# "Simple Linear Regression"

→ Data Science Project life cycle (Crisp-dm) → Data mining.

Cross-Industry Standard Process For

↓ Methodolgy

1. Bussiness problem Understanding

2. Data Understanding → Data collection [Data Eng]
   → Data variables (research - domain Expert
   → Dataset understandrioay

3. Data Preprocessing → EDA
   → Data cleaning
   → Data Wrangling
   → Train / test split

4. Modelling → Apply Different Algorithms

5. Evaluation → Accuray of The Model checking (or) If not go back and start again

6. Presentation → Visulazation ( Tableau, Power BI)

---

⇒ Basic Liberaries 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 <u>relationship</u> between <u>total advertising spend</u> and <u>Sales</u> ?

- Our <u>next</u> ad <u>Campaign</u> will have a <u>total spend of</u> ~$200, how <u>many units</u> do we expect to <u>sell</u> as a <u>result of this</u> ?

Data collection

> , Comma Seperated values

```
# 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: 230.1×100 already done scaling

→ Radio :- Advertising dollars spent on Radio

→ Newspaper : Advertising dollars spent on newspaper

## 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()

| | column | Non Null Count | Dtype |
|---|---|---|---|
| 0 | TV | 200 non null | float 64 |
| 1 | Radio | 200 nan null | float 64 |
| 2 | Newspaper | 200 non null | float 64 |
| 3 | Sales | 200 non null | float 64 |

[out]

(200,4)

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

A :- We have simplified This quite a bit by Combining all Features into " total spend "

We add new column . total of 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 | 250.3 |

Step- 3.1 Exploratory Data Analysis [EDA]

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 ?
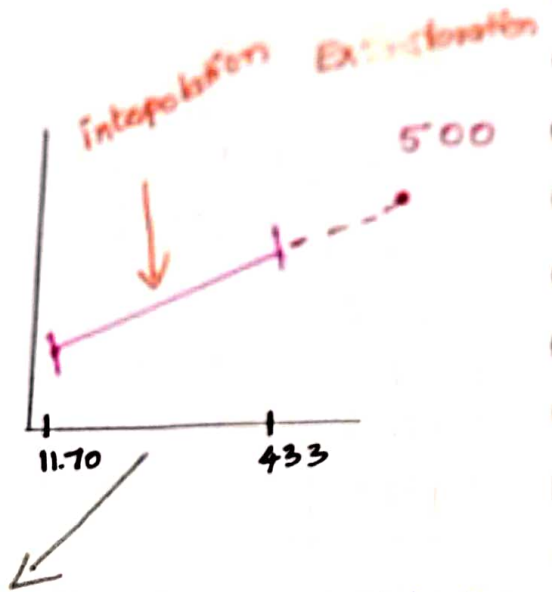
3. Given ad spending, can Sales be predicted ?

# # 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

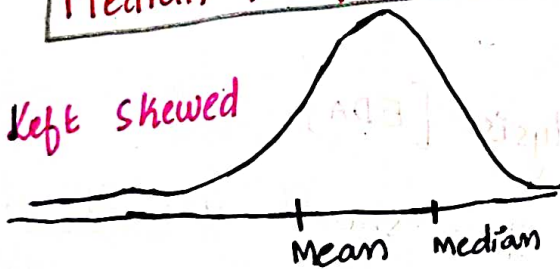Mean and median are close. They normal distribution

Interpolation   Extrapolation



500

11.70      433

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

* **Extrapolation** :- Prediction The value, Outside of known values it is called Extrapolation.
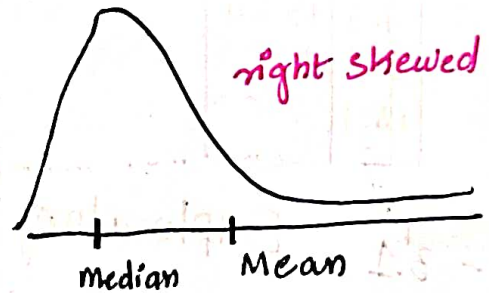
***
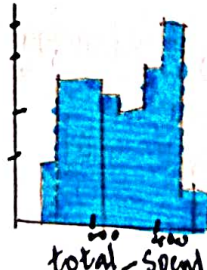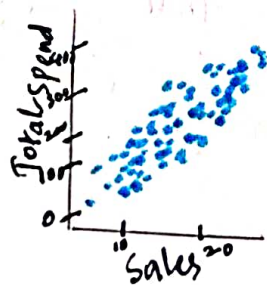⇒ Median > (greater Than) Mean

⇒ Median < [less Than] Mean

Left skewed



Mean   Median

right skewed



median   Mean

# # Sns. pairplot (df)

# # plt. show ()

Out :



1. Lineanity
2. direction
3. Correlation

x value, y value increasing

# df. corr ()

|          | Sales    | total spend |
|----------|----------|-------------|
| Sales    | 1.000000 | 0.867712    |
| total_Spend | 0.867712 | 1.000000 |

## Step 3.2   Data cleaning

