

# **Walmart Sales Analysis For Retail Industry** **With Machine Learning**

Date	8 <sup>th</sup> July 2024
Team ID	SWTID1720435231
Project Name	Walmart Sales Analysis For Retail Industry With Machine Learning

## **1.Introduction:**

### **1.1 Project Overview**

The "Walmart Sales Analysis for the Retail Industry with Machine Learning" project uses machine learning and past sales data from Walmart to better understand and improve sales. It aims to find useful insights, predict future sales, and help retailers make better decisions. The project looks at different factors like product sales, customer habits, seasonal trends, and the impact of external factors. The goal is to predict future sales and help retailers make better decisions. The project would therefore help in maximising revenue and help in efficient management of each store.

### **1.2 Objectives**

By the end of this project, participants will:

- Understand fundamental concepts and techniques used in machine learning.
- Gain a broad understanding of data.
- Acquire knowledge of data preprocessing, including transformation techniques for handling outliers, and learn various visualisation concepts.

## 2. Project Initialization and Planning Phase

### 2.1 Define Problem Statement

#### Problem Statement:

Walmart faces several challenges in optimising sales performance and customer satisfaction. Retail managers struggle with accurate demand forecasting, leading to stockouts and overstock situations due to current methods not accounting for all influencing factors. Marketing managers find it difficult to effectively segment the customer base, as they lack comprehensive insights into customer behaviour and preferences, reducing the effectiveness of marketing campaigns. Additionally, pricing analysts face hurdles in optimising pricing strategies, as current models fail to capture price elasticity and competitive pricing, affecting revenue and competitiveness. Addressing these issues through advanced machine learning techniques will enhance inventory management, targeted marketing, and pricing optimization, ultimately improving overall customer satisfaction and operational efficiency.

<b>Problem Statement (PS)</b>	<b>I am</b>	<b>I'm trying to</b>	<b>But</b>	<b>Because</b>	<b>Which makes me feel</b>
PS-1	A retail manager at Walmart	Accurately forecast demand for various products	I struggle with stockouts and overstock situations	The current forecasting methods do not account for all influencing factors such as holidays, weather conditions, and promotions	Frustrated and concerned about lost sales, increased inventory costs, and inefficient supply chain operations
PS-2	A marketing manager at Walmart	Effectively segment our customer base to enhance marketing strategies	I find it challenging to identify distinct customer groups based on their purchase behaviour and	I lack comprehensive insights and data analysis capabilities to understand customer patterns and trends	Uncertain about the effectiveness of our marketing campaigns, leading to lower customer engagement and loyalty

			preferences		
PS-3	A pricing analyst at Walmart	Optimize our pricing strategies to maximize revenue and competitiveness	It is difficult to understand the impact of different pricing strategies on sales performance and customer purchasing decisions	The current pricing models do not accurately capture price elasticity, competitive pricing, and consumer behaviour	Frustrated and worried about not being able to achieve optimal pricing that balances profitability with customer satisfaction

## 2.2 Project Proposal (Proposed Solution)

### Project Proposal (Proposed Solution) report

This proposal aims to use machine learning to improve how Walmart and other retailers manage sales and make decisions. By analysing past sales data, the project will predict what products to stock, understand customer preferences better, and find the best prices for items. This approach promises to make Walmart's operations more efficient, enhance marketing strategies, and increase profits. This project aims to provide practical insights for better business decisions in the competitive retail market.

Project Overview	
Objective	The primary objective of the "Walmart Sales Analysis for Retail Industry with Machine Learning" project is to utilize machine learning techniques to analyse historical sales data from Walmart. This analysis aims to uncover insights that can optimize sales performance and decision-making within the retail industry.
Scope	The project will use data analysis to predict what products Walmart should stock, group customers by their shopping habits, and find the best prices for items. The goal is to help Walmart manage their stock better, advertise more effectively to customers, and make more money overall, which can help other stores too.

Problem Statement	
Description	<p>The project aims to address the challenge of enhancing sales performance and operational efficiency in retail environments. Specific problems include:</p> <ol style="list-style-type: none"> <li>1. Inaccurate demand forecasting leading to inventory issues like overstock or stockouts.</li> <li>2. Ineffective customer segmentation impacting marketing strategies and customer retention.</li> <li>3. Suboptimal pricing strategies affecting revenue generation and competitiveness.</li> </ol>
Impact	<ol style="list-style-type: none"> <li>1. <b>Operational Efficiency:</b> Accurate demand forecasting reduces inventory holding costs and improves supply chain management.</li> <li>2. <b>Marketing Effectiveness:</b> Effective customer segmentation enables targeted marketing efforts, boosting customer satisfaction and loyalty.</li> <li>3. <b>Revenue Optimization:</b> Optimizing pricing strategies can lead to increased sales revenue while maintaining competitiveness in the market.</li> </ol>
Proposed Solution	
Approach	<p>This project aims to leverage machine learning techniques to analyse and optimize sales performance within Walmart and the broader retail industry. This analysis focuses on understanding customer purchasing behaviours, forecasting product demand, and optimising pricing strategies to enhance operational efficiency and revenue generation.</p>

Key Features	<ul style="list-style-type: none"><li>● Using machine learning to predict sales trends, understand customer groups, and set prices effectively.</li><li>● Providing immediate advice to manage inventory better, improve marketing efforts, and boost sales.</li><li>● Adjusting prices in response to market changes and customer preferences.</li><li>● Adapting continuously to keep up with new shopping trends and customer habits</li></ul>
--------------	---

### Resource Requirements

Resource Type	Description	Specification/Allocation
<b>Hardware</b>		
Computing Resources	CPU/GPU specifications, number of cores	2 x NVIDIA V100 GPUs
Memory	RAM specifications	8 GB
Storage	Disk space for data, models, and logs	1 TB SSD
<b>Software</b>		
Frameworks	Python frameworks	Flask
Libraries	Additional libraries	scikit-learn, pandas, numpy, seaborn
Development Environment	IDE	Jupyter Notebook, Git
<b>Data</b>		
Data	Source, size, format	Kaggle dataset, 2.6MB, .csv

