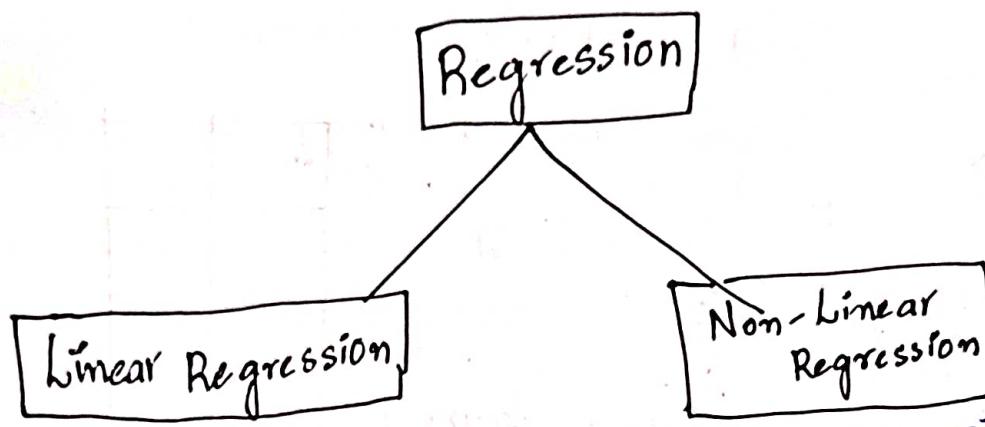


7/4/22
9:30pm

Regression

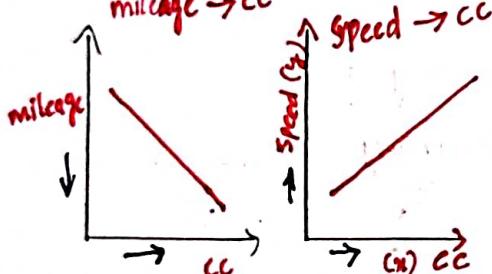
Predicting Output which is Continuous



(Assuming the relation between x & y)

is Linear")

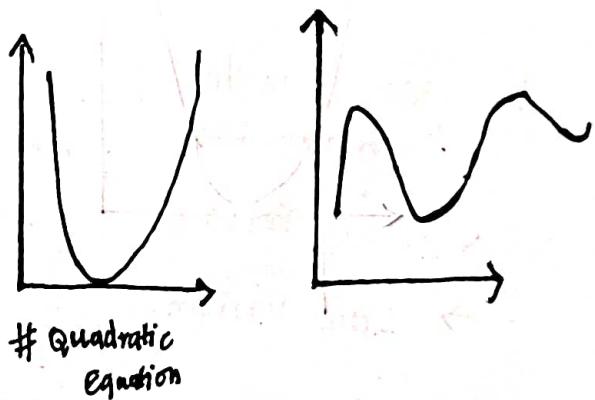
Ex:- In The Form of $(y = mx + c)$



* mileage decreases ↑ * cc increases ↑
* cc increases ↑ * speed increases ↑

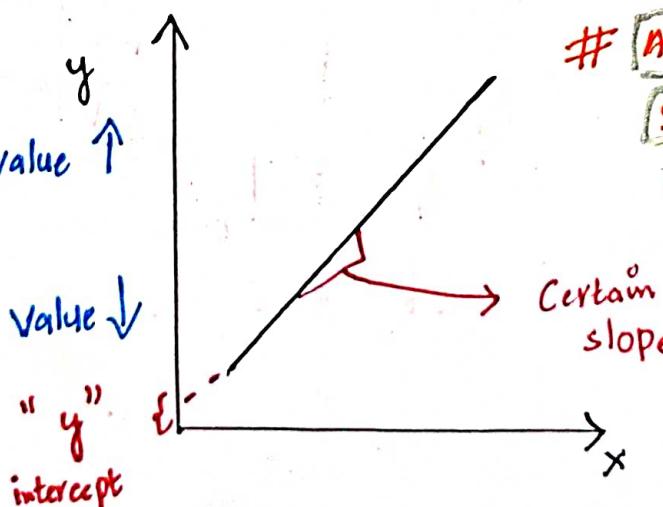
Non-Linear Regression

Any line other than straight line is Non-Linear



x value ↑, y value ↑
(+ve slope)

x value ↑, y value ↓
(-ve slope)



Any relation in terms of Straight Line, it is called as "Linear Regression"

Ex:-

X	Y
1	15
2	25
3	35
4	45
5	55
6	? 35

In Early Ages, They took "Avg" values.

Like

$$\frac{15 + 25 + 35 + 45 + 55}{5} = 35$$

7 35

8 35

→ They say, Avg Line is Best Fit Line;

* Best Fittants

tips (waiter)

- 1. 10 \$
- 2. 1 \$
- 3. 4 \$
- 4. 3.6 \$
- 5. 1.4 \$
- 6 → Avg 4 \$

Bill amount
(not given)

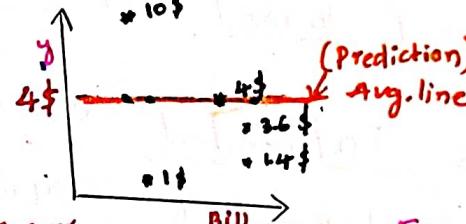
so, For only One
Variable / column - (Avg)

For Single
variable

Best Fit line

Average line

7 → 4 \$ Every time Answer
is 4 \$ (because of
Avg line)



For x value [Predict
Va
c

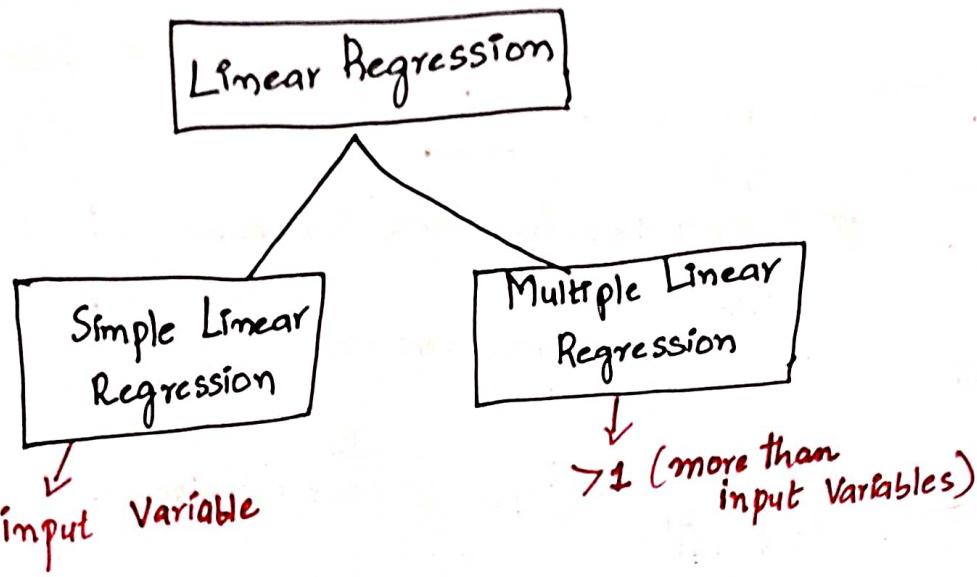
H₀: Average line is best fit line

H₁: Regression line > Avg. line.

Y	avg.line Y-Prediction	Error σ	Error σ^2
10	4	6	36
1	4	-3	9
4	4	0	0
3.6	4	-0.4	0.16
1.4	4	-2.6	6.76

0

51.92



Simple Linear Regression :

- In Simple Linear regression, We predict the value of One variable "Y" based on Other variable "X".
- "X" is independent variable (or) input variable (or) exploratory variable
"y" is dependent variable (or) output variable (or) response variable

Ex:-

Weight	CC	Speed	Mileage
indepn dvar	in. var	in. var	→ O/P (it depends on weight, CC, Speed)

Que :- Why it is Simple ?

Because, it Examines

Relationship

between

two

Ans :-

Variables Only.

X	Y

$$y = x$$

Que :- Why linear ?

independent

Variable increases

(or Decreases)

Ans :- When the

dependent Variable

increases (or Decreases) in

The dependent Variable

a linear Fashion.

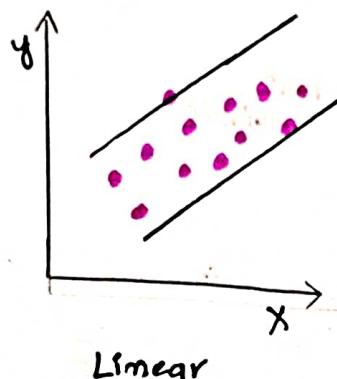
Que: How to check Linear Regression ?

(or) not

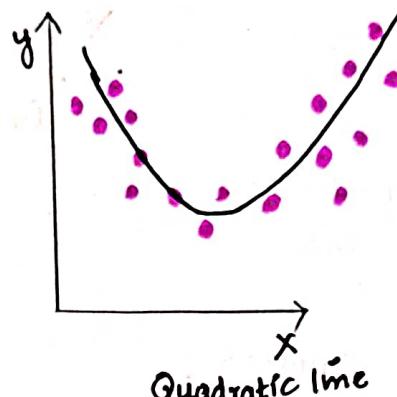
Ans:-

By scatter plot

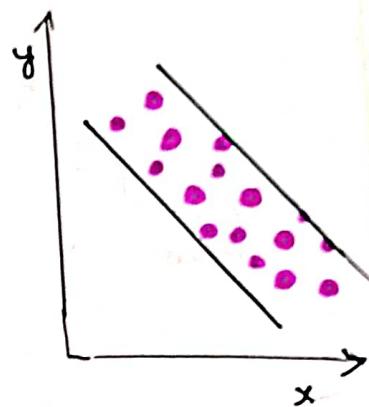
Linearity
direction



Linear



Quadratic line



Linear but
downward
direction

As "x" value increasing
"y" value also increasing

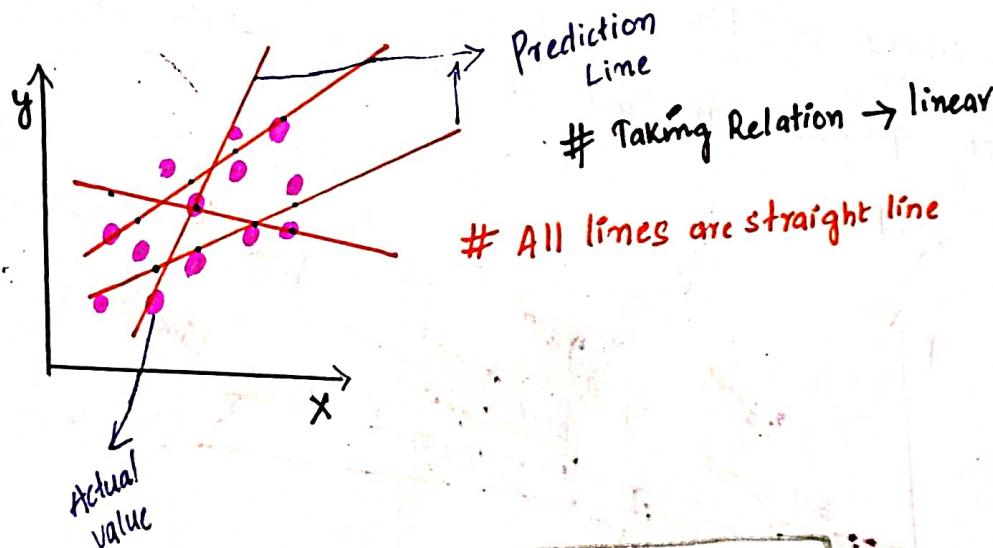
"Positive Direction"

As "x" value increasing
"y" value also decreases

"Negative direction"

Que:- Which Line to be Considered ?

Ans:-



Que:- Different notations ? But Same Concept ?

Ans:-

$$y = mx + c$$

slope

$$y = an + b$$

"y" intercept

$$y = b_0 + b_1 n$$

slope

"y" intercept

y intercept

$$y = \beta_0 + \beta_1 n$$

slope

"y" intercept

slope

Ques: What is Standard notation?

Ans :-

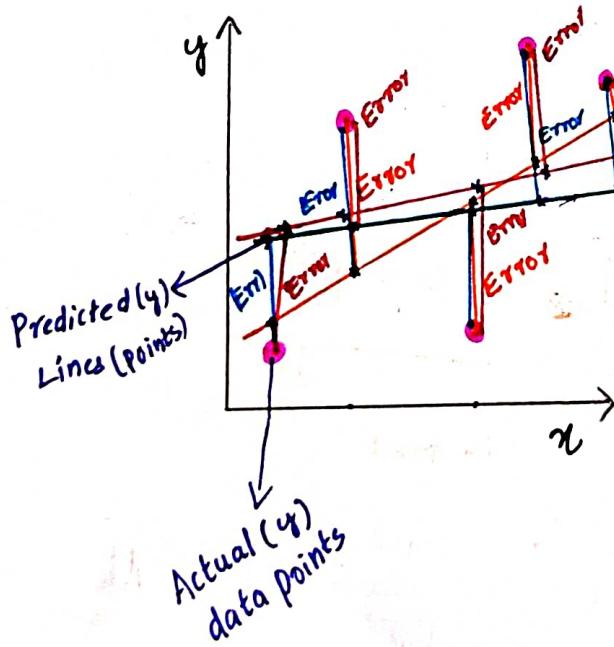
$$y = \beta_0 + \beta_1 x_1$$

Annotations:

- (dependent variable)
- "y intercept" value
- Constant coefficient
- Independent variable
- slope value

Ordinary Least squares :

- Least squares fitting is a way to find The best fit Curve
- Least squares fitting is a way to find The best fit Curve
 - Line for a set of points
 - Sum of Squares of Residuals Errors are used to estimate the best fit curve (or) Line
 - Least squares method is used to obtain The coefficient of "m" and "b"



best fit line?
Where The Sum of Squares of Error is Minimum

distance b/w Actual data point - Predicted data point [Error]

Square The errors [each error do square]
after Squaring Error, Then Sum Error

Que:- What is SSE?

