

# Sales Data Analysis

This project, titled “Sales Data Analysis,” focuses on analyzing Walmart’s sales dataset to derive actionable insights that can support better decision-making. The dataset includes information such as store details, weekly sales, holiday indicators, temperature, fuel price, and economic factors like CPI and unemployment rate. By performing detailed data cleaning, preprocessing, and exploratory analysis, the project aims to identify sales patterns, seasonal fluctuations, and the impact of external variables on sales performance.

Through the use of Python (VS Code) the project transforms raw data into meaningful visualizations and metrics. The insights gained from this analysis help highlight top-performing stores and products, understand seasonal demand, and identify business opportunities for improving inventory management and sales forecasting.

## Objective or Problem Statement

The goal of this project is to analyze and visualize Walmart sales data to uncover meaningful business insights. It aims to identify sales trends, seasonal patterns, and product performance across different stores and categories.

The project seeks to solve key business challenges such as recognizing top and low-performing products, understanding factors that influence sales (like holidays and discounts), and improving inventory management and sales forecasting through data-driven insights.

## Tools and Technologies Used

**Python (VS Code Environment):** Used for data cleaning, preprocessing, and exploratory data analysis (EDA).

- **Libraries:**
  - **Pandas** – for data manipulation and analysis
  - **Matplotlib / Seaborn** – for creating visualizations and identifying sales trends

## Database Used

The dataset used in this project is the Walmart Sales Data, obtained from Kaggle. It contains detailed sales information from multiple Walmart stores.

The dataset includes approximately 10,000+ records of sales transactions collected over a specific time period.

It consists of various fields such as:

- Store ID
- Date of Sale
- Weekly Sales
- Holiday Flag
- Temperature
- Fuel Price
- CPI (Consumer Price Index)
- Unemployment Rate

This dataset provides both numerical and categorical data, enabling analysis of sales trends, store performance, and external factors affecting sales.

## Methodology

### 1. Data Collection

- The Walmart Sales Dataset was imported from Kaggle.
- The dataset was loaded into Python (VS Code) using the Pandas library for analysis.
- Data integrity was verified by checking for duplicates and missing values.

### 2. Data Preprocessing

- Missing values were handled appropriately, and duplicate records were removed.
- Date columns were converted to a standard format for time-series analysis.
- Data types were standardized to ensure consistency for visualization and calculations.

### **3. Exploratory Data Analysis (EDA)**

- Conducted exploratory analysis using Pandas, Matplotlib, and Seaborn to identify:
  - Sales trends over time
  - Seasonal and regional variations
  - High-performing and low-performing stores and products
- Visualizations such as line charts, bar graphs, and heatmaps were created to interpret patterns.

### **4. Data Visualization**

- Total Sales by Store and Category
- Impact of holidays and temperature on weekly sales
- Correlations between economic indicators and sales performance

### **5. Insights and Recommendations**

- Key insights were summarized to help management understand:
  - Which stores and products drive the most revenue
  - Seasonal demand fluctuations
  - The impact of external factors like holidays and unemployment on sales

## **Features Implemented/ Key Visuals**

### **1. Data Cleaning and Preprocessing**

- Removed duplicate and null values to ensure data accuracy.
- Standardized date and numerical formats for consistent analysis.

### **2. Trend Analysis**

- Created time-series visualizations to show overall sales performance over months and years.
- Identified seasonal patterns and holiday impacts on weekly sales.

### **3. Product and Store Performance**

- Built charts to highlight top-performing stores and product categories.
- Compared sales across different regions and time periods using bar and line graphs.

## **4. Correlation Analysis**

- Used visual tools to examine the relationship between economic indicators (e.g., CPI, fuel price, unemployment) and sales trends.

## **6. Insights Presentation**

- Combined all findings into a comprehensive visual summary to support business decision-making in forecasting, inventory planning, and promotional strategies.

