# Technical Case Study

Will Sparrow

# Overview

- I recently conducted an exploratory and predictive analysis on a multi-location retail dataset to identify patterns in weekly sales, evaluate store segmentation, and explore forecasting approaches influenced by seasonal/economic factors.

- This project was completed as part of a technical case study.

- Dataset structure and task framing have been modified and anonymized.

# Table of Contents

- Summary
- A Look at the Data
- Weekly Sales Insights & Predictive Approach
- Store Classification & Efficiency Analysis
- Improving Models with External Inputs
- Python Script

Adapted from a technical case prompt. Names, values, and variables have been modified.

# The Data – Weekly Sales

- ~40 Total Stores
  - ~60 departments per store
- Time Range: 2010 - 2013
- Interval: Weekly
- H_Period shows which weeks are Holidays

| Location_ID | Dept | Date | Weekly_Sales | H_Period |
|---|---|---|---|---|
| 1 | 1 | 2/5/10 | 24924.5 | FALSE |
| 1 | 1 | 2/12/10 | 46039.49 | TRUE |
| 1 | 1 | 2/19/10 | 41595.55 | FALSE |
| 1 | 1 | 2/26/10 | 19403.54 | FALSE |
| 1 | 1 | 3/5/10 | 21827.9 | FALSE |
| 1 | 1 | 3/12/10 | 21043.39 | FALSE |
| 1 | 1 | 3/19/10 | 22136.64 | FALSE |
| 1 | 1 | 3/26/10 | 26229.21 | FALSE |
| 1 | 1 | 4/2/10 | 57258.43 | FALSE |

Weekly Sales

Adapted from a technical case prompt. Names, values, and variables have been modified.

# The Data – Factors & Store Type

- Economic and Regional Factors
  - For each week, displays the reported measures of four economic and/or regional variables: **temperature**, **gas price**, the **Consumer Price Index (CPI)**, and **unemployment** rate.

| Location_ID | Date | Avg_Reg_Temp | Gas_Price | CPI_Index | Unemp_Pct |
|---|---|---|---|---|---|
| 1 | 2/5/10 | 42.31 | 2.57 | 211.10 | 8.11 |
| 1 | 2/12/10 | 38.51 | 2.55 | 211.24 | 8.11 |
| 1 | 2/19/10 | 39.93 | 2.51 | 211.29 | 8.11 |
| 1 | 2/26/10 | 46.63 | 2.56 | 211.32 | 8.11 |
| 1 | 3/5/10 | 46.5 | 2.63 | 211.35 | 8.11 |
| 1 | 3/12/10 | 57.79 | 2.67 | 211.38 | 8.11 |
| 1 | 3/19/10 | 54.58 | 2.72 | 211.22 | 8.11 |
| 1 | 3/26/10 | 51.45 | 2.73 | 211.02 | 8.11 |
| 1 | 4/2/10 | 62.27 | 2.72 | 210.82 | 7.81 |

Economic and Regional Factors

- Store Type
  - Displays the type and size of store

| Location_ID | Type | Size |
|---|---|---|
| 1 | A | 151315 |
| 2 | A | 202307 |
| 3 | B | 37392 |
| 4 | A | 205863 |
| 5 | B | 34875 |

Store Type

Adapted from a technical case prompt. Names, values, and variables have been modified.
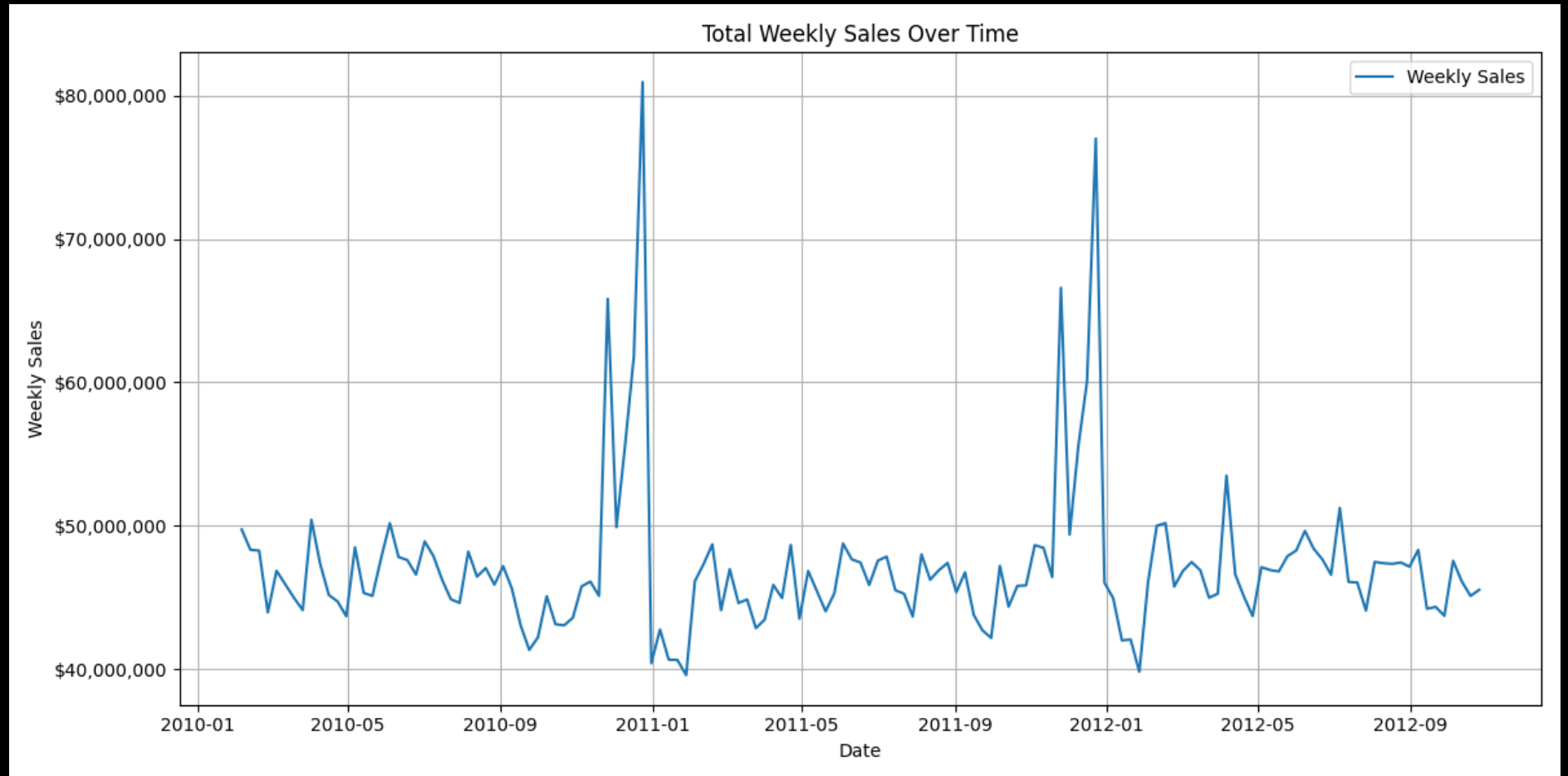
# Weekly Sales Insights & Predictive Approach

