

A hand holding a pen is pointing at a line graph on a screen. The graph shows a fluctuating line, likely representing stock prices. The background is slightly blurred, focusing attention on the hand and the graph.

# **Supervised ML- Regression Capstone Project**

## **Yes Bank-Stock Price Prediction**

Abhishek Kumar  
Amitha Sridhar  
Mohita Rathour  
Mukesh Sablani  
Sanjay Paul

# Introduction

*“All there is to investing is picking good stocks at good times and staying with them as long as they remain good companies.”*

*-Warren Buffett*

The stock market is known for being volatile, dynamic, and nonlinear. Accurate stock price prediction is extremely challenging because of multiple (macro and micro) factors, such as politics, global economic conditions, unexpected events, a company's financial performance, and so on.

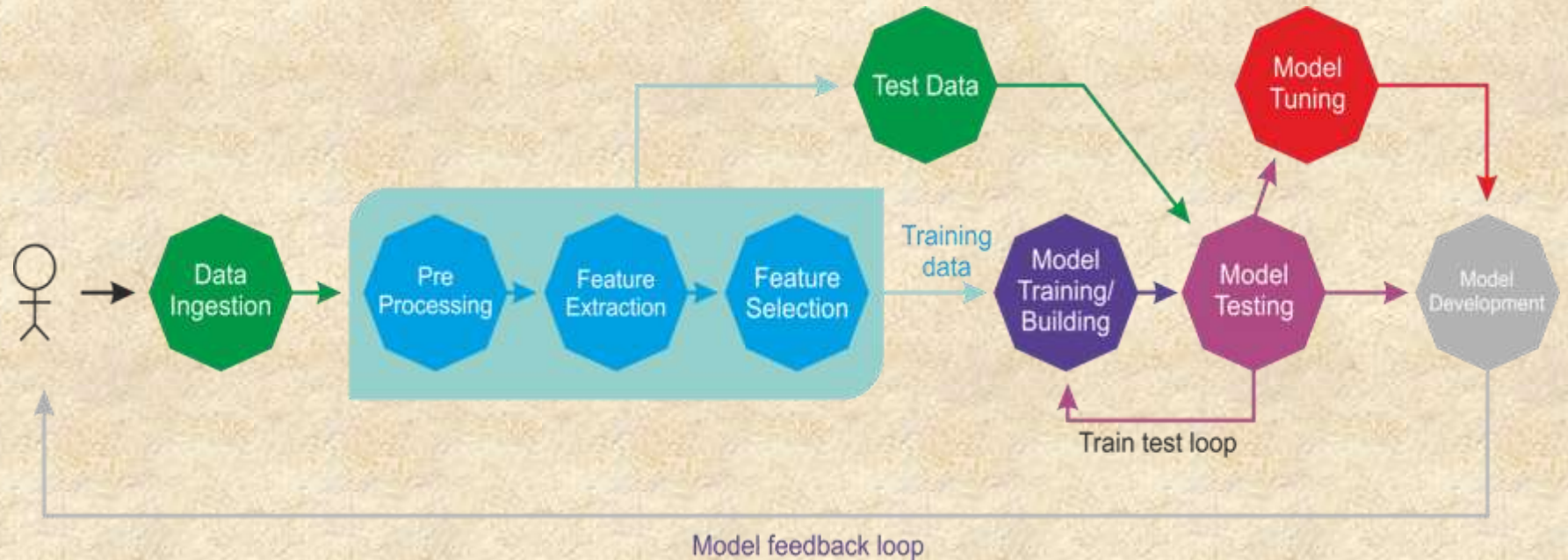
In the era of big data, deep learning for predicting stock market prices and trends has become even more popular than before.

Yes Bank is an Indian bank headquartered in [Mumbai](#), India and was founded by [Rana Kapoor](#) and [Ashok Kapoor](#) in 2004.

In 2018 the Enforcement Directorate (ED) has alleged that Yes Bank co-founder Rana Kapoor and Dewan Housing Finance Limited (DHFL) promoters Kapil and Dheeraj Wadhawan siphoned off funds worth ₹ 5,050 crore through suspicious transactions.

In this project we will be studying the stock price pattern from the inception of the Yes Bank Stock and the effects of the alleged 2018 fraud case of Rana Kapoor on the stock price.

# Workflow



# Features

This data has 185 rows and 5 columns

- > **Date:-**A trade date refers to the month, day, and year that an order is executed in the market.
- > **Open:-**It is the price at which the financial security opens in the market when trading begins.
- > **High:-** the highest price at which a stock traded during the course of the trading day
- > **Low:-**the lowest price that a stock trades in that day.
- > **Close:-**the last level at which it was traded on any given day.

# Year wise Study Open/Close



It's seen that the Opening value of the stock Price has a steady growth since its inception till 2018. A total growth of around 900%. Then after the fraud case there is steep fall post 2018.

