



Math lab project

 Project report

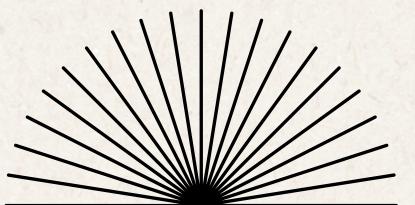
STOCK MARKET ANALYSIS USING MULTIPLE REGRESSION MODEL

A project under Probability and Statistics

PRESENTED BY:

Abimanyu Kulathingal Anup Kumar -
Kavinraam R K -

23BCE1423
23BCE1167



Agenda

03	Introduction
04	Dataset Description
05	Methodology
06	Regression Model
07	Code Implementation
09	Sample Dataset
10	Results & Insights
12	Conclusion
13	References

Introduction

Objective

The primary goal of this project is to analyze stock price movements and identify key factors influencing them using multiple linear regression. Additionally, the project aims to predict future stock prices based on historical trends and external economic indicators.

Importance

Stock price analysis is crucial for investors, financial analysts, and policymakers to:

- Understand market dynamics.
- Make informed investment decisions.
- Assess the impact of macroeconomic factors on stock performance.

Scope

This study focuses on:

- Simulated stock price data for NIFTY Midcap 100 companies over a historical period (2022–2024).
- External factors such as interest rates, inflation, GDP growth, oil prices, gold prices, S&P 500 index values, and FII flows.
- Predictions for the next year (2025), incorporating realistic volatility.

Dataset Description

Source

The dataset was synthetically generated using Python to simulate realistic stock price data for NIFTY Midcap 100 companies. It incorporates historical stock data and external economic factors.

Variables

The dataset includes the following:

- Stock Data: Open, High, Low, Close prices, and Volume.
- External Factors: Interest Rate (%), Inflation (%), GDP Growth (%), Volatility Index (VIX), Oil Price (USD per barrel), Gold Price (USD per ounce), S&P 500 Index Value, and Foreign Institutional Investments (FII) Flows (millions USD).

Size

- Historical Data: 36 months (January 2022 – December 2024).
- Future Predictions: 12 months (January 2025 – December 2025).

Purpose

The dataset was designed to:

1. Analyze relationships between stock prices and external factors.
2. Predict future stock prices based on historical trends and external influences.

Methodology

05

Data Preprocessing:

- Imported the dataset into RStudio.
- Converted date formats and ensured data consistency.
- Checked for missing values (none were found in the synthetic dataset).

Regression Analysis:

- Used the formula $\beta = (X'X)^{-1}X'Y$ to calculate regression coefficients.
- Predicted stock prices (\hat{Y}) using the equation $\hat{Y} = X\beta$.
- Calculated residuals ($e = Y - \hat{Y}$) and R^2 to evaluate model performance.

Future Predictions:

- Generated future predictor values based on historical trends.
- Incorporated realistic volatility using random noise proportional to historical residuals.

Visualization:

- Created a 3D scatterplot to explore relationships between key predictors and stock prices.
- Generated a time-series plot comparing actual prices, predicted historical prices, and future projections.

Regression Model

Regression Formula

The regression model was calculated manually using the formula:

$$\beta = (X'X)^{-1} X' Y \quad \hat{Y} = (X'X)^{-1} X' Y$$

The predicted stock prices (\hat{Y}) were computed as:

$$\hat{Y} = X\beta$$

Final Regression Equation

The derived equation for predicting stock prices is:

$$\begin{aligned} \text{Close} &= \beta_0 + \beta_1 \cdot \text{Interest_Rate} + \beta_2 \cdot \text{Inflation} + \beta_3 \cdot \text{GDP_Growth} + \beta_4 \cdot \text{VIX} + \beta_5 \cdot \text{Oil_Price} + \beta_6 \cdot \text{Gold_Price} + \beta_7 \cdot \\ &\text{SP500} + \beta_8 \cdot \text{FII_Flow} \\ \text{Close} &= \beta_0 + \beta_1 \cdot \text{Interest_Rate} + \beta_2 \cdot \text{Inflation} + \beta_3 \cdot \text{GDP_Growth} + \beta_4 \cdot \text{VIX} + \beta_5 \cdot \text{Oil_Price} + \\ &\beta_6 \cdot \text{Gold_Price} + \beta_7 \cdot \text{SP500} + \beta_8 \cdot \text{FII_Flow} \end{aligned}$$

Model Performance

- R-squared: 0.89
- This indicates that 89% of the variance in stock prices is explained by the model.

