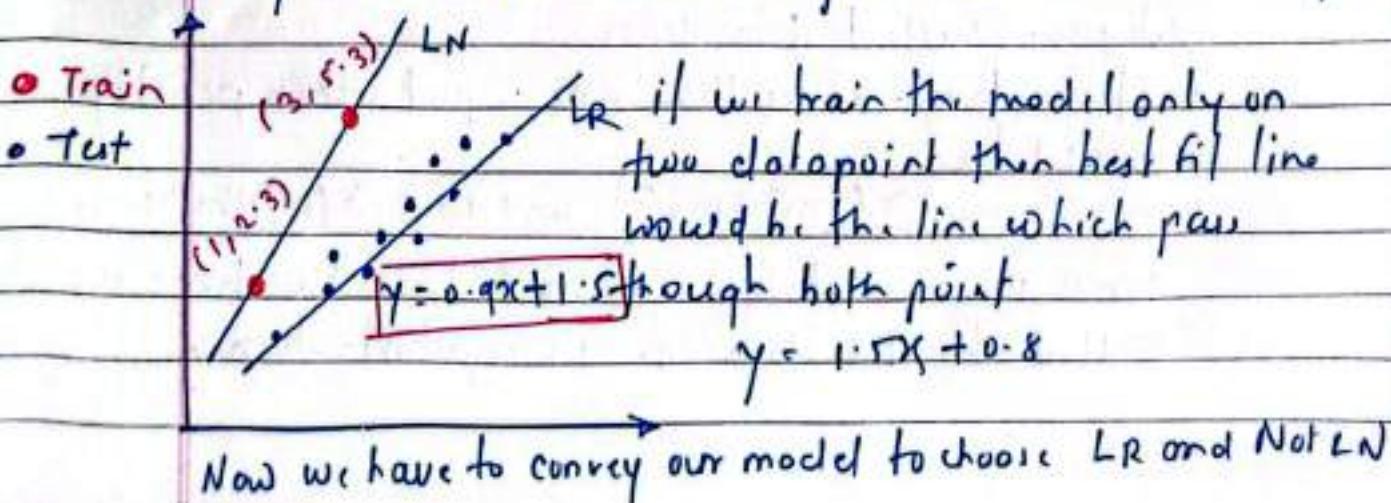
- 1) Ridge Regression (L2 regularization)
- 2) Lasso Regression (L1 regularization)
- 3) Elastic Net (Combo of Ridge and Lasso)

Ridge Regression:-

Overfitting:- means your train acc is too high but test accuracy is very low.

$$y = mx + b$$

m or slope which is coefficient of x is defining change in y , so to reduce overfitting means to reduce slope.



Now we have to convey our model to choose LR and Not LN

$$L = \sum_{i=1}^n (\hat{y}_i - y_i)^2 + \lambda (m^2) \quad \text{--- (1)}$$

Now we add penalty before we will only reducing which will add help this term which is for the first term nothing but Error / γ = hyper-parameter Residual m = slope.

Now we will calculate loss for both line for $\lambda = 1$ with eqn (1)

Since line is passing through both points perfectly that's why 1st term will be zero

$$L = 0 + 1(1.5)^2 \\ = 2.25$$

$$(2.3 - 0.9 - 1.5)^2 + \\ (5.3 - 2.7 - 1.5)^2 + (0.9)^2 \\ = (0.1)^2 + (1.1)^2 + (0.9)^2$$

$$= 2.03$$

here we getting significant reduction in loss for the new line

As our model can see this change it will select 2nd model although it will give bad accuracy on training since Variance is significantly reduced although bias increased.

Why we called L² since

- if we have more than one input then penalty would be

$$\lambda(m_1^2 + m_2^2) \text{ and for } 3 \lambda(m_1^2 + m_2^2 + m_3^2)$$

since we are doing square (^2) all the terms we called it L² Norm or L² regularization.

Lasso Regression (L1 Regularization)

- this is also help to reduce overfitting.
- In Ridge regression we have seen:

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \|w\|^2$$

MSE penalty term
 $(w_1^2 + w_2^2 + \dots + w_n^2)$
Weight in mLR

Lasso is just another variation of Ridge

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \|w\|$$

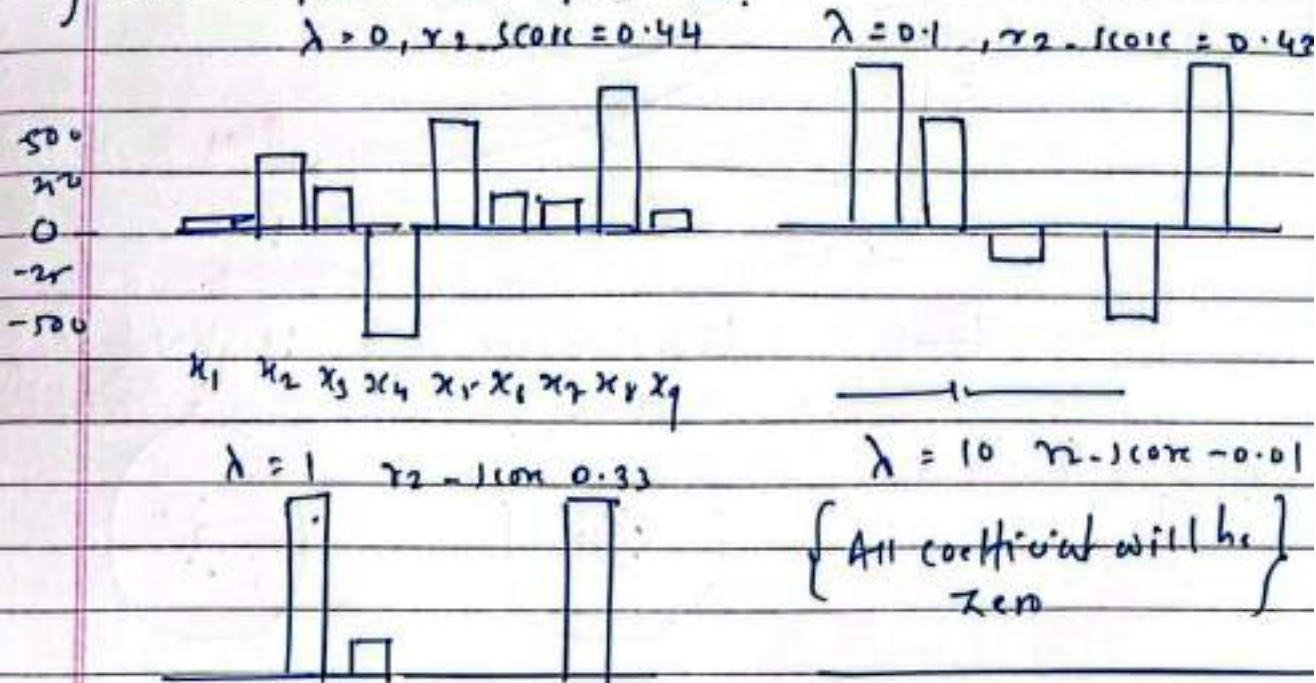
OR

$$= \text{MSE} + \lambda [|w_1| + |w_2| + |w_3| + \dots + |w_n|]$$

- In ridge regression for any value of λ there were always some value for coefficient of input feature but in Lasso, if you continuously increase value of λ at certain point you will zero value for some of coefficient which are not important so here we unknowingly doing feature selection and it is advantage of Lasso.
- so when you are working on high dimensional data and some feature are not imp we should prefer Lasso over Ridge.

there are some key points need to discuss about Lasso.

- How coefficients are affected?



- Higher coefficients are affected more

- Generally as you increased λ coefficient value will be decreased gradually toward zero
- usually higher coefficients are affected first/rapidly
- for it we should be cautious for selecting optimum value to do proper feature selection

- Impact on Bias and Variance

- As we know if λ increases, overfitting decreases (\downarrow) which lead to increase in bias (\uparrow)

High Bias



Underfitting

High Variance



Overfitting

fig a.

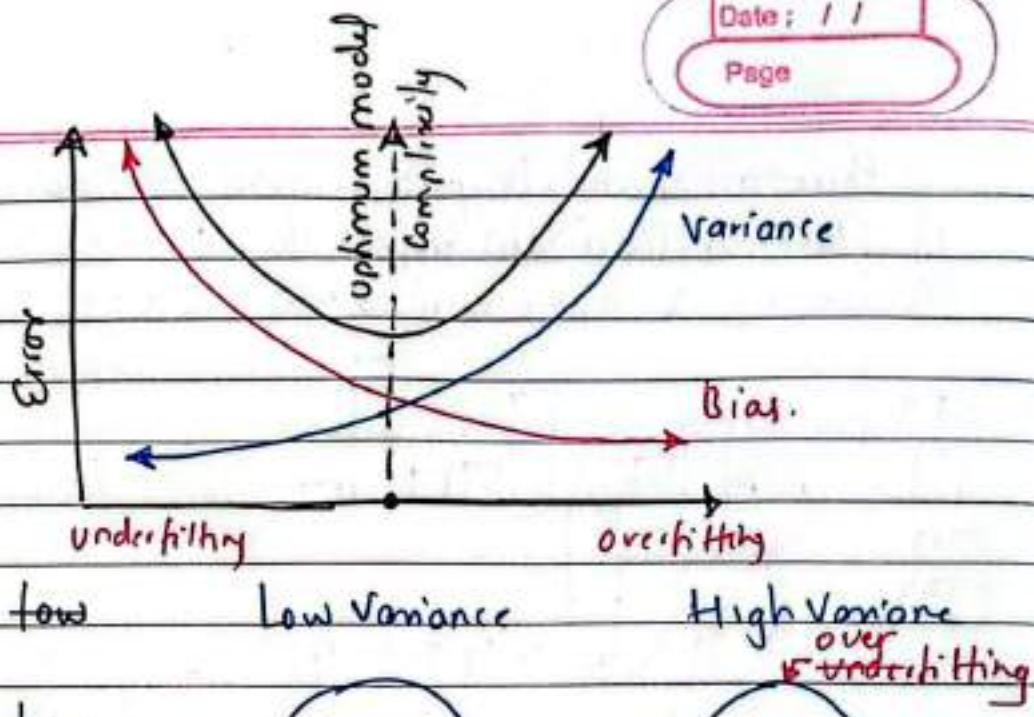
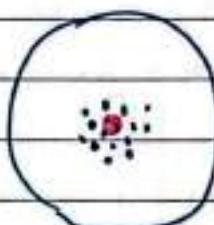


fig b

low bias



High bias

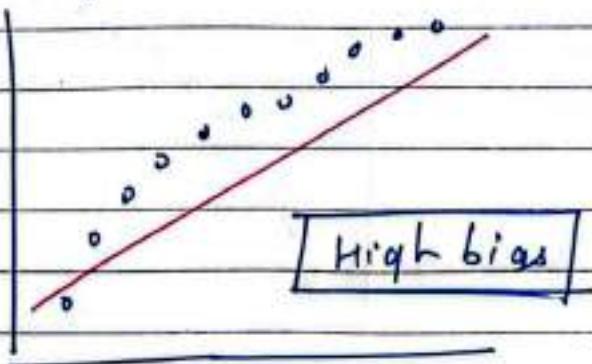


Bias - Variance Tradeoff .

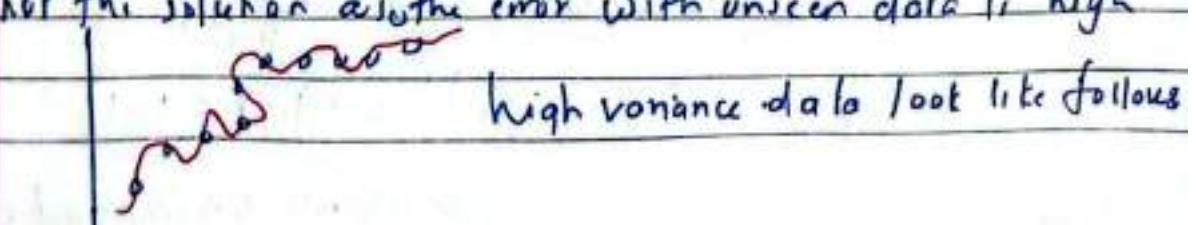
- It is important to understand prediction errors (bias and variance) when it comes to accuracy in any ML Algorithm
- There is a tradeoff between model's ability to minimize bias and variance which is referred to as best solution for selecting a value of regularization constant
- Proper understanding of these errors would help to avoid the overfitting and underfitting of a dataset while training algorithm.

Bias: Bias is known as difference between prediction values by ML model and correct values. Being high in biasing will give large error in training as well as testing. & that's why it is always recommended that algorithm should always be low biased to avoid underfitting.

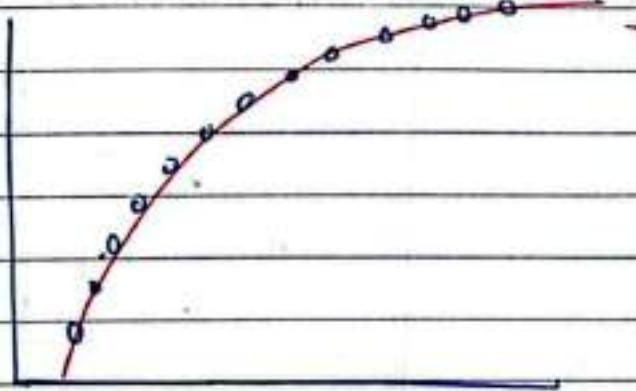
- if high bias data is predicted in straight line format, thus not fitting accurately in the data in the dataset such fitting called as underfitting.
- This happens the hypothesis is too simple or linear in nature



- Variance**
- The variability of model prediction for given data point which tells us spread of our data is called Variance of model.
 - The model with high variance has a very complex fit to training data and thus not able to fit on the test data or data hasn't seen earlier. Such model works very well on training data but has high error rates on test data.
 - When model is high on variance it is said to be overfitting of data. Overfitting is fitting the training set accurately via complex curve and high order hypothesis but is not the solution as the error with unseen data is high.



- Bias-Variance trade-off:- if the algorithm is too simple (hypothesis with linear eqⁿ) then it may be on high bias and low variance and thus it is error prone.
- if error fit too complex (hypothesis with high degree equation) then it may be on high variance and low bias, in the latter condition the new entries will not perform well. there is something between both of these condition known as tradeoff or bias-variance tradeoff



4) Effect of Regularization on loss function

loss function:-

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

- it measures how far an estimated value from its true value
- if we are training on different models LR, DT, RF to know which model performs better and which parameters are better loss function is useful.

Elastic Net :- $(L_1 + L_2)$: it is combination of both regularization techniques

$$L_{\text{Reg}} = \underbrace{\sum_n (y_i - \hat{y}_i)^2}_{\text{penalty which is imposed on normal}} + \underbrace{\lambda (|m| + |c|)}_{\text{hyperparameter}}$$

$$L_2 \text{ reg} = \frac{\sum (y_i - \hat{y}_i)^2}{n} + \underbrace{\lambda (m^2 + c^2)}_{\text{penalty}}$$

Elastic ($L_1 + L_2$) :-

$$\frac{\sum (y_i - \hat{y}_i)^2}{n} + \lambda [(|m| + |c|) + (m^2 + c^2)]$$

adding combination percentage of Lasso & Ridge

$$\frac{\sum (y_i - \hat{y}_i)^2}{n} + \lambda \left[\underbrace{c(|m| + |c|)}_{\text{Lasso}} + \underbrace{(1-c)(m^2 + c^2)}_{\text{Ridge}} \right]$$

c is between 0 to 1

$c=1$ Lasso

$c=0$ RIDGE

$c=0.5 = 50\% \text{ Ridge} \& 50\% \text{ Lasso}$

Summary :- Ridge is majority is going to focus on regularization but Lasso is going to focus on feature selection as well

- if my job is only feature selection i will go for Lasso and if i want regularization then i will go for Ridge. and if i want both then i would go for Elastic Net.

Time Series forecasting.

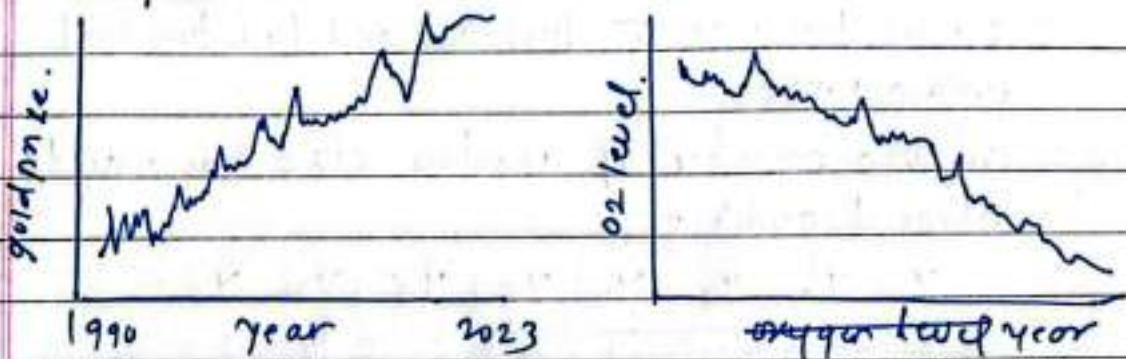
time series :- it is data which is index by time.

- the most important part of time series is sequence.

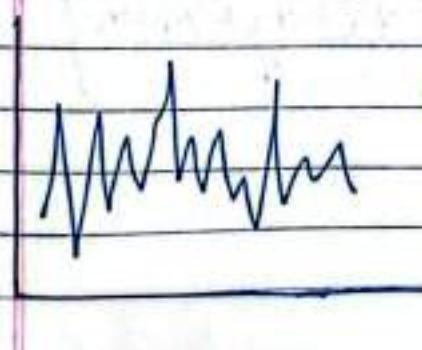
what we do with time series :- we capture the past data based on that we predict for future this entire procedure is called forecasting

- (a) plotting time series
- (b) component of time series
- (c) forecasting in time series
 - (i) Data based prediction
 - (ii) model based prediction

A) line plot:-



- line plot gives holistic view or overall view.
- line plot gives having understanding about long term i.e. for long term what happened to my data.
- after time line plot are difficult to understand.