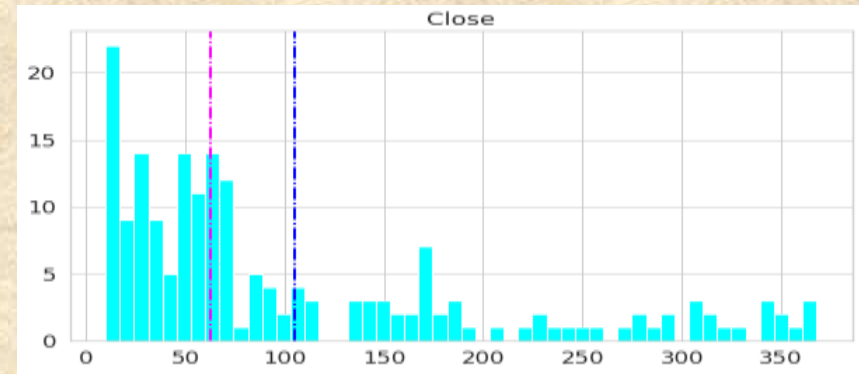
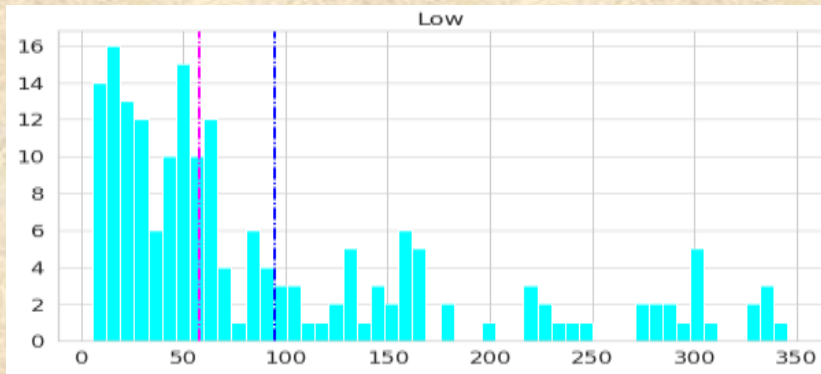
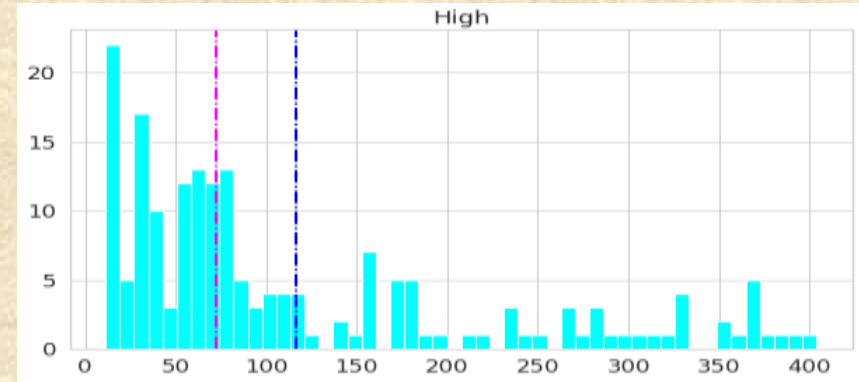
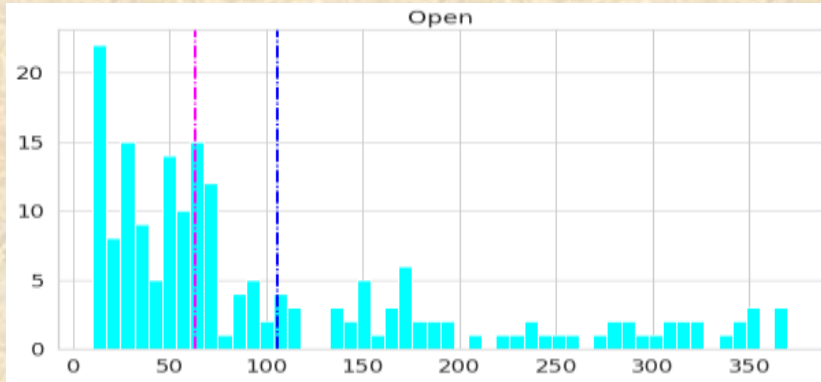
# Year wise Study High/Low



The Growth pattern shadows that of the Open/Close pattern. A steady growth since its inception till 2018 and then a steep fall post the alleged fraud case.

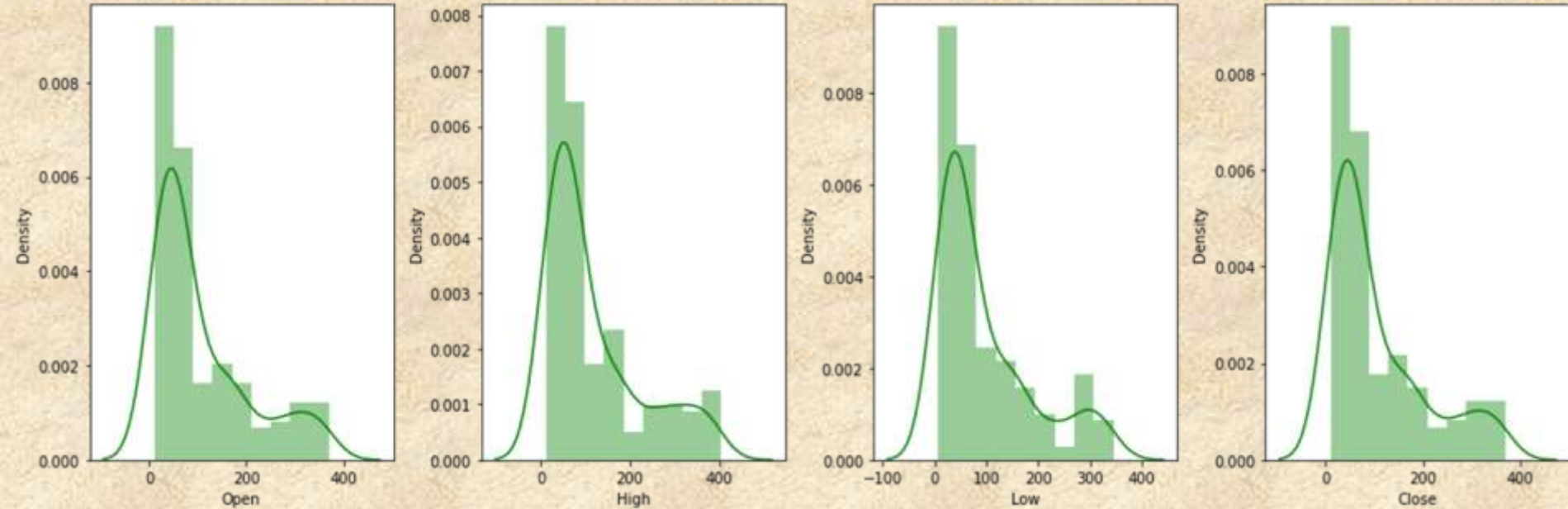


# Mean and Median Study



Data isn't evenly distributed. The difference between the mean and the median is high. Thus normalization of data is needed.

