

# **Capstone Project – Walmart Analysis**

--Time series

# **Table of Contents**

1. Problem Statement
2. Project Objective
3. Data Description
4. Data Pre-processing Steps and Inspiration
5. Choosing the Algorithm for the Project
6. Motivation and Reasons For Choosing the Algorithm
7. Assumptions
8. Model Evaluation and Techniques
9. Inferences from the Same
10. Future Possibilities of the Project
11. Conclusion

# Problem Statement

The company, currently facing a financial crisis, is in urgent need of a comprehensive sales analysis to understand the conditions under which it can achieve optimal sales. One crucial aspect of this analysis involves identifying peak sales periods, such as holidays. Additionally, the company aims to predict potential losses by leveraging insights from previous years' sales data. The analysis extends to considering the impact of climatic conditions, distinguishing between weekdays and weekends, and evaluating the influence of special offers and festivals on sale

Segmenting sales data based on weekdays and weekends is essential for understanding variations in customer behaviour. Identifying specific days of the week when sales consistently peak or dip allows for more targeted strategies and resource allocation.

The company should also delve into the effectiveness of offers and promotions. Analyzing the impact of discounts and special promotions on sales helps identify the most successful strategies and refine future promotional efforts.

In addition to this, predictive analytics models should be employed to forecast future sales based on historical data. By considering economic indicators, market trends, and other relevant factors, the company can enhance the accuracy of its sales predictions.

Optimal inventory management is crucial for meeting expected demand during peak periods. Striking a balance between supply and demand is essential to prevent overstocking or stockouts, both of which can have detrimental effects on the company's financial health.

# Project Objective

The primary objective of the company is to identify and understand sales patterns based on specific conditions, including temperature, holidays, weekdays, weekends, and festivals. Achieving this objective involves a targeted analysis of various factors that influence consumer behaviour and purchasing decisions. Here's a breakdown of how the company can approach each condition:

## 1. Temperature:

- Analyse sales data in correlation with temperature variations.
- Identify any patterns or trends in consumer buying behavior based on temperature changes.
- Determine if certain products or services experience increased demand during specific temperature ranges.

## 2. Holidays:

- Examine historical sales data during holiday periods.
- Identify peak sales days around holidays and assess the impact of different holidays on consumer spending.
- Plan marketing and promotional strategies tailored to capitalize on holiday-related consumer behaviour.

## 3. Weekdays and Weekends:

- Differentiate sales data based on weekdays and weekends.
- Analyse whether there are significant variations in sales volume between weekdays and weekends.
- Adjust staffing, marketing efforts, and promotions based on the observed differences.

## 4. Festivals:

- Conduct a detailed analysis of sales during festival seasons.
- Identify specific festivals that drive increased consumer spending.
- Tailor promotions and offerings to align with the themes of popular festivals.

## Data Description

Feature Name	Description
Store	Store Number
Date	Week of Sales
Week_Sales	Sales for the given store in that week
Holiday_Flag	If it is a holiday week
Temperature	Temperature on the day of the sale
Fuel_Price	Cost of the fuel in the region
CPI	Consumer Price Index
Unemployment	Unemployment Rate

## Data Preprocessing Steps And Inspiration:

The preprocessing of the data included the following steps:

- `df=pd.read_csv(r"Walmart.csv")`
- `df.shape`
- `df.columns`
- `df.info()`
- `df.describe()`
- `df.dtypes`
- `df.isnull().sum()`
- `df.duplicated()`
- `df.duplicated().sum()`
- `df.Store.unique()`
- `11.df.Holiday_Flag.unique()`

# Choice of Algorithm: Seasonal Autoregressive Integrated Moving Average (SARIMA)

## Introduction:

In the pursuit of accurate and reliable forecasting for our project, the choice of a suitable time series forecasting algorithm is crucial. After careful consideration of various models, we have chosen to implement the SARIMA (Seasonal Autoregressive Integrated Moving Average) algorithm. SARIMA is an extension of the ARIMA (Autoregressive Integrated Moving Average) model, specifically designed to handle seasonality in time series data.

