

Linear Regression

```
1 1 - It is a Supervised Learning predictive model used for finding the linear
2   relationship between dependent variable(Y) and independent variable(s)(X)
3 2 - Linear Regression is used when output variable has continues values.
4 3 - Equaion of linear regression :  $Y = mx + c$ 
5
6 Two Types:
7 1. Simple Linear Regression(Only one independent Variable)
8 2. Multiple Linear Regreesion (Two or more independent Variables)
```

Required Liabraries for Linear Regression

In []:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.preprocessing import LabelEncoder,OneHotEncoder
6 from sklearn.model_selection import train_test_split
7 from sklearn.linear_model import LinearRegression
8 from sklearn.metrics import mean_absolute_error,mean_squared_error,r2_scores
```

Assumption of Linear Regression

1 - Linearity

```
1 1-Indipendent variable(X) and Dependent variable(Y) should be linear to each other.
2 2-i.e. relationship between x and y is linear.
```

2 - Independence

```
1 All observations(Indipendent variable) are indipendent to each other.
```

3 - No Multicollinearity

```
1 There is no strong co-relation between indipendent variables.
```

4 - Normality

```
1 1 - Residuals should have normally distribution.
2 2 - Normal distribution also known as Gaussian Distribution.
3 3 - Gaussian distribution is a probablity Distribution that is
4   symmetric about the mean and it should be in bell shape curve.
```

5 - Homoscedasticity

- 1 1 - Homo means same & scedasticity means variance , so it means same variance.
- 2 2 - If the variance in the residual error is constant regardless of the dependent
- 3 variable(x) then it is homoscedasticity.

Project Steps

In []:

- 1 1. Problem Statement
- 2 2. Data Gathering(JSON, CSV, Excel, PDF, Images, Videos, Text, etc)
- 3 3. Exploratory Data Analysis(Pandas, Matplotlib, Seaborn) (60 %)
- 4 4. Feature Engineering (Scaling, Handling Outliers, Encoding, Log Transform, Binning)
- 5 5. Feature Selection (Required Feature to train the model)
- 6 6. Model Building (LR, Logistic Regression, DT, RF, etc)
- 7 7. Model Evaluation (MSE, AMSE, Classification Report, Confusion matrix, etc)
- 8 8. Model Deployment (AWS, GCP, Azure, Heroku)

Coefficient of Correlation(R)

- 1 1 -It is also called as Karl Pearson Correlation Coefficient/ R-value.
- 2 2 -It gives the correlation between IV(X1) and IV(X2).
- 3 3 -It gives correlation between IV(X) and DV(Y).
- 4 4 -It gives the correlation between all the variables present in the Dataset.
- 5 5 -Range of R-value is -1 to +1
- 6
- 7 $R > 0.7$ or $R < -0.7$ ---> Good predictors or strongly correlated.
- 8 $R = -0.3$ to 0.3 ---> Bad predictors.
- 9 $R = 0$ ---> No relation
- 10
- 11 $R_{xy} = \text{Covariance} / \text{Std dev of}(x,y)$

Covariance

- 1 It is measure of association between X and Y.
- 2 1 - If Y increase with increasing X then it is +ve covariance.
- 3 2 - If Y decrease with incresing X then it is -ve covariance.
- 4 3 - If there is no linear tendency for Y with change to X then it is zero(0)

Gradient Descent Algorithm (GD)

In []:

- 1 1. Use to find best fit line.
- 2 2. Best value of m and c.
- 3 3. To find m and c, GD uses partial derivatives which is use reduce cost function(MSE).
- 4 4. Cost function = Mean squared Error(MSE)
- 5 5. It uses the X_train and Y_train.

Best Fit Line(BFL)

In []:

- 1 1. It passes to maximum number of datapoints
- 2 2. Line which has lowest Sum of Errors (MSE)
- 3 3. Gradient Descent Algorithm is used to find BFL
- 4 3. Algorithm finds one BFL from infinite number of possibilities.
- 5 4. BFL always passes from Xmean and Ymean.
- 6

Mean Square Error(MSE)

In []:

- 1 1. It is define as the average of square of difference between Yactual and Ypredicted.
- 2 2. $MSE = \frac{\sum(Y_{act} - Y_{pred})^2}{N}$

Cost Function (j)

In []:

- 1 1. Cost function for 'm'
- 2 $Dj/Dm = -2/N \sum(Y_{act} - Y_{pred}) * X_{act}$
- 3
- 4 2. Cost function for 'c'
- 5 $Dj/Dc = -2/N \sum(Y_{act} - Y_{pred})$

Model Evaluation

Prediction

- 1 `y_pred = linear_model.predict(x_test)` >> we get predicted values of dependant
- 2 variable Y

1. SSE(sum of square error)

- 1 1. Called as Residual or Schochastic Error
- 2 2. Residual Error = (y-actual - y-pred)
- 3 3. Sum of Square Error :
- 4 $SSE = \sum(y_{-actual} - y_{-pred})^2$

2. SSR(Sum of Square Due to Regression)

- 1 1. Called as Regression or Deterministic Error
- 2 2. Regression Error = (y-pred - ymean)
- 3 3. Sum of Square Due to Regression :
- 4 $SSR = \sum (y\text{-pred} - y\text{mean})^2$

3. SST(Sum of Squares Total)

- 1 1. Square difference between Dependent variable()Y and its mean.
- 2 2. $SST = \sum (y\text{-actual} - y\text{-mean})^2$
- 3 3. $SST = SSE + SSR$

4. mean_squared_error

- 1 `mse = mean_squared_error(y_test,y_pred)`
- 2 `print('Mean Squared Error is :',mse)`

5. Root Mean Squared Error

- 1 `rmse = np.sqrt(mse)`
- 2 `print('Root Mean Squared Error is :',rmse)`

6. Mean Absolute Error

- 1 `mae = mean_absolute_error(y_test,y_pred)`
- 2 `print('Mean absolute Error is :',mae)`

7. r2_score (Coefficient of determination)

In []:

- 1 1. It is to check goodness of best fit line
- 2 2. For Good Correlated features R2 will increase more
- 3 3. For Bad Predictors r2 will very low (0)
- 4 4. R2 can be negative
- 5 5. $R2 = 1 - (SSE/SST)$

Disadvantage of r2_score

In []:

```

1 1. R2 will never decrease
2 2. When we add more features R2 score will increase
3   (For Correlated and non correlated features)
4 3. ex. if R value is low for x3 then again r2 will increase
5
6 for Testing Data : r2_score(y_test,y_pred)
7 for Training Data : r2_score(y_train,y_pred_train)

```

8. Adjusted r2_Score

In []:

```

1 1. When we add more features R2 score will increase ,this is disadvantage of
2   R2_score so we use Adjusted R2 score.
3 2. There is no built in function for adjusted R2 score, we need to calculate it.
4
5 Adjusted r2_score = R2 - [(k-1/n-k)*(1-R2)^2]
6 k - no. of parameters
7 n - total no. of size

```

Advantages of Linear Regression

In []:

```

1 1. Perform exceptionally well for linearly separated data
2 2. Easy to implement
3 3. It can handle overfitting

```

Disdvantages of Linear Regression

In []:

```

1 1. Model Fails,If the relation between independent variables and dependent variab
2 2. If the independent variables are correlated, then it may affect performance
3 3. Impact of missing Values / Sensitive to missing Values
4 4. Impact of Outliers / Sensitive to outliers

```

Applications

In []:

```

1 1. Price
2 2. Population
3 3. Age
4 4. Any Contineous Data

```