# Data Distribution Graph



Data is Positively Skewed

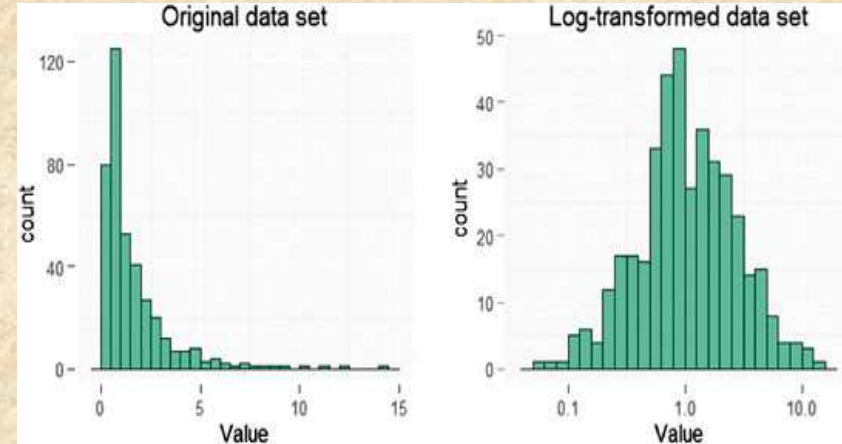


# Data Transformation

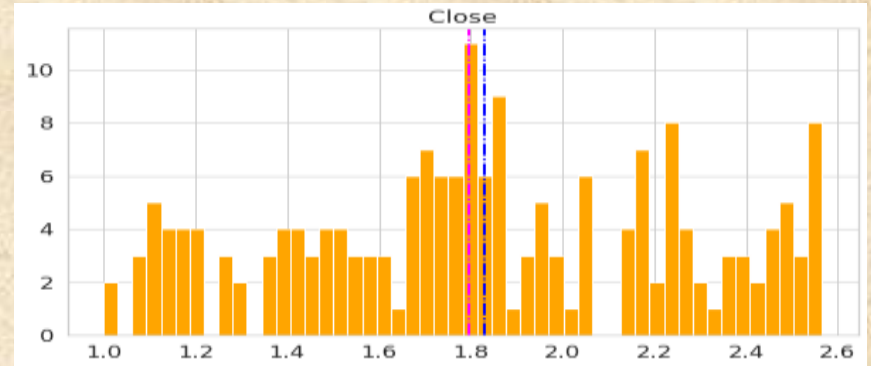
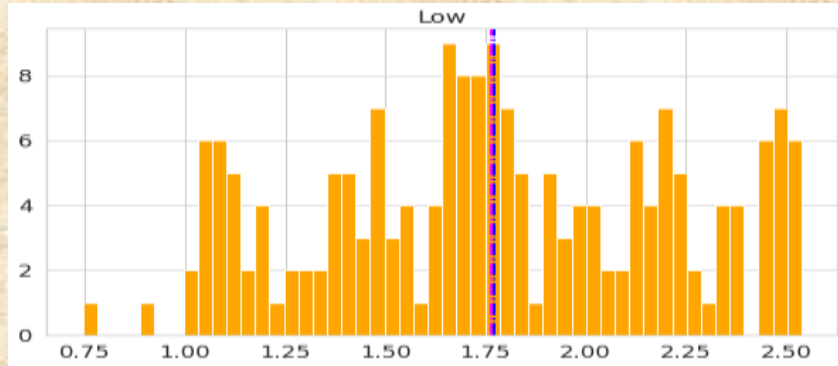
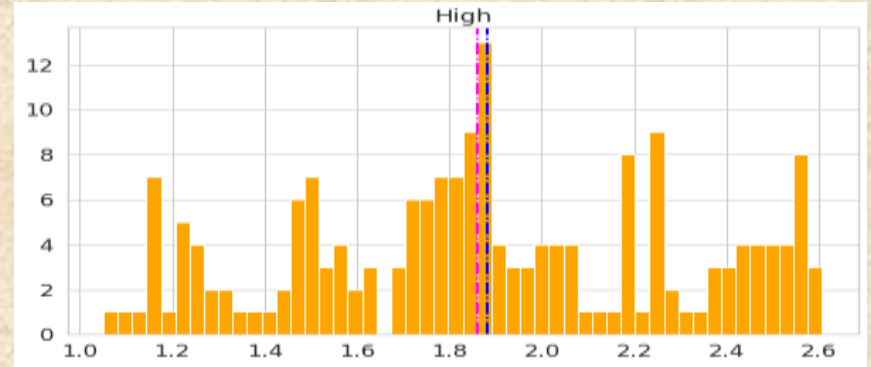
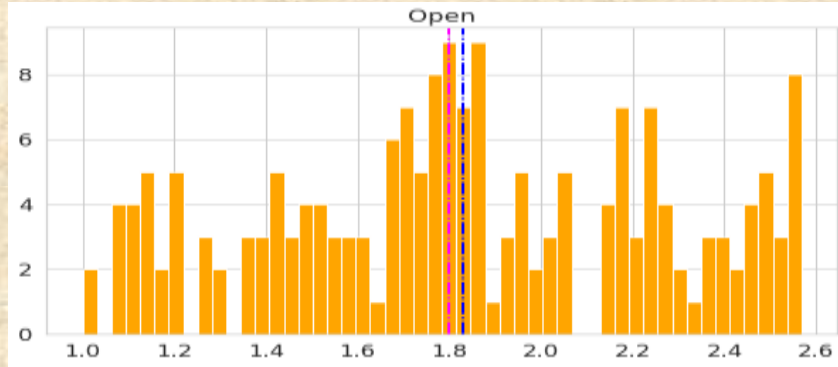
Data transformation is the process of converting, **cleansing**, and structuring data into a usable format that can be analyzed to support decision making processes.

Log transformation is a data transformation method in which it replaces each variable  $x$  with a  $\log(x)$ .

It is primarily used to convert a **skewed distribution** to a normal distribution/less-skewed distribution.

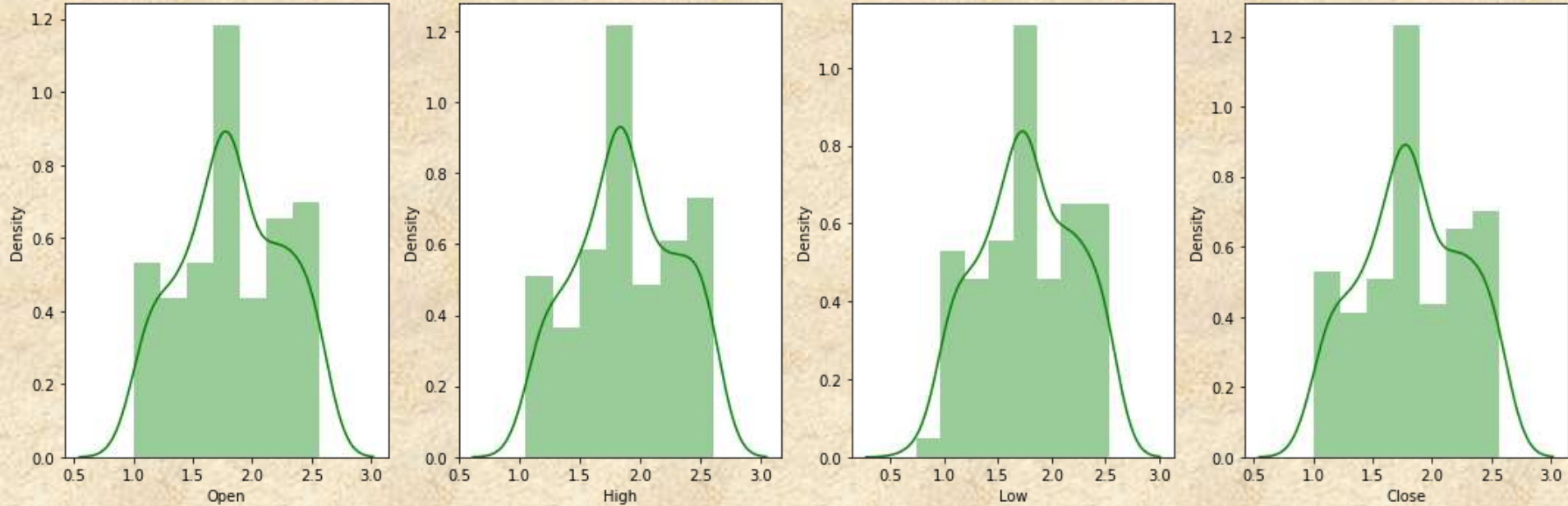


# Mean and Median after Log Transformation



After applying Log Transformation, the difference between Mean-Median is minimalized.