Ans:- SSE \Rightarrow Some of squares of Error

Denoted by $\sum [y - \hat{y}]^2$ squaring them
= SSE_{min}

Sum
Actual value
Predicted value

$\sum [y - \hat{y}]^2$ $\underset{\text{min}}{\rightarrow}$ This value Should be Minimum
* you Have identify the line such a way $\sum [y - \hat{y}]^2$ min

Que :- How to identify The Exact line?

Ans:- Explanation:- OF Calculus

Minimum Point : $y = x^n$ (what is Least "y" value)

$y = 0$ # if 0.00001 also it get squared
so, "0" is Least y value

how it is identified

$\frac{dy}{dx} \text{ at } x=0$ \Rightarrow $\frac{dy}{dx} \text{ at } x=0 \Rightarrow \frac{d(x^n)}{dx} \text{ at } x=0 \Rightarrow 2x \Rightarrow 2(0) = 0$

Ex:- $y = 4x^2 + 5 \rightarrow 4x^2 = 8$

$\frac{dy}{dx} \text{ at } x=0 \Rightarrow 8x + 0 \Rightarrow 8(0) + 0 \Rightarrow 0$

* $\frac{d}{dx}(4x^2) = 2 \times 4 = [8x]$

* $\frac{d}{dx}(5) = 0$

$$\# \text{ Minimum Value} = \frac{dy}{dx} \text{ at } x=0$$

$$\# SSE_{\min} = \mathbb{E}[y - \hat{y}]^2$$

So, in order to get minimum value for S.S.E

calculate

$$\frac{d[\mathbb{E}[y - \hat{y}]^2]}{dx} \text{ at } x=0$$

$$\hat{y} = \beta_0 + \beta_1 x$$

$$\frac{d}{dx} \mathbb{E}(y - \beta_0 + \beta_1 x)$$

Expand This Equation ↑

$$\beta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

"y" Mean "x" Mean

Example :

$H_0: \text{Avg. line} \leq \text{Regression Line}$
 $H_1: \text{Regression line} < \text{Avg. line}$

X	Y
4	11
6	4
7	6
8	5
12	8

identifying best fit line for these values

By using

$$y = \beta_0 + \beta_1 x \quad [\text{Linear line}]$$

not in $y = \log x, e^x, ax^n, \beta x^{-n}$

For Average Fit line

For $2+5x$ Equation

For $5+4x$ Equation

X	Y	y-Pred avg	Error	Error ²
4	11	6.8	4.2	17.64
6	4	6.8	-2.8	7.84
7	6	6.8	-0.8	0.64
8	5	6.8	-1.8	3.24
12	8	6.8	1.2	1.44
		30.8		

Random Value
For $2+5x$ (Eq)

SSE

X	Y	y-Pred	Error	Error ²
4	11	22	-11	121
6	4	32	-28	784
7	6	37	-31	961
8	5	42	-37	1369
12	8	62	-54	2916
		6151		

For $5+4x$ Equation

SSE

X	Y	y-Pred	Error	Error ²
4	11	21	-10	100
6	4	29	-25	625
7	6	33	-27	729
8	5	37	-32	1024
12	8	53	-45	2025
		4503		

SSE

$$\text{* Avg} = \frac{11+4+6+5+8}{5} = 6.8$$

$$\begin{aligned} \text{* Error} &\Rightarrow 11 - 6.8 = 4.2 \\ (\text{Actual} - \text{Predicted}) &= 4 - 6.8 = -2.8 \\ y &\quad y_{\text{pred}} = 6 - 6.8 = -0.8 \\ &\quad 5 - 6.8 = -1.8 \\ &\quad 8 - 6.8 = 1.2 \end{aligned}$$

$$\begin{aligned} \text{* Error}^2 &= (4.2)^2 = 17.64 \\ (2.8)^2 &= 7.84 \\ (0.8)^2 &= 0.64 \\ (1.8)^2 &= 3.24 \\ (1.2)^2 &= 1.44 \end{aligned}$$

$$\begin{aligned} \text{* Some of square of Error} &= 17.64 + 7.84 + 0.64 + 3.24 + 1.44 \\ &= 30.8 \end{aligned}$$

For, $2+5x$

SSE

$$\begin{aligned} \Rightarrow y_{\text{pred}} &\Rightarrow 2+5(4) = 22 \\ 2+5(6) &= 32 \\ 2+5(7) &= 37 \\ 2+5(8) &= 42 \\ 2+12(12) &= 62 \end{aligned}$$

⇒ Error

$$(y - \hat{y}) = 11 - 22 = -11$$

$$4 - 32 = -28$$

$$6 - 37 = -31$$

$$5 - 42 = -37$$

$$8 - 62 = -54$$

$$\begin{aligned} \Rightarrow \text{Error}^2 &= (-11)^2 = 121 \\ (-28)^2 &= 784 \\ (-31)^2 &= 961 \\ (-37)^2 &= 1369 \\ (-54)^2 &= 2916 \end{aligned}$$

$$= 6151$$

$$\begin{aligned} y_{\text{pred}} &\Rightarrow 5+4(4) = 29 \\ 5+4(6) &= 33 \\ 5+4(7) &= 37 \\ 5+4(8) &= 37 \\ 5+4(12) &= 53 \end{aligned}$$

$$\begin{aligned} \text{* Error} &= 21 - 29 \\ y - y_{\text{pred}} &= 4 - 29 \\ &= 6 - 33 \\ &= 5 - 37 \\ &= 8 - 53 \end{aligned}$$

$$\begin{aligned} \text{* Error}^2 &= (10)^2 = 100 \\ (25)^2 &= 625 \\ (27)^2 &= 729 \\ (32)^2 &= 1024 \\ (45)^2 &= 2025 \end{aligned}$$

$$\begin{aligned} \text{SSE} &= 100 + 625 + 729 + 1024 + 2025 \\ &= 2025 \end{aligned}$$

$$\begin{aligned} \Rightarrow \text{SSE} &= 121 + 784 + 961 + 1369 + 2916 \\ &= 6151 \end{aligned}$$

For Regression line

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

X	Y	$x - \bar{x}_{\text{mean}}$	$y - \bar{y}$	$\bar{x} * \bar{y}$	$(x - \bar{x})^2$
4	11	-3.4	4.2	-14.28	11.56
6	4	-1.4	-2.8	3.92	1.96
7	6	-0.4	-0.8	0.32	0.16
8	5	0.6	-1.8	-1.08	0.36
12	8	4.6	1.2	5.52	21.16
				$\Sigma E = -5.6$	$E(35.2)$

X Avg: 7.4

Y Avg: 6.8

* Avg. of X = $\frac{4+6+7+8+12}{5} = 7.4$

* Avg. of Y = $\frac{11+4+6+5+8}{5} = 6.8$

* $x - \bar{x}_{(\text{mean})} = 4 - 7.4 = -3.4$
 $6 - 7.4 = -1.4$
 $7 - 7.4 = -0.4$
 $8 - 7.4 = 0.6$
 $12 - 7.4 = 4.6$

* $y - \bar{y}_{(\text{mean})} = 11 - 6.8 = 4.2$
 $4 - 6.8 = -2.8$
 $6 - 6.8 = -0.8$
 $5 - 6.8 = -1.8$
 $8 - 6.8 = 1.2$

* $\bar{x}_{\text{mean}} * \bar{y}_{\text{mean}} =$
 $-3.4 * 4.2 = -14.28$
 $-1.4 * -2.8 = 3.92$
 $-0.4 * -0.8 = 0.32$
 $+0.6 * -1.8 = -1.08$
 $4.6 * 1.2 = 5.52$

$$\therefore \beta_0 = \bar{y} - \beta_1 (\bar{x})$$

$$\beta_1 = \frac{\Sigma (x - \bar{x})(y - \bar{y})}{\Sigma (x - \bar{x})^2}$$

* $\bar{x}_{\text{mean}} * \bar{y}_{\text{mean}}$

$$= -14.28 + 3.92 + 0.32 - 1.08 + 5.52 = -5.6$$

* $(x - \bar{x})^2$

$$\begin{aligned} (3.4)^2 &= 11.56 \\ (-1.4)^2 &= 1.96 \\ (-0.4)^2 &= 0.16 \\ (0.6)^2 &= 0.36 \\ (4.6)^2 &= 21.16 \end{aligned}$$

* $\Sigma (x - \bar{x})^2$

$$= 11.56 + 1.96 + 0.16 + 0.36 + 21.16 = 35.2$$

$$\beta_1 = \frac{\Sigma (x - \bar{x})(y - \bar{y})}{\Sigma (x - \bar{x})^2} \Rightarrow \frac{-5.6}{35.2} = 0.15909$$

Slope

$$\beta_0 = \bar{y} - \beta_1 (\bar{x}) = 6.8 - (0.159)(7.4)$$

$6.8 - (0.159) = 7.977$
Now, The equation is

Intercept

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

0.15909 [n]

$$y = 7.977$$

Here

$$y = 7.977273 - 0.159[x]$$

Equation.

X	Y	y _{pred}	Err	Error ²
4	11	7.34	3.66	13.39
6	4	7.02	-3.02	9.12
7	6	6.86	-0.86	0.73
8	5	6.70	-1.7	2.89
12	8	6.06	1.94	3.76
			29.89	SSE

$$\text{Equation: } 7.977 - 0.159[x]$$

$$7.977 - 0.159(4) = 7.34$$

$$7.977 - 0.159(6) = 7.02$$

$$7.977 - 0.159(7) = 6.86$$

$$7.977 - 0.159(8) = 6.70$$

$$7.977 - 0.159(12) = 6.06$$

Comparing with
Avg. line and
Other Equations

$$y = 7.977 - 0.159[x]$$

Best fit line.

[SSE] min

Error

$$11 - 7.34 = 3.66$$

$$4 - 7.02 = -3.02$$

$$6 - 6.86 = -0.86$$

$$5 - 6.70 = -1.7$$

$$8 - 6.06 = 1.94$$

Error²

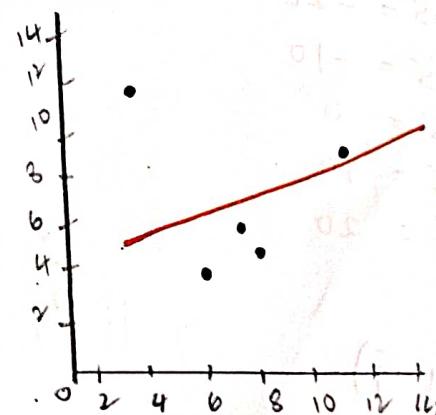
$$(3.66)^2 = 13.39$$

$$(-3.02)^2 = 9.12$$

$$(-0.86)^2 = 0.73$$

$$(-1.7)^2 = 2.89$$

$$(1.94)^2 = 3.76$$



$$\begin{aligned} \text{SSE} &= 13.39 + 9.12 + 0.73 + 2.89 + 3.76 \\ &= 29.89 \end{aligned}$$

For, Average \rightarrow

For, $10x+5$

X	Y	\bar{y}	Error	c_{error}
1	15	35	+20	400
2	25	35	+10	100
3	35	35	0	0
4	45	35	-10	100
5	55	35	-20	400
				$SST = 1000$

EX [For Intercept]

$$\therefore R^2 = 1 - \frac{0_{\text{reg}}}{1000_{\text{tot}}}$$

$$R^2 = 1$$

Maximum value

X	Y	$x-\bar{x}$	$y-\bar{y}$	$(x-\bar{x})(y-\bar{y})$	$(x-\bar{x})^2$	\bar{y}	Error	Error
1	15	-2	-20	40	-4	15	0	0
2	25	-1	-10	10	-1	25	0	0
3	35	0	0	0	0	35	0	0
4	45	1	10	10	1	45	0	0
5	55	2	20	40	4	55	0	0
				100	0			
						$SSE = 0$		

$$\bar{x} = 3 \quad \bar{y} = 35$$

$$X \text{ Avg} = \frac{1+2+3+4+5}{5} = 3$$

$$Y \text{ Avg} = \frac{15+25+35+45+55}{5} = 35$$

$$* x - \bar{x} = 1 - 3 = -2$$

$$2 - 3 = -1$$

$$3 - 3 = 0$$

$$4 - 3 = 1$$

$$5 - 3 = 2$$

$$* y - \bar{y} = 15 - 35 = -20$$

$$25 - 35 = -10$$

$$35 - 35 = 0$$

$$45 - 35 = 10$$

$$55 - 35 = 20$$

$$* (x - \bar{x})(y - \bar{y})$$

$$= -2 \times (-20) = 40$$

$$-1 \times (-10) = 10$$

$$1 \times 10 = 10$$

$$2 \times 20 = 40$$

$$* \varepsilon(x - \bar{x})(y - \bar{y})$$

$$= 40 + 10 + 0 + 10 + 40 = 100$$

$$* (x - \bar{x})^2 \# \text{ For squaring don't consider } (+, -)$$

$$= (-2)^2 = +4$$

$$(-1)^2 = +1$$

$$(1)^2 = 1$$

$$(2)^2 = 4$$

$$* \varepsilon(x - \bar{x})^2 = -4 - 1 + 1 + 4 = 10$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

$$\beta_1 = \frac{\varepsilon(x - \bar{x})(y - \bar{y})}{\varepsilon(x - \bar{x})^2}$$

$$\# \beta_1 = \frac{100}{10} = 10$$

slope

$$\beta_0 = 35 - 10 \cdot 3 = 5$$

y intercept

Equation is

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

$$y = 5 + 10x$$

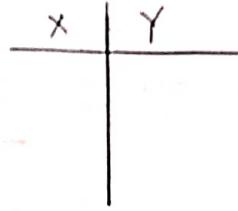
✓
8/4/22
3:30 am

8/4/22
3:00 PM

* BIVARIATE STATISTICS *

⇒ Simple linear regression :-

it is **2 Sample "T" test**



H_0 : Avg. line \geq Regression line
 H_1 : Regression line $>$ Avg. line.