To overcome this we have technique called Smoothing.

how to do smoothing:-

moving average:- Ex i have rainfall data

Month	M_1	M_2	M_3	M_{120}
data	y_1	y_2	y_3	y_{120}

where M_1 - January

M_2 - february

M_3 - March.

for window size 3

$$y_2 = \frac{y_1 + y_2 + y_3}{3} \quad y_3 = \frac{y_2 + y_3 + y_4}{3} \text{ and so on}$$

- but it would not be same for first and last data point since they will not have one datapoint so their value will remain as it is.
- for window size 5 first two and last two will remain same.
- for even number of window size we wouldn't get balanced number

$$y_1 \ y_2 \ y_3 \ y_4 \ y_5 \ y_6 \ y_7 \ y_8$$

$$y_3 = \frac{y_2 + y_3 + y_4 + y_5}{4} \{ \text{giving high priority to future data} \}$$

$$y_5 = \frac{y_1 + y_2 + y_3 + y_4}{4} \{ \text{giving high priority to past data} \}$$

Solution :-

$$= \frac{1}{2} \left[\frac{y_1 + y_2 + y_3 + y_4}{4} \right] + \frac{1}{2} \left[\frac{y_2 + y_3 + y_4 + y_5}{4} \right]$$

- How to figure out best window size?

1. Way of domain knowledge.

Eg. sales of AC is based on season like in summer it sold out rapidly whereas in winter is not seem that like expected.

so here we can assume or try window size 3/4

2. when we don't have domain knowledge another way
it by calculating size for all and check wherever it
become J mouth.

Components of time series:

it generally has three components:

1) Trend - holistic view - Moving average.

2) Seasonality - periodic change

c) Resid - Noise or Error - less the better.

Decomposition of T.S [Trend, Seasonality, Noise]

- Trend: - find optimum window size

- To find seasonality & Noise subtract trend from original data.

Let's say we have data 2001 to 2010

Month M₁ M₂ M₃ M₄ M₁₂₀

data $y_1 \ y_2 \ y_3 \ y_4 \ \dots \ \dots \ \dots \ y_{120}$

$$\text{Trend} : T_1 = M_1, T_2 = \frac{M_1 + M_2 + M_3}{3}, T_3 = \frac{M_2 + M_3 + M_4}{3}$$

$$\text{Jesomlhy : } S_1 = T_{\text{an}} = \frac{T_1 + T_{13} + T_{25} + \dots + T_{109}}{10}.$$

$$S_2 = f_{cb} = \frac{T_2 + T_{14} + T_{26} + \dots + T_{110}}{10}$$

$$\text{Noise : } N_1 = y_1 - T_1 - S_1, \quad N_2 = y_2 - T_2 - S_2, \quad \dots \quad N_{120} = y_{120} - T_{120} - S_{120}$$

forecasting in time series.

Data based forecasting :-

- ① Simple Exponential :- it is applicable to those time series where there is no trend no seasonality.
- ② Double Exponential :- it is only performed good for those time series where it only have trend but not seasonality.
- ③ Holt winter's model - (triple exponential) it is performed well where there is both trend as well as seasonality (when you don't know whether time series have trend or seasonality then Holt winter is best)

→ Simple Exponential :-

Day	D ₁	D ₂	D ₃	D ₁₉₉	D ₂₀₀	D ₂₀₁
Temp	t ₁	t ₂	t ₃	t ₁₉₉	t ₂₀₀	Today → (predict)

what will be the easiest way to predict for tomorrow's temp when we have data upto today?

① taking Avg :- $\frac{t_1 + t_2 + t_3 + \dots + t_{200}}{200} = x$.

but this not good method since temp varies a lot throughout year summer to winter

ii) another way is whatever is today's temp it will same for tomorrow.

$$t_{\text{day 200}} = 25.7^\circ$$

$$t_{\text{day 201}} = 25.7^\circ$$

Simple Exponential :- instead of taking all 200 data into consideration, we will take recent data for prediction.

- Simple Exponential try to calculate the component in Local Average.

Days	0 ₁	0 ₂	0 ₃	0 ₄	D ₁₉₉	D ₂₀₀
temp	t ₁	t ₂	t ₃	t ₄			t ₁₉₉	t ₂₀₀
	L ₁	L ₂	L ₃	L ₄			L ₁₉₉	L ₂₀₀

- Local Average at Day 4 i.e. L₄ will be combination of two component yesterday's Local Avg., and today temp.

$$L_4 = \frac{1}{2} L_3 + \frac{1}{2} t_4$$

Note: L₁ → t₁ since there is no previous record.

Simple Exponential is applicable, second point onwards.

Mathematical Understanding :-

$$L_{200} = \frac{1}{2} L_{199} + \frac{1}{2} t_{200} \quad \textcircled{1}$$

$$L_{199} = \frac{1}{2} L_{198} + \frac{1}{2} t_{199} \quad \textcircled{2}$$

$$L_{200} = \frac{1}{2} \left[\frac{1}{2} L_{198} + \frac{1}{2} t_{199} \right] + \frac{1}{2} t_{200} \dots \text{2 in } \textcircled{1}$$

$$= \frac{1}{4} L_{198} + \frac{1}{4} t_{199} + \frac{1}{2} t_{200} \quad \textcircled{3}$$

$$L_{198} = \frac{1}{2} L_{197} + \frac{1}{2} t_{198} \quad \textcircled{4}$$

put $\textcircled{4}$ in $\textcircled{3}$

$$\frac{1}{4} \left[\frac{1}{2} L_{197} + \frac{1}{2} t_{198} \right] + \frac{1}{4} t_{199} + \frac{1}{2} t_{200}$$

$$\frac{1}{8} L_{197} + \frac{1}{8} t_{198} + \frac{1}{4} t_{199} + \frac{1}{2} t_{200}$$

If I rearrange this, we will get

$$L_{200} = \frac{1}{2} t_{200} + \frac{1}{4} t_{199} + \frac{1}{8} t_{198} + \frac{1}{16} t_{197} \dots$$

$$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$$

$$(\frac{1}{2})^1 \quad (\frac{1}{2})^2 \quad (\frac{1}{2})^3 \quad (\frac{1}{2})^4$$

from current datapoint weightage of coefficient reduce exponentially whereas in simple average every coefficient were getting equal average.

Let understand with short Example.

Day	1	2	3	4	5	6
Local	1	2	1	2	1	3
	$L_1 \rightarrow L_2 \rightarrow L_3 \rightarrow L_4 \rightarrow L_5 \rightarrow L_6$					

$$L_1 = 1, L_2 = \frac{L_1 + S_2}{2} = \frac{1+2}{2} = 1.5$$

$$L_3 = \frac{L_2 + S_3 - 1.5 + 1}{2} \quad L_4 = \frac{L_3 + S_4 - 1.25 + 2}{2} = \frac{3.25}{2}$$

$$= \frac{2.75}{2} = 1.375$$

$$= 1.625$$

$$L_5 = 1.3125, L_6 = 2.1562$$

All future prediction going to be last available local Average.

Double Exponential :- it contain trend but no seasonality.

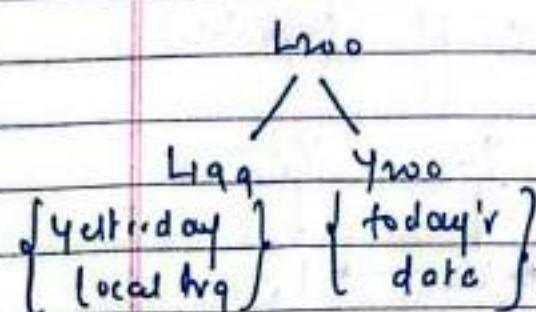
- we can't always rely on local Average because for ex. If chocolate price increase by Rs 1 each day

Day	1	2	3	4	5	6
Price	1	2	3	4	5	6
						L_6

L_6 = around 4.2 which is not true.

Hence we need consider a trend as well, that's why we need both trend as well as Local Average.

for Local Avg



∴ but Local Avg is also expected to change by some amount.

there is little modification in Local Average

$$L_{200} = \frac{1}{2} L_{199} + \frac{1}{2} Y_{200}$$

$$L_{200} = \frac{1}{2} (L_{199} + T_{199}) + \frac{1}{2} Y_{200}$$

↑
trend is added.

$$T_{200} = \frac{1}{2} T_{199} + \underbrace{\frac{1}{2} (L_{200} - L_{199})}_{\text{change in Local Average from yesterday to today}}$$

↑
(today's)
trend

Raw Data $y_1 \quad y_2 \quad y_3 \quad y_4 \quad \dots \quad y_{20}$

Local Avg $L_1 \quad L_2 \quad L_3 \quad L_4 \quad \dots \quad L_{20}$

Trend $T_1 \quad T_2 \quad T_3 \quad T_4 \quad \dots \quad T_{20}$

so from raw data we will calculate Local Avg & trend.

$$L_1 = 0 \quad L_2 = \frac{1}{2} (L_1 + T_1) + \frac{1}{2} (y_2)$$

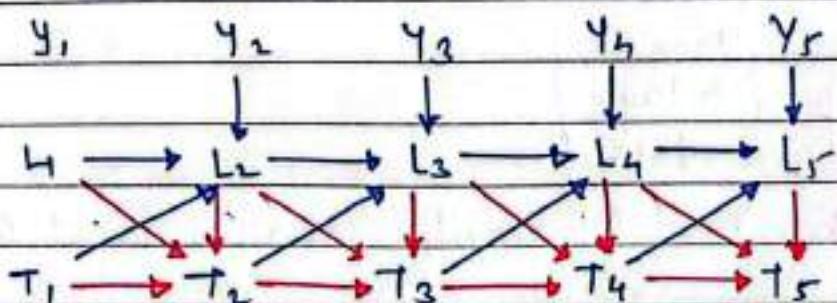
$$T_1 = 0$$

$$T_2 = \frac{1}{2} (T_1) + \frac{1}{2} (L_2 - L_1)$$

Let us understand with Example, we will solve problem, both Simple Exponential & Double Exponential.

Day	1	2	3	4	5	$\frac{L_1 + L_2}{2} = 1.5$
price (y)	1	2	3	4	5	
$L(S)$	1	1.5	2.25	3.125	4.0625	$\frac{1.5 + 3}{2} = \frac{4.5}{2} = 2.25$
$L(D)$	1	1.5	2.375	3.46875	4.6484375	

Now we will calculate with Double Exponential.



$$L_1 = \frac{1}{2}(L_1 + T_1) + \frac{1}{2}(y_2)$$

$$T_2 = \frac{1}{2}(T_1) + \frac{1}{2}(L_2 - L_1)$$

$$L_2 = \frac{1}{2}(L_1 + T_1) + \frac{1}{2}(y_2) = \frac{1}{2} + 1 = 1.5$$

$$T_2 = \frac{1}{2}(0) + \frac{1}{2}(1.5 - 1) = 0 + \frac{0.5}{2} = 0.25$$

$$L_3 = \frac{1}{2}(L_2 + T_2) + \frac{1}{2}(y_3)$$

$$= \frac{1}{2}(1.5 + 0.25) + \frac{1}{2}(3) = \frac{1}{2}(1.75) + 1.5$$

$$= 0.875 + 1.5$$

$$= 2.375$$

$$T_3 = \frac{1}{2}T_2 + \frac{1}{2}(L_3 - L_2)$$

$$= \frac{1}{2}(0.25) + \frac{1}{2}(2.375 - 1.5)$$

$$= 0.125 + 0.4375$$

$$= 0.5625$$

$$\begin{aligned}
 L_4 &= \frac{1}{2}(L_3 + T_3) + \frac{1}{2}(Y_4) \\
 &= \frac{1}{2}(2.375 + 0.5625) + \frac{1}{2}(5) \\
 &= \frac{2.9375}{2} + 2 \\
 &= 1.46875 + 2 = 3.46875
 \end{aligned}$$

$$\begin{aligned}
 T_4 &= \frac{1}{2}(T_3) + \frac{1}{2}(L_4 - L_3) \\
 &= \frac{1}{2}(0.5625) + \frac{1}{2}(3.46875 - 2.375) \\
 &= 0.28125 + 0.5(1.09375) \\
 &= 0.28125 + 0.546875 \\
 &= 0.828125
 \end{aligned}$$

$$\begin{aligned}
 L_5 &= \frac{1}{2}(L_4 + T_4) + \frac{1}{2}(Y_5) \\
 &= \frac{1}{2}(3.46875 + 0.828125) + \frac{1}{2}(5) \\
 &= 2.1484375 + 2.5 \\
 &= 4.6484375 \\
 T_5 &= \frac{1}{2}T_4 + \frac{1}{2}(L_5 - L_3) \\
 &= \frac{1}{2}(0.828125) + \frac{1}{2}(4.6484375 - 2.46875) \\
 &= 0.4140625 + 0.5(1.1797) \\
 &= 0.4140625 + 0.589875 \\
 &= 1.0039375
 \end{aligned}$$

$T_1 \quad T_2 \quad T_3 \quad T_4 \quad T_5$ } Trend is increased.
 0 0.25 0.56 0.82 1.0039

	Actual	Predicted
Simpl. Exponential	4.0625	5
W/ Double Exponential	4.6484375	5

Holt's winter model :- which contain both trend & seasonality.

- Here along with local Average we will find out trend as well as seasonality.

time series data.

y_1	y_2	y_3	y_{200}
A	L_1	L_2	L_{200}
trend	T_1	T_2	T_{200}
seasonality	s_1	s_2	s_{200}

for 200 Months \rightarrow 16 years - 7 months

for seasonality - periodic hm. of 12 months

Let first convert it into double exponential problem.

by subtracting seasonality from time series

$$y_1 - s_1, y_2 - s_2, y_3 - s_3 \dots \dots \dots y_{200} - s_{200}$$

$$L_{200} = \frac{1}{2}(L_{199} + T_{199}) + \frac{1}{2}(y_{200} - s_{200}) \quad (1)$$

{ Note Since $y_{200} = y_{200} - s_{200}$ }

$$T_{200} = \frac{1}{2}(T_{199}) + \frac{1}{2}(L_{200} - L_{199}) \quad (2)$$

$$s_{200} = \frac{1}{2}(s_{200} - 12) + \frac{1}{2}(s_{200})$$

from (1) we can get to know to get

L_{200} we subtract s_{200} from y_{200} then
now s_{200} we can do $y_{200} - L_{200}$.

$$s_{200} = \frac{1}{2}(s_{200} - 12) + \frac{1}{2}(y_{200} - L_{200}) \quad (3)$$

If you are carefully observe then you find there is loop between ③ equations.

for L_{200} we need s_{200} for s_{200} we need L_{200}

Since seasonality is periodic

$$s_{200} = s_{188}$$

$$\text{in eq } ① \quad s_{200} = s_{188}$$

$$L_{200} = \frac{1}{2} (L_{99} + T_{199}) + \frac{1}{2} (Y_{200} - s_{188})$$

$$T_{200} = \frac{1}{2} (T_{199}) + \frac{1}{2} (L_{200} - L_{199})$$

$$s_{200} = \frac{1}{2} (s_{200} - l_2) + \frac{1}{2} (Y_{200} - L_{200})$$

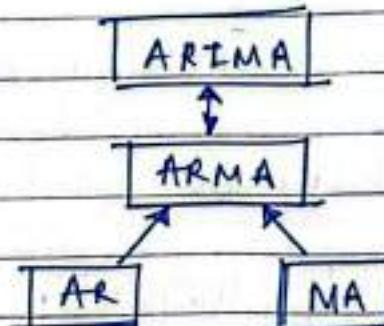
Summary table.

Model	when we can apply	what should be future prediction	hyperparameters
Simple	No trend	(not possible, local Avg = 1)	
Exponential	No seasonality	all future prediction	
Double	only trend	last possible local Avg &	2
Exponential	No seasonality	change of trend	
Triple	both trend and local Average + trend +		3.
Exponential	Seasonality		

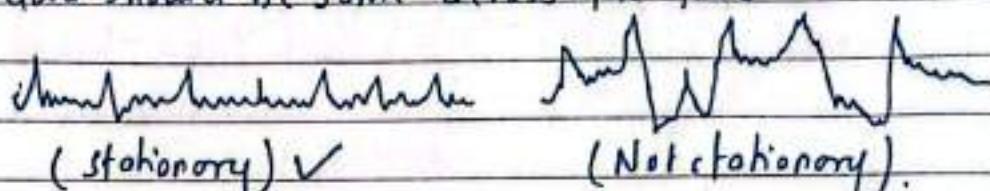
Model Based forecasting :-

Basically there are three models.

- 1) AR
- 2) MA
- 3) ARMA
- 4) ARIMA



Here data must be stationary. It means distribution of data should be same across the time



If we have stationary data then go for model based otherwise databased since there is no any condition

AR (Auto Regressive) : AR is completely depend upon past data and nothing else for future prediction

MA: Moving Average : future data is only depend only external outsiders factors.

ARMA :- future data only going to deal with both past data data and external factors.

AR model :-

$$\text{In regression } y = \alpha x + \beta t + \epsilon$$

$$\text{future } y_t = \alpha y_{t-1} + \beta t + \epsilon$$

↑ ↑ ↑ ↑
past note time Error

Here future data is only depend on past data
if it would depend on past two data

$$Y_t = \alpha_1 Y_{t-1} + \alpha_2 Y_{t-2} + \beta + \epsilon \quad AR(2)$$

$$Y_t = \alpha_1 Y_{t-1} + \alpha_2 Y_{t-2} + \alpha_3 Y_{t-3} + \beta + \epsilon \dots AR(3)$$

How to select Best possible model. AR1, AR2, AR3
ARt ~~will~~ based on test data accuracy whichever will
perform good will finalize that.