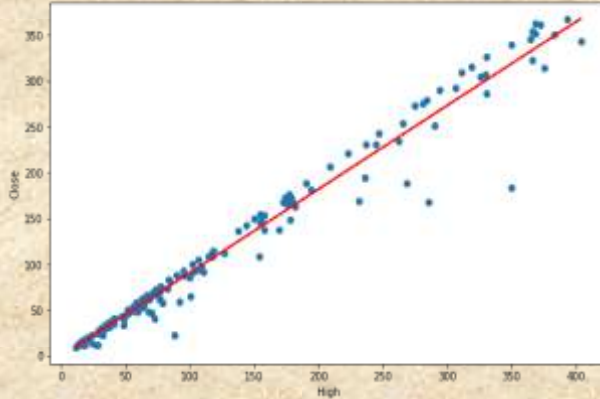
# Data distribution graph after Log Transformation



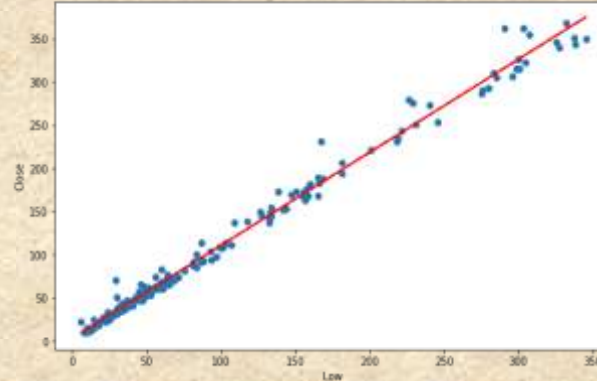
Post Log Transformation, data is uniformly distributed. Normal Bell curve achieved.

# Correlation Graphs

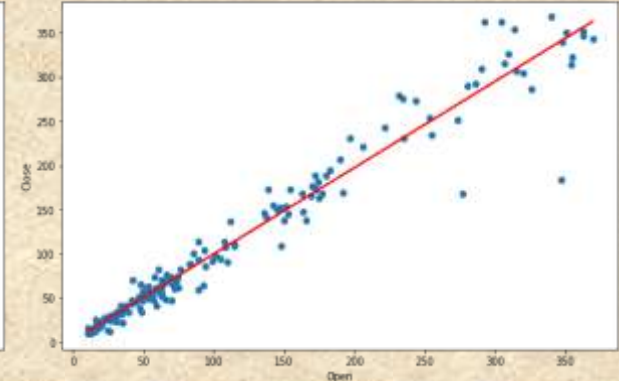
Price Vs High\_correlation:0.9850513315779623



Price Vs Low\_correlation:0.9953579476474373



Price Vs Open\_correlation:0.9779710062230934



Independent Features:- High , Low, Open

Dependent Feature:- Close

High Correlation between the dependent and independent Variables.

# Multicollinearity And VIF

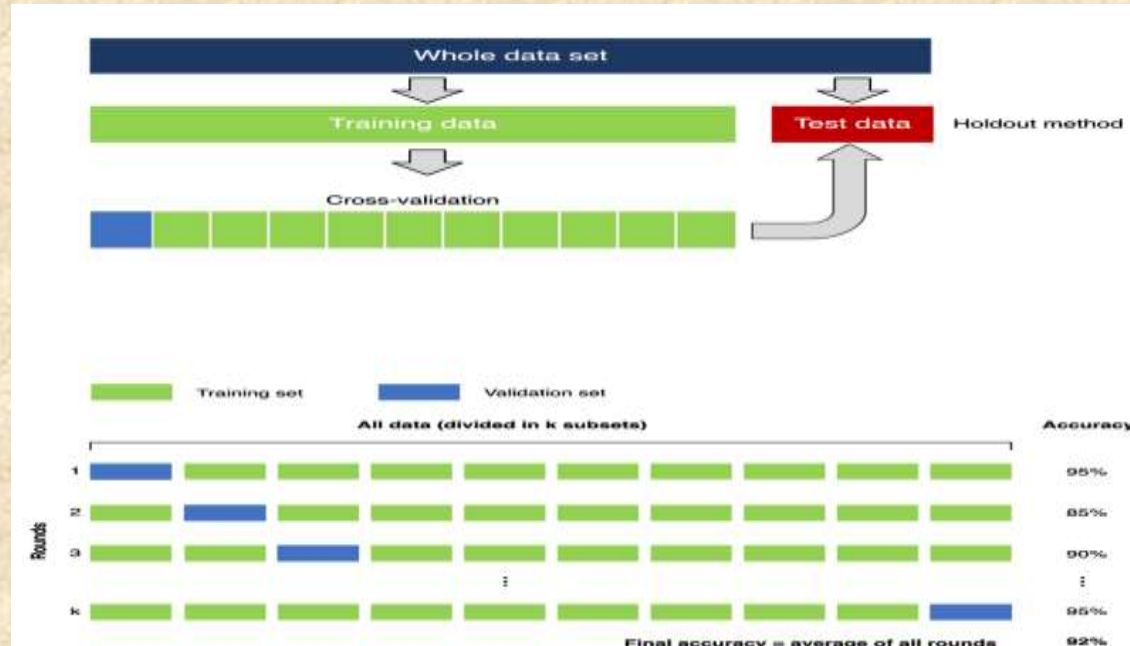
Multicollinearity is the occurrence of high intercorrelations among two or more independent variables in a regression model.

A variance inflation factor (VIF) is a measure of the amount of multicollinearity in regression analysis.

$$VIF_i = \frac{1}{1 - R_i^2}$$

Feature	VIF
Open	175.1857041
High	167.0575232
Low	71.574137

# Train Test Model



**Train-Test Split** is a procedure that allows to simulate how a model would perform on given dataset.

**Cross-Validation** is a method of evaluating and comparing learning algorithms by dividing data into two segments: one used to learn or train a model and the other used to validate the model.

**Hyperparameter tuning** consists of finding a set of optimal hyperparameter values for a learning algorithm while applying this optimized algorithm to any data set.