- On the next few slides, I investigate sales trends through exploratory analysis, focusing on whether holiday periods and regional or economic factors meaningfully influence performance

- I also outline a potential forecasting model that could be used to predict future sales, including the general approach and evaluation considerations.
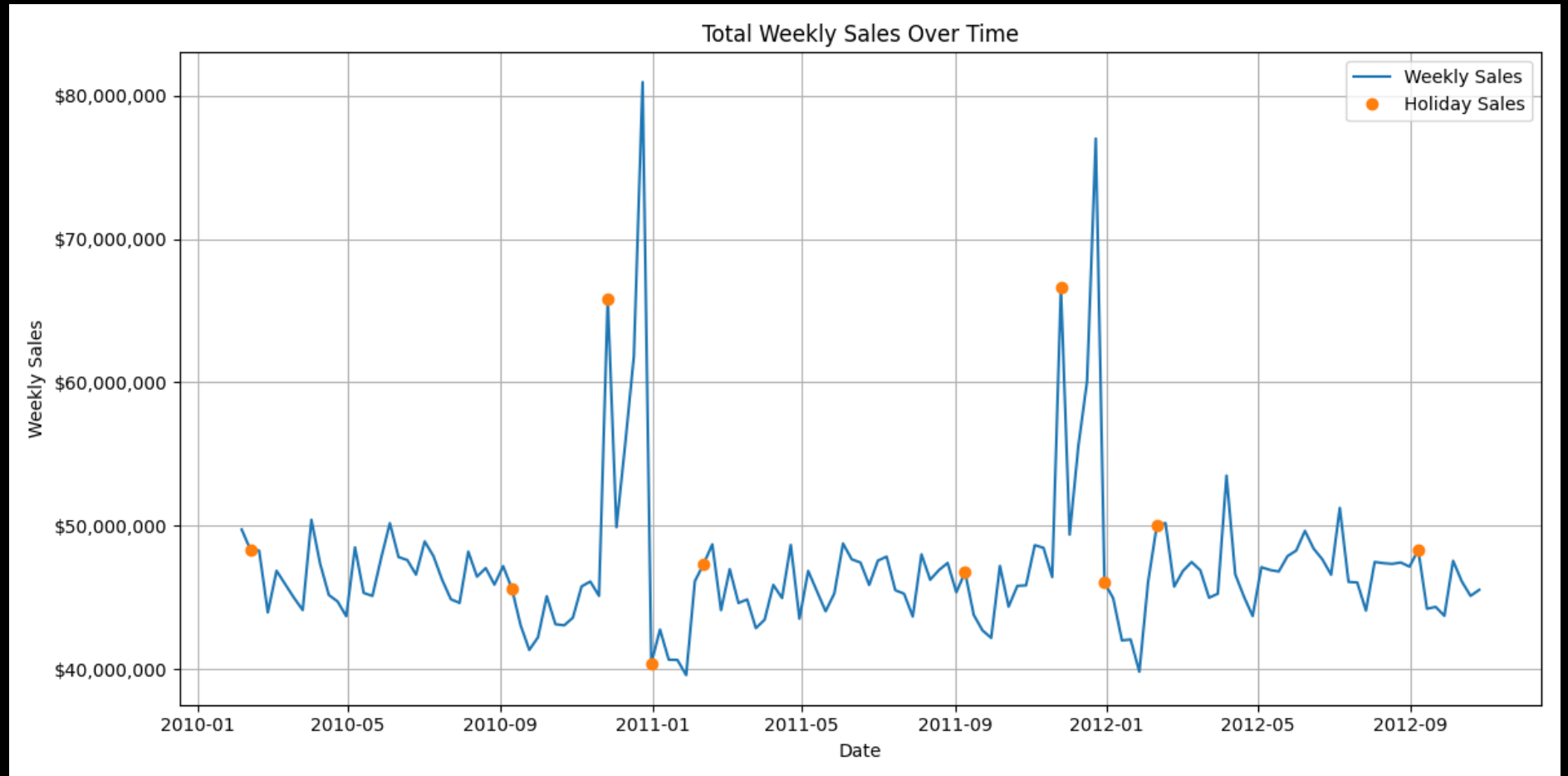
# Weekly Sales Insights & Predictive Approach

- Weekly sales over time across ALL stores/deps



Adapted from a technical case prompt. Names, values, and variables have been modified.