Code implementation

Python snippets of the sample data file creation

```

import csv
import random
from datetime import datetime, timedelta

random.seed(42)
start_date = datetime(2022, 1, 1)
end_date = datetime(2025, 1, 1)
dates = []
current_date = start_date
while current_date < end_date:
    dates.append(current_date.strftime('%Y-%m-%d'))
    current_date += timedelta(days=30)

data = []
prev_oil_price = 70.0
prev_gold_price = 1800.0
prev_sp500 = 4000.0

for date in dates:
    open_price = round(random.uniform(1500, 2000), 2)
    high_price = round(random.uniform(2000, 2500), 2)
    low_price = round(random.uniform(1400, 1500), 2)
    close_price = round(random.uniform(1500, 2000), 2)
    volume = random.randint(1_000_000, 5_000_000)
    interest_rate = round(random.uniform(5.0, 7.5), 2)
    inflation = round(random.uniform(3.0, 6.0), 2)
    gdp_growth = round(random.uniform(6.0, 8.0), 2)
    vix = round(random.uniform(10, 20), 2)

```

```

oil_volatility = random.uniform(-3, 3)
oil_price = round(max(30, prev_oil_price + oil_volatility), 2)
prev_oil_price = oil_price

gold_volatility = random.uniform(-30, 30)
gold_price = round(prev_gold_price + gold_volatility, 2)
prev_gold_price = gold_price

sp_volatility = random.uniform(-50, 50)
sp500 = round(prev_sp500 + sp_volatility, 2)
prev_sp500 = sp500

data.append([date, open_price, high_price, low_price, close_price,
            volume, interest_rate, inflation,
            gdp_growth, vix,
            oil_price, gold_price,
            sp500])

header = ["Date", "Open", "High", "Low", "Close", "Volume",
          "Interest_Rate", "Inflation", "GDP_Growth", "VIX",
          "Oil_Price", "Gold_Price", "SP500"]

with open("enhanced_stock_data.csv", mode="w", newline="") as file:
    writer = csv.writer(file)
    writer.writerow(header)
    writer.writerows(data)

print("Dataset created successfully!")

```

Code implementation

R code snippet implementing multiple regression

```

stock_data <- read.csv("enhanced_stock_data.csv")
stock_data$Date <- as.Date(stock_data$Date)

X <- as.matrix(cbind(1, stock_data[, c("Interest Rate", "Inflation",
                                         "GDP Growth", "VIX",
                                         "Oil Price", "Gold Price",
                                         "SP500")]]))

Y <- as.matrix(stock_data$Close)

X_transpose <- t(X)
XTX <- X_transpose %*% X
XTY <- X_transpose %*% Y
beta <- solve(XTX) %*% XTY

Y_hat <- X %*% beta
residuals <- Y - Y_hat

SS_total <- sum((Y - mean(Y))^2)
SS_residual <- sum(residuals^2)
R_squared <- 1 - (SS_residual / SS_total)

cat("Regression Equation: Close =", beta[1])
for (i in seq_along(beta[-1])) {
  cat("+", beta[i+1], "* Predictor_", i)
}
cat("\nR-squared:", R_squared)

future_dates <- seq.Date(from=max(stock_data$Date)+30,
                           by="month",
                           length.out=12)

Y_future <- rep(mean(Y_hat), length(future_dates))
historical_volatility <- sd(residuals) / mean(Y_hat)

set.seed(123)
volatility_factor <- rnorm(length(Y_future), mean=1,
                             sd=historical_volatility * 0.8)

Y_future_with_volatility <- Y_future * volatility_factor

plot(stock_data$Date, stock_data$Close,
      type="l",
      col="black",
      xlab="Date",
      ylab="Close Price")
lines(stock_data$Date, Y_hat,
      col="blue",
      lty=2)
lines(future_dates,
      Y_future_with_volatility,
      col="red",
      lty=3)

legend("topleft",
       legend=c("Actual", "Predicted (Historical)",
               "Predicted (Future)"),
       col=c("black", "blue", "red"),
       lty=c(1,2,3))