MA : Moving Average Here my future data is
only depend upon external factors.

Ex. rate of some production which depends upon
seasons (grain or fruit) due to natural
calamities.

data	y_1	y_2	y_3	.	.	y_t
External factor	e_1	e_2	e_3	.	.	e_t

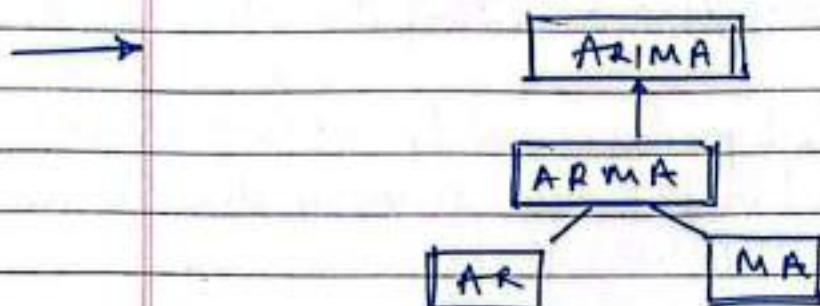
$$Y_t = \alpha e_{t-1} + \beta + \epsilon \dots MA(1)$$

this year
data | last year
data \ constant

$$Y_t = \alpha_1 e_{t-1} + \alpha_2 e_{t-2} + \beta + \epsilon \dots MA(2)$$

Model based forecasting.

- ① AR
- ② MA
- ③ ARMA (AR+MA)
- ④ ARIMA (which is combination of ARMA itself)



AR → future prediction are completely depend on past data, it is just copy of regression model.

$$\text{target } \gamma = \underbrace{\alpha x + \beta}_{\text{prediction}} + \varepsilon \quad \text{error.}$$

variable.

$$y_t = \underbrace{\alpha y_{t-1} + \beta}_{\text{past data}} + \varepsilon \quad \text{AR(1)}$$

it is depend on one past data

$$y_t = \alpha_1 y_{t-1} + \alpha_2 y_{t-2} + \alpha_3 y_{t-3} + \beta + \varepsilon$$

tomorrow ↑ today ↑ yesterday ↑ day before yesterday

it is similar to SLR

$$\text{SLR: } \gamma = \alpha x + \beta + \varepsilon$$

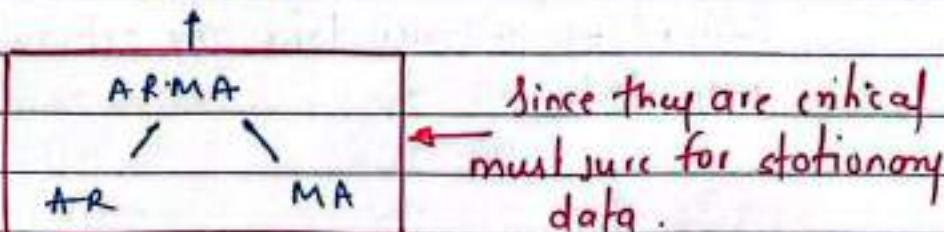
$$\text{AR(1)}: \gamma_t = \alpha y_{t-1} + \beta + \varepsilon_t$$

$$\text{n features } y = \alpha_1 x_1 + \alpha_2 x_2 + \beta + \varepsilon$$

$$\text{AR(2)}: \gamma_t = \alpha_1 y_{t-1} + \alpha_2 y_{t-2} + \beta + \varepsilon_t$$

- one of the necessary condition for apply these data is data should be stationary i.e. it means distribution of data should be same across all the data point.
- Since among all four model ARMA, AR and MA are very critical for that we must sure that data should be stationary.

ARIMA



MA → Moving Average :- here future prediction are completely depend upon external factors.

production of rice	y_1	y_2	y_3	y_4	y_5	y_6	y_7
Error due to external factor	ϵ_1	ϵ_2	ϵ_3	ϵ_4	ϵ_5	ϵ_6	ϵ_7

- when we say our data is only depending on last time point for example for 8th year

$$y_8 = \alpha \epsilon_7 + \beta + \epsilon_8$$

so general formula is

$$y_t = \alpha y_{t-1} + \beta + \epsilon_t \quad \text{Moving Avg (MA)} \\ \text{depend on External Error factors.}$$

$$y_t = \alpha y_{t-1} + \beta + \epsilon_t \quad \text{Auto Regressio (AR)} \\ \text{depend on past data.}$$

for MA(2) and MA(3) equation would be .

$$MA(2) : Y_t = \alpha_1 \epsilon_{t-1} + \alpha_2 \epsilon_{t-2} + \beta + \epsilon$$

$$MA(3) : Y_t = \alpha_1 \epsilon_{t-1} + \alpha_2 \epsilon_{t-2} + \alpha_3 \epsilon_{t-3} + \beta + \epsilon.$$

- it is exactly similar to AR model only difference is instead of previous data here external factors are introduced.

ARMA : it is combination of AR and MA model.

$$\boxed{ARMA(1,1)} \\ \boxed{AR(1) + MA(1)}$$

$$ARMA(p,q) = AR(p) + MA(q)$$

$$Y_t = \alpha_1 Y_{t-1} + \beta_1 \epsilon_{t-1} + \beta_0 + \epsilon_t .$$

if $\beta_1 \epsilon_{t-1} = 0$ then it will be AR(1)

$$Y_t = \underline{\alpha_1 t_{t-1}} + \beta_1 \epsilon_{t-1} + \beta_0 + \epsilon_t .$$

i. $\alpha_1 t_{t-1} = 0$ then it will be MA(1)

- that's why you don't need to put MA and AR model separately .
- if you think your model mostly inclined to MA then we can put AR component zero if model inclined to AR then we can put MA component zero.

- thing we need to take care for these three models.
- i) stationary :: if it is stationary we gonna directly used ARMA instead of using AR and MA Model separately.
- y. ARMA model has two hyperparameter (P, q)
 P corresponds to AR and q corresponds to MA

how to choose best value of p and q .

by creating multiple model in range of (0 to 5)
 for each p and q .

$$p : \{0, 1, 2, 3, 4, 5\} = 6 \quad 6 \times 6 = 36$$

$$q : \{0, 1, 2, 3, 4, 5\} = 6$$

since 0,0 is not possible, pair we have $36 - 1 = 35$

ARIMA : ARIMA doesn't care about data is stationary or Not.

this not a new model but extension of ARMA

ARIMA - Auto Regressive integrated moving average.

- if transformed data from original data follows ARMA model then original data will must follow ARIMA.

$$\begin{array}{ccccccc} y_1 & y_2 & y_3 & y_4 & \dots & y_n \\ \downarrow & \downarrow & & & \downarrow & \\ z_1 & z_2 & z_3 & z_4 & \dots & z_n \end{array}$$

transformed data

if this $\{z_1, z_2, \dots, z_n\}$ follows ARMA then
 (y_1, y_2, \dots, y_n) will must follow ARIMA.

for ex.

Bank balance of Every month.

$$y_1 \ y_2 \ y_3 \ y_4 \ \dots \ y_t \leftarrow (\text{today})$$

If less money is added every month not due to new
added but due to interest

what will interest added each month

$$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\ z_1 \ z_2 \ z_3 \ z_4 \ z_t$$

where $z_1 = y_2 - y_1$ } this is newly constructed data
 $z_2 = y_3 - y_2$ } if this follow ARMA then parat
 $z_3 = y_4 - y_3$ } data will follow ARIMA

ARMA has two component (P, q) whereas

ARIMA has three component (P, q, d)

^{degree of subtraction}

d - for subtraction of consecutive numbers.

$$d = 1$$

$$d = 2$$

$$d = 3$$

$$z_1 = y_2 - y_1$$

$$z_1 = y_3 - y_1$$

$$z_1 = y_4 - y_1$$

$$z_1 = y_3 - y_2$$

$$z_2 = y_4 - y_2$$

$$z_3 = y_5 - y_2$$

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

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

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

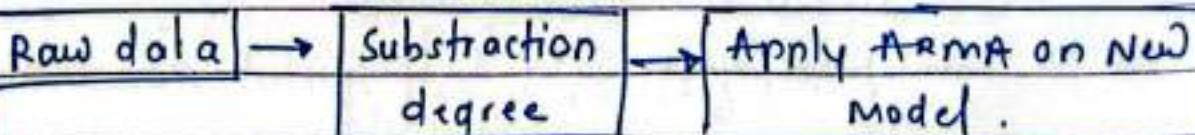
$$z_{100} \ y_{101} \ y_{100}$$

$$z_{100} = y_{102} - y_{100}$$

$$z_{100} = y_{103} - y_{100}$$

what is lifecycle of ARIMA

Transformation



Summary table.

Model	No of parameter	Conditions	depends on
AR	1	stationary	past
MA	1	stationary	External factor
ARMA	2	stationary	past + External factor
ARIMA	3	Noneed - since I'm transforming the data.	

forecasting

Data Based

there is no previous condition

Simple Exponential Double Exponential Holt winters

No trend No seasonality. trend ✓ but No seasonality. both trend & seasonality. { when we are not sure }

AR MA ARMA ARIMA

{ stationary data }

Model Based

{ whenever we don't know which to use we can use directly Holt winter }

Evaluation of time series Model.

there are two matrix we gonna used.

- ① Mean Square Error
- ② MAPE score.

Mean Square Error.

Let say i have so many future data.

Actual	Prediction
y_t	\hat{y}_t
y_{t+1}	\hat{y}_{t+1}
y_{t+2}	\hat{y}_{t+2}

$$MSE = \frac{(y_t - \hat{y}_t)^2 + (y_{t+1} - \hat{y}_{t+1})^2 + (y_{t+2} - \hat{y}_{t+2})^2}{n}$$

- ② MAPE Score. (Mean Absolute Percentage Error)

Let compare 2 scenarios.

	s1	s2
Actual	3	1000
Predicted	2	999

which is more accurate?

{ if you try MSE for both it would be same }
 it means, sometimes we are not interested in amount of error but percentage of error

formula for MAPE scores is

$$\frac{|y_A - y_P|}{y_A} \quad \text{this is value of percentage}$$

$$\text{first case: } \frac{|3-2|}{3} = \frac{1}{3} = 33.3\%$$

$$\text{Second case: } \left| \frac{1000 - 999}{100} \right| = 0.001 \\ 0.01\%$$

Since we are not only concerned about Absolute error always but absolute percentage error also.

Exercises :- calculate the MSE and MAPE score

Actual	Predicted
5	10
20	15
30	25

$$\text{MSE} = \frac{(5-10)^2 + (20-15)^2 + (30-25)^2}{3} \\ = \frac{(-5)^2 + (5)^2 + (5)^2}{3} = 25$$

MAPE score :

$$\left| \frac{y_A - y_P}{y_A} \right| = \left| \frac{5 - 10}{5} \right| = -5/5 = -1 = 100\%$$

$$\therefore \quad = \left| \frac{20 - 15}{20} \right| = 25\%$$

$$\therefore \quad = \left| \frac{30 - 25}{25} \right| = 16.6\%$$

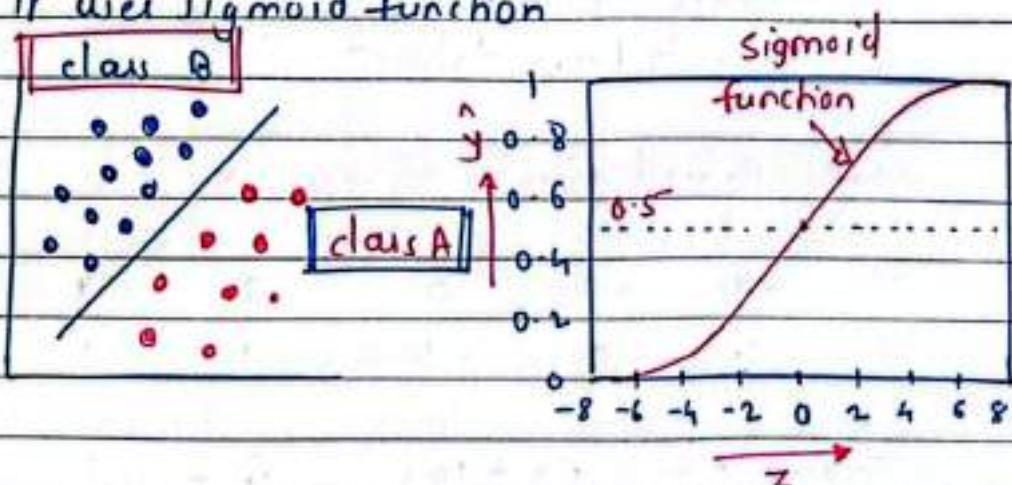
$$\text{Mean abs error} = \frac{100 + 25 + 16.6}{3} = \underline{\underline{47.2}}$$

Classification Algorithm

- 1) Logistic Regression 2) SVM
- 3) KNN
- 4) Decision Tree 5) Random forest
- .

1) Logistic Regression:-

1. it is supervised learning model.
2. it is classification model and best for binary classification.
3. it uses sigmoid function



$$\text{Sigmoid function} = \hat{y} = \frac{1}{1 + e^{-z}}$$

where $z = w_1x_1 + w_2x_2 + b$.

$(w_1x_1 + w_2x_2 + b)$ - eqn of line

slope $m \rightarrow$ weight w

intercept $c \rightarrow$ bias b

• \hat{y} = probability that ($y=1$)

$\hat{y} = p(y=1|x) \dots \{ \text{probability of } y \text{ being 1 for given value of } x \}$

x = input feature

w = weights (it will be in the form of)

{ number of weight equal to number of feature in the dataset }

b = bias

$\hat{y} = \sigma(z)$

- Advantages :- 1) Easy to implement
 2) perform well on data with linear relationship
 3) less prone to overfitting for low dimensional dataset.

Disadvantages :- 1) High dimensional dataset causes overfitting.

- 2) difficult to capture complex relationship in dataset
- 3) sensitive to outliers
- 4) Need large dataset.

Math Behind Logistic Regression.

X	-9	-8	0	8	9
Y	0	0	1	1	1

$$\text{Assume } z = 5x + 10 \quad \hat{y} = \frac{1}{1 + e^{-z}}$$

$x = -9$	$x = -8$	$x = 0$	$x = 8$	$x = 9$
----------	----------	---------	---------	---------

$$\begin{aligned} z &= 5(-9) + 10 & z &= 5(-8) + 10 & z &= 5(0) + 10 & z &= 5(8) + 10 & z &= 5(9) + 10 \\ &= -35 & & = -30 & & = 10 & & = 50 & & = 55 \end{aligned}$$

$$\begin{aligned} \hat{y} &= \frac{1}{1 + e^{-35}} & \hat{y} &= \frac{1}{1 + e^{-30}} & \hat{y} &= \frac{1}{1 + e^{10}} & \hat{y} &= \frac{1}{1 + e^{50}} & \hat{y} &= \frac{1}{1 + e^{55}} \end{aligned}$$

$$\begin{aligned} \hat{y} &= 0 & \hat{y} &= 0 & \hat{y} &= 1 & \hat{y} &= 1 & \hat{y} &= 1 \end{aligned}$$

- Advantages :- 1) Easy to implement
 2) perform well on data with linear relationship
 3) less prone to overfitting for low dimensional dataset.

Disadvantages :- 1) High dimensional dataset causes overfitting.

- 2) difficult to capture complex relationship in dataset
- 3) sensitive to outliers
- 4) Need large dataset.

Math Behind Logistic Regression.

X	-9	-8	0	8	9
Y	0	0	1	1	1

Assume $Z = 5x + 10$ $\hat{y} = \frac{1}{1 + e^{-z}}$

$x = -9$	$x = -8$	$x = 0$	$x = 8$	$x = 9$
----------	----------	---------	---------	---------

$Z = 5(-9) + 10$ = -35	$Z = 5(-8) + 10$ = -30	$Z = 5(0) + 10$ = 10	$Z = 5(8) + 10$ = 50	$Z = 5(9) + 10$ = 55
---------------------------	---------------------------	-------------------------	-------------------------	-------------------------

$\hat{y} = \frac{1}{1 + e^{-35}}$	$\hat{y} = \frac{1}{1 + e^{-30}}$	$\hat{y} = \frac{1}{1 + e^{10}}$	$\hat{y} = \frac{1}{1 + e^{50}}$	$\hat{y} = \frac{1}{1 + e^{55}}$
-----------------------------------	-----------------------------------	----------------------------------	----------------------------------	----------------------------------

$\hat{y} = 0$	$\hat{y} = 0$	$\hat{y} = 1$	$\hat{y} = 1$	$\hat{y} = 1$
---------------	---------------	---------------	---------------	---------------

Inference : if z value is large positive number,

$$\hat{y} = \frac{1}{1+e^{-z}} \approx \hat{y} = 1.$$

if z is large negative number,

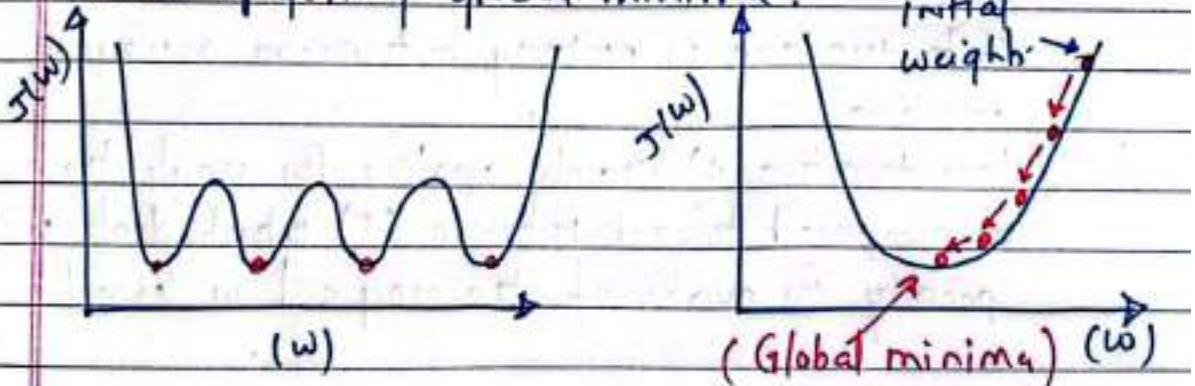
$$\hat{y} = \frac{1}{1+e^{z}} \approx \hat{y} = 0.$$

Loss function & cost function for Logistic Regression.

