

# Understanding and Implementing the ARIMA Model

## Introduction

This assignment focused on studying and implementing the ARIMA (AutoRegressive Integrated Moving Average) model for time series forecasting using historical stock price data. The primary aim was not only to obtain forecasts, but also to understand the assumptions behind the model, the reasoning involved in choosing its parameters, and the practical challenges that arise when applying statistical models to real-world data. Through this exercise, I gained a clearer understanding of how theoretical concepts such as stationarity, autocorrelation, and residual analysis translate into concrete steps during model implementation.

## Overview of the ARIMA Model

ARIMA is a classical statistical approach designed for univariate time series forecasting, where future values are predicted based solely on past observations of the same variable. The model combines three distinct components:

- **AutoRegressive (AR)** terms, which describe how the current value of the series depends on its own previous values. This captures persistence or momentum in the data.
- **Integrated (I)** terms, which refer to differencing operations applied to the data in order to remove trends and make the series stationary.
- **Moving Average (MA)** terms, which model the influence of past forecast errors on current observations.

These components are summarized using the notation ARIMA(p, d, q). Choosing appropriate values of p, d, and q is crucial, as overly complex models can overfit the data, while overly simple models may fail to capture important dynamics. This assignment emphasized the importance of balancing interpretability, simplicity, and predictive performance.

## Key Concepts and Their Interpretation

### Stationarity

Stationarity is a core assumption underlying ARIMA models. A stationary time series has statistical properties—such as mean, variance, and autocorrelation—that do not change over time. When the original stock price series was plotted, it clearly showed an upward trend and

changing variability, which are typical features of financial price data. This visual observation suggested non-stationarity.

To confirm this statistically, the Augmented Dickey-Fuller (ADF) test was applied. The resulting p-value indicated that the null hypothesis of non-stationarity could not be rejected. This reinforced the conclusion that the raw series was unsuitable for direct ARIMA modeling and required transformation.

## Differencing

To address non-stationarity, first-order differencing was applied to the time series. Differencing transforms the data by computing changes between consecutive observations rather than using absolute values. This process effectively removes long-term trends and stabilizes the mean of the series.

After differencing, the ADF test was performed again, and the p-value indicated that the differenced series was stationary. This confirmed that a single level of differencing ( $d = 1$ ) was sufficient. This step highlighted an important practical insight: even simple transformations can significantly change the statistical properties of a time series and determine whether a model is valid or invalid.

## ACF and PACF Analysis

The AutoCorrelation Function (ACF) and Partial AutoCorrelation Function (PACF) plots were used to guide the selection of the AR and MA parameters. These plots visually represent how current values relate to past values at different lags.

The PACF plot showed a sharp drop after the first lag, suggesting that the direct influence of earlier observations diminishes quickly, which is characteristic of an AR(1) process. Similarly, the ACF plot exhibited a significant spike at lag one followed by a rapid decline, indicating the presence of an MA(1) component. Based on these observations, an ARIMA(1,1,1) model was selected. This choice was guided by data-driven diagnostics rather than arbitrary selection.

## Implementation Experience and Learnings

One of the most important lessons from the coding exercise was the importance of careful data preprocessing. Even small mistakes—such as incorrect file loading or improperly handled missing values—can lead to misleading errors or invalid results. Ensuring that the time index was correctly parsed and that missing values were addressed was essential before applying any statistical tests or models.

Another key learning was the value of following a structured modeling pipeline. Each step—visualization, stationarity testing, differencing, parameter selection, and evaluation—served a specific purpose. Skipping or misordering these steps would have made

the results difficult to interpret. This process helped demystify ARIMA and reinforced that it is not a black-box method, but rather a model grounded in well-defined statistical assumptions.

## Interpretation of Plots and Results

Several diagnostic and evaluation plots were generated during the analysis:

- **Time Series Plot:** The original plot revealed a clear upward trend, supporting the conclusion that the series was non-stationary and unsuitable for direct modeling.
- **Forecast vs Actual Plot:** The forecasted values generally followed the direction of the test data, indicating that the model captured the underlying trend. However, during periods of sharp price movements, the forecasts lagged behind, highlighting the limitations of linear models in volatile markets.
- **Residual Plot:** The residuals fluctuated randomly around zero, suggesting that most systematic patterns in the data were captured by the model.
- **Residual Density Plot:** The distribution of residuals was centered near zero, indicating that the model's forecasts were largely unbiased.
- **ACF and PACF of the Differenced Series:** These plots showed minimal remaining autocorrelation, providing evidence that the chosen ARIMA order was appropriate.

## Evaluation and Model Performance

Model performance was evaluated using standard error metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE). These metrics provided a quantitative measure of forecast accuracy. While the error values were not negligible, they were reasonable given the inherent volatility of stock price data.

Additionally, walk-forward validation was performed to simulate real-time forecasting. This approach demonstrated that the model's performance remained relatively stable over time, suggesting that it was not overly sensitive to specific segments of the dataset.

## Limitations and Scope for Improvement

Despite its usefulness, the ARIMA model has inherent limitations. It assumes linear relationships and relies solely on past values of the series, making it unsuitable for capturing nonlinear dynamics, regime changes, or sudden market shocks. Furthermore, it does not incorporate external information such as macroeconomic indicators, company news, or market sentiment.

Future improvements could include extending the model to **SARIMA** to capture seasonal patterns or using **ARIMAX** to incorporate exogenous variables. Comparing ARIMA with machine learning approaches, such as recurrent neural networks or LSTM models, could also provide valuable insights into its relative strengths and weaknesses.

## **Conclusion**

This assignment provided a thorough understanding of the ARIMA model from both a theoretical and practical perspective. By working through data preprocessing, stationarity testing, parameter selection, model fitting, and evaluation, I developed a clearer appreciation of how time series forecasting models are constructed and assessed. While ARIMA has limitations, it remains a strong foundational model and an important benchmark for more advanced forecasting techniques.