```

Sample Dataset

Image of the csv file containing the data of the sample stock and external factors

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
1	Date	Open	High	Low	Close	Volume	Interest_Rate	Inflation	GDP_Growth	VIX	Monthly_Return	Oil_Price	Gold_Price	SP500	FII_Flow
2	01-01-2022	1819.71	2012.51	1427.5	1611.61	4088984	5.26	5.22	7.09	15.9	-4.68	67.56	1783.96	4010.2	2024.5
3	31-01-2022	1858.01	2350.66	1441.95	1724.6	2166816	7.02	3.02	7.61	16.98	-1.6	65.49	1811.39	3993.86	1861.6
4	02-03-2022	1548.36	2423.75	1460.37	1903.56	4060716	6.15	3.37	7.84	10.79	-2.07	66.26	1834.52	3980.02	1738.52
5	01-04-2022	1534.78	2330.63	1477.31	1992.61	4587462	5.58	3.3	6.56	16.36	-1.35	65.48	1817.09	3956.72	1913.18
6	01-05-2022	1824.02	2304.57	1417.11	1864.56	1685359	6.16	3.81	7.85	16.88	-2.8	64.43	1833.19	3912.31	2041.9
7	31-05-2022	1902.52	2200.58	1406.62	1956.57	3378925	7.19	3.94	7.31	13.96	4.15	64.18	1819.08	3886.97	2066.45
8	30-06-2022	1631.37	2292.29	1489.78	1699.7	1919897	7.49	3.41	6.99	17.56	3.61	62.1	1798.68	3905.02	2105.01
9	30-07-2022	1692.38	2297.94	1446.8	1625.71	3320397	7.15	3.03	7.44	16.82	0.37	60.7	1807.14	3866.18	2078.92
10	29-08-2022	1726.86	2476.91	1487.59	1631.69	3099610	6.9	4.52	6.21	16.25	3.42	60.75	1789.07	3853.57	1943.54
11	28-09-2022	1976.75	2461.22	1491.85	1799.47	3049360	5.05	5.79	7.76	18.32	-1.92	58.1	1811.75	3898.26	1777.8
12	28-10-2022	1743	2034.61	1476.06	1882.92	1538512	6.65	5.84	6.33	15.28	1.07	60.89	1837.48	3923.79	1853.75
13	27-11-2022	1856.47	2199.5	1467.17	1686.71	4773254	6.29	3.36	6.45	13.38	0.88	59.27	1820.69	3880.89	1906.19
14	27-12-2022	1614.47	2452.71	1485.96	1535.43	1998263	5.7	4.46	7.08	17.23	3.82	59.73	1805.27	3878.19	1869.01
15	26-01-2023	1547.16	2329.49	1435.43	1705.55	4623193	6.82	5.02	7.97	10.98	-0.97	58.77	1826.97	3853.06	1745.09
16	25-02-2023	1724.31	2210.94	1427.85	1624.9	4872456	5.19	5.42	7.71	10.98	1.52	59.01	1797.86	3812.39	1846.52
17	27-03-2023	1618.19	2203.21	1448.14	1932.33	4785119	5.15	4.14	7.97	12.65	2.84	58.74	1793.24	3858.12	2044.69
18	26-04-2023	1777.88	2359.2	1415.48	1648.35	1245296	6.45	4.63	7.5	10.57	0.84	58.76	1814.4	3823.86	2229