- loss function measures how far an estimated value is from true value.

$$\text{Loss function for linear regression} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_{pred})^2$$

if we use this function we will get many local minima instead of getting global minima.



- Binary cross entropy loss function (or) log loss.

$$L(y, \hat{y}) = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$$

Here,

$$y \rightarrow 0 \text{ or } 1$$

But $\hat{y} \rightarrow 0 \text{ to } 1$ (probability could be continuous)

when $y=1$

$$L(1, \hat{y}) = -(1 \log \hat{y} + (1-1) \log(1-\hat{y})) \\ = -\log \hat{y}$$

- Since we always want smaller loss function value hence \hat{y} should be very large (from 0 to 1) if it is the $-\log \hat{y}$ will be very large negative number or very small number.

when $y=0$

$$L(1, \hat{y}) = -(0 \log \hat{y} + (1-0) \log(1-\hat{y})) \\ = -\log(1-\hat{y}).$$

- Since we want smaller loss function value, hence \hat{y} should be very small the automatically $(1-\hat{y})$ will be very large thus $-\log(1-\hat{y})$ will be large negative number or very small number.

- Cost function is nothing but mean average of loss function

- Loss function (L) mainly applies for single training set as compared to cost function (J) which deals with a penalty for number of training set or complete batch.

Loss function:

$$L(y, \hat{y}) = -(y \log \hat{y} + (1-y) \log(1-\hat{y})) \quad \text{for single}$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m (L(y^{(i)}, \hat{y}^{(i)})) =$$

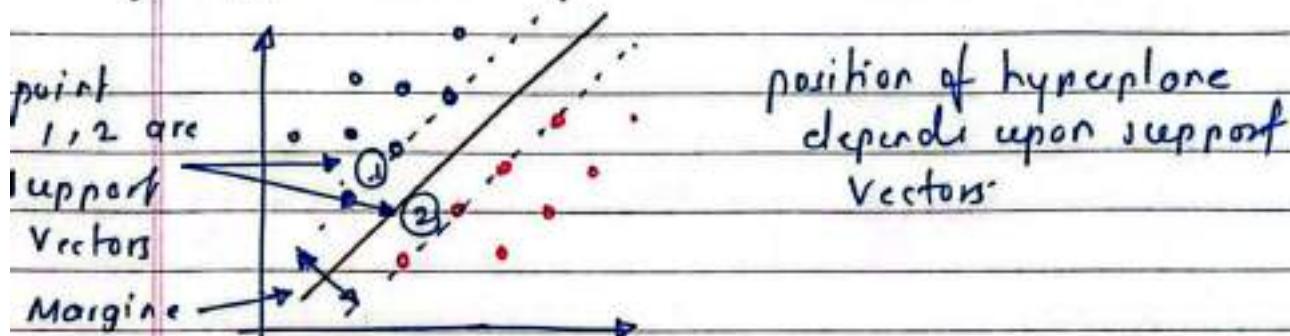
$$-\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)}))$$

{'m' denotes number of data-points in the }
training set

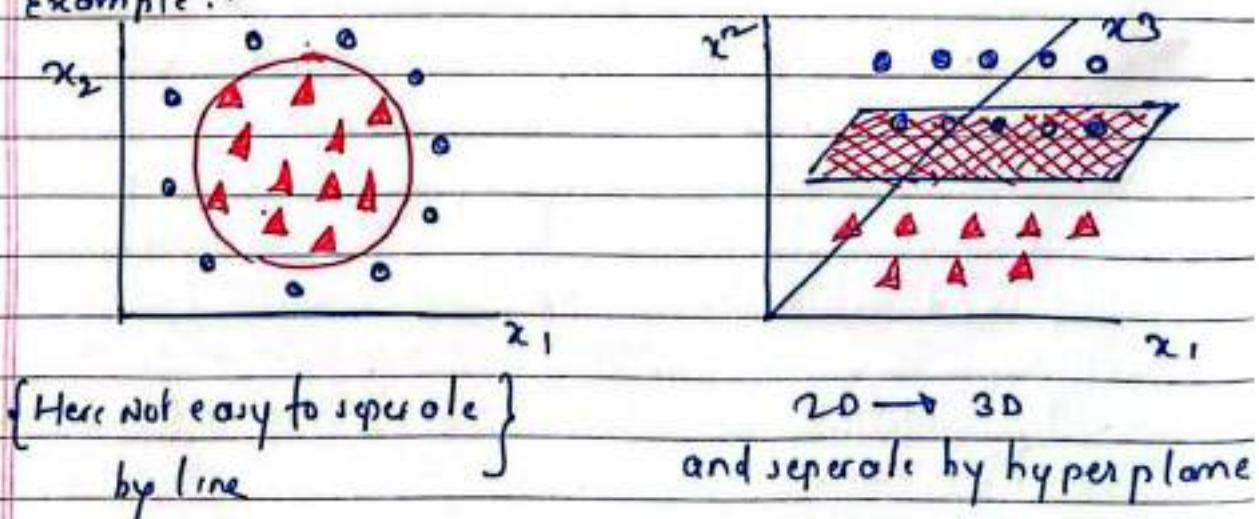
Support Vector Machine (SVM)

Basic about SVM.

- 1) it is supervised ML model.
- 2) it can be used for both classification as well as regression but it is predominantly used for binary classification.
- 3) Hyperplane.
- 4) Support vectors



- for 2D data it is easy to draw hyperplane but where data point can't be separated by line need to convert into 3D where we can separate the datapoint by hyperplane
- Example :-



Hyperplane:- Hyperplane is line (in 2D) or plane that separates the data point into two classes

Support Vectors :- these are the datapoints which are nearest to hyperplane if these datapoints changes position of hyperplane changes.

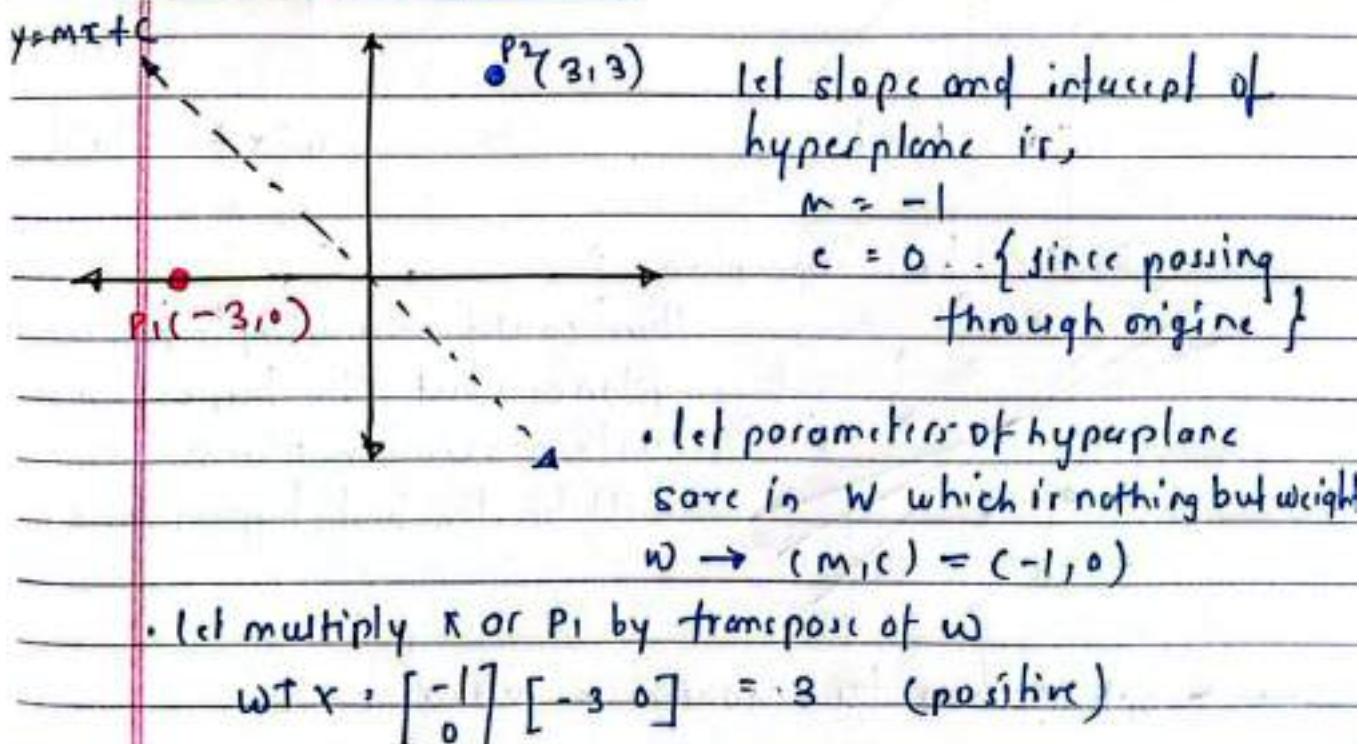
Advantages of SVM

- 1) works fine with smaller dataset
- 2) works fine or efficiently where there is clear margin of separation
- 3) works well with high dimensional data

Disadvantages

- 1) Not suitable for large dataset as training time would take very large.
- 2) Not suitable for noisy (outlier) dataset with overlapping classes.

Math Behind SVM



[Note : why transpose ? \rightarrow for matrix multiplication no. of columns of 1st Matrix must be equal to no. of rows of second matrix]

- positive value indicates all points of hyperplane will be positive class.

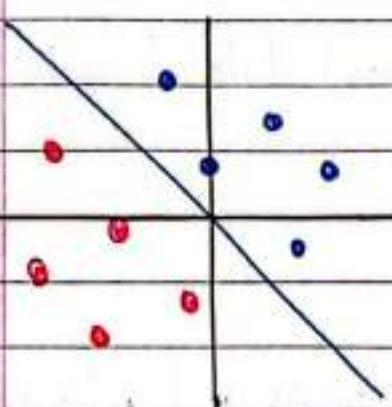
for $P_2(3, 3)$

$$w^T x = \begin{bmatrix} -1 \\ 0 \end{bmatrix} \begin{bmatrix} 3 \\ 3 \end{bmatrix}$$

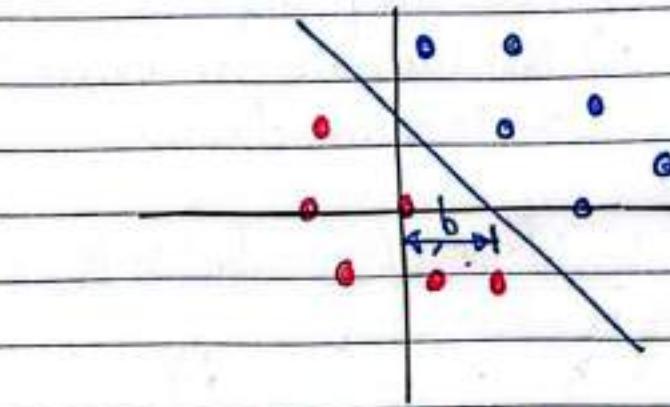
$$= -3 \text{ (Negative)}$$

Here for all the points which lie on the right side of hyperplane will be belong to negative class.

But Not all the time hyperplane will pass through origin.

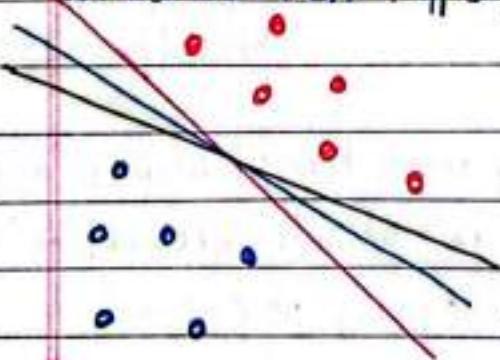


$$w^T x = \text{label}$$



$$w^T x + b = \text{label}$$

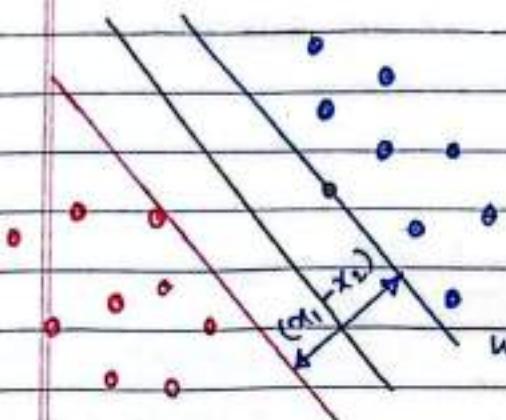
which is best hyperplane?



there could be multiple hyperplane, but the hyperplane with maximum margin size will be the best hyperplane.

→ optimization for maximum margin

$$w^T x + b = 1 \text{ or } -1$$



Equation of point or
blue support vector &
its output value any
negative value.

$w^T x + b = 1 \Rightarrow$ this is equation of point or red
support vector and its output value could be any
positive value

to get margin let subtract one from another.

$$w^T x_1 + b = 1$$

$$(-) w^T x_2 + b = -1$$

$$w^T (x_1 - x_2) = 2$$

$$w^T (x_1 - x_2) = 2$$

divide both sides by $\|w\|$

$$\frac{w^T (x_1 - x_2)}{\|w\|} = \frac{2}{\|w\|}$$

$$(x_1 - x_2) \cdot \frac{2}{\|w\|} \leftarrow \text{this is nothing but magnitude of vector.}$$

and

$$y_i = \begin{cases} -1 & w^T x_i + b \leq -1 \\ 1 & w^T x_i + b \geq 1 \end{cases} \quad (\text{label})$$

To max $\left(\frac{2}{\|w\|} \right)$ such that.

$$y_i = \begin{cases} -1 & w^T x_i + b \leq -1 \\ 1 & w^T x_i + b \geq 1 \end{cases}$$

Max instead of using $\frac{1}{2} \left(\frac{w^T w}{||w||^2} \right)$ we can also try Min which make better sense

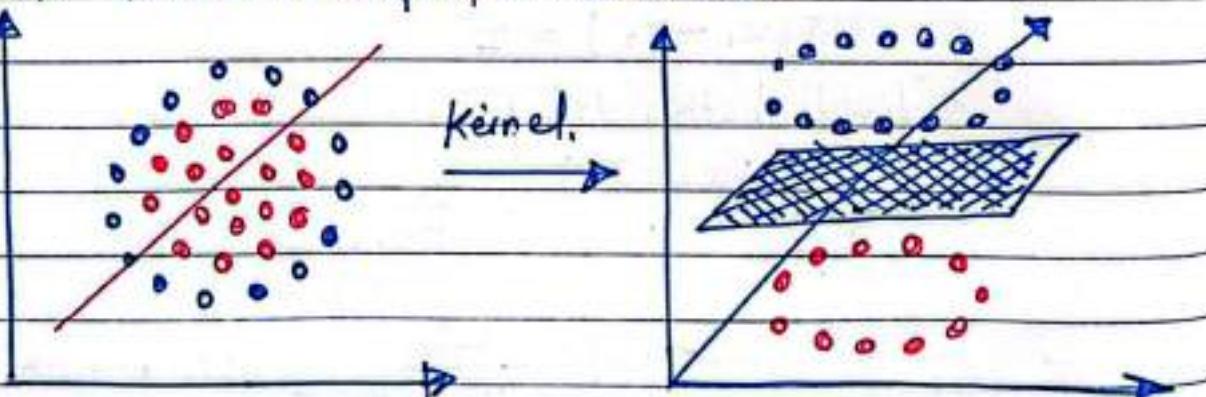
$$\min \left(\frac{||w||^2}{2} \right) + c \times \sum \epsilon_i$$

c : Number of error

ϵ_i : Error magnitude

(we all model to train with some error to avoid overfitting (i.e. it will be good and train and will bad for test data))

Kernels' in SVM : Generally function of the kernel is to transform the training set of data so that non-linear decision surface can be transformed to a linear equation in higher number of dimension space. It return the inner product between two points in standard feature dimension.



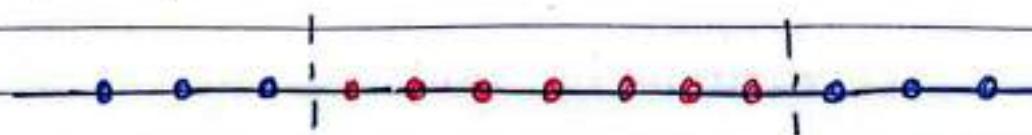
Type of SVM Kernel.

- 1) Linear
- 2) polynomial
- 3) Radial Basis function. (rbf)
- 4) sigmoid.

