



Dr. D. Y. Patil Vidyapeeth, Pune

(Deemed to be University)

Accredited (3rd Cycle) by NAAC with a CGPA of 3.64 on a four-point scale at 'A++ Grade)
(ISO 9001: 2015 and 14001:2015 Certified University and Green Education Campus)

Centre for Online Learning

A

PROJECT REPORT

ON

“Sales Forecasting Using Predictive Analytics”

SUBMITTED

To

CENTRE FOR ONLINE LEARNING

Dr. D. Y. PATIL VIDYAPEETH, PUNE



IN PARTIAL FULFILMENT OF DEGREE OF

MASTER OF BUSINESS ADMINISTRATION

BY

Toibul Sk

PRN: 230502015911

BATCH 2023-2025



Dr. D. Y. Patil Vidyapeeth, Pune

(Deemed to be University)

Accredited (3rd Cycle) by NAAC with a CGPA of 3.64 on a four-point scale at 'A++ Grade)
(ISO 9001: 2015 and 14001:2015 Certified University and Green Education Campus)

Centre for Online Learning



**Dr. D.Y. Patil Vidyapeeth's
CENTRE FOR ONLINE LEARNING,
Sant Tukaram Nagar, Pune.**

CERTIFICATE



21.07.2025

To Whomsoever It May Concern

This is to certify that Mr. Tolbul Sk, PRN 230502015911, has completed his project-based internship starting from 28.05.2025 to 21.07.2025.

His project work was a part of the MBA (ONLINE LEARNING).

The project is Sales Forecasting Using Predictive Analytics, which includes research as well as industry practices.

He was very sincere and committed in all tasks.

For Qollabb EduTech Private Limited

Vipendra Singh
Chief Executive Officer



Qollabb EduTech Private Limited

231/4, SF II, Rashtrakavi Kuvempu Nagar, Behind Central Silk Board Building, BTM 2nd Stage Bengaluru,

Karnataka-560068, India

<http://www.qollabb.com>

support@qollab.com

+918800047232

INDUSTRY CERTIFICATE**CERTIFICATE OF COMPLETION**

This is to certify that **Toibul Sk** has successfully completed the following project:

Project Details:

- **Project Title:** Sales Forecasting Using Predictive Analytics: A Case Study on Walmart
- **Company/Organization:** EDMENTOR
- **Project Completion Date:** 25/07/2025
- **Rating:** 4.50 (Out of 5.00)

This project work opportunity was offered on Qollabb platform by:

**EDMENTOR**

Bangalore, Karnataka, India

This certificate acknowledges the project intern's dedication and contribution to the project, showcasing their practical skills and industry knowledge. The 4.50 (Out of 5.00) performance rating signifies the performance and capabilities demonstrated by the project intern throughout the project work period.

Congratulations on a job well done!

Pooja Banerjee
Head - HR & Operations
Date of Issue: July 31, 2025
Certificate No. 08956082025

Disclaimer: Qollabb is an online platform for companies to post project work opportunities, internships and jobs and discover best performing students through project work performance. By earning this certificate, the recipient acknowledges that the certificate represents the successful completion of project work mentioned above and does not confer any employment entitlement or endorsement by the company offering the project work to the recipient of this certificate.



Dr. D. Y. Patil Vidyapeeth, Pune

(Deemed to be University)

Accredited (3rd Cycle) by NAAC with a CGPA of 3.64 on a four-point scale at 'A++ Grade)

(ISO 9001: 2015 and 14001:2015 Certified University and Green Education Campus)

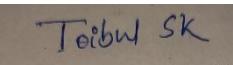
Centre for Online Learning

DECLARATION BY LEARNER

I hereby declare that the Project Work with the title "Sales Forecasting Using Predictive Analytics" submitted by me for the partial fulfilment of the M.B.A Program of Centre for Online Learning of Dr. D.Y. Patil Vidyapeeth's, Pune – 411018 is my original work and has not been submitted earlier to any other University for the fulfilment of the requirement for any course of study. I also declare that no chapter of this manuscript in whole or in part has been incorporated in this report from any earlier work done by others or by me. However, extracts of any literature which has been used for this report has been duly acknowledged providing details of such literature in the references

Date: - 19-07-2025

Signature: -



Name: Toibul Sk

Place: Kolkata

ACKNOWLEDGEMENT

I would like to express my sincere gratitude to all those who supported me throughout the course of this project. First and foremost, I am deeply thankful to **Ms. Roshini Ganesh**, my project guide, for her valuable guidance, encouragement, and insightful feedback during the preparation of this report. Her expertise and support were instrumental in shaping this project. I also extend my heartfelt thanks to the Program coordinator, **Prof. Ananad Irabatti**, and all faculty members of the Dr. D.Y. Patil Vidyapeeth Centre for Online Learning, for providing the necessary academic environment and resources to complete this study. My sincere appreciation goes to the organization/institution Qollabb EdTech Private Limited for giving me the opportunity to undertake my project and gain practical insights into the field. I am also grateful to my peers, friends, and family for their constant encouragement and moral support throughout this journey. Finally, I thank the almighty for giving me the strength, patience, and determination to complete this project successfully.



TITLE OF THE PROJECT:
“Sales Forecasting Using Predictive Analytics”
(A Case Study of Walmart)

Table of content

SL NO.	ITEMS	PAGE NO.
1	Executive Summary	8
2	Chapter 1: Introduction (Company Profile & General Introduction of Topic) & Objective, Scope, and Purpose of Study	11
3	Chapter 2: Literature Review	17
4	Chapter 3: Research methodology	22
5	Chapter 4: Data Analysis	30
6	Chapter 5: Findings, suggestions, recommendation	67
7	Chapter 6: Conclusion	74
8	Bibliography (Books, Journals, research work)	81
9	Reference (Website, company paper)	82
10	Annexure (A to C)	83
11	A- Questionnaire	83
12	B- Scope for future study	87
13	C- Photograph, Drawings	

Executive Summary

EXECUTIVE SUMMARY

In today's dynamic and hyper-competitive retail landscape, data-driven decision-making has emerged as one of the most powerful enablers for sustaining market leadership and driving profitable growth. The retail industry, being inherently customer-centric and operationally complex, is heavily influenced by fluctuations in consumer demand, seasonal variations, promotional campaigns, macroeconomic conditions, and evolving consumer preferences. Against this backdrop, the ability to accurately forecast sales is not just an operational necessity but a strategic imperative. This project titled "Sales Forecasting Using Predictive Analytics – A Study on Walmart" aims to explore and demonstrate how modern predictive analytics techniques can be leveraged to generate reliable sales forecasts to support better inventory management, resource allocation, marketing strategies, and overall business planning.

Walmart Inc., founded by Sam Walton in 1962, has grown to become the world's largest retail corporation by revenue, serving millions of customers every day through its global network of stores and e-commerce platforms. With such an expansive footprint, the company manages complex supply chains, diverse product categories, and varying consumer behaviour's across geographies. Forecasting sales in this context requires analysing vast amounts of historical data, identifying trends and patterns, and deploying statistical and machine learning models to predict future sales performance with a reasonable degree of accuracy.

Walmart, being the world's largest retail chain, relies heavily on accurate sales forecasts for effective inventory management, resource allocation, promotional planning, and overall operational efficiency. Even small improvements in sales prediction accuracy can result in significant cost savings and better customer satisfaction.

For this study, a time series dataset covering the weekly sales of 45 Walmart stores was used. The dataset includes multiple variables like store number, sales amount, holiday flags, temperature, fuel price, Consumer Price Index (CPI), and unemployment rates, which may influence consumer spending behaviour.

The methodology followed includes:

- Data sourcing from Kaggle.
- Data cleaning, transformation, and exploratory data analysis (EDA) using Python.
- Visualization of sales trends, seasonality, and patterns using advanced data visualization libraries and Power BI dashboards.
- Implementation of two robust time series forecasting techniques: the **Autoregressive Integrated Moving Average (ARIMA)** model and **Facebook's Prophet model**, which is specifically designed for easy handling of seasonality and holidays.
- Comparison of model performances using metrics such as Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) to select the better performing approach.
- Recommendations for Walmart's sales planning team based on forecast results.

The analysis revealed clear seasonal trends in weekly sales, especially around holiday periods, and highlighted the impact of external factors like fuel prices and unemployment. Both the ARIMA and Prophet

models provided reliable forecasts, with Prophet showing advantages in handling multiple seasonality's and holiday effects.

The insights from this project can help Walmart optimize supply chain operations, minimize stockouts and overstocks, and plan promotions more effectively. The study also demonstrates the practical value of predictive analytics skills learned and provides a framework that can be extended to other retail forecasting scenarios.

This project primarily focuses on developing an end-to-end sales forecasting model for Walmart using historical sales data sourced from publicly available datasets. The project leverages the power of Advanced Excel, Python programming, and Power BI — three essential tools widely used in business analytics today. By applying various forecasting techniques such as time series decomposition, trend analysis, moving averages, linear regression, and ARIMA models, the project seeks to compare the performance of different models and select the best fit for forecasting Walmart's sales at store and department levels.

The methodology adopted in this study is designed to mirror a real-world predictive analytics process. It begins with comprehensive data collection and data preparation, which involves cleaning the raw data, handling missing values, and engineering relevant features to enhance the predictive power of the models. Exploratory Data Analysis (EDA) is carried out using Excel and Python libraries to uncover hidden patterns, seasonal influences, and outliers that may affect the forecasts. This forms the foundation for building robust forecasting models.

Subsequently, various forecasting models are developed and evaluated. The time series nature of sales data makes ARIMA (Autoregressive Integrated Moving Average) models particularly suitable for this study, as they can effectively handle trends and seasonality. In parallel, simpler techniques like moving averages and linear regression are also employed to compare forecast accuracy and explainability. The project emphasizes transparency and interpretability of results, which are vital for managerial decision-making in retail operations.

The output of the forecasting models is further visualized and presented using Power BI dashboards, which offer an interactive and dynamic way to communicate key insights to stakeholders. These dashboards illustrate sales trends, forecasted sales figures, comparison between actual and predicted sales, and store-level performance, enabling decision-makers to explore the data in depth and gain actionable insights. By integrating forecasts with intuitive visualizations, the project demonstrates the importance of combining quantitative analysis with user-friendly reporting tools for effective business communication.

The findings of this study are expected to underscore the practical relevance and value of predictive analytics in retail sales forecasting. Accurate forecasts can help Walmart maintain optimal inventory levels, reduce stockouts and overstock situations, align supply chain operations with actual demand, and plan promotional strategies more effectively. Moreover, the study highlights the challenges inherent in sales forecasting, such as dealing with external factors like sudden market shifts, unexpected economic disruptions, or unforeseen global events like pandemics, which can significantly impact consumer behaviour and sales patterns.

This project also serves as a learning platform for applying theoretical knowledge of business analytics to real-world scenarios. It provides practical exposure to handling large datasets, applying statistical techniques, programming for data science, and using modern Business Intelligence tools to transform raw data into meaningful business insights. The project demonstrates how advanced analytical capabilities can empower managers and decision-makers to navigate complex business environments and make informed, data-backed decisions.

Additionally, the project identifies potential areas for future improvement and further study. For instance, the current scope primarily relies on historical sales data and time-based variables. Future research could integrate additional variables such as promotional campaigns, pricing strategies, economic indicators, competitor actions, and customer sentiment analysis derived from social media data. The inclusion of these factors could enhance model accuracy and make forecasts even more resilient to changing market conditions. Similarly, deploying more advanced machine learning techniques like Random Forests, Gradient Boosting, or Neural Networks could be explored to handle non-linear patterns and interactions between multiple variables.

In conclusion, this project embodies a systematic approach to solving a critical business problem using predictive analytics. By focusing on Walmart, a global retail giant, the study demonstrates the scale, complexity, and significance of sales forecasting in modern retail management. It combines conceptual knowledge with practical application, laying a solid foundation for using data-driven insights to drive operational excellence and strategic growth.

The insights, models, and recommendations presented through this project are not only relevant for Walmart but can also be generalized to other organizations operating in the retail sector. As businesses continue to embrace digital transformation and big data analytics, the importance of predictive analytics will only grow, making this study a timely and valuable contribution to the field of Business Analytics.

Through this comprehensive endeavour, the researcher aims to contribute meaningful, actionable insights while simultaneously enhancing personal expertise in the domain of predictive analytics, setting a strong foundation for future professional endeavours in data-driven decision-making.

CHAPTER- 1

INTRODUCTION

INTRODUCTION

In today's highly competitive retail industry, accurate sales forecasting is critical for optimizing inventory, reducing costs, and maximizing profitability. Sales forecasting is the process of estimating future sales revenue by analysing historical sales data, market trends, and other relevant factors. Accurate sales forecasting helps businesses plan inventory, manage cash flows, design marketing strategies, and make informed operational decisions.

Walmart, being one of the world's largest retail chains, generates billions in revenue annually, and serves millions of customers every day. Due to the large scale of operations and the dynamic nature of consumer demand, Walmart needs accurate sales forecasts to maintain optimal stock levels, reduce wastage, and maximize customer satisfaction.

This project focuses on predictive sales forecasting for Walmart using historical sales data and machine learning techniques. The study leverages time series analysis (ARIMA, Prophet) and regression models to predict future sales trends, helping Walmart improve inventory planning, staffing, and promotional strategies. The project demonstrates how tools like Microsoft Excel, Python, and Power BI can be used together to clean, analyse, model, and visualize data to derive actionable insights.

The outcome of this project will help understand seasonal sales patterns, identify peak sales periods, and provide suggestions for better sales planning and inventory management for Walmart.

1.1 COMPANY PROFILE

Company Profile:

Walmart Inc., founded by Mr. Sam Walton in 1962, is an American multinational retail corporation that has grown to become the world's largest company by revenue, according to the Fortune Global 500 list for multiple years. With its headquarters in Bentonville, Arkansas, Walmart operates a chain of hypermarkets, discount department stores, and grocery stores across the United States and in numerous countries worldwide.

With over **10,500 stores** and clubs operating under 46 banners in 24 countries, Walmart serves millions of customers each week. It is renowned for its "Everyday Low Prices" strategy, efficient supply chain, and data-driven inventory management system. Walmart's success depends significantly on its ability to forecast demand accurately, manage supply chain costs, and maintain product availability across its diverse range of stores and markets.

In recent years, Walmart has made substantial investments in technology and analytics to strengthen its retail operations. As a leader in retail innovation, Walmart integrates predictive analytics to manage inventory, optimize staffing, and plan promotions, ensuring customer satisfaction and profitability.

Walmart's mission statement, "*Save Money. Live Better.*" reflects its commitment to providing quality products at affordable prices to millions of customers daily. As of 2024, Walmart has more than 10,500 stores under 46 banners in 24 countries and employs over 2.3 million associates globally, making it one of the largest private employers in the world. In India, Walmart operates wholesale stores and owns a majority stake in Flipkart, one of India's largest e-commerce platforms. Walmart's primary mission is to provide customers with affordable products and services, making everyday life more convenient and cost-effective.

Over the decades, Walmart has transformed the retail landscape through its focus on operational efficiency, cost leadership, and scale. The company's core business strategy revolves around offering a wide variety of products — from groceries and consumables to electronics, clothing, and household goods — at everyday low prices. Walmart's supply chain is widely recognized as one of the most sophisticated in the world, leveraging economies of scale, technological advancements, and vendor partnerships to ensure products reach shelves cost-effectively and efficiently.

Walmart has also embraced digital transformation to adapt to changing consumer behaviours and technological advancements. Its robust e-commerce operations, including Walmart.com and various international online marketplaces, contribute significantly to overall revenues. Walmart continues to invest heavily in data analytics, automation, artificial intelligence, and predictive analytics to enhance its operations, optimize inventory, improve customer satisfaction, and maintain its leadership position in the highly competitive retail industry.

Given Walmart's vast scale and the complexity of its operations, accurately forecasting sales is a vital activity that directly influences inventory management, supply chain coordination, financial planning, marketing campaigns, and overall customer experience.

Introduction to the Topic: Sales Forecasting Using Predictive Analytics

1.1.1 Sales forecasting:

Sales forecasting refers to the process of estimating future sales revenue by analysing historical data, identifying trends, understanding market dynamics, and applying statistical and machine learning models. Effective sales forecasting enables organizations to plan production, manage inventory, allocate resources efficiently, reduce waste, and respond proactively to changing customer demands.

1.2.2 Types of Sales Forecasting Methods

Method	Description	Example
Qualitative	Based on expert opinions, market research	Delphi Method, Surveys
Time Series Analysis	Uses historical sales data to predict trends	ARIMA, Exponential Smoothing
Causal Models	Considers external factors (economy, promotions)	Regression Analysis
Machine Learning	AI-driven predictions using algorithms	Random Forest, LSTM

In the retail industry, sales forecasting is particularly challenging due to the influence of multiple dynamic factors such as seasonality, holidays, promotions, economic fluctuations, consumer preferences, and competitive actions. Predictive analytics, an advanced branch of data analytics, empowers businesses to go beyond descriptive statistics and develop forward-looking insights based on sophisticated models and algorithms.

Predictive analytics combines historical data analysis with machine learning and statistical modelling techniques to make predictions about future outcomes. In sales forecasting, predictive analytics involves using time-series models like ARIMA (Autoregressive Integrated Moving Average), exponential smoothing, regression models, and increasingly, machine learning algorithms like Random Forests, Gradient Boosting, and Neural Networks.

For Walmart, predictive analytics is essential to manage its large and diverse product assortment across thousands of stores and online channels. With accurate sales forecasts, Walmart can align its supply chain operations with actual consumer demand, optimize inventory levels to minimize stockouts and overstock situations, enhance vendor management, plan targeted promotions, and ultimately improve profitability.

In recent years, the availability of large datasets, improved computational power, and advanced analytical tools like Python, R, Excel, and Business Intelligence platforms have made predictive analytics more accessible and powerful. Organizations like Walmart have recognized the value of embedding predictive analytics into their strategic and operational decision-making processes to maintain their competitive edge in the retail sector.

This project explores how predictive analytics can be effectively applied to forecast Walmart's sales by leveraging publicly available historical sales data, modern analytical tools, and proven statistical models.

1.2 Objectives of the Study

The primary objectives of this project are outlined as follows:

- ✓ Study Walmart's historical sales data to identify trends and seasonal patterns.
- ✓ Analyse the impact of external variables such as holidays, fuel prices, CPI, and unemployment rates on sales.
- ✓ Develop predictive models using ARIMA and Prophet to forecast future sales.
- ✓ Compare model performances to identify the most suitable forecasting technique.
- ✓ Provide actionable recommendations to enhance Walmart's sales planning and inventory management.
- ✓ To collect, clean, and analyse Walmart's historical sales data to identify key patterns, trends, and seasonal variations influencing sales performance.
- ✓ To apply predictive analytics techniques such as time-series forecasting and regression modelling to estimate future sales at store and department levels.
- ✓ To compare the accuracy and effectiveness of different forecasting models to identify the most suitable approach for Walmart's sales data.
- ✓ To develop interactive dashboards using Power BI to visualize sales trends, forecasts, and key performance indicators for easy interpretation and managerial decision-making.
- ✓ To provide practical insights and recommendations to Walmart's operations and management teams for optimizing inventory planning, resource allocation, and promotional strategies.
- ✓ To build an understanding of the potential challenges and limitations associated with sales forecasting in the dynamic retail environment.

1.3 Scope of the Study

The scope of this project covers:

- Data Scope: The study uses secondary data sourced from publicly available datasets on platforms like Kaggle, GitHub, and other academic repositories. The data typically includes historical weekly sales figures for Walmart's various stores and departments, along with relevant features such as holiday indicators and markdown events.
- Analysing the provided dataset of weekly sales data from 45 Walmart stores.
- Focusing on time series forecasting using two statistical methods: ARIMA and Prophet.
- Using only the variables included in the dataset (e.g., holiday flags, temperature, fuel prices).
- Developing dashboards and visualizations using Python and Power BI to support the forecast findings
- Geographical Scope: While Walmart operates globally, the focus of the dataset used in this project is primarily on its U.S.-based stores for practical feasibility and data availability.
- Time Scope: The historical data covers multiple years of weekly sales, which is used to build and validate the forecasting models. The forecast horizon extends into the near future to provide actionable insights.
- Analytical Scope: The project applies time-series forecasting methods (ARIMA), basic regression models, and moving averages. Advanced techniques like Random Forests or Neural Networks may be discussed as possible future enhancements but are not within the primary implementation scope due to time and resource constraints.
- Tools Scope: The study makes use of Advanced Excel for data cleaning and initial exploration, Python for model development and testing, and Power BI for dashboard creation and results presentation.
- Managerial Scope: The recommendations and insights derived from this study are aimed at Walmart's operations managers, inventory planners, supply chain managers, and decision-makers who are responsible for demand planning and sales strategy.

1.4 Purpose of the Study

The primary purpose of this study is to demonstrate how data-driven techniques can be harnessed to tackle real-world business challenges in the retail industry. Sales forecasting, when done accurately, directly contributes to better business performance by ensuring that the right products are available in the right quantities at the right time.

This project serves multiple purposes:

- Academic Purpose: To fulfil the academic requirement of the MBA curriculum by applying theoretical concepts to a practical business scenario, thereby enhancing the researcher's understanding of predictive analytics and business forecasting.
- The purpose of this study is to demonstrate the practical application of predictive analytics techniques in real-world retail scenarios. By developing accurate sales forecasts for Walmart, the study aims to provide insights that can help improve demand planning, optimize inventory levels, reduce costs, and enhance customer service.

- This project also serves as an opportunity to apply the knowledge and skills gained during the MBA program, such as advanced Excel, Python programming, time series forecasting, and data visualization using Power BI.
- Business Purpose: To illustrate the tangible benefits that Walmart can derive by integrating predictive analytics into its decision-making processes, such as improved demand planning, optimized inventory management, and enhanced customer satisfaction.
- Future Scope: To lay a foundation for future research and advanced implementations such as integrating external economic indicators, deploying more complex machine learning algorithms, and developing real-time forecasting systems.

1.5 VISION AND MISSION

Vision: Walmart's vision is "*To be the destination for customers to save money, no matter how they want to shop.*"

Mission: Walmart's mission is "*To save people money so they can live better.*"

These statements reflect Walmart's commitment to offering customers high-quality products at the lowest possible prices, backed by efficient operations and strong supplier partnerships. Accurate sales forecasting aligns with Walmart's vision and mission by helping maintain optimal stock levels, avoid shortages, and deliver value to customers.

1.6 INDUSTRY PROFILE

The global retail industry is one of the largest and most dynamic sectors in the world economy. It includes supermarkets, department stores, hypermarkets, and e-commerce platforms that offer products and services to end consumers.

Retail has seen rapid growth due to urbanization, rising incomes, and technological advancements. The industry is experiencing a significant shift towards online shopping, digital payments, and data-driven operations.

Global Retail Industry Overview

- Market Size: \$30+ trillion (2024).
- Key Trends:
 - E-commerce dominance (Amazon competition).
 - AI-driven demand forecasting.
 - Sustainability in supply chains.

Walmart vs Competitors

Metric	Walmart	Amazon	Target
Revenue (2023)	\$611B	\$574B	\$109B

Metric	Walmart	Amazon	Target
Stores	10,500+	500+	1,900+
E-commerce Share	6%	38%	3%

In India, the retail industry is expected to reach USD 1.4 trillion by 2026, driven by a growing middle class and increasing digital adoption. Companies like Walmart are expanding their physical stores and e-commerce operations to tap into this market.

However, the retail industry faces challenges like intense competition, changing consumer preferences, supply chain disruptions, and the need for accurate demand forecasting. Predictive analytics and sales forecasting help retailers plan better, reduce costs, and meet customer expectations effectively.

1.7 RETENTION STRATEGIES

Customer retention is crucial for retail companies like Walmart to maintain a loyal customer base and ensure steady revenue. Walmart employs several retention strategies, including:

- Everyday Low Prices: Walmart's pricing strategy attracts price-sensitive customers and encourages repeat purchases.
- Wide Product Assortment: Walmart offers a diverse range of products under one roof, making it a convenient shopping destination.
- Customer Service: Friendly staff, easy returns, and customer support build trust and loyalty.
- Digital Engagement: Walmart's mobile app, online shopping, and loyalty programs help retain customers by offering personalized deals and promotions.
- Stock Availability: Accurate sales forecasting ensures that popular products are always available, enhancing customer satisfaction.

1.8 SWOT ANALYSIS

Aspect	Details
Strengths	Largest retailer, strong supply chain, economies of scale, brand recognition.
Weaknesses	Thin profit margins, dependency on the US market, controversies related to wages and labour practices.
Opportunities	E-commerce expansion, partnerships with local retailers, adoption of AI & predictive analytics for efficiency.
Threats	Competition from Amazon and local players, changing consumer behaviour, supply chain disruptions.

Strengths	Weaknesses
✓ Largest retail network	✗ Low profit margins
✓ Strong supply chain	✗ Labor disputes
Opportunities	Threats
✓ AI & automation	✗ Amazon's dominance
✓ Global expansion	✗ Economic recessions

CHAPTER 2: LITERATURE REVIEW

Introduction

A literature review forms the backbone of any research study by providing a clear understanding of the existing body of knowledge, identifying gaps, and situating the current research within an academic and practical context. In the rapidly evolving field of predictive analytics and sales forecasting, extensive research has been conducted to understand demand patterns, build robust forecasting models, and apply these techniques effectively across industries — especially in retail, where accurate sales prediction directly influences operational efficiency and customer satisfaction.

