

Analyzing the Stock Market Price of the S&P 500

Rakesh Rathod
Anushree Gupta
Abhishek Mathur
Gautami Mudaliar
Padmasini Krishnan Venkat



Introduction to the S&P 500

The S&P 500 is a stock market index that tracks the performance of 500 large-cap companies listed on US stock exchanges. Its performance is a benchmark for the US stock market.



Objective

Is to develop and evaluate a GARCH model to forecast the volatility of S&P 500 stock market prices.



Methodology of Analysis

Is to calculate the actual volatility, then volatility using GARCH and forecast them. Then, use various tests to figure out the best GARCH model, cross correlation between actual volatility and GARCH volatility etc.



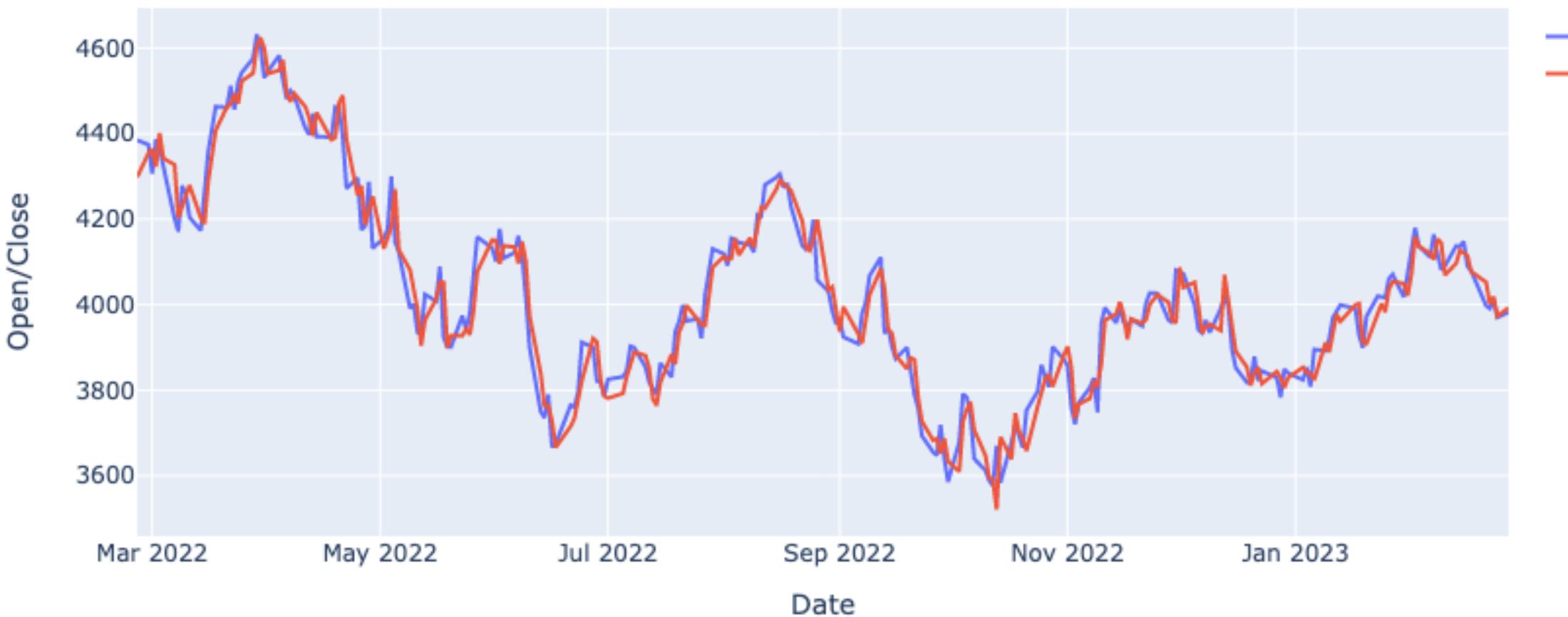


DATASET

Includes the {Date, opening price, closing price, Highest price of the day, Lowest price of the day} of S&P500 for past 10 years.