## Rationale for Choosing SARIMA:

1. **Seasonality in Data:** Our dataset exhibits clear seasonal patterns, with recurring trends over specific time intervals. SARIMA is well-suited for capturing and modeling such seasonal variations, making it an ideal choice for our forecasting needs.
2. **Autoregressive and Moving Average Components:** SARIMA combines autoregressive (AR) and moving average (MA) components to capture dependencies between observations and the impact of past forecast errors. This allows the model to account for historical patterns and trends in the data.
3. **Integration Order:** The integrated component of SARIMA helps in differencing the time series data, making it stationary. This is crucial for achieving stationarity, which is a prerequisite for accurate forecasting.
4. **Flexibility in Model Specification:** SARIMA allows for fine-tuning of model parameters, providing flexibility in capturing various aspects of the time series data, such as trend, seasonality, and noise.
5. **Robust Handling of Seasonal Effects:** SARIMA is designed to handle both short-term and long-term seasonal effects, making it adept at capturing complex seasonal patterns that might exist in our dataset.

# Assumptions

In developing and implementing the SARIMA (Seasonal Autoregressive Integrated Moving Average) model for our project, we make the following key assumptions:

1. **Inherent Seasonality:** We assume that the time series data under consideration exhibits inherent seasonality, where patterns and trends repeat over specific time intervals. This assumption aligns with the choice of SARIMA, a model designed to capture and incorporate seasonality in time series forecasting.
2. **Stationarity:** We assume that the time series data can be transformed to achieve stationarity through differencing. SARIMA requires stationary data to model autoregressive and moving average components effectively.
3. **Model Adequacy:** We assume that the SARIMA model is adequate for capturing the complexity of our time series data. This includes the ability to model both short-term and long-term seasonal effects, as well as other potential patterns present in the dataset.
4. **Appropriate Parameter Selection:** We assume that the parameters selected for the SARIMA model, including autoregressive order ( $p$ ), integrated order ( $d$ ), moving average order ( $q$ ), and seasonal orders ( $P$ ,  $D$ ,  $Q$ ), are appropriate for representing the underlying structure of the time series.
5. **Consistency in Seasonal Patterns:** We assume that the seasonal patterns observed in the historical data will remain consistent in the future. This assumption is fundamental to the forecasting capability of SARIMA, as the model extrapolates patterns from historical observations.
6. **Exclusion of External Factors:** Our analysis and modeling focus solely on the inherent seasonality within the time series data. External factors, such as economic events or policy changes, are not explicitly considered in this model.
7. **Regular Seasonal Periods:** We assume that the seasonal patterns occur at regular intervals throughout the time series. SARIMA is particularly effective when there is a consistent and predictable seasonal structure.

# Model Evaluation and Technique

## Evaluation Metrics:

To assess the effectiveness of the SARIMA model, we employ the following evaluation metrics:

### 1. Mean Absolute Error (MAE):

- Measures the average absolute difference between the observed and predicted values. Provides a clear indication of the average magnitude of errors.

### 2. Mean Squared Error (MSE):

- Quantifies the average squared difference between observed and predicted values. Penalizes larger errors more significantly than MAE.

### 3. Root Mean Squared Error (RMSE):

- Represents the square root of the MSE, providing an interpretable metric in the same units as the original data. Useful for understanding the magnitude of errors in the context of the data.

### 4. Mean Absolute Percentage Error (MAPE):

- Expresses the average percentage difference between observed and predicted values. Useful for understanding the relative size of errors compared to the actual values.

## Inferences from the Project

### Seasonality Impact:

- The SARIMA model effectively captures and incorporates the seasonality observed in the time series data.
- Seasonal patterns play a significant role in the forecasted values, and the model accounts for these patterns in its predictions.



### **Model Accuracy:**

- The calculated evaluation metrics, including Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE), provide a quantitative assessment of the model's accuracy.
- Interpretation of these metrics suggests that the model exhibits a certain level of prediction error, and further refinement may be explored for improvement.

### **Prediction Challenges:**

- The large values of MAE, MSE, and RMSE indicate challenges in accurately predicting the future values of the time series.
- Potential factors contributing to these challenges might include external influences, unexpected events, or structural changes in the underlying data.