Suppose feature (x)

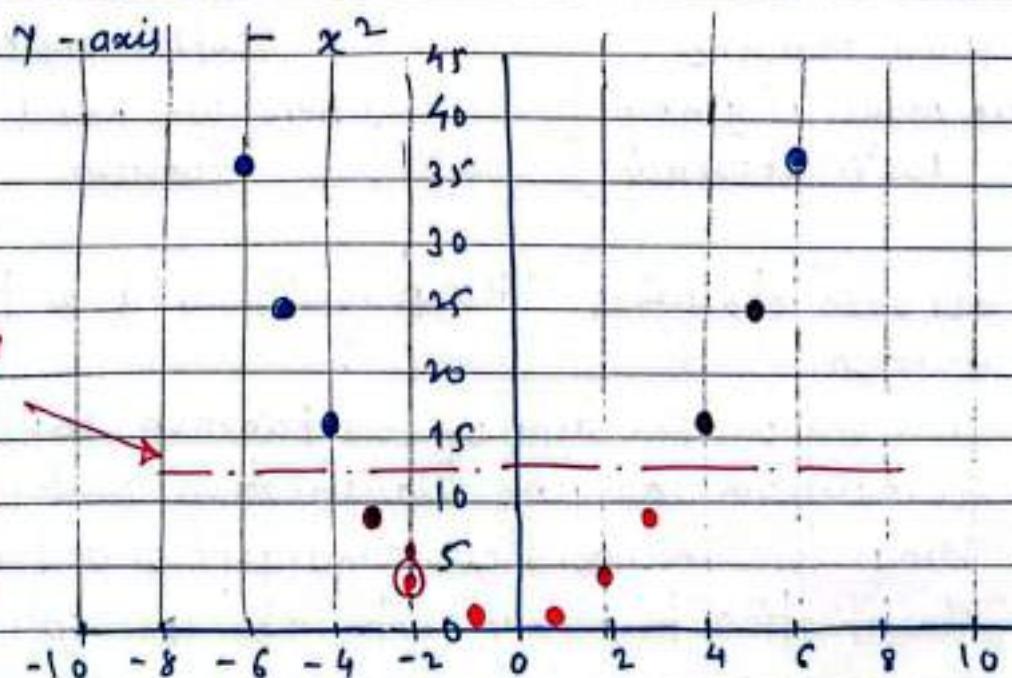
x	-6	-5	-4	-3	-2	-1	0	1	2	3	4	5	6
x^2	36	25	16	9	4	1	0	1	4	9	16	25	36

If you try to plot x on this 1D line.



- we can see none of the line could separate the two class perfectly.
- that's why we add another feature which is function of x i.e x^2

x -axis = x



1) Linear kernel :- $K(x_1, x_2) = x_1^T x_2$

{ best suitable for having too many features }

2) polynomial kernel.

$$K(x_1, x_2) = (x_1^T x_2 + r)^d$$

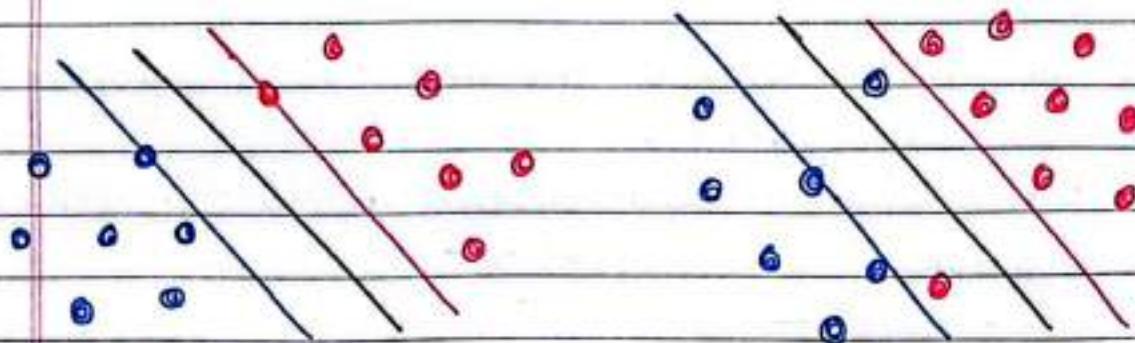
3.) radial basis function. (rbf kernel)

$$k(x_1, x_2) = \exp(-\gamma \cdot \|x_1 - x_2\|^2)$$

4. Sigmoid function

$$k(x_1, x_2) = \tanh(\gamma \cdot x_1 \cdot x_2 + \gamma).$$

Loss function for SVM classifier.



Hard Margin

{Here model is giving
100% accuracy}

soft margin

{Here we need loss
function}

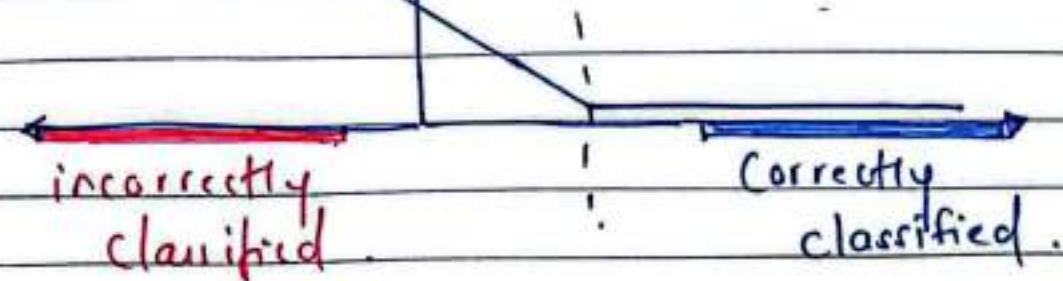
- for SVM classifier "Hinge loss" is used as loss function.
- it is one of the type of loss function mainly used for maximum margin classification model.
- Hinge loss incorporates a margin or distance from classification boundary into the loss calculation. Even if new observation classified correctly they can incur penalty if the margin from decision boundary is not large enough.

$$L = \max(0, 1 - y_i(\omega^T x + b))$$

0 - for correct classification

1 - for wrong classification.

~~Hinge loss~~



Let's calculate for misclassification.

$$y_i = 1, \quad \hat{y}_i = -1$$

$$\begin{aligned} L &:= (1 - 1)(-1) \\ &= 1 + 1 \\ &= 2 \end{aligned}$$

$$y_i = -1 \quad \hat{y} = 1$$

$$\begin{aligned} L &:= (1 - (-1))(1) \\ &= 1 + 1 \\ &= 2 \end{aligned}$$

{ both are high loss value }

Now let for correct classification.

$$y_i = 1 \quad \hat{y}_i = 1$$

$$(0 - (1))(1)$$

$$0 - 1$$

$$-1$$

$$y_i = -1 \quad \hat{y}_i = -1$$

$$(0 - (-1))(-1)$$

$$0 - 1$$

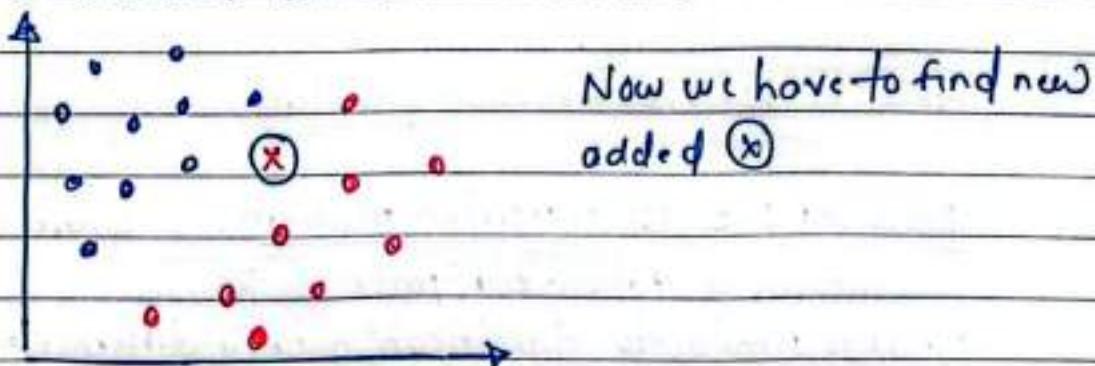
$$-1$$

{ both are low loss value }

KNN - K-Nearest Neighbour.

- The abbreviation KNN stands for "K-Nearest Neighbor". It is supervised machine learning algorithm. This algorithm can be used to solve both classification as well as regression problem.
- The number of nearest neighbor to a new unknown variable that has to be predicted or classified denoted by the symbol ' k '.
- Whenever new data will come if new data is close to 1 then prediction will be class 1 otherwise 0 for binary classification. Some principle for multi-classification problem.
- In general k is odd number.

Let understand with some example



Working of KNN Algorithm.

Step 1: Loading training as well as test data.

Step 2: Next we choose value of k (hyperparameter). The nearest datapoint i.e. k can be any integer.

Step 3: For each data point (new) test data do the following.

3.1) Calculate distance between test data and each row of training data with the help of distance calculating method like Euclidean, Manhattan or Hamming distance. The most commonly method to calculate is Euclidean.

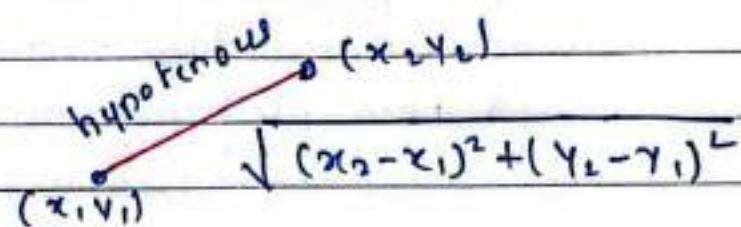
3.2) Now based on its distance value, sort them in ascending order.

3.3) Next it will choose the top k rows from sorted array.

3.4) Now it will assign a class to a test point based on most frequent class of those rows.

Step 4 : End.

a) Euclidean distance :-



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

b) Manhattan distance



$$M.\text{distance} = x + y$$

this is for classification problem.

Now let see for regression problem. : where output is continuous data. e.g. price of house.

It is similar to classification only difference is last step where we were considering most frequent class here from k number of output (nearest 'k') we calculate their mean, that's it.

Limitation : Not applicable to huge dataset since calculating distance would consume lot of time.

- Sensitive to outliers
- Sensitive to missing values

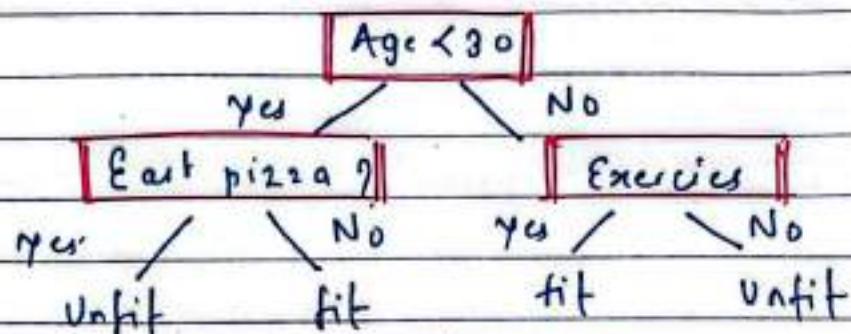
Design Tree.

Show decision Tree.

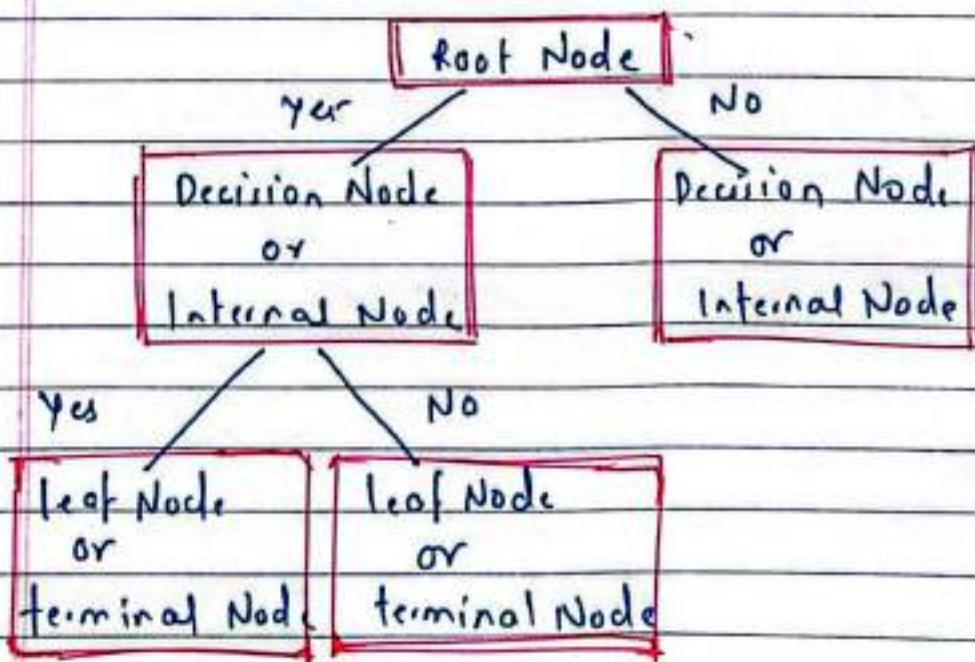
- 1) It is supervised ML Model
- 2) Used both classification & Regression.
- 3) Build Decision Nodes at each step.
- 4) Basis of Tree based model.

Let understand with Example.

Is person fit or Not?



Structure & terminology of DT :-



Advantages :-

- 1) Can be used for both classification & Regression
- 2) Easy to interpret
- 3) No need for Normalization or scaling
- 4) Not sensitive to outliers.

Disadvantages:

- 1) Overfitting issue
- 2) small changes in the data alter the tree structure causing instability.
- 3) training time is relatively high.

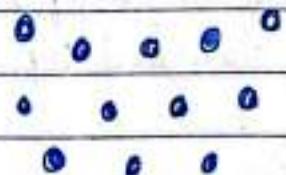
Some concepts in DT :-

- a) Entropy
- b) Information Gain
- c) Gini Impurity

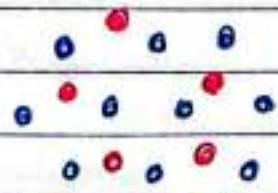
Entropy : High Information gain : low Gini Impurity : High	Entropy : low Information gain : High Gini Impurity : low
--	---

{ Entropy and gini impurity are inversely proportional to each other }

Entropy :- In ML Entropy is the quantitative measure of randomness of information being processed.



Low entropy



High entropy

- A high value of entropy means that randomness in the system is high and thus making accurate prediction is tough.
- A low value of Entropy means that randomness in the system is low and thus making accurate prediction is easier.

$$\text{Entropy} = \sum_{i=1}^c -p_i \log_2 p_i \quad c = \text{number of classes}$$

$p_i = \text{probability of } i\text{th class}$

- Information Gain :- once we find entropy to find which feature to be selected as root Node or internal Node we use information gain:
- it is measure of how much information a feature provides about class low entropy leads to increased Information Gain and high entropy lead low information gain
- information gain computes th. difference between entropy before split and average entropy after split of the dataset based on given value

$$\text{Info gain}(T, f) = \text{Entropy}(T) - \sum_{\text{feature } f} \frac{|T_f|}{|T|} \cdot \text{Entropy}(f)$$

Target ↑
 feature

- Gini Impurity :- it is measure of impurity at Node
- The split made in decision tree is said to be pure if all the data point are accurately separated into different class.
- it measures the likelihood that randomly selected data point would be incorrectly classified by specific Node.

$$\text{formula} = 1 - (P_Y^2 + P_N^2)$$

P_Y = probability of class Y

P_N = probability of class N

Decision Tree for regression.

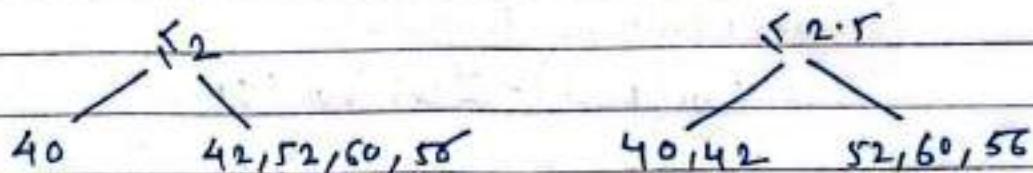
Let understand with example

freq	Gap	Salary (k)
2	yes	40
2.5	yes	42
3	No	52
4	No	60
4.5	yes	<u>56</u>
		$\bar{y} = 50 \leftarrow \text{Average.}$

(+) take experience at root node

{Note:- Since exp is continuous data OT arranging it in ascending order}

- Now for comparison we will take two node example



- Now to decide which split is suitable we used one concept called "Variance reduction."

$$\text{Variance} = \frac{1}{n} \sum_{i=1}^n (y - \bar{y})^2 \leftarrow (\text{MSE formula})$$

where \bar{y} = Average output.

- Now we have to calculate Variance at each node.
- 1st we will calculate variance at root.

$$\begin{aligned}
 \text{Variance (Root)} &= \frac{1}{5} \left[(40-50)^2 + (42-50)^2 + (52-50)^2 \right. \\
 &\quad \left. + (60-50)^2 + (56-50)^2 \right] \\
 &= \frac{1}{5} [100 + 64 + 4 + 100 + 36] \\
 &= 60.8
 \end{aligned}$$

Now we will calculate variance of each internal node or decision node

$$\text{Variance (IN1)} = \frac{1}{2} [(40 - 50)^2]$$

$$= 100$$

$$\text{Variance (IN2)} = \frac{1}{4} [(42 - 50)^2 + (52 - 50)^2 + (60 - 50)^2 + (50 - 50)^2]$$

$$= \frac{1}{4} [(-8)^2 + (2)^2 + (10)^2 + (0)^2]$$

$$= \frac{1}{4} [64 + 4 + 100 + 36]$$

$$= 51$$

Variance reduction formula:-

$$= \text{Var(Root)} - \sum w_i \text{Var(IN)}$$

$$= 60 \cdot 8 - \left[\frac{1}{5}(100) + \frac{4}{5}(51) \right]$$

$$= 60 \cdot 8 - 20 - 40 \cdot 8$$

$$= 0$$

Some we will calculate for second condition ≤ 2.5
whoever have large variance reduction we will
finalize it for splitting.

$$\text{Var(IN1)} = \frac{1}{2} [(40 - 50)^2 + (42 - 50)^2]$$

$$= \frac{1}{2} [100 + 64]$$

$$= \frac{\sqrt{168}}{2} = 82$$

$$\begin{aligned}\text{Var}_{IN2} &= \frac{1}{2} [(52-50)^2 + (60-50)^2 + (56-50)^2] \\ &= \frac{1}{3} [14 + 100 + 36] \\ &= \frac{140}{3} = 46.66\end{aligned}$$

Variance reduction for next split i.e. S2-S

$$\therefore \text{Var}(\text{Root}) - \sum w_i \text{Var}(IN)$$

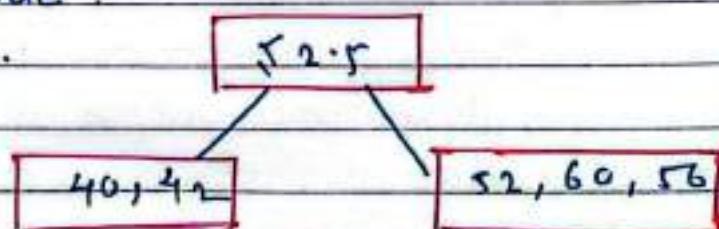
$$\begin{aligned}&= 60 \cdot 8 - \left[\frac{2}{5} (82) + \frac{3}{5} (46.66) \right] \\ &= 0.304\end{aligned}$$

$\text{Var}(\text{split 2}) > \text{Var}(\text{split 1})$ that's why we will select second split.