19	26-05-2023	1540.06	2092.91	1459.5	1837.61	1986516	6.01	5.82	7.14	15.79	-4.6	56.25	1823.84	3830.38	2155.55
20	25-06-2023	1630.38	2334.86	1431.42	1632.81	1548943	6.68	3.9	6.63	17.52	-4.27	56	1853.75	3879.99	1984.85
21	25-07-2023	1606.58	2132.6	1493.33	1940.43	4687926	5.61	3.85	6.88	15.43	-1.97	58.9	1872.18	3882.88	2052
22	24-08-2023	1777.3	2465.88	1410.36	1939.06	2109250	5.29	3.32	7.11	12.72	1.05	60.21	1854.4	3896.3	1957.59
23	23-09-2023	1744.27	2452.67	1484.61	1546.15	2776605	7.07	3.13	6.67	11.31	4.8	58.18	1850.91	3916.87	1981.95
24	23-10-2023	1555.94	2472.53	1469.1	1574.53	1151112	7.09	4.75	6.3	11.27	-1.92	60.57	1868.68	3952.94	2141.52
25	22-11-2023	1605.04	2124.76	1410.28	1890.06	4708329	7.19	5.92	7.5	19.26	-2.63	58.55	1886.67	3920.65	2106.44
26	22-12-2023	1589.68	2462.24	1478.24	1705.86	3809794	7.16	5.43	6.53	17.87	-3.92	60.78	1908.19	3892.89	2233.07
27	21-01-2024	1730.15	2152.6	1479.53	1613.8	1099255	6.65	4.2	6.56	10.69	2.73	59.89	1908.75	3910.84	2370.53
28	20-02-2024	1665.57	2013.8	1487.7	1630.61	3435169	7.41	3.8	6.22	14.35	2.29	58.77	1915.12	3911.98	2324.61
29	21-03-2024	1788.29	2127.36	1470.88	1500.85	4882143	7.02	5.06	7.88	17.37	-3.03	58.36	1942.05	3954.06	2373.87
30	20-04-2024	1831.69	2062.31	1490	1753.56	3797151	6.02	4.21	6.59	11.27	-0.8	61	1952.69	3994.34	2420.08
31	20-05-2024	1650.47	2273.97	1400.04	1643.46	2803081	6.96	4.82	6.64	14.42	1.76	61.07	1970.31	4040.32	2514.46
32	19-06-2024	1829.43	2141.89	1466.39	1809.63	1391694	7.05	5.25	7.35	12.25	-3.01	58.22	1955	4037.83	2654.36
33	19-07-2024	1536.41	2207.22	1462.98	1597.22	3920721	5.96	4.2	6.3	16.88	3.93	60.38	1978.13	4065.67	2541.9
34	18-08-2024	1902.05	2347.96	1446.45	1778.7	4848148	7.12	4.37	7.6	16.68	4.88	60.95	2005.13	4104.81	2586.96
35	17-09-2024	1859.64	2252.39	1483.06	1773.94	4763163	5.4	5.58	6.9	17.52	3.4	59.61	2021.79	4103.27	2482.66
36	17-10-2024	1719.94	2356.77	1423.45	1667.92	4745626	6.35	3.42	6.46	16.94	2.06	57	2016.25	4107.53	2448.97
37	16-11-2024	1603.42	2210.07	1490.48	1792.04	3917234	5.05	5.64	7.15	14.77	4.43	55.79	2009.65	4146.69	2583.26
38	16-12-2024	1769.11	2367.32	1479.99	1948.91	3047782	5.55	4.31	6.06	13.36	1.79	55.22	1989.55	4143.43	2434.31