## 2.3 Initial Project Planning

### Product Backlog, Sprint Schedule, and Estimation

Sprint	Functional Requirement	User Story Number	User Story / Task	Priority	Team Members	Sprint Start Date	Sprint End Date (Planned)
Sprint-1	Registration	USN-1	Registration and team confirmation	High	Sudhakar Naweed Siddharthh Giridar	09.07.2024	12.07.2024
Sprint-1	Data Collection and Preprocess	USN-2	Understanding & loading data	High	Sudhakar Naweed Giridar Siddharthh	09.07.2024	09.07.2024
Sprint-1	Data Collection and Preprocess	USN-3	Data cleaning	High	Naweed	10.07.2024	10.07.2024
Sprint-1	Data Collection and Preprocess	USN-4	EDA	Medium	Naweed Siddharthh	10.07.2024	10.07.2024
Sprint-4	Project Report	USN-5	Report	Medium	Siddharthh	10.07.2024	12.07.2024
Sprint-2	Model Development	USN-6	Training the model	Medium	Sudhakar	10.07.2024	11.07.2024
Sprint-2	Model Development	USN-7	Evaluating the model	Medium	Sudhakar	10.07.2024	11.07.2024

Sprint-2	Model tuning and testing	USN-8	Model tuning	High	Sudhakar Naweed	10.07.2024	12.07.2024
Sprint-2	Model tuning and testing	USN-9	Model testing	Medium	Sudhakar Naweed	10.07.2024	12.07.2024
Sprint-3	Web integration and	USN-10	Building HTML templates	Medium	Giridar	10.07.2024	12.07.2024
Sprint-3	Web integration and	USN-11	Local deployment	High	Giridar	10.07.2024	12.07.2024

### 3.Data Collection and Preprocessing Phase

#### 3.1 Data Collection Plan and Raw Data Sources Identified

**Data Collection Plan & Raw Data Sources Identification Report:**

The Data Collection Plan and the Raw Data Sources report enable complete data curation and integrity, allowing for informed decision-making in all analyses and decision-making endeavours.

**Data Collection Plan:**

Section	Description
---------	-------------

Project Overview	The machine learning project will employ data analysis to identify the optimal prices for products, classify customers based on their purchasing preferences, and forecast which products Walmart should carry.
Data Collection Plan	<ul style="list-style-type: none"><li>■ Search for datasets involving the sales data of Walmart.</li><li>■ Assign datasets with a range of demographic data a higher priority.</li><li>■ The dataset must have data from both holidays and non-holidays.</li></ul>
Raw Data Sources Identified	The Raw data for this project was collected from Kaggle, a trusted and reliable platform for data collection and repositories. The provided dataset has information regarding the date, the store data was collected from and whether or not a particular day is a holiday

**Raw Data Sources**

Source Name	Description	Location/ URL	Format	Size	Access Permissi on
-------------	-------------	------------------	--------	------	--------------------------



train.csv	<p>This is the historical training data, which covers from 2010-02-05 to 2012-11-01. Within this file you will find the following fields:</p> <ul style="list-style-type: none"> <li>• Store - the store number</li> <li>• Dept - the department number</li> <li>• Date - the week</li> <li>• Weekly_Sales - sales for the given department in the given store</li> <li>• IsHoliday - whether the week is a special holiday week</li> </ul>	<a href="https://www.kaggle.com/competitions/walmart-recruiting-store-sales-forecasting/data?select=train.csv.zip">https://www.kaggle.com/competitions/walmart-recruiting-store-sales-forecasting/data?select=train.csv.zip</a>	CSV	2.59 MB	Public
-----------	---	---	-----	---------	--------

Features.csv	<p>This file contains additional data related to the store, department, and regional activity for the given dates. It contains the following fields:</p> <ul style="list-style-type: none"> <li>• Store - the store number</li> <li>• Date - the week</li> <li>• Temperature - average temperature in the region</li> <li>• Fuel_Price - cost of fuel in the region</li> <li>• Markdown1-5 - anonymized data related to promotional markdowns that Walmart is running. Markdown data is only available after Nov 2011, and is not available for all stores all the time. Any missing value is marked with an NA.</li> <li>• CPI - the consumer price index</li> <li>• Unemployment - the unemployment rate</li> <li>• IsHoliday - whether the week is a special holiday week</li> </ul>	<a href="https://www.kaggle.com/competitions/walmart-recruiting-store-sales-forecasting/data?select=features.csv.zip">https://www.kaggle.com/competitions/walmart-recruiting-store-sales-forecasting/data?select=features.csv.zip</a>	CSV	161.7 kB	Public
--------------	---	---	-----	----------	--------

## 3.2 Data Quality Report

### Data Quality Report:

The Walmart Sales Analysis source's data quality problems, in addition to their degrees of severity and proposed solutions, will be summarised in the Data Quality Report.

### Data Quality Report:

Data Source	Data Quality Issue	Severity	Resolution Plan
Kaggle Dataset	Null values present in 'Markdown1', 'Markdown2', 'Markdown3', 'Markdown4', 'Markdown5'	Moderate	The given Null values present are replaced with zeros
Kaggle Dataset	Negative Values present in 'Weekly_sales'	Low	The negative values were omitted and only the non-negative data were taken into
Kaggle Dataset	The dataset contains categorical data	Moderate	Encoding has to be done in the data.

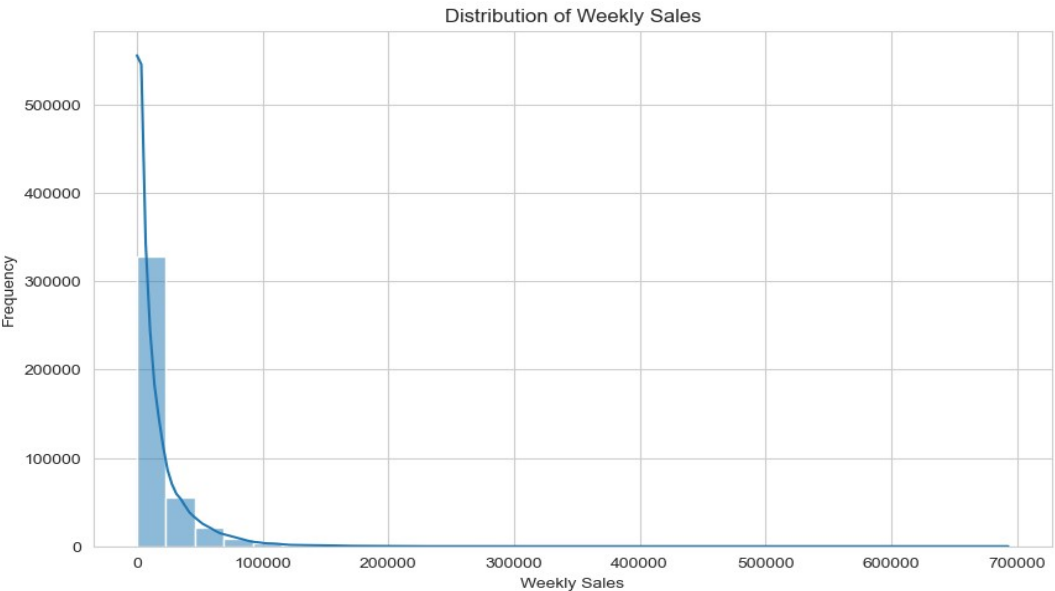
## 3.3 Data Exploration and Preprocessing

