

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.

Coefficient:

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

dependent variable

Independent Variables.

⇒ Example :

TV	Bad to	nP	Sales

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

TV Radio news
paper

y-intercept :

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]

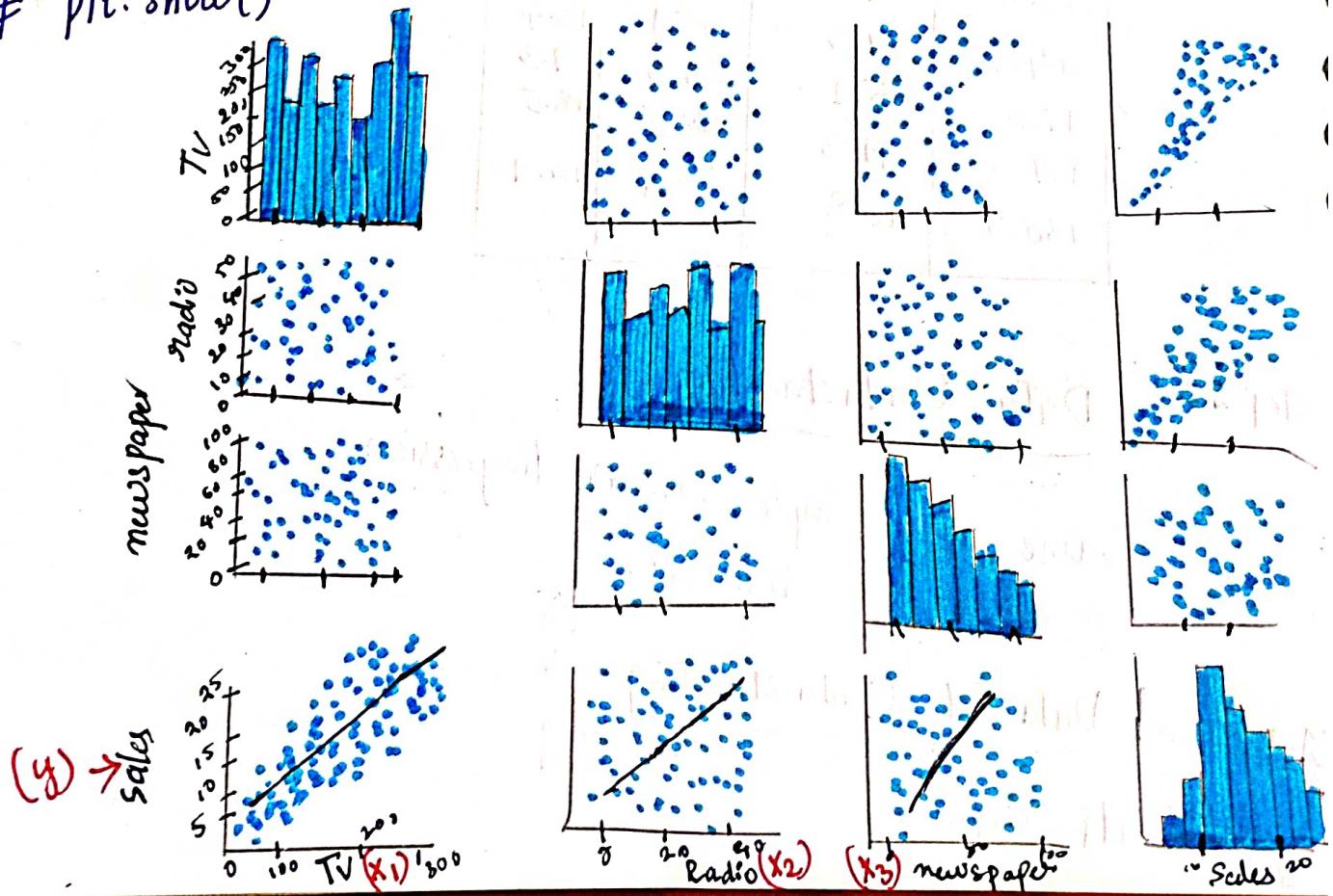
df.describe()