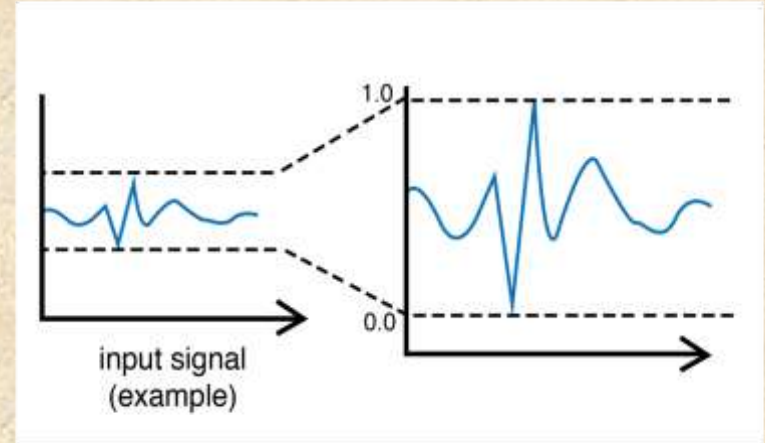


# Feature Scaling

Feature Scaling is the process of scaling or converting all the values in our dataset to a given scale.

The MinMaxscaler is a type of scaler that scales the minimum and maximum values to be between 0 and 1 respectively.

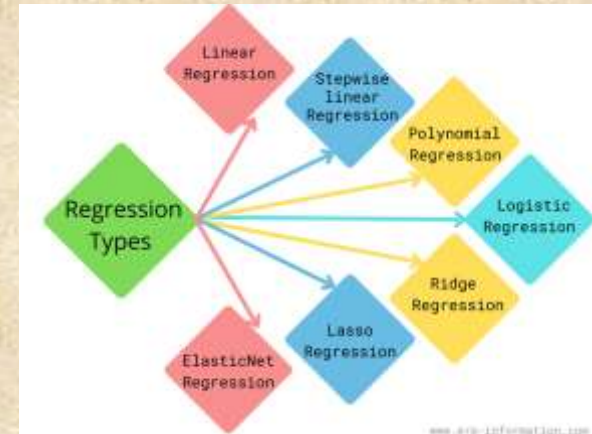
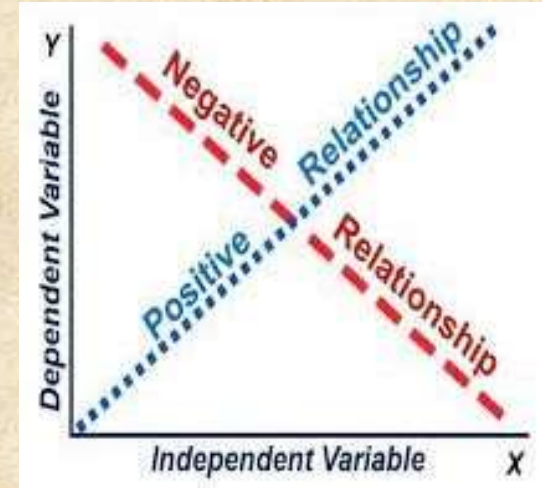
It makes the flow of gradient descent smooth and helps algorithms quickly reach the minima of the cost function. Without scaling features, the algorithm may be biased toward the feature which has values higher in magnitude.



# Model Training

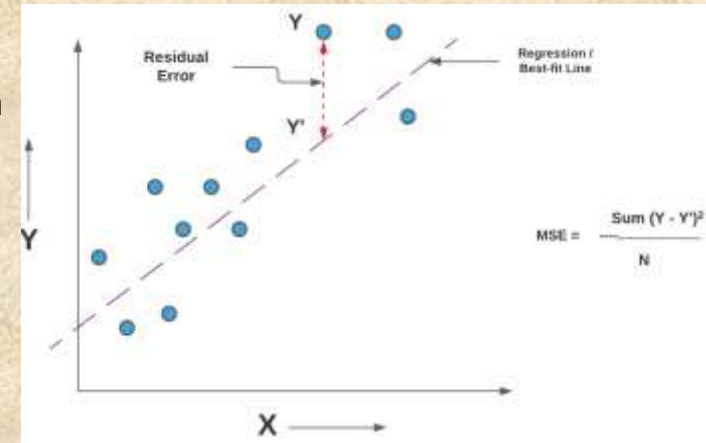
Regression is the statistical method used in **finance** and other disciplines to determine the strength and character of the relationship between **one dependent variable** (Close) and a series of **other independent variables** (High,Low,Open).

Model training in machine language is the process of feeding an ML algorithm with data to help identify and learn good values for all **attributes involved**. It is the primary step in machine learning, resulting in a working model that can then be validated, tested and deployed.



# Metric Comparison Features

→ **MSE**:- Mean squared error (MSE) measures the amount of error in statistical models. It provides the average squared difference between the observed and predicted values.

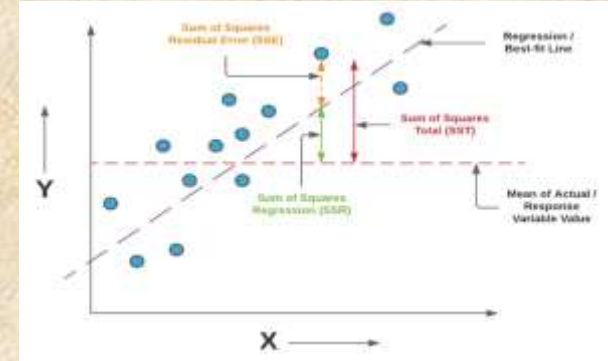


→ **RMSE**:- Root-mean-square Error (RMSE) is used to measure the differences between values predicted by a model and the values observed.

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

# Metric Comparison Features

- **R<sup>2</sup>:** R-squared (R<sup>2</sup>) represents the proportion of the variance for a dependent variable that is predicted from an independent variable.
- **Adjusted R<sup>2</sup>:** Adjusted R<sup>2</sup> measures the proportion of variation explained by only those independent variables that really help in explaining the dependent variable.
- **MAE:** Mean Absolute Error (MAE) is the magnitude of absolute difference between the prediction of an observation and the true value of that observation.



$$R^2 = 1 - \frac{SS_{residuals}}{SS_{total}}$$

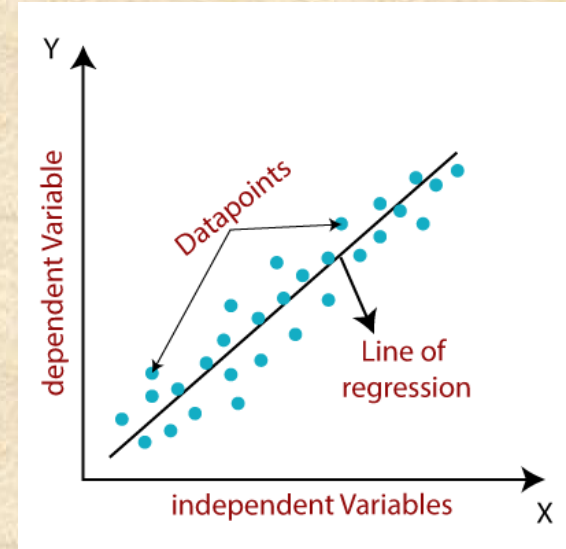
$$\text{Adjusted } R^2 = 1 - \frac{SS_{residuals} / (n - K)}{SS_{total} / (n - 1)}$$

$$MAE = \frac{1}{n} \sum \left| y - \hat{y} \right|$$

# Linear Regression

Linear regression is used to model the relationship between two variables and estimate the value of a response by using a line-of-best-fit.