EXPLORING DATASET

Date vs Open/Close of F.Y. 2022-23



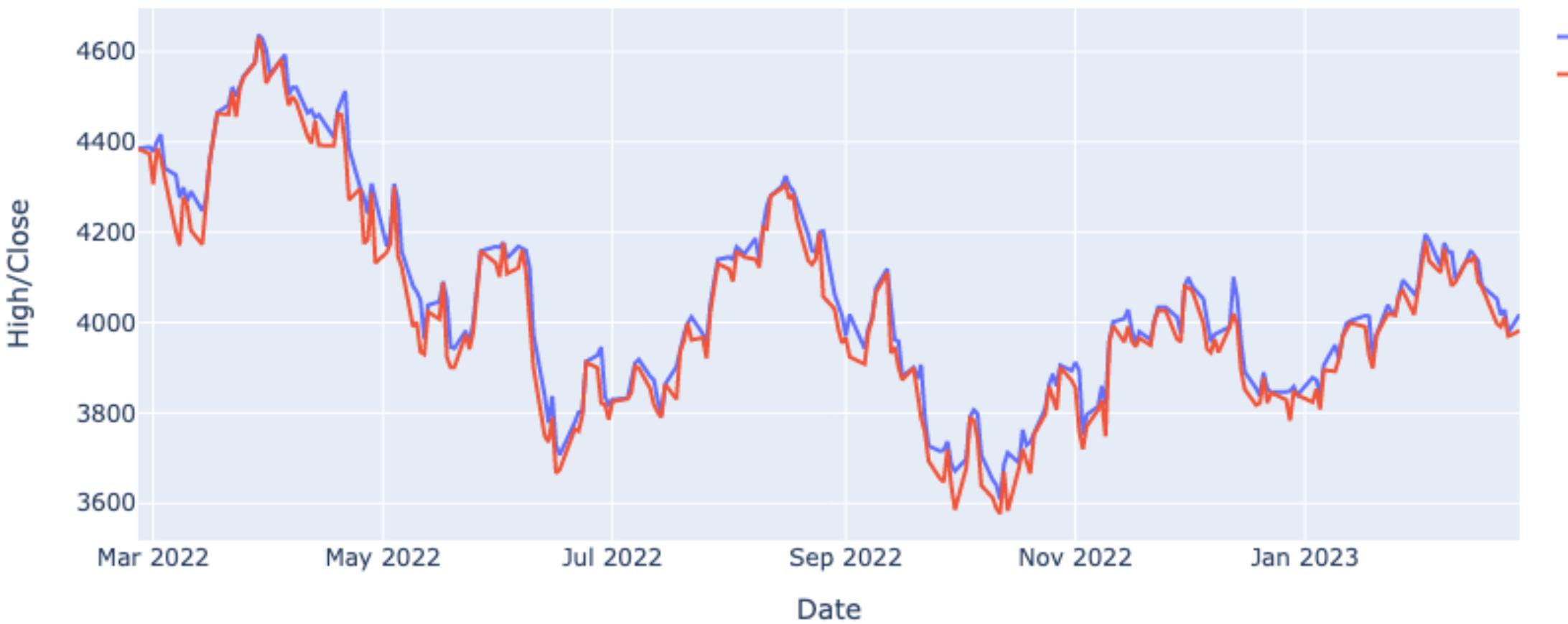
EXPLORING DATASET

Date vs Low/Open of F.Y. 2022-23



EXPLORING DATASET

Date vs High/Close of F.Y. 2022-23



Actual Volatility

can be calculated as follows:

$$AV = \sqrt{n} * \text{stdev}(R)$$

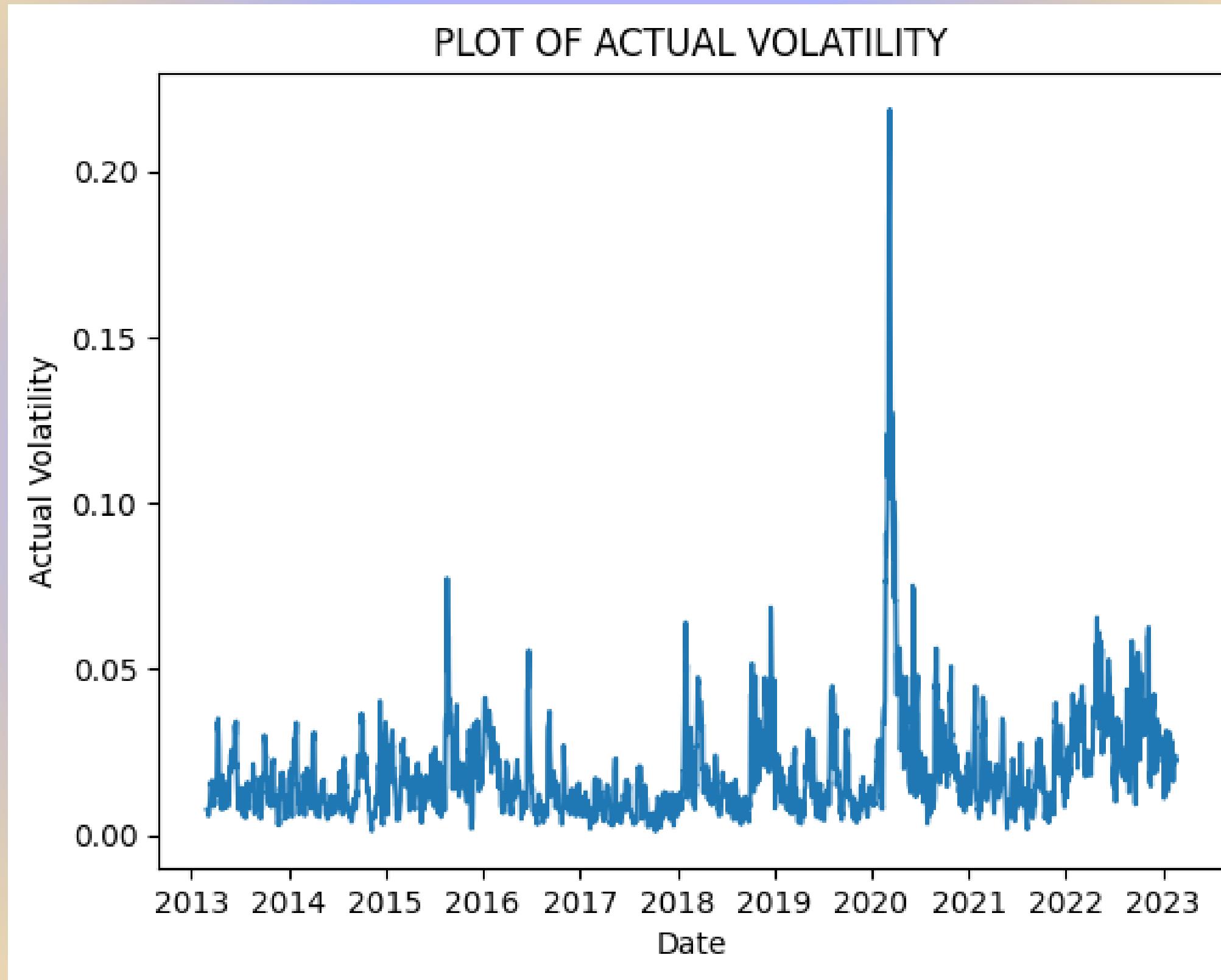
Where:

"n" is the rolling window

" $\text{stdev}(R)$ " is the standard deviation of the returns over the rolling window.



Actual Volatility



GARCH Model

The Generalized Autoregressive Conditional Heteroskedasticity (GARCH) model is a statistical approach used to analyze financial data. It is a time series model that is designed to capture the volatility clustering in financial data, which means that large movements in asset prices tend to be followed by more large movements, and small movements tend to be followed by more small movements.



GARCH Model and its parameters

can be calculated as follows:

$$\sigma_t^2 = \omega + \sum_{i=1}^p \alpha_i \varepsilon_{t-i}^2 + \sum_{j=1}^q \beta_j \sigma_{t-j}^2$$

ε_t is the standardized residual at time t .

σ_t is the conditional standard deviation of the return at time t .

ω , α , and β are parameters to be estimated.



GARCH Model and its assumptions

1. Stationarity: The statistical properties of the data remain constant over time.

Eg: Any major event like covid lockdown might affect stationarity

2. Conditional heteroscedasticity: The conditional variance of the data is time-varying and depends on the past values of the data.

3. Normality: The GARCH model assumes that the distribution of the data is normal, which implies that the data has a bell-shaped distribution.