# Weekly Sales Insights & Predictive Approach

- Holidays highlighted by orange dots

- These points are some of the highest weekly sales numbers

- Often a sharp dip following holiday



Total Weekly Sales Over Time

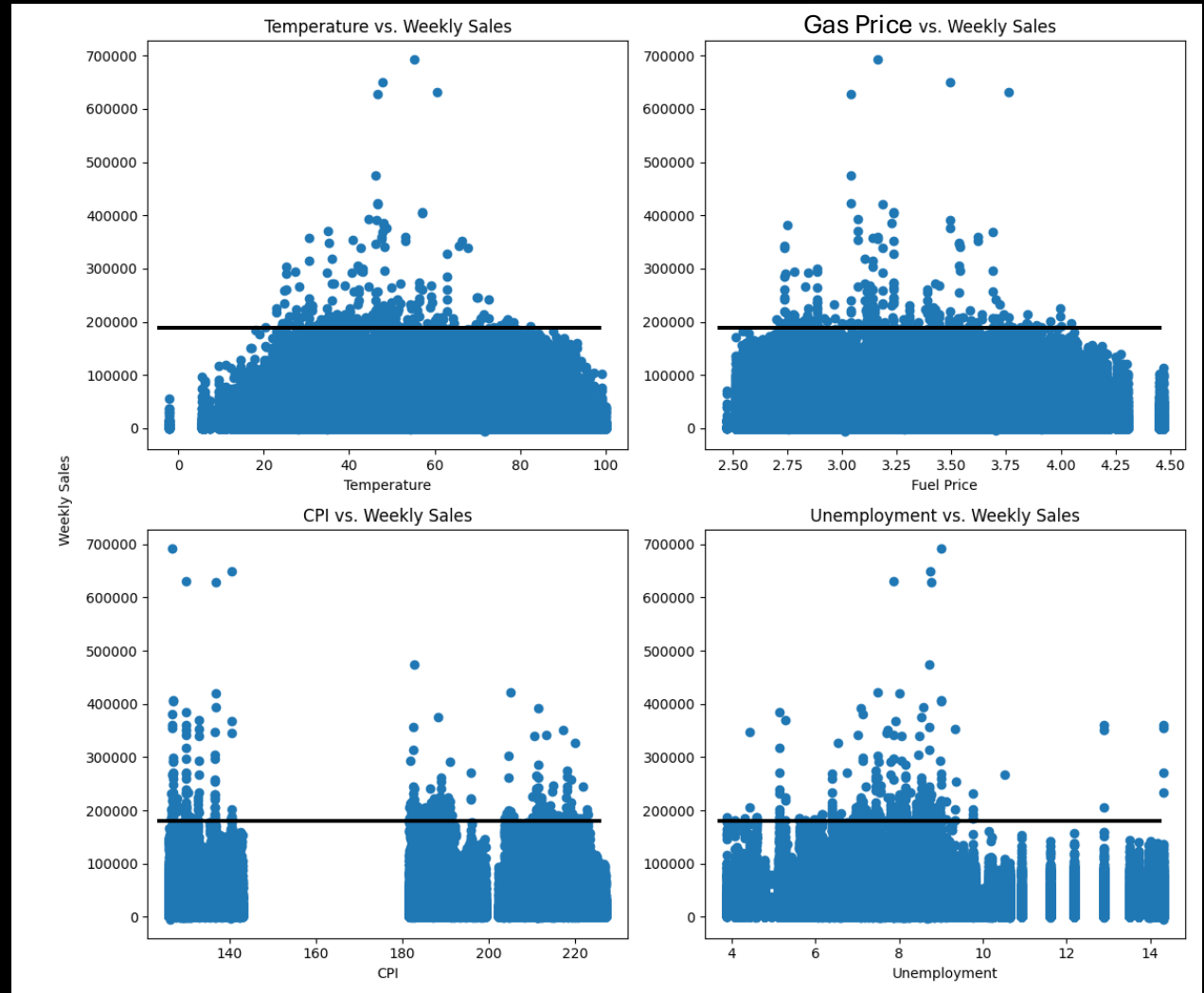# Weekly Sales Insights & Predictive Approach

- Visually, it appears that holidays have an impact, but let's run a quick test as well.

- After comparing the average weekly sales for non-holidays versus holidays, we get:

```
H_Period
False        15901.445069
True         17035.823187
```

During holiday weeks, sales are ~7% higher

# Weekly Sales Insights & Predictive Approach

- All four regional/economic factors plotted against weekly sales shows no correlation

- Temperature, gas prices, CPI, and unemployment rate therefore DO NOT affect sales across all stores

# Weekly Sales Insights & Predictive Approach

The below model is a way in which we can predict future sales to account for seasonality of data:

- Methodology: ARIMA model
  - We can use an ARIMA model which takes into account seasonality to forecast weekly sales. This model looks at past data and can use patterns such as the sales output during holiday weeks, to help predict future sales.

- Validation Techniques: Test-Train-Split
  - Because the data is time based, we must sort the data chronologically by the date field. Then take the first 80% of data to train the model, tune the various parameters of the model, and use the last 20% to test the model.

- Metrics to Evaluate Model: Mean Absolute Error
  - After using the test set to generate predictions, take the absolute value between each prediction and the actual weekly sales value  and divide that by the number of predictions made. This will give you on average, how far each prediction is away from its actual sales value.