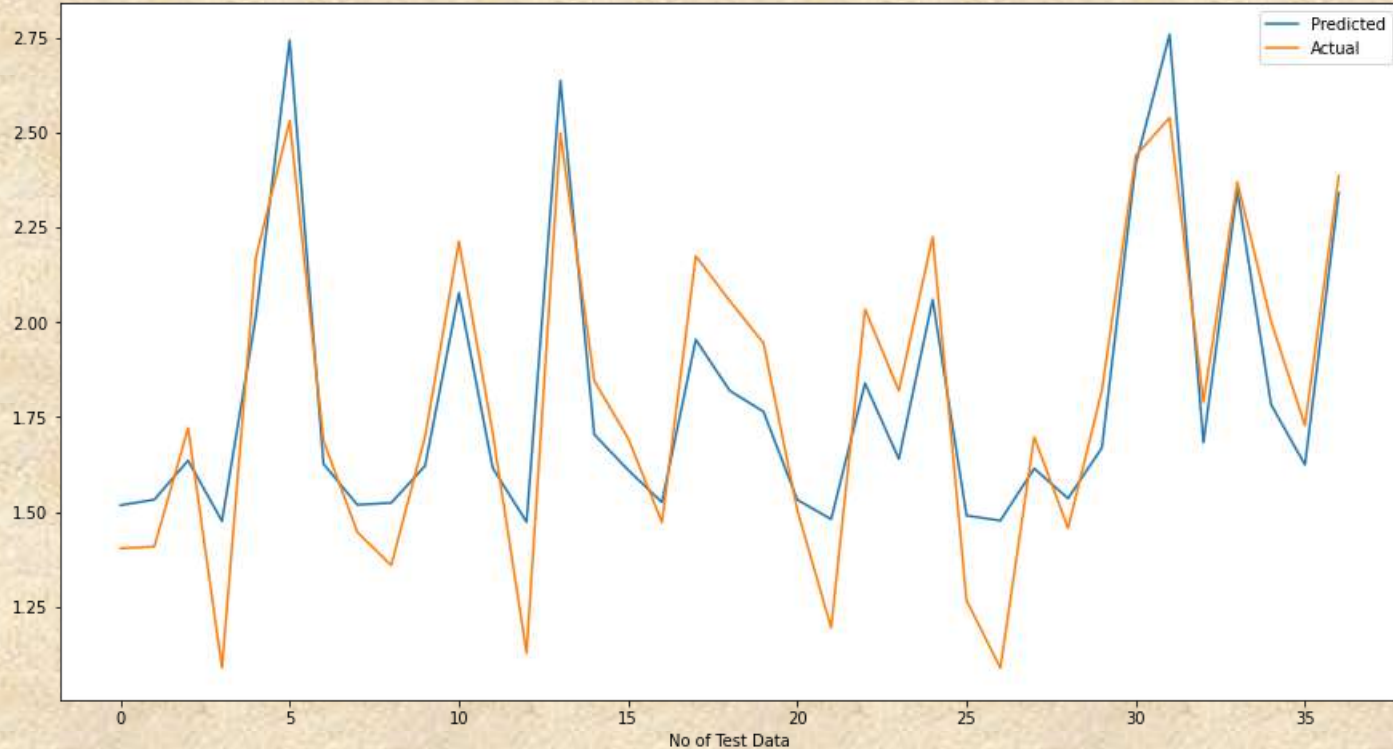
This model is built for simple linear regression, where we study the independent features Open, High and Low and how the dependent feature Close response to it.





# Linear Regression graph

Linear Regression-Actual Vs Predicted



**MSE:-0.031583**

**MAE:-0.151285**

**RMSE:-0.177715**

**R2:-0.822570**

**Adjusted R2:-0.806440**

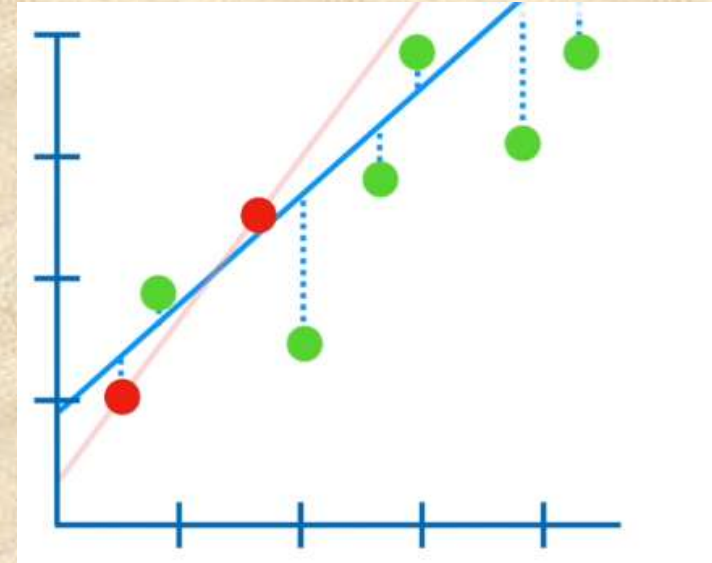


# Ridge Regression

Ridge **regression** is a model tuning method that is used to analyse any data that suffers from multicollinearity.

