

Dimensionality Reduction :-

Ex:-

Roll No Name Mobile Number

Date - 20/10/22 UNIT-3 Model Development

Regression :-

- Regression analysis is a statistical method to model the relationship between a dependent and independent variable with one or more independent variables.
- Regression analysis reveals average relationship between two variables and this makes possible estimation / Prediction.
- Regression used for prediction, forecasting, time series modeling, and determining the causal-effect relationship between variables.
- Ex:- Prediction of rain using temp. and other factors
Determining Market trends
Prediction of road accidents due to rash driving.

Simple Linear Regression :-

It is a type of Regression algorithms that models the relationships between a dependent variable and a single independent variable.

↳ The key point of Simple Linear Regression is that the dependent variable must be a continuous real value.

↳ ex:- Model relationship between income and expenditure.

Forecasting new observation.

Least Square Method

$$\sum (Y - Y_c)^2$$

Regression equation of Y on X , $[Y = a + bx]$

Formula:

$$\begin{cases} \sum Y = n a + b \sum X \\ \sum XY = a \sum X + b \sum X^2 \end{cases}$$

n = no. of observation.

Example-1

(in) (X)	(in) (Y)	XY	X^2
1	2	2	1
2	5	10	4
3	8	24	9
4	12	48	16
5	14	70	25
6	18	108	36
$\sum X = 51$		$\sum Y = 59$	$\sum X^2 = 91$
$\sum XY = 262$		$\sum X^2 = 91$	

$$\begin{aligned} \text{Eq } & y = na + bx \\ \Rightarrow & 59 = 6x_1 + b x_2 \\ \Rightarrow & 59 = 6a + 21b \quad \text{---(1)} \\ \text{Multiply } 7 \text{ on eqn (1)} & \\ 7(6a + 21b) &= 7 \times 59 \\ \Rightarrow 42a + 147b &= 413 \quad \text{---(2)} \end{aligned}$$

$$\begin{aligned} \text{Eq } & xy = ax + bx^2 \\ \Rightarrow & 262 = ax_1 + b x_1^2 \\ \Rightarrow & 262 = 21a + 91b \quad \text{---(3)} \\ \text{multiply (2) on eqn (3)} & \\ 2(21a + 91b) &= 2 \times 262 \\ \Rightarrow 42a + 182b &= 524 \quad \text{---(4)} \end{aligned}$$

Subtract eqn (3) from eqn (4), we get

$$42a + 182b - 42a - 147b = 524 - 413$$

$$\Rightarrow 35b = 111$$

$$\Rightarrow b = \frac{111}{35}$$

$$\Rightarrow \boxed{b = 3.17}$$

$$y = a + bx$$

$$\Rightarrow \boxed{y = -1.26 + 3.17x}$$

Put this on eqn (1)

$$6a + 21b = 59$$

$$\Rightarrow 6x_1 + 21 \times 3.17 = 59$$

$$\Rightarrow 6a + 66.57 = 59$$

$$\Rightarrow 6a = -7.57$$

$$\Rightarrow a = -\frac{7.57}{6} = -1.26.$$

<u>Q</u>	(x) over	(y) Runs	<u>xy</u>	<u>y²</u>
1		2		
2		5	10	25
3		8	24	64
4		12	48	144
5		14	90	196
6		18	108	324
<u>21</u>	<u>59</u>		<u>262</u>	<u>757</u>