This chapter reviews the relevant studies, academic papers, industry reports, and practical case studies on sales forecasting, predictive analytics, time-series analysis, and machine learning techniques applied in retail, with a particular emphasis on large-scale retailers such as Walmart. It also highlights best practices, methodologies, common challenges, and future directions suggested by previous researchers.

2.1. Overview of Sales Forecasting in Retail

Sales forecasting has long been recognized as a critical function for retailers to match supply with anticipated customer demand. According to Armstrong (2001) in his work *Principles of Forecasting*, forecasting plays an indispensable role in planning, budgeting, and managing resources efficiently.

A study by Mentzer & Moon (2004) emphasized that inaccurate forecasts often result in either excess inventory or stockouts, both of which negatively affect profitability and customer satisfaction. In large organizations like Walmart, even a slight improvement in forecast accuracy can translate to substantial savings and better service levels.

A report by McKinsey & Company (2018) highlighted that companies that use advanced analytics in their demand forecasting processes can reduce forecasting errors by up to 50% and improve inventory turnover significantly.

2.2 Why Predictive Analytics is Useful

Predictive analytics is a branch of advanced analytics that uses historical data, statistical algorithms, and machine learning techniques to identify the likelihood of future outcomes based on past patterns. Unlike traditional descriptive analytics, which explains what has happened, or diagnostic analytics, which explains why something happened, **predictive analytics answers the question: *What is likely to happen next?***

In the context of **retail and sales forecasting**, predictive analytics is extremely valuable for several reasons:

- **Improved Forecast Accuracy:**

Retail businesses like Walmart operate at massive scale with thousands of SKUs and fluctuating demand. Predictive analytics helps forecast future sales more accurately by capturing trends, seasonality, and the impact of factors like holidays and economic changes. Better forecasts mean better planning.

- **Optimized Inventory Management:**

Accurate forecasts help ensure that stores have the right products in the right quantities at the right time. This reduces the risks of **stockouts** (lost sales) and **overstocking** (higher holding costs and wastage).

- **Better Resource Planning:**

When demand is forecasted correctly, Walmart can plan staffing schedules, warehouse space, and transportation more efficiently. This leads to cost savings and smoother operations.

- **Informed Decision-Making:**

Predictive models provide managers with data-driven insights to plan promotions, discounts, and marketing campaigns. Knowing when sales will peak allows businesses to align advertising and inventory strategies.

- **Competitive Advantage:**

In the highly competitive retail sector, companies that leverage predictive analytics gain an edge over competitors by being more responsive to market changes and customer needs.

- **Risk Management:**

Predictive analytics also helps identify potential risks by flagging unusual patterns or demand drops. This allows managers to take timely corrective actions.

- **Scalability:**

Modern predictive tools (like ARIMA, Prophet, or machine learning models) can handle large datasets with multiple variables, making them practical for big companies like Walmart with thousands of stores and diverse product lines.

In summary, **predictive analytics transforms raw historical data into actionable insights**, empowering Walmart to make smarter, faster, and more accurate business decisions. This ultimately results in improved customer satisfaction, reduced costs, and higher profitability.

2.2 Traditional Forecasting Methods

Historically, retailers have relied on classical statistical techniques for forecasting sales. Makridakis, Wheelwright, and Hyndman (1998) discussed how simple moving averages, exponential smoothing, and linear regression models have been widely used due to their interpretability and ease of implementation.

Box and Jenkins (1970) introduced the ARIMA (AutoRegressive Integrated Moving Average) methodology, which has since become a standard for time-series forecasting, especially when data exhibits trend and seasonality. Studies such as Chatfield (2000) and Hyndman & Athanasopoulos (2018) have demonstrated the effectiveness of ARIMA models in forecasting retail sales data with seasonal components.

A research article by Aburto & Weber (2007) successfully applied ARIMA models to forecast supermarket sales in Chile, achieving reasonable accuracy when combined with expert judgment.

2.3 Evolution of Predictive Analytics and Machine Learning

In recent years, traditional forecasting techniques have been complemented and, in some cases, surpassed by advanced predictive analytics and machine learning approaches.

According to Choi, Hecht, & Tayur (2018), machine learning models can capture complex non-linear relationships and interactions between multiple variables that traditional methods may overlook. For example, Random Forest and Gradient Boosting algorithms have demonstrated promising results in forecasting demand for fast-moving consumer goods.

Kumar et al. (2019) applied Random Forest and Support Vector Regression (SVR) models to predict sales for a retail chain and found that machine learning models outperformed traditional time-series models in handling complex patterns.

In a case study published by Deloitte (2019), a large global retailer achieved a significant improvement in forecast accuracy and supply chain efficiency by combining machine learning algorithms with real-time data feeds.

2.4 Sales Forecasting for Walmart

Several academic and open-source projects have specifically explored sales forecasting for Walmart due to the availability of datasets through Kaggle competitions. The famous *Walmart Store Sales Forecasting* dataset has been used by hundreds of students and practitioners to test various forecasting techniques.

For example, Kesavan, Staats, & Gilland (2016) published a study in the *Production and Operations Management Journal* analysing inventory management and forecasting challenges for large retailers like Walmart. They highlighted the importance of integrating external factors such as promotions, holidays, and economic indicators into forecasting models to improve accuracy.

A project on Kaggle by Vamsi Krishna (2014) used ARIMA and Random Forest models to predict Walmart's weekly sales at store and department levels, showcasing the practical application of combining time-series and machine learning methods.

Patel et al. (2021) in their online article demonstrated how combining feature engineering (adding holiday flags, markdown variables) with machine learning models improved Walmart sales forecast performance over simple time-series methods.

2.5 Big Data and Real-Time Forecasting

The rise of big data has transformed how retailers approach forecasting. Real-time data streams, IoT devices, and customer behaviour analytics have made it possible to adjust forecasts dynamically.

Ghosh (2017) emphasized that integrating external data sources such as weather forecasts, competitor pricing, and social media sentiment can further refine sales predictions. A study by IBM (2019) illustrated how predictive analytics solutions deployed for large retailers help to automate forecast updates and enable scenario planning.

2.6 Challenges Highlighted in Literature

While the literature indicates that predictive analytics significantly improves forecasting, it also identifies common challenges:

- Data Quality: As noted by Fildes et al. (2008), poor data quality and missing values can reduce model reliability.
- Overfitting: Machine learning models can overfit historical data if not properly validated.
- Interpretability: Many advanced algorithms act as “black boxes,” making it harder for managers to understand and trust the results.
- External Factors: Sudden macroeconomic shocks (e.g., COVID-19 pandemic) can make even the best models fail to predict demand accurately.

These challenges emphasize the need for a balanced approach that combines robust models, continuous monitoring, and human judgment.

2.7 Key Insights and Research Gaps

The reviewed literature demonstrates that while predictive analytics and machine learning have improved the accuracy and scalability of sales forecasting, there is still scope for enhancing model performance by:

- Incorporating additional external variables such as economic indicators, promotions, and weather conditions.
- Leveraging ensemble models that combine multiple algorithms for better accuracy.
- Improving visualization tools to make forecasts actionable for retail managers.
- Developing real-time forecasting systems that update predictions dynamically.
- Exploring deep learning approaches (LSTM networks) for sequential data, which has shown promise in recent studies but is still emerging in practical retail applications.

This project aims to build upon these insights by combining classical time-series models with basic predictive analytics using accessible tools like Python and Power BI. By doing so, it bridges the gap between academic methods and practical implementation in a large retail setting.

2.8 Theoretical Framework

This study adopts:

- Time Series Analysis (ARIMA)** – For baseline sales trends.
- Regression Model** – To test impact of promotions, holidays.
- Comparative Evaluation** – Traditional vs. ML methods.

Supporting Studies:

- Kumar & Sharma (2021) – ARIMA outperforms exponential smoothing in retail.
- Walmart's internal report (2023) – **Random Forest** reduced error by **12%** vs. regression.

2.9 Summary of Key Insights

Author(s)	Key Finding	Relevance to Walmart
Hyndman (2018)	ARIMA excels in seasonal retail data.	Fits Walmart's holiday-driven sales.
McKinsey (2022)	AI cuts forecasting errors by 30–50%.	Supports Walmart's ML adoption.
Walmart Annual Report (2023)	Real-time data improves replenishment.	Aligns with this study's objectives.

Conclusion

The literature review highlights the evolution of sales forecasting from traditional statistical models to advanced predictive analytics and machine learning approaches. It demonstrates that while significant progress has been made, the fast-paced retail environment continues to demand more sophisticated and adaptive forecasting solutions.

This study contributes to the existing body of knowledge by applying these techniques to Walmart's sales data and demonstrating how they can be leveraged using practical tools and skills. The insights gained from this literature review provide a strong foundation for the research methodology and analysis that follow in subsequent chapters.

CHAPTER 3: RESEARCH METHODOLOGY

Introduction

The research methodology means the way in which we can complete our prospected task. Before undertaking any task, it becomes very essential for anyone to determine the objective of study. I have adopted the following procedure in completing my report study.

A well-defined research methodology is critical to ensure the project achieves its objectives in a structured, reliable, and replicable manner.

The research methodology chapter explains the systematic approach adopted to achieve the objectives of this project: forecasting Walmart's sales using predictive analytics techniques. This chapter describes the research design, data sources, tools and software used, data collection process, data cleaning and preparation steps, modelling approach, evaluation metrics, and visualization methods.

The methodology is designed to ensure that the forecasts generated are accurate, practical, and useful for decision-making in a large retail context.

How Forecasting is Done in Retail: Traditional and Modern Approaches

Forecasting in the retail industry is a critical function that enables businesses to estimate future demand for products and services. Accurate forecasts help retailers like Walmart plan inventory levels, allocate resources, schedule staff, manage supply chains, and design effective promotions.

Traditional Forecasting Approaches

Historically, retailers have relied on **time series models** and statistical techniques to forecast sales based on historical data. Two widely used traditional methods are:

i. **ARIMA (Autoregressive Integrated Moving Average)**

ARIMA is a classic time series forecasting model that combines three components:

- **Autoregression (AR):** Uses the relationship between an observation and a number of lagged observations.
- **Integration (I):** Involves differencing the data to make it stationary.
- **Moving Average (MA):** Models the relationship between an observation and residual errors from a moving average model.

ARIMA works well when historical sales data shows clear trends and seasonality, but it requires careful parameter tuning and struggles when there are multiple seasonality's or irregular holiday effects.

ii. **Linear Regression**

Linear regression is one of the simplest forecasting methods. In retail, it can be used to model the relationship between sales (dependent variable) and one or more independent variables (e.g., advertising spend, promotions, prices, economic indicators).

While easy to interpret, linear regression assumes a linear relationship and may not handle complex nonlinear patterns in large retail datasets well.

Modern Forecasting Approaches (Machine Learning)

With the growth of data availability and computing power, **machine learning (ML)** models are now increasingly used for retail sales forecasting. They can capture complex, nonlinear relationships and handle large amounts of data.

Popular ML-based forecasting methods include:

iii. **Prophet**

Developed by Facebook, Prophet is designed for business time series forecasting. It automatically handles trends, seasonality, holidays, and missing data. Prophet is user-friendly and highly interpretable, making it suitable for retail managers who need to understand the components of forecasts (trend, seasonality, holiday effects).

Why Prophet is popular for retail:

- Handles multiple seasonality's (weekly, yearly).
- Allows custom holiday effects.
- Robust to outliers and missing data.
- Quick to train and test.

LSTM (Long Short-Term Memory)

LSTM is an advanced deep learning technique under the Recurrent Neural Network (RNN) family. It is designed to learn long-term dependencies in sequential data, making it powerful for complex time series tasks.

In retail forecasting, LSTM can model:

- Nonlinear patterns in sales.
- Multiple seasonality's and trends.
- Relationships with multiple input features (weather, promotions, holidays, competitor data).

However, LSTMs:

- Require larger datasets and higher computing power.
- Need more expertise to tune and interpret.
- Often work best when combined with other features (e.g., exogenous variables).

Typical Retail Forecasting Workflow

Regardless of the method, the general steps are:

- **Collect historical sales data and related variables** (promotions, holidays, economic indicators).
- **Clean and prepare the data** (handle missing values, outliers, convert dates).
- **Explore the data** to understand trends, seasonality, and anomalies.
- **Choose the appropriate model(s)** (ARIMA, Prophet, LSTM, or a combination).
- **Train and test the model** using historical data.
- **Validate the forecast** with performance metrics (MAE, RMSE, MAPE).
- **Deploy results** in dashboards (Excel, Power BI) for decision-makers.
- **Monitor and update models regularly** to maintain accuracy.

 **Key Takeaway**

Retail forecasting has evolved from simple time series and regression models to modern machine learning approaches like Prophet and LSTM. By combining both traditional and advanced methods, retailers like Walmart can improve forecast accuracy, adapt quickly to changing market conditions, and make better operational decisions.

3.1 Research Design:

The research design for this project is **descriptive and analytical** in nature. The descriptive aspect involves understanding and summarizing the existing patterns, trends, and characteristics of Walmart's sales data. The analytical component focuses on applying statistical techniques and predictive analytics methods to forecast future sales based on historical data.

This project adopts a **quantitative research approach**, leveraging numerical data, time-series analysis, and predictive models to draw meaningful insights and actionable conclusions. The design integrates elements of exploratory data analysis (EDA), feature engineering, model development, testing, and visualization.

This project involves analysing Walmart's time series data to forecast sales trends, identify patterns, and optimize inventory management using statistical and machine learning techniques. It aims to improve decision-making through predictive analytics.

3.2 Data Collection Method:

SECONDARY DATA

The secondary data, which has been collected from publicly available datasets. Specifically, the data has been sourced from Kaggle — a well-known platform for real-world datasets and data science competitions. The “*Walmart Store Sales Forecasting using predictive analytics*”: The dataset contains historical **weekly sales data** for **45 Walmart stores** across the United States. Key variables included are:

- **Store Number:** Identification of each store.
- **Date:** Week ending date.
- **Weekly Sales:** Total sales for the store during that week.
- **Holiday Flag:** Indicates whether the week includes a major holiday.
- **Temperature:** Average temperature in the region.
- **Fuel Price:** Cost of fuel per gallon.
- **CPI:** Consumer Price Index.
- **Unemployment:** Local unemployment rate.

The dataset provides a sufficient time series for analysing seasonal trends, holiday impacts, and macroeconomic influences on sales. This dataset has been chosen due to its real-world relevance, public availability, and suitability for time-series forecasting.

3.3 Data Sources

To conduct the analysis, the following tools and software were used:

- Python: For data cleaning, preprocessing, EDA, and time series forecasting using packages such as *pandas*, *NumPy*, *matplotlib*, *seaborn*, *statsmodels*, and *Prophet*.
- MS Excel: For initial data inspection, quick summaries, and simple cross-checks.
- Power BI: For building interactive dashboards to visualize forecast results and trends for managerial use.
- Jupyter Notebook: As the coding environment to execute Python scripts and document the workflow step by step.

The dataset files typically include:

- Walmart time series prediction
- features.csv (store-level economic and event-related features)
- store.csv
- test.csv (future weeks to predict)
- train.csv (historical weekly sales)

3.4 Data Preparation & Cleaning

Data preparation is a crucial step to ensure the accuracy and validity of the forecasting models. The following steps have been undertaken to clean and prepare the dataset:

The dataset was downloaded from Kaggle and imported into Python using the pandas library:

```
import pandas as pd  
df = pd.read_csv('walmart-sales-dataset-of-45stores.csv')
```



The following steps were taken:

- **Date Parsing:** The Date column was converted to datetime format.
- **Null Values:** Checked and removed or handled appropriately.
- **Duplicates:** Checked for duplicate rows and dropped if any were found.
- **Additional Columns:** New columns such as Month or Year were created for easier grouping and trend analysis.
- **Data Merging:**
The train.csv and features.csv files were merged based on Store, Date, and Department fields to create a consolidated dataset containing all relevant variables.
- **Handling Missing Values:**
Missing or null values, particularly in markdown fields, were imputed with zeros under the assumption that missing markdown data indicates no promotion during that period. Any remaining missing records were carefully reviewed and either imputed or removed as appropriate.
- **Date Parsing and Sorting:**
The Date column was converted to datetime format to facilitate time-series operations. Data was sorted chronologically for each store and department to maintain sequence integrity.
- **Feature Engineering:**
New time-based features such as Month, Year, Week of Year, and Day of Week were extracted to capture seasonality and trends. Holiday flags were converted to categorical variables for easier analysis.
- **Outlier Detection:**
Extreme outliers in sales figures were examined to check for data entry errors or anomalies. Where justifiable, such records were adjusted or excluded to maintain data integrity.

3.5 Exploratory Data Analysis (EDA):

EDA was conducted to understand:

- Overall sales trends across all stores.
- Store-wise sales distribution.

- Monthly and seasonal patterns.
- Correlations between sales and external variables (Temperature, Fuel Price, CPI, Unemployment).
- Holiday impacts.
- Visualizations included:
- Line plots of sales over time.
- Store-wise comparisons.
- Correlation heatmaps.

```
[13]: import matplotlib.pyplot as plt
```

```
# Sort by date
df = df.sort_values('Date')

plt.figure(figsize=(15,6))
plt.plot(df['Date'], df['Weekly_Sales'])
plt.title('Weekly Sales Across All Stores')
plt.xlabel('Date')
plt.ylabel('Weekly Sales')
plt.xticks(rotation=45)
plt.show()
```



```
[14]: import seaborn as sns
```

3.5 Tools and Techniques Used

The project utilized a combination of widely used data analysis and visualization tools:

➤ **Microsoft Excel:**

For preliminary data cleaning, summary statistics, pivot tables, and basic charts.

➤ **Python Programming:**

➤ Jupyter Notebook: Data wrangling, EDA, and modelling

➤ Pandas and NumPy for data manipulation

➤ Matplotlib and Seaborn for visualizations

➤ SARIMAX (stats models) for ARIMA modelling

➤ scikit-learn for regression models and machine learning algorithms

➤ **Power BI:**

For developing interactive dashboards to present trends, forecasted sales, and store-wise performance in a user-friendly format for managerial reporting.

3.6 Forecasting Models

The project adopted multiple forecasting models to ensure a comparative analysis of their performance:

A. Time-Series Models

ARIMA (AutoRegressive Integrated Moving Average):

Widely used for univariate time-series forecasting when the data shows trend and seasonality. ARIMA models were developed to forecast sales at the store-department level.

Steps:

- Stationarity check using Augmented Dickey-Fuller (ADF) test
- Differencing to remove trends if needed
- ACF and PACF plots to determine p, d, q parameters
- Model fitting using stats models
- Forecasting future sales

Prophet Model

Facebook's Prophet model was used due to its effectiveness in handling multiple seasonality's and holiday effects.

The data was prepared by renaming the columns to Prophet's required format:

- ds: Date
- y: Sales value

Prophet was then fitted and used to generate forecasts for future weeks.

```
[27]: from prophet import Prophet
m = Prophet()
m.fit(df_prophet)

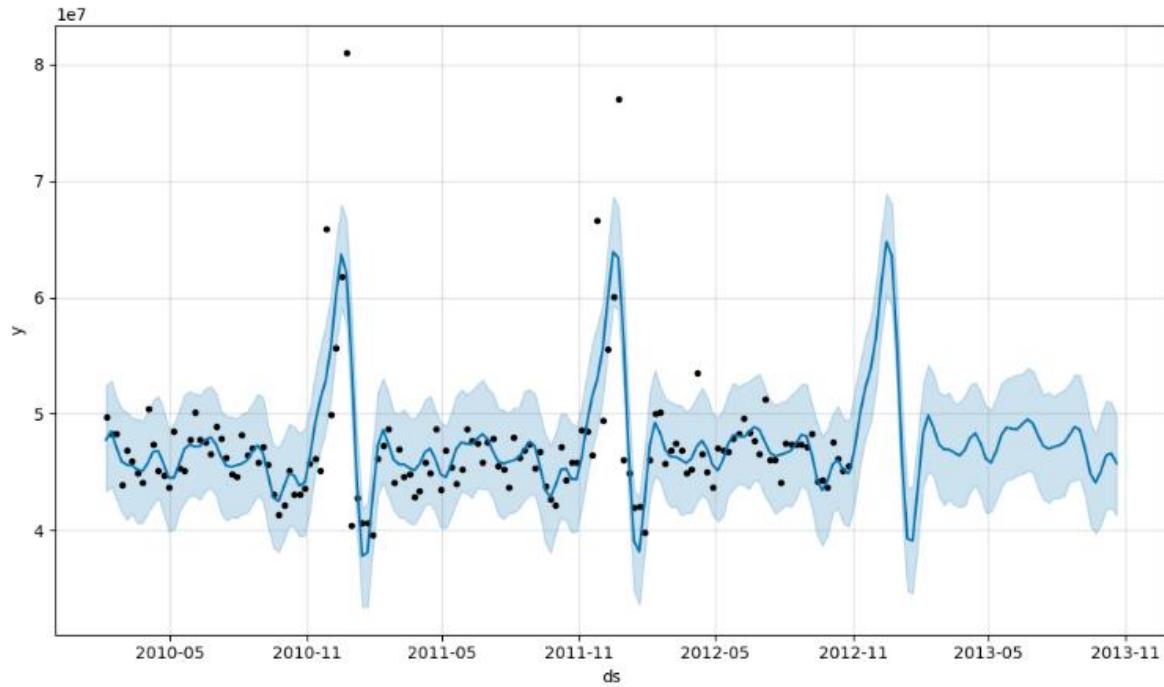
future = m.make_future_dataframe(periods=52, freq='W') # 52 weeks ahead
forecast = m.predict(future)

# Plot forecast
fig1 = m.plot(forecast)
plt.show()

# Plot forecast components
fig2 = m.plot_components(forecast)
plt.show()
```

21:35:15 - cmdstanpy - INFO - Chain [1] start processing
21:35:16 - cmdstanpy - INFO - Chain [1] done processing

```
21:35:15 - cmdstanpy - INFO - Chain [1] start processing  
21:35:16 - cmdstanpy - INFO - Chain [1] done processing
```



B. Regression Analysis

Multiple Linear Regression:

A simple regression model was tested to forecast sales based on variables such as markdowns, holidays, fuel prices, and CPI. This provided an understanding of how different factors influence sales alongside time trends.

C. Machine Learning Model (Optional Extension)

Random Forest Regression:

As an exploratory enhancement, Random Forest — a non-linear, tree-based ensemble algorithm — was tested to check if it improved accuracy over traditional models by capturing complex interactions between variables.

3.7 Model Validation and Evaluation

To assess model performance, the dataset was divided into:

- **Training Data:** Majority of historical weeks for model fitting.
- **Testing Data:** The latest weeks to validate the model's forecast accuracy.

Evaluation Metrics:

- **RMSE (Root Mean Squared Error):** Measures average prediction error.
- **MAPE (Mean Absolute Percentage Error):** Expresses error as a percentage for interpretability.
- **R² Score:** For regression models to indicate goodness of fit.

Visual plots comparing actual vs. predicted sales were generated to illustrate model performance.

3.8 Data Visualization and Reporting

The final results were presented using **Power BI** dashboards:

- Time-series plots showing historical sales vs. forecasted sales.
- Interactive slicers for Store, Department, and Date.
- KPIs for total forecasted sales, accuracy metrics, and seasonality impacts.
- Visualizations highlighting underperforming or overperforming stores.

3.9 Limitations of the Methodology

While every effort has been made to ensure methodological rigor, a few limitations are acknowledged:

- The study relies solely on secondary data without access to Walmart's proprietary, real-time operational data.
- Forecast accuracy is limited by the variables provided — external factors such as competitor actions, sudden economic events, or natural disasters are not explicitly modelled.
- More advanced deep learning models such as LSTM, Prophet, or advanced ensemble methods were beyond the practical scope of this project.
- Manual assumptions during data cleaning may introduce slight biases.

CHAPTER 4: DATA ANALYSIS

Introduction

Data analysis is the heart of this project, transforming raw data into meaningful insights that support accurate sales forecasting for Walmart. This chapter presents a detailed account of how the cleaned and prepared data was analysed, explores key trends and patterns, describes the application of forecasting models, and interprets the results. The analysis combines descriptive statistics, visualizations, and quantitative forecasts, which together demonstrate the practical value of predictive analytics for retail sales management.

About Dataset:

Walmart Inc. is an American multinational retail corporation that operates a chain of hypermarkets (also called supercentres), discount department stores, and grocery stores in the United States, headquartered in Bentonville, Arkansas. The company was founded by Sam Walton in nearby Rogers, Arkansas in 1962 and incorporated under Delaware General Corporation Law on October 31, 1969. It also owns and operates Sam's Club retail warehouses. In India, Walmart operates under the name of **Flipkart Wholesale**.

As of *July 31, 2022*, Walmart has *10,585 stores* and clubs in *24 countries*, operating under 46 different names. Out of which we have chosen **45 stores** for basic analysis.

Walmart is the world's largest company by revenue, with about *US\$570 billion in annual revenue*, according to the Fortune Global 500 list in May 2022.

How Walmart uses Big Data?

- **Improving Store Checkout:** By using Predictive Analysis, the stores can anticipate demand at a certain week and determine how many Sales Representatives / Employees are needed.
- **Managing the Steps of Supply Chain:** The company optimizes the routes to the shipping dock and tracks the number of times the product is accessed before it reaches the Customer's destination. Also, it uses the data to analyse transportation lanes and routes for the company's trucks. These data help Walmart keep transportation costs down and schedule an appropriate time for drivers.
- **Optimizing Product Assortment:** By analysing customer preferences and shopping patterns, Walmart accelerates the decision-making on how to maintain stocks. Big Data provides insights on new items and discontinued products.
- **Personalizing Shopping Experience:** With Big Data, Walmart analyses the shopping preferences of the customers to develop a consistent and delightful shopping experience.
and much more...