# **Results/ Key Insights**

## **1. Sales Performance Trends**

- Steady growth was observed during specific months, indicating seasonal demand spikes, particularly around holiday periods.
- Significant sales dips occurred during off-season months, suggesting potential opportunities for promotions during those times.

## **2. Store and Regional Insights**

- A few stores consistently contributed to the highest revenue, showing strong customer demand and operational efficiency.
- Some regions displayed lower sales figures, pointing to the need for localized marketing efforts or revised pricing strategies.

## **3. Product-Level Insights**

- Top-performing product categories accounted for a large share of total revenue.
- Low-selling products were identified, helping management focus on optimizing inventory and discontinuing underperforming items.

## **4. Impact of External Factors**

- External variables such as holidays, temperature changes, and unemployment rate had a noticeable impact on weekly sales.
- Holiday weeks showed a clear increase in sales volume, confirming the importance of seasonal campaigns.

## 5. Business Recommendations

- Focus marketing and stock efforts on high-demand months and holiday periods.
- Strengthen supply chain and logistics for top-performing stores.
- Use predictive analytics for better sales forecasting and demand planning in future operations.

## Challenges Faced

### 1. Data Quality Issues

**Challenge:** The dataset contained missing values, duplicate records, and inconsistent date formats, which affected data accuracy.

**Solution:** Used Python (Pandas) functions to clean the data duplicates were removed, missing values were handled using imputation techniques, and date formats were standardized for time-series analysis.

### 2. Large and Unstructured Data

**Challenge:** Managing and analyzing large amounts of raw sales data made it difficult to quickly generate insights.

**Solution:** Utilized data filtering and grouping techniques in Python to organize and summarize the dataset efficiently.

### 3. Identifying Seasonal Trends

**Challenge:** Detecting clear seasonal sales patterns was complex due to overlapping product categories and time periods.

**Solution:** Applied line charts and time-series visualizations in Matplotlib to clearly observe and interpret monthly and seasonal variations.

#### **4. Interpretation of Correlation Factors**

**Challenge:** Understanding the influence of external factors like fuel price and CPI on sales performance required careful statistical interpretation.

**Solution:** Used correlation heatmaps and pair plots to visually examine relationships and draw accurate conclusions.

### **Conclusion**

The Sales Data Analysis project successfully provided meaningful insights into the sales performance of Walmart stores using data-driven techniques. Through data cleaning, exploration, and visualization, the project identified key trends, high-performing products, and the impact of external factors such as holidays and economic conditions on weekly sales.

The findings revealed that a small number of stores and products contribute to the majority of total sales, highlighting areas of strong business performance. Seasonal and holiday periods were found to significantly boost revenue, confirming the importance of targeted marketing and inventory planning during these times. Overall, the project demonstrates how effective data analysis and visualization can play a crucial role in understanding business performance and supporting data-driven decision-making in the retail sector.

# Screenshots

```
walmart_analysis.py > ...
1  import os
2  import warnings
3  warnings.filterwarnings("ignore")
4
5  import matplotlib
6  import sys
7
8  # Try to force an interactive backend commonly available.
9  # If this fails, we'll still try to show plots and provide a helpful message.
10 for preferred in ("TkAgg", "Qt5Agg", "MacOSX"):
11     try:
12         matplotlib.use(preferred, force=True)
13         break
14     except Exception:
15         pass
16
17 import pandas as pd
18 import numpy as np
19 import matplotlib.pyplot as plt
20 import seaborn as sns
21 from sklearn.ensemble import RandomForestRegressor
22 from sklearn.metrics import mean_squared_error
23 from statsmodels.tsa.seasonal import seasonal_decompose
24 import joblib
25
26 DATA_PATH = "Walmart.csv" # put your file here
27
28
29 def load_data(path):
30     if not os.path.exists(path):
31         raise FileNotFoundError(f"{path} not found.")
32     sample = pd.read_csv(path, nrows=5)
33     date_cols = [c for c in sample.columns if "date" in c.lower() or "week" in c.lower()]
34     if date_cols:
35         df = pd.read_csv(path, parse_dates=[date_cols[0]])
36         df.rename(columns={date_cols[0]: "Date"}, inplace=True)
37     else:
38         df = pd.read_csv(path)
39         df.rename(columns={df.columns[0]: "Date"}, inplace=True)
40         df["Date"] = pd.to_datetime(df["Date"], errors="coerce")
41     return df
42
43
44 def clean_data(df):
45     df = df.copy().drop_duplicates()
46     sales_col = next((c for c in df.columns if "sale" in c.lower()), None)
47     if not sales_col:
48         raise ValueError("No column with 'sales' in the name was found.")
```