Q: How can this be proved statistically?

A: Where

(SSE)_{Regression line}

< (SSE)_{Average line}

Ex:- 120

30 <

⇒ Multiple linear regression :-

x_1	x_2	x_3	x_4	y

Where we have multiple "x" values and only single "y". We use **"ANOVA"** Statistical Test ✓ applied when output variable is "Continuous".

Q: What is the hypothesis Test apply for Regression?

A: **Anova Test**

(> 2 samples) Apply for More than 2 samples

A:

if we apply Only for 2 samples

it is **"2 sample T Test"**

* **"Correlation"** is high

gives "good" results

→ **Linear Regression**

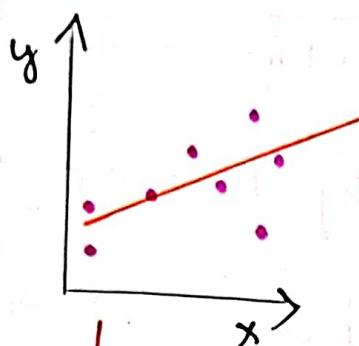
" Relation between two variables
Modulus of "r" any value
 $|r| = \pm 1$

$|r| \geq 0.8$ (strong) correlation
 $0.5 < |r| < 0.8$ (moderate)
 $|r| < 0.5$ (weak)

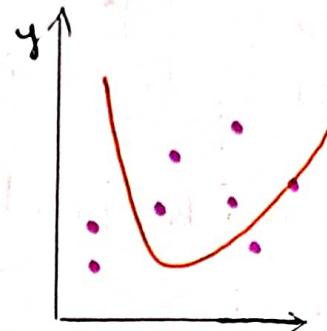
The value of "One Variable", is a function of "Other Variable"

The value of y , is a function of x :

$$y = f(x)$$



$$\therefore y = mx + c$$



$$y = ax^r + bx + c$$

The value of dependent variable, is function of independent variable

$$\Rightarrow E(y) = \beta_0 + \beta_1 x$$

"slope"
"y" intercept
Estimated value

slope β_1 is 0

β_0

slope β_1 is +

β_0

slope β_1 is -

β_0

$$E(y) = \beta_0 + 0(x)$$

$$E(y) = \beta_0 + \beta_1 x$$

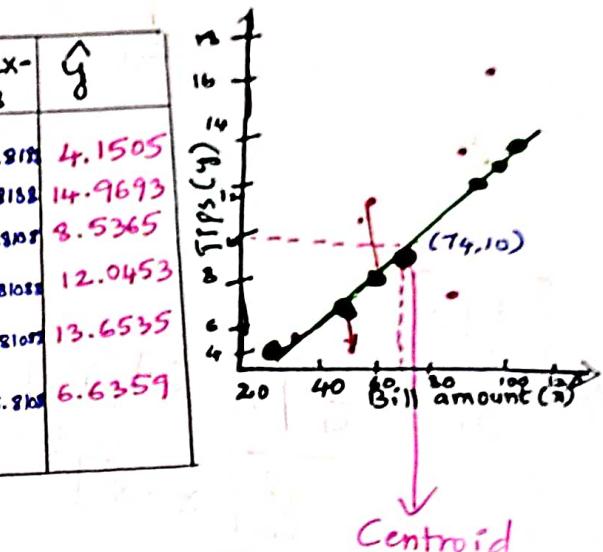
$$E(y) = \beta_0 - \beta_1 x$$

* Calculations

X	Y	$X - \bar{X}$	$Y - \bar{Y}$	$(X - \bar{X})(Y - \bar{Y})$	$(X - \bar{X})^2$	$y = 0.1462x - 0.8188$	\hat{y}
34	5	-40	-5	200	1600	$0.1462[34] - 0.8188$	4.1505
108	17	34	7	238	1156	$0.1462[108] - 0.8188$	14.9693
64	11	-10	1	-10	100	$0.1462[64] - 0.8188$	3.5365
88	8	14	-2	-28	196	$0.1462[88] - 0.8188$	12.0453
99	14	25	4	100	625	$0.1462[99] - 0.8188$	13.6535
51	5	-23	-5	115	529	$0.1462[51] - 0.8188$	6.6359
$\bar{X} = 74$		$\bar{Y} = 10$		$\sum E = 615$	$\sum F = 4206$		

$$y = 0.1462[x] - 0.8188$$

Equation



The line should pass through centroid.

$$\text{Avg. of } X := \frac{34 + 108 + 64 + 88 + 99 + 51}{6} = \bar{X} = 74$$

$$\text{Avg of } Y := \frac{5 + 17 + 11 + 8 + 14 + 5}{6} = \bar{Y} = 10$$

How This Equation is formed?

Descriptive Statistics / Centroid

$$\bar{X}, \bar{Y} = 74, 10$$

$$\text{Variance } (X) = \frac{\sum (x - \bar{x})^2}{(n-1)} \rightarrow s^2$$

* Line $[\hat{y} = b_0 + b_1 x]$ min

Variance(X) can be written as.

$$b_1 = \frac{\sum [x_i - \bar{x}][y_i - \bar{y}]}{\sum [x_i - \bar{x}]^2}$$

$$\frac{\sum (x - \bar{x})(y - \bar{y})}{(n-1)}$$

covariance (X, Y) written as

$$\frac{\sum (x - \bar{x})(y - \bar{y})}{(n-1)}$$

Substituting & Equating

Slope =

$$\frac{\sum [x - \bar{x}][y - \bar{y}]}{\sum [x - \bar{x}]^2}$$

$$\text{cov}(X, Y) = \frac{\sum [x - \bar{x}][y - \bar{y}]}{n-1}$$

(S.S.E)

* Sum of Squares of Error = $\sum [y - \hat{y}]^2$

$$= \sum [y - (\beta_0 + \beta_1 x)]^2$$

$$= \sum [y - \bar{y} - \beta_1 x]^2$$

↓ should be minimum

Slope $b_1 = \frac{\sum [(x - \bar{x})(y - \bar{y})]}{\sum [x - \bar{x}]^2}$

$$b_1 = \frac{615}{4206}$$

$$b_1 = 0.1462$$

y intercept $\Rightarrow b_0 = \bar{y} - b_1 \bar{x}$

$$b_0 = 10 - 0.1462 [74]$$

$$b_0 = -0.8188$$

Final Equation $\Rightarrow y = b_0 + b_1 x$

$$y = -0.8188 + 0.1462 x$$

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

(or)

$$y = 0.1462 x + (-0.8188)$$

*** "Accuracy" called as
"Coefficient of Determination"

$$R^2 = 1 - \frac{SSE_{\text{reg}}}{SST_{\text{Avg}}}$$

"Determining The variation of "y"
with respect To variation of "x"
→ As "x" value changes, how
"y" value changes.

(or)
"y" variable Explained by
"x" variable.

Ques: What is Max. value of R^2

Ans: $R^2 = 1$ [Max. value]

$$(SSE)_{\text{Reg}} = 0$$

$$R^2 = -\infty$$

$\therefore R^2$ (W.R.T mod)

**** Range of R^2

$$[\alpha, 1]$$

Infinite

"Regression Squared Error"

\hat{y}	Error $y - \hat{y}$	Square Error $(y - \hat{y})^2$
4.1505	(5.4 - 4.1505)	(0.8495) ²
14.9693	2.0307	4.1237
8.5365	2.4635	6.0688
12.0453	-4.0453	16.3645
13.6535	0.3465	0.1201
6.6359	-1.6359	2.6762

$$SSE \Rightarrow \sum = 30.075$$

- * Avg. line = 120
 - * Regression line = 30.075
- difference between these are $\frac{30}{120} \Rightarrow 75\%$.

Coefficient of Determination

(Or)

$$R^2 = 1 - \frac{SSE(\text{Reg})}{SST(\text{Avg})}$$

$\sum [y - \hat{y}]^2$
SSE [Sum of squares of Error (Regression line)]
SST [Sum of squares of Error (Average line)]

$$R^2 = 1 - \frac{\sum [y - \hat{y}]^2_{\text{reg}}}{\sum [y - \bar{y}]^2_{\text{Avg}}}$$

$$R^2 = 1 - \frac{30}{120^4}$$

$$R^2 = \frac{4-1}{4} = \frac{3}{4} \Rightarrow 0.75 \Rightarrow 75\% \text{ accuracy}$$

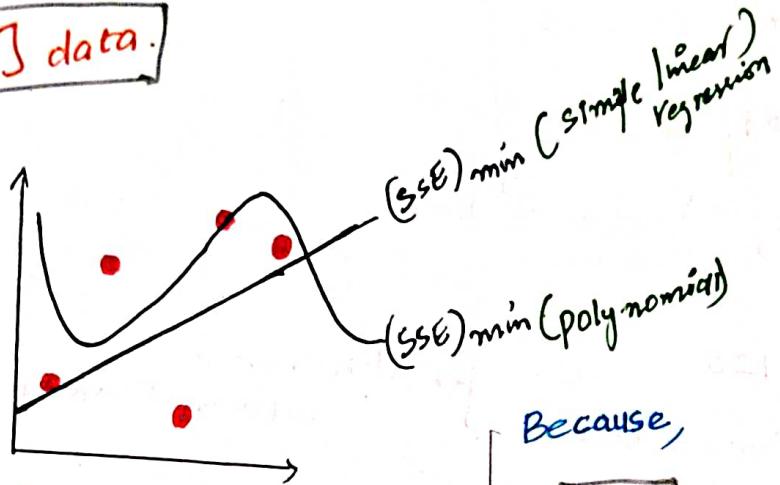
enf
8/4/22
4:30 PM

Dr. - 9/4/22
3:00 PM

Evaluation Metrics in Regression

After Model Fitting, we would like to assess The Performance of model by Comparing model Predictions to actual

[True] data:



1. **MAE** [Mean Absolute Error]

bad

\downarrow
differentiation
 \downarrow
squaring of units

$$\frac{\sum |y - \hat{y}|}{n}$$
2. **MSE** [Mean Squared Error]

$$\frac{\sum [y - \hat{y}]^2}{n}$$

Because,

$$\text{MAE} \rightarrow \frac{d[|x|]}{dx}$$

$$\frac{d[|x|]}{dx} \underset{\text{at } x=0}{\rightarrow} \frac{x}{|x|} = \frac{0}{0} \Rightarrow \infty$$

↓
Infinite

$$\text{MSE} \rightarrow \frac{dx^2}{dx}$$

$$\frac{dx^2}{dx} \underset{\text{at } x=0}{=} 2x$$

Example milage of a car

y	\hat{y}	y^2
18.2	18.9	0.7

Here
18.2 Kmph
18.9 Kmph

$$\text{Error} = (0.7)^2 = 0.49 \frac{\text{Km}^2}{\text{h}^2}$$

If we square it

Squaring
Solves the
problem of
units

$$\frac{d[\]}{dx} \text{ at } x=0$$

x is only notation it's not $\frac{dx}{dy}$

Ques: What are Evaluation Metrics used for Regression?

Ans: We have 1. Mean Absolute Error (MAE)
and 2. Mean square error [MSE].

Ques \Rightarrow Mean absolute Error (MAE) Where, we are going to use?

Ans: MAE is "Absolute Values of Mean" is called mean absolute Error.

But, problem with $\frac{d}{dx}$ of $(\text{Mod})|x|$ is infinite

is $\frac{x}{|x|}$ and at $x=0$, The answer is "MAE"

So, we are not going to consider

Ques \Rightarrow Mean Square Error (MSE) with calculating of $\frac{d}{dx}$ of x^n

Ans: Problem with MSE is

we will get some value of $2x$. But, by squaring the Error. Units are also squaring up

it gives wrong prediction.

Consider "MSE".

And Solution For This is :-

$$\text{RMSE} = \sqrt{\frac{\sum (y - \hat{y})^2}{N}}$$

Same as in Variance

We take Standard deviation

"Root Mean Square Error"

Q4c :- What is "Root Mean Square Error" RMSE?

A :- "RMSE" represents standard deviation of residuals [Errors] difference between model predictions and True values (Training data)

* RMSE can be easily interpreted compared To MSE because RMSE Units match units of Output.

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n [y - \hat{y}]^2}$$

$$(or) \quad \sqrt{\frac{\sum [y - \hat{y}]^2}{n}}$$

Example

Squaring of Errors

Converting negative values to positive values by adding 11

X	\hat{Y}	Error	MAE	MSE
1	1.4	0.4	0.4	0.16
2	2.2	0.2	0.2	0.04
3	3.8	0.8	0.8	0.64
4	4.1	0.1	0.1	0.01
5	5.6	0.6	0.6	0.36

$$\frac{\sum E}{n} = 0.42 \quad \frac{\sum E^2}{n} = 0.242$$

$$\text{Root Mean Square Error} = \sqrt{0.242}$$

⇒ Mean percentage Error [MPE]

$$MPE = \frac{100\%}{n} \sum_{i=1}^n (y_i - \hat{y}_i) / y_i$$

In this we calculate Percentage of Error. But we don't consider "Positive" or "Negative" value, without absolute operation.