$$\Sigma x = na + b \Sigma y$$

$$\Rightarrow 21 = 6a + 59b \quad \textcircled{1}$$

$$\Sigma xy = a \Sigma y + b \Sigma y^2$$

$$\Rightarrow 262 = 59a + 757b \quad \textcircled{2}$$

$$59(6a + 59b) = 21 \times 59$$

$$\Rightarrow 354a + 3481b = 1239 \quad \textcircled{3}$$

$$6(59a + 757b) = 6 \times 262$$

$$\Rightarrow 854a + 4547b = 1572 \quad \textcircled{4}$$

$$35/a + 4547b - 384953481b = 1572 - 1239$$

$$\therefore 1066b = 333$$

$$\therefore b = \frac{333}{1066}$$

$$\boxed{b = 0.31}$$

$$6a + 59b = 21$$

$$\therefore 6a + 59 \times 0.31 = 21$$

$$\therefore 6a + 18.29 = 21$$

$$\therefore 6a = 2.71$$

$$\therefore a = \frac{2.71}{6}$$

$$\boxed{a = 0.451}$$

$$X = a + b y$$

$$x = 0.41 + 0.$$

$$\therefore X = 0.41 + 0.31 y$$

$$X = 0.41 + 0.31 \times 36$$

$$\therefore X = 0.41 +$$

$$\boxed{X = 11.57}$$

Utility of Regression :-

- 1) Nature of Relationship
- 2) Estimation of Relationship.
- 3) Prediction
- 4) Useful in economic and Business Research.

Multiple Regression :-

→ Considering the values of the available multiple independent variables and predicting the value of one dependent variable.

→ The variables considered for the model should be

$$\Sigma Y =$$

Date - 27/10/22

Q) ~~Structure~~

$$\Sigma Y = n/a + b_1 \Sigma X_1 + b_2 \Sigma X_2$$

$$\Sigma Y X_1 = a \Sigma X_1 + b_1 \Sigma X_1^2 + b_2 \Sigma X_1 X_2$$

$$\Sigma Y X_2 = a \Sigma X_2 + b_1 \Sigma X_1 X_2 + b_2 \Sigma X_2^2$$

} Multiple Regression

stu-name	Marks	Liveclass	Book
A	8	3	4
B	9	4	5
C	7	3	3
D	10	5	5
E	6	2	3

Name	y	x_1	x_2	xx_{x_1}	xxx_2	x_1xx_2	y^2	x_1^2	x_2^2
A	8	3	4	24	32	12	64	9	16
B	9	4	5	36	45	20	81	16	25
C	7	3	3	21	21	9	49	9	9
D	10	5	5	50	50	25	100	25	25
E	6	2	3	12	18	6	36	4	9
	$\Sigma = 40$	$\Sigma = 17$		$\Sigma = 143$	$\Sigma = 166$	$\Sigma = 72$	$\Sigma = 3$	$\Sigma = 84$	