Adapted from a technical case prompt. Names, values, and variables have been modified.
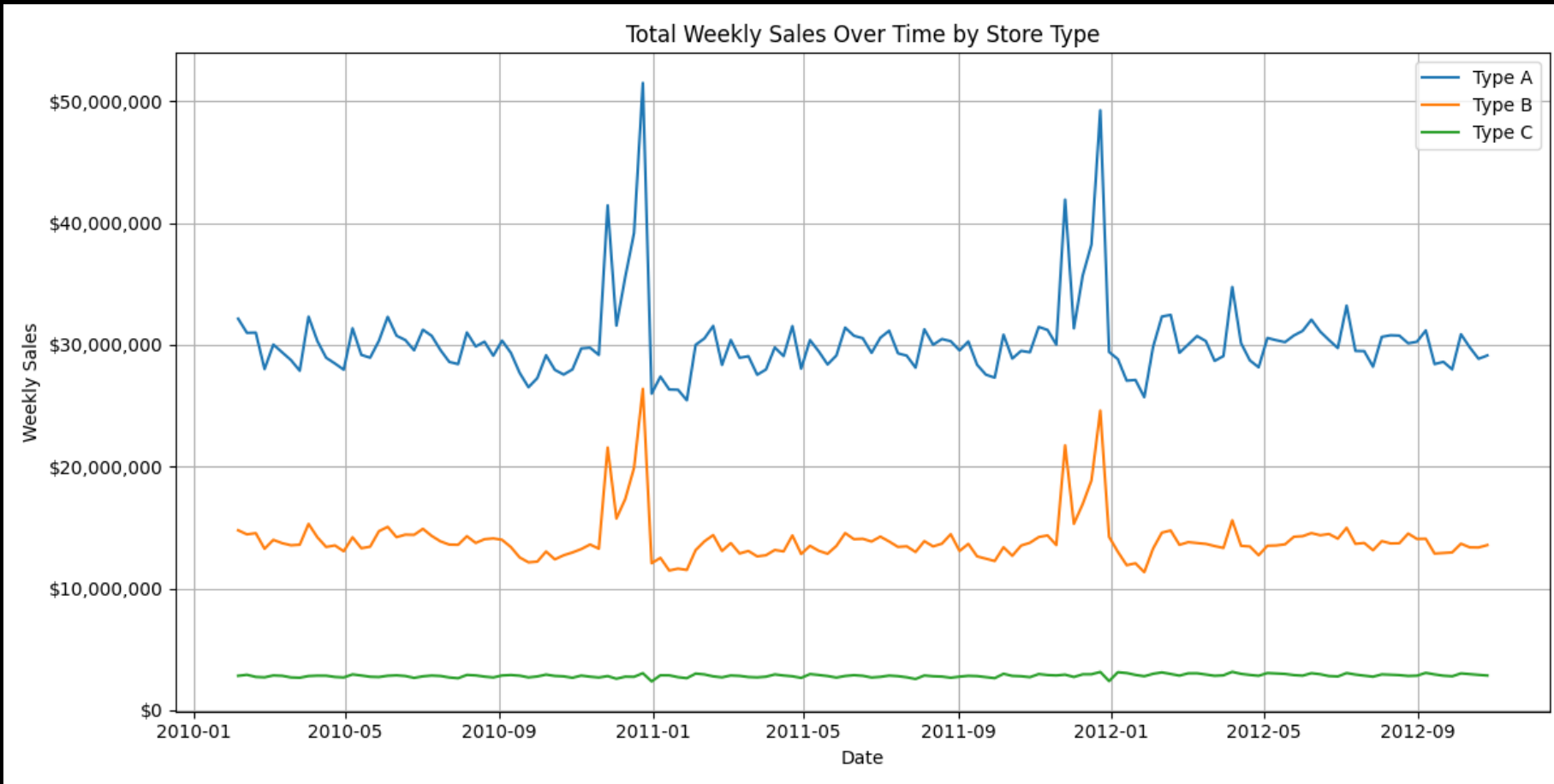
# Store Classification & Efficiency Analysis

- On the following slides, I segment the stores by size, classification, and regional characteristics, then analyze sales performance across each group to identify the factors most associated with differences in outcomes.

