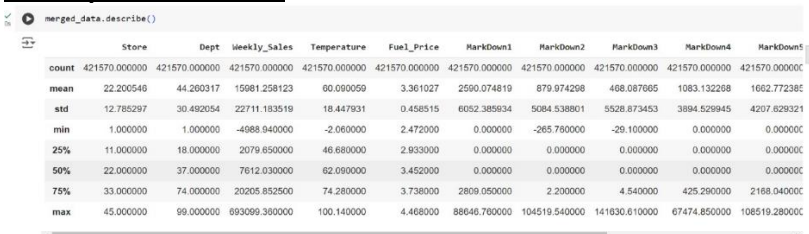
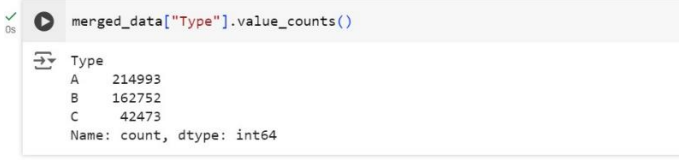


Data Collection and Preprocessing Phase

Date	4 June,2024
Team ID	SWTID1719938571
Project Title	Walmart Sales Analysis for Retail Industry with Machine Learning
Maximum Marks	6 Marks

Data Exploration and Preprocessing

In this phase, we will statistically analyze dataset variables to identify patterns and outliers. Python will be used for preprocessing tasks, including normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring the quality of data for subsequent analysis and modeling. This process will form a strong foundation for generating reliable insights and accurate predictions.

Section	Description
Data Overview	<p><u>Dimensions:</u> 421570 rows × 17 columns</p> <p><u>Descriptive Statistics:</u></p> <pre>merged_data.describe()</pre> 
Univariate Analysis	<p>Univariate Analysis</p> <pre>merged_data["Type"].value_counts()</pre> 

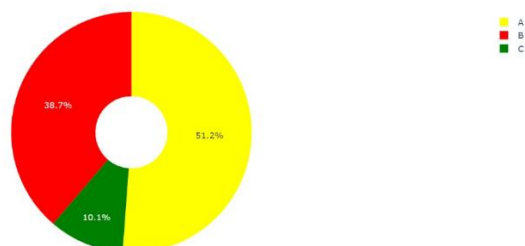
```
[29] print("Maximum Weekly Sales:", merged_data["Weekly_Sales"].max())
      print("Minimum Weekly Sales:", merged_data["Weekly_Sales"].min())
      print("Average Weekly Sales:", merged_data["Weekly_Sales"].mean())
```

```
Maximum Weekly Sales: 693099.36
Minimum Weekly Sales: 0.0
Average Weekly Sales: 16031.557675777807
```

```
print(f"Maximum Temperature: {merged_data['Temperature'].max():.2f}°C")
print(f"Minimum Temperature: {merged_data['Temperature'].min():.2f}°C")
print(f"Average Temperature: {merged_data['Temperature'].mean():.2f}°C")
```

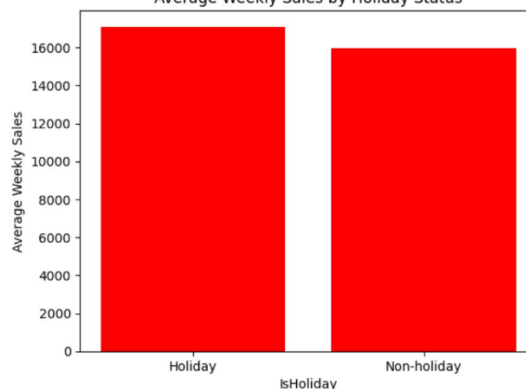
```
Maximum Temperature: 100.14°C
Minimum Temperature: 5.54°C
Average Temperature: 60.10°C
```