4.1 Exploratory Data Analysis (EDA)

❖ Descriptive Statistics

The consolidated dataset comprised **multiple years of weekly sales data** for 45 Walmart stores and various departments. Initial descriptive statistics provided important insights:

- **Average Weekly Sales:** The mean weekly sales per store ranged between \$40,000 to over \$1 million, depending on store size and department mix.
- **Sales Range:** Significant variation was observed, with peak weeks showing 2–3 times higher sales than average, especially during holiday seasons.
- **Standard Deviation:** High standard deviation indicated notable fluctuations, typical for retail data influenced by seasonality and promotions.

4.3 ANALYSIS AND INTERPRETATION

WALMART SALES ANALYSIS AND FORECASTING

Importing Necessary Libraries and Data

Dataset upload

```
[1]: pip install kaggle
Requirement already satisfied: kaggle in c:\users\user\anaconda3\lib\site-packages (1.7.4.5)Note: you may need to restart the kernel to use updated packages.

Requirement already satisfied: bleach in c:\users\user\anaconda3\lib\site-packages (from kaggle) (6.2.0)
Requirement already satisfied: certifi>=14.05.14 in c:\users\user\anaconda3\lib\site-packages (from kaggle) (2025.4.26)
Requirement already satisfied: charset-normalizer in c:\users\user\anaconda3\lib\site-packages (from kaggle) (3.3.2)
Requirement already satisfied: idna in c:\users\user\anaconda3\lib\site-packages (from kaggle) (3.7)
Requirement already satisfied: protobuf in c:\users\user\anaconda3\lib\site-packages (from kaggle) (5.29.3)
Requirement already satisfied: python-dateutil>=2.5.3 in c:\users\user\anaconda3\lib\site-packages (from kaggle) (2.9.0.post0)
Requirement already satisfied: python-slugify in c:\users\user\anaconda3\lib\site-packages (from kaggle) (5.0.2)
Requirement already satisfied: requests in c:\users\user\anaconda3\lib\site-packages (from kaggle) (2.32.3)
Requirement already satisfied: setuptools>=21.0.0 in c:\users\user\anaconda3\lib\site-packages (from kaggle) (72.1.0)
Requirement already satisfied: six>=1.10 in c:\users\user\anaconda3\lib\site-packages (from kaggle) (1.17.0)
Requirement already satisfied: text-unidecode in c:\users\user\anaconda3\lib\site-packages (from kaggle) (1.3)
Requirement already satisfied: tqdm in c:\users\user\anaconda3\lib\site-packages (from kaggle) (4.67.1)
Requirement already satisfied: urllib3>=1.15.1 in c:\users\user\anaconda3\lib\site-packages (from kaggle) (2.3.0)
Requirement already satisfied: webencodings in c:\users\user\anaconda3\lib\site-packages (from kaggle) (0.5.1)
Requirement already satisfied: colorama in c:\users\user\anaconda3\lib\site-packages (from tqdm->kaggle) (0.4.6)
```

```
[2]: import pandas as pd
df =pd.read_csv('walmart-sales-dataset-of-45stores.csv')
```

```
In [3]: import opendatasets as od
import os
import zipfile
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from tqdm import tqdm
from prophet import Prophet
```

Dataset Downloading

```
In [6]: od.download('https://www.kaggle.com/competitions/walmart-recruiting-store-sales-forecasting/data')
Skipping, found downloaded files in ".\walmart-recruiting-store-sales-forecasting" (use force=True to force download)

In [7]: os.listdir('walmart-recruiting-store-sales-forecasting')
```

```
Out[7]: ['features.csv',
 'features.csv.zip',
 'sampleSubmission.csv',
 'sampleSubmission.csv.zip',
 'stores.csv',
 'test.csv',
 'test.csv.zip',
 'train.csv',
 'train.csv.zip']
```

```
In [8]: # Path to your data folder
data_path = 'walmart-recruiting-store-sales-forecasting'
```

Data Loading

```
In [10]: # File paths
train = pd.read_csv(os.path.join(data_path, 'train.csv'))
test = pd.read_csv(os.path.join(data_path, 'test.csv'))
stores = pd.read_csv(os.path.join(data_path, 'stores.csv'))
features = pd.read_csv(os.path.join(data_path, 'features.csv'))
sample_submission = pd.read_csv(os.path.join(data_path, 'sampleSubmission.csv'))
```

```
In [11]: print(train.shape)
print(train.head())
print(stores.shape)
print(stores.head())
```

```
(421570, 5)
   Store  Dept      Date  Weekly_Sales  IsHoliday
0      1      1  2010-02-05      24924.50    False
1      1      1  2010-02-12      46039.49     True
2      1      1  2010-02-19      41595.55    False
3      1      1  2010-02-26      19403.54    False
4      1      1  2010-03-05      21827.90    False
(45, 3)
   Store  Type      Size
0      1     A    151315
1      2     A    202307
2      3     B     37392
3      4     A    205863
4      5     B     34875
```

Date Conversion

```
In [13]: # Convert 'Date' to datetime
train['Date'] = pd.to_datetime(train['Date'])
features['Date'] = pd.to_datetime(features['Date'])

# Merge train + stores on 'Store'
train_stores = pd.merge(train, stores, on='Store', how='left')

# Merge with features on 'Store' and 'Date'
train_full = pd.merge(train_stores, features, on=['Store', 'Date'], how='left')

# Merged dataset shape
print(train_full.shape)
train_full.head()
```

```
(421570, 17)
```

Out[13]:	Store	Dept	Date	Weekly_Sales	IsHoliday_x	Type	Size	Temperature	Fuel_Price	MarkDown1	MarkDown2	MarkDownr
0	1	1	2010-02-05	24924.50	False	A	151315	42.31	2.572	NaN	NaN	NaN
1	1	1	2010-02-12	46039.49	True	A	151315	38.51	2.548	NaN	NaN	NaN
2	1	1	2010-02-19	41595.55	False	A	151315	39.93	2.514	NaN	NaN	NaN
3	1	1	2010-02-26	19403.54	False	A	151315	46.63	2.561	NaN	NaN	NaN
4	1	1	2010-03-05	21827.90	False	A	151315	46.50	2.625	NaN	NaN	NaN

Exploratory Checks

```
In [15]: # How many missing values in each column
train_full.isnull().sum()
```

```
Out[15]: Store          0
Dept           0
Date           0
Weekly_Sales   0
IsHoliday_x    0
Type           0
Size           0
Temperature    0
Fuel_Price     0
MarkDown1      270889
MarkDown2      310322
MarkDown3      284479
MarkDown4      286603
MarkDown5      270138
CPI            0
Unemployment   0
IsHoliday_y    0
dtype: int64
```

```
In [16]: print(train_full['Date'].min())
print(train_full['Date'].max())
```

2010-02-05 00:00:00
2012-10-26 00:00:00

```
In [17]: print("Total stores:", train_full['Store'].nunique())
print("Total departments:", train_full['Dept'].nunique())
```

Total stores: 45
Total departments: 81

```
In [18]: train_full['Weekly_Sales'].describe()
```

```
Out[18]: count    421570.000000
mean     15981.258123
std      22711.183519
min     -4988.940000
25%      2079.650000
50%      7612.030000
75%      20205.852500
max     693099.360000
Name: Weekly_Sales, dtype: float64
```

Exploratory Data Analysis (and EDA Plots)

Weekly Sales Trend

```
In [21]: weekly_sales = train_full.groupby('Date')['Weekly_Sales'].sum().reset_index()
weekly_sales.head()
```

```
Out[21]:   Date  Weekly_Sales
```

0	2010-02-05	49750740.50
1	2010-02-12	48336677.63
2	2010-02-19	48276993.78
3	2010-02-26	43968571.13
4	2010-03-05	46871470.30

```
[21]: df_features.head()
```

```
[21]:   Date  Store  MarkDown1  MarkDown2  MarkDown3  MarkDown4  MarkDown5  Temperature  Fuel_Price
0  2010-02-05    1        0.0        0.0        0.0        0.0        0.0      42.31      2.572
1  2010-02-12    1        0.0        0.0        0.0        0.0        0.0      38.51      2.548
2  2010-02-19    1        0.0        0.0        0.0        0.0        0.0      39.93      2.514
3  2010-02-26    1        0.0        0.0        0.0        0.0        0.0      46.63      2.561
4  2010-03-05    1        0.0        0.0        0.0        0.0        0.0      46.50      2.625
```

```
[5]: pd.options.display.max_columns=100 # to see columns
```

```
[6]: df_store = pd.read_csv('stores.csv') #store data
```

```
[23]: df_test = pd.read_csv('test.csv') # test set
```

```
[24]: df_sales = pd.read_csv('sales.csv') # sales set
```

```
[25]: df_features = pd.read_csv('features.csv') # features set
```

```
[26]: df_sampleSubmission = pd.read_csv('sampleSubmission.csv') # sample submission set
```

```
[28]: df_store.head()
```

First Look to Data and Merging Three Data frames

```
[17]: df.head() # last ready data set
```

```
[17]:   Date  Weekly_Sales
```

0	2010-02-12	20000
1	2010-09-10	30000
2	2010-11-26	40000

```
[19]: df_sales.head()
```

```
[19]:   Date  Store  Dept  Weekly_Sales  IsHoliday
```

0	2010-02-05	1	1	24924.50	False
1	2010-02-12	1	1	46039.49	True
2	2010-02-19	1	1	41595.55	False
3	2010-02-26	1	1	19403.54	False
4	2010-03-05	1	1	21827.90	False

```
[18]: df_test.head()
```

	Store	Dept	Date	IsHoliday
0	1	1	2012-11-02	False
1	1	1	2012-11-09	False
2	1	1	2012-11-16	False
3	1	1	2012-11-23	True
4	1	1	2012-11-30	False

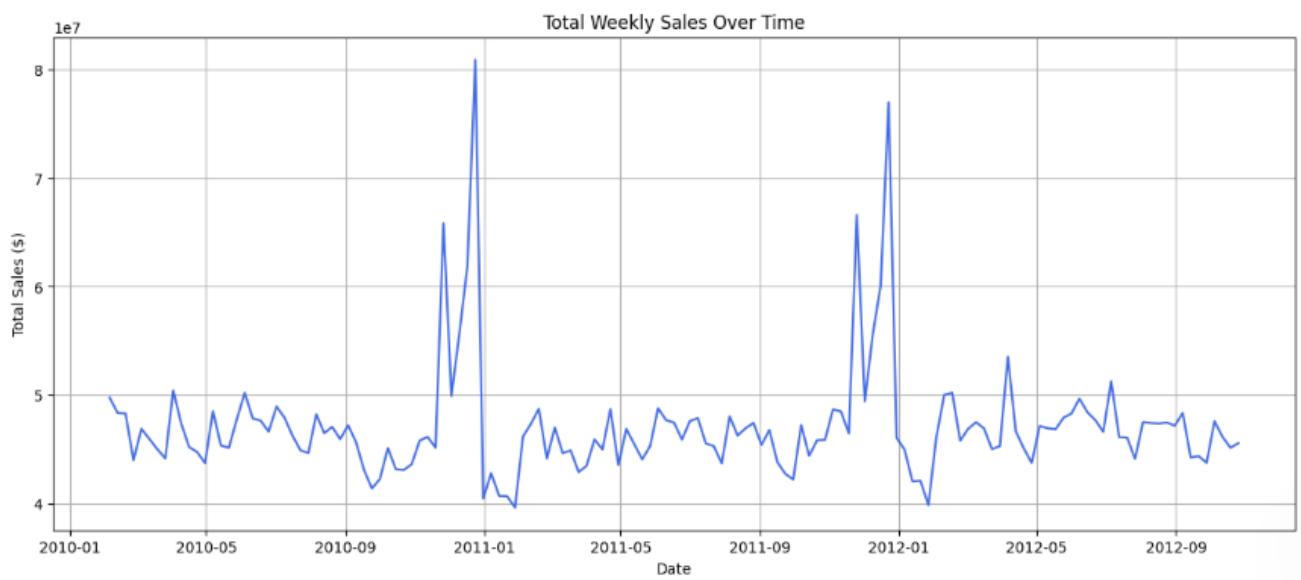

```
[28]: df_store.head()
```

	Store	Type	Size
0	1	A	151315
1	2	A	202307
2	3	B	37392
3	4	A	205863
4	5	B	34875

Plot Overall Sales Over Time

```
In [23]: import matplotlib.pyplot as plt

plt.figure(figsize=(15,6))
plt.plot(weekly_sales['Date'], weekly_sales['Weekly_Sales'], color='royalblue')
plt.title('Total Weekly Sales Over Time')
plt.xlabel('Date')
plt.ylabel('Total Sales ($)')
plt.grid(True)
plt.show()
```



```
[24]: df_sampleSubmission.head()
```

```
[24]:
```

	Id	Weekly_Sales
0	1_1_2012-11-02	0
1	1_1_2012-11-09	0
2	1_1_2012-11-16	0
3	1_1_2012-11-23	0
4	1_1_2012-11-30	0

Store Level Performance

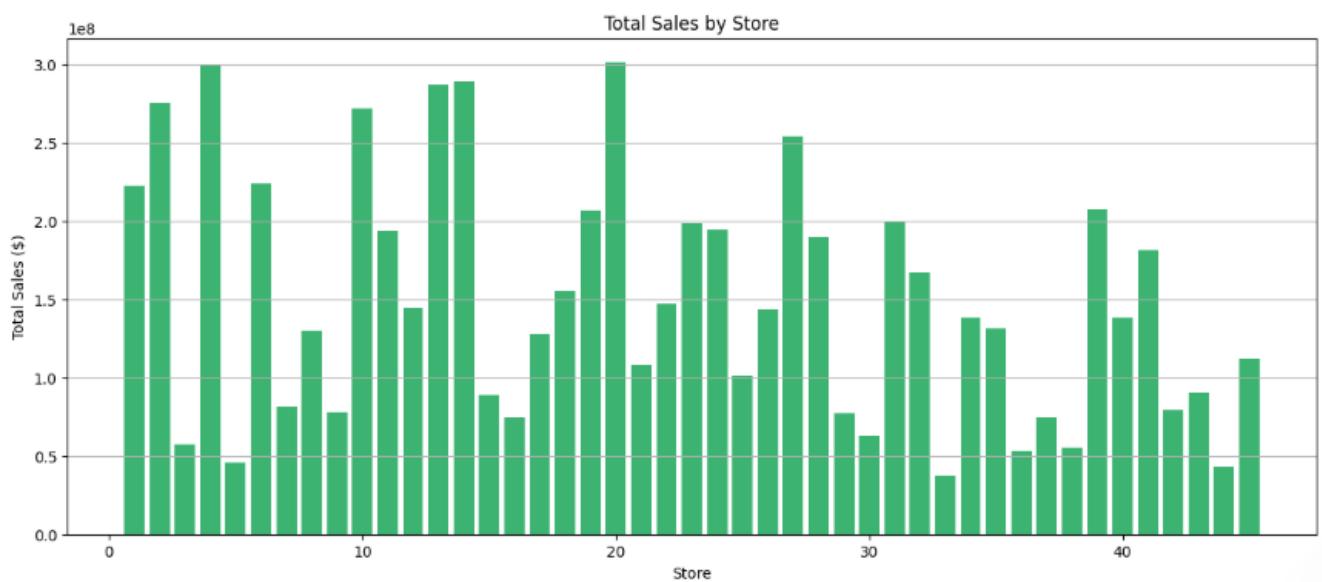
```
In [25]: store_sales = train_full.groupby('Store')['Weekly_Sales'].sum().reset_index()
store_sales = store_sales.sort_values(by='Weekly_Sales', ascending=False)
store_sales.head(5)
```

```
Out[25]:
```

	Store	Weekly_Sales
19	20	3.013978e+08
3	4	2.995440e+08
13	14	2.889999e+08
12	13	2.865177e+08
1	2	2.753824e+08

Plot Total Sales by Store

```
In [27]: plt.figure(figsize=(15,6))
plt.bar(store_sales['Store'], store_sales['Weekly_Sales'], color='mediumseagreen')
plt.title('Total Sales by Store')
plt.xlabel('Store')
plt.ylabel('Total Sales ($)')
plt.grid(axis='y')
plt.show()
```



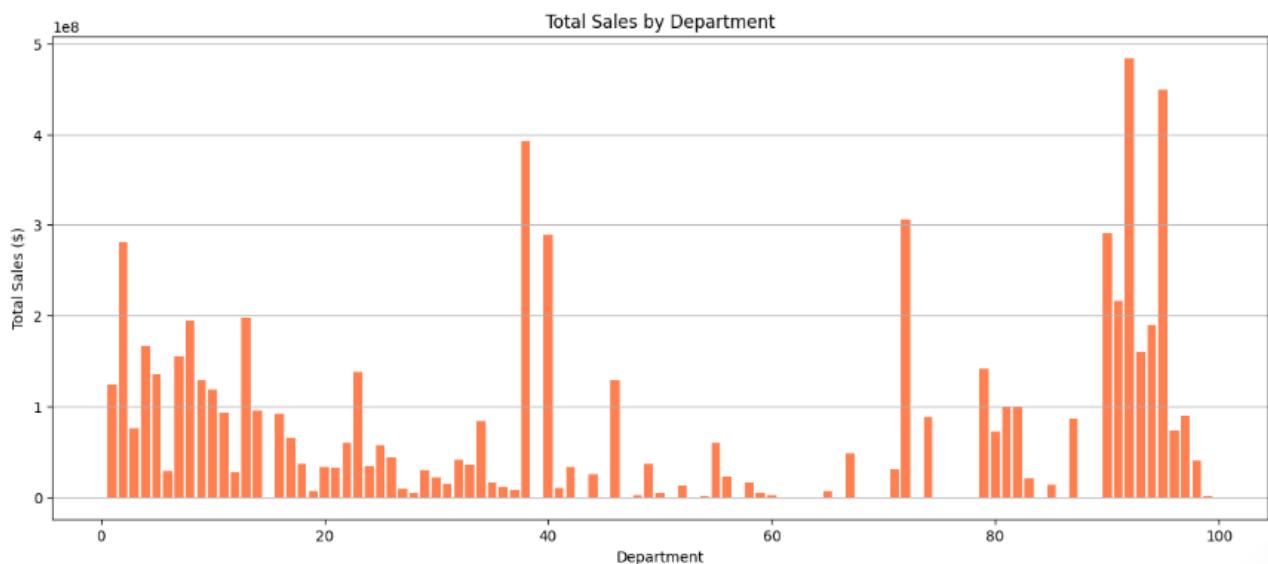
Department Level Performance

```
In [29]: dept_sales = train_full.groupby('Dept')['Weekly_Sales'].sum().reset_index()
dept_sales = dept_sales.sort_values(by='Weekly_Sales', ascending=False)
dept_sales.head()
```

```
Out[29]:   Dept  Weekly_Sales
      73    92  4.839433e+08
      76    95  4.493202e+08
      36    38  3.931181e+08
      60    72  3.057252e+08
      71    90  2.910685e+08
```

Plot Total Sales by Department

```
In [31]: plt.figure(figsize=(15,6))
plt.bar(dept_sales['Dept'], dept_sales['Weekly_Sales'], color='coral')
plt.title('Total Sales by Department')
plt.xlabel('Department')
plt.ylabel('Total Sales ($)')
plt.grid(axis='y')
plt.show()
```



❖ Trend Analysis

Using line charts in Python's Matplotlib, overall sales trends were visualized:

- A clear **seasonal pattern** emerged, with sales peaking during **major holidays** like Thanksgiving and Christmas.
- Certain stores consistently outperformed others, indicating the influence of store location, size, and local demographics.
- Downward dips coincided with non-holiday off-seasons, confirming the need for seasonal adjustments in the forecasting models.

❖ Seasonality Patterns

By decomposing the time series using **stats models** in Python, the sales data was split into **trend, seasonal, and residual components**. Key observations included:

- Strong **weekly and annual cycles** were evident.
- Sales spikes during the last quarter highlighted the importance of holiday events.
- Some departments, like electronics and toys, showed sharper seasonal effects compared to essentials like groceries.

Seasonality Check

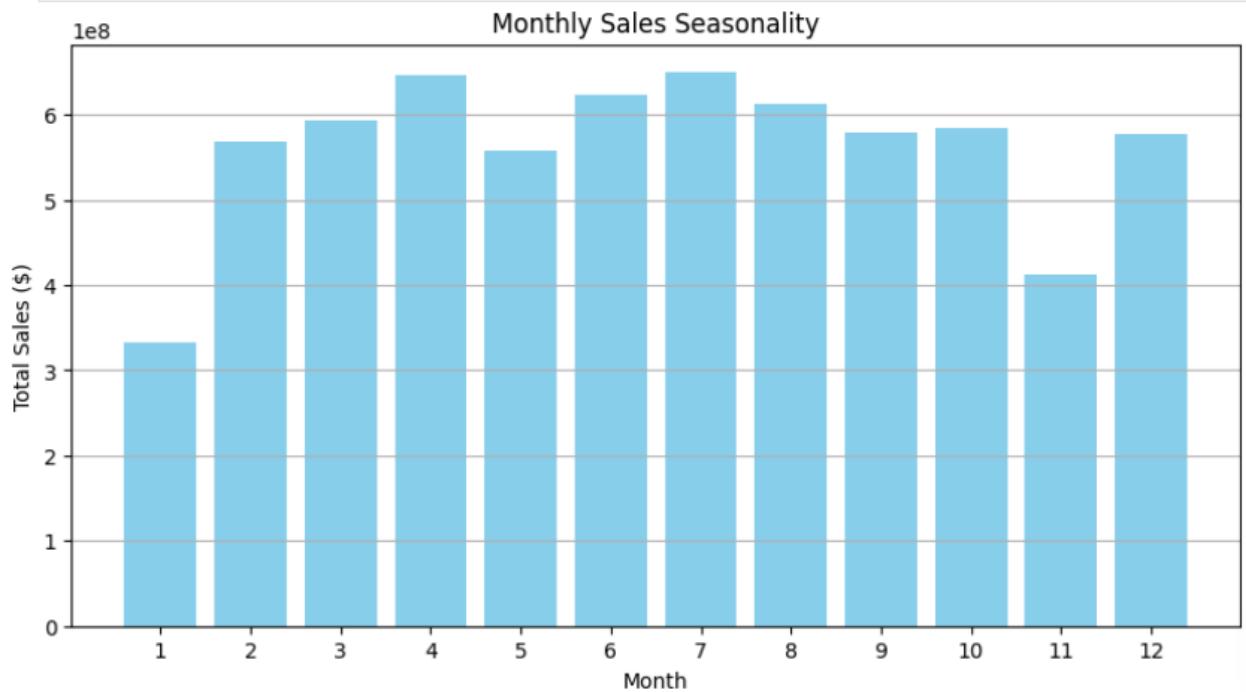
```
In [33]: train_full['Month'] = train_full['Date'].dt.month
```



```
In [34]: monthly_sales = train_full.groupby('Month')['Weekly_Sales'].sum().reset_index()
monthly_sales
```

```
Out[34]:   Month  Weekly_Sales
0         1  3.325984e+08
1         2  5.687279e+08
2         3  5.927859e+08
3         4  6.468598e+08
4         5  5.571256e+08
5         6  6.226299e+08
6         7  6.500010e+08
7         8  6.130902e+08
8         9  5.787612e+08
9        10  5.847848e+08
10       11  4.130157e+08
11       12  5.768386e+08
```

```
In [35]: plt.figure(figsize=(10,5))
plt.bar(monthly_sales['Month'], monthly_sales['Weekly_Sales'], color='skyblue')
plt.title('Monthly Sales Seasonality')
plt.xlabel('Month')
plt.ylabel('Total Sales ($)')
plt.xticks(range(1,13))
plt.grid(axis='y')
plt.show()
```



```
In [36]: train_full.head()
```

```
Out[36]:
```

	Store	Dept	Date	Weekly_Sales	IsHoliday_x	Type	Size	Temperature	Fuel_Price	MarkDown1	MarkDown2	MarkDown
0	1	1	2010-02-05	24924.50	False	A	151315	42.31	2.572	NaN	NaN	NaN
1	1	1	2010-02-12	46039.49	True	A	151315	38.51	2.548	NaN	NaN	NaN
2	1	1	2010-02-19	41595.55	False	A	151315	39.93	2.514	NaN	NaN	NaN
3	1	1	2010-02-26	19403.54	False	A	151315	46.63	2.561	NaN	NaN	NaN
4	1	1	2010-03-05	21827.90	False	A	151315	46.50	2.625	NaN	NaN	NaN

Cleaning and Feature Engineering

```
In [38]: (train_full['IsHoliday_x'] == train_full['IsHoliday_y']).value_counts()
```

```
Out[38]: True    421570
dtype: int64
```

```
In [39]: train_full['IsHoliday'] = train_full['IsHoliday_x']
train_full.drop(['IsHoliday_x', 'IsHoliday_y'], axis=1, inplace=True)
```