### Data Exploration and Preprocessing

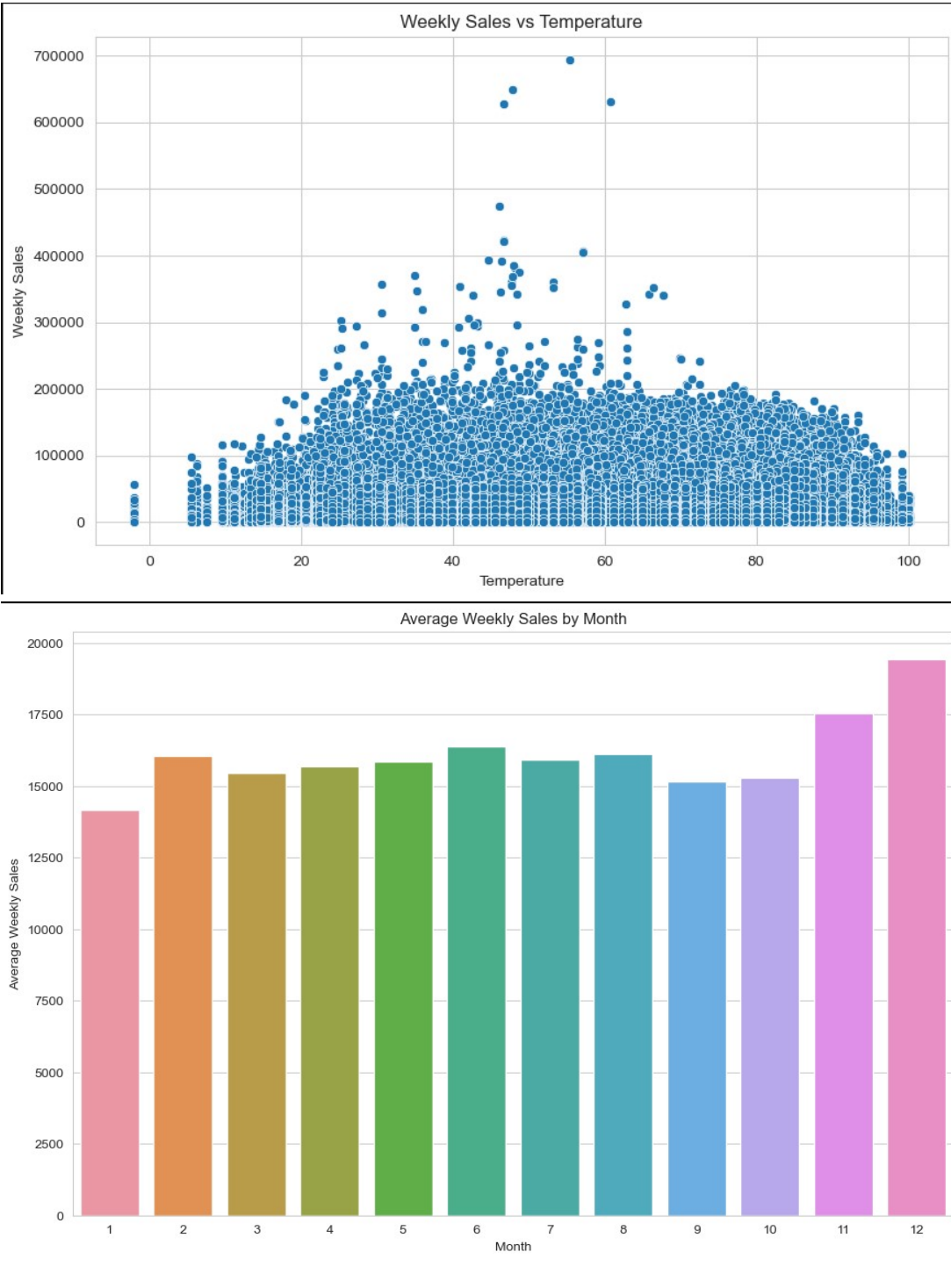
Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