4. Independence:
there is no correlation between the residuals at different time periods.



An Example of model estimation

Suppose we have daily stock returns for a company over the past 10 days:

0.01%, 0.03%, -0.02%, 0.05%, 0.01%, -0.03%, -0.01%, -0.02%, 0.02%, -0.01%

Estimate α , β and ω by fitting the data to model,
Say, $\omega = 0.0003$, $\alpha = 0.1$, $\beta = 0.8$

conditional variance for the 11th day:

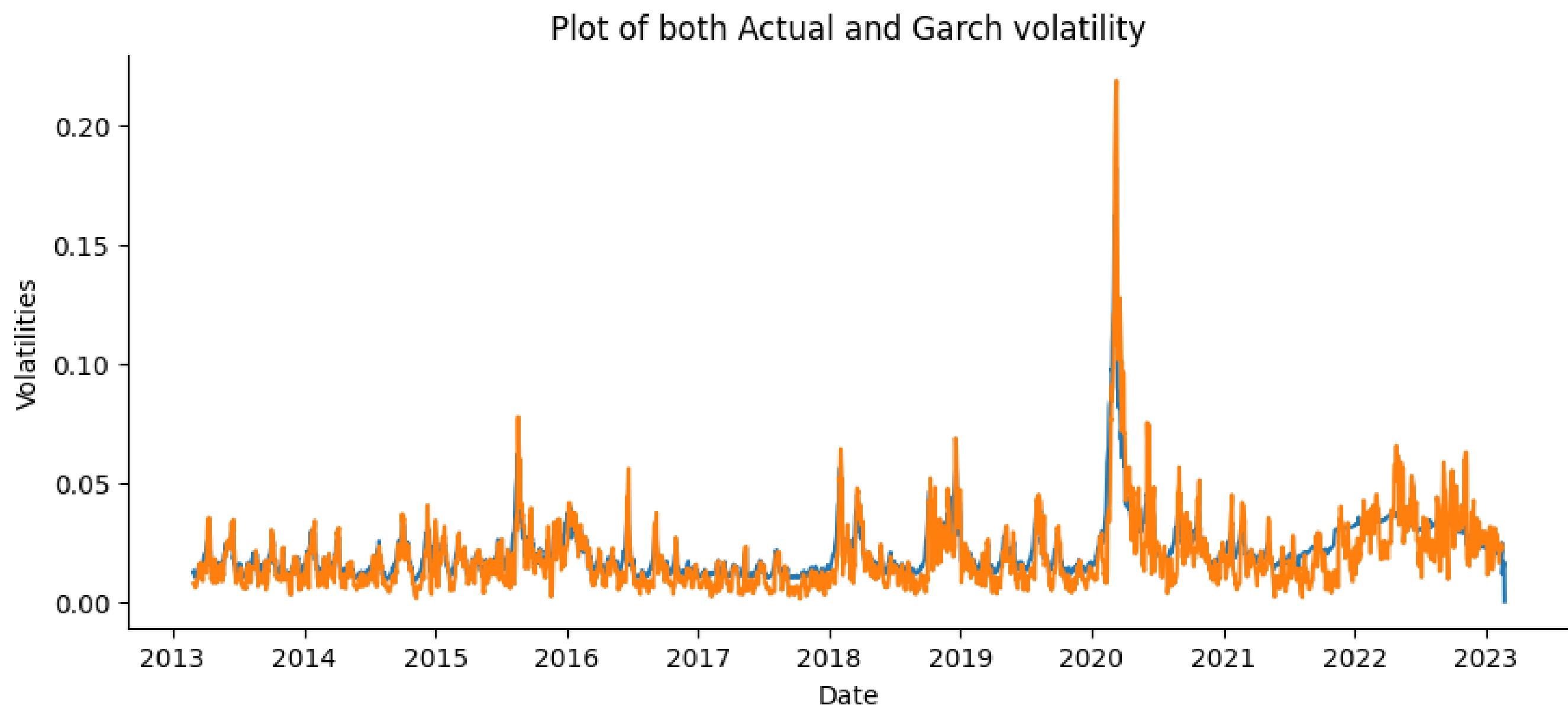
For $p=1, q=1$
 $\sigma_{11}^2 = 0.0003 + (0.1 \times 0.0001) + (0.8 \times 0.0002) = 0.00046$

