

```
In [1]: 1 # Import necessary libraries
        2 import pandas as pd
        3 import seaborn as sns
        4 import numpy as np
        5 import matplotlib.pyplot as plt
        6 from matplotlib import dates
        7 from datetime import datetime
```

Business Understanding ¶

Walmart is an American retail corporation that operates a chain of hypermarkets, discount department stores, and grocery stores.

In this project, we focused to answer the following questions:

1. Which store has minimum and maximum sales?
2. Which store has maximum standard deviation i.e., the sales vary a lot. Also, find out the coefficient of mean to standard deviation
3. Which store/s has good quarterly growth rate in Q3'2012
4. Some holidays have a negative impact on sales. Find out holidays which have higher sales than the mean sales in non-holiday season for all stores together
5. Provide a monthly and semester view of sales in units and give insights
6. Build prediction to forecast demand.

Data Understanding

The data contains these features:

- Store - the store number
- Date - the week of sales
- Weekly_Sales - sales for the given store
- Holiday_Flag - whether the week is a special holiday week 1 – Holiday week 0 – Non-holiday week
- Temperature - Temperature on the day of sale
- Fuel_Price - Cost of fuel in the region
- CPI – Prevailing consumer price index
- Unemployment - Prevailing unemployment rate

In [3]:

```

1 # Load dataset
2 data = pd.read_csv('Walmart_Store_sales.csv')
3 data

```

Out[3]:

	Store	Date	Weekly_Sales	Holiday_Flag	Temperature	Fuel_Price	CPI	Unemployo
0	1	05-02-2010	1643690.90	0	42.31	2.572	211.096358	8
1	1	12-02-2010	1641957.44	1	38.51	2.548	211.242170	8
2	1	19-02-2010	1611968.17	0	39.93	2.514	211.289143	8
3	1	26-02-2010	1409727.59	0	46.63	2.561	211.319643	8
4	1	05-03-2010	1554806.68	0	46.50	2.625	211.350143	8
...
6430	45	28-09-2012	713173.95	0	64.88	3.997	192.013558	8
6431	45	05-10-2012	733455.07	0	64.89	3.985	192.170412	8
6432	45	12-10-2012	734464.36	0	54.47	4.000	192.327265	8
6433	45	19-10-2012	718125.53	0	56.47	3.969	192.330854	8
6434	45	26-10-2012	760281.43	0	58.85	3.882	192.308899	8

6435 rows × 8 columns



Data Preparation

```
In [4]: 1 # Convert date to datetime format and show dataset information
        2 data['Date'] = pd.to_datetime(data['Date'])
        3 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6435 entries, 0 to 6434
Data columns (total 8 columns):
#   Column          Non-Null Count  Dtype
---  -
0   Store            6435 non-null   int64
1   Date             6435 non-null   datetime64[ns]
2   Weekly_Sales     6435 non-null   float64
3   Holiday_Flag     6435 non-null   int64
4   Temperature      6435 non-null   float64
5   Fuel_Price       6435 non-null   float64
6   CPI              6435 non-null   float64
7   Unemployment     6435 non-null   float64
dtypes: datetime64[ns](1), float64(5), int64(2)
memory usage: 402.3 KB
```

C:\Users\Parag\AppData\Local\Temp\ipykernel_10960\236554556.py:2: UserWarning: Parsing dates in DD/MM/YYYY format when dayfirst=False (the default) was specified. This may lead to inconsistently parsed dates! Specify a format to ensure consistent parsing.

```
data['Date'] = pd.to_datetime(data['Date'])
```

```
In [5]: 1 # checking for missing values
        2 data.isnull().sum()
        3
```

```
Out[5]: Store            0
        Date            0
        Weekly_Sales     0
        Holiday_Flag     0
        Temperature      0
        Fuel_Price       0
        CPI              0
        Unemployment     0
        dtype: int64
```

In [6]:

```
1 # Splitting Date and create new columns (Day, Month, and Year)
2 data["Day"] = pd.DatetimeIndex(data['Date']).day
3 data['Month'] = pd.DatetimeIndex(data['Date']).month
4 data['Year'] = pd.DatetimeIndex(data['Date']).year
5 data
```

Out[6]:

	Store	Date	Weekly_Sales	Holiday_Flag	Temperature	Fuel_Price	CPI	Unemploy
0	1	2010-05-02	1643690.90	0	42.31	2.572	211.096358	{
1	1	2010-12-02	1641957.44	1	38.51	2.548	211.242170	{
2	1	2010-02-19	1611968.17	0	39.93	2.514	211.289143	{
3	1	2010-02-26	1409727.59	0	46.63	2.561	211.319643	{
4	1	2010-05-03	1554806.68	0	46.50	2.625	211.350143	{
...	
6430	45	2012-09-28	713173.95	0	64.88	3.997	192.013558	{
6431	45	2012-05-10	733455.07	0	64.89	3.985	192.170412	{
6432	45	2012-12-10	734464.36	0	54.47	4.000	192.327265	{
6433	45	2012-10-19	718125.53	0	56.47	3.969	192.330854	{
6434	45	2012-10-26	760281.43	0	58.85	3.882	192.308899	{