Pie Chart for Type of Store



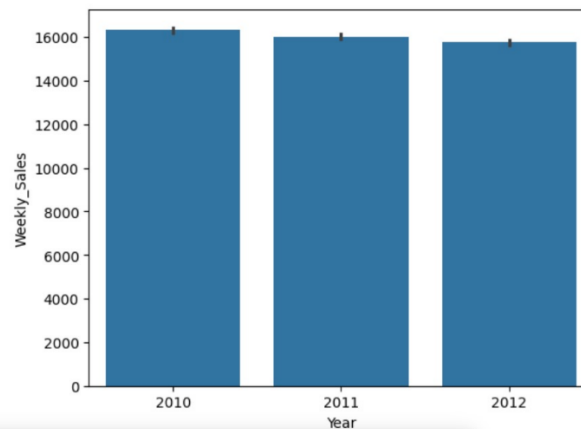
Bivariate Analysis

Average Weekly Sales by Holiday Status

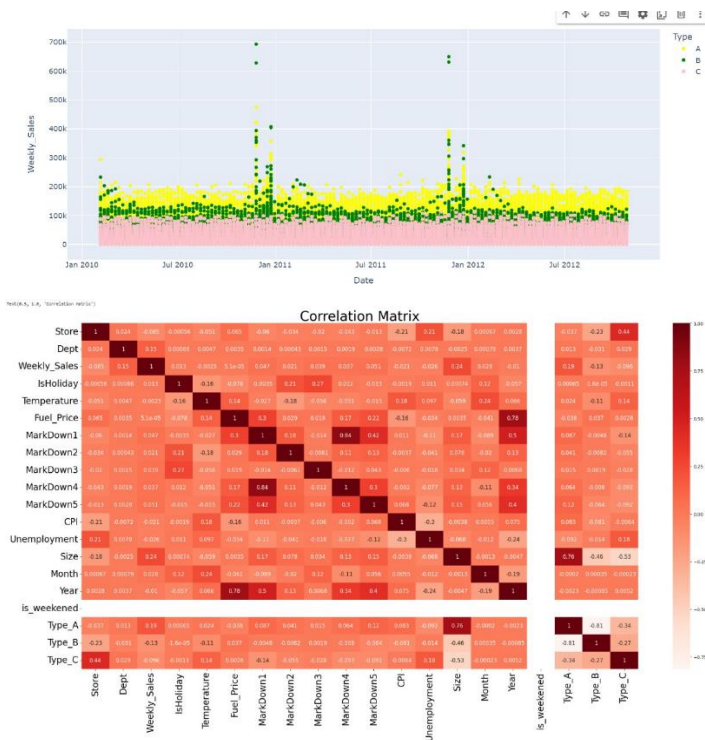
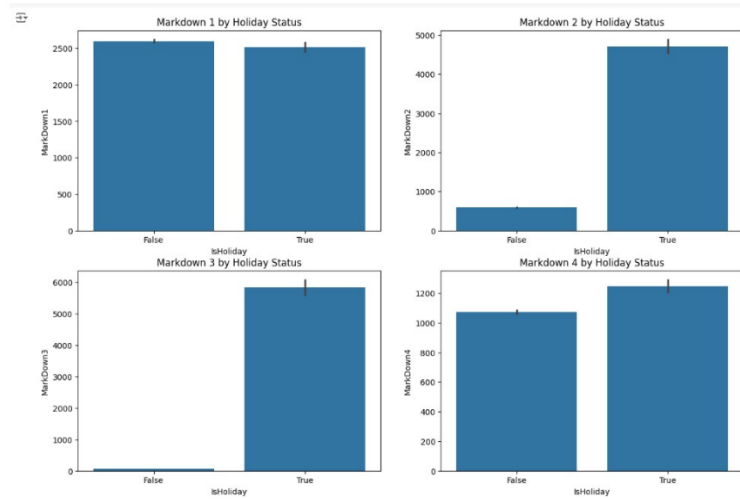


```
# bar graph
sns.barplot(x=merged_data["Year"], y=merged_data["Weekly_Sales"])
```