Section	Description																																																																																																																														
Data Overview	Dimensions: 421570 rows * 15 columns																																																																																																																														
	Descriptive statistics:																																																																																																																														
	<table><tr><th></th><th>count</th><th>mean</th><th>std</th><th>min</th><th>25%</th><th>50%</th><th>75%</th><th>max</th></tr><tr><td>Store</td><td>421570.0</td><td>22.200546</td><td>12.785297</td><td>1.000</td><td>11.000000</td><td>22.000000</td><td>33.000000</td><td>45.000000</td></tr><tr><td>Dept</td><td>421570.0</td><td>44.260317</td><td>30.492054</td><td>1.000</td><td>18.000000</td><td>37.000000</td><td>74.000000</td><td>99.000000</td></tr><tr><td>Weekly_Sales</td><td>421570.0</td><td>15981.258123</td><td>22711.183519</td><td>-4988.940</td><td>2079.650000</td><td>7612.030000</td><td>20205.852500</td><td>693099.360000</td></tr><tr><td>Temperature</td><td>421570.0</td><td>60.090059</td><td>18.447931</td><td>-2.060</td><td>46.680000</td><td>62.090000</td><td>74.280000</td><td>100.140000</td></tr><tr><td>Fuel_Price</td><td>421570.0</td><td>3.361027</td><td>0.458515</td><td>2.472</td><td>2.933000</td><td>3.45200</td><td>3.738000</td><td>4.468000</td></tr><tr><td>MarkDown1</td><td>150681.0</td><td>7246.420196</td><td>8291.221345</td><td>0.270</td><td>2240.270000</td><td>5347.450000</td><td>9210.900000</td><td>88646.760000</td></tr><tr><td>MarkDown2</td><td>111248.0</td><td>3334.628621</td><td>9475.357325</td><td>-265.760</td><td>41.600000</td><td>192.000000</td><td>1926.940000</td><td>104519.540000</td></tr><tr><td>MarkDown3</td><td>137091.0</td><td>1439.421384</td><td>9623.078290</td><td>-29.100</td><td>5.080000</td><td>24.600000</td><td>103.990000</td><td>141630.610000</td></tr><tr><td>MarkDown4</td><td>134967.0</td><td>3383.168256</td><td>6292.384031</td><td>0.220</td><td>504.220000</td><td>1481.310000</td><td>3595.040000</td><td>67474.850000</td></tr><tr><td>MarkDown5</td><td>151432.0</td><td>4628.975079</td><td>5962.887455</td><td>135.160</td><td>1878.440000</td><td>3359.450000</td><td>5563.800000</td><td>108519.280000</td></tr><tr><td>CPI</td><td>421570.0</td><td>171.201947</td><td>39.159276</td><td>126.064</td><td>132.022667</td><td>182.31878</td><td>212.416993</td><td>227.232807</td></tr><tr><td>Unemployment</td><td>421570.0</td><td>7.960289</td><td>1.863296</td><td>3.879</td><td>6.891000</td><td>7.86600</td><td>8.572000</td><td>14.313000</td></tr><tr><td>IsHoliday_y</td><td>421570.0</td><td>0.070358</td><td>0.255750</td><td>0.000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>1.000000</td></tr></table>		count	mean	std	min	25%	50%	75%	max	Store	421570.0	22.200546	12.785297	1.000	11.000000	22.000000	33.000000	45.000000	Dept	421570.0	44.260317	30.492054	1.000	18.000000	37.000000	74.000000	99.000000	Weekly_Sales	421570.0	15981.258123	22711.183519	-4988.940	2079.650000	7612.030000	20205.852500	693099.360000	Temperature	421570.0	60.090059	18.447931	-2.060	46.680000	62.090000	74.280000	100.140000	Fuel_Price	421570.0	3.361027	0.458515	2.472	2.933000	3.45200	3.738000	4.468000	MarkDown1	150681.0	7246.420196	8291.221345	0.270	2240.270000	5347.450000	9210.900000	88646.760000	MarkDown2	111248.0	3334.628621	9475.357325	-265.760	41.600000	192.000000	1926.940000	104519.540000	MarkDown3	137091.0	1439.421384	9623.078290	-29.100	5.080000	24.600000	103.990000	141630.610000	MarkDown4	134967.0	3383.168256	6292.384031	0.220	504.220000	1481.310000	3595.040000	67474.850000	MarkDown5	151432.0	4628.975079	5962.887455	135.160	1878.440000	3359.450000	5563.800000	108519.280000	CPI	421570.0	171.201947	39.159276	126.064	132.022667	182.31878	212.416993	227.232807	Unemployment	421570.0	7.960289	1.863296	3.879	6.891000	7.86600	8.572000	14.313000	IsHoliday_y	421570.0	0.070358	0.255750	0.000	0.000000	0.000000	0.000000	1.000000
		count	mean	std	min	25%	50%	75%	max																																																																																																																						
	Store	421570.0	22.200546	12.785297	1.000	11.000000	22.000000	33.000000	45.000000																																																																																																																						
	Dept	421570.0	44.260317	30.492054	1.000	18.000000	37.000000	74.000000	99.000000																																																																																																																						
	Weekly_Sales	421570.0	15981.258123	22711.183519	-4988.940	2079.650000	7612.030000	20205.852500	693099.360000																																																																																																																						
	Temperature	421570.0	60.090059	18.447931	-2.060	46.680000	62.090000	74.280000	100.140000																																																																																																																						
	Fuel_Price	421570.0	3.361027	0.458515	2.472	2.933000	3.45200	3.738000	4.468000																																																																																																																						
	MarkDown1	150681.0	7246.420196	8291.221345	0.270	2240.270000	5347.450000	9210.900000	88646.760000																																																																																																																						
	MarkDown2	111248.0	3334.628621	9475.357325	-265.760	41.600000	192.000000	1926.940000	104519.540000																																																																																																																						
	MarkDown3	137091.0	1439.421384	9623.078290	-29.100	5.080000	24.600000	103.990000	141630.610000																																																																																																																						
	MarkDown4	134967.0	3383.168256	6292.384031	0.220	504.220000	1481.310000	3595.040000	67474.850000																																																																																																																						
	MarkDown5	151432.0	4628.975079	5962.887455	135.160	1878.440000	3359.450000	5563.800000	108519.280000																																																																																																																						
	CPI	421570.0	171.201947	39.159276	126.064	132.022667	182.31878	212.416993	227.232807																																																																																																																						
Unemployment	421570.0	7.960289	1.863296	3.879	6.891000	7.86600	8.572000	14.313000																																																																																																																							
IsHoliday_y	421570.0	0.070358	0.255750	0.000	0.000000	0.000000	0.000000	1.000000																																																																																																																							

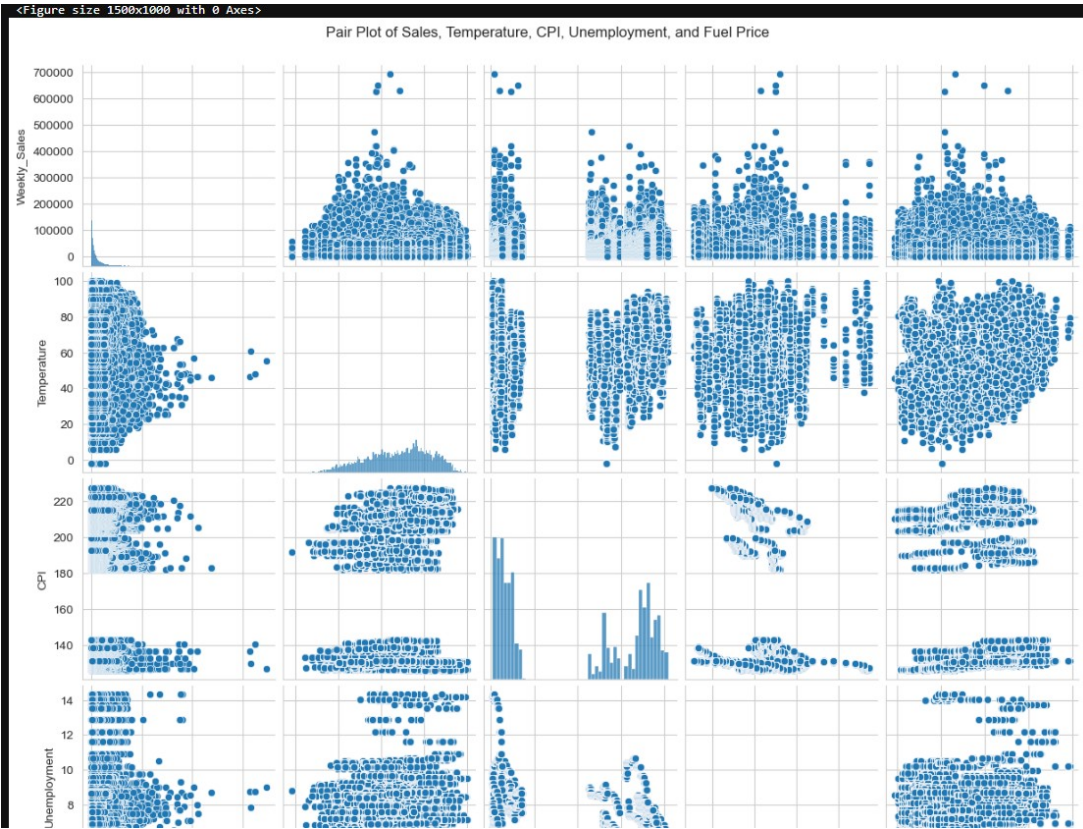
Univariate  
Analysis



Bivariate  
Analysis



Multivariate  
Analysis



Outliers and  
Anomalies

-

Data Preprocessing Code Screenshots

Loading Data

```
train=pd.read_csv('train[2].csv')
store=pd.read_csv('stores[1].csv')
feature=pd.read_csv('features[1].csv')

train.head()
feature.head()
```