# Store Classification & Efficiency Analysis

- We will be looking at four different variables:
  - Store Size
    - For simplicity's sake, I have assumed Store Type and Store Size are the same variable
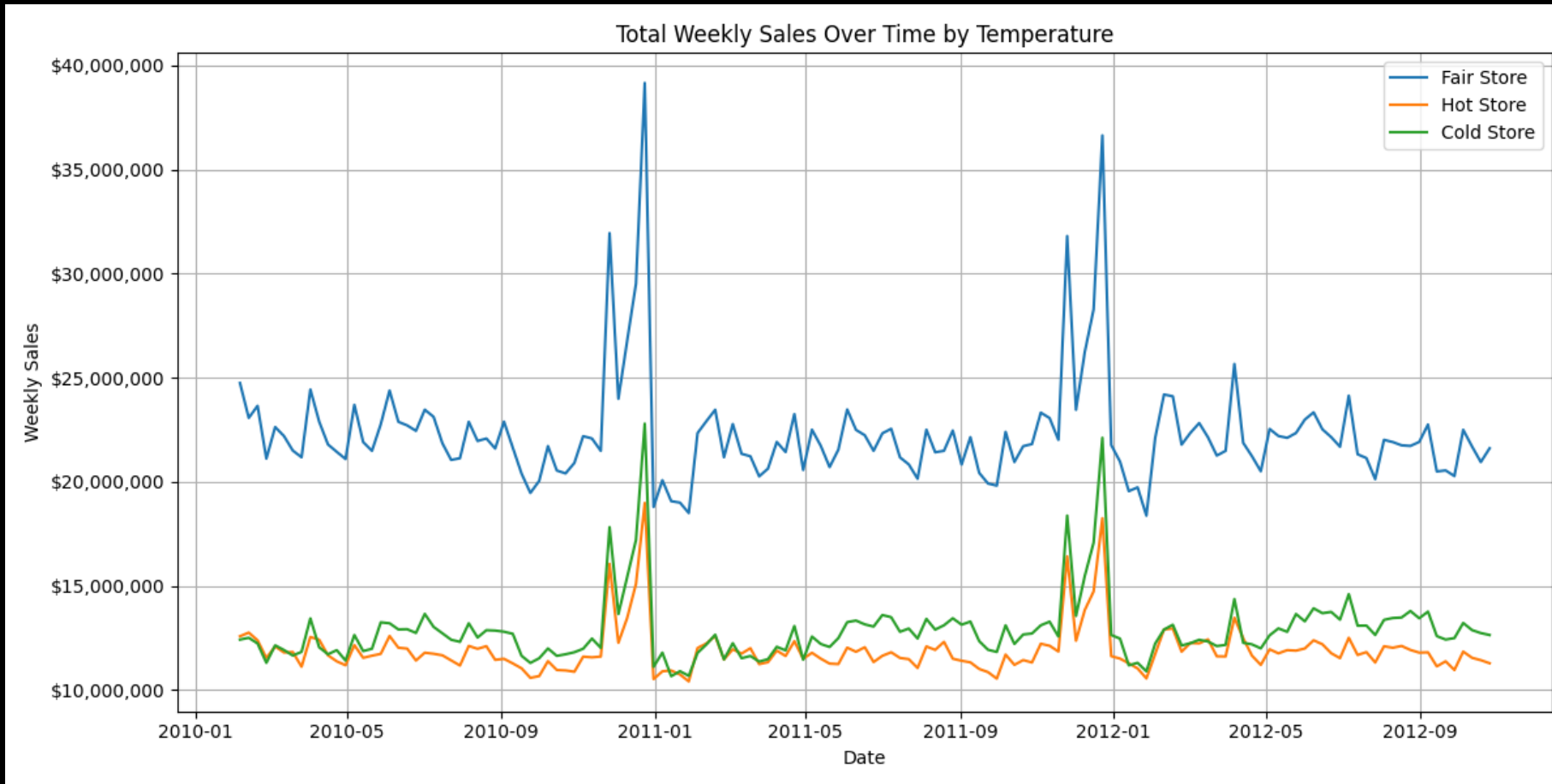  - Temperature
  - Gas Prices
  - CPI

# Store Classification & Efficiency Analysis



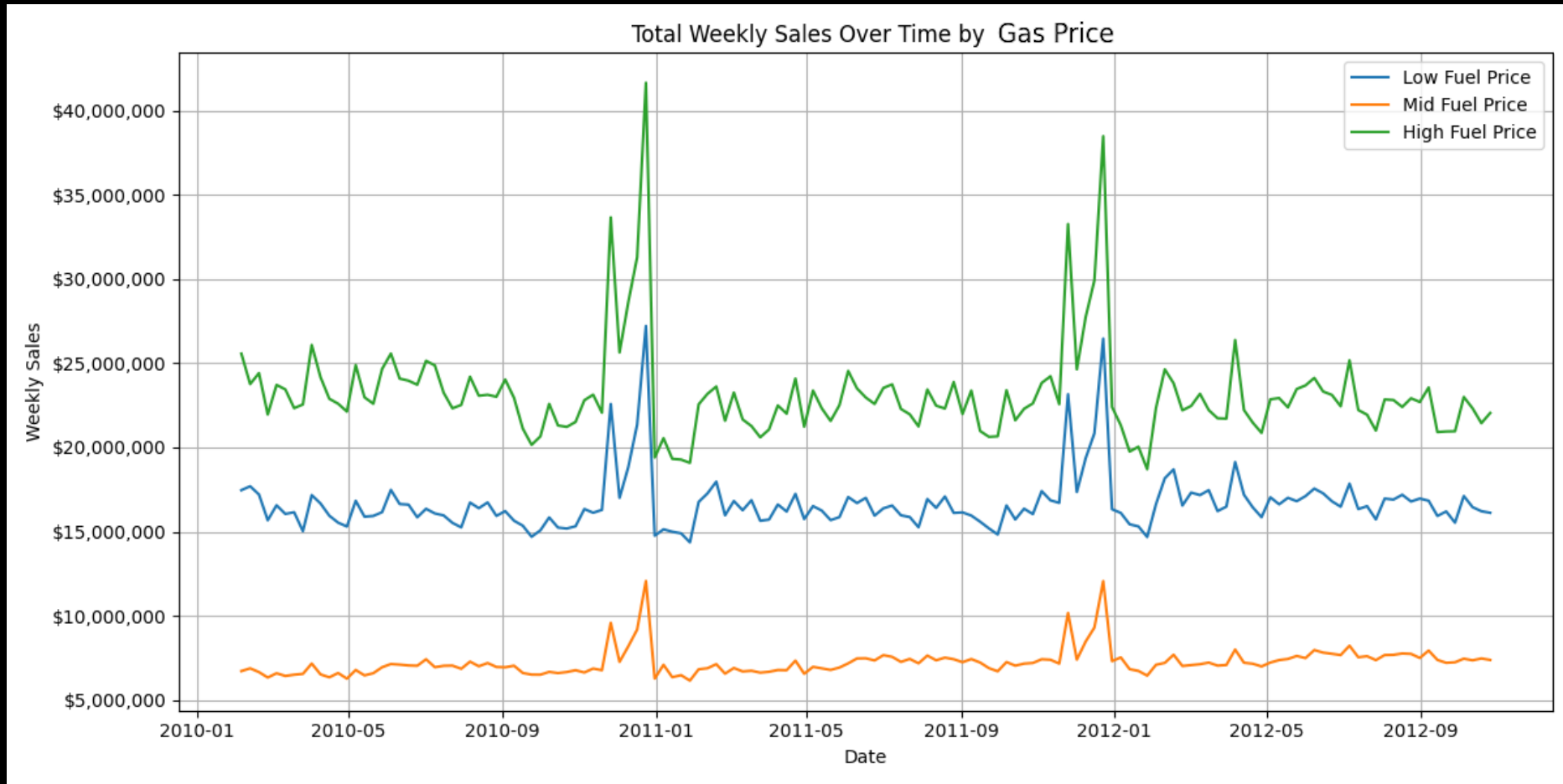Total Weekly Sales Over Time by Store Type

- Size
  - Type A stores are greater than 150,000
  - Type B are stores between 50,000 and 150,000
  - Type C are stores less than 50,000
  - Larger the store the better the sales

Adapted from a technical case prompt. Names, values, and variables have been modified.