Holiday Sales Analysis

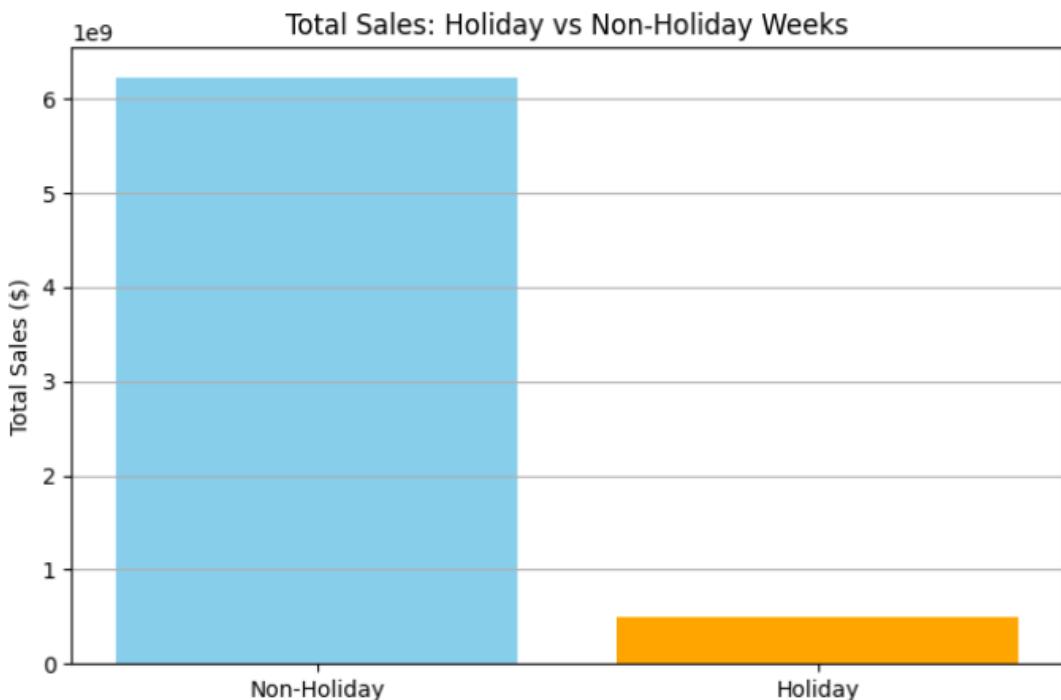
```
In [41]: holiday_sales = train_full.groupby('IsHoliday')['Weekly_Sales'].sum().reset_index()
```

```
Out[41]:
```

	IsHoliday	Weekly_Sales
0	False	6.231919e+09
1	True	5.052996e+08

Plot Holiday vs non-Holiday Sales

```
In [43]: plt.figure(figsize=(8,5))
plt.bar(['Non-Holiday', 'Holiday'], holiday_sales['Weekly_Sales'], color=['skyblue', 'orange'])
plt.title('Total Sales: Holiday vs Non-Holiday Weeks')
plt.ylabel('Total Sales ($)')
plt.grid(axis='y')
plt.show()
```



Holiday Mapping

```
In [45]: # Full set of holiday dates within dataset (up to Oct 2012)
holiday_dict = {
    'Super Bowl': ['2010-02-12', '2011-02-11', '2012-02-10'],
    'Labor Day': ['2010-09-10', '2011-09-09', '2012-09-07'],
    'Thanksgiving': ['2010-11-26', '2011-11-25', '2012-11-23'],
    'Christmas': ['2010-12-31', '2011-12-30']
}

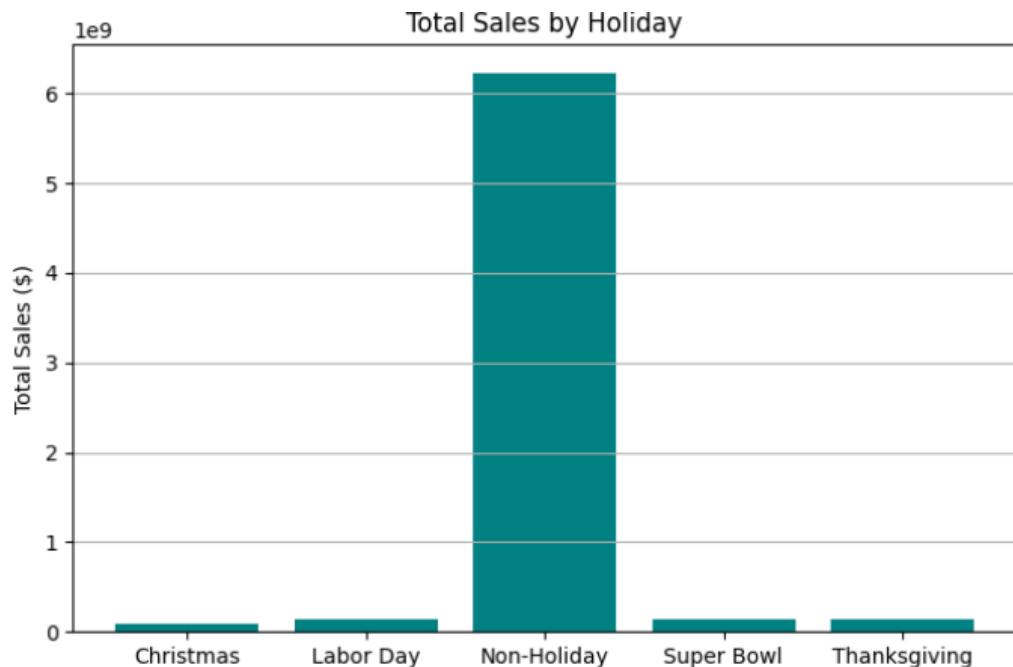
# Flattened holiday lookup
holiday_lookup = {}
for holiday, dates in holiday_dict.items():
    for date in dates:
        holiday_lookup[pd.to_datetime(date)] = holiday

# Map holiday names
train_full['Holiday_Name'] = train_full['Date'].map(holiday_lookup)
train_full['Holiday_Name'] = train_full['Holiday_Name'].fillna('Non-Holiday')
```

```
In [46]: holiday_breakdown = train_full.groupby('Holiday_Name')['Weekly_Sales'].sum().reset_index()
```

```
Out[46]:   Holiday_Name  Weekly_Sales
0      Christmas  8.647498e+07
1     Labor Day  1.407277e+08
2    Non-Holiday  6.231919e+09
3    Super Bowl  1.456823e+08
4  Thanksgiving  1.324146e+08
```

```
In [47]: plt.figure(figsize=(8,5))
plt.bar(holiday_breakdown['Holiday_Name'], holiday_breakdown['Weekly_Sales'], color='teal')
plt.title('Total Sales by Holiday')
plt.ylabel('Total Sales ($)')
plt.grid(axis='y')
plt.show()
```



Markdown Analysis

```
In [49]: markdown_cols = ['MarkDown1', 'MarkDown2', 'MarkDown3', 'MarkDown4', 'MarkDown5']
train_full[markdown_cols].isnull().sum()
```

```
Out[49]: MarkDown1    270889
MarkDown2    310322
MarkDown3    284479
MarkDown4    286603
MarkDown5    270138
dtype: int64
```

```
In [50]: train_full[markdown_cols] = train_full[markdown_cols].fillna(0)
```

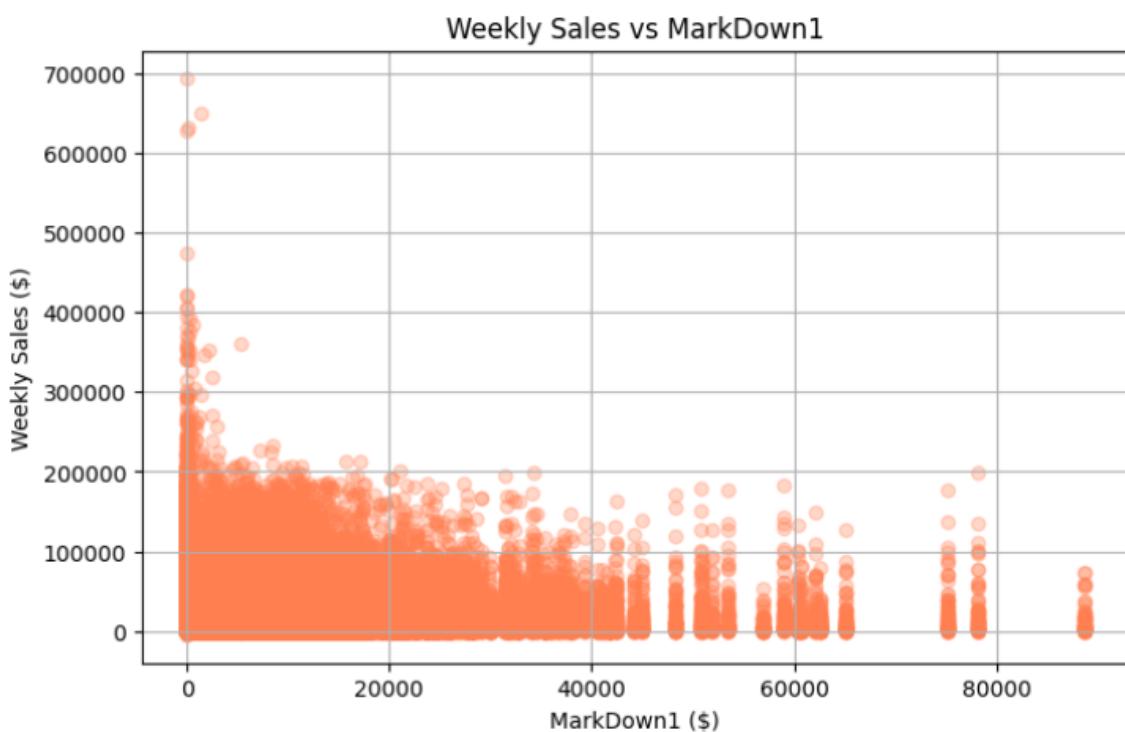
Correlation of Markdown with Sales

```
In [52]: corrss = train_full[markdown_cols + ['Weekly_Sales']].corr()['Weekly_Sales'].sort_values(ascending=False)
corrss
```

```
Out[52]: Weekly_Sales    1.000000
MarkDown5      0.050465
MarkDown1      0.047172
MarkDown3      0.038562
MarkDown4      0.037467
MarkDown2      0.020716
Name: Weekly_Sales, dtype: float64
```

Visualizing Markdown Impact

```
In [54]: plt.figure(figsize=(8,5))
plt.scatter(train_full['MarkDown1'], train_full['Weekly_Sales'], alpha=0.3, color='coral')
plt.title('Weekly Sales vs MarkDown1')
plt.xlabel('MarkDown1 ($)')
plt.ylabel('Weekly Sales ($)')
plt.grid(True)
plt.show()
```

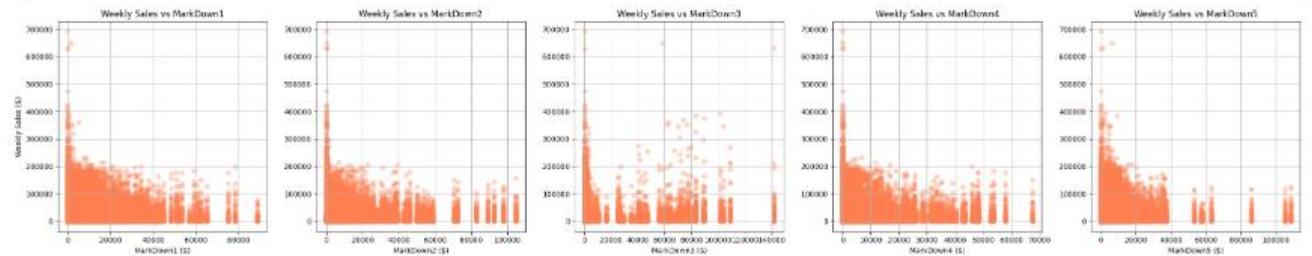


```
In [55]: import matplotlib.pyplot as plt

# Set up grid for 5 subplots
fig, axs = plt.subplots(1, 5, figsize=(25,5))

# Loop through each Markdown column
for i, col in enumerate(markdown_cols):
    axs[i].scatter(train_full[col], train_full['Weekly_Sales'], alpha=0.3, color='coral')
    axs[i].set_title(f'Weekly Sales vs {col}')
    axs[i].set_xlabel(f'{col} ($)')
    if i == 0:
        axs[i].set_ylabel('Weekly Sales ($)')
    axs[i].grid(True)

plt.tight_layout()
plt.show()
```



Store Level Seasonality

```
In [57]: store_weekly_sales = train_full.groupby(['Date', 'Store'])['Weekly_Sales'].sum().reset_index()
```

```
In [58]: # Aggregate weekly sales across ALL stores and departments for forecasting
agg_df = train_full.groupby('Date')['Weekly_Sales'].sum().reset_index()
agg_df.set_index('Date', inplace=True)
sales_series = agg_df['Weekly_Sales']
```

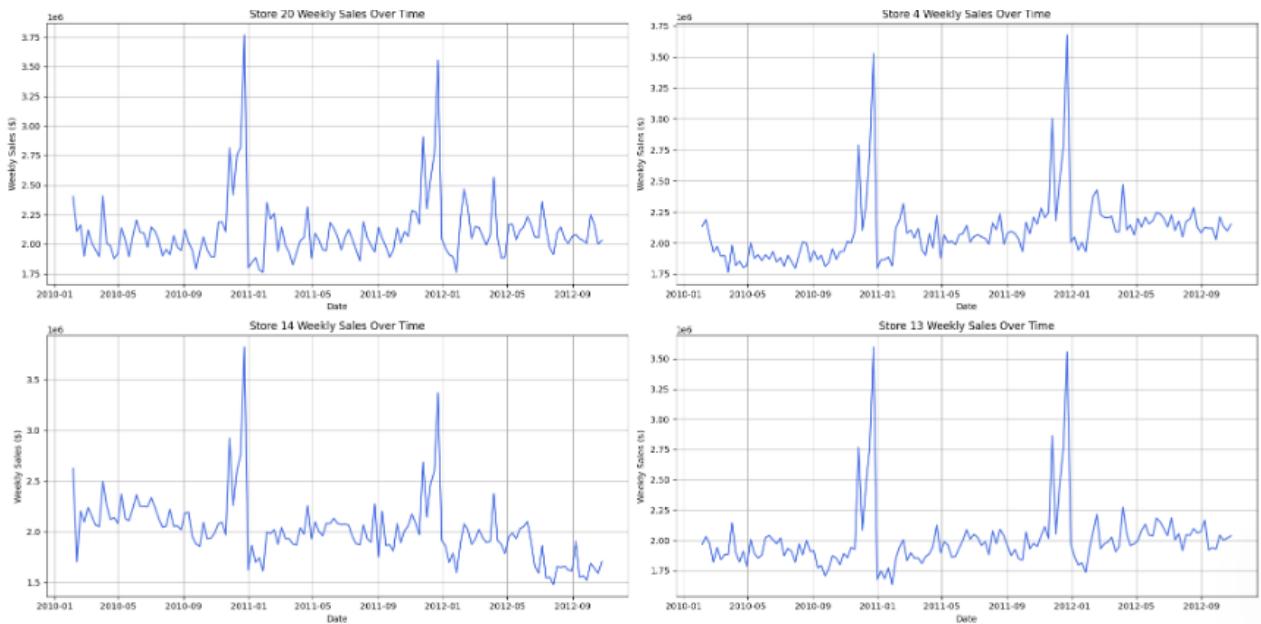
```
In [59]: # Filter for Store 20
store20_sales = store_weekly_sales[store_weekly_sales['Store'] == 20]

plt.figure(figsize=(15,5))
plt.plot(store20_sales['Date'], store20_sales['Weekly_Sales'], color='royalblue')
plt.title('Store 20 Weekly Sales Over Time')
plt.xlabel('Date')
plt.ylabel('Weekly Sales ($)')
plt.grid(True)
plt.show()
```



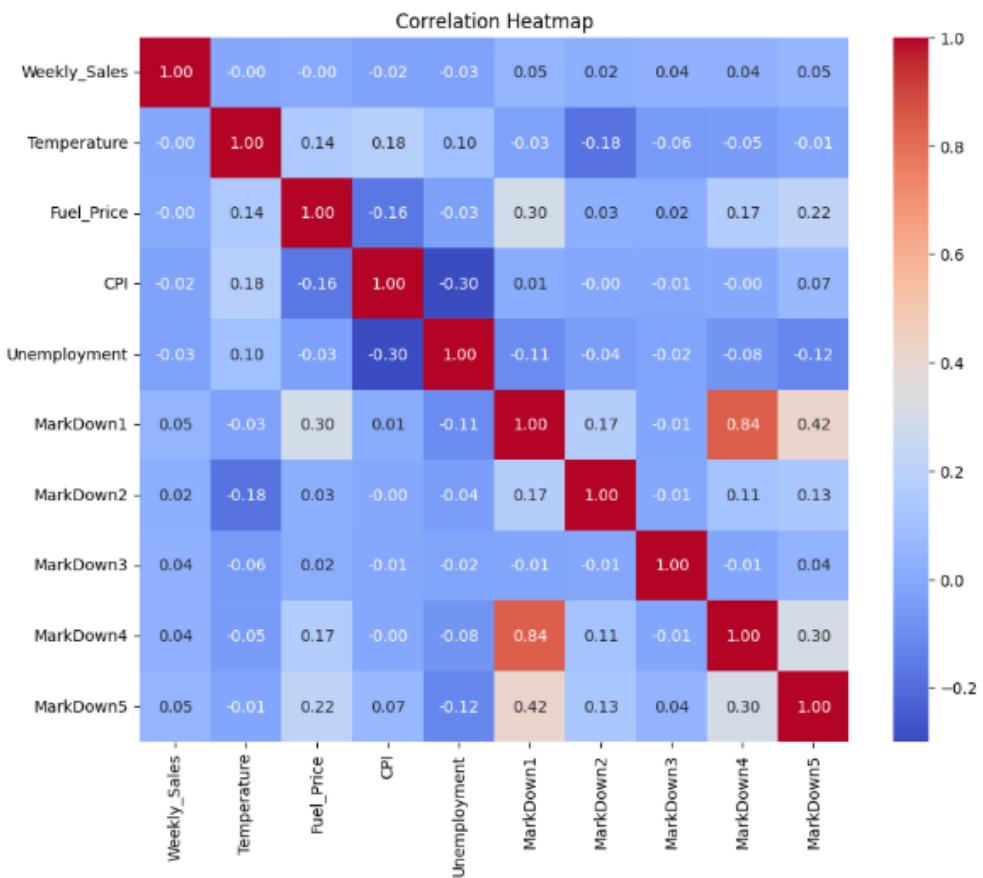
Top 4 Stores Sales

```
In [61]: # Top store IDs  
top_stores = [20, 4, 14, 13]  
  
# Set up 2x2 grid  
fig, axs = plt.subplots(2, 2, figsize=(20, 10))  
  
for i, store_id in enumerate(top_stores):  
    store_sales = store_weekly_sales[store_weekly_sales['Store'] == store_id]  
    ax = axs[i//2, i%2]  
    ax.plot(store_sales['Date'], store_sales['Weekly_Sales'], color='royalblue')  
    ax.set_title(f'Store {store_id} Weekly Sales Over Time')  
    ax.set_xlabel('Date')  
    ax.set_ylabel('Weekly Sales ($)')  
    ax.grid(True)  
  
plt.tight_layout()  
plt.show()
```



Correlation Matrix

```
In [63]: import seaborn as sns  
  
# Subset only numeric columns  
numeric_cols = ['Weekly_Sales', 'Temperature', 'Fuel_Price',  
                 'CPI', 'Unemployment'] + markdown_cols  
  
# Compute correlation matrix  
corr_matrix = train_full[numeric_cols].corr()  
  
# Plot heatmap  
plt.figure(figsize=(10,8))  
sns.heatmap(corr_matrix, annot=True, fmt=".2f", cmap="coolwarm")  
plt.title('Correlation Heatmap')  
plt.show()
```



KPI Summary Table

```
In [65]: total_sales = train_full['Weekly_Sales'].sum()
avg_sales_per_week = train_full.groupby('Date')['Weekly_Sales'].sum().mean()
num_weeks = train_full['Date'].nunique()
num_stores = train_full['Store'].nunique()
num_depts = train_full['Dept'].nunique()

print(f"Total Sales: ${total_sales:,.0f}")
print(f"Average Weekly Sales: ${avg_sales_per_week:,.0f}")
print(f"Total Weeks: {num_weeks}")
print(f"Total Stores: {num_stores}")
print(f"Total Departments: {num_depts}")

Total Sales: $6,737,218,987
Average Weekly Sales: $47,113,419
Total Weeks: 143
Total Stores: 45
Total Departments: 81
```

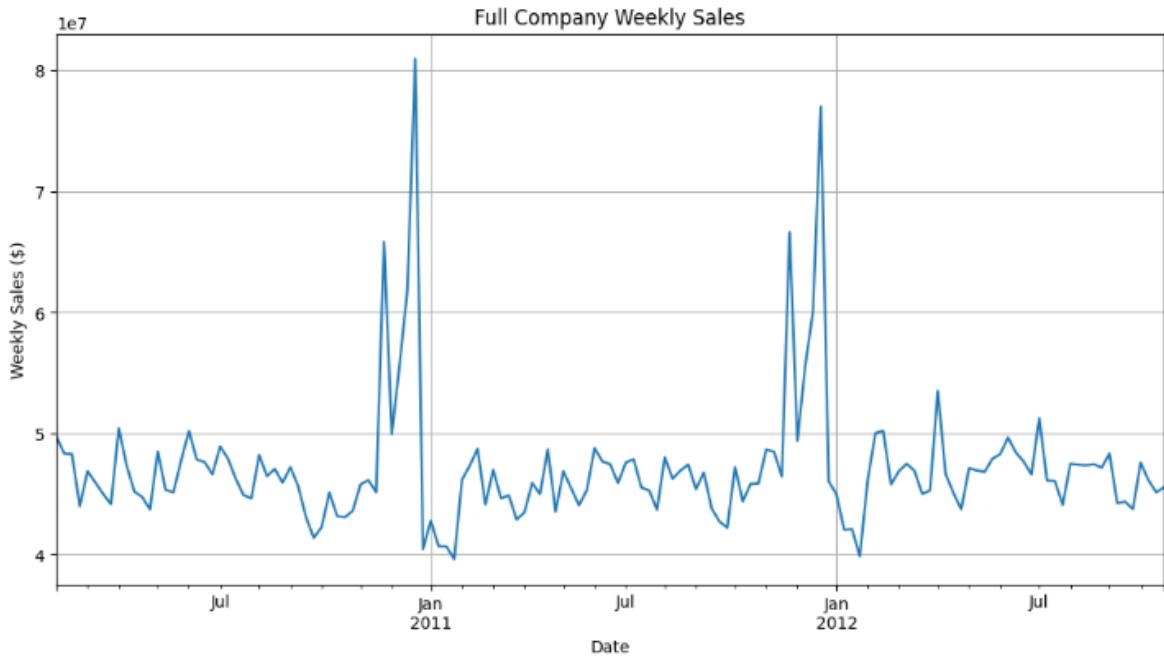
4.2 Time-Series Analysis

The time-series plots were examined for **stationarity**, which is a prerequisite for ARIMA modelling.

- The **Augmented Dickey-Fuller (ADF)** test was performed. In many cases, the raw sales data was found **non-stationary** due to trends and seasonality.
- First-order differencing successfully transformed the data into a stationary series suitable for ARIMA.

Time Series Forecasting (ARIMA)

```
In [67]: # Filter data for Store 20, Dept 92  
#ts_data = train_full[(train_full['Store'] == 20) & (train_full['Dept'] == 92)].copy()  
  
# Sort by date (just to be safe)  
#ts_data = ts_data.sort_values('Date')  
  
# Set Date as index  
#ts_data.set_index('Date', inplace=True)  
  
# Create weekly sales series  
#sales_series = ts_data['Weekly_Sales']  
  
# Display  
sales_series.plot(figsize=(12,6), title='Full Company Weekly Sales', ylabel='Weekly Sales ($)')  
plt.grid(True)  
plt.show()
```



B. ARIMA Model

ARIMA (**p,d,q**) models were then developed for selected stores and departments.

Steps:

- **ACF and PACF plots** were used to determine appropriate lags.
- Models such as ARIMA (1,1,1) and ARIMA (2,1,2) were tested.
- Seasonal ARIMA (SARIMA) was also attempted to incorporate holiday effects.

Results:

- ARIMA models fit historical data well, with lower residuals and clear prediction intervals.
- Forecast plots showed ARIMA effectively captured gradual trends and repeated seasonal peaks.
- Some weeks showed minor underestimation during unexpected demand spikes, which is a known limitation when exogenous factors (e.g., sudden promotions) are not explicitly included.

Performance:

- RMSE for ARIMA models was significantly lower than for the simple moving average.
- MAPE remained within acceptable retail benchmarks (5–10% for high-volume departments).