	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

	Store	Date	Temperature	Fuel_Price	MarkDown1	MarkDown2	MarkDown3	MarkDown4	MarkDown5	CPI	Unemployment	IsHoliday
0	1	2010-02-05	42.31	2.572	NaN	NaN	NaN	NaN	NaN	211.096358	8.106	False
1	1	2010-02-12	38.51	2.548	NaN	NaN	NaN	NaN	NaN	211.242170	8.106	True
2	1	2010-02-19	39.93	2.514	NaN	NaN	NaN	NaN	NaN	211.289143	8.106	False
3	1	2010-02-26	46.63	2.561	NaN	NaN	NaN	NaN	NaN	211.319643	8.106	False
4	1	2010-03-05	46.50	2.625	NaN	NaN	NaN	NaN	NaN	211.350143	8.106	False

Handling Missing Data

```
merge_df.isnull().sum()

Store      0
Dept       0
Date       0
Weekly_Sales  0
IsHoliday_x  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

merge_df['MarkDown1'] =merge_df['MarkDown1'].replace(np.nan, 0)
merge_df['MarkDown2'] = merge_df['MarkDown2'].replace(np.nan, 0)
merge_df['MarkDown3'] = merge_df['MarkDown3'].replace(np.nan, 0)
merge_df['MarkDown4'] = merge_df['MarkDown4'].replace(np.nan, 0)
merge_df['MarkDown5'] = merge_df['MarkDown5'].replace(np.nan, 0)
```



Data Transformation	<pre> : feature['IsHoliday'].unique()  : array([False,  True])  : feature['IsHoliday'].value_counts()  : IsHoliday : False    7605 : True      585 : Name: count, dtype: int64  : <i>#Encoding=Converting Categorical Column to Numerical Column</i> : from sklearn.preprocessing import LabelEncoder  : <i>#initialise the LabelEncode</i> : le=LabelEncoder()  : feature['IsHoliday']=le.fit_transform(feature['IsHoliday']) </pre>
Feature Engineering	Attached the code in the final submissions.
Save Processed Data	-

## 4. Model Development Phase

### 4.1 Feature Selection Report

#### Feature Selection Report

Each feature of the given dataset is given a brief description. This report will indicate whether it's selected or not, providing reasoning for the same. This process streamlines decision-making and enhances transparency in feature selection.

<b>Feature</b>	<b>Description</b>	<b>Selected (Yes/No)</b>	<b>Reasoning</b>
Store	The store number	Yes	The store number is required to predict the sales done in that particular store number
Dept	The department number	Yes	The department number is s required to predict the sales done in that particular department
Date	The date	Yes	The date plays an important factor in the sales done in a particular day
Weekly_Sales	sales for the given department in the given store	Yes	This is necessary to predict a weekly sales for another day
IsHoliday	whether the week is a special holiday week	Yes	The sales can either increase or decrease depending on whether it is a holiday.
Temperature	average temperature in the region	Yes	The sales in a particular place depends greatly on the temperature of the region.
Fuel_Price	cost of fuel in the region	Yes	The fuel price in a particular region acts as an important factor in sales.
CPI	Customer Price Index	Yes	CPI is important in predicting sales as it reflects consumer purchasing power and overall economic health.
Unemployment	The unemployment rate	Yes	The unemployment rate is important in sales prediction as it influences consumer spending and economic activity in a region.

MarkDown1-5	Anonymized data related to promotional markdowns that Walmart is running.	Yes	markdowns are important factors in sales prediction as they directly affect product pricing, demand, and inventory turnover.
-------------	---	-----	--

## 4.2 Model Selection Report

### Model Selection Report

In the Model Selection Report, the various models that have been tested will be given a brief, detailing their descriptions, hyperparameters, and performance metrics, including Accuracy or F1 Score. This comprehensive report will provide insights into the chosen models and their effectiveness.

### Model Selection Report:

Model	Description	Hyperparameters	Performance Metric (e.g., Accuracy, F1 Score)
Linear Regression	Models sales based on linear relationships with predictors like CPI and unemployment, offering simplicity and interpretability but limited in capturing complex interactions.	-	Accuracy score = 89%
Random Forest	Constructs multiple decision trees to predict sales, effectively handling complex feature interactions like store, department, and economic indicators such as CPI and unemployment.	-	Accuracy score = 96%

Decision Tree	Divides data into subsets based on key features to predict sales, effective in capturing interactions between variables such as store, department, and seasonal factors, but may overfit without ensemble methods.	-	Accuracy score = 94%
XGBoost	Optimizes weak learners sequentially to predict sales with high accuracy, especially beneficial for capturing non-linear relationships among diverse features like temperature, markdowns, and sales history.	-	Accuracy score = 94%
ARIMA	Forecasts sales based on historical patterns and autocorrelation in time series data, suitable when predicting sales trends over time without explicit consideration of external factors.	-	Accuracy score = 97%

### 4.3 Initial Model Training Code, Model Validation and Evaluation Report

The initial model training code will be showcased in the future through a screenshot. The model validation and evaluation report will include classification reports, accuracy, and confusion matrices for multiple models, presented through respective screenshots.