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

$x_{\text{train}}, x_{\text{test}}, y_{\text{train}}, y_{\text{test}} = \text{train-test-split}(x, y, \text{test_size} = 0.3, \text{Random-state} = 29)$

Step-4

Modelling

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

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

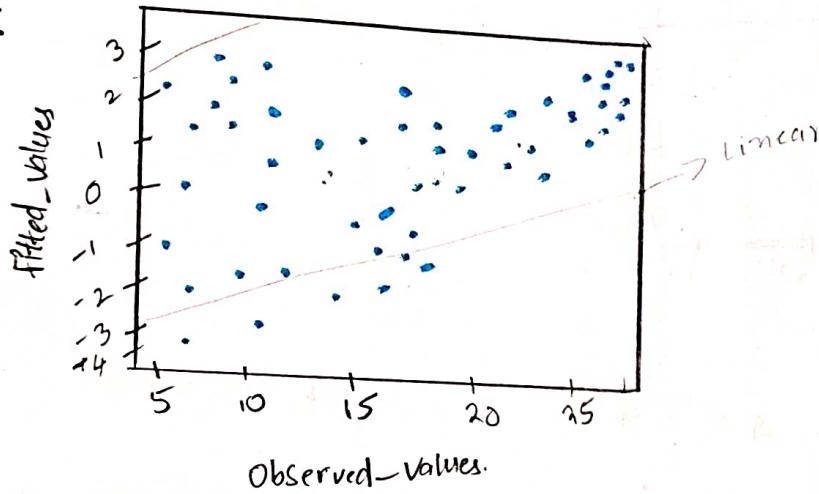
plt. Scatter (y-test, test-res)

plt. xlabel ("Observed-values")

plt. ylabel ("fitted-values")

plt. show()

out:



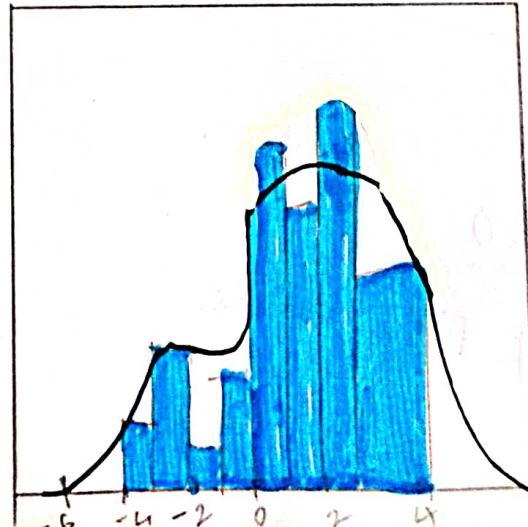
* Normality of errors

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

plt.show()

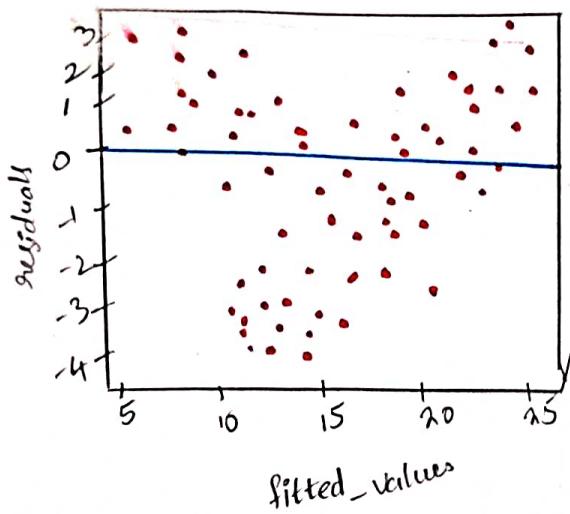
out:

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 $\star \star \star \star$ condition. Than, 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)