```
# df. is null ( ). Sum ( )
```

out

```
Sales       0
total-spend 0
```

## Step 3.3   Data Wrangling

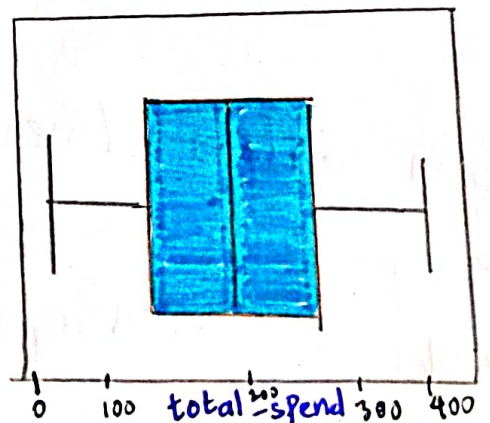*. <u>Outliers</u>
— x —

```
# 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 import train_test_split

# 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 "Baseline model (or) raw model.
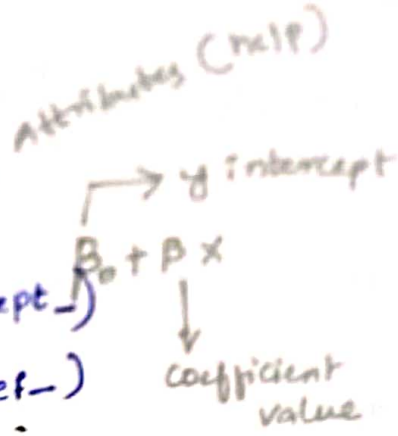
# from sklearn.linear_model import LinearRegression

storing
# model = LinearRegression()

fitting the data

```
# model. fit (x.train, y_train)
```
Out :- Linear Regression ()

Attributes (help)

```
# Print ("Model Intercept : ", model. intercept_)
  Print 1(" Model Coefficient: ", model. coef_)
```

→ y intercept

$\beta_0 + \beta x$ → coefficient value

out, :- Model Inter cept : 4.33176
       Model Coefficient : 0.0480

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

⇒ <u>Predictions</u>
   — * —

```
# 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]])
```
→ 2 dimensional

input should given in

out    array ([13.93860])

Predicting on X train, X test

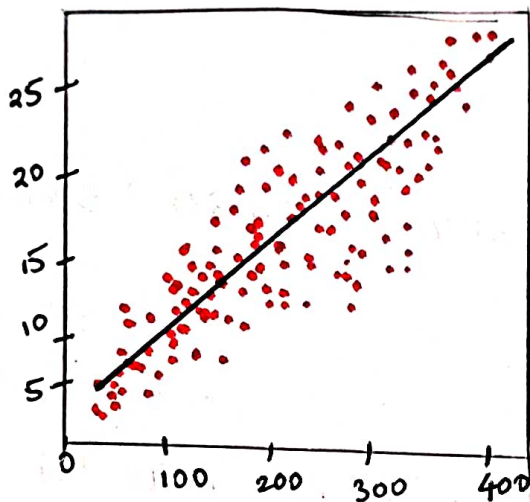| | X | y | ŷ | |
|---|---|---|---|---|
| X Train | 1 | 10 | 14 | → Train accuracy |
| | 2 | 15 | 12 | |
| | 3 | 20 | 13 | |
| | 4 | 40 | 33 | |
| X Test | 5 | 50 | 42 | → Test accuracy |

```
# 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 ()



**Steps**

⇒ Evaluation metrics
  ——— * ———

from sklearn . metrics **Import** mean_ absolute_error

# 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 : 1.97768
  MAE For Test data : 1. 90653