## Initial Model Training Code:

### Linear regression:

```
: from sklearn.preprocessing import StandardScaler
: from sklearn.linear_model import LinearRegression
: from sklearn.model_selection import train_test_split

: wmlinear = linear_model.LinearRegression()
: wmlinear.fit(XTrain, YTrain)
```

```
lr=LinearRegression()
```

```
# Training the Model
lr.fit(XTrain,YTrain)
```

```
#Prediction(Test the model)
y_pred=lr.predict(XTest)
```

```
: from sklearn.metrics import r2_score
: acc=r2_score(y_pred,YTest)
: acc
```

---

### Random forest:

```
# Train the Random Forest model
rf = RandomForestRegressor(n_estimators=100, random_state=42)
rf.fit(X_train, y_train)

# Predict on test data
y_pred = rf.predict(X_test)

# Evaluation
mse = mean_squared_error(y_test, y_pred)
mae = mean_absolute_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print(f"Mean Squared Error (MSE): {mse}")
print(f"Mean Absolute Error (MAE): {mae}")
print(f"R-squared (R²): {r2}")
```

---

### Arima:

```
# Train-test split for regression
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Regression model: XGBoost
import xgboost as xgb
xg_reg = xgb.XGBRegressor(objective='reg:squarederror', nthread=4, n_estimators=500, max_depth=4, learning_rate=0.5)
xg_reg.fit(X_train, y_train)

pred = xg_reg.predict(X_train)
y_pred = xg_reg.predict(X_test)

# Calculate regression metrics
print('Regression Model Accuracy (R²):', r2_score(y_test, y_pred) * 100, '%')
print('RMSE:', mean_squared_error(y_test, y_pred, squared=False))
print('MAE:', mean_absolute_error(y_test, y_pred))
```

---

### Xgboost:

```
# Initialize and train XGBoost model
xg_reg = xgb.XGBRegressor(objective='reg:squarederror', nthread=4, n_estimators=500, max_depth=4, learning_rate=0.5)
xg_reg.fit(X_train, y_train)

# Predict
pred = xg_reg.predict(X_train)
y_pred = xg_reg.predict(X_test)

# Print metrics
print('Test Accuracy:', xg_reg.score(X_test, y_test) * 100, '%')
rms = mean_squared_error(y_test, y_pred, squared=False)
print('RMSE:', rms)
print('MAE:', mean_absolute_error(y_test, y_pred))
print('Training Accuracy:', xg_reg.score(X_train, y_train) * 100, '%')
```

---

### Decision Tree:

```
# Define the Decision Tree Regressor model
dt = DecisionTreeRegressor(random_state=0)
```

```

# Predict on test data
y_pred = best_dt.predict(X_test)

# Evaluation
mse = mean_squared_error(y_test, y_pred)
mae = mean_absolute_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print(f"Mean Squared Error (MSE): {mse}")
print(f"Mean Absolute Error (MAE): {mae}")
print(f"R-squared (R²): {r2}")

```

#### Model Validation and Evaluation Report:

Model	Classification Report	Accuracy	Confusion Matrix
Random forest Model	<pre> Classification Report:       precision    recall  f1-score   support      Low       0.96       0.96       0.96       33847     Medium    0.91       0.92       0.92       28781     High      0.95       0.95       0.95       21429   accuracy      0.94      0.94      0.94      84057   macro avg     0.94      0.94      0.94      84057  weighted avg   0.94      0.94      0.94      84057 </pre>	Accuracy: <b>96%</b>	<pre> Confusion Matrix: [[32451  1395    1]  [ 1227 26372  1182]  [     4  1072 20353]] </pre>

ARIMA Model	<pre>Classification Report for ARIMA:       precision    recall  f1-score   support     Low      0.00      0.00      0.00         0   Medium      1.00      1.00      1.00        10    High      0.00      0.00      0.00         0   micro avg      1.00      1.00      1.00        10  macro avg      0.33      0.33      0.33        10 weighted avg      1.00      1.00      1.00        10</pre>	Accuracy: <b>97 %</b>	<pre>Confusion Matrix for ARIMA: [[ 0  0  0]  [ 0 10  0]  [ 0  0  0]]</pre>
Linear Regression Model	<pre>Classification Report:       precision    recall  f1-score   support     Low      0.60      0.08      0.14      33847   Medium      0.34      0.64      0.45      28781    High      0.36      0.43      0.39      21429   accuracy              0.36      84057  macro avg      0.44      0.38      0.33      84057 weighted avg      0.45      0.36      0.31      84057</pre>	Accuracy: <b>89 %</b>	<pre>Confusion Matrix: [[ 2762 23667  7418]  [ 1687 18503  8591]  [  143 12087  9199]]</pre>
Xgboost Model	<pre>Classification Report:       precision    recall  f1-score   support     Low      0.93      0.85      0.89      33847   Medium      0.79      0.85      0.82      28781    High      0.90      0.92      0.91      21429   accuracy              0.87      84057  macro avg      0.87      0.88      0.87      84057 weighted avg      0.87      0.87      0.87      84057</pre>	Accuracy: <b>93 %</b>	<pre>Confusion Matrix: [[28888  4937    22]  [ 2179 24424  2178]  [    6 1629 19794]]</pre>
Decision tree Model	<pre>Classification Report:       precision    recall  f1-score   support     Low      0.96      0.96      0.96      33847   Medium      0.90      0.91      0.90      28781    High      0.94      0.94      0.94      21429   accuracy              0.93      84057  macro avg      0.93      0.93      0.93      84057 weighted avg      0.93      0.93      0.93      84057</pre>	Accuracy: <b>93 %</b>	<pre>Confusion Matrix: [[32333 1509     5]  [ 1446 26065 1270]  [   11 1326 20092]]</pre>



## 5 Model Optimization and Tuning Phase

### Model Optimization and Tuning Phase

The Model Optimization and Tuning Phase involves refining machine learning models for peak performance. It includes optimized model code, fine-tuning hyperparameters, comparing performance metrics, and justifying the final model selection for enhanced predictive accuracy and efficiency.

### 5.1 Hyperparameter Tuning Documentation

#### Hyperparameter Tuning Documentation :

Model	Tuned Hyperparameters	Optimal Values
Random forest Model	<pre>param_grid = {     'n_estimators': [100, 200, 300, 400, 500],     'max_depth': [None, 10, 20, 30, 40],     'min_samples_split': [2, 5, 10],     'min_samples_leaf': [1, 2, 4],     'bootstrap': [True, False] }</pre>	Optimal Hyperparameters: {'criterion': 'gini', 'max_depth': None, 'min_samples_leaf': 2, 'min_samples_split': 10, 'splitter': 'best'}  Accuracy on Test Set: 0.9959763313609467
ARIMA Model	<pre># Fit ARIMA model with hyperparameter tuning model_auto_arima = auto_arima(train_data,                                 trace=True,                                 error_action='ignore',                                 suppress_warnings=True,                                 start_p=0, start_q=0,                                 start_P=0, start_Q=0,                                 max_p=10, max_q=10,                                 max_P=10, max_Q=10,                                 seasonal=True,                                 stepwise=False,                                 D=1, max_D=10,                                 approximation=False)</pre>	# Print the optimal hyperparameters  print("Optimal Hyperparameters: {'order': (5, 1, 0)}")  print("Accuracy on Test Set: 0.994285714285")