$$\Sigma y = n a + b_1 \Sigma x_1 + b_2 \Sigma x_2$$

$$y_0 = 5a + 17b_1 + 20b_2 \quad \textcircled{1}$$

$$\Sigma y x_1 = a \Sigma x_1 + b_1 \Sigma x_1^2 + b_2 \Sigma x_1 x_2$$

$$\therefore 143 = 17a + 63b_1 + 72b_2 \quad \textcircled{2}$$

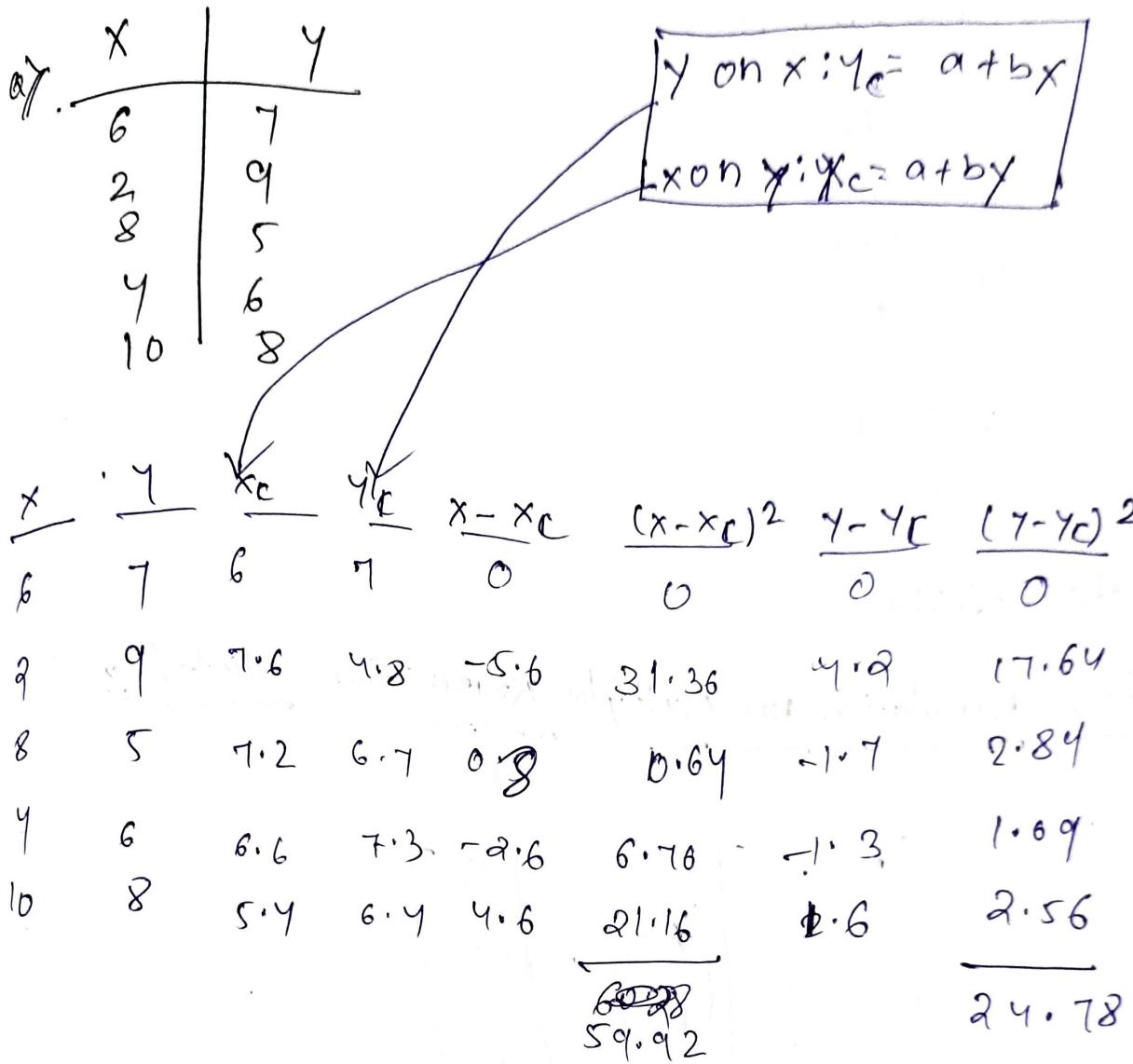
$$\Sigma y x_2 = a \Sigma x_2 + b_1 \Sigma x_1 x_2 + b_2 \Sigma x_2^2$$