```python
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).
np.sqrt ↗

```python
# 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

```python
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

(14)

Another way for R²

# model score (x_train, y_train)

Out: 0.76534

# model Score (x_test, y_test)

Out: 0.740017

# Cross_validation → 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]

→ 0.74337 → should be Equal to Test square
                      Then it is good model.

, Linear Regression assumptions
                              * ~~~~ *
                                    Errors

    L  →  Linearity
    I  →  Independent
    N  →  Normality
    E  →  Equal Variance

10/4/22
3:30 Am

# Assumptions of Linear regression

1. check whether model has overfitting (or) underfitting problem

2. Is test Accuracy = Cross Validation Score

3. check assumptions (if it is Linear Regression)

Purely Assumption of The Researcher

4. check model meets bussiness problem requirements

5. Finally, save the model and share to deployment Team.

---

1 Que :- Is model has Overfitting (or) underfitting problem ?

Ans :- it's good model.

2 Que  Is test accuracy = Cross Validation Score ?

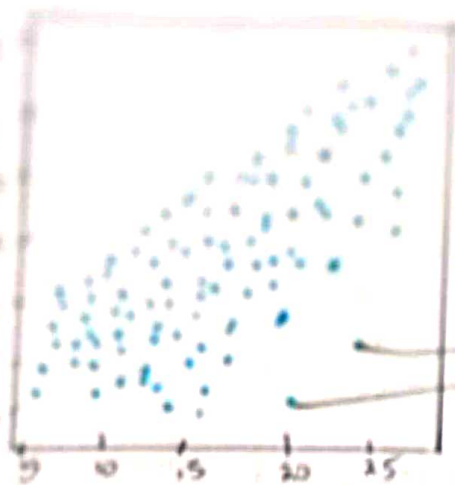Ans :- Applied K-Fold cross validation and it is Equal in test accuracy and CV Score.

3 Que :- check Assumptions (if it is linear regression)

Ans :- Line Assumptions.

1 Linearity of Errors

# test_res = y_test - test-Predictions.

↑ Saved in TRI.

# plt. scatter (y_test, test_res)
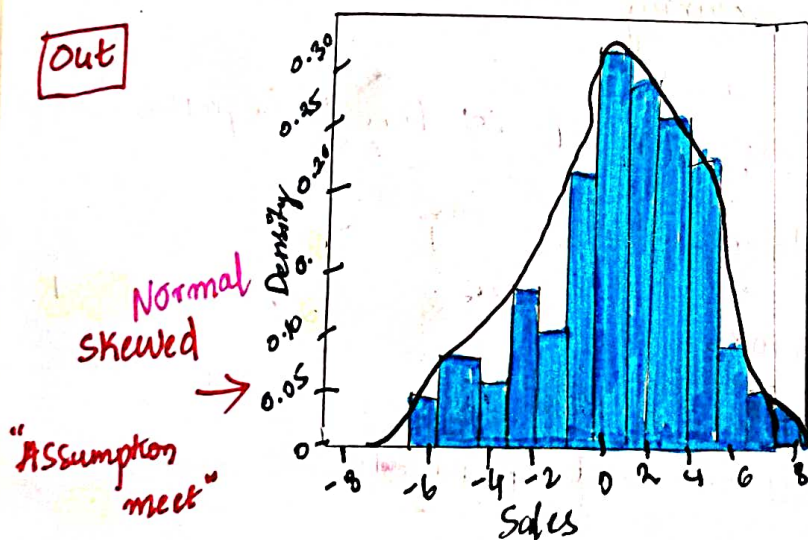plt. XLabel ("Observed_values")
plt. YLabel ("fitted_values")
plt. show()