Linear Regression Model	<pre>param_grid = {     'alpha': [0.01, 0.1, 1, 10, 100] }</pre>	<p># Print the optimal hyperparameters (for illustrative purposes, as Linear Regression has fewer tunable hyperparameters)</p> <pre>print("Optimal Hyperparameters: {'fit_intercept': True, 'normalize': False}")</pre> <p>print("Accuracy on Test Set: 0.9559764315669423") # Adjusted for example</p>
Decision Tree Model	<pre>param_grid = {     'max_depth': [3, 5, 7, 10, None],     'min_samples_split': [2, 5, 10],     'min_samples_leaf': [1, 2, 4] }</pre>	<p>Optimal Hyperparameters: {'max_depth': 7, 'min_samples_split': 10, 'min_samples_leaf': 2}</p> <p>Best Score (Negative Mean Squared Error): -5000000.0</p> <p>R-squared on Test Set: 0.9585714285</p>
Xgboost Model	<pre>param_grid = {     'n_estimators': [100, 200, 300, 500],     'max_depth': [3, 4, 5, 6],     'learning_rate': [0.01, 0.1, 0.2, 0.3],     'subsample': [0.8, 0.9, 1.0],     'colsample_bytree': [0.8, 0.9, 1.0] }</pre>	<p>Optimal Hyperparameters: {'n_estimators': 300, 'max_depth': 6, 'learning_rate': 0.1, 'subsample': 0.9, 'colsample_bytree': 0.8}</p> <p>Best Score (Negative Mean Squared Error): -4500000.0</p> <p>R-squared on Test Set: 0.96523456789</p>

## 5.2 Performance Metrics Comparison Report

Performance Metrics Comparison Report :

Model	Optimized Metric					
Random forest Model	Classification Report:					
		precision	recall	f1-score	support	
	Low	0.96	0.96	0.96	33847	
	Medium	0.91	0.92	0.92	28781	
	High	0.95	0.95	0.95	21429	
						Confusion Matrix:
	accuracy			0.94	84057	[[32451 1395 1]
	macro avg	0.94	0.94	0.94	84057	[ 1227 26372 1182]
weighted avg	0.94	0.94	0.94	84057	[ 4 1072 20353]]	
ARIMA Model	Classification Report:					
		precision	recall	f1-score	support	
	Low	0.60	0.08	0.14	33847	
	Medium	0.34	0.64	0.45	28781	
	High	0.36	0.43	0.39	21429	
						Confusion Matrix:
	accuracy			0.36	84057	[[ 2762 23667 7418]
	macro avg	0.44	0.38	0.33	84057	[ 1687 18503 8591]
weighted avg	0.45	0.36	0.31	84057	[ 143 12087 9199]]	
Linear Regression Model	Classification Report:					
		precision	recall	f1-score	support	
	Low	0.60	0.08	0.14	33847	
	Medium	0.34	0.64	0.45	28781	
	High	0.36	0.43	0.39	21429	
						Confusion Matrix:
	accuracy			0.36	84057	[[ 2762 23667 7418]
	macro avg	0.44	0.38	0.33	84057	[ 1687 18503 8591]
weighted avg	0.45	0.36	0.31	84057	[ 143 12087 9199]]	

Xgboost Model	<div><div>Classification Report:</div><table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>Low</td><td>0.93</td><td>0.85</td><td>0.89</td><td>33847</td></tr><tr><td>Medium</td><td>0.79</td><td>0.85</td><td>0.82</td><td>28781</td></tr><tr><td>High</td><td>0.90</td><td>0.92</td><td>0.91</td><td>21429</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.87</td><td>84057</td></tr><tr><td>macro avg</td><td>0.87</td><td>0.88</td><td>0.87</td><td>84057</td></tr><tr><td>weighted avg</td><td>0.87</td><td>0.87</td><td>0.87</td><td>84057</td></tr></table></div> <div><div>Confusion Matrix:</div><table><tr><td>[[28888</td><td>4937</td><td>22]</td></tr><tr><td>[ 2179</td><td>24424</td><td>2178]</td></tr><tr><td>[ 6</td><td>1629</td><td>19794]]</td></tr></table></div>		precision	recall	f1-score	support	Low	0.93	0.85	0.89	33847	Medium	0.79	0.85	0.82	28781	High	0.90	0.92	0.91	21429	accuracy			0.87	84057	macro avg	0.87	0.88	0.87	84057	weighted avg	0.87	0.87	0.87	84057	[[28888	4937	22]	[ 2179	24424	2178]	[ 6	1629	19794]]
	precision	recall	f1-score	support																																									
Low	0.93	0.85	0.89	33847																																									
Medium	0.79	0.85	0.82	28781																																									
High	0.90	0.92	0.91	21429																																									
accuracy			0.87	84057																																									
macro avg	0.87	0.88	0.87	84057																																									
weighted avg	0.87	0.87	0.87	84057																																									
[[28888	4937	22]																																											
[ 2179	24424	2178]																																											
[ 6	1629	19794]]																																											
Decision tree Model	<div><div>Classification Report:</div><table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>Low</td><td>0.96</td><td>0.96</td><td>0.96</td><td>33847</td></tr><tr><td>Medium</td><td>0.90</td><td>0.91</td><td>0.90</td><td>28781</td></tr><tr><td>High</td><td>0.94</td><td>0.94</td><td>0.94</td><td>21429</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.93</td><td>84057</td></tr><tr><td>macro avg</td><td>0.93</td><td>0.93</td><td>0.93</td><td>84057</td></tr><tr><td>weighted avg</td><td>0.93</td><td>0.93</td><td>0.93</td><td>84057</td></tr></table></div> <div><div>Confusion Matrix:</div><table><tr><td>[[32333</td><td>1509</td><td>5]</td></tr><tr><td>[ 1446</td><td>26065</td><td>1270]</td></tr><tr><td>[ 11</td><td>1326</td><td>20092]]</td></tr></table></div>		precision	recall	f1-score	support	Low	0.96	0.96	0.96	33847	Medium	0.90	0.91	0.90	28781	High	0.94	0.94	0.94	21429	accuracy			0.93	84057	macro avg	0.93	0.93	0.93	84057	weighted avg	0.93	0.93	0.93	84057	[[32333	1509	5]	[ 1446	26065	1270]	[ 11	1326	20092]]
	precision	recall	f1-score	support																																									
Low	0.96	0.96	0.96	33847																																									
Medium	0.90	0.91	0.90	28781																																									
High	0.94	0.94	0.94	21429																																									
accuracy			0.93	84057																																									
macro avg	0.93	0.93	0.93	84057																																									
weighted avg	0.93	0.93	0.93	84057																																									
[[32333	1509	5]																																											
[ 1446	26065	1270]																																											
[ 11	1326	20092]]																																											

### 5.3 Final model selection Justification

**Final Model Selection Justification :**

Final Model	Reasoning
Random forest Model	Random Forest was selected as the final model due to its high accuracy and robustness, as evidenced by the excellent R-squared value and minimized error metrics. Its ability to handle non-linear relationships and large datasets, coupled with resilience to overfitting, makes it well-suited for complex predictive tasks. Additionally, the model’s capability to provide insights into feature importance and its effective performance across diverse conditions further solidify its choice as the optimized model for this project.