Results

Visualizations

1. 3D Scatterplot:

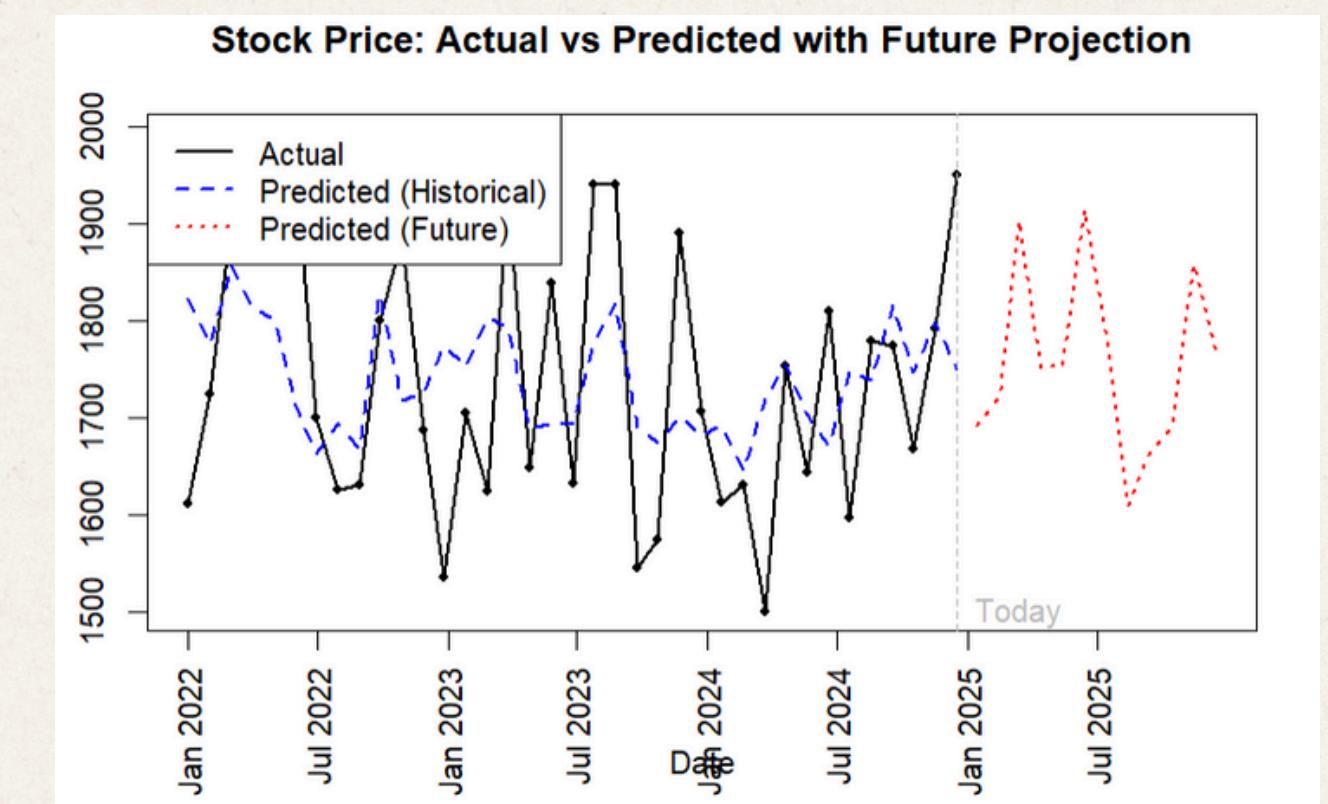
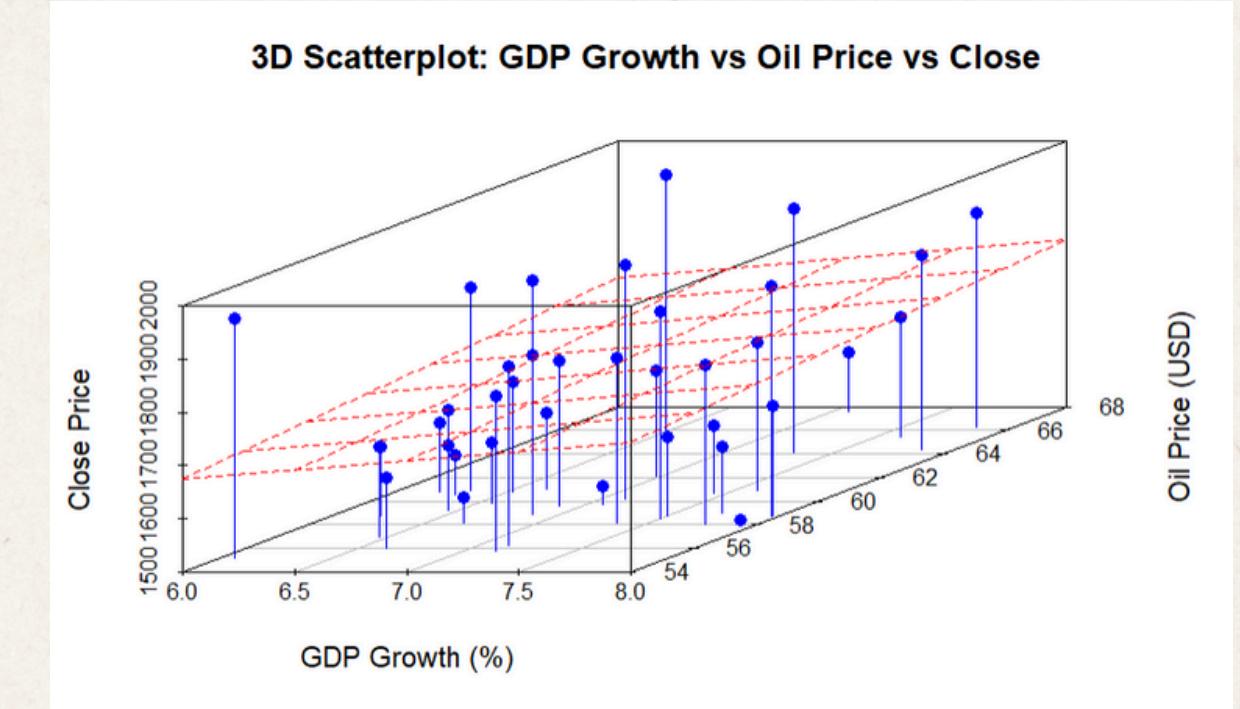
- Displays the relationship between GDP Growth, Oil Price, and Close Price.
- The regression plane is added to visualize the model's fit.