For $p=2, q=1, \alpha_1 = 0.1, \alpha_2 = 0.2, \beta = 0.7$
 $\sigma_{16}^2 = 0.0003 + (0.1 \times 0.0001) + (0.2 \times 0.0001) + (0.7 \times 0.0002) = 0.00047$

ω = parameter represents the long-term average variance
 α_1, α_2 , and β = parameters represent the short-term and long-term impact of past squared errors and conditional variances on the current conditional variance



Actual and forecasted GARCH Volatility





Tests

1. Residual/Error Analysis -> Ljung-box, Heteroskedasticity
2. Model Evaluation - > AIC, BIC
3. Test on predicted vs actual -> Jarque Bera, Cross correlation



LJUNG BOX TEST

The Ljung-Box test is a statistical test used to determine whether there is evidence of autocorrelation in a time series. It is commonly used in time series analysis to check for autocorrelation in the residuals of a fitted model. If there is evidence of autocorrelation in the residuals, it suggests that the model is not capturing all the information in the data, and additional terms or a different model should be considered.

Results

Lags	lb_stat	lb_pvalue
1	49.736124	1.758777e-12
2	70.191153	5.730399e-16
3	70.562995	3.233686e-15
4	76.945698	7.722403e-16
5	83.606342	1.474958e-16
6	115.022002	1.806598e-22
7	171.684024	1.108268e-33
8	215.521938	3.399261e-42
9	264.741066	7.662683e-52
10	274.312172	4.123589e-53

AIC & BIC TEST

AIC : AIC (Akaike Information Criterion) is a statistical measure used to compare different models that have been fit to the same dataset. It is a tool used in model selection to balance the trade-off between model complexity and goodness of fit. AIC is calculated as follows:

$$AIC = -2\log(L) + 2k$$

where L is the likelihood function of the model, k is the number of parameters in the model, and the negative sign ensures that smaller values of AIC are better.

BIC : BIC (Bayesian Information Criterion) is a statistical measure used to compare different models that have been fit to the same dataset. Like AIC, BIC is a tool used in model selection to balance the trade-off between model complexity and goodness of fit.

BIC is calculated as follows:

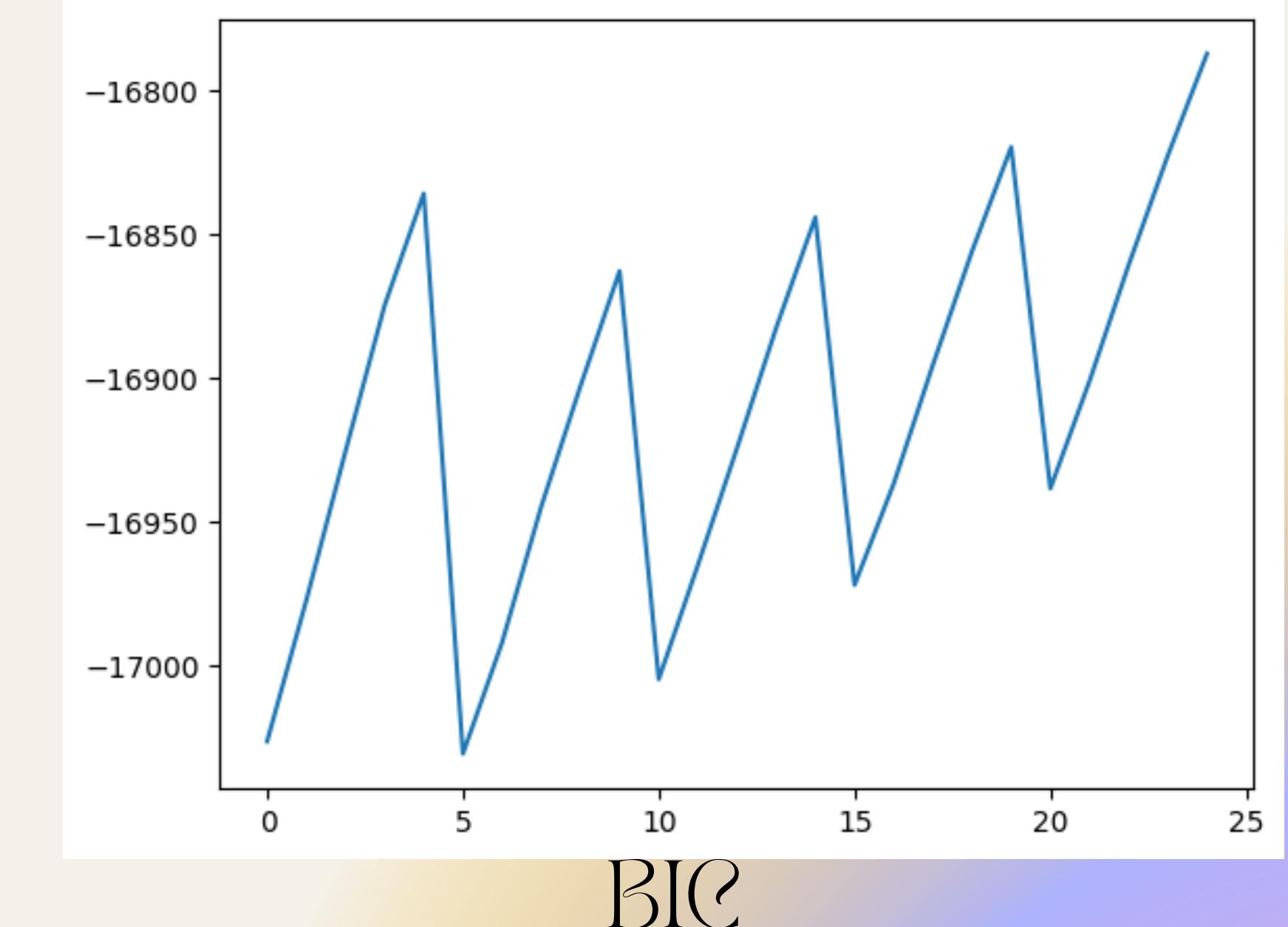
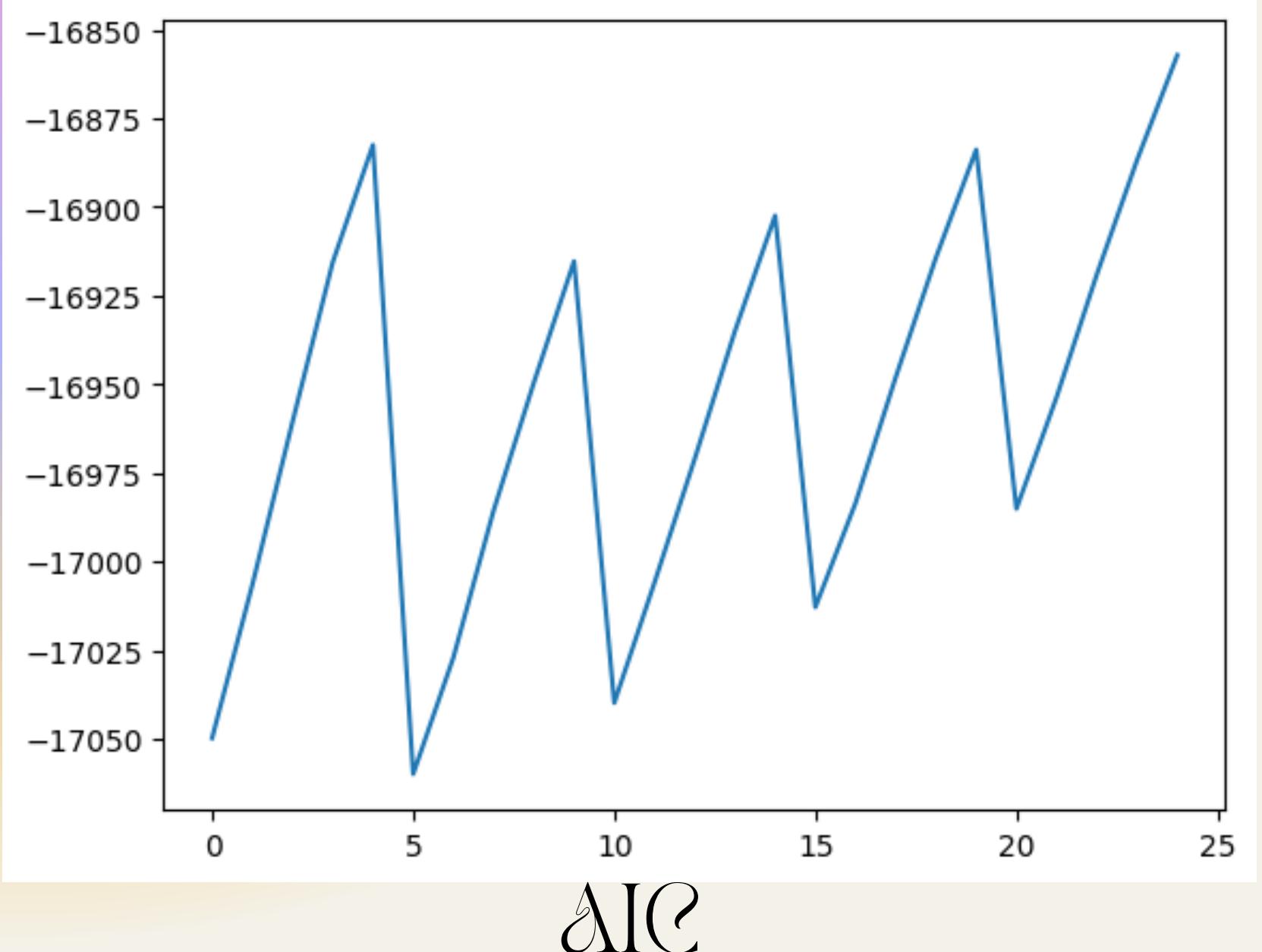
$$BIC = -2\log(L) + k\log(n)$$