How to calculate o/p for the test data

- whichever leaf node your test data reached to take avg of all the numbers present in the same leaf Node.

for Ex.



if test data reached to 1st leaf then

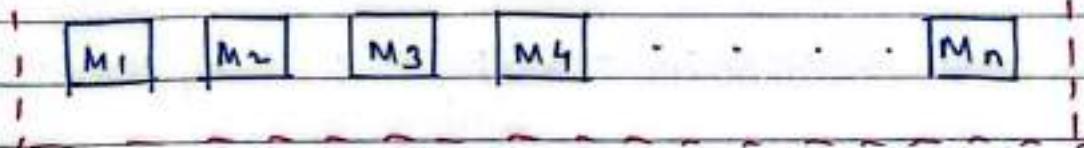
$$\frac{40+41}{2} = 41$$

and if it reached to second leaf then

$$\frac{52+60+56}{3} = 56$$

Introduction to Ensemble Learning or Techniques.

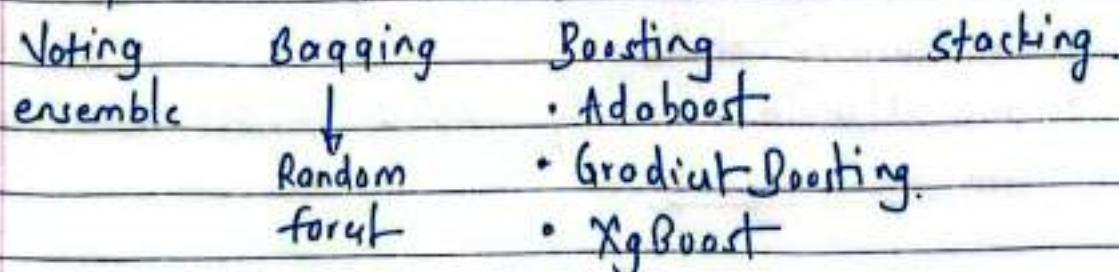
- Ensemble technique is a machine learning technique that combines several basic models in order to produce one optimum predictive model.



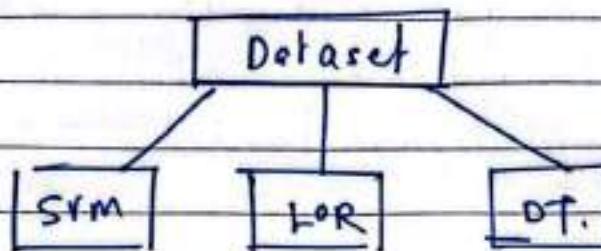
- Collectively these models called as Ensemble model. provided model should be different from each other.
- for making the models different we can adopt two way
 - ① Algorithm chosen for the model should be different
Ex. LR, SVM, NB, KNN like
 - ② Or if we select same algorithm for all the model then choose different data for all the model by some sampling method.

- In classification problem output is decided on majority count for the particular class predicted by all the model.
- In regression mean or average of all the output predicted by model.

Types of Ensemble.



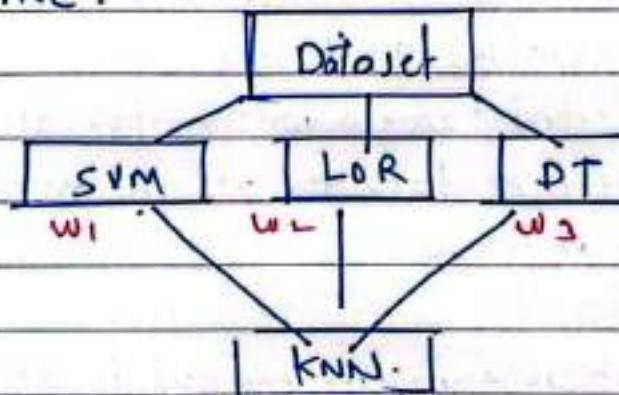
Voting Ensemble :- this depends on model should be different or particularly different algorithm.



for classification :- output is given by majority of count.

for regression :- output is calculated by mean of output given by all the algorithm.

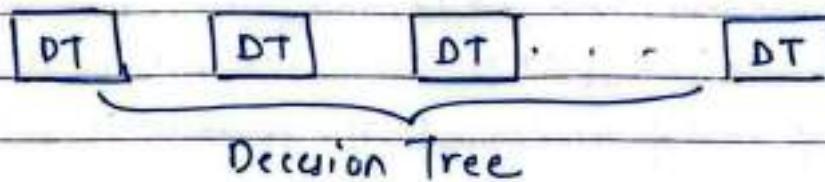
stacking :- unlike voting here we give additional weightage to the models based on their performance.



- Here we add one more additional algorithm which is trained on the output of above models or ensemble model.
- this newly added algorithm gives weightage to each model based on their performance based on the training data.
- and after this newly added model will predict the final result.

Bagging :- Bootstrap Aggregation

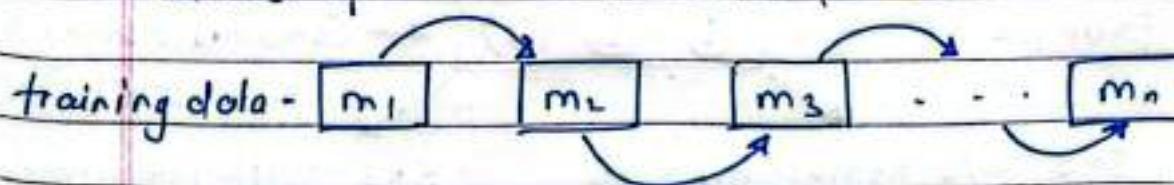
Here we use same algorithm for all the models just dataset is different for each model.



Random forest is special case of Bagging where base model is Decision tree. but we can other algorithm as well.

Boosting :- Boosting is an ensembling modeling

technique that attempts to build a strong classifier from the number of weak classifier . it is done by building model by using weak model in series firstly a model is built from training data . Then second model is built which tries to correct the error present in the first model . This procedure is continued and the model are added until either the complete training data set predicted correctly or the maximum number of models are added .



- Benefits :-
- 1) Improvement in performance.
- 2) It helps to achieve low bias and low variance.
- 3) Robust to variance.

When to use :- always.

Disadvantages :- Here we have to train multiple models which increase computational complexity.

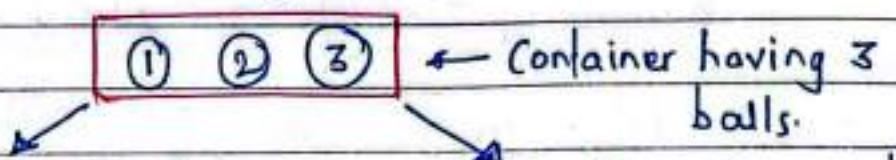
1) Let understand bagging and Random forest deeply.

- A Bagging classifier is ensemble method that fits base classifier each on random subsets of the original dataset and then aggregate their individual prediction (either by voting or averaging) to form final prediction
- Let understand bagging with Decision tree.
- Here we will create multiple Decision Tree, no. of decision tree is hyperparameter. Let say '10'
- and suppose we have 1000 rows dataset. Now how to distribute among the 10 models?
- if we split it in 10 each model will get only 100 rows. Since it is less data it will create bad model.
- So each model should get enough data and also different data.

there are two techniques for sampling :-

- 1) Simple random sampling with replacement
- 2) Simple random sampling without replacement.

Example



SRS w/o Replacement

① ②
① ③
② ③

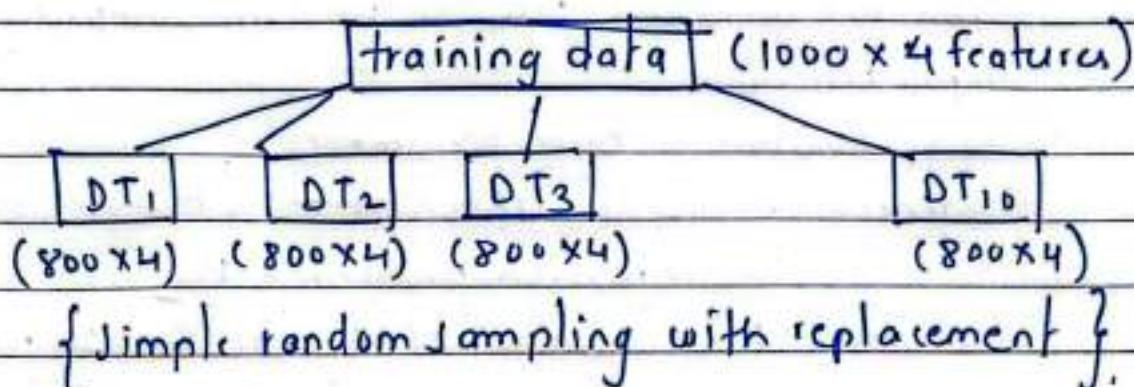
{Each row/column }
should not
repeated twice

① ① ① ②
② ② ① ③
③ ③ ② ③

{each row or column
can repeat
multiple times }.

so to select how many rows randomly for each model it depends, or user generally people prefer 80% of the dataset.

- so earlier due to splitting each model might get 100 only but by this technique each model will get 800 different dataset for training.



- Now its time to aggregate result by combining method or Voting.

Voting: if to predict from class 1 or class 0

$$DT_1 \rightarrow 1$$

$$DT_2 \rightarrow 0$$

$$DT_3 \rightarrow 0$$

.

$$DT_{10} \rightarrow 1$$

out of 10

$$8 DT \rightarrow 1$$

$$2 DT 0$$

{ so in this case due to majority input or test query belong to class 1 }

hyperparameter :- i) Number of DT's

ii) How much data you want to provide for each DT model.

Random forest :- if number of feature increases

then complexity of creating decision tree also increases to avoid that this technique is used.

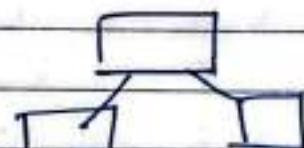
- instead of selecting random number of rows we select random number of columns.
- the each model will get lesser number of columns for this we will use simple random sampling w/o replacement.
- no of column is hyperparameter.
generally for classification = \sqrt{P}
regression problem = $P/2$
where P is number of features in original dataset
- so for the columns we will do SRS w/o Replacement
but for the rows we will do SRS with Replacement.
- why it is called random sampling?
→ because everytime we select the features randomly.
- hyperparameters : 1) No. of trees
2) No of features
3) No of datapoints.

Boosting :- 1) Adaboost

2) Gradient boost

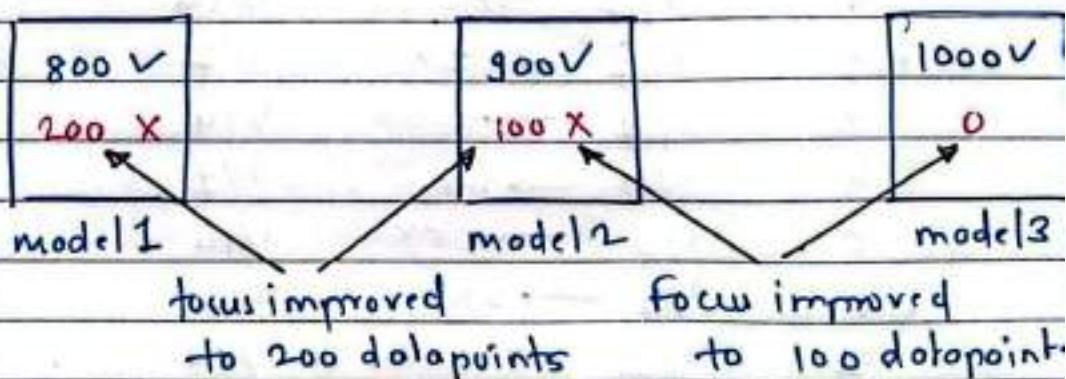
3) Xgboost.

1) **Adaboost**: - also called as Adaptive boosting is a technique in Machine learning used as an Ensemble method. The most common estimator used with Adaboost is decision trees with one level which means Decision tree, with only 1 split. These tree are also called Decision stumps.



Boosting is nothing but whenever my model performing bad need improve focus there and wherever it is performing good need to focus less there.

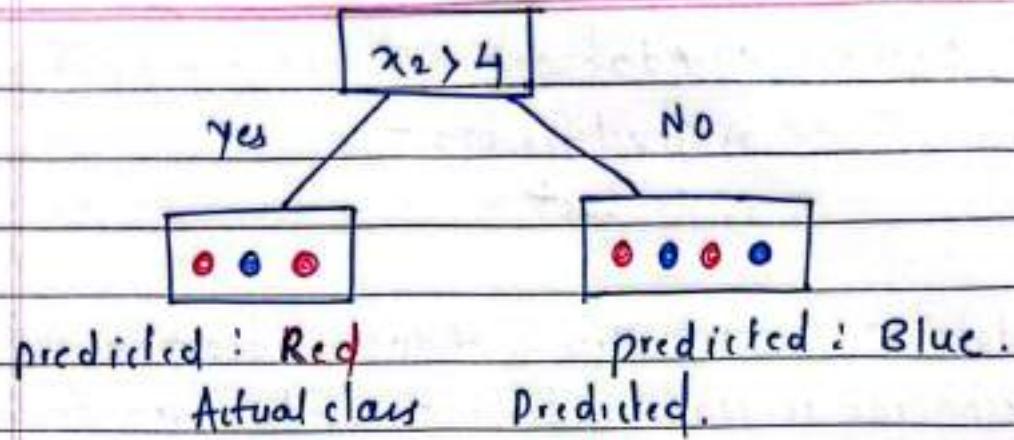
Workflow :→



- how many model need to create that user can decide i.e it is hyperparameter.

Let understand
with
Example.

x_1	x_2	y	$P(n-1) = 2(7-1) = 2 \times 6 = 12$
1	7	●	$P = \text{no. of features}$
2	6	○	$n = \text{no. of rows}$
3	5	●	so with total 12 condition we can split DT for stump.
4	4	●	
5	3	○	
6	2	●	Let say $x_2 > 4$.
7	1	○	



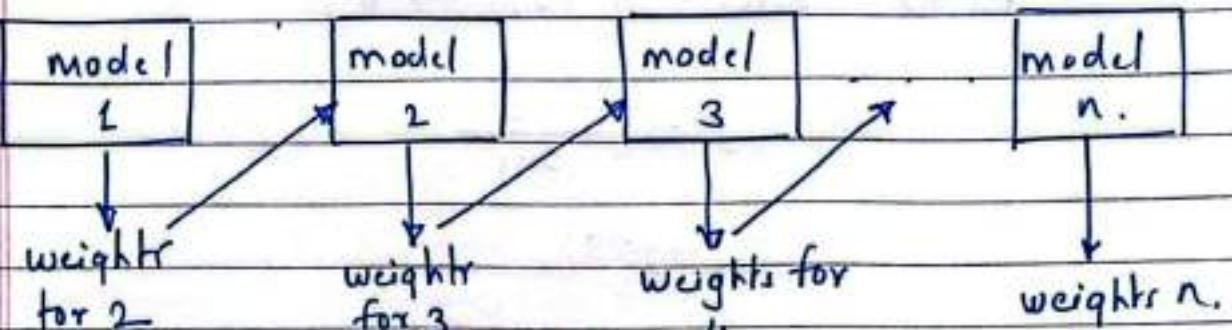
1	Red	Red
• 2	blue	Red → wrong prediction
3	Red	Red
• 4	Red	blue → wrong prediction
5	blue	blue
• 6	Red	blue → wrong prediction
7	blue	blue

Now let understand my distribution of focus.

	earlier focus	New focus	updated focus
1	$1/7 \rightarrow 0.5x$	$1/14$	$1/16$
• 2	$1/7 \rightarrow 2x$	$2/7$	$4/16$
3	$1/7 \rightarrow 0.5x$	$1/14$	$1/16$
• 4	$1/7 \rightarrow 2x$	$2/7$	$4/16$
5	$1/7 \rightarrow 0.5x$	$1/14$	$1/16$
• 6	$1/7 \rightarrow 2x$	$2/7$	$4/16$
7	$1/7 \rightarrow 0.5x$	$1/14$	$1/16$
		$\sum = 1.14$	$\sum = 1$

- Since 2, 4, 6 are wrong prediction need improve the focus. and rest 1, 3, 5, 7 all correct will reduce the focus or we can also called it as weights.
- for ex. we will improve our focus by 2 for incorrect prediction and reduce by 2

- sum of updated focus should be not greater than 1
since new focus was $16/16 = 1$. We try to normalize it.
- to do so we divide New focus by $16/16$ and we get new updated focus whose sum is equal to 1
- thus whichever point get higher focus/wt will get high priority in next model and those point with lesser focus/wt will get lesser priority.