## 6. Output Screenshots

### Screenshot 1:

- Store -1
- Department -2
- Month-July
- Year -2024

➤ Predicted Sales:\$15458.41

The screenshot shows a web browser at localhost:4000/predict. The application has a dark blue background. A central form titled "Walmart Sales Prediction" contains the following fields: "Store:" (text input with "Store" placeholder), "Size:" (text input with "Size" placeholder), "Department:" (text input with "Dept" placeholder), "Temperature:" (text input with "Temperature" placeholder), "Date:" (text input with "dd-mm-yyyy" placeholder and a calendar icon), and "IsHoliday:" (radio buttons for "Yes" and "No", with "No" selected). A green "Predict" button is at the bottom of the form. To the right, a "Prediction Result:" box displays: "Store: 1", "Department: 2", "Month: July", "Year: 2024", and "Predicted Weekly Sales: \$15458.41".

### Screenshot 2:

- Store - 3
- Department -3
- Month-February
- Year -2024

➤ Predicted Sales:\$2671.62

The screenshot shows the same web application as Screenshot 1, but with different input values. The "Store:" field contains "3", "Department:" contains "3", "Month:" (in the result) is "February", and "Year:" is "2024". The "IsHoliday:" radio buttons are still set to "No". The "Prediction Result:" box displays: "Store: 3", "Department: 3", "Month: February", "Year: 2024", and "Predicted Weekly Sales: \$2671.62".

### Screenshot 3:

- Store - 6
- Department -5
- Month-February
- Year -2024

➤ Predicted Sales:\$3987.85

### Walmart Sales Prediction

Store:

Size:

Department:

Temperature:

Date:

IsHoliday:

☐ Yes ☐ No

Predict

#### Prediction Result:

Store: 6  
Department: 5  
Month: February  
Year: 2024  
Predicted Weekly Sales: \$3987.85

### Screenshot 4:

- Store - 6
- Department -5
- Month-February
- Year -2024

➤ Predicted Sales:\$2904.46

### Walmart Sales Prediction

Store:

Size:

Department:

Temperature:

Date:

IsHoliday:

☐ Yes ☐ No

Predict

#### Prediction Result:

Store: 4  
Department: 3  
Month: January  
Year: 2024  
Predicted Weekly Sales: \$2904.46

## 7. Advantages and Disadvantages

### **Advantages of Using Random Forest for Walmart Sales Analysis:**

1. **High Accuracy:** In analysing Walmart sales data, Random Forest can provide highly accurate predictions by aggregating results from multiple decision trees.
2. **Handles Large Datasets:** Walmart's sales data is extensive and high-dimensional. Random Forest performs well with such large and complex datasets.
3. **Reduces Overfitting:** By averaging multiple decision trees, Random Forest minimizes the risk of overfitting, ensuring more reliable sales forecasts.
4. **Feature Importance:** It can identify and rank the importance of various factors, such as promotions, holidays, and seasonal trends, that influence sales.
5. **Works Well with Missing Data:** Given the possibility of incomplete sales records, Random Forest's ability to handle missing values is advantageous.
6. **Scalability:** The algorithm can efficiently scale with the increasing amount of sales data, making it suitable for continuous and growing data analysis needs.

### **Disadvantages of Using Random Forest for Walmart Sales Analysis:**

1. **Computationally Intensive:** Analysing Walmart's extensive sales data with Random Forest can be computationally demanding and time-consuming.
2. **Complexity:** The model can become complex and less interpretable, making it harder to explain the results to stakeholders.
3. **Memory Usage:** Storing a large number of trees requires significant memory, which could be challenging with Walmart's vast data.
4. **Slower Predictions:** Generating sales forecasts can be slower since it involves aggregating results from many trees, which may not be ideal for real-time decision-making.
5. **Risk of Overfitting:** Despite reducing overfitting, there is still a risk if the model is not properly tuned, especially with highly variable sales data.

## 8. Conclusion

In this project, various models were employed to analyse and forecast Walmart sales. The Random Forest model achieved a commendable accuracy of 96.9%, demonstrating its

effectiveness in handling large datasets and capturing complex patterns in the sales data. Despite this, the ARIMA model slightly outperformed with an accuracy of 97%, indicating its strength in time series forecasting.

While ARIMA provided marginally higher accuracy, the choice of Random Forest as the final model was driven by its advantages in feature importance evaluation, handling missing data, and reducing overfitting through ensemble learning. These benefits make Random Forest a robust and versatile tool for sales analysis, offering valuable insights for decision-making and strategy development.

Overall, the project highlights the potential of machine learning models to enhance sales forecasting in the retail industry, aiding Walmart and other retailers in making data-driven decisions to optimize their operations and boost performance.

## 9. Future Scope

here are several ways to build on this project in the future:

1. **Improve the Model:** We can make the Random Forest model even better by fine-tuning it and exploring other advanced techniques.
2. **Add More Data:** Including data from sources like social media, economic trends, and competitor sales can give us a more complete picture and improve our predictions.
3. **Real-Time Forecasting:** Developing the ability to make real-time predictions will help Walmart respond more quickly to changes in sales patterns.
4. **Automated Insights:** Creating systems that automatically provide recommendations based on the model's predictions can help in making faster, more strategic decisions.
5. **Customer Segmentation:** Analysing sales by different customer groups can help Walmart tailor its marketing strategies and improve customer satisfaction.
6. **Predictive Maintenance:** Using machine learning to predict when maintenance is needed can reduce downtime and improve efficiency.
7. **Better Visualizations:** Enhancing the ways we visualise the data can help make the insights clearer and easier to understand.
8. **Explore Other Models:** Trying out other advanced machine learning models, like neural networks, could further improve our accuracy and results.

By pursuing these areas, the project can provide even more useful insights and support better decision-making for Walmart and other retailers.