n is the number of observations in the dataset

Results

Parameter p	Parameter q	AIC value	BIC value
1	1	-17049.808266838936	-17026.45016451211
1	2	-17006.60403090591	-16977.40640299738
1	3	-16960.590735749363	-16925.55358225913
2	1	-17059.963042949985	-17030.76541504146
2	2	-17027.073008546886	-16992.03585505665
2	3	-16985.59158445888	-16944.71490538695
3	1	-17039.88964620231	-17004.85249271208
3	2	-17005.943840430587	-16965.06716135865
3	3	-16971.02878460011	-16924.31257994647

RESULTS



Based on the results of the ΔIC test, the best value for our parameters is $p=2$ and $q=1$.

JARQUE-BERA TEST



The Jarque-Bera test is a statistical test used to check whether a given sample of data has a normal distribution.

The test is based on the skewness and kurtosis of the data, and the null hypothesis is that the data has a normal distribution.

The test statistic is calculated by:

$$JB = n/6 * (g_1^2 + 1/4 * (g_2 - 3)^2)$$

where n is the sample size.

g_1 and g_2 are the skewness and kurtosis respectively.

If the p-value of the test is less than the chosen significance level, it is concluded that the data does not have a normal distribution. The test is commonly used in finance and economics to assess the normality assumption in financial returns and can be useful in identifying non-normal data distributions in other fields.

Results

Test statistic = 36250.37437102552
p-value = 0.0.

The very small p-value suggests that the data significantly deviate from a normal distribution.



Heteroscedasticity

Heteroscedasticity refers to the variation in the errors of a regression model that is not consistent across different values of the independent variable.

This phenomenon can be problematic in financial analysis when modeling stock prices or other financial data.

To address heteroskedasticity, different methods can be used, such as weighted least squares regression and robust standard errors. These techniques help to adjust the estimates of the model parameters and improve the accuracy of the model.

Results



Test statistic - 965.1415547176754.

P-value - 6.023021775011075e-201.

Degrees of freedom - 155.3993868249614.

Significance level - 3.42092223877066e-254.