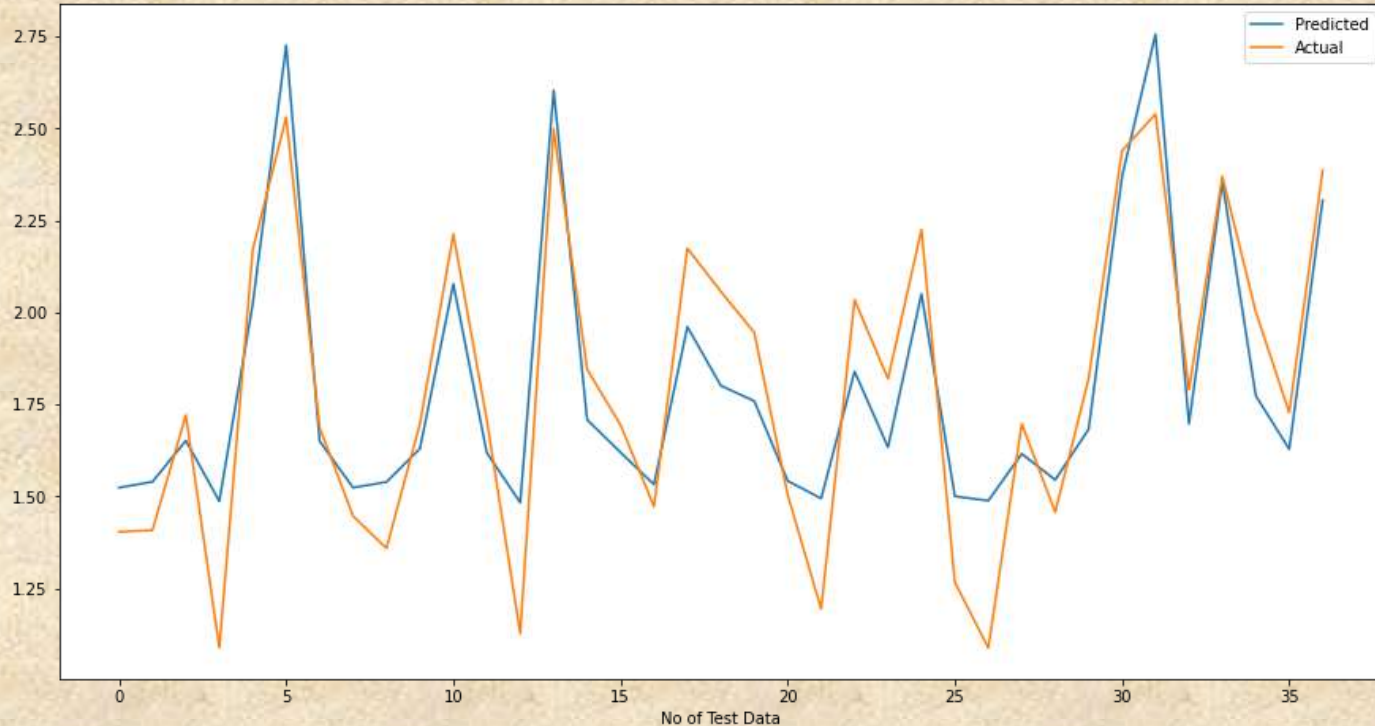
Ridge regression is one of the types of linear regression in which a small amount of bias is introduced so that we can get better long-term predictions.

It is also called as **L2 regularization** as it disperse the error terms in all the weights and leads to more accurate customized final models.



# Ridge Regression graph

Ridge Regression-Actual Vs Predicted



**MSE:- 0.031685**

**MAE:- 0.151477**

**RMSE:- 0.178001**

**R2:- 0.821997**

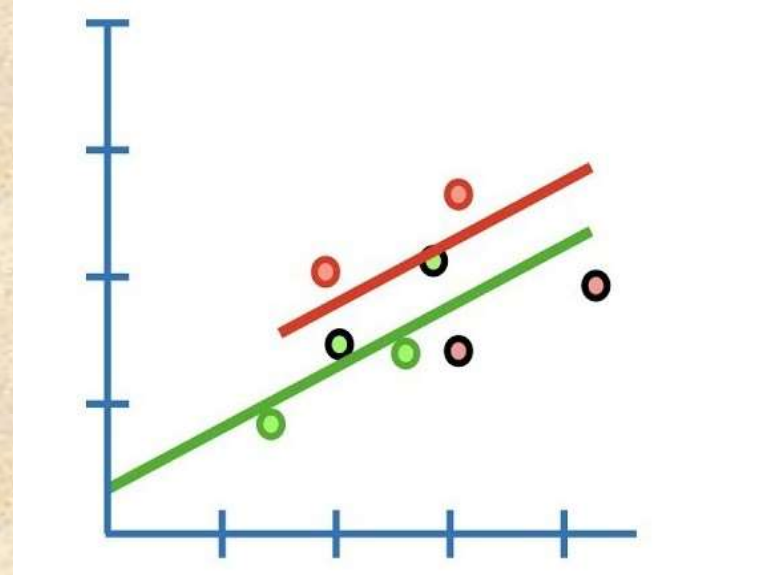
**Adjusted R2:-0.805815**

# Lasso Regression

Lasso Regression (Least Absolute Shrinkage and Selection Operator) is a regression analysis method that performs both variable selection and regularization in order to enhance the prediction accuracy and interpretability of the model.

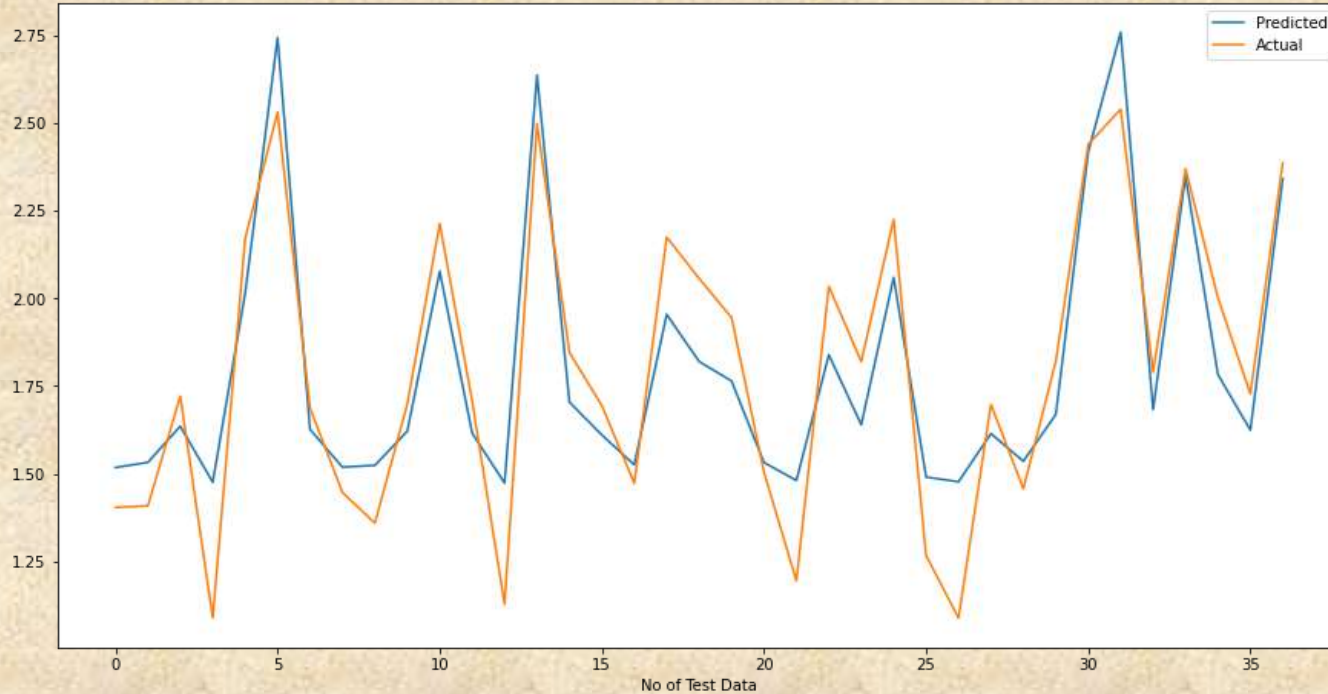
This model uses shrinkage. Shrinkage is where data values are shrunk towards a central point as the mean.