$$\therefore 166 = 20a + 92b_1 + 84b_2 \quad \textcircled{3}$$

$$Y = 2.52 + 0.833x_1 + 0.66x_2$$

$$\begin{aligned} a &= 2.52 \\ b_1 &= 0.833 \\ b_2 &= 0.66 \end{aligned}$$

Regression & Correlation analysis :-



$$\text{Reg. eqn. } y \text{ on } x \quad y = 7.9 - 0.15x$$

$$\text{Reg. eqn. } x \text{ on } y \quad x = 10.2 - 0.6y$$

$$y_c = 7.9 - 0.15x$$

$$S_{ny} = \sqrt{\frac{\sum (x - \bar{x}_c)^2}{n}}$$

$$\Rightarrow y_c = 7.9 - 0.15 \times 6$$

$$= \sqrt{12.05611 \cdot 984}$$

$$\Rightarrow y_c = 7$$

$$= 3.47 \quad 3.46$$

$$S_{yn} = \sqrt{\frac{\sum (y - y_c)^2}{n}}$$

$$= \sqrt{\frac{24.78}{5}}$$

$$= \sqrt{4.946}$$

$$= 2.223$$

Date - 28/10/22

Matrix Formulation for Multiple Regression Model :-

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$Y = X\beta + \epsilon$$

find the coefficient of Regression in matrix form
from the data below

y	9	10	13	14	16
x_1	1	3	4	6	7
x_2	10	14	15	18	20

$$\beta_0 = 1$$

$$y = \begin{bmatrix} 9 \\ 10 \\ 13 \\ 14 \\ 16 \end{bmatrix}, x = \begin{bmatrix} 1 & 1 & 10 \\ 1 & 3 & 14 \\ 1 & 4 & 15 \\ 1 & 6 & 18 \\ 1 & 7 & 20 \end{bmatrix}$$

$$\beta = (x^T x)^{-1} x^T y = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}$$

$$x^T = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 3 & 4 & 6 & 7 \\ 10 & 14 & 15 & 18 & 20 \end{bmatrix}$$

$$x^T x = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 3 & 4 & 6 & 7 \\ 10 & 14 & 15 & 18 & 20 \end{bmatrix} \begin{bmatrix} 1 & 1 & 10 \\ 1 & 3 & 14 \\ 1 & 4 & 15 \\ 1 & 6 & 18 \\ 1 & 7 & 20 \end{bmatrix}$$

$$\begin{aligned} & \begin{bmatrix} 1+1+1+1+1 \\ 1+3+4+6+7 \\ 10+14+15+18+20 \end{bmatrix} \\ & \begin{bmatrix} 1+3+4+6+7 & 10+14+15+18+20 \\ 1+9+16+36+49 & 10+12+60+108+140 \\ 10+142+60+108+140 & 100+196+225+324+400 \end{bmatrix} \end{aligned}$$

$$= \begin{bmatrix} 5 & 21 & 77 \\ 21 & 111 & 360 \\ 77 & 360 & 1245 \end{bmatrix}$$

$$x^{-T} y = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 3 & 4 & 6 & 7 \\ 10 & 14 & 15 & 18 & 20 \end{bmatrix} \begin{bmatrix} 9 \\ 10 \\ 13 \\ 14 \\ 16 \end{bmatrix}$$

$$\therefore \begin{bmatrix} 9+10+13+14+16 \\ 9+30+52+84+112 \\ 90+140+195+252+820 \end{bmatrix}$$

$$= \begin{bmatrix} 162 \\ 287 \\ 997 \end{bmatrix}$$

$$\boxed{A^{-1} = \frac{1}{|A|} \cdot \text{adj } A}$$

$$\det(x^T x) = \begin{bmatrix} 5 & 21 & 77 \\ 21 & 111 & 360 \\ 77 & 360 & 1245 \end{bmatrix}$$

$$= 5 \left[111 \times 1245 + 360 \times 360 \right] - 21 \left[\dots \right]$$

$$= 5 (13875 + 129600) - 21 (26145 + 27720) + \frac{97560}{49541}$$

$$71 \cancel{4375} = 1131105 + 1240239$$

$$S_1 = \det(x^T x)$$

adj of $(x^T x) \Rightarrow$ coefficient of $(x^T x)$

$$\left[\begin{bmatrix} 111 & 360 \\ 300 & 1245 \end{bmatrix} - \begin{bmatrix} 21 & 360 \\ 77 & 1245 \end{bmatrix} + \begin{bmatrix} 21 & 111 \\ 77 & 360 \end{bmatrix} \right] \\ \left[\begin{bmatrix} 111 & 360 \\ 360 & 1245 \end{bmatrix} + \begin{bmatrix} 21 & 360 \\ 77 & 1245 \end{bmatrix} - \begin{bmatrix} 21 & 111 \\ 77 & 360 \end{bmatrix} \right] \\ \left[\begin{bmatrix} 111 & 360 \\ 360 & 1245 \end{bmatrix} + \begin{bmatrix} 21 & 360 \\ 77 & 1245 \end{bmatrix} - \begin{bmatrix} 21 & 111 \\ 77 & 360 \end{bmatrix} \right]$$