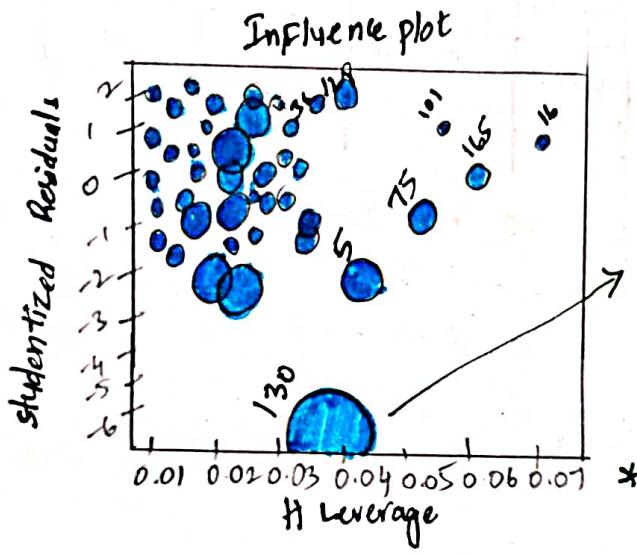
Why?

$$\beta_3 = 0,$$

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

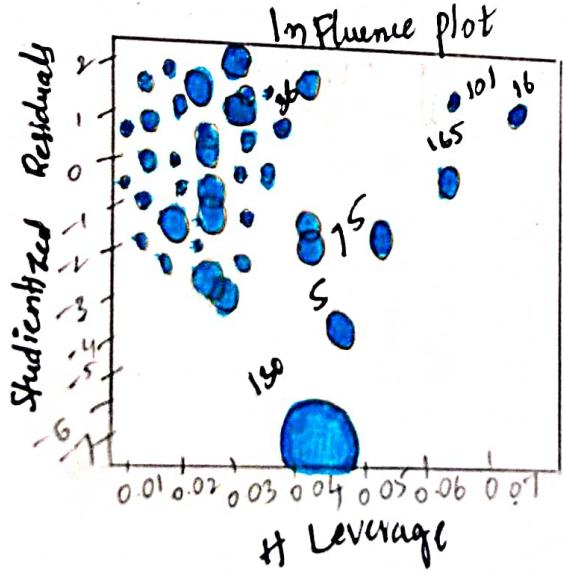
By dropping, we are losing $\frac{1}{3}$ of data = 33%.
In that case, we drop some records, that are influencing more.
 $Ex: - \frac{1}{200} = 0.05\%$.

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]

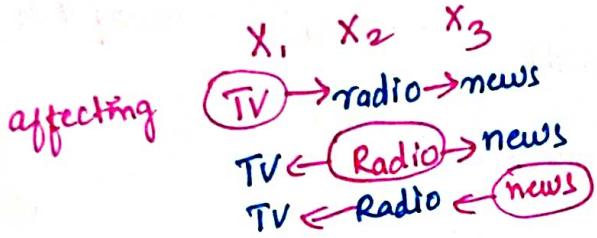
Variance inflation Factor [VIF] measures ratio between the variance for a given regression coefficient with only that variable in the model versus the variance for given regression coefficient with all variables in the model.

1 independent variable

influence on other

independent variable is called as

(VIF)



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

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

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

$$[\text{SSE}]_{\min} \quad [\mathbb{E}(y - \hat{y})^2]_{\min}$$

$$\frac{\partial}{\partial x} \text{ at } x=0$$

#

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?

$$\frac{\partial \left[\mathbb{E}[y - \beta_0 - \beta_1 x_1]^2 \right]}{\partial x_1} \text{ at } x=0 \quad \begin{array}{l} \text{He we have only} \\ \text{One variable.} \end{array}$$

In Multiple linear Regression?