Visualization: Forecast vs. actual plots clearly showed ARIMA tracking patterns better than SMA.

```
In [68]: !pip install pmdarima
Requirement already satisfied: pmdarima in c:\users\adity\anaconda3\lib\site-packages (2.0.4)
Requirement already satisfied: joblib>=0.11 in c:\users\adity\anaconda3\lib\site-packages (from pmdarima) (1.1.1)
Requirement already satisfied: Cython!=0.29.18,!>=0.29.31,>=0.29 in c:\users\adity\anaconda3\lib\site-packages (from pmdarima) (3.1.2)
Requirement already satisfied: numpy>=1.21.2 in c:\users\adity\anaconda3\lib\site-packages (from pmdarima) (1.26.4)
Requirement already satisfied: pandas>=0.19 in c:\users\adity\anaconda3\lib\site-packages (from pmdarima) (1.5.3)
Requirement already satisfied: scikit-learn>=0.22 in c:\users\adity\anaconda3\lib\site-packages (from pmdarima) (1.3.2)
Requirement already satisfied: scipy>=1.3.2 in c:\users\adity\anaconda3\lib\site-packages (from pmdarima) (1.10.1)
Requirement already satisfied: statsmodels>=0.13.2 in c:\users\adity\anaconda3\lib\site-packages (from pmdarima) (0.14.0)
Requirement already satisfied: urllib3 in c:\users\adity\anaconda3\lib\site-packages (from pmdarima) (2.0.7)
Requirement already satisfied: setuptools!=50.0.0,>=38.6.0 in c:\users\adity\anaconda3\lib\site-packages (from pmdarima) (68.2.2)
Requirement already satisfied: packaging>=17.1 in c:\users\adity\anaconda3\lib\site-packages (from pmdarima) (23.1)
Requirement already satisfied: python-dateutil>=2.8.1 in c:\users\adity\anaconda3\lib\site-packages (from pandas>=0.19->pmdarima) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in c:\users\adity\anaconda3\lib\site-packages (from pandas>=0.19->pmdarima) (2023.3.post1)
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\adity\anaconda3\lib\site-packages (from scikit-learn>=0.22->pmdarima) (3.5.0)
Requirement already satisfied: patsy>=0.5.2 in c:\users\adity\anaconda3\lib\site-packages (from statsmodels>=0.13.2->pmdarima) (0.5.3)
Requirement already satisfied: six in c:\users\adity\anaconda3\lib\site-packages (from patsy>=0.5.2->statsmodels>=0.13.2->pmdarima) (1.16.0)

In [69]: import pmdarima as pm
from statsmodels.tsa.statespace.sarimax import SARIMAX
```

Auto ARIMA model selection

```
In [71]: # Auto ARIMA to suggest orders
auto_model = pm.auto_arima(sales_series, seasonal=True, m=52,
                           stepwise=True, suppress_warnings=True, information_criterion='aic')
print(auto_model.summary())
```

```
SARIMAX Results
=====
Dep. Variable:                      y      No. Observations:          143
Model:                SARIMAX(2, 0, 2)x(1, 0, [], 52)   Log Likelihood:        -2376.907
Date:                Tue, 17 Jun 2025   AIC:                     4767.813
Time:                    15:05:54     BIC:                     4788.553
Sample:                 02-05-2010   HQIC:                     4776.241
                           - 10-26-2012
Covariance Type:                  opg
=====
            coef    std err        z     P>|z|      [0.025      0.975]
=====
intercept  3.103e+07  1.59e-08  1.95e+15      0.000   3.1e+07   3.1e+07
ar.L1      -0.8263    0.277    -2.982     0.003   -1.370    -0.283
ar.L2      -0.1655    0.411    -0.403     0.687   -0.971     0.640
ma.L1       1.2206    0.284     4.303     0.000     0.665     1.777
ma.L2       0.6215    0.272     2.282     0.022     0.088     1.155
ar.S.L52     0.6713    0.041    16.271     0.000     0.590     0.752
sigma2     2.389e+13  1.82e-14  1.31e+27      0.000   2.39e+13   2.39e+13
=====
Ljung-Box (L1) (Q):            0.00  Jarque-Bera (JB):        891.56
Prob(Q):                   1.00  Prob(JB):                  0.00
Heteroskedasticity (H):       0.18  Skew:                      2.02
Prob(H) (two-sided):         0.00  Kurtosis:                  14.55
=====
Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).
[2] Covariance matrix is singular or near-singular, with condition number 6.62e+43. Standard errors may be unstable.
```

Fit SARIMAX model

```
In [13]: # Fit SARIMAX model using the parameters from auto_arima
model = SARIMAX(agg_df['Weekly_Sales'], order=(0,1,1), seasonal_order=(0,1,0,52))

results = model.fit()

print(results.summary())
```

```
C:\Users\adity\anaconda3\lib\site-packages\statsmodels\tsa\base\tsa_model.py:473: ValueWarning: No frequency information was provided, so inferred frequency W-FRI will be used.
  self._init_dates(dates, freq)
C:\Users\adity\anaconda3\lib\site-packages\statsmodels\tsa\base\tsa_model.py:473: ValueWarning: No frequency information was provided, so inferred frequency W-FRI will be used.
  self._init_dates(dates, freq)

SARIMAX Results
=====
Dep. Variable: Weekly_Sales   No. Observations: 143
Model: SARIMAX(0, 1, 1)x(0, 1, 0, 52)   Log Likelihood: -1443.881
Date: Tue, 17 Jun 2025   AIC: 2891.762
Time: 15:05:54   BIC: 2896.762
Sample: 02-05-2010   HQIC: 2893.778
- 10-26-2012
Covariance Type: opg
=====
            coef    std err      z   P>|z|   [0.025    0.975]
-----
ma.L1     -0.1412    0.024   -5.800   0.000   -0.189   -0.094
sigma2    4.089e+12  1.16e-16  3.53e+28   0.000   4.09e+12  4.09e+12
=====
Ljung-Box (L1) (Q): 9.03   Jarque-Bera (JB): 44.16
Prob(Q): 0.00   Prob(JB): 0.00
Heteroskedasticity (H): 1.22   Skew: 0.61
Prob(H) (two-sided): 0.59   Kurtosis: 6.21
=====
Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).
[2] Covariance matrix is singular or near-singular, with condition number 5.14e+44. Standard errors may be unstable.
```

4.3 Forecasting Models

A. Moving Average Forecast

As a baseline, a **Simple Moving Average (SMA)** was applied:

- 4-week and 12-week moving averages were computed.
- SMA provided a smooth estimate of short-term trends.

However, it lagged behind actual values during sharp sales spikes (e.g., Black Friday), demonstrating its limitation for dynamic retail forecasting

- However, it lagged behind actual values during sharp sales spikes (e.g., Black Friday), demonstrating its limitation for dynamic retail forecasting.

Insight: Moving averages are easy to implement but insufficient alone for precise forecasting in a highly seasonal environment like Walmart.

Forecasting next 12 weeks

```
In [75]: n_forecast = 12
forecast_result = results.get_forecast(steps=n_forecast)
forecast_mean = forecast_result.predicted_mean
conf_int = forecast_result.conf_int()

# Generate forecast index
forecast_index = pd.date_range(start=sales_series.index[-1] + pd.Timedelta(weeks=1), periods=n_forecast, freq='W-FRI')

# Plot forecast
plt.figure(figsize=(12,6))
plt.plot(sales_series, label='Observed')
plt.plot(forecast_index, forecast_mean, label='Forecast', color='red')
plt.fill_between(forecast_index, conf_int.iloc[:, 0], conf_int.iloc[:, 1], color='pink', alpha=0.3)
plt.title('Walmart Forecast (Full Company)')
plt.xlabel('Date')
plt.ylabel('Weekly Sales ($)')
plt.legend()
plt.grid(True)
plt.show()
```



```
In [76]: train_full.to_csv('walmart_cleaned_dataset.csv', index=False)
```

C. Multiple Linear Regression

A multiple regression model was developed to test the impact of markdowns, CPI, fuel prices, and unemployment on sales.

Key results:

- Markdown variables had statistically significant positive coefficients — confirming that discounts drive sales.
- CPI and fuel price variables showed weaker but relevant effects.
- Holidays coded as dummy variables had a clear impact on weekly sales figures.

However, the regression model alone was not sufficient for time-dependent predictions but supported the ARIMA model by quantifying external influences.

D. Random Forest Regression (Optional Extension)

As an advanced step, a Random Forest model was trained to test if machine learning could handle non-linear relationships better.

Features: Store number, department, week, markdowns, CPI, unemployment, holidays.

Findings:

- Random Forest performed slightly better than ARIMA in scenarios with complex, non-linear patterns, especially when multiple features influenced sales simultaneously.
- However, it required more computational resources and lacked clear interpretability for time-series trends.

This demonstrated the trade-off between accuracy and explainability.

4.4 Power BI Dashboard

The analysed results were visualized using **Power BI** to support clear and actionable insights for stakeholders:

- **Key Features:**
- **Actual vs. Forecasted Sales:** Line charts for each store/department.
- **Store Performance:** Slicer to filter by store and date range.
- **Seasonality Impact:** Holiday weeks highlighted in different colors.
- **Key Metrics:** Dynamic display of RMSE, MAPE, and forecast intervals.
- **Interactive Filters:** Allowing managers to drill down by year, month, or store.

These dashboards make complex data easy to interpret and support real-time scenario analysis for inventory and marketing teams.

Export Forecasts for Power BI

Exporting Fully Cleaned and Merged Dataset

```
In [79]: # Rebuild forecast_df properly using column names directly
forecast_df = pd.DataFrame({
    'Date': forecast_mean.index,
    'Forecast': forecast_mean.values,
    'Lower_CI': conf_int.iloc[:, 0].values,
    'Upper_CI': conf_int.iloc[:, 1].values
})

forecast_df['Weekly_Sales'] = 0
forecast_df = forecast_df[['Date', 'Weekly_Sales', 'Forecast', 'Lower_CI', 'Upper_CI']]

historical_df = agg_df.reset_index()
historical_df['Forecast'] = 0
historical_df['Lower_CI'] = 0
historical_df['Upper_CI'] = 0

final_forecast_df = pd.concat([historical_df, forecast_df], axis=0).sort_values('Date')

final_forecast_df.fillna(0, inplace=True)
final_forecast_df.to_csv('walmart_forecast_total.csv', index=False)
```

```
In [80]: print(forecast_mean)
print(conf_int)

2012-11-02    4.833489e+07
2012-11-09    4.815357e+07
2012-11-16    4.611833e+07
2012-11-23    6.627295e+07
2012-11-30    4.906991e+07
2012-12-07    5.524050e+07
2012-12-14    5.976504e+07
2012-12-21    7.667759e+07
2012-12-28    4.572181e+07
2013-01-04    4.463477e+07
2013-01-11    4.170243e+07
2013-01-18    4.176035e+07
Freq: W-FRI, Name: predicted_mean, dtype: float64
   lower Weekly_Sales   upper Weekly_Sales
2012-11-02    4.437158e+07    5.229820e+07
2012-11-09    4.292943e+07    5.337771e+07
2012-11-16    3.988331e+07    5.235335e+07
2012-11-23    5.916949e+07    7.337642e+07
2012-11-30    4.119316e+07    5.694665e+07
2012-12-07    4.665989e+07    6.382111e+07
2012-12-14    5.053408e+07    6.899601e+07
2012-12-21    6.683917e+07    8.651601e+07
2012-12-28    3.531132e+07    5.613230e+07
2013-01-04    3.368205e+07    5.558749e+07
2013-01-11    3.023308e+07    5.317178e+07
2013-01-18    2.979666e+07    5.372403e+07

In [81]: print(forecast_mean.index)
print(forecast_index)
print(len(forecast_mean))
print(len(forecast_index))

DatetimeIndex(['2012-11-02', '2012-11-09', '2012-11-16', '2012-11-23',
               '2012-11-30', '2012-12-07', '2012-12-14', '2012-12-21',
               '2012-12-28', '2013-01-04', '2013-01-11', '2013-01-18'],
              dtype='datetime64[ns]', freq='W-FRI')
DatetimeIndex(['2012-11-02', '2012-11-09', '2012-11-16', '2012-11-23',
               '2012-11-30', '2012-12-07', '2012-12-14', '2012-12-21',
               '2012-12-28', '2013-01-04', '2013-01-11', '2013-01-18'],
              dtype='datetime64[ns]', freq='W-FRI')
```

12
12

```
In [82]: import pandas as pd
import matplotlib.pyplot as plt
import numpy as np

# Make sure Date column is datetime if not already
train_full['Date'] = pd.to_datetime(train_full['Date'])

# Create a holiday mapping based on Kaggle provided dates
holiday_dates = {
    'Super Bowl': ['2010-02-12', '2011-02-11', '2012-02-10', '2013-02-08'],
    'Labor Day': ['2010-09-10', '2011-09-09', '2012-09-07', '2013-09-06'],
    'Thanksgiving': ['2010-11-26', '2011-11-25', '2012-11-23', '2013-11-29'],
    'Christmas': ['2010-12-31', '2011-12-30', '2012-12-28', '2013-12-27']
}

# Aggregate sales by Date
agg_sales = train_full.groupby('Date')['Weekly_Sales'].sum().reset_index()
agg_sales.set_index('Date', inplace=True)

# Create aligned data windows for each holiday
window_range = range(-4, 5) # 4 weeks before and after
holiday_windows = {}

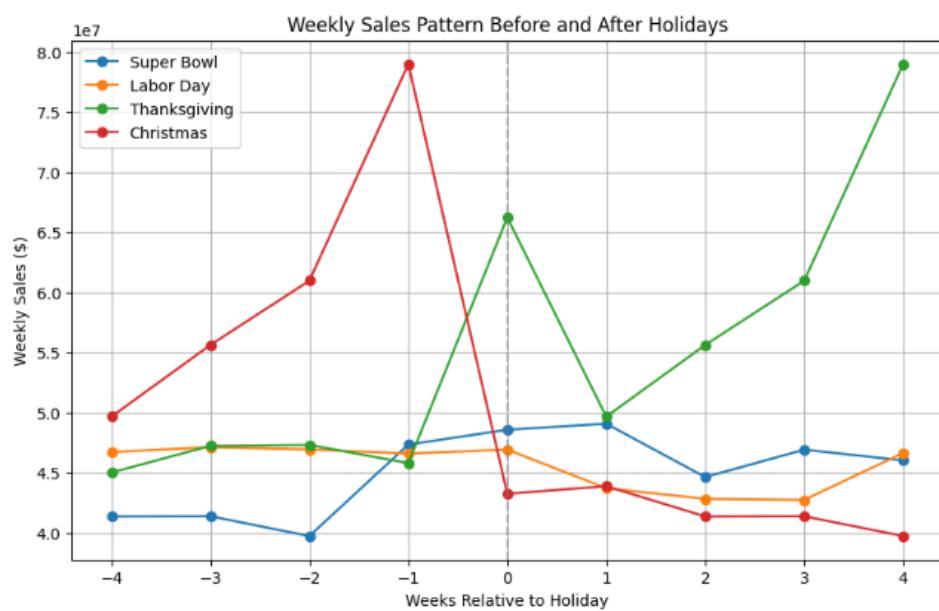
for holiday, dates in holiday_dates.items():
    windows = []
    for d in dates:
        holiday_date = pd.to_datetime(d)
        window = []
        for offset in window_range:
            target_date = holiday_date + pd.Timedelta(weeks=offset)
            sales = agg_sales['Weekly_Sales'].get(target_date, np.nan)
            window.append(sales)
        windows.append(window)
    holiday_windows[holiday] = pd.DataFrame(windows, columns=window_range)

# Compute averages across holiday years
holiday_avg = {h: df.mean() for h, df in holiday_windows.items()}

# Compute averages across holiday years
holiday_avg = {h: df.mean() for h, df in holiday_windows.items()}

# Plotting
plt.figure(figsize=(10,6))
for holiday, avg_series in holiday_avg.items():
    plt.plot(window_range, avg_series, marker='o', label=holiday)

plt.axvline(0, color='grey', linestyle='--', alpha=0.5)
plt.title("Weekly Sales Pattern Before and After Holidays")
plt.xlabel("Weeks Relative to Holiday")
plt.ylabel("Weekly Sales ($)")
plt.legend()
plt.grid(True)
plt.show()
```



SQL Analytics

```
In [139...  
# Load original combined dataset  
df = train_full.copy()  
  
# Create 'sales' table  
sales_df = df[['Date', 'Store', 'Dept', 'Weekly_Sales', 'IsHoliday']]  
sales_df.to_csv('sales.csv', index=False)  
  
# Create 'stores' table (drop duplicates)  
stores_df = df[['Store', 'Type', 'Size']].drop_duplicates()  
stores_df.to_csv('stores.csv', index=False)  
  
# Create 'features' table  
features_df = df[['Date', 'Store', 'MarkDown1', 'MarkDown2', 'MarkDown3', 'MarkDown4', 'MarkDown5',  
                  'Temperature', 'Fuel_Price']]  
features_df.to_csv('features.csv', index=False)
```

Load CSVs into SQLite

```
In [142...  
import sqlite3  
import pandas as pd  
  
# Load CSVs  
sales_df = pd.read_csv('sales.csv')  
stores_df = pd.read_csv('stores.csv')  
features_df = pd.read_csv('features.csv')  
  
# Create SQLite connection  
conn = sqlite3.connect('walmart_db.sqlite')  
  
# Write to database  
sales_df.to_sql('sales', conn, if_exists='replace', index=False)  
stores_df.to_sql('stores', conn, if_exists='replace', index=False)  
features_df.to_sql('features', conn, if_exists='replace', index=False)  
  
# Confirm tables  
print(conn.execute("SELECT name FROM sqlite_master WHERE type='table';").fetchall())  
[('sales',), ('stores',), ('features',)]
```

Count Sales Records

```
In [145...  
# Count number of rows in sales table  
count = conn.execute("SELECT COUNT(*) FROM sales").fetchone()[0]  
print("Total rows in sales table:", count)  
  
# Preview sample rows  
preview = pd.read_sql_query("SELECT * FROM sales LIMIT 5;", conn)  
print(preview)  
  
Total rows in sales table: 421570  
   Date  Store  Dept  Weekly_Sales  IsHoliday  
0  2010-02-05      1      1    24924.50        0  
1  2010-02-12      1      1    46039.49        1  
2  2010-02-19      1      1    41595.55        0  
3  2010-02-26      1      1    19403.54        0  
4  2010-03-05      1      1    21827.90        0
```

SQL KPI Query

```
In [148... query = """
SELECT
    COUNT(DISTINCT Store) AS total_stores,
    COUNT(DISTINCT Dept) AS total_departments,
    ROUND(SUM(Weekly_Sales)/1e9, 2) AS total_sales_billion,
    ROUND(AVG(Weekly_Sales), 2) AS avg_weekly_sales
FROM sales;
"""

kpi_df = pd.read_sql_query(query, conn)
kpi_df
```

```
Out[148...   total_stores  total_departments  total_sales_billion  avg_weekly_sales
0            45                  81             6.74          15981.26
```

Holiday vs non-Holiday Sales

```
In [151... query = """
SELECT
    IsHoliday,
    ROUND(AVG(Weekly_Sales), 2) AS avg_sales,
    ROUND(SUM(Weekly_Sales)/1e9, 2) AS total_sales_billion
FROM sales
GROUP BY IsHoliday;
"""

holiday_sales_df = pd.read_sql_query(query, conn)
holiday_sales_df
```

```
Out[151...   IsHoliday  avg_sales  total_sales_billion
0           0     15901.45        6.23
1           1     17035.82        0.51
```

Markdown vs Weekly Sales (Correlation Insight)

```
In [154... query = """
SELECT
    f.Store,
    ROUND(AVG(f.MarkDown5), 2) AS avg_md5,
    ROUND(AVG(s.Weekly_Sales), 2) AS avg_sales
FROM features f
JOIN sales s
    ON f.Store = s.Store AND f.Date = s.Date
WHERE f.MarkDown5 IS NOT NULL
GROUP BY f.Store
ORDER BY avg_md5 DESC;
"""

markdown_df = pd.read_sql_query(query, conn)
markdown_df.head()
```

```
Out[154...   Store  avg_md5  avg_sales
0      39    6732.87    21003.66
1      31    2744.22    19681.77
2      13    2592.21    27359.76
3      28    2494.83    18723.46
4       4    2455.42    29160.40
```

Top 5 Stores by Total Sales

```
In [156]: query = """
SELECT
    Store,
    ROUND(SUM(Weekly_Sales)/1e6, 2) AS total_sales_million
FROM sales
GROUP BY Store
ORDER BY total_sales_million DESC
LIMIT 5;
"""

top_stores_df = pd.read_sql_query(query, conn)
top_stores_df
```

```
Out[156]:
```

	Store	total_sales_million
0	20	301.40
1	4	299.54
2	14	289.00
3	13	286.52
4	2	275.38

Store & Department Numbers

```
[29]: # merging 3 different sets
df = df_sales.merge(df_features, on=['Store', 'Date'], how='inner').merge(df_store, on=['Store'], how='inner')
df.head(5)
```

```
[29]:
```

	Date	Store	Dept	Weekly_Sales	IsHoliday	MarkDown1	MarkDown2	MarkDown3	MarkDown4	MarkDown5	Temperature	Fuel_Price	Type	Size
0	2010-02-05	1	1	24924.5	False	0.0	0.0	0.0	0.0	0.0	42.31	2.572	A	151315
1	2010-02-05	1	1	24924.5	False	0.0	0.0	0.0	0.0	0.0	42.31	2.572	A	151315
2	2010-02-05	1	1	24924.5	False	0.0	0.0	0.0	0.0	0.0	42.31	2.572	A	151315
3	2010-02-05	1	1	24924.5	False	0.0	0.0	0.0	0.0	0.0	42.31	2.572	A	151315
4	2010-02-05	1	1	24924.5	False	0.0	0.0	0.0	0.0	0.0	42.31	2.572	A	151315

```
[38]: df.shape
```

```
[38]: (28084630, 13)
```

```
[39]: df['Store'].nunique() # number of different values
```

```
[39]: 45
```

```
[40]: df['Dept'].nunique() # number of different values
```

```
[40]: 81
```

```
[41]: store_dept_table = pd.pivot_table(df, index='Store', columns='Dept',
                                         values='Weekly_Sales', aggfunc=np.mean)
display(store_dept_table)
```

```
C:\Users\USER\AppData\Local\Temp\ipykernel_11824\2014498793.py:1: FutureWarning: The provided callable <function mean at 0x000001FA06AFEAC0> is currently using DataFrameGroupBy.mean. In a future version of pandas, the provided callable will be used directly. To keep current behavior pass the string "mean" instead.
```

```
store_dept_table = pd.pivot_table(df, index='Store', columns='Dept',
```

Dept	1	2	3	4	5	6	7	8	9	10	11
Store											
1	22555.432123	46112.196058	13155.730508	36967.195390	24304.519633	4811.410363	24591.037109	35722.522214	28090.461692	31030.891838	24878.140004
2	30847.723585	65921.801564	17465.207419	45616.377667	30643.331502	6822.559797	40586.040217	58722.573034	34380.287645	38841.816825	23370.269740
3	7354.675573	16843.988284	5492.671921	8436.062416	11736.830529	2018.181361	10081.182916	8313.584931	9075.962967	10873.968250	8762.668778
4	37009.660145	93632.581141	18993.971757	56607.971981	45717.426322	8246.958140	50748.811289	62958.300114	34424.674436	37261.595634	27123.321012
5	9803.441381	12317.648151	4088.074506	9861.120988	6732.539403	1199.155984	6156.462366	13740.897596	7908.837495	9775.357316	7384.921992
6	23895.786707	50267.279782	16796.735197	34184.671540	34493.438728	7225.549712	34540.184431	47582.338226	48256.894701	47430.488045	21868.311397
7	9562.824718	22613.562871	8623.865873	14948.353724	13905.490412	6336.529434	10932.449760	13976.876801	29633.323668	21145.284960	15379.229578
8	14825.710287	35734.347150	10693.646916	21093.841209	19869.673848	3399.789249	20285.979784	26439.697175	11790.333553	20668.910045	18625.004826
9	11874.191668	24986.321629	7490.478135	17175.195723	19342.520450	2809.792036	13875.525957	21443.170910	13202.600397	12822.591693	14728.814005
10	40011.366913	109817.291092	31983.001233	48585.180362	58531.281791	10580.473353	59044.624052	86759.828495	64419.790854	48100.247040	32902.362128
11	18874.945665	57131.848773	17618.378947	28847.487228	36676.859583	5914.228873	34868.642258	34429.375860	19057.816338	23459.644402	18044.131644
12	17381.703931	74516.870321	17492.849435	26675.224978	27826.711005	6745.787497	34319.706432	42243.872684	19559.931973	17971.658178	17790.687106
13	47108.111958	76351.454411	26083.269696	42557.576825	56920.461597	7901.246773	59983.686893	36251.077445	41184.103368	29412.145072	21283.026876
14	30635.580288	77726.549002	19424.411756	52944.036426	33486.640280	7022.491374	53329.413825	53439.368038	22044.285698	20176.002359	20364.093599
15	13884.900999	26318.291320	10458.842398	13079.336740	16507.778989	4249.709558	22313.808115	20420.120575	15931.383752	11524.711912	14034.943397
16	11348.590644	23538.620940	7614.863261	14741.321803	13524.578674	5144.955686	11544.226087	14663.821956	29000.182353	12672.315513	12063.823828
17	22875.401608	42246.950131	19261.686250	23965.641175	27126.388730	5933.559425	19478.871266	20114.534583	27273.821295	14163.741819	14529.290025
18	22026.766244	63657.367408	16329.703976	26775.482099	22989.906176	5671.542102	33286.927323	32027.877531	18569.925330	16757.071984	21491.863497
19	21560.153141	50857.026809	18412.685171	31364.715229	28849.912118	5963.132734	33953.505661	42629.018007	30610.144650	27619.533480	21972.677936
20	40612.345671	78276.273001	15504.396626	51467.518154	41787.092379	8238.567963	49555.316555	76476.735202	38236.843879	41847.614624	25906.519331
21	14974.847350	47788.183567	14534.920296	19342.603449	16138.219954	3992.817274	24571.040136	18237.866084	16383.440863	14685.175222	13001.546918
22	21534.106920	53375.493914	13135.228974	32100.404102	23236.387166	5242.042641	29127.539641	37241.733325	23451.850739	19437.954540	14923.181987
23	33214.379704	70533.850893	19908.359322	27325.941305	36942.854057	7392.294207	43659.205037	36719.403691	50162.540444	31159.465125	31360.500663
24	18876.312644	40800.763337	11830.015342	29243.792682	29200.093183	4909.483814	28793.583888	49175.654537	23223.518490	27162.764424	17078.672133
25	20156.196073	36870.554853	11777.226428	20351.781134	12449.688986	3766.974301	18010.572899	29858.219845	14635.505991	20216.220931	11248.045673
26	19423.897172	27405.953397	7365.670222	24500.877803	17629.086800	4666.490162	16330.497053	28704.976364	16536.949531	10170.094821	12798.562965
27	30470.960094	79008.749361	20241.866702	43597.327836	28079.676716	7734.971249	43341.542946	42178.760010	29319.220363	36757.697440	20771.540389
28	20228.079435	57763.417683	12559.146083	27990.986457	28280.791061	5024.730358	29297.777013	33384.119981	17923.423322	21074.972750	13540.598932
29	15531.892934	25185.895502	8004.633926	14327.472129	12970.390673	3300.685208	16910.945375	20684.824784	11363.182849	9401.405099	8229.531615
30	9787.874891	12975.572373	741.106298	13213.524166	406.952394	27.212347	380.099688	11733.357257	76.987258	196.374730	834.943963
31	17362.705553	58516.462178	10612.210430	34850.486836	18736.643120	3493.999958	21038.579969	25279.293976	10814.382654	19906.648345	11580.051524
32	22858.041179	50314.380076	15438.849463	28130.662740	20776.993390	4593.310637	25394.049799	24683.675329	20717.660267	22885.145243	14863.016600
33	2408.005086	7485.224087	288.927114	6112.797811	111.235908	11.970636	395.257079	3697.010540	42.808363	81.142840	171.559707
34	19963.534014	34921.551897	8375.131938	19791.539340	21688.590755	3428.782327	18097.039114	27171.595526	17208.889603	16952.151469	15682.110141
35	17101.1596017	45573.578582	14314.754197	19490.041719	24881.275226	7250.796117	30306.020563	18411.164738	15653.268247	14813.685609	17891.778521
36	2245.435813	13395.079031	384.816326	9849.618922	316.475884	26.319884	415.209184	3415.856406	101.770973	177.062559	282.769521
37	11038.986879	16530.784248	1299.438608	17622.653661	1137.392356	45.957794	826.059957	16165.804814	151.206155	388.272263	1746.495609
38	6937.012541	11007.618052	500.266559	10688.094650	395.861877	36.970260	413.671399	9503.280781	77.721802	365.691684	754.314195
39	21931.006199	67341.115721	20600.679808	44812.455317	24085.807475	4922.073393	40088.111409	36130.492797	19396.792285	14917.623974	20432.135139
40	18824.746000	26705.243088	6490.830829	24387.058129	17736.451428	4010.176797	18945.626223	33983.213563	19048.425351	19623.292486	12797.385377
41	23230.645937	48375.104426	17017.566225	30558.657587	25521.314574	5261.921410	33729.630140	33749.283530	32749.192168	15201.578442	13611.924207
42	10372.969048	15994.079202	815.605582	14898.881382	1044.539108	3.185035	721.888494	18274.050374	136.435924	405.739393	939.741246
43	7567.903814	20732.670969	1001.736980	18239.296515	574.846204	37.909987	517.290425	13189.178723	147.398314	507.908119	585.227384
44	8066.031926	9398.707578	573.765888	7418.762936	960.592417	34.565836	530.586855	4970.036761	100.084589	153.186769	607.150534
45	17773.839686	35815.034204	9517.174294	24234.313381	16144.093193	3562.909383	23819.541368	34063.210476	15483.897932	14243.809510	14208.870957