$$\left[(138195 - 129600) - (26145 - 27720) + (7560 - 8547) \right] \\ \left[-(138195 - 129600) + (26145 - 27720) - (7560 - 8547) \right]$$

$$\begin{bmatrix} 8595 & 1575 & -987 \\ 1575 & 298 & -183 \\ -987 & -183 & 114 \end{bmatrix} = \text{adj } A$$

$$\beta = (x^T x)^{-1} x^T y = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}$$

$$\Rightarrow \frac{1}{51} \begin{bmatrix} 8595 & 1575 & -987 \\ 1575 & 296 & 183 \\ -987 & 183 & 114 \end{bmatrix} * \begin{bmatrix} 63 \\ 287 \\ 947 \end{bmatrix} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}$$

$$\Rightarrow \frac{1}{51} \begin{bmatrix} 532890 + 452025 - 984039 \\ 97050 + 84952 - 18245 \\ -61194 - 52521 + 113658 \end{bmatrix} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}$$

$$\Rightarrow \frac{1}{51} \begin{bmatrix} 876 \\ 151 \\ -57 \end{bmatrix} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 17.17 \\ 2.96 \\ -1.11 \end{bmatrix} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}$$

$$y_i = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

$$y_i = 17.17 + 2.96 x_1 - 1.11 x_2$$

Date - 29/10/22

Model Evaluation and Selection :-

- ▷ Metrics for performance Evaluation
- ▷ Methods for performance Evaluation
- ▷ Methods for Model comparison
- ▷ Model Selection.

Confusion Matrix :-

A table that is often used to describe the performance of a classification model (or "classifier") on a set of test data for which the true values are known.

ix

		Predicted class		
		Yes	No	
Actual Yes Values	Yes	40	20	
	No	TP	Fn	60 (P)
No	Yes	10	90	100 (N)
	FP	TN		
				$P = \frac{TP + TN}{TP + TN + FP + FN}$ = 160

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

Date - 31/10/22

true positive (TP)

↳ Positive tuples that were correctly level by the classifier.

↳ negative tuples that were correctly leveled by the classifier.

false positive (FP)

↳ negative tuples that were incorrectly level by the positive.

↳ Positive tuples that were mislabeled as negative.

Performance matrices :-

1. Accuracy :-

↳ It is also called as Recognition Rate.

↳ Percentage of test set tuples that are correctly classified.

$$\frac{TP + TN}{P + N}$$

2. Error Rate :-

→ 1 - Accuracy

→ It is also called as Missclassification Rate

$$\text{1. } \frac{FP + FN}{P+N} \Rightarrow \frac{10 + 20}{160} = 0.19$$

calculate sensitivity :-

- ↳ It also refers to the TP recognition rate.
- ↳ The proportion of positive tuples that are currently identified.

$$\text{sensitivity} = \frac{TP}{P}$$

sensitivity :-

- ↳ It refers to the TN recognition rate.
- ↳ The proportion of negative tuples that are currently identified.

$$\text{specificity} = \frac{TN}{N}$$

Precision :-

- ↳ It is called measure of exactness. It means what % of tuples is truly positive.
- ↳ It lies between 0 to 1.

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

6. Recall :-

↳ It is a measure of completeness.

↳ what % of positive tuples are labeled as Positive.

$$\hookrightarrow \frac{TP}{TP + FN}$$

7. F-Score / F-Measure :-

$$2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{1.056}{1.46} = 0.72$$

↳ It is the harmonic mean of precision & recall.

↳ It gives same weightage to precision & recall.

↳ It is also called as F_1 score.

Date - 05/11/22 Measures of insample Evaluation :-

Loss Function

cost function

A function that calculate loss of 1 data point is known as Loss function.

$$\rightarrow (y_i^o - \hat{y}_i)^2$$

y_i^o = Actual Value

\hat{y}_i = Predicted Value,

A function that calculate loss of entire data is known as Cost function.

$$\text{By Mean sq. error} = \frac{1}{n} \sum_{i=1}^n (y_i^o - \hat{y}_i)^2$$

n = sample size.