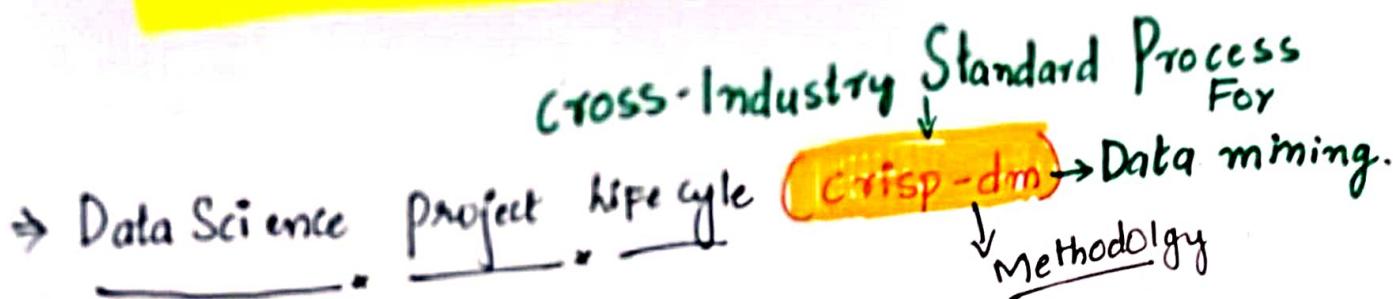
⇒ Mean absolute Percentage Error [MAPE]

In This "MAPE" we convert Every value To positive by applying "absolute" to the Equation.

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

9/4/22
5:10pm

* "Simple Linear Regression" *



1. Business problem Understanding

2. Data Understanding

3. Data Preprocessing

4. Modelling

5. Evaluation

6. Presentation

Data collection [Data Eng.]

Data variables (Research
- domain Expert)

Dataset understanding

EDA

Data cleaning

Data Wrangling

Train/test split

Apply Different Algorithms

Accuracy checking

(or) If not go back and
start again

Visualization

(Tableau, Power BI)

⇒ Basic Libraries for all Machine Learning

Import numpy as np

Import Pandas as pd

Import matplotlib.pyplot as plt
%matplotlib inline

Import Seaborn as sns

Step 1: Business Problem Understanding

- * Is there a relationship between total advertising Spend and Sales?
- * Our next ad campaign will have a total spend of \$200k, how many units do we expect to sell as a result of this?

Step 2.1 Data collection

df = pd.read_csv ("Advertising.csv")

df.head()

Out

TV	Radio	Newspaper	Sales
230.1	37.8	69.2	22.1
44.5	39.3	45.1	10.4
17.2	45.9	69.3	9.3
151.5	41.3	59.5	18.5
180.8	10.8	58.4	12.9

Step 2.2 Data Understanding

The sample data displays Sales (in thousands of units) for a particular product as a function of advertising budgets (in thousand of dollars) for TV, radio and Newspaper media.

Independent Variables

- TV : Advertising dollars spent on TV for a Single product in a Given Market (in thousands of dollars) → EX: 230.1×1000 already done scaling
- Radio : Advertising dollars Spent on Radio
- Newspaper : Advertising dollars Spent on news paper

Target Variable

- Sales : Sales of a Single product in a given market (in thousands of widgets) → EX: 22.1×1000 = 22100 products sold

Step - 2.3 Dataset Understanding

df.info()

Out	column
0	TV
1	Radio
2	Newspaper
3	Sales

	Non Null Count	Dtype
0	200 non null	float 64
1	200 non null	float 64
2	200 non null	float 64
3	200 non null	float 64

(200,4)

Que :- If some one was to spend a total of \$200, what would the expected sales be?

X :- We have simplified all Features

This quite entails bit by Combining into "total spend"

We add new column, total or 3 column creating new one

df[total-spend] = df["TV"] + df["radio"] + df["newspaper"]

df.head()

Out	TV	Radio	newspaper	Sales	total-Spend
0	230.1	37.8	69.2	22.1	337.1
1	44.5	39.3	45.1	10.4	128.9
2	17.2	45.9	69.3	9.3	132.4
3	151.5	41.3	58.5	18.5	251.3
4	180.8	10.8	58.4	12.9	250.0

We drop 3 columns, consider only two columns For Calculating

df.drop (columns = ["TV", "radio", "newspaper"], inplace=True)

df.head()

Out

Sales	Total-Spend
22.1	337.1
10.4	128.9
9.3	132.4
18.5	251.3
12.9	250.0

Exploratory Data Analysis [EDA]

Step - 3.1

On The basis of this data. How should you spend advertising money in future? These general questions might lead you to more specific questions:

1. Is there a relationship between ads and Sales?

2. How strong is that correlation?

3. Given ad spending, can Sales be predicted?

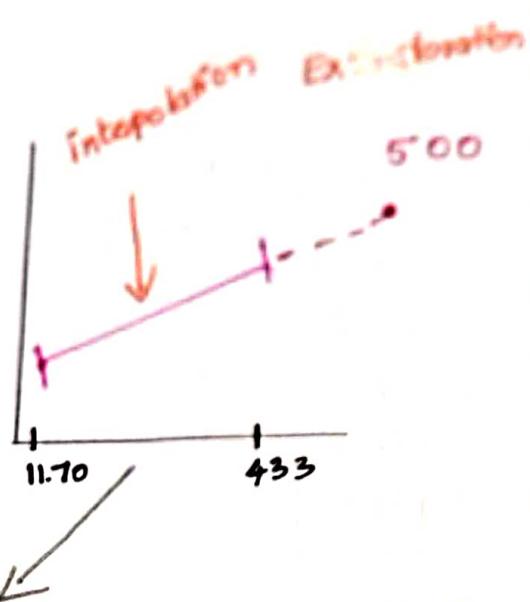
df.describe()

[Out]:

	Sales	total-spend
Count	200.000000	200.000000
Mean	14.022500	200.860500
Std	5.217457	92.985181
min	1.600000	11.700000
25%	10.375000	123.550000
50%	12.900000	207.350000
75%	17.400000	281.125000
max	27.000000	433.600000

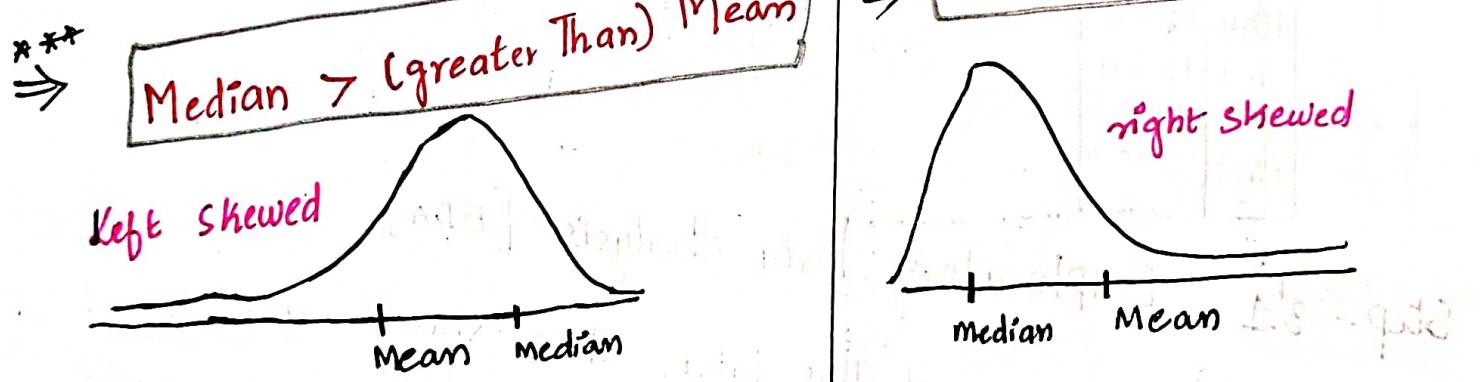
How close
The values
in dataset
To the mean

Mean and median
are close.
They normal
distribution



* Interpolation : Estimation of value with in two known values in sequence of values.

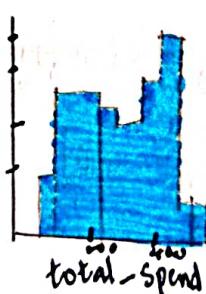
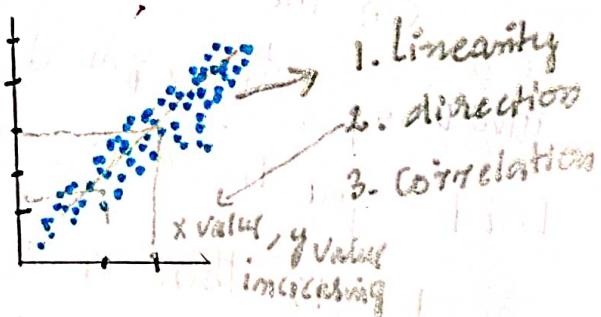
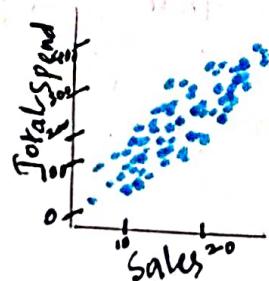
* Extrapolation : Prediction of value outside of known values it is called Extrapolation.



sns.pairplot(df)

plt.show()

[Out]:



df.corr()

out	Sales	total_spend
Sales	1.000000	0.867712
total-spend	0.867712	1.000000

Step 3.2 Data cleaning

df.isnull().sum()

out

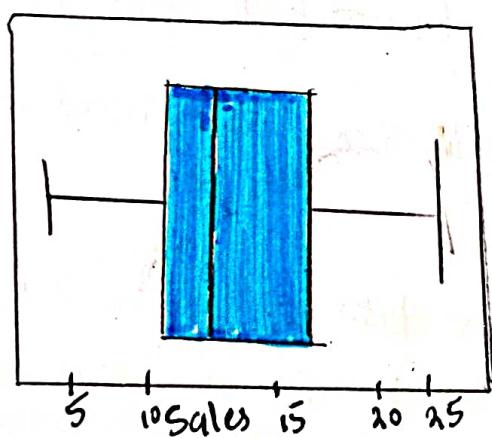
Sales	0
total-spend	0

Step 3.3 Data Wrangling

* Outliers

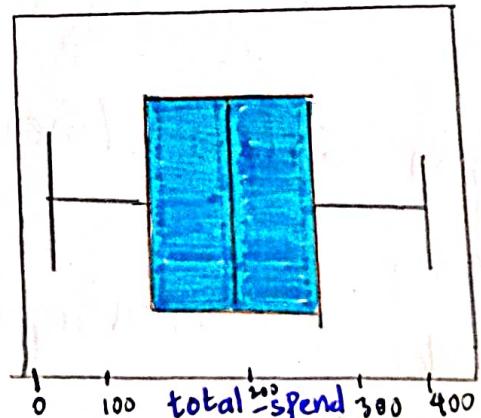
sns.boxplot(df.sales)

plt.show()



sns.boxplot(df.total_spend)

plt.show()



* Feature transformation

df.sales.skew()

[Out]: 0.407

df.total_spend.skew()

[Out]: 0.049

Step

3.

Train Test split

x = df.drop(columns = "sale")

y = df["sales"]

from sklearn.model_selection

x-train, x-test, y-train, y-test = train-test-split(x,y,
test-size = 0.3, random_state =)

Step-4 Modelling :- If doing with data preprocessing directly
working on modelling is "Base line model" (0%) raw model.

from sklearn.linear_model

Storing
model = LinearRegression()

model.fit(x-train, y-train)

[Out]: Linear Regression()

Print("Model Intercept : ", model.intercept_)

print ("Model Coefficient : ", model.coef_)

[Out]: Model Intercept : 4.33176
Model Coefficient : 0.0480

$$\begin{array}{c} \text{Attributes (x)} \\ \rightarrow y \text{ intercept} \\ B_0 + B_1 x \\ \downarrow \\ \text{coefficient value} \end{array}$$

$$\hat{y} = 4.33176 + 0.0480(x)$$

→ Predictions

spend = 200

predicted - sales = 4.33176 + 0.0480 * spend

Predicted - sales

[Out]: 13.93860

If we have multiple equations, we can write this

model.predict([[200]])

input should
given in
2 dimensional

[Out]: array([13.93860])

Predicting on x-train, x-test

	x	y	\hat{y}
x train	1	10	14
	2	15	12
	3	20	13
	4	30	14
x test	5	50	33

→ Train accuracy

→ Test accuracy

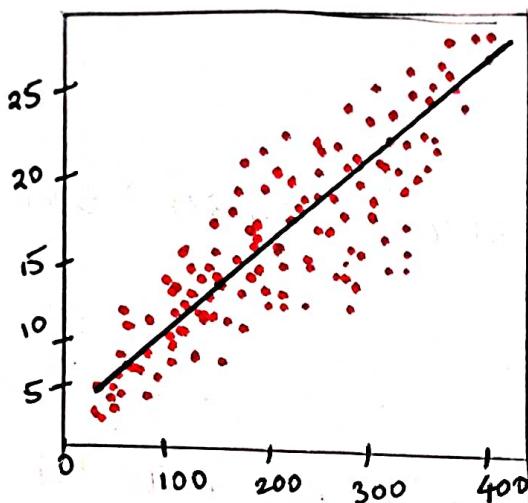
train_predictions = model.predict(x-train)

test_predictions = model.predict(x-test)

Only stores No output

→ plotting the Least Squares line

```
# plt.scatter(x-train, y-train, color = "red")
# train-predictions = model.predict(x-train)
# plt.plot(x-train, train-predictions, color = "black")
# plt.show()
```



Step 5

⇒ Evaluation metrics

```
# print("MAE For Test data:", mean_absolute_error(y-test, test-predictions))
# print("MAE For Train data:", mean_absolute_error(y-train, train-predictions))
```

Out: MAE For Train data : 1.97168
MAE For Test data : 1.90653

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

```
# print ("MSE for test data : ", mean_squared_error(y-test,  
test_predictions))
```

```
# print ("MSE for train data : ", mean_squared_error(y-train,  
train_predictions))
```

out: MSE for Test data : 6.82220

MSE for Train data : 6.40720

RMSE (Just write `np.sqrt`). \rightarrow `np.sqrt`

```
# print ("RMSE for test data", np.sqrt(mean_squared_error  
(y-test, test_predictions)))
```

```
# print ("RMSE for train data", np.sqrt(mean_squared_error  
(y-train, train_predictions)))
```

out: RMSE for Test data : 2.611935

RMSE for Train data : 2.53124

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

```
# print ("R2 for test data", r2_score(y-test, test-  
predictions))
```

```
# print ("R2 for train data", r2_score(y-train, train-  
predictions))
```

out: R2 for Test data : 0.74001

R2 for Train data : 0.76534

Another way for R²
model.score(x_train, y_train)

[Out]: 0.76534

model.score(x_test, y_test)

[Out]: 0.740017

cross-validation \Rightarrow K-Fold Cross Validation
from sklearn.model_selection import cross_val_score

scores = cross_val_score(model, x, y, cv=5)
Here K = 5

print(scores)

scores.mean()

[Out]: [0.74964192, 0.79455226, 0.76417134, 0.74872042,
0.65980565]

$\Rightarrow 0.74337 \rightarrow$ should be equal to test square
Then it is good model.