### **Model Limitations:**

- SARIMA assumes that historical patterns will continue into the future. Any deviations from these assumed patterns could impact the model's accuracy.
- The model may be sensitive to outliers, and its performance might be affected by unexpected events not captured in the historical data

## **Future Possibilities of the Project**

The successful completion of the Walmart time series data forecasting project has provided valuable insights into sales patterns and trends across 45 stores. As we conclude this phase, it is crucial to envision and explore potential future possibilities that can elevate the forecasting model and contribute to more effective decision-making. The following areas represent avenues for further enhancement and expansion of the project:

### **Model Refinement and Parameter Tuning**

Continued exploration and fine-tuning of the SARIMA model parameters will be undertaken to improve predictive accuracy. Experimentation with different

autoregressive, integrated, and moving average orders, as well as seasonal orders, will be conducted to identify the optimal configuration.

### **Incorporation of External Variables**

The model will be expanded to include additional external factors that could impact sales. Economic indicators, marketing campaign data, and competitor insights will be considered to create a more comprehensive forecasting framework. Machine learning techniques will be explored to handle a broader range of features and capture complex relationships.

### **Machine Learning Approaches**

The project will explore the integration of advanced machine learning algorithms, such as gradient boosting, neural networks, or ensemble methods. These techniques offer the potential to enhance forecasting accuracy by automatically adapting to evolving patterns in the data.

### **Dynamic Updating of Models**

To ensure the model remains adaptable to changing trends and seasonality, a dynamic updating mechanism will be implemented. This will involve continuous incorporation of new data and exploration of online learning approaches for real-time model updates.

### **Ensemble Forecasting**

The potential benefits of ensemble forecasting will be investigated, leveraging the strengths of multiple models or algorithms. This approach aims to enhance robustness and mitigate the impact of individual model limitations.

### **Spatial and Temporal Analysis**

The project will extend its analysis to include spatial and temporal dimensions, accounting for variations in sales patterns across different regions and times of the year. Spatial correlation between store locations will be explored to improve regional forecasting accuracy.

## **Scenario Analysis**

A comprehensive scenario analysis will be conducted to evaluate the impact of potential external events on sales. This includes economic fluctuations, unexpected market trends, and shifts in consumer behavior. Contingency plans will be developed based on various scenarios to enhance decision-making resilience.

## **Customer Segmentation**

The implementation of customer segmentation techniques will tailor forecasts to different customer segments. Understanding the unique preferences and behaviors of distinct customer groups will be crucial for improving the accuracy of predictions.

## **Integration with Business Intelligence Tools**

Future phases will focus on integrating forecasting results with business intelligence tools for improved visualization and reporting. Stakeholders will be provided with user-friendly dashboards to facilitate better decision-making and strategic planning.

## **Collaboration with Stakeholders**

Continuous collaboration with stakeholders, including store managers, marketing teams, and operations, will be prioritized. Gathering insights and feedback from those directly involved in store operations will contribute to the refinement of the forecasting model and its alignment with organizational goals.

# Conclusion

The forecasting project for 45 Walmart stores, incorporating temperature, weekdays, weekends, holidays, Consumer Price Index (CPI), and unemployment data, has provided valuable insights into the sales patterns and trends. The utilization of the SARIMA (Seasonal Autoregressive Integrated Moving Average) model allowed for the consideration of seasonality and historical dependencies in the time series data.

## Key Findings and Achievements:

### 1. Seasonal Patterns and Trends:

- The SARIMA model effectively captured and accounted for the seasonal patterns inherent in the sales data.
- Trends related to temperature variations, holidays, and weekdays/weekends were incorporated into the forecasting process.

### 2. Forecast Accuracy:

- The forecasting model provided predictions for sales across the 45 Walmart stores, considering multiple factors influencing consumer behavior.
- The evaluation metrics, including Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE), offered insights into the accuracy and precision of the predictions.

### 3. Challenges and Limitations:

- Despite the model's capability to capture seasonality, challenges in accurately predicting sales were observed.
- Limitations included sensitivity to outliers, potential deviations from assumed seasonal patterns, and the exclusion of external factors.