Mean Absolute Error

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

lower the MAE

higher the accuracy

Arithmatic avg of absolute Errors.

Mean Bias Error

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

A positive bias means the error from the data is overestimated

A negative bias means the error is underestim.

Relative Absolute Errors

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

where, $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \rightarrow$ Mean of actual values.

Result - 0 - 1

0 → good model.

Mean absolute percentage error

$$MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| * 100\%$$

Root Mean Square error

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

RMSE = 0 → Perfect Model

Relative squared error :-

$$RSE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i^0 - \bar{y})^2}$$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i^0 \rightarrow \text{mean of actual values.}$$

Relative Root Mean Squared Error :-

$$RRMSE = \sqrt{\frac{\frac{1}{n} \sum_{i=1}^n (y_i^0 - \hat{y}_i)^2}{\sum_{i=1}^n (y_i^0)^2}}$$

excellent when RRMSE < 10%.
 good 10% to 20%.
 fair 20 to 30%.
 poor > 30%.
 Error %

Root Mean Squared Logarithmic Error :-

$$RMSE = \sqrt{\log(y_i^0 + 1) - \log(\hat{y}_i + 1)^2}$$

Coefficient of Determination (R squared) :-

It is used to analyze how differences one variable can be explained by a difference in a second variable.

Its ranges from 0 to 1 (0% to 100%).

If $R^2 = 0$, then dependent variable can't be predicted from the independent variable.

If $R^2 = 1$, then the dependent variable can be predicted from the independent variable.

If R^2 between 0 and 1, then the dependent variable can be predicted.

$$r = \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{\sqrt{[n\Sigma x^2 - (\Sigma x)^2][n\Sigma y^2 - (\Sigma y)^2]}}$$

$$R^2 = \frac{\text{variable can be predicted formula - 2}}{1 - (RSS/TSS)}$$

Residual sum square Total sum square

Date - 07/11/22

Polynomial Regression :-

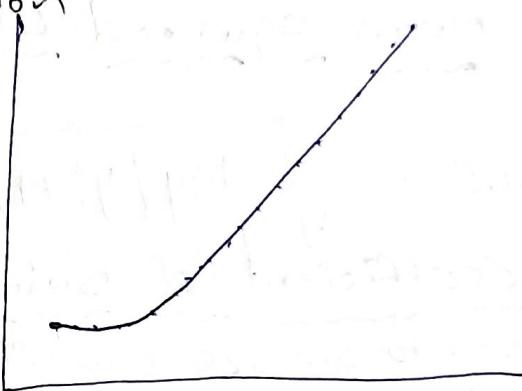
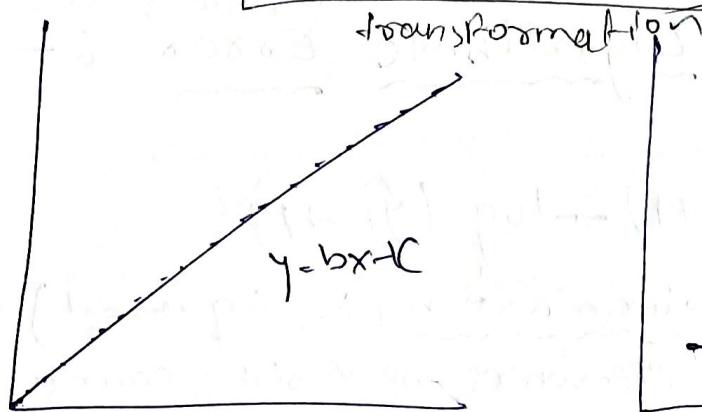
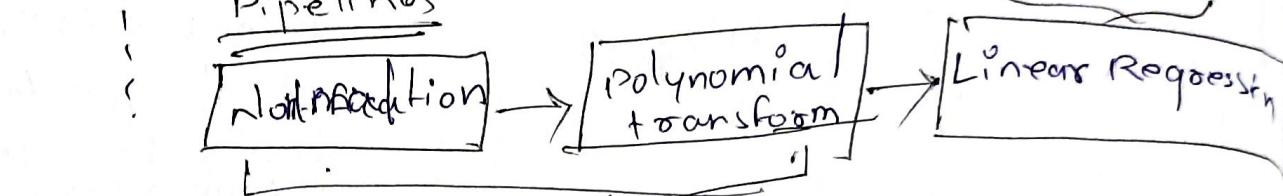
- Polynomial Regression is a Regression algorithm that models a relationship b/w a dependent(y) and independent variable(x) as nth degree polynomial
if $y = \text{constant} - 0\text{-degree}$.