Linear Regression assumptions
L * I * N * E Errors

- L \rightarrow Linearity
- I \rightarrow Independent
- N \rightarrow Normality
- E \rightarrow Equal Variance

any
10/14/22
3:30 AM

→ CHECK LIST ←

Assumptions of linear regression

1. check whether model has overfitting (if) underfitting problem
2. Is test Accuracy = Cross Validation Score
3. check assumptions (if it is linear Regression)
check model meets business problem requirements
4. Finally, save the model and share to deployment team.
5. purely Assumption of The Researcher

Ques:- Is model has overfitting (if) underfitting problem?

Ans :- It's good model.

Ques :- Is test accuracy = Cross Validation Score?

Ans :- Applied K-fold cross validation and it is equal in test accuracy and CV score.

Ques :- Check Assumptions (if it is linear regression)

Ans :- Line Assumptions.

1. Linearity * of Errors

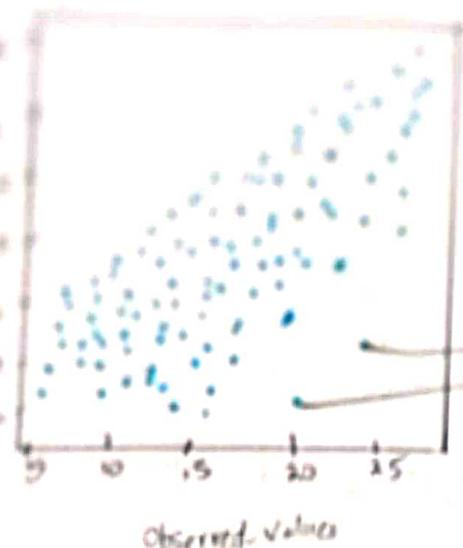
test_res = $y_{\text{test}} - \text{test_Predictions}$.
(↑ Saved in TM)

plt.scatter(y_test, test_res)

plt.xlabel("Observed-values")

plt.ylabel("Fitted-values")

plt.show()

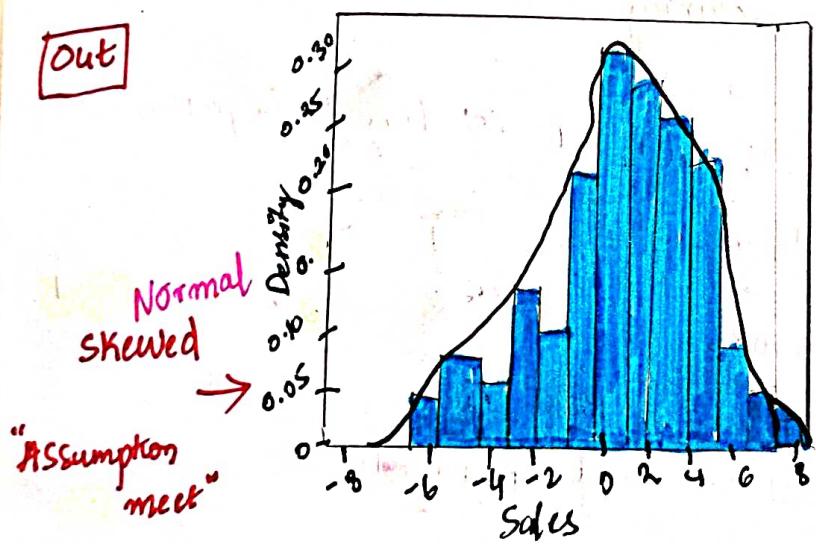


We can ignore
The Edge points
upto $\pm 5\%$ error

2. Normality of Errors

`sns.distplot(test_res, bins=15, kde=True)`
`plt.show()`

[Out]



Normal
Skewed
→
"Assumption
meet"

↓ p. of positive Errors
↓ p. of negative Errors
↓ ≈ 0.02
color = red

(Homoscedasticity)

3. Equal variance of Error

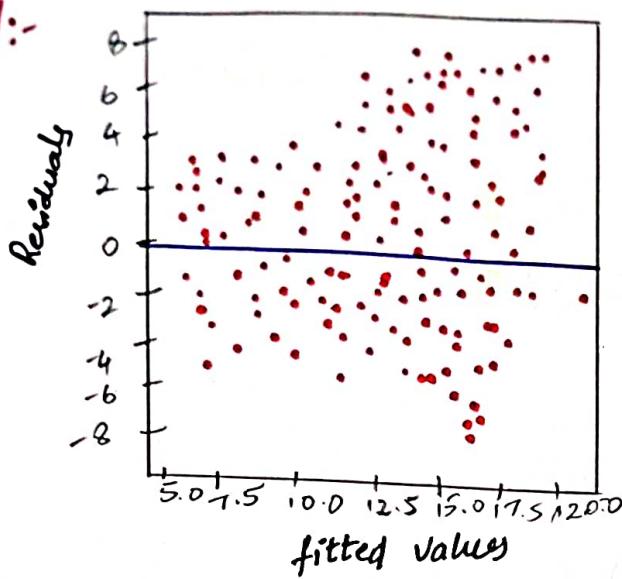
`plt.scatter(test_predictions, test_res, c="blue")`

`plt.axhline(y=0, color = "blue")`

↓
ax = axis → horizontal line

```
# plt.xlabel ("fitted values")
# plt.ylabel ("Residuals")
# plt.show()
```

[out]:-



For model

H_0 : Augline best fit line
 H_1 : regression line best line

Anova test
 $P \leq \alpha$
 $P \geq \beta$

4. Variable significance

```
# import statsmodels.formula.api as Smf
```

m = Smf.ols ("y~x", data=df).fit ()

m. summary ()

OLS Regression Results

[out]

Dep. Variable : y

R-squared : 0.753

Adj. R-squared : 0.752

F. static : 603.4

prob.(F.static) : 5.06×10^{-62}

p t	0.025	0.975
0.000	3.378	5.908
0.000	0.045	0.053

Confidence Interval
 $\text{Std.Error} = 6/5n$

Model : OLS

method : Least squares

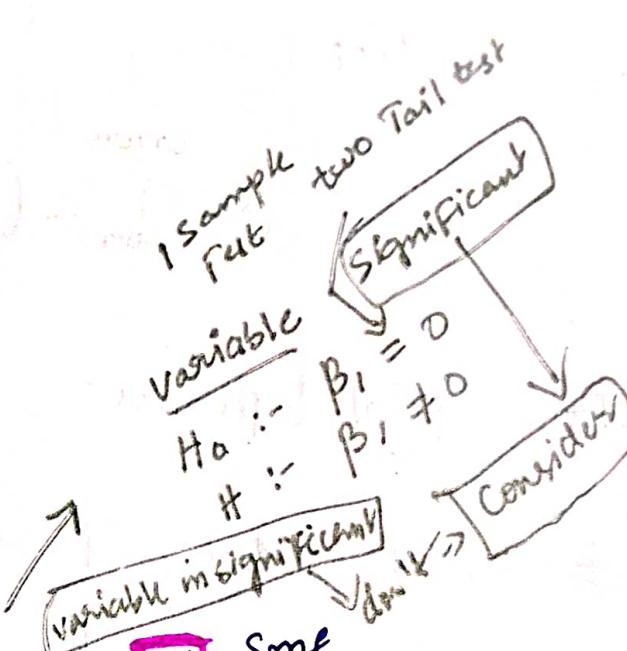
coffs std.error t

intercept

4.923 0.049 9.676
0.0487 0.002 24.564

1-sample T Test
Two-tail

t-test
Central limit
Theorem



When doing \hat{y} going to fail in Regression?
When sum of square Coefficient Average Line
When sum of square the sum of square of regression
↳ **Last Point** the sum of square of regression
 $(SST)_{avg} < (SSE)_{reg}$ it going to fail in this situation.

$$(SST)_{avg} < (SSE)_{reg}$$

Step G Final Inferences

```
# model.predict([[200]])
```

Out array ([13.9898])

→ Save a Model
from joblib import dump
"Sales-model.joblib")

```
# dump(model,
```

Out: ["sales-model.joblib"]

→ saves in working directory

→ Load a Model → from client side

```
# from joblib import load
```

```
# loaded-model = load ("sales-model.joblib")
```

```
# loaded-model.predict ([[200]])
```

[out] array ([13.989])

```
# loaded-model.predict ([[500]]) # check with multiple values.
```

[out]: array ([28.3488])

~~27/4/22~~
19/4/22
5.00pm

Date :- 11/04/22
10:30 pm

Multiple linear Regression

Multiple linear Regression :-

Examines Relationship between

(or) more than two

Variables

- * Each independent Variable has its own corresponding coefficient.

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + \dots + b_n * x_n$$

dependent variable

Independent Variables.

⇒ Example :

TV	Bad	NP	Sales

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

TV Radio news paper Sales

↓ ↓ ↓ ↓

TV News paper news paper

y-intercept : \hat{y}

Coefficient : $\beta_1, \beta_2, \beta_3$

Same Example : But different method

Step-1

Bussiness Problem Understanding.

- * What is The relation b/w each advertising channel [TV, Radio, Newspaper] and Sales? $y = f(x_1, x_2, x_3)$
- * Previously, We Explored is There a relationship between Total advertising Spend and Sales? as well as predicting the total sales for some value of Total spend.

Step: 2.1

Data collection

```
# df = pd.read_csv("Advertising.csv")
# df.head()
```

Out

TV	Radio	Newspaper	Sales
230.1	37.8	69.2	22.1
44.5	39.3	45.1	10.4
17.2	45.9	69.3	9.3
151.5	41.3	58.5	18.5
180.8	10.8	58.4	12.9

Step 2.2

Data Understanding

Same as Simple linear Regression Dataset.

Step 2.3

Data Set Understanding

```
# df.info()
```

Step-3.2

Exploratory Data Analysis [EDA]

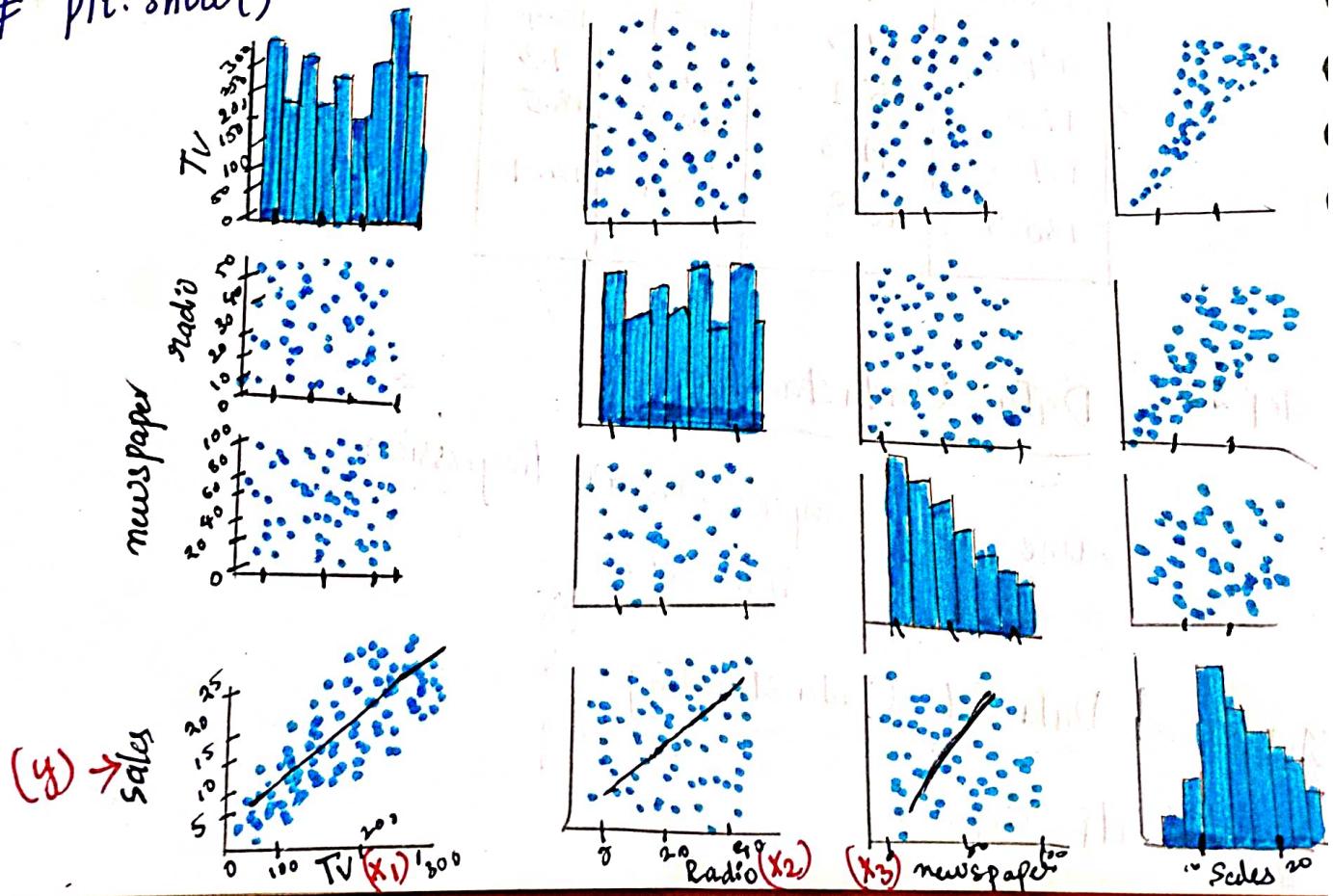
df.describe()

out:

	TV	Radio	newspaper	Sales
Count	200.000000	200.000000	200.000000	200.000000
Mean	147.042500	23.264000	30.554000	14.024500
Std	85.854236	14.846809	21.778621	5.217457
min	0.700000	0.000000	0.300000	1.600000
25%	74.375000	9.975000	12.750000	10.375000
50%	149.750000	22.900000	25.750000	12.900000
75%	218.825000	36.525000	45.100000	17.400000
max%	296.400000	49.600000	114.000000	27.000000

sns.pairplot(df)

plt.show()



→ By Observing The Scatter plot, we made an assumption of relation between y and $[x_1 + x_2 + x_3]$ is Linear (symmetrical Matrix)

df. corr(C) →

- i) y (vs) x_1
- ii) y (vs) x_2
- iii) y (vs) x_3

relation → strong → High Accuracy

x_1 (vs) x_2 → low
 x_1 (vs) x_3 ↓
 x_2 (vs) x_3 ↓

If strong, collinearity problem

	TV	Radio	News Paper	Sales
TV	1.000000	0.054809	0.056648	0.782224
Radio	0.054809	1.000000	0.354104	0.576223
News paper	0.056648	0.354104	1.000000	0.228299
Sales	0.782224	0.576223	0.228299	1.000000

dependent
Independent variable

- ⇒ The relation between "y" and "x" be high.
- ⇒ The higher The value - stronger The correlation
- ⇒ The relation between any two independent variables should be "Low"

**** if the correlation between any two (2) independent variables is strong. then it is called "collinearity problem".

Should be "weak"

Example

* independent Variable

* dependent Variable

[student 2]

[Teacher]

should be weak

should be strong

Surf 12/4/22
3:30 PM

Dt: 12/4/21
1:00PM

Step: 3.2

Data Cleaning

Same as Simple Linear

Step 3.3

Data Wrangling

not Same

Step 3.4

Train-Test Split

$x = df.\text{drop}[\text{columns} = \text{"Sales"}]$

$y = df[\text{"Sales"}]$

from sklearn.model_selection import train-test-split

$x_train, x_test, y_train, y_test = \text{train-test-split}(x, y, \text{test_size} = 0.3, \text{Random-state} = 29)$

Step-4

Modelling

Here, $y = \beta_0 + \beta_1 \times [TV] + \beta_2 \times [\text{radio}] + \beta_3 \times [\text{Newspaper}]$

from sklearn.linear_model import Linear Regression

model = Linear Regression()

model.fit(x_train, y_train)

: Linear Regression()

Out

```
# model.coef
```

```
[Out]: array([0.04422917, 0.18181641, 0.0075874])  
           $\beta_1$        $\beta_0$        $\beta_2$ 
```

```
# model.intercept
```

```
[Out]: 2.974097357  
         $\beta_0$ 
```

→ Predictions → Predicting On "X"

```
# train-predictions = model.predict(x-train)
```

```
# test-predictions = model.predict(x-test)
```

Step: 5

Evaluations → on "y"

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

```
# test_RMSE = np.sqrt(mean_squared_error(y-test, test-predictions))
```

```
# train_RMSE = np.sqrt(mean_squared_error(y-train, train-predictions))
```

```
# print(train_RMSE, test_RMSE)
```

```
[Out]: 1.64082, 1.7805
```

model.score (xtrain, y-train) (train Rⁿ)

[out]: 0.88817

model.score (x-test, y-test) (test Rⁿ)

[out]: 0.905258

⇒ checklist

Que: Is model has Underfitting (or) Overfitting problem?

Ans:- Good model

Que: Is Test Accuracy = Cross Validation Score

Ans:

from sklearn.model_selection import cross_val_score

Scores = cross_val_score (model, x, y, cv=5)

print (Scores)

Scores.mean()

[0.87865198, 0.9176312, 0.92933, 0.81443, 0.895]

[out]

0.887106349

compare with xtest,ytest accuracy.

3. Check Assumptions

* Linearity of errors

$\text{test_res} = y_{\text{test}} - \text{test_predictions}$

xtest

Residuals

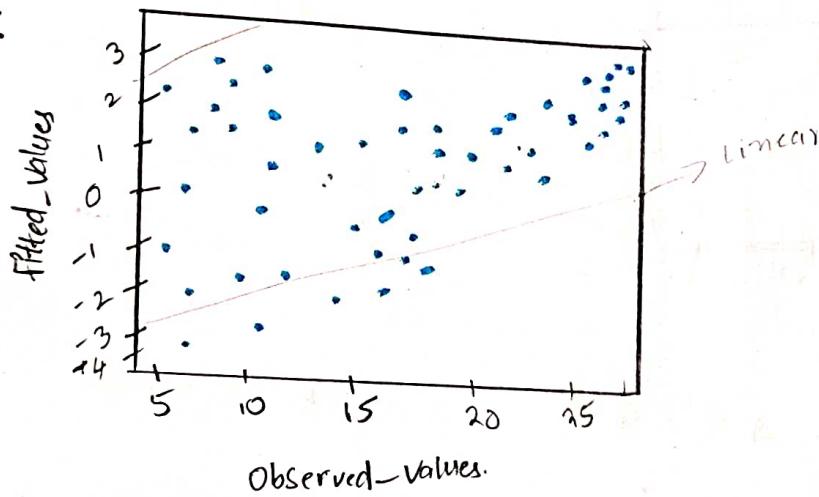
plt. Scatter (y-test, test-res)

plt. xlabel ("Observed-values")

plt. ylabel ("fitted-values")

plt. show()

out:



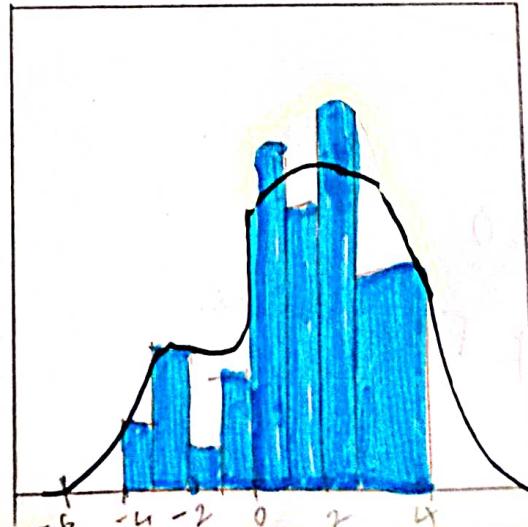
* Normality of errors

sns.distplot(test-res, bins=15, kde=True)

plt.show()

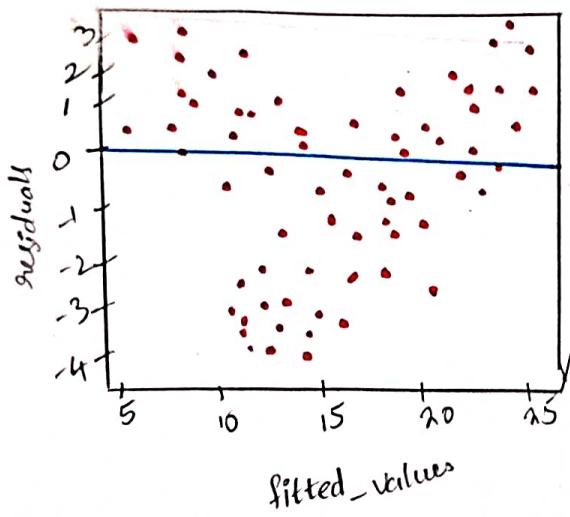
out:

Left skewed



3. Equal Variance of Errors (Homoscedasticity) X

```
# plt.scatter(test_predictions, test_res, c="r")  
# plt.axhline(y=0, color="blue")  
# plt.xlabel("fitted-values")  
# plt.ylabel("residuals")  
# plt.show()
```



* ASSumptions Failed, $R^2 = 86\%$. [we reject the model]

⇒ Every checklist should satisfy the condition. Then, Only we consider the model.

4. Variable Significance

$$H_0 : \beta_1 = 0$$

$$H_1 : \beta_1 \neq 0$$

⇒ Hypothesis Testing For Variables.

```
import statsmodels.formula.api smf
```

```
# model 2 = Smf. OLS("y ~ x", data=df).fit()
```

```
# model 2 . summary ()
```

[out]: OLS Regression Results.

For Model: Accept H_0 :
↓
Reject H_0 .

Dependent Variable : y

R-squared: 0.897

Model : OLS

Adj. R-squared: 0.896

method : least squares

F statistic: 570.3

Prob (F-statistic): 1.58×10^{-96}

	coeff	std err	t	p > t	[0.025]	0.975
Intercept	2.9389	0.312	9.422	0.000	2.324	3.554
β_1 x(0)	0.0458	0.001	32.809	0.000	0.043	0.049
β_2 x(1)	0.1885	0.009	21.893	0.000	0.472	0.206
β_3 x(2)	-0.0010	0.006	-0.177	0.860	-0.013	0.011

Variable 1

$H_0: \beta_1 = 0$

$H_1: \beta_1 \neq 0$

Variable wise

Variable 2

$H_0: \beta_2 = 0$

$H_1: \beta_2 \neq 0$

Variable 3

$H_0: \beta_3 = 0$

$H_1: \beta_3 \neq 0$

$P < 0.05 \rightarrow P_{\text{low}}$

Null, go \rightarrow Reject H_0

$P > 0.05$

⇒ Checking whether data has any influence value using influence index plots

1) * Influential index plots, For Making "P" value equal < 0.05

import statsmodels.api as sm

sm. graphics.influence_plot(model1)

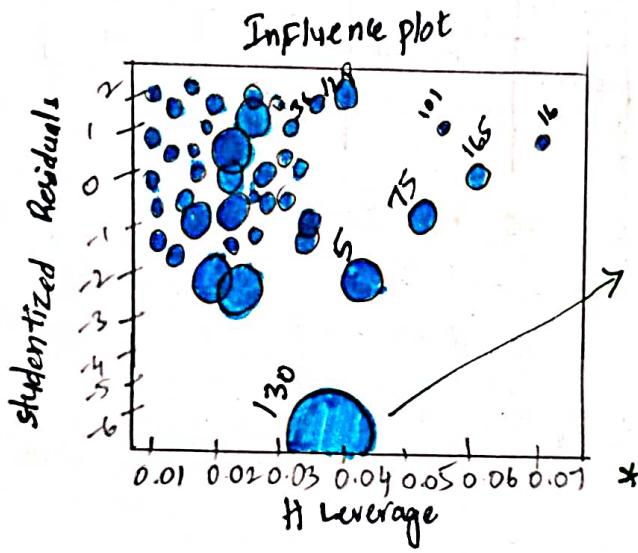
Why?

$$\beta_3 = 0,$$

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

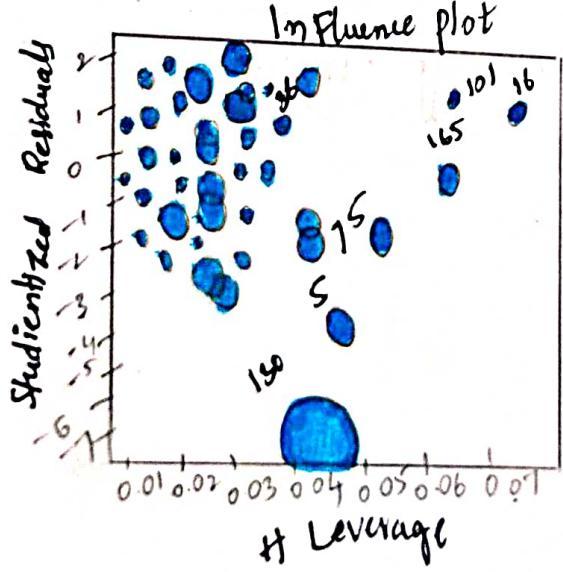
By dropping, we are losing $\frac{1}{3}$ of data = 33%.
In that case, we drop some records, that are influencing more.
 $Ex: - \frac{1}{200} = 0.05\%$.

But:



Index 130 is showing high influence. So, we exclude that entire row.

studentized residuals = $\frac{\text{Residuals}}{\text{standard deviation of residuals}}$



df_new = df.drop(df.index[[130]], axis=0)

df_new.