# Store Classification & Efficiency Analysis



Total Weekly Sales Over Time by Temperature

- Temperature
  - Stores in milder climates perform the best
  - Both hot/cold stores perform about the same

# Store Classification & Efficiency Analysis
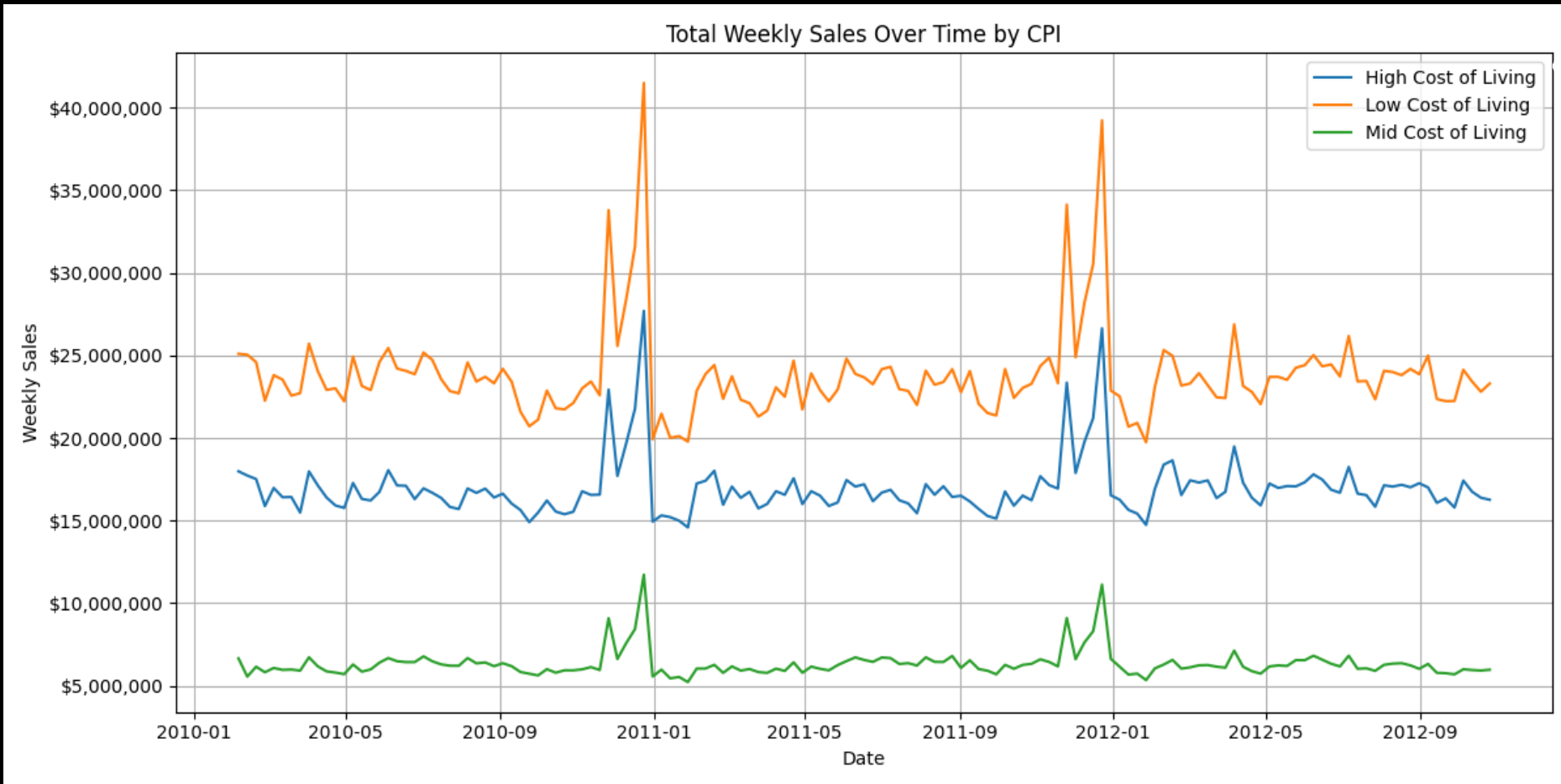


Total Weekly Sales Over Time by Gas Price

## Gas Price
- Stores in regions with high gas prices perform the best
- Stores in. regions with low gas prices perform second best
- Stores in regions with mid gas prices perform last

# Store Classification & Efficiency Analysis



Total Weekly Sales Over Time by CPI

CPI
- Stores in LCOL regions perform best
- Stores in HCOL regions perform second best

Adapted from a technical case prompt. Names, values, and variables have been modified.

# Store Classification & Efficiency Analysis

| Variable | Effect | Potential Explanation |
|---|---|---|
| Store Size | The bigger the store size, the more weekly sales. | Assuming this is measured in square footage, the bigger the store is as it relates to area, the more inventory you can hold, and shoppers can ultimately buy. |
| Temperature | Stores in "fair" regions (milder temps) do much better than both hot/cold stores. | Shoppers are able to do more shopping year round if the weather is fair. They are less likely to go out and purchase things in extremely hot/cold weather. |
| Fuel Prices | Stores in area with high gas prices perform the best, then stores with low gas prices, and finally stores with mid gas prices. | Shoppers in an area with higher gas prices may have more absolute income (e.g., California residents in major city). |
| CPI | Stores in lower cost of living areas perform best, then stores in higher cost of living, and finally stores in mid cost of living. | The types of good this store sells, could be less luxury based, and cater more to lower income families. |

Adapted from a technical case prompt. Names, values, and variables have been modified.