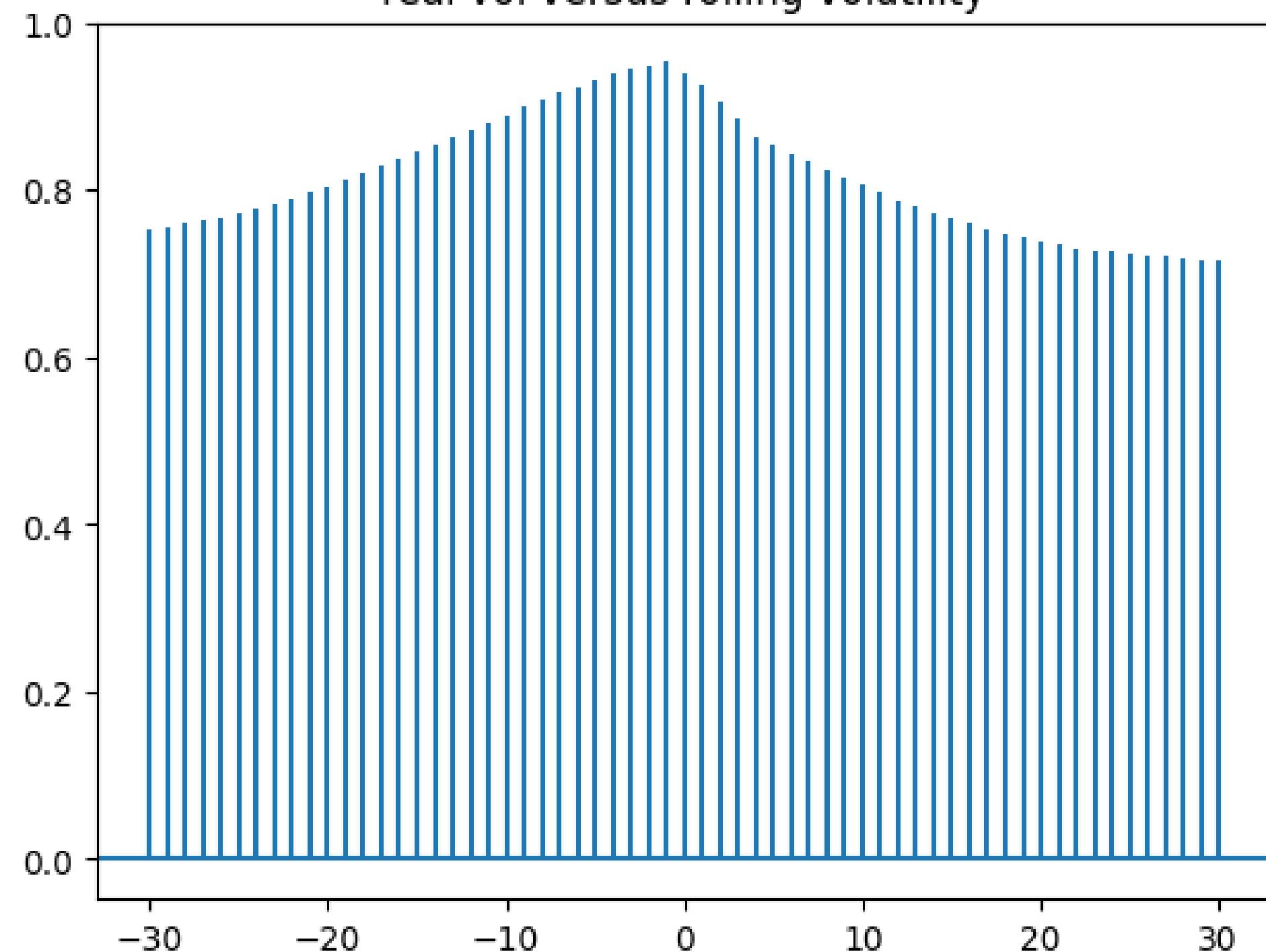
Based on these results, we can conclude that there is strong evidence of heteroscedasticity in the residual data of the stock market price dataset.

CROSS CORRELATION TEST



Cross-correlation is a statistical technique used to measure the correlation between two time series at different time lags. The cross-correlation test is used to determine whether there is a significant correlation between two time series at a given time lag. The test can be used to analyze relationships between financial time series, as well as in other fields to analyze relationships between different time series. The test involves computing the cross-correlation function and testing the significance of the correlation at each time lag using a hypothesis test such as the t-test or F-test.

real vol versus rolling volatility



Results

This plot presents the cross correlation between 'Actual Volatility' and 'forecasted GARCH Volatility' with 0 to 30 lags.

CONCLUSION

Based on the above results, we have come across that the data that the model is fit on isn't very reliable for stock prices prediction since it can lead to various misconceptions and erroneous results.



Thank you