→ we can ignore
The Edge points
upto ± 5% error

Observed values

2. **Normality of Errors**

\# Sns. distplot (test_res, bins = 15, kde=True)

plt. show ()

Out



Normal
skewed
→

"Assumption meet"

Sales

% of positive Errors
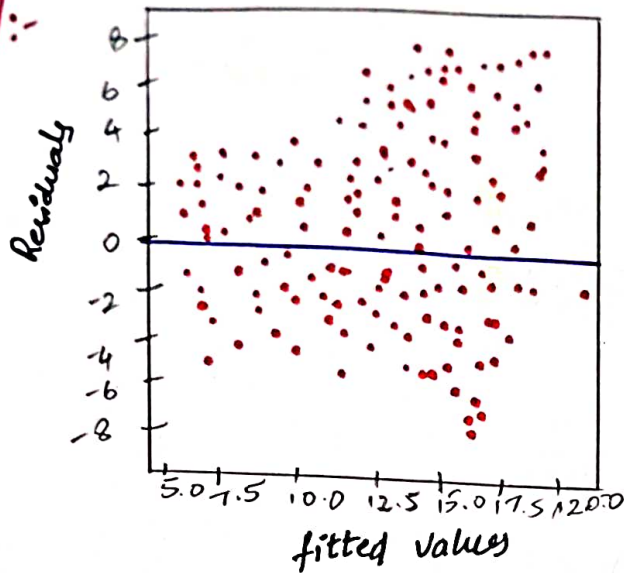% of negative Errors
color = red

(Homoscadesicity)

3. **Equal variance of Error**

\# plt. scatter (test_predictions, test_res, c = "r")

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

ax = axis → h= horizontal line

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

Out:-



For model

$H_0$ :- Avg line best fit line

$H_1$ : regression line best line

↓ Anova Test

$p : \leq \alpha$

$p \geq \alpha$

1 Sample two Tail test

variable → Significant

$H_0$ :- $\beta_1 = 0$

$H$ :- $\beta_1 \neq 0$ → Consider

variable insignificant don't → 

4. Variable , Significance

```python
#  import Statsmodels.formula.api as Smf

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

#  m.Summary()
```

"y" upon "x"
y multiple column
y~x+x+x+x
ordinary least square (SSE) min

out   OLS Regression Results.

Train

$p = 5.06 \times 10^{-62}$
$p < 0.05$
Reject $H_0$

Dep. Variable : Y          R-squared : 0.753

Model : OLS          Adj. R-Squared : 0.752

method : Least squares      F. static : 603.4

Prob.(F.static) : 5.06 $\times 10^{-62}$

Confidence Interval
Std.Error = $\frac{\sigma}{\sqrt{n}}$

$\bar{x} \pm Z_{1-\alpha} \frac{\sigma}{\sqrt{n}}$
Central limit Theorem

one sample
t Test applied

| | coeffs | Std.Err | t | P>\|t\| | [0.825  0.915] |
|---|---|---|---|---|---|
| intercept | 4.9730 | 0.499 | 9.676 | 0.000 | 3.378  5.508 |
| x | 0.0487 | 0.002 | 24.564 | 0.000 | 0.045  0.053 |

1 Sample T Test

Two Tail   95%   2.5  Error

When Average line is going to Fail in Regression?

When Sum of Square (SSE)req Average line is Less than the Sum of Square of regression line (SSE)reg it going To Fail in this Situation.

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

Errors

## Step 6    Final Inferences

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

[Out] array ([13.9898 ])

⇒ <u>Save</u> a <u>Model</u>

from joblib import dump

\# dump ( model , "Sales - model . Joblib")

[Out] : ["Sales - model. Joblib"] → saves in working directory.

⇒ <u>Load</u> a <u>Model</u> → From client side

\# from Joblib import load

```python
# loaded_model = load("sales_model.joblib")
```

```python
# loaded_model.Predict([[200]])
```

out  array([13.989])

```python
# loaded_model.Predict([[500]]) # check with
                                 #  multiple values.
```

out ∴ array([28.3488])

with
74% Accuracy.

19/4/22
5:00pm