# Improving Models with External Inputs

- On the following slide, I explore how incorporating external data sources could potentially strengthen the predictive models by adding additional context and improving forecast accuracy. I also outline the practical challenges and limitations involved in integrating these data sources with the existing dataset.

# Improving Models with External Inputs

| Additional Data Sources | How Can Data Add Context or Improve Forecast/Operational Efficiency? | Potential Limitations of Additional Data Sources |
|---|---|---|
| Demographic Data | Allows us to understand the characteristics (income, age, etc.) of the surrounding population around our stores. Households and/or individuals with higher income may spend more. | Census data is not always accurate, up to date, and certain measures you might be interested in may not be reported. |
| Marketing Data | Helps us understand the impact of marketing campaigns on store sales. A recently released commercial or email campaign may result in a spike in weekly sales. | Difficult to quantify if marketing campaigns had direct effect on consumers. Additionally, dollars spent on marketing might not yield dollars spent on goods until much later. |
| Social Media Data | Helps us understand what the current sentiment of the company is and how it is trending. If the sentiment is negative, then you might expect weekly sales to decline (e.g., boycott of companies). | Sentiment analysis can be very noisy and unreliable. Additionally, general negativity (or positivity) for a company might not be able to be measured regionally. |

Adapted from a technical case prompt. Names, values, and variables have been modified.

# Python Script

- The following slides will show Python code I used to conduct the analyses.

```python
import pandas as pd
import matplotlib.pyplot as plt
import matplotlib.ticker as ticker

# Read in data
holiday_df = pd.read_excel("Data Science - Holiday.xlsx")
macro_factors_df = pd.read_excel("Data Science - Macro Factors.xlsx")
store_sales_df = pd.read_excel("Data Science - Store Sales.xlsx")
store_type_df = pd.read_excel("Data Science - Store Type.xlsx")

# Combine store_sales_df with macro_factors_df
combined_df = pd.merge(
    store_sales_df,
    macro_factors_df.drop("H_Period", axis=1),  # keep H_Period col from store_sales_df
    on=["Location_ID", "Date"],
    how="left"
)

# Add store type and size to combined df
combined_df = pd.merge(
    combined_df,
    store_type_df,
    on="Location_ID",
    how="left"
)

# Reorder cols so that it reads better for me and save to csv so I can open file and explore
combined_df = combined_df[["Location_ID", "Dept", "Type","Size", "Date",
                           "Weekly_Sales", "H_Period", "Avg_Reg_Temp", "Gas_Price", "CPI_Index", "Unemp_Pct"]]
combined_df.to_csv("Final Combined Sales.csv", index=None)
```

Adapted from a technical case prompt. Names, values, and variables have been modified.

```
32   ## Conduct exploratory data analysis to understand sales trends and the impact of holidays and regional economic factors.
33   # Filter for only holiday rows
34   holidays = combined_df.loc[combined_df["H_Period"] == True]
35
36   # Drop duplicates to keep only unique dates
37   unique_holidays = holidays[["Date", "H_Period"]].drop_duplicates()  # Assumes that all locations/depts observe same holiday schedule
38
39   # Convert to dictionary
40   holiday_dict = dict(zip(unique_holidays["Date"], unique_holidays["H_Period"]))
41
42   # Group by and plot weekly sales across all locations/departments
43   grouped_by_week_sales = combined_df.groupby("Date")["Weekly_Sales"].sum()
44
45   # Add holiday sales which just returns the weekly_sales number if that row's date is a holiday week
46   weekly_sales_df = grouped_by_week_sales.copy()
47   weekly_sales_df = weekly_sales_df.to_frame(name="Weekly_Sales").reset_index()
48   weekly_sales_df["Holiday Sales"] = weekly_sales_df["Weekly_Sales"][
49       weekly_sales_df["Date"].apply(lambda d: d in holiday_dict)
50   ]
51
52   # Plot weekly sales over time across all locations/depts and shows which weeks are holidays
53   fig1, ax1 = plt.subplots(figsize=(12, 6))
54   ax1.plot(weekly_sales_df["Date"], weekly_sales_df["Weekly_Sales"], label="Weekly Sales")
55   ax1.plot(weekly_sales_df["Date"], weekly_sales_df["Holiday Sales"], "o", label="Holiday Sales")
56   ax1.set_title("Total Weekly Sales Over Time")
57   ax1.set_xlabel("Date")
58   ax1.set_ylabel("Weekly Sales")
59   # Format Y-axis as dollars
60   ax1.yaxis.set_major_formatter(ticker.FuncFormatter(lambda x, pos: f"${x:,.0f}"))
61   ax1.grid(True)
62   ax1.legend()
63
64   fig1.tight_layout()
65   fig1.savefig("weekly_sales_and_holidays.png")  # Save the figure
66   # plt.show()
67   plt.close(fig1)
```