walmart\_analysis.py > explore\_and\_show

```
44 def clean_data(df):
48     raise ValueError("No column with 'sales' in the name was found.")
49     df.rename(columns={sales_col: "Weekly_Sales"}, inplace=True)
50     df["Weekly_Sales"] = pd.to_numeric(df["Weekly_Sales"], errors="coerce")
51     df["Weekly_Sales"].fillna(df["Weekly_Sales"].median(), inplace=True)
52
53     for col in list(df.columns):
54         if "store" in col.lower() and col != "Store":
55             df.rename(columns={col: "Store"}, inplace=True)
56         if "dept" in col.lower() and col != "Dept":
57             df.rename(columns={col: "Dept"}, inplace=True)
58
59     df.sort_values(by=["Date"], inplace=True)
60     df.reset_index(drop=True, inplace=True)
61     return df
62
63
64 def explore_and_show(df):
65     print("Matplotlib backend:", matplotlib.get_backend())
66     df["Date"] = pd.to_datetime(df["Date"], errors="coerce")
67     ts = df.set_index("Date").resample("W")["Weekly_Sales"].sum()
68
69     # 1) Total sales over time
70     fig, ax = plt.subplots(figsize=(12, 5))
71     ax.plot(ts.index, ts.values, linewidth=1.5)
72     ax.set_title("Total Weekly Sales Over Time")
73     ax.set_xlabel("Date")
74     ax.set_ylabel("Weekly Sales")
75     plt.show(block=True) # <--- blocking: window appears and stays until closed
76
77     # 2) Top 10 stores (if available)
78     if "Store" in df.columns:
79         top_stores = df.groupby("Store")["Weekly_Sales"].sum().nlargest(10)
80         fig, ax = plt.subplots(figsize=(10, 5))
81         sns.barplot(x=top_stores.index.astype(str), y=top_stores.values, ax=ax)
82         ax.set_title("Top 10 Stores by Total Sales")
83         ax.set_xlabel("Store")
84         ax.set_ylabel("Total Sales")
85         plt.setp(ax.get_xticklabels(), rotation=45)
86         plt.show(block=True)
87
88     # 3) Seasonal decomposition if enough points
89     if len(ts) >= 104:
90         try:
91             result = seasonal_decompose(ts, model="additive", period=52, extrapolate_trend='freq')
92             # result.plot() creates subplots; show them
93             result.plot()
94             plt.show(block=True)
```