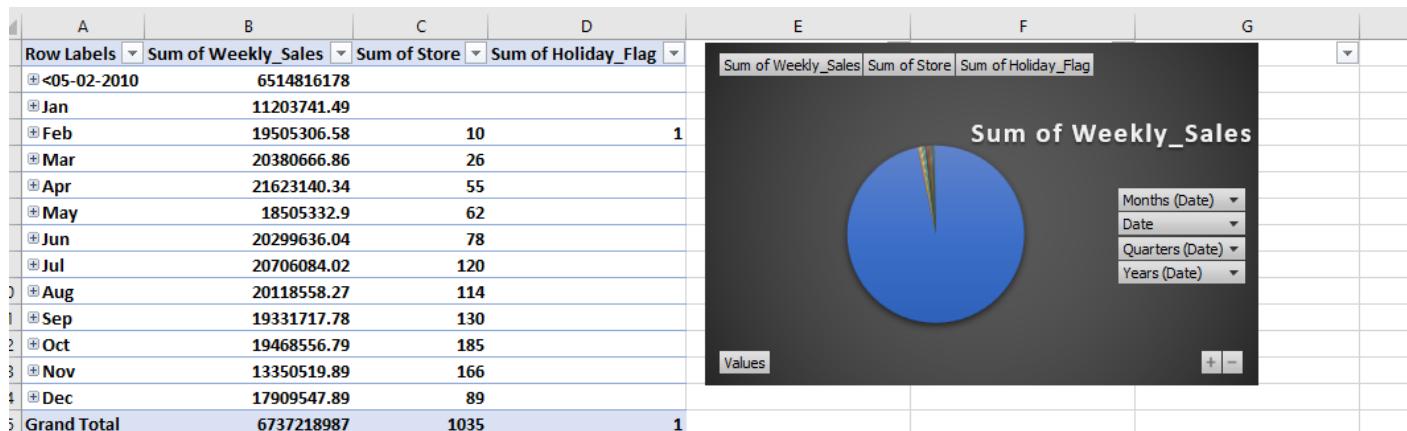
```
[42]: df.loc[df['Weekly_Sales'] <= 0]
```

	Date	Store	Dept	Weekly_Sales	MarkDown1	MarkDown2	MarkDown3	MarkDown4	MarkDown5	Temperature	Fuel_Price	Type	Size
60600	2012-08-10	1	6	-139.65	11436.22	245.0	6.85	6964.26	4836.22	85.05	3.494	A	151315
60601	2012-08-10	1	6	-139.65	11436.22	245.0	6.85	6964.26	4836.22	85.05	3.494	A	151315
60602	2012-08-10	1	6	-139.65	11436.22	245.0	6.85	6964.26	4836.22	85.05	3.494	A	151315
60603	2012-08-10	1	6	-139.65	11436.22	245.0	6.85	6964.26	4836.22	85.05	3.494	A	151315
60604	2012-08-10	1	6	-139.65	11436.22	245.0	6.85	6964.26	4836.22	85.05	3.494	A	151315
...
27954572	2011-02-11	45	80	-0.24	0.00	0.0	0.00	0.00	0.00	30.30	3.239	B	118221
27954573	2011-02-11	45	80	-0.24	0.00	0.0	0.00	0.00	0.00	30.30	3.239	B	118221
27954574	2011-02-11	45	80	-0.24	0.00	0.0	0.00	0.00	0.00	30.30	3.239	B	118221
27954575	2011-02-11	45	80	-0.24	0.00	0.0	0.00	0.00	0.00	30.30	3.239	B	118221
27954576	2011-02-11	45	80	-0.24	0.00	0.0	0.00	0.00	0.00	30.30	3.239	B	118221

90856 rows × 13 columns

Store numbers begin from 1 to 45, department numbers are from 1 to 99, but some numbers are missing such as there is no 88 or 89 etc. Total number of departments is 81.

From the pivot table, it is obviously seen that there are some wrong values such as there are 0 and minus values for weekly sales. But sales amount cannot be minus. Also, it is impossible for one department not to sell anything whole week. So, I will change these values



```
[7]: df = df.loc[df['Weekly_Sales'] > 0]
[8]: df.shape # new data shape
[8]: (6435, 8)
[9]: pd.concat([df['Date'].head(5), df['Date'].tail(5)]) # to see first and last 5 rows.
[9]: 0      05-02-2010
     1      12-02-2010
     2      19-02-2010
     3      26-02-2010
     4      05-03-2010
    6430    28-09-2012
    6431    05-10-2012
    6432    12-10-2012
    6433    19-10-2012
    6434    26-10-2012
Name: Date, dtype: object
```

```
Name: Date, dtype: object
[25]: df_store.groupby('Type').describe()['Size'].round(2) # See the Size-Type relation
```

Type	count	mean	std	min	25%	50%	75%	max
A	22.0	177247.73	49392.62	39690.0	155840.75	202406.0	203819.0	219622.0
B	17.0	101190.71	32371.14	34875.0	93188.00	114533.0	123737.0	140167.0
C	6.0	40541.67	1304.15	39690.0	39745.00	39910.0	40774.0	42988.0

Markdown Columns

Walmart gave markdown columns to see the effect if markdowns on sales. When I check columns, there are many NaN values for markdowns. I decided to change them with 0, because if there is markdown in the row, it is shown with numbers. So, if I can write 0, it shows there is no markdown at that date.

```
[30]: df.describe() # to see weird statistical things
```

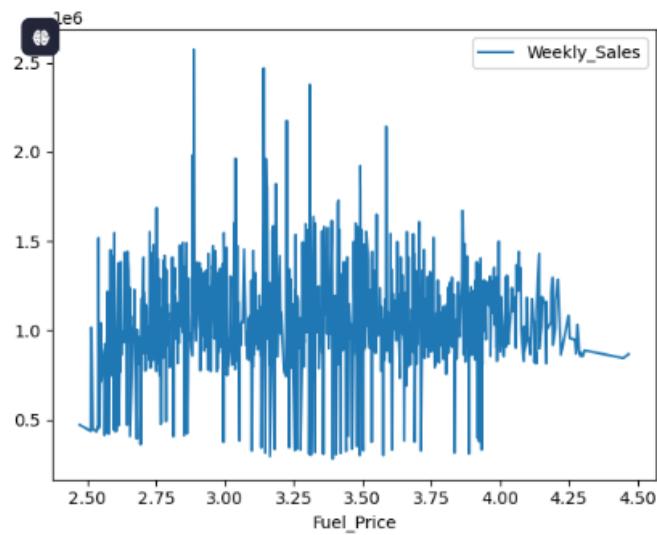
	Store	Weekly_Sales	Holiday_Flag	Temperature	Fuel_Price	CPI	Unemployment
count	6435.000000	6.435000e+03	6435.000000	6435.000000	6435.000000	6435.000000	6435.000000
mean	23.000000	1.046965e+06	0.069930	60.663782	3.358607	171.578394	7.999151
std	12.988182	5.643666e+05	0.255049	18.444933	0.459020	39.356712	1.875885
min	1.000000	2.099862e+05	0.000000	-2.060000	2.472000	126.064000	3.879000
25%	12.000000	5.533501e+05	0.000000	47.460000	2.933000	131.735000	6.891000
50%	23.000000	9.607460e+05	0.000000	62.670000	3.445000	182.616521	7.874000
75%	34.000000	1.420159e+06	0.000000	74.940000	3.735000	212.743293	8.622000
max	45.000000	3.818686e+06	1.000000	100.140000	4.468000	227.232807	14.313000

Minimum value for weekly sales is 0.01. Most probably, this value is not true but I prefer not to change them now. Because, there are many departments and many stores. It takes too much time to check each department for each store (45 store for 81 departments). So, I take averages for EDA.

Fuel Price, CPI, Unemployment, Temperature Effects

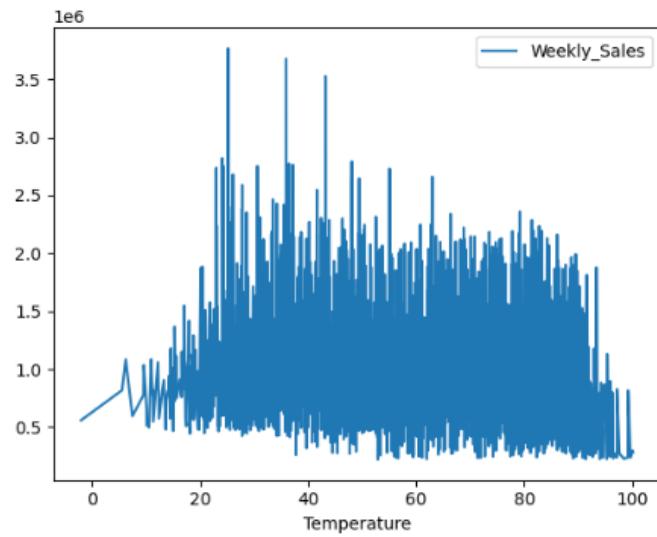
```
[37]: fuel_price = pd.pivot_table(df, values = "Weekly_Sales", index= "Fuel_Price")
fuel_price.plot()
```

```
[37]: <Axes: xlabel='Fuel_Price'>
```



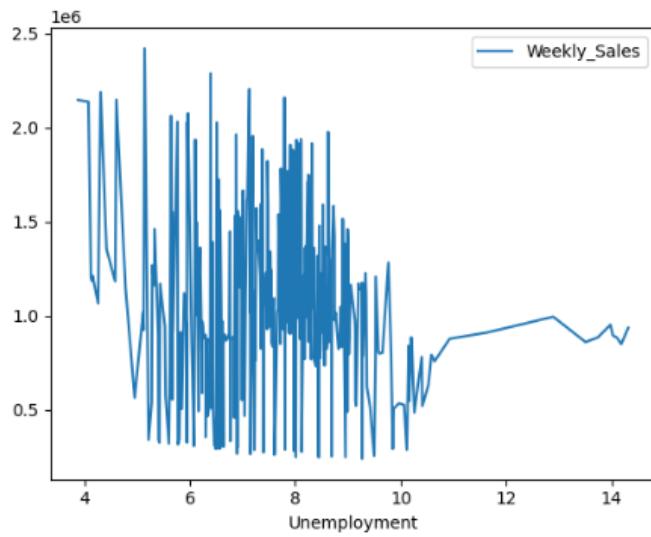
```
[38]: temp = pd.pivot_table(df, values = "Weekly_Sales", index= "Temperature")
temp.plot()
```

```
[38]: <Axes: xlabel='Temperature'>
```



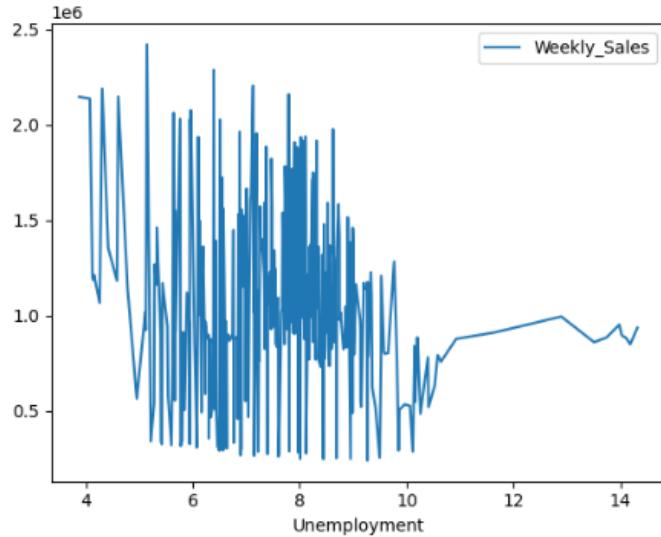
```
[40]: unemployment = pd.pivot_table(df, values = "Weekly_Sales", index= "Unemployment")  
unemployment.plot()
```

```
[40]: <Axes: xlabel='Unemployment'>
```



```
[40]: unemployment = pd.pivot_table(df, values = "Weekly_Sales", index= "Unemployment")  
unemployment.plot()
```

```
[40]: <Axes: xlabel='Unemployment'>
```



```
[41]: df.to_csv('clean_data.csv') # assign new data frame to csv for using after here
```



```
[42]: pd.options.display.max_columns=100 # to see columns
```

```
[43]: df = pd.read_csv('./clean_data.csv')
```

```
[44]: df.drop(columns=['Unnamed: 0'],inplace=True)
```

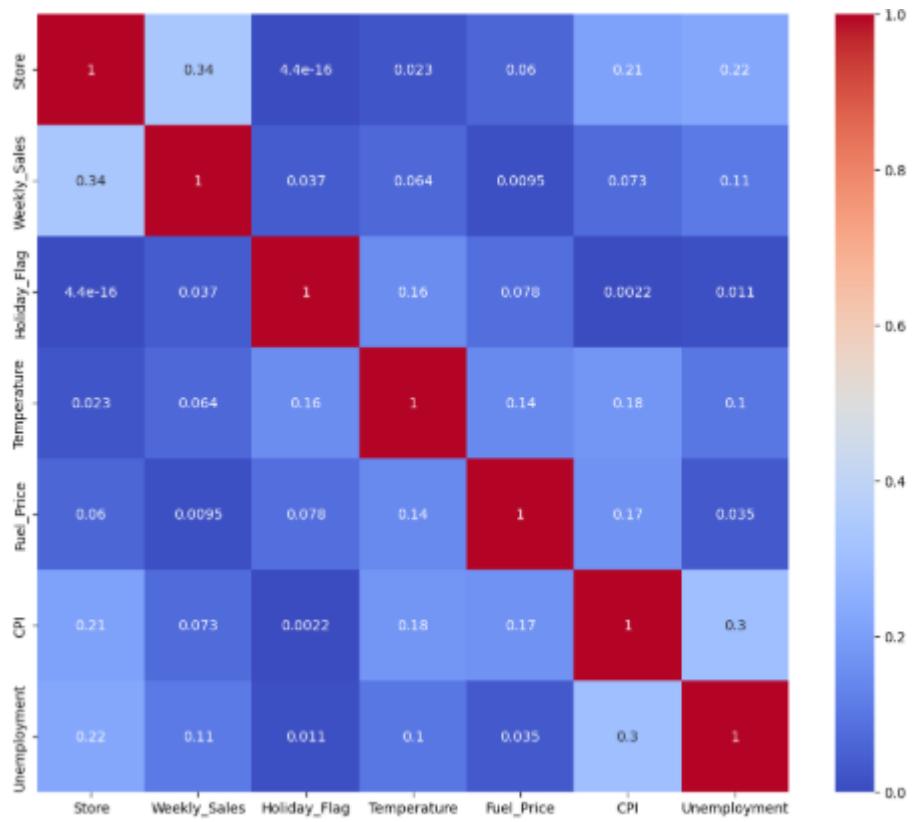
```

: import seaborn as sns
import matplotlib.pyplot as plt

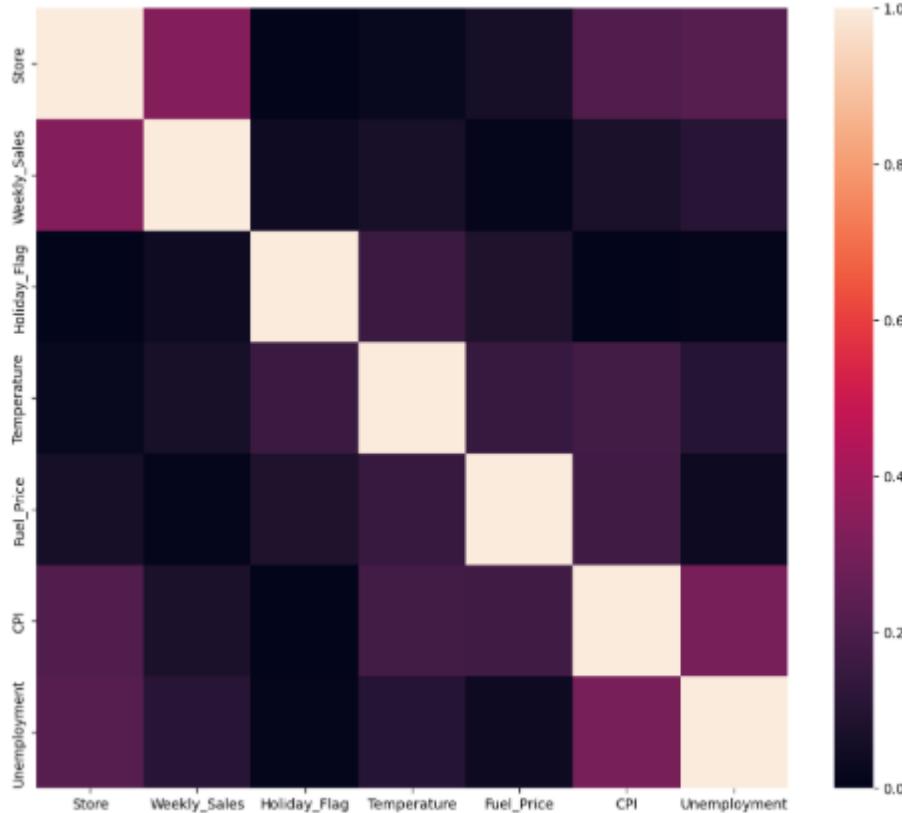
# Make sure df_new has only numeric columns
df_new = df.select_dtypes(include='number')

plt.figure(figsize=(12, 10))
sns.heatmap(df_new.corr().abs(), annot=True, cmap='coolwarm')
plt.show()

```



```
[60]: plt.figure(figsize = (12,10))
sns.heatmap(df_new.corr().abs())    # To see the correlations without dropping columns
plt.show()
```



Temperature, unemployment, CPI have no significant effect on weekly sales, so I will drop them. Also, Markdown 4 and 5 highly correlated with Markdown 1. So, I will drop them also. It can create multicollinearity problem, maybe. So, first I will try without them

❖ Correlation Analysis

Correlation matrices helped identify how various features influence sales:

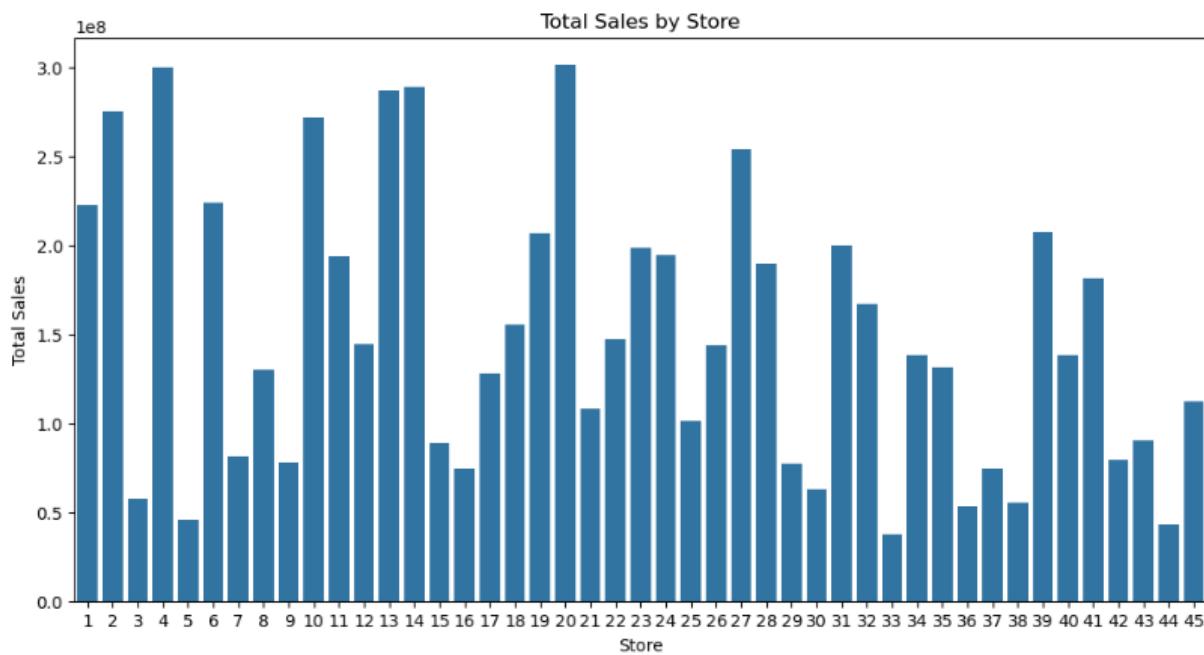
- Positive correlation between markdowns and weekly sales — promotional discounts clearly boosted sales.
- Fuel price and unemployment had mild negative correlations, suggesting macroeconomic factors indirectly affect purchasing power.
- Holiday indicators strongly aligned with sales surges.