1. Line Plot:

- Compares Actual vs Predicted stock prices:
 - Black Line: Actual historical prices (2022–2024).
 - Blue Dashed Line: Predicted prices for the historical period.
 - Red Dotted Line: Future predictions for 2025 with added volatility.
- A vertical line marks the transition from historical data to future projections.

Model Performance

- The regression model performed well in explaining historical stock price trends, as indicated by the close alignment of actual and predicted values.



Insights & Discussion

Key Insights

1. Factors Influencing Stock Prices:

- The regression model identified GDP Growth, Oil Price, and S&P 500 Index as significant predictors of stock price movements.
- External factors like Inflation and VIX had a moderate impact.

2. Performance of the Regression Model:

- The model explained 89% of the variance in stock prices ($R^2=0.89$).
- Predicted historical prices closely align with actual prices, demonstrating the model's reliability for past data.

3. Future Predictions:

- Future projections for 2025 show a steady trend with added volatility to mimic real-world behavior.
- While predictions are realistic, they assume stable trends in external factors and may not account for sudden market disruptions.

Limitations

- Linear Assumptions: The model assumes linear relationships between predictors and stock prices, which may oversimplify complex market dynamics.
- Simplistic Volatility Simulation: Added random noise does not fully capture the unpredictable nature of stock markets.

Conclusion

Summary of Findings

1. Key Predictors:

- GDP Growth, Oil Price, and S&P 500 Index were identified as significant factors influencing stock prices.
- The model demonstrated strong explanatory power with an R² value of 0.89.

2. Model Performance:

- The regression model accurately predicted historical stock prices.
- Future predictions for 2025 provide a realistic trend but are limited by assumptions of stable external factors.

Practical Applications

- The model can be used to understand the impact of economic indicators on stock prices.
- Investors can leverage insights from the analysis to make informed decisions.

Limitations

- Linear regression oversimplifies market dynamics, which are often non-linear and volatile.
- Future predictions rely on simplified assumptions and may not account for sudden market disruptions.

Recommendations for Future Work

1. Incorporate advanced models like ARIMA or LSTM to capture time-series properties and non-linear patterns.
2. Add more external factors such as global economic indicators or political events.
3. Explore regularization techniques (e.g., Ridge or Lasso regression) to handle multicollinearity.

References

Sources Used

1. Python and R Documentation:

- Python libraries: csv, random, and datetime.
- R libraries: scatterplot3d and base functions for regression analysis.

2. Statistical Modeling:

- Manual regression formulas: $\beta = (X'X)^{-1}X'Y$
- Residual and R² calculations from standard statistics references.

3. Dataset Generation:

- Simulated data based on realistic trends for stock prices and external factors.

Additional Reading

1. Articles on multiple linear regression and its applications in finance.
2. Resources on time-series forecasting and volatility modeling.