walmart\_analysis.py > ...

```
64 def explore_and_show(df):
95     except Exception as e:
96         print("Seasonal decomposition error:", e)
97     else:
98         print("Not enough data for seasonal decomposition (need >=104 weekly points).")
99
100
101 def forecasting_and_show(df):
102     ts = df.set_index("Date").resample("W")["Weekly_Sales"].sum()
103     if len(ts) < 10:
104         print("Too few data points for forecasting.")
105         return
106
107     def create_lags(series, lags=[1, 2, 3, 4, 52]):
108         df_lag = pd.DataFrame({'y': series})
109         for lag in lags:
110             df_lag[f'lag_{lag}'] = df_lag['y'].shift(lag)
111             df_lag["rolling_4"] = df_lag['y'].shift(1).rolling(4).mean()
112         return df_lag.dropna()
113
114     df_feat = create_lags(ts)
115     X = df_feat.drop("y", axis=1)
116     y = df_feat["y"]
117     split = int(len(X) * 0.8)
118     X_train, X_test = X.iloc[:split], X.iloc[split:]
119     y_train, y_test = y.iloc[:split], y.iloc[split:]
120
121     model = RandomForestRegressor(n_estimators=100, random_state=42)
122     model.fit(X_train, y_train)
123     preds = model.predict(X_test)
124     preds_series = pd.Series(preds, index=y_test.index)