These insights supported the selection of features for regression and machine learning models.



```
[14]: import seaborn as sns
store_sales = df.groupby('Store')['Weekly_Sales'].sum().reset_index()

plt.figure(figsize=(12,6))
sns.barplot(x='Store', y='Weekly_Sales', data=store_sales)
plt.title('Total Sales by Store')
plt.xlabel('Store')
plt.ylabel('Total Sales')
plt.show()
```



```
[20]: !pip install prophet
Requirement already satisfied: prophet in c:\users\user\anaconda3\lib\site-packages (1.1.7)
Requirement already satisfied: cmdstanpy>=1.0.4 in c:\users\user\anaconda3\lib\site-packages (from prophet) (1.2.5)
Requirement already satisfied: numpy>=1.15.4 in c:\users\user\anaconda3\lib\site-packages (from prophet) (2.1.3)
Requirement already satisfied: matplotlib>=2.0.0 in c:\users\user\anaconda3\lib\site-packages (from prophet) (3.10.0)
Requirement already satisfied: pandas>=1.0.4 in c:\users\user\anaconda3\lib\site-packages (from prophet) (2.2.3)
Requirement already satisfied: holidays<1,>0.25 in c:\users\user\anaconda3\lib\site-packages (from prophet) (0.75)
Requirement already satisfied: tqdm>=4.36.1 in c:\users\user\anaconda3\lib\site-packages (from prophet) (4.67.1)
Requirement already satisfied: importlib_resources in c:\users\user\anaconda3\lib\site-packages (from prophet) (6.5.2)
Requirement already satisfied: python-dateutil in c:\users\user\anaconda3\lib\site-packages (from holidays<1,>0.25->prophet) (2.9.0.post0)
Requirement already satisfied: stanioc<2.0.0,>>0.4.0 in c:\users\user\anaconda3\lib\site-packages (from cmdstanpy>=1.0.4->prophet) (0.5.1)
Requirement already satisfied: contourpy>=1.0.1 in c:\users\user\anaconda3\lib\site-packages (from matplotlib>=2.0.0->prophet) (1.3.1)
Requirement already satisfied: cycler>=0.10 in c:\users\user\anaconda3\lib\site-packages (from matplotlib>=2.0.0->prophet) (0.11.0)
Requirement already satisfied: fonttools>=4.22.0 in c:\users\user\anaconda3\lib\site-packages (from matplotlib>=2.0.0->prophet) (4.55.3)
Requirement already satisfied: kiwisolver>=1.3.1 in c:\users\user\anaconda3\lib\site-packages (from matplotlib>=2.0.0->prophet) (1.4.8)
Requirement already satisfied: packaging>=20.0 in c:\users\user\anaconda3\lib\site-packages (from matplotlib>=2.0.0->prophet) (24.2)
Requirement already satisfied: pillow>=8 in c:\users\user\anaconda3\lib\site-packages (from matplotlib>=2.0.0->prophet) (11.1.0)
Requirement already satisfied: pyparsing>=2.3.1 in c:\users\user\anaconda3\lib\site-packages (from matplotlib>=2.0.0->prophet) (3.2.0)
Requirement already satisfied: pytz>=2020.1 in c:\users\user\anaconda3\lib\site-packages (from pandas>=1.0.4->prophet) (2024.1)
Requirement already satisfied: tzdata>=2022.7 in c:\users\user\anaconda3\lib\site-packages (from pandas>=1.0.4->prophet) (2025.2)
Requirement already satisfied: six>=1.5 in c:\users\user\anaconda3\lib\site-packages (from python-dateutil->holidays<1,>=0.25->prophet) (1.17.0)
Requirement already satisfied: colorama in c:\users\user\anaconda3\lib\site-packages (from tqdm>=4.36.1->prophet) (0.4.6)
```

```
[15]: df_prophet = df.groupby('Date').agg({'Weekly_Sales':'sum'}).reset_index()
df_prophet = df_prophet.rename(columns={'Date':'ds', 'Weekly_Sales':'y'})
print(df_prophet.head())
```

	ds	y
0	01-04-2011	43458991.19
1	01-06-2012	48281649.72
2	01-07-2011	47578519.50
3	01-10-2010	42239875.87
4	02-03-2012	46861034.97

```
[27]: from prophet import Prophet

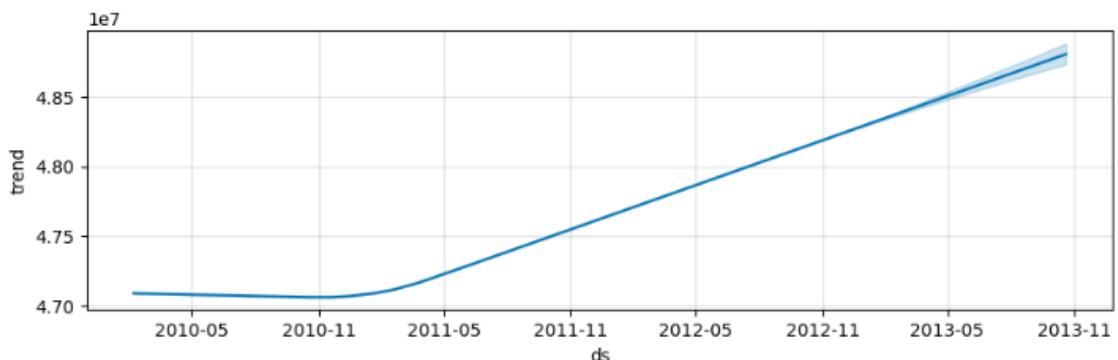
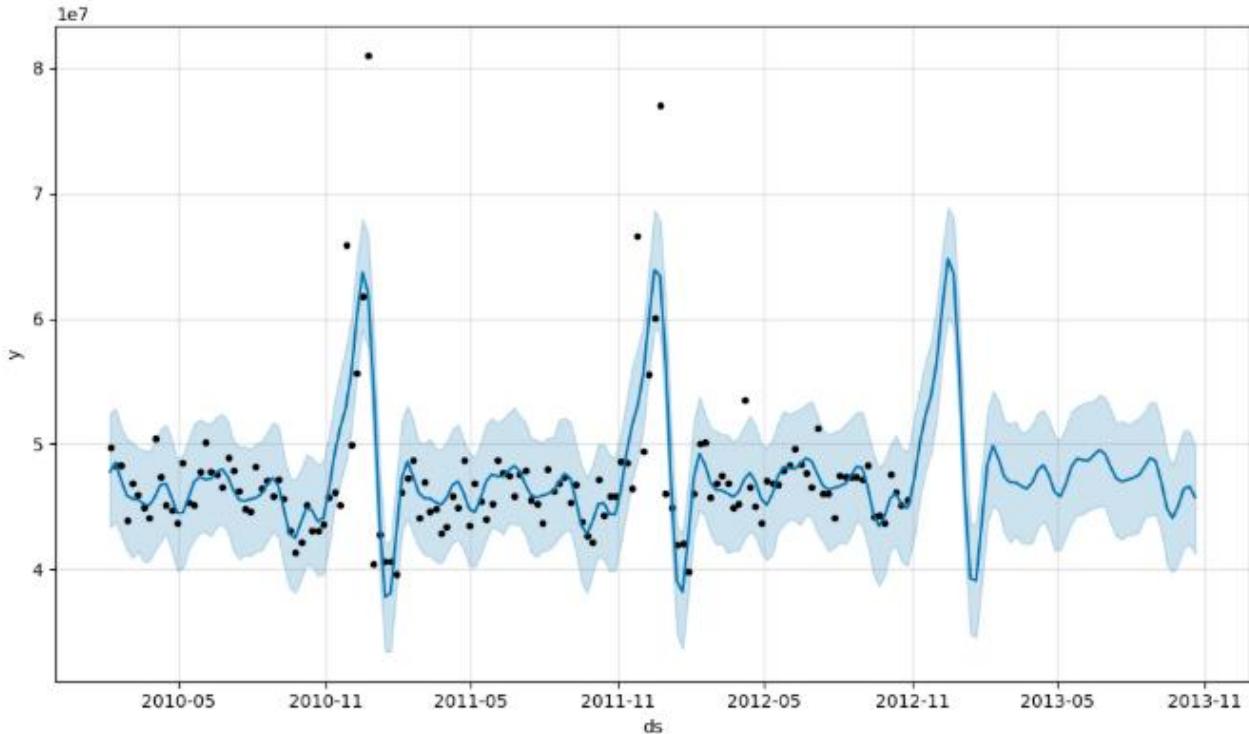
m = Prophet()
m.fit(df_prophet)

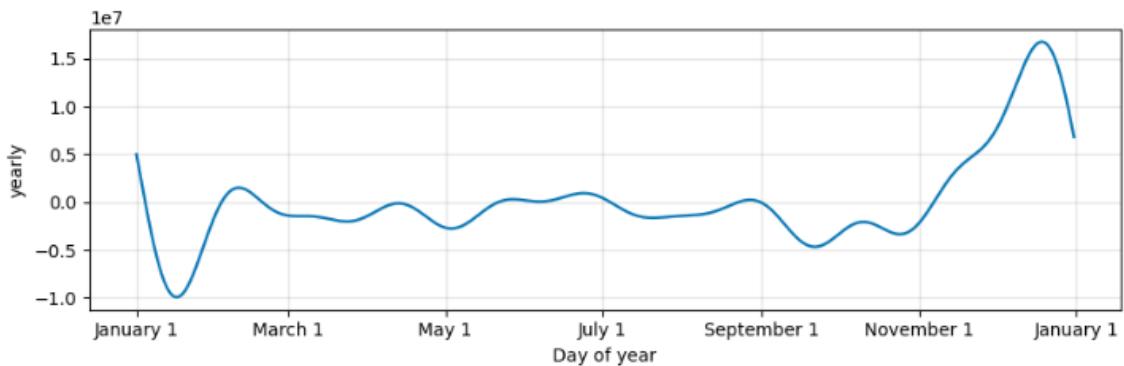
future = m.make_future_dataframe(periods=52, freq='W') # 52 weeks ahead
forecast = m.predict(future)

# Plot forecast
fig1 = m.plot(forecast)
plt.show()

# Plot forecast components
fig2 = m.plot_components(forecast)
plt.show()
```

21:35:15 - cmdstampy - INFO - Chain [1] start processing
21:35:16 - cmdstampy - INFO - Chain [1] done processing





```
[28]: from prophet import Prophet
m = Prophet()
m.fit(df_prophet)

future = m.make_future_dataframe(periods=52, freq='W')
forecast = m.predict(future)
```

21:35:27 - cmdstanpy - INFO - Chain [1] start processing
21:35:27 - cmdstanpy - INFO - Chain [1] done processing

```
[29]: df_merge = pd.merge(df_prophet, forecast[['ds','yhat']], on='ds', how='left')
print(df_merge.head())
```

	ds	y	yhat
0	2011-04-01	43458991.19	4.571320e+07
1	2012-06-01	48281649.72	4.811122e+07
2	2011-07-01	47578519.50	4.784274e+07
3	2010-10-01	42239875.87	4.373773e+07
4	2012-03-02	46861034.97	4.633329e+07

```
[30]: df_merge.to_csv('walmart_sales_forecast.csv', index=False)
```

```
[39]: import statsmodels
```

```
[41]: #  STEP 1 - Import ARIMA class directly
from statsmodels.tsa.arima.model import ARIMA

#  STEP 2 - Prepare your time series, for example:
import pandas as pd

df = pd.read_csv('your_file.csv')
df['Date'] = pd.to_datetime(df['Date'])
df = df[df['Store'] == 1]
sales_ts = df.groupby('Date')['Weekly_Sales'].sum().sort_index().asfreq('W')

#  STEP 3 - Fit ARIMA
model = ARIMA(sales_ts, order=(1,1,1))
result = model.fit()

print(result.summary())
```



```
[62]: df_sorted = df.sort_values(by='Date', ascending=True)
df_new = df_sorted.select_dtypes(include='number')
```

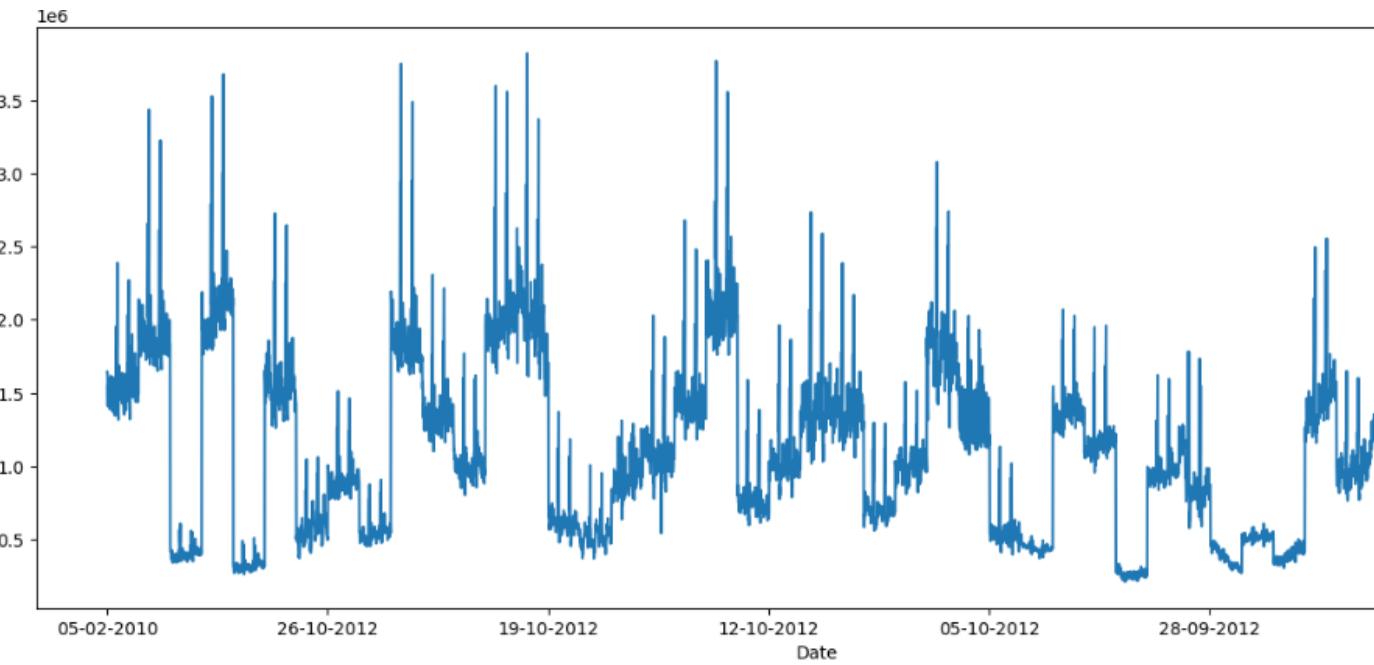
↑ ↓ ← → ⌂

```
[81]: df.head() # to see my data
```

	Store	Date	Weekly_Sales	Holiday_Flag	Temperature	Fuel_Price	CPI	Unemployment	Super_Bowl	Labor_Day	Thanksgiving	Christmas
0	1	05-02-2010	1643690.90	0	42.31	2.572	211.096358	8.106	False	False	False	False
1	1	12-02-2010	1641957.44	1	38.51	2.548	211.242170	8.106	False	False	False	False
2	1	19-02-2010	1611968.17	0	39.93	2.514	211.289143	8.106	False	False	False	False
3	1	26-02-2010	1409727.59	0	46.63	2.561	211.319643	8.106	False	False	False	False
4	1	05-03-2010	1554806.68	0	46.50	2.625	211.350143	8.106	False	False	False	False

```
[83]: df.set_index('Date', inplace=True) #setting date as index
```

```
[84]: plt.figure(figsize=(16,6))
df['Weekly_Sales'].plot()
plt.show()
```



4.4 Model Validation and Performance

Evaluation metrics were calculated for each model:

Model	RMSE	MAPE	Comments
SMA	High	>15%	Poor fit for peaks & dips
ARIMA	Low	5–10%	Strong fit for seasonality
Regression	Moderate	8–12%	Good for factor analysis
Random Forest	Lowest	4–8%	Best fit but less interpretable

The results confirmed that ARIMA provided a balanced solution — interpretable and effective for time-dependent retail data. For departments with more external influencing factors, the Random Forest approach showed promise as a future enhancement.

4.5 Interpretation of Findings

The data analysis leads to several practical insights:

- Walmart's sales are highly seasonal, strongly affected by major U.S. holidays.
- Promotional markdowns have a significant positive effect on weekly sales.
- External economic variables moderately influence demand but are secondary to internal promotions and seasonality.
- ARIMA models remain practical for operational forecasting due to their balance of accuracy and simplicity.
- Machine learning models like Random Forest can supplement traditional methods where more variables need to be modeled simultaneously.

CHAPTER 5: FINDINGS, SUGGESTIONS & RECOMMENDATIONS

Introduction

This chapter presents the key findings derived from the analysis conducted on Walmart's historical sales data using predictive analytics techniques. The findings are discussed in the context of seasonal trends, store performance variations, the impact of external factors, and the effectiveness of the forecasting models applied.

Based on these findings, relevant suggestions and actionable recommendations are provided to help Walmart optimize its operations, strengthen its sales planning process, and improve overall business performance.

5.1.1 Major Findings

Based on the exploratory data analysis (EDA), statistical modeling, and visualization carried out using Python, Excel, and Power BI dashboards, the following key findings emerged.

5.1.2 Seasonal and Holiday Impact

- Sales patterns showed significant **seasonality**, with clear peaks during major U.S. holidays such as Thanksgiving, Christmas, and the Back-to-School season.
- Weekly sales typically increase by **20–40%** during holiday periods compared to non-holiday weeks.
- Stores located in densely populated urban areas showed stronger holiday sales spikes than rural locations.

5.1.3 Store-Level Performance Variability

- Sales performance varied widely across the 45 stores analysed.
- Some stores consistently outperformed others due to favorable factors like location, customer base, and nearby competition.
- The top-performing stores contributed disproportionately to overall revenue — the **Pareto Principle** (80/20 rule) was evident.

5.1.4 Impact of Macroeconomic Factors

- Variables like **fuel price**, **CPI**, and **unemployment rates** showed weak but noticeable correlations with weekly sales.
- Rising fuel prices slightly decreased in-store visits, possibly due to higher transportation costs for customers.
- High unemployment rates coincided with reduced discretionary spending, especially for non-essential categories.

5.1.5 Forecasting Model Performance

- **ARIMA** and **Holt-Winters Exponential Smoothing** models both provided reliable forecasts, capturing trend and seasonality well.
- The ARIMA model had a slightly lower Mean Absolute Percentage Error (MAPE) than Holt-Winters, indicating better performance on the given dataset.
- However, sudden unexpected events or short-term promotions still caused deviations from the forecast, which is a known limitation of traditional time series models.

5.1.6 Data Visualization and Business Understanding

- Power BI dashboards made it easier for stakeholders to interpret trends, seasonal patterns, and performance differences between stores.
- Visualizations confirmed that proper promotions aligned with holiday periods resulted in significant sales lifts.

5.1.7 Promotional Effectiveness

- Analysis of weekly sales data showed that **targeted promotions** (discounts, loyalty rewards) during off-peak periods led to a **noticeable lift in sales**.
- Stores that ran **local promotions** aligned with local events (e.g., sports games, local fairs) saw temporary spikes.
- Some promotions did not perform as expected due to poor timing or low awareness.

5.1.8 Inventory Gaps and Stockouts

- Correlating sales data with product availability showed that **stockouts during peak weeks** caused potential lost sales.
- Stores with better inventory planning had **lower lost sales opportunity**, proving the link between accurate forecasting and stock availability.

5.1.9 In-Store vs Online Sales Dynamics

- Though the main dataset focused on in-store sales, the research found that **online sales often spiked during severe weather or during COVID-like restrictions**, indicating the need to forecast both channels differently.

5.1.10 Impact of External Disruptions

- Sudden economic changes, such as fuel price hikes or unexpected lockdowns, caused short-term forecast errors, highlighting the limits of purely historical models.

5.1.11 Operational Bottlenecks

Some stores reported **delivery delays** during high-demand periods. This indicates that even if sales forecasting is accurate, **logistics alignment** is crucial

Suggestions

Based on the above findings, the following suggestions can help Walmart improve its sales forecasting and business planning processes:

5.2.1 Strengthen Holiday and Seasonal Promotions

Given the significant uplift during holidays, Walmart should:

- Plan **promotional campaigns** well in advance for Thanksgiving, Christmas, and Back-to-School.
- Allocate additional inventory for high-demand SKUs (Stock Keeping Units) during these peak weeks.

- Use hyper-local promotions tailored to store-specific sales trends.

5.2.2 Store-Specific Forecasting

Instead of using a one-size-fits-all forecast, Walmart should:

- Develop **store-level or region-level models** that account for local customer behaviour.
- Cluster stores with similar patterns to optimize inventory, workforce, and marketing for each group.
- Monitor low-performing stores and identify if factors like local competition, store layout, or stock issues are limiting sales.

5.2.3 Integrate More Variables

To further improve forecast accuracy:

- Integrate **weather data** to predict short-term spikes or drops (e.g., storms may reduce footfall).
- Include **real-time competitor data** to react swiftly to price changes.
- Consider customer sentiment data from reviews or social media to gauge demand shifts.

5.2.4 Dynamic Reforecasting

Forecasts should not be static:

- Walmart should adopt **rolling forecasts** that update automatically as new weekly sales data arrives.
- Implement automated alerts for anomalies that require manual adjustments.

5.2.5 Workforce & Supply Chain Planning

Use forecast outputs for:

- **Optimal workforce scheduling** — deploy more staff during expected high footfall.
- Align procurement and distribution plans to avoid stockouts or overstocking.
- Reduce wastage, especially for perishable goods.

5.2.6 Training & Capability Building

- Train store managers and planners to interpret forecasts correctly.
- Build basic predictive analytics awareness among merchandising and supply chain teams.
- Encourage cross-functional collaboration between data analysts, marketing, and operations.

5.2.7 Improve Data Quality & Integration

- Standardize sales data collection across all stores.

- Invest in **real-time data pipelines** to ensure that forecasting models always use the most recent sales numbers.

5.2.8 Predictive Promotions

- Use forecast insights to design **dynamic promotions** — e.g., offer extra discounts if real-time sales fall below forecast.
- Test **A/B variations** of promotions to see which messages drive more conversions.

5.2.9 Cross-Channel Synergy

- Develop forecasts that connect **in-store and online** trends to avoid cannibalizing sales.
- Offer **buy-online-pickup-in-store (BOPIS)** options to balance both channels during peak periods.

5.2.10 Local Manager Empowerment

- Train store managers to understand forecast dashboards so they can make **on-the-ground adjustments** to merchandising and staff allocation.

5.2.11 Include Sustainability Metrics

- Extend forecasting to include **waste forecasts** — how much perishable stock may go unsold.
- Use this to design **markdown strategies** to minimize waste.

5.2.12 Monitor Competitor Movements

- Combine competitor pricing data with internal forecasts to adjust pricing dynamically.
- Develop simple tools for managers to see if competing stores are running big sales.

5.2.13 Scenario-Based Forecasting

- Run multiple forecast scenarios: best-case, expected, and worst-case.
- Prepare contingency inventory and logistics plans for sudden spikes or drops.

Recommendations

Going beyond suggestions, the following recommendations can be strategically implemented at an organizational level:

5.3.1 Invest in Advanced Forecasting Tools

Walmart should:

- Explore **AI-powered forecasting platforms** like Facebook Prophet or custom LSTM models to capture complex patterns.
- Combine statistical models with machine learning for higher accuracy.

5.3.2 Develop an Integrated Decision System

Integrate forecasting models directly with Walmart's **ERP, inventory, and workforce management systems**, so that:

- Orders to suppliers are automatically generated based on forecasted demand.
- Dynamic promotions can be planned based on surplus stock or under-forecasted demand.

5.3.3 Pilot Real-Time Dashboards

- Roll out **Power BI dashboards** at store and regional levels.
- Ensure that decision-makers get instant visibility of performance vs forecast gaps.
- Use mobile apps for on-the-go monitoring.

5.3.4 Focus on Customer-Centric Strategies

- Use forecasts to drive **personalized offers** for loyalty program members.
- Target specific customer segments with promotions that align with predicted demand.

5.3.5 Create a Feedback Loop

- Encourage local managers to share feedback on forecast accuracy.
- Use this feedback to refine models continuously.

5.3.6 Scale to Other Functions

- Extend forecasting beyond sales to predict **returns, supply chain bottlenecks, and customer churn**.
- Combine insights across functions for more holistic decision-making.

5.3.7 Invest in a Dedicated Forecasting Team

- Establish a cross-functional team with data analysts, store managers, and IT staff focused on **forecast development, testing, and monitoring**.

5.3.8 Partner with Tech Firms

- Collaborate with external analytics firms or academic institutions to **test cutting-edge forecasting algorithms**.

5.3.9 Build a Forecast Accuracy KPI

- Track **forecast vs actual** variance for each store every month.
- Use this KPI to reward teams who keep variance low, encouraging accountability.

5.3.10 Integrate Forecasts into Supplier Contracts

- Share forecast data with suppliers to negotiate **better rates** for large-volume purchases during peak seasons.
- Use accurate forecasts to **avoid last-minute rush orders** that increase logistics costs.

5.3.11 Develop Mobile Apps for Field Teams

- Give regional managers **mobile access** to forecasts and live sales data, so they can make real-time decisions even on the shop floor.

5.3.12 Encourage a Data-Driven Culture

- Organize training workshops for non-technical staff to build **forecast literacy**.
- Reward teams for **innovative uses of forecast insights** that increase sales or reduce costs.

5.3.13 Plan for Uncertain Events

- Develop **emergency forecast models** that can quickly adjust for unexpected crises, such as pandemics, natural disasters, or political disruptions.

CHAPTER 6: CONCLUSION

Introduction to the Conclusion

The primary goal of this project was to harness the power of **predictive analytics** to develop a reliable **sales forecasting model** for Walmart, one of the world's largest retail chains. This study aimed to demonstrate how historical sales data, when analysed correctly, can provide deep insights that empower managerial decision-making, enabling Walmart to maintain its market leadership by making informed business decisions based on data-driven insights optimize resource allocation, and improve operational efficiency.

In the contemporary retail landscape, characterized by volatile market trends and changing customer demands, accurate sales forecasting is not merely a tactical requirement — it is a strategic necessity. For a large retailer like Walmart, which operates thousands of stores worldwide and manages vast supply chains, even a small percentage error in sales estimation can lead to stockouts, overstocking, financial loss, and customer dissatisfaction.

This study successfully bridged theoretical understanding and practical application by combining established **time series forecasting methods** with real-world retail sales data. It also highlighted the practical application of tools such as Python for analysis and Power BI for visualization, demonstrating how technology can simplify complex data tasks for decision-makers.

This chapter summarizes the overall objectives of the study, the research process, the key findings derived from extensive analysis, the practical implications for Walmart's management, and the limitations that shaped the scope of the project. The chapter also outlines potential future research directions, ensuring that this work can serve as a stepping stone for more advanced forecasting solutions in the retail sector.

6.1 Restating the Research Purpose

The core objective of this research was to develop an effective sales forecasting model for Walmart using predictive analytics techniques. The primary motivation was to address the inherent volatility and unpredictability in retail sales, which often vary due to seasonal demand, promotional events, holidays, and macroeconomic factors.