Adapted from a technical case prompt. Names, values, and variables have been modified.

```python
71    # Avg holiday versus non-holiday weekly sales
72    holiday_v_nonholiday_sales = combined_df.groupby("H_Period")["Weekly_Sales"].mean()
73    # print(holiday_v_nonholiday_sales) #~7% higher sales numbers when its holiday week
74
75    # Create 2x2 scatter plot grid
76    fig2, axes = plt.subplots(2, 2, figsize=(12, 10))  # 2 rows, 2 columns
77
78    axes[0, 0].scatter(combined_df["Avg_Reg_Temp"], combined_df["Weekly_Sales"])
79    axes[0, 0].set_title("Avg Reg Temp vs. Weekly Sales")
80    axes[0, 0].set_xlabel("Avg Reg Temp")
81
82    axes[0, 1].scatter(combined_df["Gas_Price"], combined_df["Weekly_Sales"])
83    axes[0, 1].set_title("Gas Price vs. Weekly Sales")
84    axes[0, 1].set_xlabel("Gas Price")
85
86    axes[1, 0].scatter(combined_df["CPI_Index"], combined_df["Weekly_Sales"])
87    axes[1, 0].set_title("CPI Index vs. Weekly Sales")
88    axes[1, 0].set_xlabel("CPI Index")
89
90    axes[1, 1].scatter(combined_df["Unemp_Pct"], combined_df["Weekly_Sales"])
91    axes[1, 1].set_title("Unemp Pct vs. Weekly Sales")
92    axes[1, 1].set_xlabel("Unemp Pct")
93
94    # Shared y-axis label
95    fig2.text(0.04, 0.5, "Weekly Sales", va="center", rotation="vertical")
96
97    fig2.tight_layout(rect=[0.05, 0, 1, 1])
98    fig2.savefig("economic_factors_vs_weekly_sales.png")  # Save the figure
99    # plt.show()
100   plt.close(fig2)
```

Adapted from a technical case prompt. Names, values, and variables have been modified.

```python
102    ## Segment the stores based on type, size, and regional characteristics.
103    # Multi line plot
104    def multi_line_plot(df, segment):
105        fig, ax = plt.subplots(figsize=(12, 6))
106        unique_values = df[segment].unique()
107
108        for val in unique_values:
109            sub_df = df[df[segment] == val]
110            sales_by_date = sub_df.groupby("Date")["Weekly_Sales"].sum().sort_index()
111            ax.plot(sales_by_date.index, sales_by_date.values, label=f"{val}")
112
113        if segment.endswith("_Classification_"):
114            segment = segment.replace("_Classification_", "")
115
116        ax.set_title(f"Total Weekly Sales Over Time by {segment}")
117        ax.set_xlabel("Date")
118        ax.set_ylabel("Weekly Sales")
119        ax.legend()
120        ax.grid(True)
121        ax.yaxis.set_major_formatter(ticker.FuncFormatter(lambda x, pos: f"${x:,.0f}"))
122
123        fig.tight_layout()
124        fig.savefig(f"weekly_sales_by_{segment}.png")
```

Adapted from a technical case prompt. Names, values, and variables have been modified.

```python
127    # Split by regional factors (excluding unemployment because it is national/macro factor)
128    regional_factors_agg_df = combined_df.groupby("Location_ID")[["Avg_Reg_Temp", "Gas_Price", "CPI_Index"]].agg(
129        ["min", "median", "max"]
130    )
131    # Add a column to classify locations based on their average temp
132    def classify_temperature(temp):
133        if temp < 55:
134            return "Cold Store"
135        elif 55 <= temp <= 70:
136            return "Fair Store"
137        else:
138            return "Hot Store"
139
140
141    # Add a column to classify locations based on their average CPI
142    def classify_cpi(cpi):
143        if cpi < 140:
144            return "Low Cost of Living"
145        elif 140 <= cpi <= 200:
146            return "Mid Cost of Living"
147        else:
148            return "High Cost of Living"
149
150
151    # Add a column to classify locations based on their average gas price
152    def classify_fuel_price(fuel_price):
153        if fuel_price < 3.3:
154            return "Low Fuel Price"
155        elif 3.3 <= fuel_price <= 3.5:
156            return "Mid Fuel Price"
157        else:
158            return "High Fuel Price"
159
160
161    # Apply classification
162    regional_factors_agg_df["Avg_Reg_Temp_Classification"] = regional_factors_agg_df[("Avg_Reg_Temp", "median")].apply(
163        classify_temperature
164    )
165    regional_factors_agg_df["CPI_Index_Classification"] = regional_factors_agg_df[("CPI_Index", "median")].apply(
166        classify_cpi
167    )
168    regional_factors_agg_df["Gas_Price_Classification"] = regional_factors_agg_df[("Gas_Price", "median")].apply(
169        classify_fuel_price
170    )
```

Adapted from a technical case prompt. Names, values, and variables have been modified.

```
172    # Count unique entries for each classification column
173    temp_counts = regional_factors_agg_df["Avg_Reg_Temp_Classification"].value_counts()
174    cpi_counts = regional_factors_agg_df["CPI_Index_Classification"].value_counts()
175    fuel_counts = regional_factors_agg_df["Gas_Price_Classification"].value_counts()
176
177    print("Avg Reg Temp Classification Counts:\n", temp_counts)
178    print("\nCPI Index Classification Counts:\n", cpi_counts)
179    print("\nGas Price Classification Counts:\n", fuel_counts)
180
181    # Flatten the data
182    regional_factors_agg_df.columns = ["_".join(col).strip() if isinstance(col, tuple) else col for col in regional_factors_agg_df.columns]
183
184    # Now the DataFrame has singular column names
185    regional_factors_agg_df = regional_factors_agg_df.reset_index()
186
187    # Select relevant cols and merge
188    combined_df = pd.merge(
189        combined_df,
190        regional_factors_agg_df[
191            [
192                "Location_ID",
193                "Avg_Reg_Temp_Classification_",
194                "CPI_Index_Classification_",
195                "Gas_Price_Classification_",
196            ]
197        ],
198        on="Location_ID",
199        how="left",
200    )
201
202    multi_line_plot(combined_df, "Avg_Reg_Temp_Classification_")
203    multi_line_plot(combined_df, "CPI_Index_Classification_")
204    multi_line_plot(combined_df, "Gas_Price_Classification_")
205    multi_line_plot(combined_df, "Type")
```

Adapted from a technical case prompt. Names, values, and variables have been modified.