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

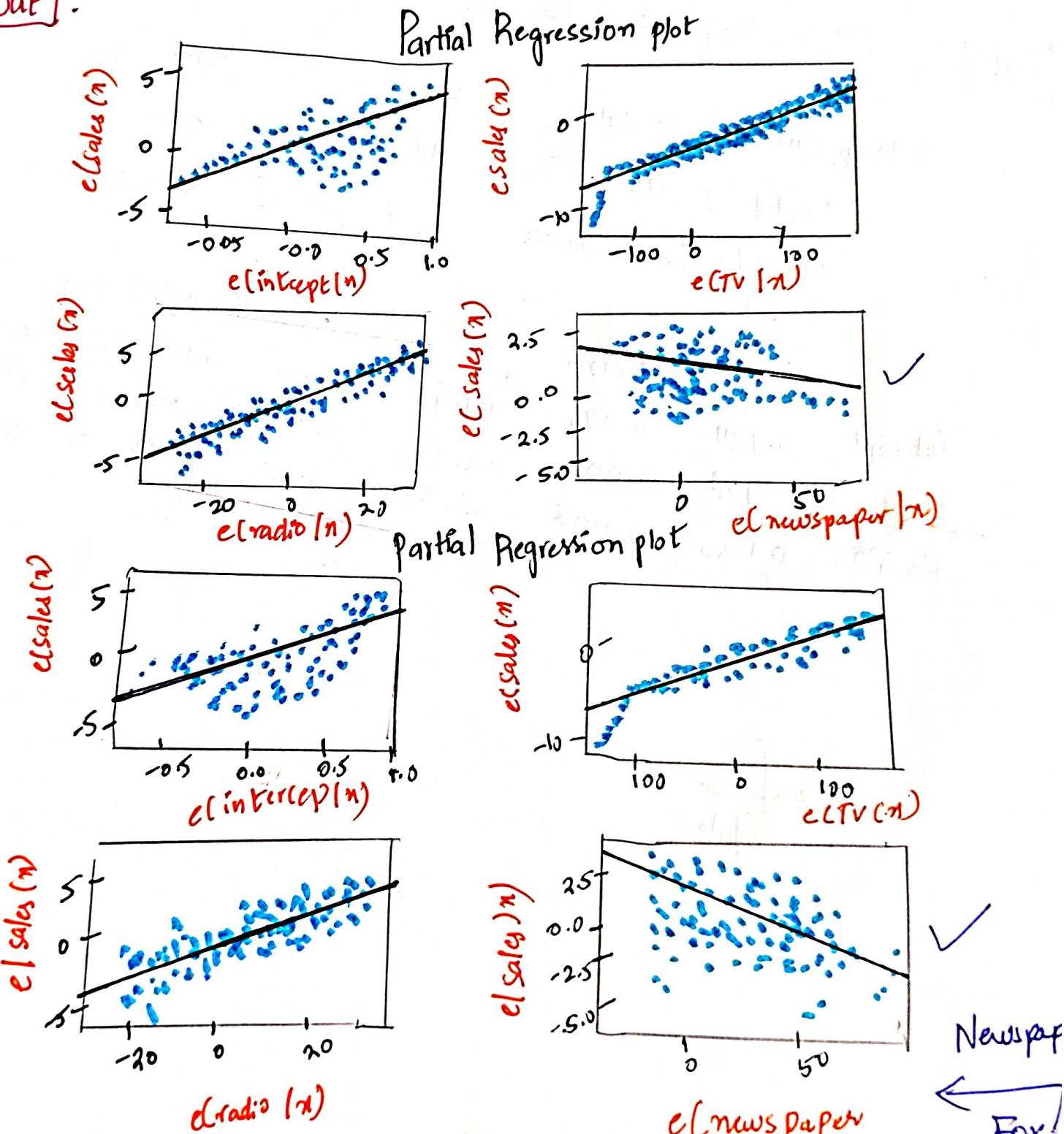
constant calculate constant
constant constant constant

↓
when we calculate with $\frac{\partial}{\partial x_1} \frac{\partial}{\partial x_1}$.
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