$$y = b_0 x^0 + c - 0\text{-degree}$$

$$y = a_0 x^2 + b_0 x + c - 2\text{-degree polynomial}$$

Pipelines

Prediction



- Polynomial Regression fits a non linear relationships.

- Here the data points are arranged non-linear fashion.

NOTE
 If we apply a linear model on simple RQ in a linear data set, then it provides us a good result but if we apply the same model without any modification on a non linear data set, then it will produce a disastrous output. due to which loss function will increase, the error rate will be high and accuracy will be decreased. so for such cases we need the polynomial regression model.

Example :-

<u>x</u>	<u>y</u>	<u>$y_i x_i$</u>	x^2	x^3	x^4
3	9.5	28.5	9	27	
4	3.2	12.8	16	64	
5	3.8	19	25	125	
6	6.5	39	36	216	
7	12.5	87.5	49	343	4659
28	28	168.8	911.7	135	9659

$$\Sigma y_i = n a_0 + a_1 \Sigma x_i + a_2 (\Sigma x_i^2)$$

$$\Sigma y_i x_i = a_0 (\Sigma x_i) + a_1 (\Sigma x_i^2) + a_2 (\Sigma x_i^3)$$

$$\Sigma y_i x_i^2 = a_0 (\Sigma x_i^2) + a_1 (\Sigma x_i^3) + a_2 (\Sigma x_i^4)$$

$$275 = 5 a_0 + 25 a_1 + \frac{135}{625} a_2$$

$$158.5 = 25 a_0 + 135 a_1 + 775 a_2$$

$$965.2 = 135 a_0 + 775 a_1 + 4659 a_2$$

$$a_0 = 12.42$$

$$a_1 = -5.51$$

$$a_2 = 0.76$$

$$y = a_0 + a_1 x + a_2 x^2 \dots \sim a_n x^{n-2}$$

$$y = 12.42 - 5.51x + 0.76x^2$$

Date - 10/11/22

Residual plots for regression model validation ?

A Residual is a measure of how far away a point is vertically from the regression line. Simply, it is the error between a predicted value and the observed actual value.

$$\hookrightarrow \text{Residual}(\epsilon) = y - \hat{y}$$

Residual plot

A typical residual plot has the residual values on the y-axis and the independent variable on the x-axis.

Residual plot Analysis :-

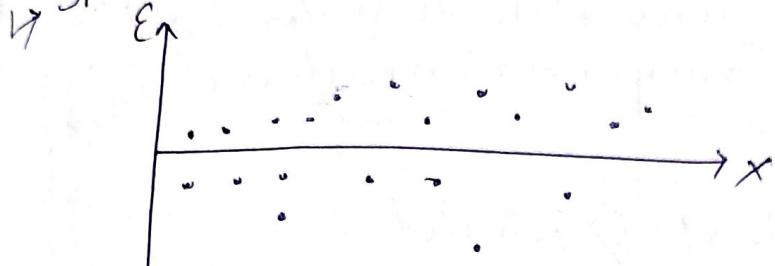
The most important assumption of a linear regression model is that the errors are independent and normally distributed.