Lasso regression performs **L1 regularization**, i.e. it adds the “absolute value of magnitude” of the coefficient as a penalty term to the loss function.



# Lasso Regression Graph

Lasso Regression-Actual Vs Predicted



**MSE:- 0.032040**

**MAE:- 0.151477**

**RMSE:- 0.178996**

**R2:- 0.820001**

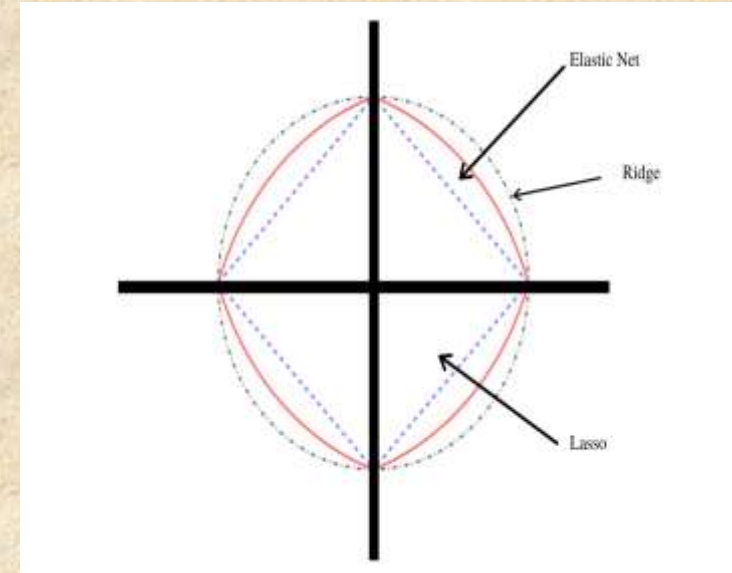
**ADJUSTED R2:- 0.803638**

# Elastic Net Regression

The technique combines both the lasso and ridge regression methods by learning from their shortcomings to improve the regularization of statistical models.

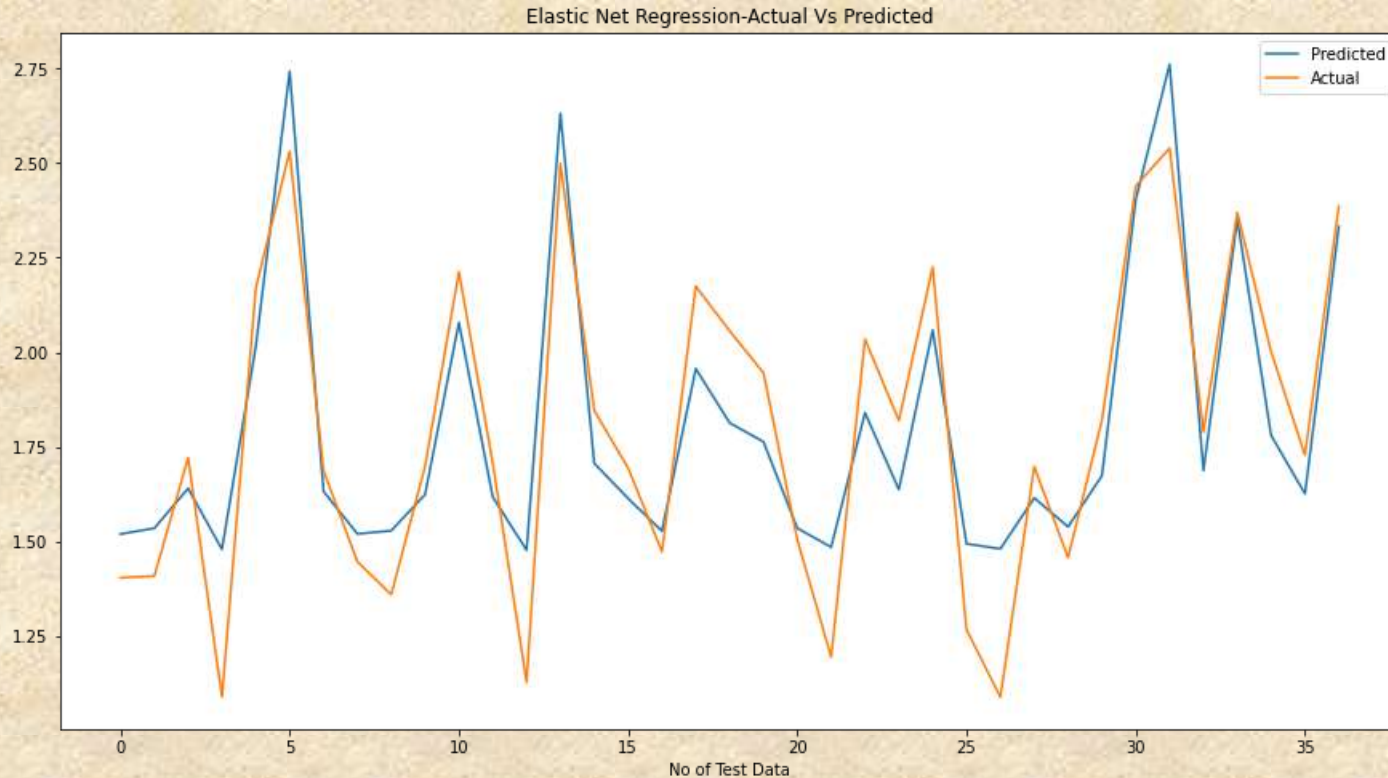
Elastic Net, a convex combination of Ridge and Lasso.

Elastic net linear regression uses the penalties from both the Lasso and Ridge techniques to regularize regression models.





# Elastic Net Regression Graph



**MSE:- 0.152095**

**MAE:- 0.031957**

**RMSE:- 0.178764**

**R2:- 0.820468**

**ADJUSTED R2:- 0.804147**



# Metric Comparison

Model Name	MSE	MAE	RMSE	R2	Adjusted R2
Linear Regression	0.031583	0.151285	0.177715	0.822570	0.806440
Elastic Net Regression	0.031589	0.152095	0.178764	0.820468	0.804147
Lasso Regression	0.032040	0.152317	0.178996	0.820001	0.803638
Ridge Regression	0.032679	0.153352	0.180773	0.816411	0.799721

From the above study we can say that Linear regression has the best R2 and Adjusted R2 score hence Linear Regression is the best model.



Thank  
You