- so we will create new dataset for next model based on weight received from last model.
- general inference is that since my first point - $1/16$ it will repeat one time and my 2nd point - $4/16$ it will repeat four times and same for $4, 6$.

x_1	x_2	new wt	Bucket range
1	7	0.0625	0 - 0.0625
2	6	0.25	0.0625 - 0.3125 (big bucket)
3	5	0.0625	0.3125 - 0.375
4	2	0.25	0.375 - 0.625 (big bucket)
5	3	0.0625	0.625 - 0.6875
6	2	0.25	0.6875 - 0.9375 (big bucket)
7	1	0.0625	0.9375 - 1

2, 4, 6 comes in big range bucket hence their chance of getting focus is more while sampling with replacement.

for this example we consider constant 2 for increasing or reducing the focus but how we decide this constant.

$$\text{constant} : e^{2k}$$

$$\text{where} : e = 2.718$$

λ = learning rate (it could be from 0 to 1)

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

for ex error = 0.107. $\rightarrow 0.1$

$$k = \frac{1}{2} \log\left(\frac{1 - 0.1}{0.1}\right)$$

$$> \frac{1}{2} \log(9)$$

$$k = \log \sqrt{9}$$

$$e^{\log \sqrt{9}} = \sqrt{9} = 3.$$

so instead of 2 for 10% error we would multiply or divide by 3 to increase or decrease our focus.

Summary:-

step :- Divide data into train and test

2 :- Create weak learner (stump)

3 :- Calculate error, find multiplication factor, value of k.

4 :- update the focus and update the data, create new model with updated data.

5 :- repeat step 2 and 4 number of this repetition equal to number of models.

6 :- Combine models for voting.

Gradient Boosting :- I'll understand with example of regression problem which is easy to understand & GBM compare classification problem.

iq	cgpa	salary (LPA)
90	8	3
100	7	4
110	6	8
120	9	6
80	5	3

We will create gradient boosting Model of 3 base model

M ₁ (mean)	M ₂ (DT)	M ₃ (DT)
--------------------------	------------------------	------------------------

Output of first model i.e. M₁ will always be mean of output or target column of dataset

→ Here we have it is 4.8 → pred 1

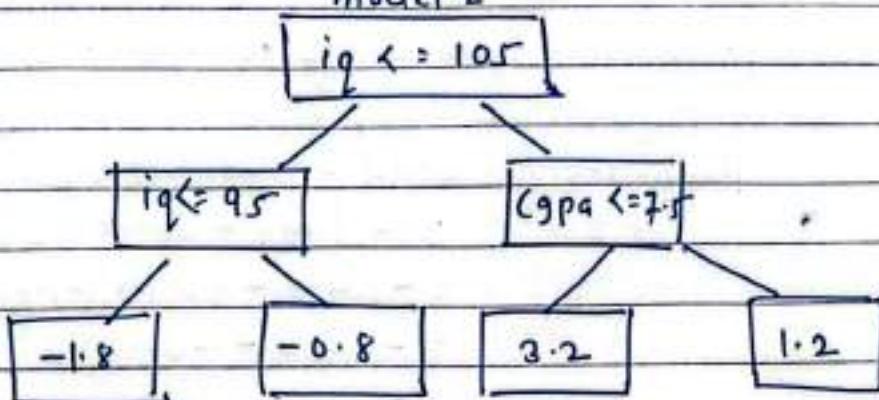
iq	cgpa	salary	pred 1	pred 01	pred 2	pred 02	pred 3
90	8	3	4.8	-1.8	-1.8	-1.62	-1.62
100	7	4	4.8	-0.8	-0.8	-0.72	-0.72
110	6	8	4.8	3.2	3.2	2.88	2.88
120	9	6	4.8	1.2	1.2	1.08	1.08
80	5	3	4.8	-1.8	-1.8	-1.62	-1.62

- In order to tell the error done by M₁ to M₂ we will calculate loss function of M₁

$$\text{actual} - \text{predicted} = \text{pseudo-residual}$$
.
- Now we will train model 2 which is DT for which input will be same (iq, cgpa) but output will be pseudo residual and not salary.

Inference:- Here we are asking to model 2 to tell us that how mistake or error is model 1 doing.

Decision Tree for
model 2



Now we will do prediction of model 2 on my data i.e pred 2

Let see how would i predict if i would have GBM of only two based model.

$$\text{prediction} = m_1 + m_2$$

Let calculate for student 1 (90, 8)

$$= 4.8 + (-1.8)$$

$$= 4.8 - 1.8 = 3$$

student 2 (100, 7)

$$= 4.8 + (-0.8) = 4.8 - 0.8 = 4$$

student 3 = 8 | student 4 = 6 | student 5 = 3.

\therefore which is all equal to actual target column, thus this problem of overfitting. like it will good for training set but will not the same for test data.

that's why in prediction we add learning rate.

$$\text{prediction} = m_1 + lr \times m_2$$

generally $lr = 0.1$

$$\begin{aligned} \text{Now student 1: } 4.8 + (0.1)(-1.8) &= 4.8 + (-0.18) \\ &= 4.62 \end{aligned}$$

with the help of learning rate we reduce the output and go closer to actual output gradually. (stepwise)

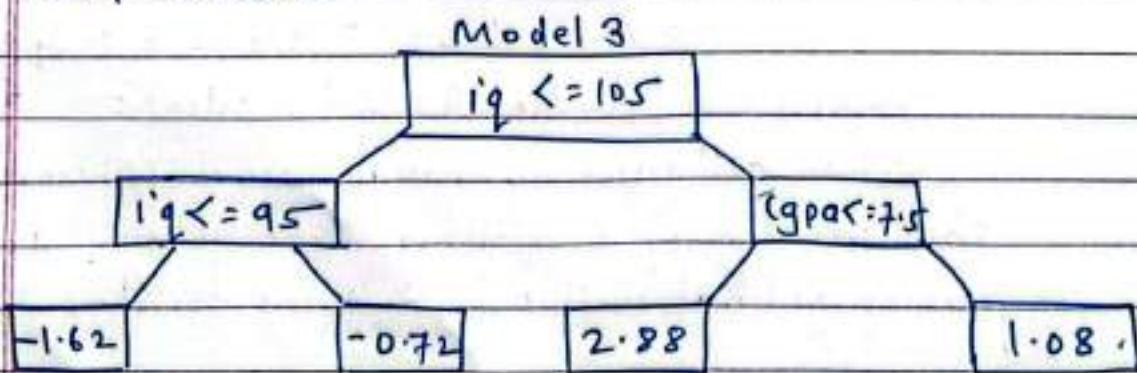
Now we will calculate model 3 for which we should know the performance of model 2 for that again we will calculate pseudo residual.

$\text{pseudo-resid} = \text{actual} - \text{predicted by } m_2$

$$\begin{aligned} \text{st 1} &= 3 - (m_1 + 0.1 \times m_2) \\ &= 3 - (4.8 + (0.1 \times (-1.8))) \\ &= 3 - (4.8 + (-0.18)) \\ &= 3 - (4.8 - 0.18) \\ &= 3 - 4.62 = -1.62 \end{aligned}$$

same for student 2/3/4/5 and it will be pseudo-residual 2 so we can clearly observe after each shift of model our residual is shift toward zero.

- Now we will create model 3 which is again DT. whose input col will be iq and cgpa but target will pseudo residual 2.
- model 3 combinedly predict the mistakes of model 1 and model 2.



- Now we will find prediction 3.

Now y -prediction by GBM will be

$$m_1 + 0.1 \times m_2 + 0.1 \times m_3$$

Let calculate student $x(60, 4.9)$

$$= 4.8 + 0.1(-1.8) + 0.1(-1.62)$$

$$= 4.8 - 0.18 - 0.162$$

$$= 4.458$$

earlier it would be 4.62 now it is 4.458

Difference between Adaboost Vs Gradient Boost

Adaboost

- Here we use decision stump. (DT whose depth is 1)
mean max leaf node could be 2

- models $y = (m_1, m_2, \dots, m_n)$
are prioritized based on their weights
 $w_1 m_1 + w_2 m_2 + w_3 m_3$

Gradient Boost

- In GBM max allowed leaf node could be from 8 to 32

- Here we use learning rate concept which is same for all the models.

Now y -prediction by GBM will be

$$m_1 + 0.1 \times m_2 + 0.1 \times m_3$$

Let calculate student $x(60, 4.9)$

$$= 4.8 + 0.1(-1.8) + 0.1(-1.62)$$

$$= 4.8 - 0.18 - 0.162$$

$$= 4.458$$

earlier it would be 4.62 now it is 4.458

Difference between Adaboost Vs Gradient Boost

Adaboost

- Here we use decision stump. (DT whose depth is 1)
mean max leaf node could be 2

- models $y = (m_1, m_2, \dots, m_n)$
are prioritized based on their weights
 $w_1 m_1 + w_2 m_2 + w_3 m_3$

Gradient Boost

- In GBM max allowed leaf node could be from 8 to 32

- Here we use learning rate concept which is same for all the models.

Gradient Boosting for classification problems

Let's understand how gradient boosting works for classification.

- when we use gradient boost for classification the initial prediction for every individual is the log(odd)

- $\log(\text{odds})$ for target : $\log \left(\frac{\text{people love the movie}}{\text{people don't love the movie}} \right)$
 $= \log \left(\frac{4}{2} \right) = 0.7$ (rounded value)
∴ this is initial prediction

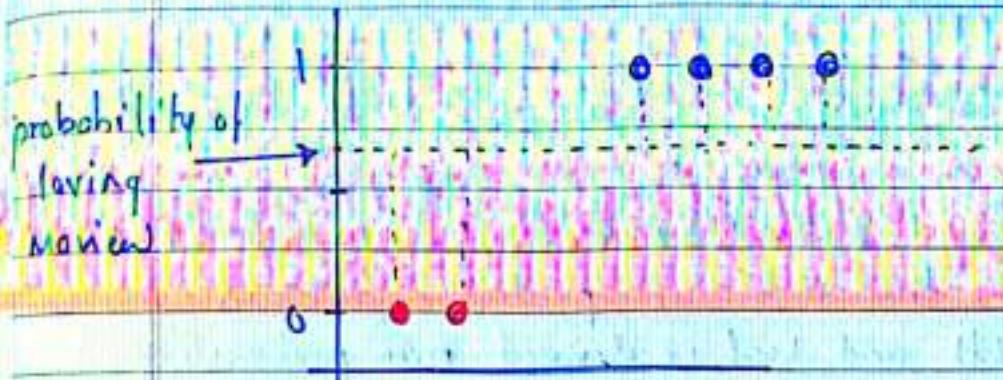
- But like with logistic regression the easiest way to use the log(odd) for classification is to convert it to a probability and we do it with logistic function.

probability of loving movie : $\frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}}$
 $= \frac{e^{0.7}}{1 + e^{0.7}} = 0.7$ (rounded value)

- Now classifying everyone in training dataset in someone who love movie is pretty lame because two of people do not love movie.

We can measure the how bad initial prediction is by calculating pseudo residuals, the difference between observed and predicted values.

$$\text{Residual} = (\text{observed} - \text{predicted})$$

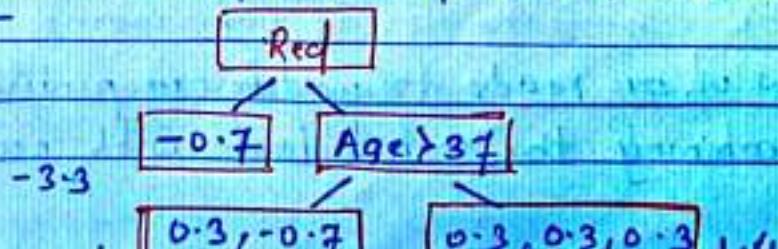


Observed value could be 0 or 1 and predicted is 0.7 if we calculate residual for initial prediction.

like popcorn	Age	favorite color	Love Movie	Residual
--------------	-----	----------------	------------	----------

yes	12	Blue	yes	0.3
yes	27	green	yes	0.3
no	44	blue	no	-0.7
yes	19	Red	no	-0.7
no	32	green	yes	0.3
no	14	blue	yes	0.3

Now we will build new tree for next model using "like popcorn and age" as input and Residual as output or target



- Note:- Just like when we used Gradient boost for regression, we are limiting the number of leaves that we will allow in the tree, here we are allowing no. of leaves to 3.

- In practice people often set the maximum number of leaves to be between 8 and 32

- In GB for regression a leaf with single residual had output value equal to that residual. In contrast when we use GB for classification the situation is little complex. this is because prediction are in terms of log(odd) and leaf's is derived from probability. so we can't just add them together to get new log(odd's) prediction without some sort of prediction
- o- for output of leaf 1st

$\frac{\sum \text{Residual}}{\sum \text{Prev. Probability} \times (1 - \text{prev probability})}$

$$\therefore \frac{-0.7}{0.7 \times (1 - 0.7)} = -3.3$$

-o- for output of leaf 2nd

$$\frac{0.3 + (-0.7)}{(0.7 \times (1 - 0.7)) + (0.7 \times (1 - 0.7))} = -1$$

-o- Output value for leaf 3.

$$\frac{0.3 + 0.3 + 0.3}{[(0.7 \times (1 - 0.7)) + (0.7 \times (1 - 0.7)) + (0.7 \times (1 - 0.7))]} = 1.4$$

Now we are ready to update our predictions by combining the initial leaf with new tree.

Prediction = initial leaf + learning rate x current leaf

- Let take learning rate is 0.8 which is very large but it for illustrative purpose however 0.1 is most common.

Now predict for each person, person L

$$0.7 + (0.8 \times 1.4) = 1.8$$

$$\text{probability} = \frac{e^{1.8}}{1+e^{1.8}} = 0.9$$

Similar we do for all the person.

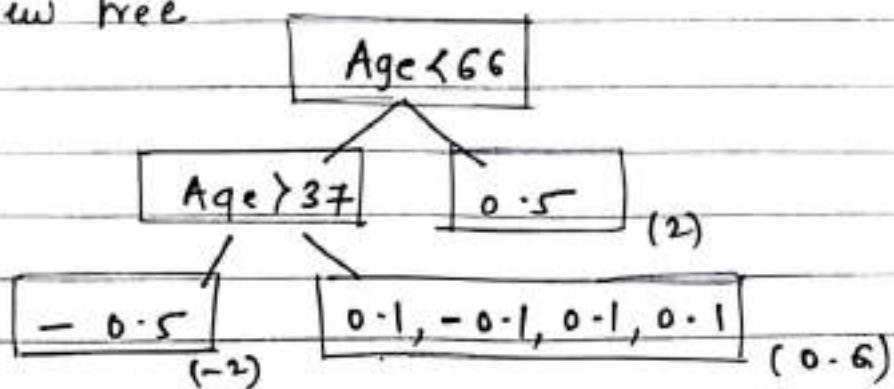
like popcorn	Age	favorite color	love movie	Residual	predicted	Residual
yes	12	Blue	yes	0.3	0.9	0.1
yes	87	green	yes	0.3	0.5	0.5
No	44	blue	No	-0.7	0.5	-0.5
yes	19	Red	No	-0.7	0.1	-0.1
No	32	green	yes	0.3	0.9	0.1
No	14	blue	yes	0.3	0.9	0.1

These new predicted probabilities are want than before & that's why we build lot of tree and not just one

- And now just like before we calculate new residuals, since Residual = observed - predicted.

$$\begin{aligned}
 \text{for person 1st: } (1-0.9) &= 0.1 \\
 \text{2nd: } (1-0.5) &= 0.5 \\
 &\vdots && \left. \right\} \text{ put in above table} \\
 \text{6th: } (1-0.9) &= 0.1
 \end{aligned}$$

Now with the help of new residual we will build new tree.



Output of Leaf 3rd: $\frac{\sum \text{Residual}}{\sum [\text{Prer prob} \times (1 - \text{Prer probability})]} = \frac{0.5}{0.5 \times (1 - 0.5)} = 2$

Output of Leaf 1st: $\frac{-0.5}{0.5 \times (1 - 0.5)} = -2$

Output of Leaf 2: $\frac{0.1 + (-0.1) + (0.1) + (0.1)}{(0.9 \times (1 - 0.9)) + (0.9 \times (1 - 0.9)) + (0.9 \times (1 - 0.9)) + (0.9 \times (1 - 0.9))} = 0.6$

Now we have calculated all of the output values for this tree we can combine it with everything else we've done so far.

This process repeats until we have made the maximum number of trees specified or the residual get super small.