[Out]

	TV	Radio	News Paper	Sales
0	230.1	37.8	69.2	22.1
1	44.5	39.3	45.1	10.4
199	232.1	8.6	8.7	13.4

199 x 4 columns

Once again Rebuild model

model 2 = smf.ols (formula = "Sales ~ TV + radio + newspaper", data = df_new).fit()

model2. summary()

[Out]: OLS Regression MODEL

	coeff	std. err	t	P> t
Intercept	3.0931	0.290	10.654	0.000
TV X(0)	0.0448	0.001	34.425	0.000
Radio X(1)	0.1939	0.008	24.130	0.000
news paper X(2)	-0.043	0.005	-0.777	0.438

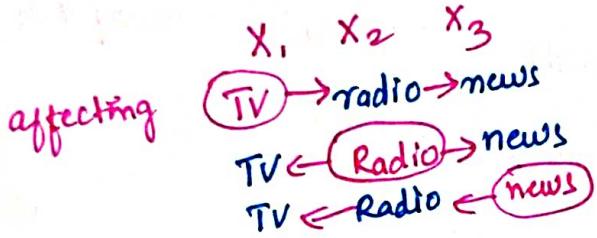
→ Reduced
From
0.86 to
0.43
By Reducing.

2) Variance inflation Factor [VIF]

Variance inflation Factor [VIF] measures ratio between the variance for a given regression coefficient with only that variable in the model versus the variance for given regression coefficient with all variables in the model.

1 independent variable

influence on other independent variable is called as ("VIF")



$$V.I.F = \frac{1}{1-R^2}$$

$r_{sq-TV} = \text{smf.ols}(\text{"TV ~ radio + newspaper"}, \text{data=df})$
• fit()

$r_{sq-TV}.\text{summary}()$

Out : OLS Regression Results

R-squared = 0.005

$$V.I.F = \frac{1}{1-0.005^2} \Rightarrow V.I.F = 1.0000$$

Calculating VIF's values of independent variables

$r_{sq-TV} = \text{smf.ols}(\text{"TV ~ radio + newspaper"}, \text{data=df}).$
• fit(), r.squared

$vif-TV = 1/(1 - r_{sq-TV})$

```
# rsq_radio = Smf.ols ("radio ~ TV + newspaper", data=df)  
    .fit().rsquared
```

```
# vif_radio = 1/(1-rsq_radio)
```

```
# rsq_newspaper = Smf.ols ("newspaper ~ radio + TV",  
    data=df).fit().rsquared
```

```
# vif_newspaper = 1/(1-rsq_newspaper)
```

Storing VIF values in a DataFrame

```
# d1 = { "variables": [ "TV", "radio", "newspaper" ],  
        "VIF": [ vif_TV, vif_radio, vif_newspaper ] }
```

```
# vif_frame = pd. Data Frame (d1)
```

```
# vif_frame
```

Out

	variable	VIF
0	TV	1.004611
1	Radio	1.144952
2	newspaper	1.145187

If The VIF model \rightarrow (greater) 4. For
any independent variable, drop
variable

2/4/22
12/04/22

4:40 pm

3) Dt: 16/9/22

⇒ AV plot [Added Variable plot)

partial differentiation instead of normal differentiation.

$$\begin{aligned} & (\text{SSE})_{\min} \\ & [\mathbb{E}(y - \hat{y})^2]_{\min} \\ & \frac{\partial}{\partial x} \text{ at } x=0 \end{aligned}$$

#

What is partial differentiation?

$\sum [y - \hat{y}]^2 \rightarrow$ it should be Minimum

$$\left\{ \left[y - [\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3] \right]^2 \right\}_{\text{minimum}}$$

$$\left\{ \left[y - \beta_0 - \beta_1 x_1 - \beta_2 x_2 + \beta_3 x_3 \right]^2 \right\}_{\text{minimum}}$$

In Simple linear Regression?

$$\frac{\partial \left[\mathbb{E}[y - \beta_0 - \beta_1 x_1]^2 \right]}{\partial x_1 \text{ at } x=0} \rightarrow \text{As we have only one variable.}$$

In Multiple linear Regression?

$$\frac{\partial}{\partial x_1} \left\{ \left[y - \beta_0 - \beta_1 x_1 - \beta_2 x_2 - \beta_3 x_3 \right]^2 \right\}$$

constant calculate constant
constant constant constant

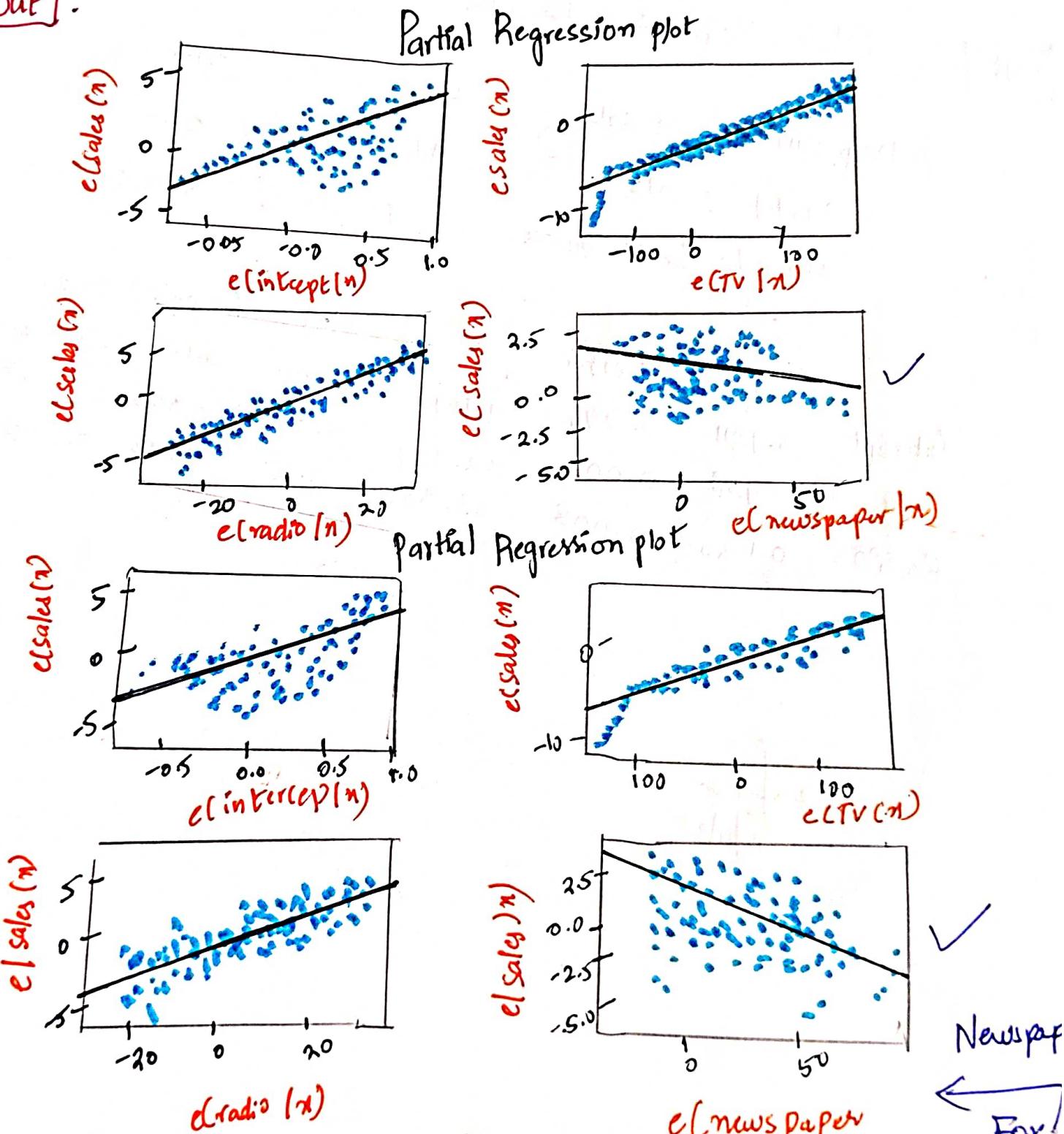
↓
when we calculate with $\frac{\partial \text{SSE}}{\partial x_i}$, $\frac{\partial}{\partial x_i}$.

When we calculate with one variable other variable will be constant.

→ it is apply Simple Linear Regression on x_1 and x_2
 Simple Linear Reg on x_3 . it is Applying individually
 S.L.R on Each and Every individual variable.

Sm. graphics . plot - partregress- grid (1m)
Partregress -

Out:



Added variable plot is not showing any significance For
 Newspaper ←

⇒ Final model including "TV" and "Radio" Only

final_model = smf.ols (formula = "Sales ~ TV + radio",
data = df). fit()

final_model.summary()

Out	OLS Regression Results		
Dep. Variable : sales	R-squared :	0.897	= 90.1.
Model : OLS	Adj. R-squared :	0.896	
Method: Least Squares	F-statistics :	859.6	
	t	p-value	
Intercept	2.9211	0.294	(0.975) 3.502
TV	0.0458	0.001	0.043
Radio	0.1880	0.008	0.172 0.204

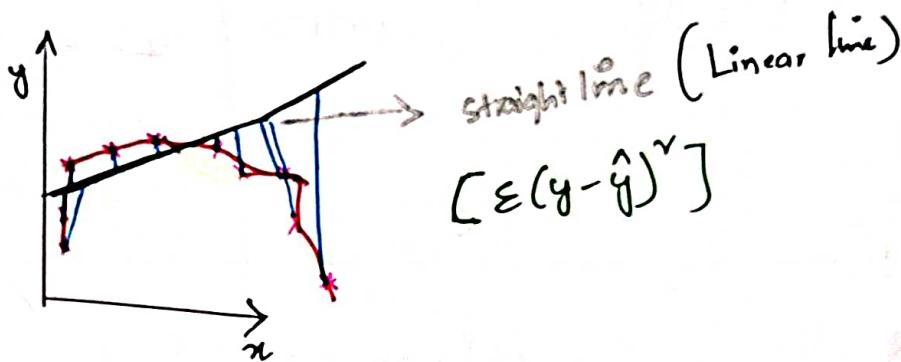
?
enf
16/11/22

Dr. Ishant 1:00 PM

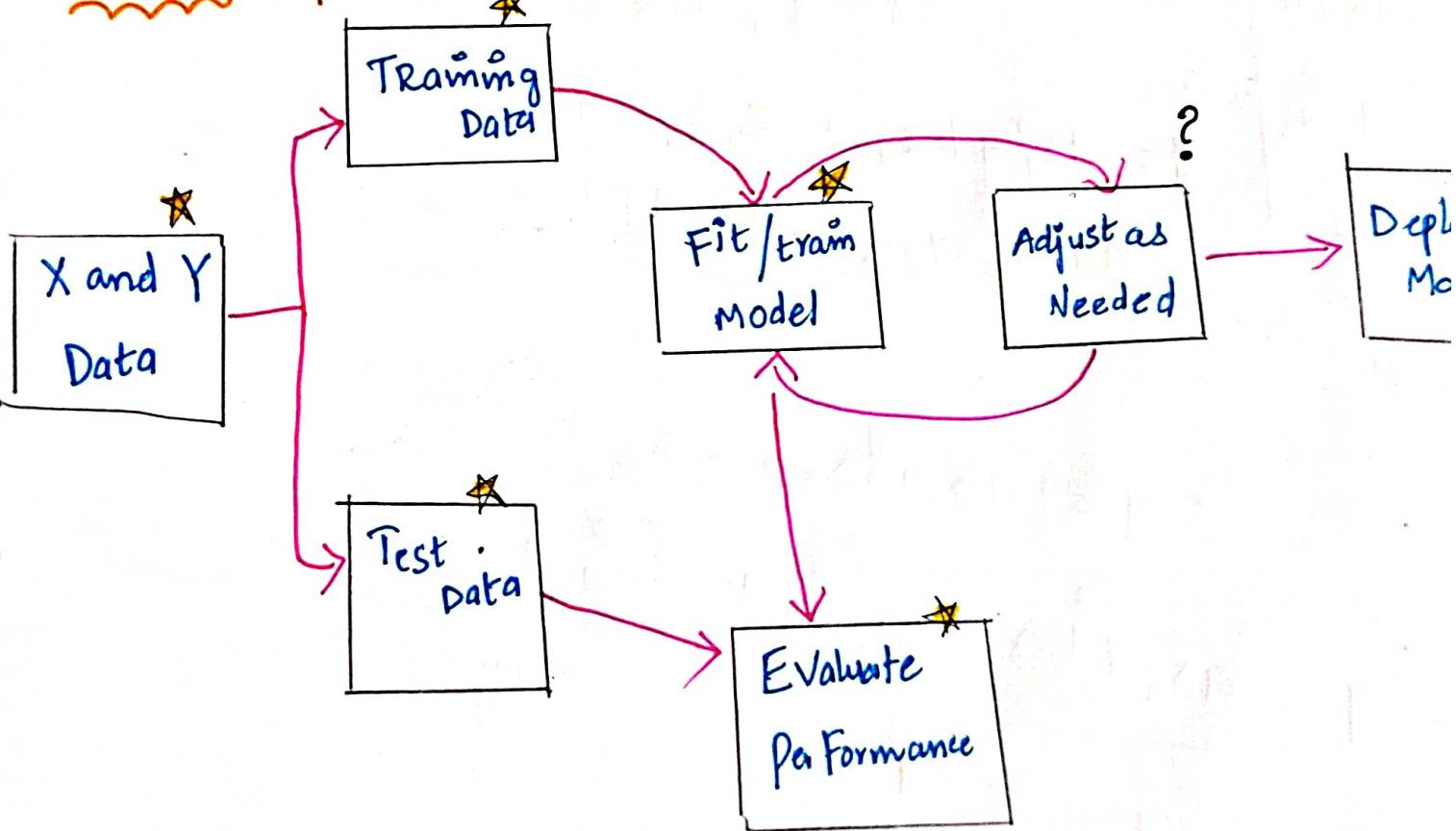
Polynomial Regression

Any Regression problem.

$$[SSE]_{\min} \Rightarrow [\varepsilon[y - \hat{y}]^2]_{\min}$$



Supervised Machine learning Process



S.L.R

M.L.R

1 Variable

$$y = ax + b$$

Linear

2 Variables

$$y = \frac{ax + b}{\text{1st variable}} + cx + d$$

$$y = ax_1 + bx_2 + c$$

it is came from joining
Two equations

$$= ax_1 + b + cx_2 + d$$

always

intercept
Have
One value Only

3 Variable

$$y = ax_1 + bx_2 + cx_3 + d$$

it is came from
Joining Three equations

$$y = ax_1 + b$$

$$y = +cx_2 + d$$

$$y = +ex_3 + f$$

Quadratic
1 Variable

$$y = ax^2 + bx + c$$

Here
we write
intercept value
as "Z"

$$y = ax^2 + bx_1 + dx_2 + ex_3 + Z$$

it is came from
Joining two
Equations

$$y = ax^2 + bx_1 + Z_1, c_1 +$$

$$+ dx_2^2 + ex_3 + Z_2, f_2$$

$$* y = ax_1 + bx_2 + cx_3 + Z$$

$$* y = ax_1 + bx_2 + cx_3 + dx_1^2 + ex_2^2 + fx_3^2 + Z$$

$$+ gx_1x_2 + hx_2x_3 + ix_3x_1$$

Combinations of variable [Iteration terms]

original data

squaring original Data

Combination (or)
Iteration terms

TV	Radio	Newspaper	Sales	TV^2	$Radio^2$	$Newspaper^2$	$TV \cdot Radio$	$Radio \cdot Newspaper$	$Newspaper^2 \cdot TV$
-	-	-	-	-	-	-	-	-	-
-	-	-	-	-	-	-	-	-	-
-	-	-	-	-	-	-	-	-	-

* Quadratic equation: $x_1^2 + x_2^2 = x_1 + x_2$?