125
126     try:
127         rmse = mean_squared_error(y_test, preds_series, squared=False)
128     except TypeError:
129         rmse = np.sqrt(mean_squared_error(y_test, preds_series))
130     print(f"Forecast RMSE: {rmse:.2f}")
131
```

```

132     # Plot actual vs predicted (show only)
133     fig, ax = plt.subplots(figsize=(10, 5))
134     ax.plot(y_test.index, y_test.values, label="Actual", linewidth=2)
135     ax.plot(preds_series.index, preds_series.values, label="Predicted", linewidth=2)
136     ax.set_title("Actual vs Predicted Sales")
137     ax.set_xlabel("Date")
138     ax.set_ylabel("Weekly Sales")
139     ax.legend()
140     ax.grid(True, linestyle="--", alpha=0.6)
141     plt.show(block=True)
142
143
144 def main():
145     print("Running show-only script. Make sure you run this in a desktop environment.")
146     df = load_data(DATA_PATH)
147     df = clean_data(df)
148     explore_and_show(df)
149     forecasting_and_show(df)
150     print("Finished. All plots were displayed (no files saved).")
151
152
153 if __name__ == "__main__":
154     main()
155

```

Figure 1

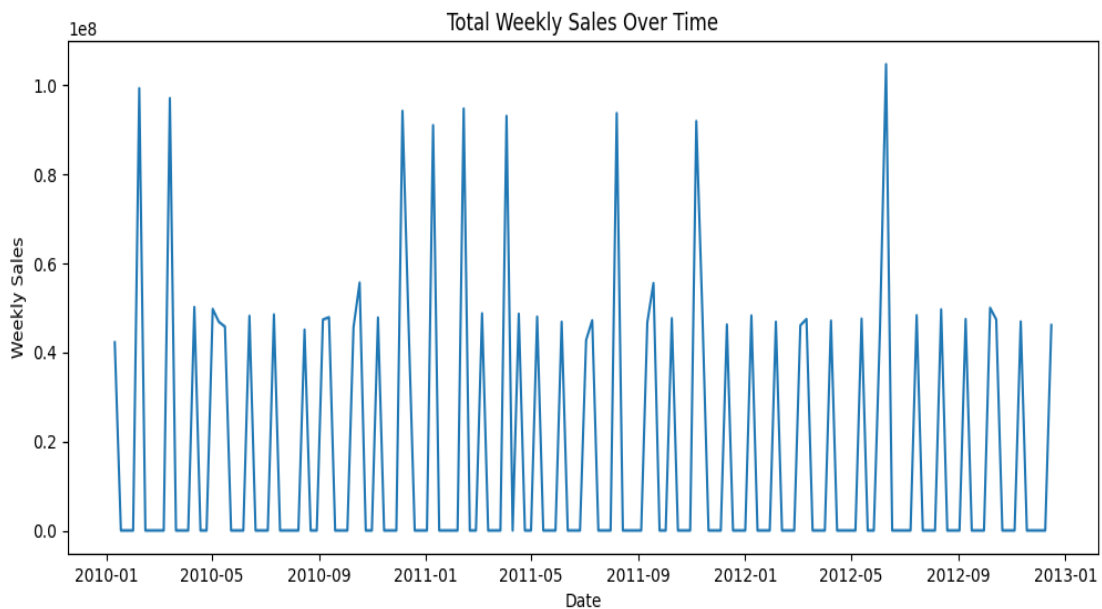
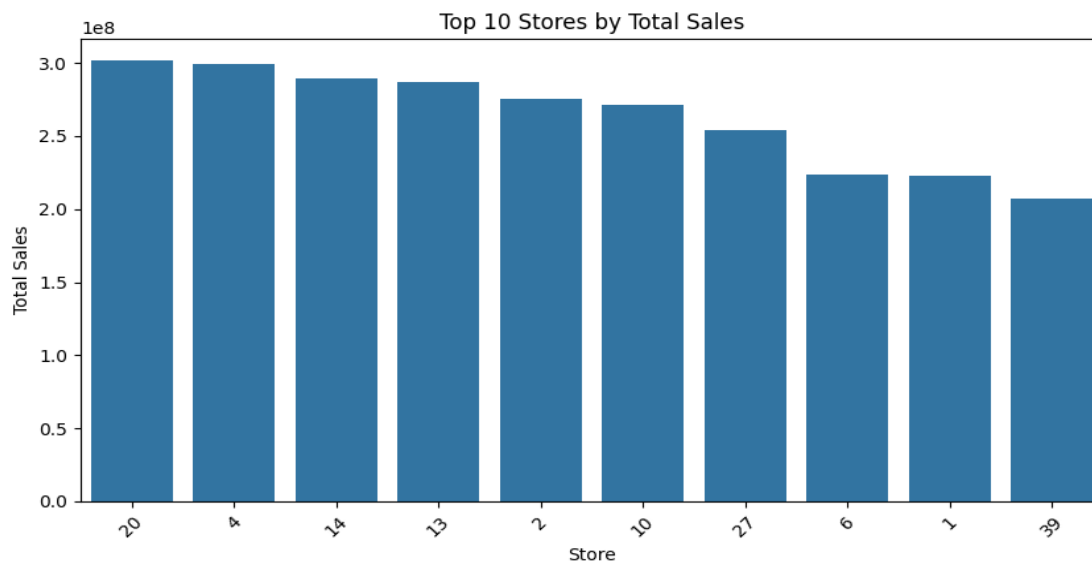