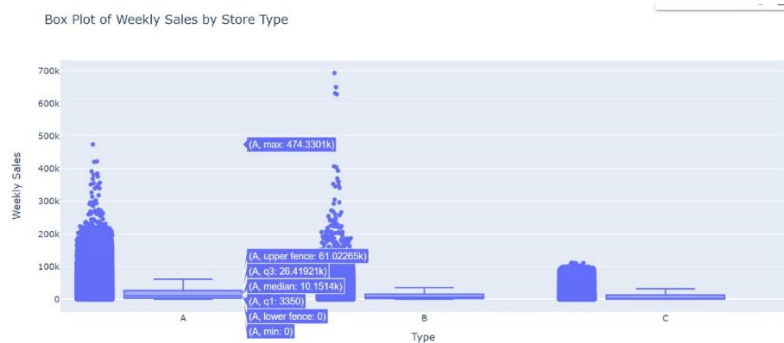
```
<Axes: xlabel='Year', ylabel='Weekly_Sales'>
```



Multivariate Analysis



Outliers and Anomalies



Data Preprocessing Code Screenshots

Loading Data

```
[2] # reading all the csv files
stores = pd.read_csv("stores.csv")
features = pd.read_csv("features.csv.zip")
train = pd.read_csv("train.csv.zip")
test = pd.read_csv("test.csv.zip")
```

```
[12] # merging all the csv files
# all the csv files have store column in common.
merged_data = train.merge(features, on="Store", how="inner").merge(stores, on="Store", how="inner")

merged_data
```

	Store	Dept	Date	Weekly_Sales	IsHoliday_x	Temperature	Fuel_Price	MarkDown1	MarkDown2	MarkDown3	MarkDown4	MarkDown5	CPI	Unemp1
0	1	1	2010-02-05	24824.50	False	42.31	2.572	NaN	NaN	NaN	NaN	NaN	211.096358	
1	1	2	2010-02-05	50605.27	False	42.31	2.572	NaN	NaN	NaN	NaN	NaN	211.096358	
2	1	3	2010-02-05	13740.12	False	42.31	2.572	NaN	NaN	NaN	NaN	NaN	211.096358	
3	1	4	2010-02-05	39954.04	False	42.31	2.572	NaN	NaN	NaN	NaN	NaN	211.096358	
4	1	5	2010-02-05	32229.38	False	42.31	2.572	NaN	NaN	NaN	NaN	NaN	211.096358	

Handling Missing Data

Code for identifying and handling missing values.

```
# Handling the null values
merged_data["MarkDown1"] = merged_data["MarkDown1"].replace(np.nan,0)
merged_data["MarkDown2"] = merged_data["MarkDown2"].replace(np.nan,0)
merged_data["MarkDown3"] = merged_data["MarkDown3"].replace(np.nan,0)
merged_data["MarkDown4"] = merged_data["MarkDown4"].replace(np.nan,0)
merged_data["MarkDown5"] = merged_data["MarkDown5"].replace(np.nan,0)
```

Data Transformation

Code for transforming variables (scaling, normalization).

```
# changing the categorical value type into numbers
merged_data = pd.get_dummies(merged_data, columns=["Type"])

merged_data["is_weekened"].replace({False:0, True:1}, inplace=True)

merged_data["IsHoliday"].replace({False:0, True:1}, inplace=True)

# Scaling the data
sc = StandardScaler()
X = sc.fit_transform(X)
print(X)
```

Feature Engineering

Code for creating new features or modifying existing ones.

```
# Date ,type and isholiday needs to be converted to numbers
merged_data["Date"] = pd.to_datetime(merged_data["Date"])
merged_data.loc[:, "DayofWeek"] = merged_data.loc[:, "Date"].dt.day_name()
merged_data.loc[:, "Month"] = merged_data.loc[:, "Date"].dt.month
merged_data.loc[:, "Year"] = merged_data.loc[:, "Date"].dt.year
```

Save Processed Data

Code to save the cleaned and processed data for future use.