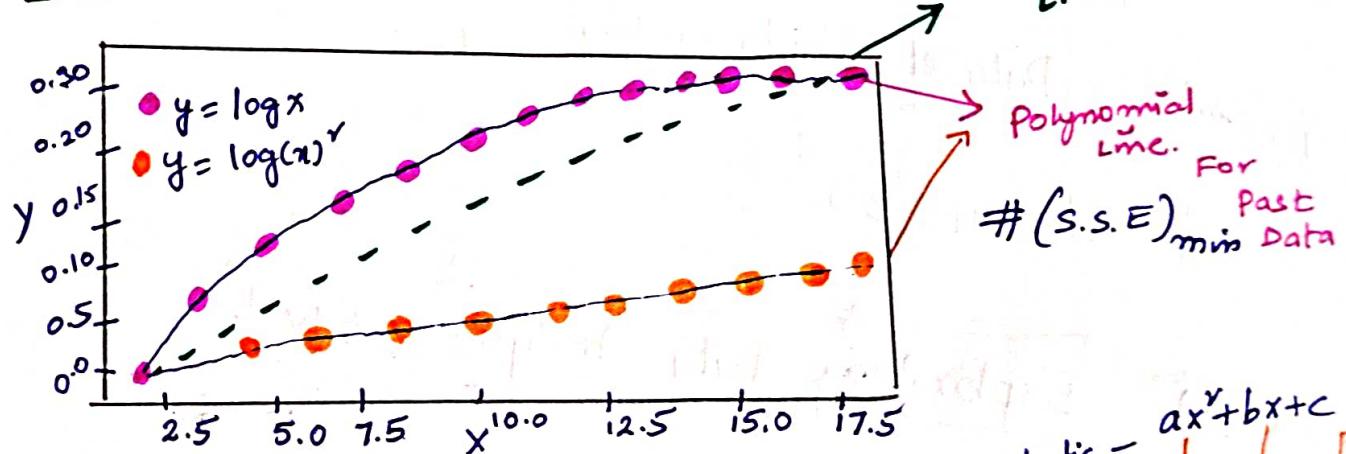
degree 2 ↗

$$a^m \cdot a^n = a^{m+n}$$

Polynomial Regression

Multiply

if we add two variables
degree 1 ↗



linear line

Polynomial line.

#(S.S.E) for Past min Data

#quadratic - $ax^2 + bx + c$

If we want to calculate combinations with squaring.
Degree 2 = $TV^2 (or) x_1^2 + x_2^2 + x_3^2 + x_1 x_2 + x_2 x_3 + x_3 x_1 + \dots$

Degree 3 = $x_1^3 + x_2^3 + x_3^3 + x_1 x_2^2 + x_2 x_3^2 + x_3 x_1^2 + x_1^2 x_2 + x_2^2 x_3 + x_3^2 x_1 + x_1 + x_2 + x_3$

cubic Equation :

$$a x^3 + b x^2 + c x + d$$

STEP 1

Problem Understanding

STEP 2

Data Collection

```
# df = pd.read_csv("Advertising.csv")  
# df.head()
```

STEP 3

Data Understanding

```
# df.shape
```

STEP 4

Exploratory Data Analysis [EDA]

```
# df.describe()
```

```
# sns. Pairplot(df)
```

STEP 5

Data Cleaning

```
# df.isnull().sum()
```

STEP 6

Data Wrangling

```
# x = df.drop("sales", axis=1)
```

```
# y = df["Sales"]
```

new Features
Adding to original data

Polynomial Regression with Scikit-learn

from sklearn.preprocessing import PolynomialFeatures

```
# polynomial_converter = PolynomialFeatures(degree=2, include_bias=False)
```

$$\begin{aligned} & ax_1^2 + bx_2^2 + cx_3^2 \\ & dx_1x_2 + ex_2x_3 + fx_3x_1 \\ & gx_1 + hx_2 + ix_3 \end{aligned}$$

don't include "z" value

```
# x_poly = polynomial_converter.fit_transform
```

↓ calculate
↓ convert
(x)

```
# x_poly.shape
```

Out (200, 9)
rows ↓
columns → changed 3 to 9

from sklearn.model_selection import train_test_split

```
# x_train, x_test, y_train, y_test = train_test_split
```

[x_poly, y, test_size=0.3,

Random_state=29]

200

140 x 9 = train
60 x 9 = test

Step 4

MODEL Fitting On Polynomial Data.

```
from sklearn.linear_model import LinearRegression
```

```
# model = LinearRegression()
```

```
# model.fit(x-train, y-train)
```

Predictions

```
# train-pred = model.predict(x-train)
```

```
# test-pred = model.predict(x-test)
```

Step 5

Evaluation

```
# model.score(x-train, y-train) [ $R^2$  Train]
```

Out 0.986

```
# model.score(x-test, y-test) [ $R^2$  Test]
```

Out 0.980

cross-validation

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

```
# scores = cross_val_score (model, X_poly, y, cv=5)
```

```
# print (scores)
```

```
# scores.mean()
```

```
[Out] [0.987, 0.989, 0.991, 0.958, 0.993]
```

```
mean: 0.984
```

```
# RMSE
```

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

```
# test-RMSE = np.sqrt (mean_squared_error (y-test, test-pred))
```

```
# train-RMSE = np.sqrt (mean_squared_error (y-train, train-pred))
```

```
# print (train-RMSE, test RMSE)
```

```
0.5950, 0.7233
```

Earlier,

Multiple Linear Regression

* RMSE - 1.94

Polynomial Regression

* RMSE - 0.72

Applying loop for knowing better Accuracy From First to Last.

```
• from sklearn.preprocessing import PolynomialFeatures  
• from sklearn.model_selection import train_test_split  
• from sklearn.linear_model import LinearRegression
```

```
# train_rmse_errors = []
```

```
# test_rmse_errors = []
```

for d in range [1, 10] :

```
# polynomial_Converter = PolynomialFeatures(degree=d, include_bias=False)
```

```
# x_poly = polynomial_Converter.fit_transform(x)
```

```
# x_train, x_test, y_train, y_test = train_test_split(x_poly, y, test_size=0.3, RandomState=29)
```

```
# model = LinearRegression()
```

```
# model.fit(x_train, y_train)
```

```
# train_pred = model.predict(x_train)
```

```
# test_pred = model.predict(x_test)
```

```
# train_RMSE = np.sqrt(mean_squared_error(y_train, train_pred))
```

```
# train_rmse_errors.append(train_RMSE)
```

```
# test_RMSE = np.sqrt(mean_squared_error(y_test, test_pred))
```

```
# test_rmse_errors.append(test_RMSE)
```

train_rmse_errors

Out: 0 [1.7345 d=1
1 0.5879 d=2
2 0.4339 d=3
3 0.3517 d=4
4 0.2509 d=5
5 0.19704 d=6
6 0.14214 d=7
7 0.14180 d=8
8 0.16654 d=9

degree

No. of bends

 = quadratic

 = cubic

 = degree of 4

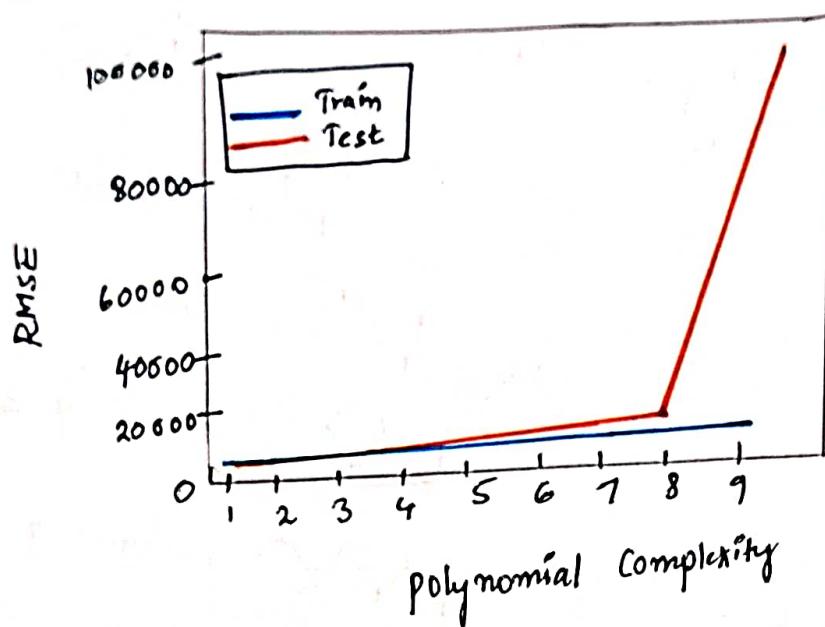
 = degree of 7

test_rmse_errors

Out: 0 1.5161 d=1
1 0.6646 d=2
2 0.5803 d=3
3 0.5077 d=4
4 2.5758 d=5
5 4.4926 d=6
6 1381.404 d=7
7 4449.599 d=8
8 9591.24 d=9

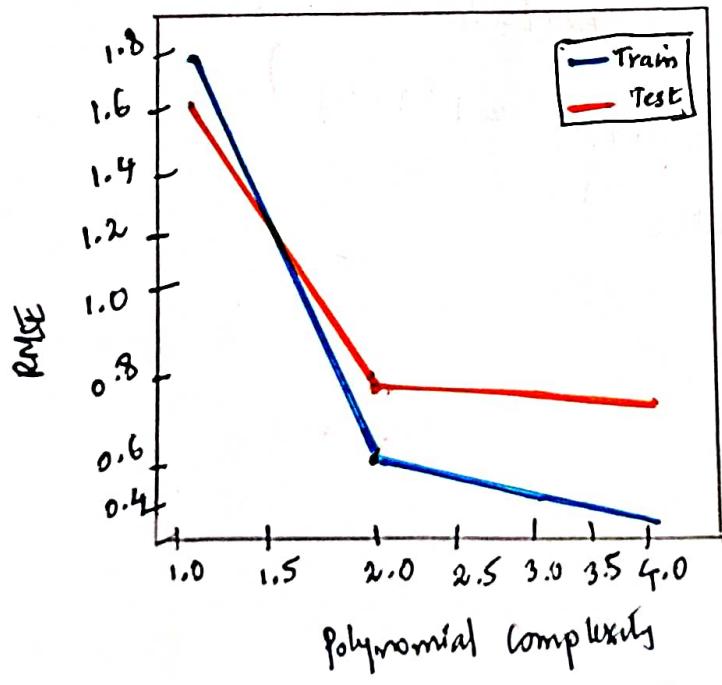
plt.plot (range(1,10), train_rmse_errors, Label = "TRAIN")
plt.plot (range(1,10), test_rmse_errors, Label = "TEST")
plt.xlabel ("Polynomial complexity")
plt.ylabel ("RMSE")
plt.legend ()
plt.show ()

Out :



```
# plt.plot(range(1,5), train_rmse_errors[:4], label = "TRAIN")
# plt.plot(range(1,5), test_rmse_errors [:4], label = "TEST")
# plt.plot(range(1,5), test_rmse_errors [:4])
# plt.xlabel ("Polynomial complexity")
# plt.ylabel ("RMSE")
# plt.legend()
# plt.show()
```

Out :



5/5

Finalizing Model choice

```
# final_poly_Converter = polynomial Features (degree = 2,  
include_bias = False)
```

```
# final_model = Linear Regression()
```

```
# final_model.fit(final_poly_Converter.fit_transform(X), y)
```

out: Linear Regression()

* Saving MODEL and CONVERTER

```
from joblib import dump
```

```
# dump(final_model, "Sales_poly_model.joblib")
```

out: Sales_poly_model.joblib

```
# dump(final_poly_Converter, "Poly_Converter.joblib")
```

out: [poly_converter.joblib]

Deployment & predictions:

Client wants to spend 149K on "TV", 22K on "radio"
12K on "Newspaper" Ads. How many units could we
expect to sell as a result?

```
from joblib import load
```

```
# loaded_poly = load ("poly-converter.joblib")
# loaded_model = load ("Sales-poly-model.joblib")
-
# Campaign_poly = loaded_poly.transform ([[149, 22, 12]])
-
# final_model.predict(campaign_poly)
[OUT]: array([14.5114])
```

anil
16/04/22
9:25 pm