Figure 1

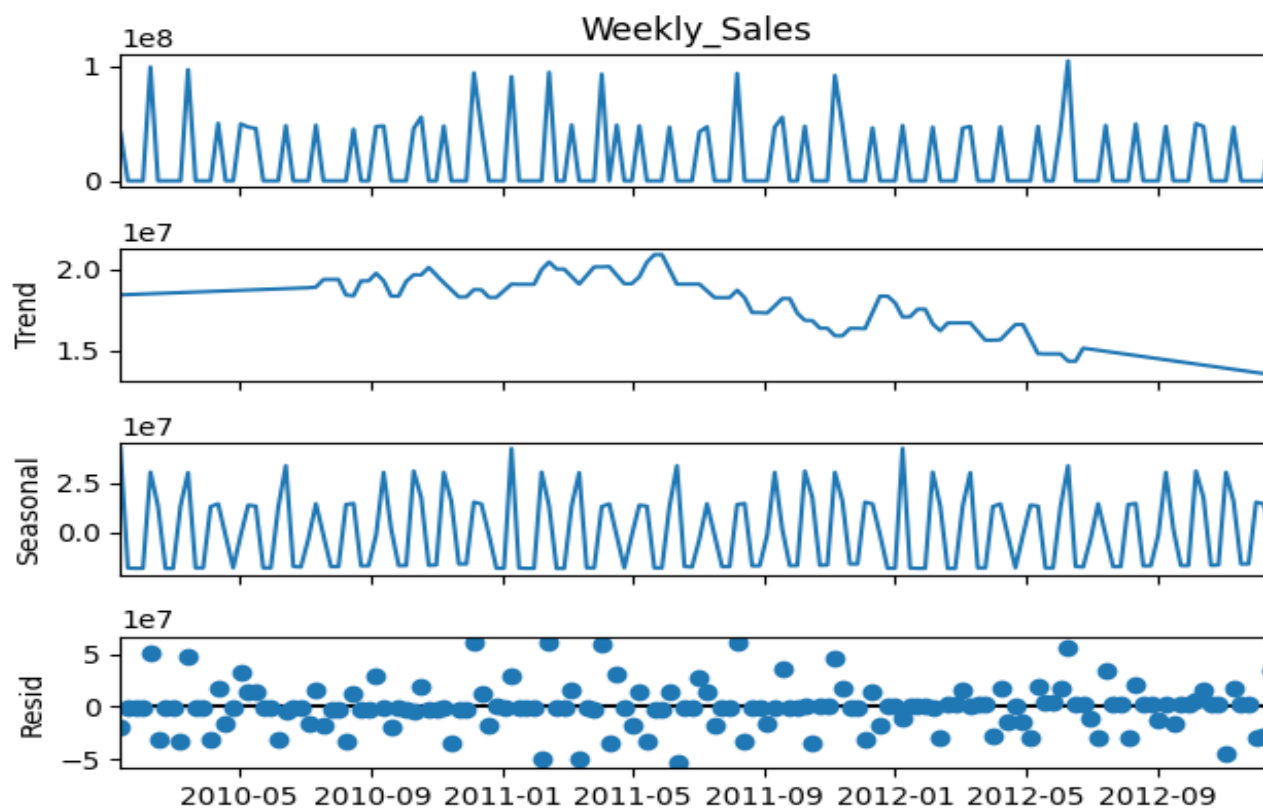
— □ ×



Home navigation icons: back, forward, search, zoom, and save.

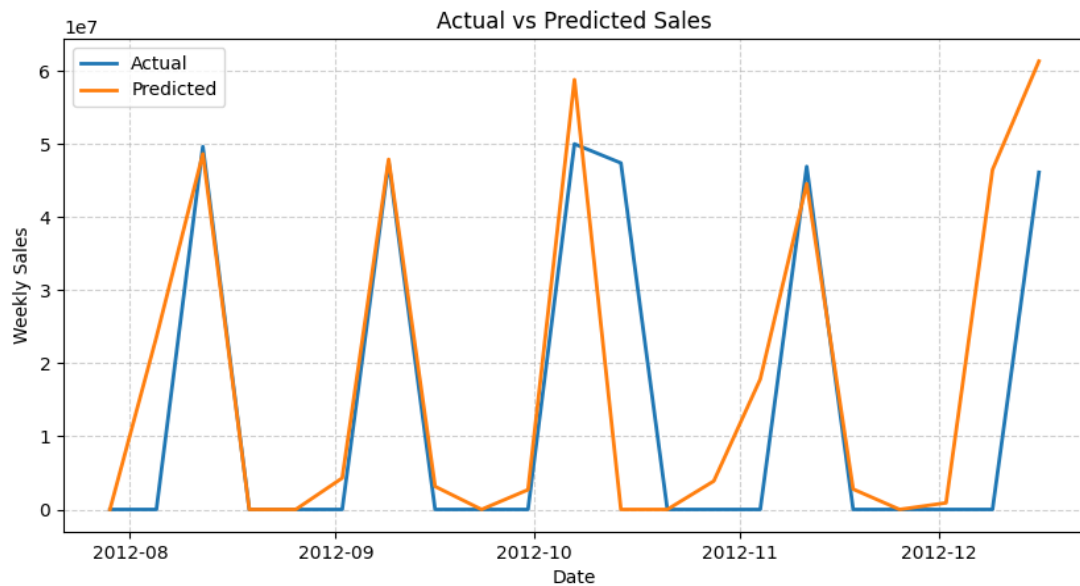
Figure 1

— □ ×



Home navigation icons: back, forward, search, zoom, and save.

Figure 1



(x, y) = (2012-11, 2.50e+07)

## Link to GitHub

GitHub: <https://github.com/Yashu-teach/Sales-Data-Analysis-Project>