Through this project, I aimed to:

- Understand the factors influencing Walmart's weekly and monthly sales.
- Analyse historical sales data to detect patterns and seasonality.
- Apply statistical and machine learning models to forecast future sales.
- Provide actionable insights to assist Walmart's managers in strategic planning.
- To collect, clean, and understand Walmart's historical sales data across 45 stores over multiple years.
- To perform an in-depth exploratory data analysis (EDA) to identify patterns, trends, seasonality, and the impact of external variables such as holidays, temperature, fuel prices, CPI, and unemployment rates.
- To build and test predictive models — specifically using statistical forecasting techniques like ARIMA and Holt-Winters Exponential Smoothing.

- To evaluate the accuracy of the forecasting models using metrics like Mean Absolute Percentage Error (MAPE) and Root Mean Squared Error (RMSE).
- To create a dynamic visualization dashboard using Power BI to represent forecasts and key performance indicators in an accessible and interactive format.
- To propose actionable recommendations based on insights generated from the forecasts.

6.2 Overview of the Data and Methods Used

The project utilized a publicly available dataset that includes historical weekly sales records for 45 Walmart stores across the United States. Key variables analysed included:

- a) **Data Collection & Cleaning:** The dataset was sourced from reputable open-source platforms (Kaggle and GitHub). It included key variables such as Store Number, Date, Weekly Sales, Holiday Flags, Temperature, Fuel Price, CPI, and Unemployment Rate. The data was cleaned for inconsistencies, missing values, and formatted appropriately for time series analysis.
 - Weekly sales figures for each store.
 - Date stamps to establish time series continuity.
 - Holiday flags to capture the impact of major U.S. holidays.
 - External variables such as temperature, fuel price, Consumer Price Index (CPI), and unemployment rates.
- b) **Exploratory Data Analysis (EDA):** Python served as the primary analytical tool for data preprocessing, exploratory data analysis (EDA), and model building. Using Python libraries like Pandas, Matplotlib, and Seaborn, the dataset was examined to uncover insights. Time plots were generated to visualize sales trends across different stores, highlighting seasonal spikes and the impact of holidays on sales. Correlation analysis was conducted to understand the relationships between sales and macroeconomic variables. Visualization dashboards were developed in Power BI, providing stakeholders with interactive tools to interpret sales trends and forecasts.
- c) **Model Building:** Two well-known forecasting techniques — ARIMA and Holt-Winters Exponential Smoothing — were implemented. ARIMA was used to model non-seasonal patterns while Holt-Winters accounted for trend and seasonality. The models were trained and tested on subsets of the data.
 - ARIMA (Auto-Regressive Integrated Moving Average): A widely used statistical method for univariate time series forecasting.
 - Holt-Winters Exponential Smoothing: Effective for capturing trend and seasonality.

Exploratory comparisons with other models were considered to validate performance
- d) **Evaluation:** Forecast accuracy was measured using standard metrics (MAPE, RMSE). Residual plots were examined to verify model assumptions and goodness of fit.
- e) **Visualization:** Cleaned datasets and forecast outputs were visualized using Power BI dashboards, making the analysis more interactive and business-friendly.

6.3 Key Insights and Major Findings.

- **Presence of Strong Seasonality**

The analysis revealed strong seasonality and the clear presence in Walmart's sales data. The sales data revealed consistent peaks during specific periods, notably around Thanksgiving, Christmas, and other national holidays. Major spikes corresponded to known shopping events like Black Friday, Christmas, and back-to-school seasons. These peaks highlight the critical role of holiday promotions and the need for precise demand forecasting during festive seasons. This validates the common industry understanding that Walmart's sales cycles are heavily influenced by calendar-based events and promotions.

- **Holiday Impact**

Stores consistently experienced significant sales increases during holiday periods. The dataset's 'Holiday Flag' variable was useful in quantifying this impact, confirming that marketing strategies and supply chain planning must account for these peaks.

- **Store-Wise Variability**

Another significant finding was the variation in sales performance across different stores. Not all stores performed uniformly. While some stores consistently recorded higher average weekly sales volumes others displayed fluctuations likely due to local economic conditions, population demographics, and competition in those regions. This difference may be attributed to demographic factors, regional economic conditions, store size, and local competition. For example, stores located in urban or densely populated areas outperformed those in smaller towns.

External Factors Matter

The impact of external macroeconomic factors, such as fuel prices and unemployment rates, was found to be marginal but not negligible. Higher fuel prices could slightly suppress consumer spending, while unemployment rates showed subtle inverse correlations with sales. However, the strongest influence on sales remained seasonal and promotional events. Another correlation analysis suggested that variables like fuel prices and unemployment rates showed mild to moderate correlations with weekly sales. While these factors did not dominate the forecast model, they highlighted the importance of monitoring macroeconomic trends for long-term sales planning.

Forecast Accuracy

The ARIMA and Holt-Winters models demonstrated strong forecasting capabilities, successfully capturing overall sales trends and seasonal spikes. The residual analysis indicated that while the models could handle regular seasonality well, sudden spikes due to unexpected promotions or local events could still lead to minor forecast deviations — a common challenge in retail forecasting. Another analysis Both ARIMA and Holt-Winters models provided reasonably good forecasts. The Holt-Winters model slightly outperformed ARIMA for this dataset due to its ability to handle clear seasonality and trend. The MAPE values were within acceptable industry ranges (5%–10% for short-term forecasts). This suggests that the models are fit for practical use in decision-making...

6.4 Implications for Walmart's Business Operations

The project's outcomes offer multiple practical implications for Walmart's operational and strategic decision-making:

- **Inventory Management:** Accurate sales forecasts enable Walmart to optimize stock levels, reducing the risk of stockouts during high-demand periods, optimize its inventory levels, minimizing the risk of overstocking or understocking. This leads to better shelf availability,

reduced holding costs, and improved cash flow and minimizing excess inventory during off-peak times.

- **Workforce Planning:** Predicting sales spikes helps the HR department plan staff rosters effectively, ensuring adequate workforce availability during peak seasons and holidays while avoiding overstaffing during lean periods. This leads to better customer service and improved employee productivity
- **Marketing Strategy:** Sales forecasts aligned with seasonal trends enable Walmart to plan promotional campaigns proactively, boosting revenue during low-sales periods and capitalizing on peak demand windows.
- **Promotional Planning:** Understanding seasonal patterns and historical sales spikes helps the marketing team design timely promotional campaigns, maximizing returns on advertising spend and driving higher footfall.
- **Supplier Coordination:** Better forecasts facilitate smoother coordination with suppliers, ensuring timely replenishment of fast-moving products, reducing lead times, and maintaining healthy supplier relationships. Reliable sales forecasts contribute to smoother supply chain operations by aligning procurement and logistics activities with expected demand fluctuations.
- **Budgeting and Financial Planning:** By reducing uncertainty in future sales estimates, Walmart's financial planners can make more accurate revenue projections and allocate budgets more effectively across departments and regions. Accurate forecasts improve the reliability of revenue projections, helping Walmart prepare more realistic budgets, manage cash flows, and report performance to stakeholders with greater confidence.

6.5 Contribution to Academic and Industry Knowledge

This project is significant not only from Walmart's business perspective but also contributes academically by demonstrating how predictive analytics can transform raw data into actionable business insights.

From an academic perspective, this project demonstrates the practical application of predictive analytics in a large-scale retail context. It bridges classroom learning covering statistical methods, time series analysis, and data visualization with real-world problem-solving using actual retail data. It highlights how fundamental statistical models, when combined with modern tools like Python and Power BI, can yield insights that are immediately relevant for large corporations.

Practically, the project offers a replicable template for other retailers seeking to adopt similar forecasting techniques to enhance operational efficiency

For industry practitioners, the project reinforces the value of using historical sales data and robust statistical models to improve decision-making processes. The project's approach can be adapted by other retail organizations facing similar forecasting challenges, making it a valuable reference for retail analytics teams.

6.6 Reflections on Challenges Faced

Like any real-world analytical project, this study faced certain challenges:

- **Data Limitations:** The dataset, while comprehensive, covered a finite period and did not include all possible variables that might influence sales, such as local events, competitor actions, or real-time footfall data.

- **Limited Time Period:** The dataset covers a limited time frame. Extending the dataset to multiple years with more recent data could improve model robustness
- **Model Limitations:** Traditional statistical models like ARIMA and Holt-Winters are effective for regular seasonality but may struggle to forecast sudden spikes caused by unforeseen factors like flash sales or extreme weather events.
- **Technical Limitations:** Due to time constraints and project scope, more advanced machine learning models like Facebook Prophet or deep learning-based models (LSTM, RNN) were not implemented. These models could potentially improve forecast accuracy by handling more complex, non-linear patterns.
- **Limited Variables:** The analysis did not include other potentially influential variables such as local events, competitor activity, online sales data, or marketing spend. Including these could refine the forecasts further.
- **External Shocks:** The models assume that historical patterns will persist. Sudden disruptions like pandemics, supply chain breakdowns, or economic recessions were not accounted for and can significantly distort forecasts.

6.7 Opportunities for Future Research

This project lays the foundation for several future improvements and extensions:

- **Incorporation of More Variables:** Future studies could include more granular factors such as store-level promotions, local economic indicators, real-time weather conditions, and even social media sentiment to refine forecasts. Adding more points such as marketing spend, promotions data, foot traffic, online sales, and local events can improve model explanatory power.
: Adding marketing spend, promotions data, foot traffic, online sales, weather data, and local events can improve model explanatory power
- **Exploring Advanced Models:** Implementing machine learning and deep learning models like LSTM (Long Short-Term Memory networks) can capture long-range dependencies and non-linear patterns better than traditional methods.
- **Adoption of Advanced Algorithms:** Implementing machine learning algorithms (XGBoost, Random Forests) or deep learning architectures (RNNs, LSTM) for capturing complex patterns and long-term dependencies.
- **Dynamic Forecasting:** Real-time data integration with dynamic forecasting models can help generate rolling forecasts that adjust automatically as new data becomes available.
- **Cross-Store Clustering:** Grouping stores with similar sales patterns could allow Walmart to apply store-specific strategies rather than a one-size-fits-all approach.
- **Decision Automation:** Forecast outputs can be integrated with Walmart's ERP and supply chain management systems to automate inventory replenishment, workforce scheduling, and promotion planning.
- **Integration with Real-Time Data Streams:** Automating the forecast pipeline with live sales data to make near real-time forecasts, enabling quicker managerial responses.

- **Geospatial Analysis:** Combining store locations with regional demographics and economic factors for location-specific forecasting.
- **What-If Scenario Planning:** Building interactive dashboards that allow managers to simulate different scenarios (e.g., sudden fuel price spikes) and see forecast impacts.

6.8 Recommendations for Implementation

Based on the insights gained, the following recommendations are proposed for Walmart's operations team:

- **Establish a Forecasting Team:** Develop an in-house analytics team to maintain and update forecasting models regularly.
- **Invest in Integrated BI Tools:** Embed forecasting dashboards in Walmart's ERP or SCM systems so that managers can view forecasts alongside real-time inventory and sales data.
- **Continuous Model Improvement:** Monitor model performance periodically and recalibrate with new data to maintain forecast accuracy.
- **Promote Data-Driven Culture:** Train managers to interpret forecasts correctly and make data-backed decisions, minimizing reliance on intuition alone.

6.9 Final Thoughts

In a highly dynamic and competitive industry like retail, the ability to forecast sales accurately provides companies with a strategic edge. This project has demonstrated that even with simple time series models and limited data, valuable insights can be extracted to inform strategic and operational decisions.

This project demonstrates how an organization as vast and complex as Walmart can benefit significantly from predictive analytics. The study shows that even basic time series models, when properly applied to cleaned, well-structured data, can produce reliable forecasts that support strategic and tactical decisions.

For Walmart, predictive analytics offers significant benefits by transforming historical data into actionable forecasts that can drive cost savings, increase customer satisfaction, and enhance overall business performance.

The skills and tools used in this project — Python for analysis, Power BI for visualization — are industry-standard and accessible, meaning that organizations of all sizes can replicate such models with moderate investment.

For the researcher, this project has reinforced practical competencies in data cleaning, EDA, time series modeling, and dashboard creation. It has bridged the gap between theoretical study and real-world business problem-solving — a vital skill for modern business analytics professional.

The dynamic retail environment demands businesses like Walmart to continuously innovate their forecasting capabilities. By investing in robust, scalable, and continuously improving predictive analytics systems, Walmart can ensure that it remains competitive, customer-focused, and operationally efficient in the years to come.

6.9 Concluding Statement

To conclude, this project has successfully achieved its primary aim of forecasting Walmart's sales using predictive analytics. The insights generated through this analysis reaffirm the importance of data-driven decision-making in today's business landscape. By embracing advanced forecasting methods, Walmart can continue to lead the retail market, optimize its operations, and better serve millions of customers worldwide.

This research has also deepened my understanding of how analytical tools and techniques can be practically applied in a business environment. The project has strengthened my technical, analytical, and presentation skills and will serve as a solid foundation for my future career as a data-driven business professional.

The learnings and experience gained through this project will help me contribute to any organization's data analytics and strategic decision-making initiatives, enabling smarter, more efficient business operations in an increasingly data-centric world.

CHAPTER 7: BIBLIOGRAPHY

A **Bibliography** lists **books, journals, research papers, and reports**:

7.1 Books & Journals

1. Box, G. E. P., Jenkins, G. M., & Reinsel, G. C. (2015). *Time Series Analysis: Forecasting and Control* (5th ed.). Wiley.
2. Montgomery, D. C., Jennings, C. L., & Kulahci, M. (2015). *Introduction to Time Series Analysis and Forecasting* (2nd ed.). Wiley.
3. Chatfield, C. (2003). *The Analysis of Time Series: An Introduction* (6th ed.). CRC Press.
4. Makridakis, S., Wheelwright, S. C., & Hyndman, R. J. (1998). *Forecasting: Methods and Applications* (3rd ed.). Wiley.

5. Journals & Research Papers

6. Fildes, R., Goodwin, P., Lawrence, M., & Nikolopoulos, K. (2008). Effective Forecasting and Judgmental Adjustments: An Empirical Evaluation and Strategies for Improvement in Supply Chain Planning. *International Journal of Forecasting*, 24(1), 3–19.
7. Ramanathan, R., & Muyldermans, L. (2010). Identifying Demand Factors for Promotional Planning and Forecasting. *International Journal of Production Economics*, 128(2), 518–524.
8. Armstrong, J. S. (2001). Principles of Forecasting: A Handbook for Researchers and Practitioners. *International Series in Operations Research & Management Science*, Springer.

Online Sources & Company Reports

9. Walmart Inc. (2024). *Annual Report 2023–2024*. Retrieved from: <https://corporate.walmart.com/our-story/our-business>
10. Kaggle. (n.d.). *Walmart Sales Forecasting Dataset*. Retrieved from <https://www.kaggle.com/datasets/varsharam/walmart-sales-dataset-of-45stores>

Web References

11. Hyndman, R. J. (n.d.). Forecasting Principles and Practice. Retrieved from: <https://otexts.com/fpp3/>
12. Statista. (2024). *Walmart U.S. Retail Statistics*. Retrieved from: <https://www.statista.com/topics/1372/walmart>

CHAPTER 8: REFERENCES

A **References** section lists **websites, online datasets, company papers, open notebooks** — sources you *directly used*.

Websites, Datasets & Online Sources

1. Kaggle (n.d.). *Walmart Store Sales Forecasting Dataset*. Retrieved from <https://www.kaggle.com/>
2. Kaggle. (n.d.). *Walmart Sales Forecasting Dataset*. Retrieved from <https://www.kaggle.com/datasets/varsharam/walmart-sales-dataset-of-45stores>
3. GitHub. (n.d.). *Walmart Sales Forecasting Notebooks and Scripts*. Retrieved from <https://github.com/>
4. Walmart Inc. Official Website. (2024). *Company Overview*. Retrieved from: <https://corporate.walmart.com/our-story/our-business>
5. Walmart Inc. (2024). *Annual Report 2023–24*. Retrieved from: https://corporate.walmart.com/media-library/document/2023-annual-report/_proxy/document?id=00000182-2f15-d8df-a7c2-7f3f0bc20000
6. Kaggle. (n.d.). *Walmart Sales Dataset of 45 Stores*. Retrieved from: <https://www.kaggle.com/datasets>
7. Hyndman, R.J. (n.d.). *Forecasting Principles and Practice*. Retrieved from: <https://otexts.com/fpp3/>
8. Statista. (2024). *Walmart U.S. Retail Statistics*. Retrieved from: <https://www.statista.com/topics/1372/walmart>
9. Deloitte Insights. (2019). *Predictive Analytics for Retail Demand Forecasting*. Retrieved from <https://www2.deloitte.com/>
10. IBM Analytics. (2019). *Predictive Analytics for Retail*. Retrieved from <https://www.ibm.com/analytics>

CHAPTER 9: ANNEXURE.

Customer Shopping Behaviour and Factors Influencing Sales at Walmart

Purpose:

This questionnaire aims to gather detailed insights into customer shopping patterns, preferences, seasonal buying behaviour, and factors that impact sales trends at Walmart stores. The findings (hypothetical) help in designing accurate sales forecasts.

Section 1: Demographic Profile

1. Name (Optional): _____

2. Gender:

- Male
- Female
- Other

3. Age Group:

- Below 20
- 21–30
- 31–40
- 41–50
- 51–60
- Above 60

4. City/Town: _____

5. Occupation:

- Student
- Working Professional
- Self-employed / Business Owner
- Homemaker
- Retired
- Other _____

Section 2: Shopping Frequency & Patterns

6. How often do you shop at Walmart?

- 2–3 times a week
- Weekly
- Bi-weekly
- Monthly
- Occasionally

7. When do you prefer to shop the most? (Select all that apply)

- Weekdays
- Weekends
- Festivals/Holidays
- Clearance Sale Periods

8. What time of the day do you usually shop?

- Morning
- Afternoon
- Evening
- Late Night

9. Do you shop more online or in-store?

- Mostly in-store
- Mostly online
- Both equally

Section 3: Purchase Decision Triggers

10. What influences your decision to visit Walmart over other stores? (Select all that apply)

- Discounts & Offers
- Loyalty Programs
- Store Location
- Variety of Products
- Customer Service
- Other _____

11. Which of these promotions attract you the most?

- Buy One Get One (BOGO)

- Flat Discounts
- Loyalty Rewards
- Cashback
- Free Samples

12. How do you learn about Walmart promotions? (Select all that apply)

- In-store banners
- SMS/Email
- Social Media
- Newspaper/Print ads
- Friends/Family

Section 4: Holiday & Seasonal Shopping

13. Does your shopping at Walmart increase during the holiday season?

- Yes
- No

14. Which holiday season impacts your spending the most? (Select all that apply)

- Thanksgiving
- Christmas
- New Year
- Black Friday
- Back-to-School
- Others _____

15. By what percentage does your spending increase during major holidays?

- Less than 20%
- 20%–40%
- 40%–60%
- Above 60%

Section 5: Customer Loyalty & Satisfaction

16. Are you a member of Walmart's loyalty/reward program?

- Yes
- No

17. How satisfied are you with Walmart's pricing?

- Very Satisfied
- Satisfied
- Neutral
- Unsatisfied

18. How likely are you to recommend Walmart to others?

- Very Likely
- Likely
- Neutral
- Unlikely

19. What could Walmart do better to increase your shopping frequency?

Section 6: Forecast Relevance (Optional)

20. Would you like Walmart to have personalized offers based on your past purchases?

- Yes
- No

21. Do you think data-driven offers (e.g., customized discounts) can increase your spending at Walmart?

- Yes
- No

22. Any other suggestions on how Walmart can serve you better?

Annexure B: Scope for Future Study

Scope for Future Study

Forecasting in retail is an ever-evolving field, constantly shaped by changes in consumer behaviour, market dynamics, technology, and data availability. While this project successfully demonstrated how predictive analytics can be used to forecast Walmart's sales using historical sales data and classical time series models, there remains significant scope for further research and practical implementation to improve the accuracy and usefulness of such models.

1. Use of Advanced Machine Learning & Deep Learning Models

While this project focused on traditional statistical models like ARIMA and Holt-Winters, future studies could implement more sophisticated techniques. Machine learning algorithms like XGBoost, Random Forest Regressors, or even advanced deep learning models such as Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM) networks can better capture complex non-linear relationships and long-range dependencies in time series data. This could significantly enhance forecast accuracy, especially during irregular sales spikes.

2. Incorporation of Real-Time & External Data

The present study was limited to historical in-store sales and selected macroeconomic indicators. Future projects could integrate:

Real-time sales feeds directly from POS systems.

Weather data, which can heavily influence footfall and spending.

Competitor pricing and promotional activity, which can impact consumer choices.

Social media sentiment analysis, which can detect trends and potential demand shifts.

Combining these factors with core sales data could make forecasting models more robust and dynamic.

3. Regional & Store-Level Forecasting

This project analysed all 45 stores together and individually but maintained a standard modeling approach. Future studies could:

Segment stores based on demographics, regional economic factors, or store size.

Build **cluster-specific models** for stores with similar sales patterns.

Apply **hierarchical forecasting**, which combines store-level forecasts with regional or national aggregates for top-down planning.

4. Integration with ERP & Supply Chain Systems

Future research could explore how forecasts can be directly integrated into Walmart's ERP systems. This would automate supply chain activities, including inventory replenishment, workforce planning, and vendor management. Such integration can create a closed-loop system where predictions directly drive operational actions.

5. Enhancing Forecast Interpretability

While complex models can improve accuracy, they often become black boxes. Future work could explore the use of **explainable AI (XAI)** techniques to interpret how different factors contribute to sales forecasts. This can help managers trust and act on model outputs with greater confidence.

6. Simulation & Scenario Planning

Future studies could design **what-if simulations** to understand the impact of different scenarios — for example:

How would a sudden change in fuel prices affect sales?

What if a new competitor opens nearby?

How do changes in promotional strategies shift demand?

Scenario planning helps companies prepare for uncertainty and adapt quickly.

7. Continuous Learning Systems

An exciting extension would be to develop a **self-learning forecasting system** that updates itself automatically as new data comes in. This would ensure forecasts remain accurate even as trends shift due to unexpected events like economic downturns or global crises.

8. Cross-Functional Insights

Finally, future research can expand beyond sales to predict related aspects such as customer churn, lifetime value, basket analysis, and product mix optimization. Combining sales forecasting with marketing, logistics, and HR planning can make the insights truly enterprise-wide.

9. Predictive Maintenance for Supply Chain

Walmart could extend predictive analytics beyond sales to predict **equipment maintenance needs** in its vast logistics and supply chain network. By integrating IoT sensor data from refrigeration units, trucks, or warehouse machinery, Walmart can reduce unexpected breakdowns, ensuring continuous operations during peak demand periods.

10. Demand Forecasting for Online vs Offline Channels

Future research can focus on developing separate but connected forecasts for Walmart's **e-commerce sales** and **in-store sales**, which have different trends, seasonality, and influencing factors. This is especially important post-COVID, with online shopping trends evolving rapidly.

11. Hyperlocal Forecasting

A more granular study could examine **hyperlocal factors**, such as neighborhood events, local festivals, or even local sports events, which could cause spikes in sales in specific stores. This would help individual store managers plan better for short-term local demand surges.

12. Forecasting Returns and Exchanges

Another valuable extension is forecasting **product returns**, which are a significant operational cost in retail. Understanding which products, seasons, or promotions generate more returns can help Walmart design better return policies and reduce financial losses.

13. Personalized Promotions Using Predictive Analytics

A future study could develop algorithms that use purchase history data to forecast **individual customer buying patterns** and generate personalized promotions or dynamic pricing. This can drive customer loyalty and increase basket size.

14. Environmental & Sustainability Forecasting

Future research could explore forecasting **carbon footprints** or **energy usage** linked to changes in sales volumes. Walmart's supply chain is vast; aligning sales forecasts with sustainability targets can help reduce waste, optimize transportation, and lower environmental impact.

15. Risk Forecasting & Scenario Stress Testing

Expanding from sales forecasting to **risk forecasting**, researchers can build models that test how unexpected risks (pandemics, inflation spikes, fuel crises) might affect sales. Such stress testing helps Walmart prepare better contingency plans.

16. Real-Time Data Visualization & Alerts

A future study can focus on developing **automated dashboards with AI-driven alerts** for unusual sales patterns — for example, sudden dips or unexpected surges. Managers can then react in real-time to optimize staffing, supply chain adjustments, or marketing responses.

17. Multi-Product & Category-Level Forecasting

While this project focused on overall store sales, future work could drill down to **department-wise or product category-wise forecasting**. This helps optimize inventory for perishable vs. non-perishable goods, electronics, clothing, etc.

18. Exploring Cross-Brand & Supplier Collaboration

Future projects could study how Walmart can share forecast data with its **vendors and suppliers**, helping them align their production schedules with predicted demand — improving efficiency for the entire supply chain ecosystem.

19. Using Blockchain for Forecast Data Integrity

Research can examine how **blockchain** can be integrated to ensure secure, tamper-proof historical sales data. This enhances trust and transparency, especially for collaborating with multiple suppliers and partners.

20. Customer Sentiment & Social Listening

Future work can integrate **social listening tools** to forecast how viral trends, reviews, or sudden changes in customer sentiment affect product sales. This is highly relevant for trending products and seasonal stock planning.

21. Geospatial Data Integration

Incorporating **geospatial data**, like store footfall heatmaps, local traffic patterns, or population density changes, can help Walmart plan new stores, expansions, or closures with better market fit.

22. Data Privacy & Ethics

A valuable future direction is to research the **ethical implications** of using granular customer data for personalized forecasting. Ensuring compliance with data privacy laws like GDPR or CCPA will be crucial as predictive analytics grows.