Response = Deterministic + Stochastic

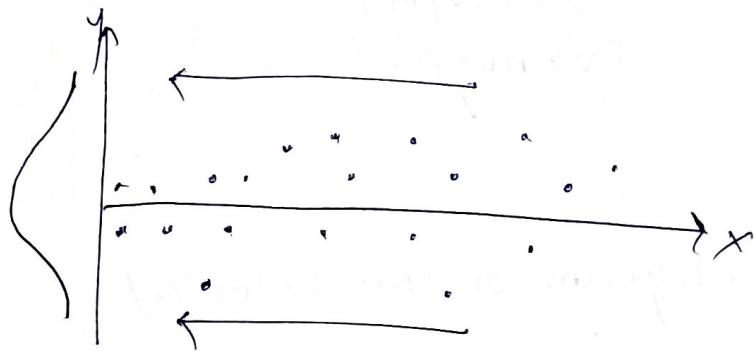
Characteristics of Good Residual plots :-

It has a high density of point close to the origin and a low density of point away from the origin.

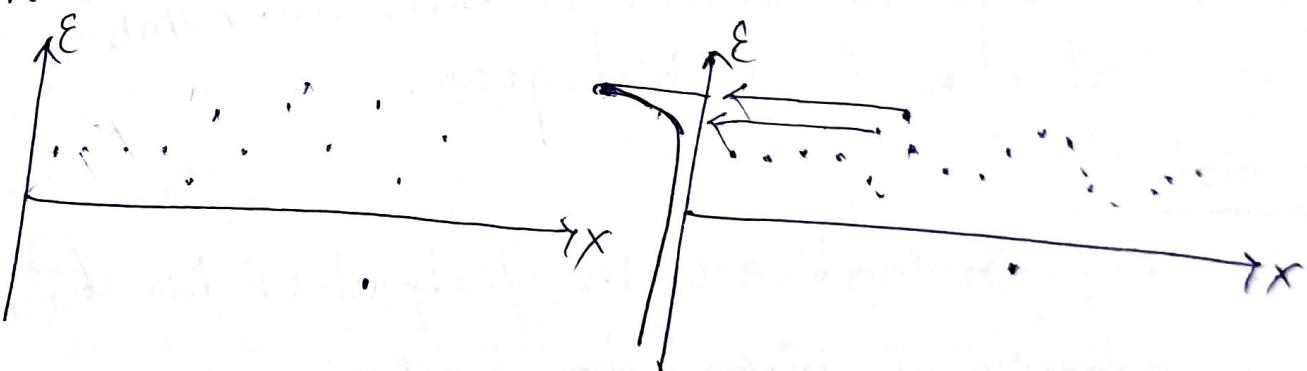
It is symmetric about the origin.



Good Residual plot



Residual errors are approximately distributed in the same manner.



Bad Residual plot

Distribution plot :-

These plot :-

↳ ~~Seaborn~~ ^{↳ seaborn is a python data visualization library based on Matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics.}

↳ These plot helps us to visualize the distribution of data. we can use these plots to understand the mean, median, range, variance, deviation, etc of the data.

↳ It is of n types of plots :- i) jointplot

ii) distplot

iii) Pairplot

iv) rugplot

a) distplot

↳ Dist plot gives us the histogram of the selected continuous variable.

↳ It is an example of a univariate analysis.

↳ We can change the number of bins i.e. number of vertical bars in a histogram.

b) joinplot :

→ It is the combination of the distplot of two vari

→ It is an example of bivariate analysis.

Data

pair plot :-

It takes all the numerical attributes of the data and plot pairwise scatter plot for two different variables and histograms from the same variables.

Rugplot :-

It draw a dash mark instead of a uniform Y-distribution as in displot.

It is an example of a univariate analysis.

Relationships :-

Prediction and Decision making :-

- ↳ Do the predicted values make sense
- ↳ Visualization
- ↳ Numerical measures for evaluation.
- ↳ comparing models

Date-12/11/22

↳ Mean square Error for a multiple Linear Regression Model will be smaller than the Mean square Error for a simple Linear Regression model, since the errors of the data will decrease when more variables are included in the model.

→ Polynomial regression will also have a smaller Mean Square Error than the linear regular regression.