If we need to classify a new person as someone who loves or not love the movie.

log(odd) prediction. first leaf + 0.8 × 2nd model + 0.8 × 3rd model
 or model

$$= 0.7 + (0.8 \times 1.4) + (0.8 \times 0.6) = \underline{\underline{2.3}}$$

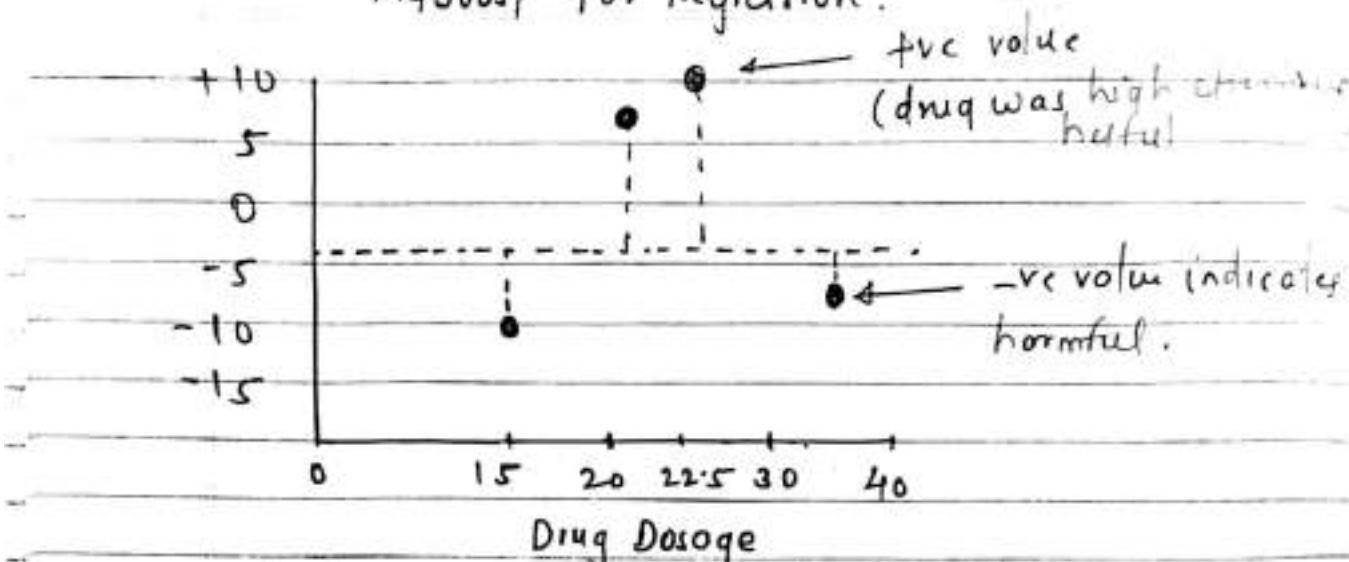
Now convert 2.3 into probability

$$\text{probability} = \frac{e^{2.3}}{1 + e^{2.3}} = 0.9$$

and thus predicted probability that this individual loves the movie is 0.9.

- Since we have kept threshold 0.5 and $0.9 > 0.5$ we classify this new person to class 1 which means he will love this movie.

XGBoost for Regression.



initial prediction is 0.5 whether it is XGboost classification or regression.

$$\text{Residual} = \text{observed} - \text{predicted}$$

Let build 1st DT with residual, the tree start with single leaf where all residual will stored.

$$[-10.5, 6.5, 7.5, -7.5] \quad \text{1st leaf}$$

We will calculate similarity score for residual

$$\text{similarity score} = \frac{[\text{sum of Resid}]}{\text{Number of Residual} + \lambda}^2$$

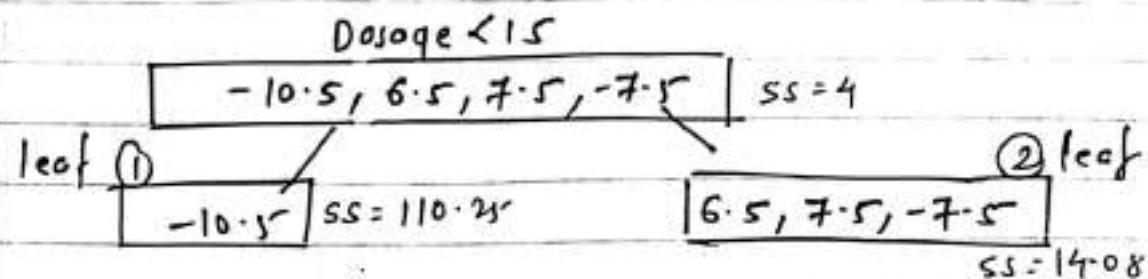
where λ = regularization parameter.

Assume $\lambda = 0$

$$\begin{aligned} & : \frac{[-10.5 + 6.5 + 7.5 + (-7.5)]}{4+0}^2 \\ & = \frac{(-4)^2}{4} = \frac{16}{4} = 4 \end{aligned}$$

Now question is whether or not we do better job of clustering similar residuals if we split them into two groups

- a] first for threshold we will consider two operation with lowest dosages and take their average for splitting



similarity score for 1st leaf : $\frac{(-10.5)^2}{1+3} = 110.25$

similarity score for 2nd leaf : $\frac{[6.5 + 7.5 + (-7.5)]^2}{3+1} = \frac{(6.5)^2}{3} = 14.08$

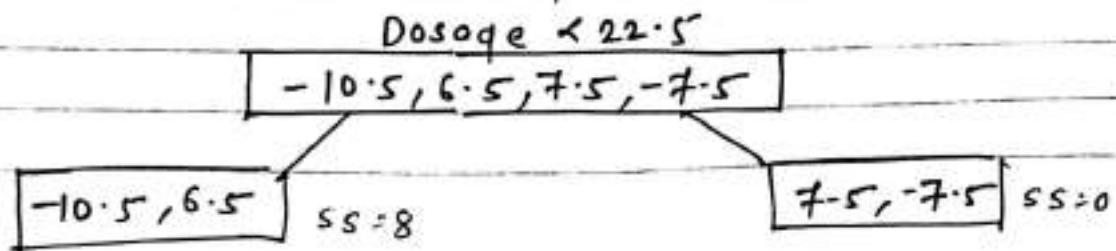
Note: when residual are very different they cancel out each other and similarity score is very low. In contrast Residual are same or only one they don't cancel each other similarity score is high.

- Now we need to quantify how much better the leaves cluster similar residuals than the root.

Gain = Left leaf similarity + Right leaf similarity - Root similarity.

Gain : $110.25 + 14.08 - 4 = \underline{\underline{120.33}}$.

- b] Now we will shift dosage threshold to 22.5



again we will calculate "Jmilarity score" for each leaf.

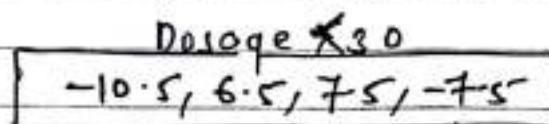
$$ss_{left} = \frac{[(-10.5) + 6.5]^2}{2+0} = 8$$

$$ss_{right} = \frac{[7.5 + (-7.5)]^2}{2+0} = 0.$$

$$\text{Gain}_{\cancel{\text{Root}}} = ss_{left} + ss_{right} - \text{Root } ss \\ = 8 + 0 - 4 = \underline{\underline{4}}$$

$\text{Gain } \cancel{\text{Root}} (\text{Dosage} > 22.5) < \text{Gain } (\text{Dosage} > 15)$

c] Now we will shift threshold value to 30.



$$ss = 4 \cdot 0.8 \quad \boxed{-10.5, 6.5, 7.5} \quad \boxed{-7.5} \quad ss = 56.25$$

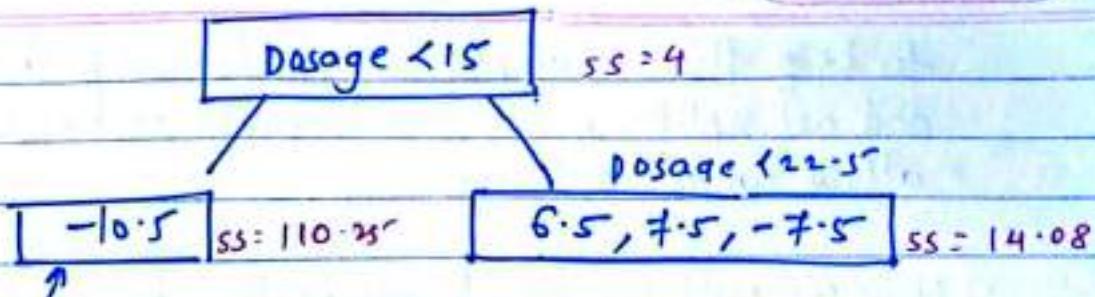
$$ss_{left} = \frac{[-10.5 + 6.5 + 7.5]^2}{3+0} = 4.08$$

$$ss_{right} = \frac{[-7.5]^2}{1} = 56.25$$

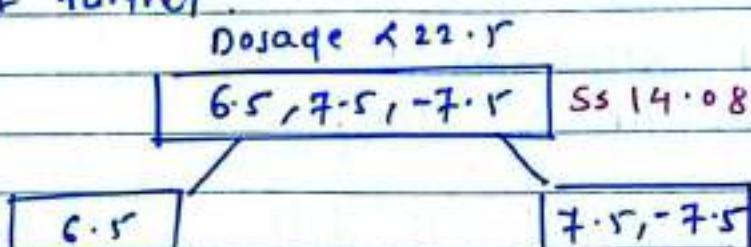
$$\begin{aligned} \text{Gain} &= ss_{left} + ss_{right} - ss_{Root} \\ &= 4.08 + 56.25 - 4 \\ &= \underline{\underline{56.33}}. \end{aligned}$$

$\text{Gain } (\text{Dosage} \cancel{> 30}) < \text{Gain } (\text{Dosage} > 30) < \text{Gain } (\text{Dosage} > 15)$

that's why we will split the root node at DT with threshold of 15 since it is giving highest Gain.



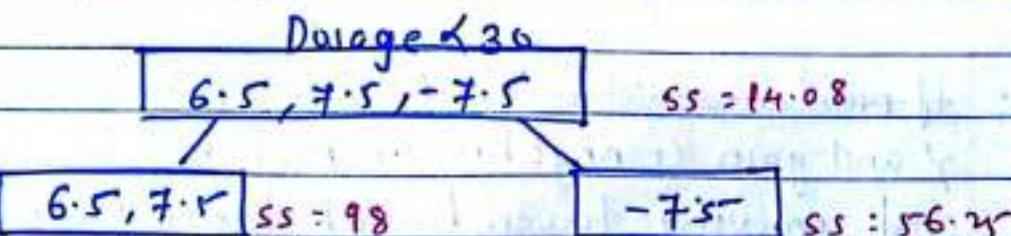
Since there is only one residual in left leaf we can't split it ahead. however we can split right leaf further.



$$\left. \begin{array}{l} \text{ss left} = 42.25 \\ \text{ss Right} = 0 \end{array} \right\} \text{See the previous formula.}$$

$$\begin{aligned} \text{Gain} &= \text{ss Left} + \text{ss Right} - \text{ss Root} \\ &= 42.25 + 0 - 14.08 \\ &= \underline{\underline{28.17}} \end{aligned}$$

Now again will shift threshold value to 30



$$\left. \begin{array}{l} \text{ss left} : 9.0 \\ \text{ss Right} : 56.25 \end{array} \right.$$

$$\text{Gain} = 9.0 \text{ ss Left} + \text{ss Right} - \text{ss Root} = \underline{\underline{140.17}}$$

$$\text{Gain}(\text{Dosage} < 22.5) < \text{Gain}(\text{Dosage} < 30)$$

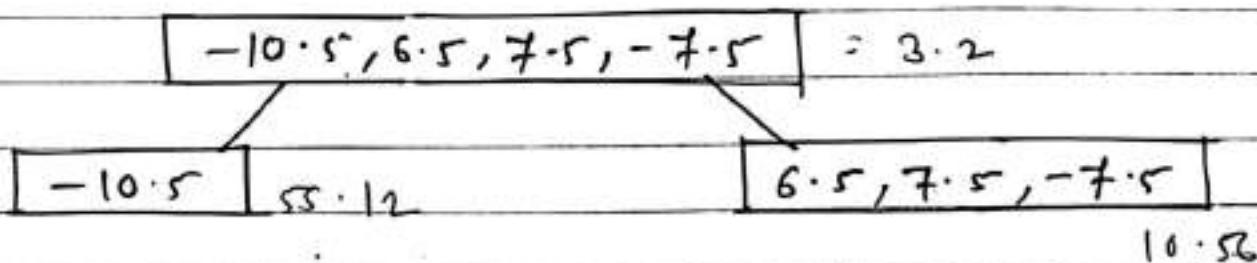
\therefore so this is perfect split.

to keep this example we are limiting it to two levels and not splitting further however in default it will allow up to 6 levels

- Now we will see how to prune this tree. In XG boost we prune based on its gain value.
- if we consider any random number like 130 this number called as γ here $\gamma = 130$.
- then we calculate the difference between γ and gain association with lowest branch of tree.
- if difference is negative we will remove the branch and if it is positive we will not remove the branch.
- To lowest branch gain $140 \cdot 17$
 ~~$\gamma = 130$~~
 $\text{Gain} - \gamma = 140 \cdot 17 - 130 = \underline{\underline{10 \cdot 17}}$
+ve value we will not remove the branch.
- if number would be 150 (i.e. $\gamma = 150$) the difference of γ and gain is negative at root but because we are not removing lowest branch hence we will not remove this as well
- if γ would be 150 then we had to remove lowest branch as well as root and whole tree would gone and it would be extreme pruning.
all this we did for $p=0$.

Now we will repeat all this step for $\lambda = 1$.

λ is regularization parameter which mean it is intended to reduce the prediction sensitivity to individual observation.



again similarity score at root

$$SS_{Root} = \frac{[-10.5 + 6.5 + 7.5 + (-7.5)]^2}{4+1} = 3.2$$

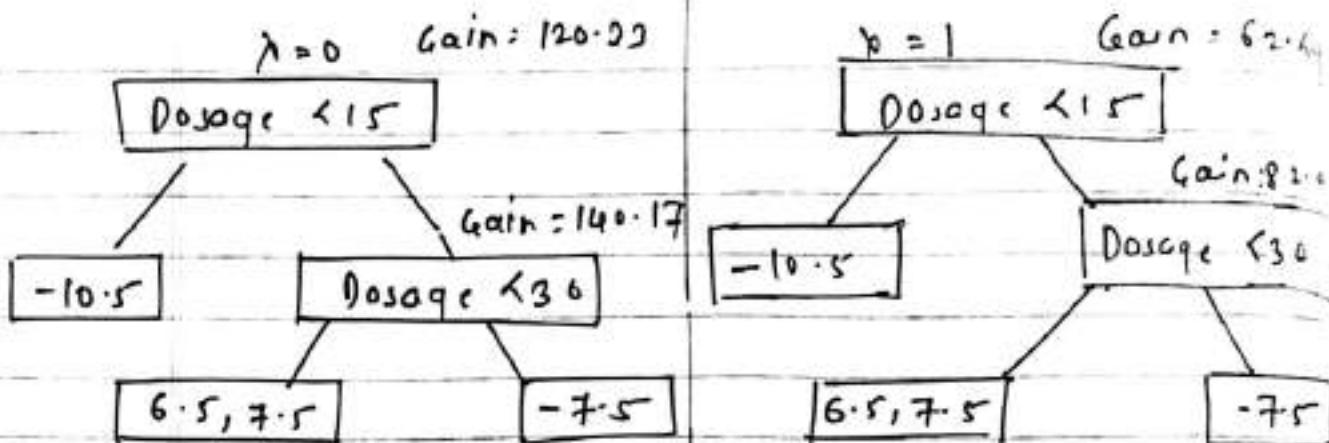
$$SS_{Left} = \frac{(-10.5)^2}{1+1} = 55.12$$

$$SS_{Right} = \frac{[6.5 + 7.5 + (-7.5)]^2}{3+1} = 10.56$$

Similarity score while $\lambda = 1$ are very less than similarity score while $\lambda = 0$, the amount of decrease is inversely proportional to the number of residual in the node.

$$\begin{aligned} \text{Gain} &= SS_{Left} + SS_{Right} - SS_{Root} \\ &= 55.12 + 10.56 - 3.2 \\ &= 62.48 \end{aligned}$$

and $\lambda = 0$ it was 120.33.



for pruning

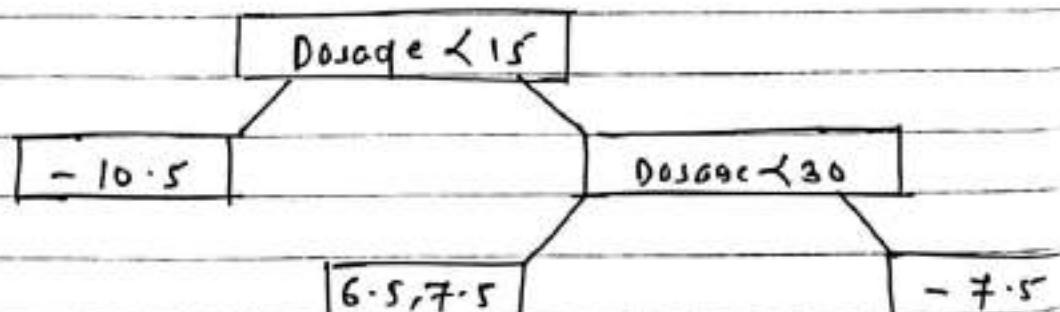
for number $\gamma = 130$
due to the difference we
didn't prune

for pruning

since lowest branch gain is
less than 130 we will remove it
and eventually prune whole
tree.

Note:- setting $\gamma = 0$ doesn't turn off pruning if
difference becomes negative it anyway will prune

on other hand $\gamma = 1$ did what it was supposed to
do , it prevent overfitting training data .



to calculate output :- [Sum of Residual]

Number of Residual + γ

$$o/p \text{ for 1st datapoint} = \frac{-10.5}{1+\gamma}$$

i) $\lambda = 0$ mean w/o regularization = -10.5

ii) $\lambda = 1$ with regularization factor = -5.25

| $\lambda > 0$ then it will reduce the amount that this individual observation adds to overall prediction.

thus λ (lambda) regularization parameter will reduce the prediction sensitivity to this individual observation

$$\text{output of leaf } \Theta : \frac{6.5 + 7.5}{2+0} = 7$$

when $\lambda=0$ o/p of value is simply avg of residual of leaf.

$$\text{output of leaf } \Theta : -7.5$$

Since we have build new tree we can make new prediction

like gradient boosting we have to use learning rate here also, default value is 0.3. Next is some

\therefore initial leaf + learning rate $\times 0.3$.

prediction for dosage 10.

$$0.5 + 0.3 (-10.5) = -2.65$$

Now we can see new residual is smaller than before this means we are taking small step in right direction

$$\text{Dosage} = 20 = 2.6$$

$$\text{Dosage} = 22.5 = -1.75 \quad \left. \right\} \text{smaller than before.}$$