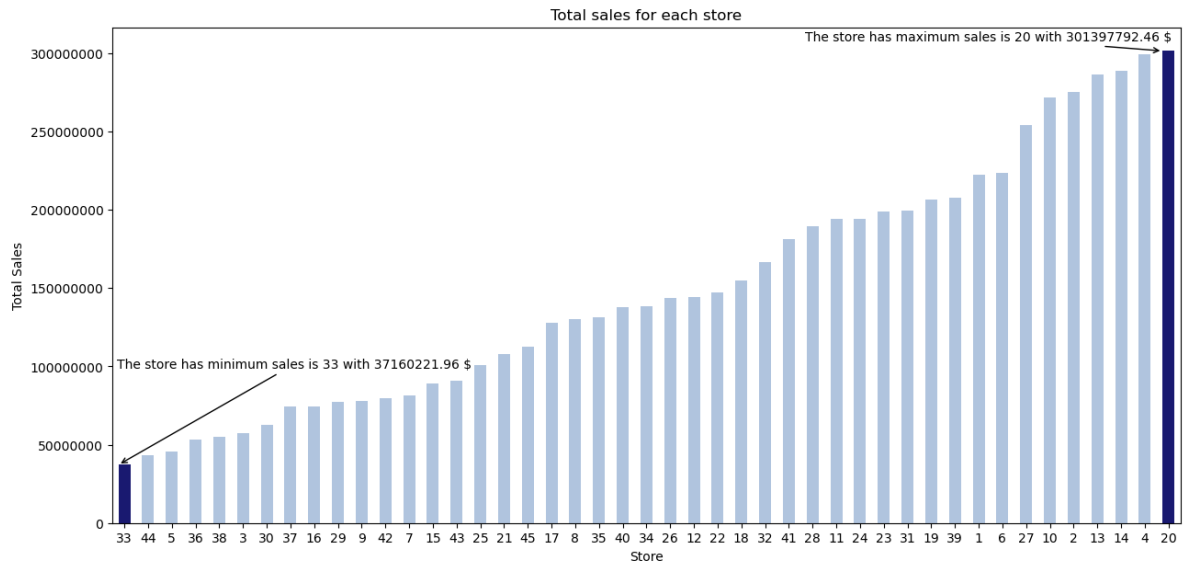
6435 rows × 11 columns



Q1: Which store has minimum and maximum sales?

```
In [7]: 1 plt.figure(figsize=(15,7))
2
3 # Sum Weekly_Sales for each store, then sortded by total sales
4 total_sales_for_each_store = data.groupby('Store')['Weekly_Sales'].sum().sort_values(ascending=False)
5 total_sales_for_each_store_array = np.array(total_sales_for_each_store) #
6
7 # Assigning a specific color for the stores have the lowest and highest sales
8 clr = ['lightsteelblue' if ((x < max(total_sales_for_each_store_array) * 0.1)) else 'lightcoral' if ((x > min(total_sales_for_each_store_array) * 0.1)) else 'lightgray' for x in total_sales_for_each_store_array]
9
10
11 ax = total_sales_for_each_store.plot(kind='bar',color=clr);
12
13 # store have minimum sales
14 p = ax.patches[0]
15 print(type(p.get_height()))
16 ax.annotate("The store has minimum sales is 33 with {0:.2f} $".format((p.get_height() * 1000000)),
17             xytext=(0.17, 0.32), textcoords='axes fraction',
18             arrowprops=dict(arrowstyle="->", connectionstyle="arc3"),
19             horizontalalignment='center', verticalalignment='center')
20
21
22 # store have maximum sales
23 p = ax.patches[44]
24 ax.annotate("The store has maximum sales is 20 with {0:.2f} $".format((p.get_height() * 1000000)),
25             xytext=(0.82, 0.98), textcoords='axes fraction',
26             arrowprops=dict(arrowstyle="->", connectionstyle="arc3"),
27             horizontalalignment='center', verticalalignment='center')
28
29
30 # plot properties
31 plt.xticks(rotation=0)
32 plt.ticklabel_format(useOffset=False, style='plain', axis='y')
33 plt.title('Total sales for each store')
34 plt.xlabel('Store')
35 plt.ylabel('Total Sales');
```

```
<class 'numpy.float64'>
```



Q2: Which store has maximum standard deviation i.e., the sales vary a lot. Also, find out the coefficient of mean to standard deviation?

In [8]:

```
1 # Which store has maximum standard deviation
2 data_std = pd.DataFrame(data.groupby('Store')['Weekly_Sales'].std().sort_
3 print("The store has maximum standard deviation is "+str(data_std.head(1).
```

The store has maximum standard deviation is 14 with 317570 \$

```
In [9]: 1 # Distribution of store has maximum standard deviation
2 plt.figure(figsize=(15,7))
3 sns.distplot(data[data['Store'] == data_std.head(1).index[0]]['Weekly_Sales'])
4 plt.title('The Sales Distribution of Store #' + str(data_std.head(1).index[0]))
```

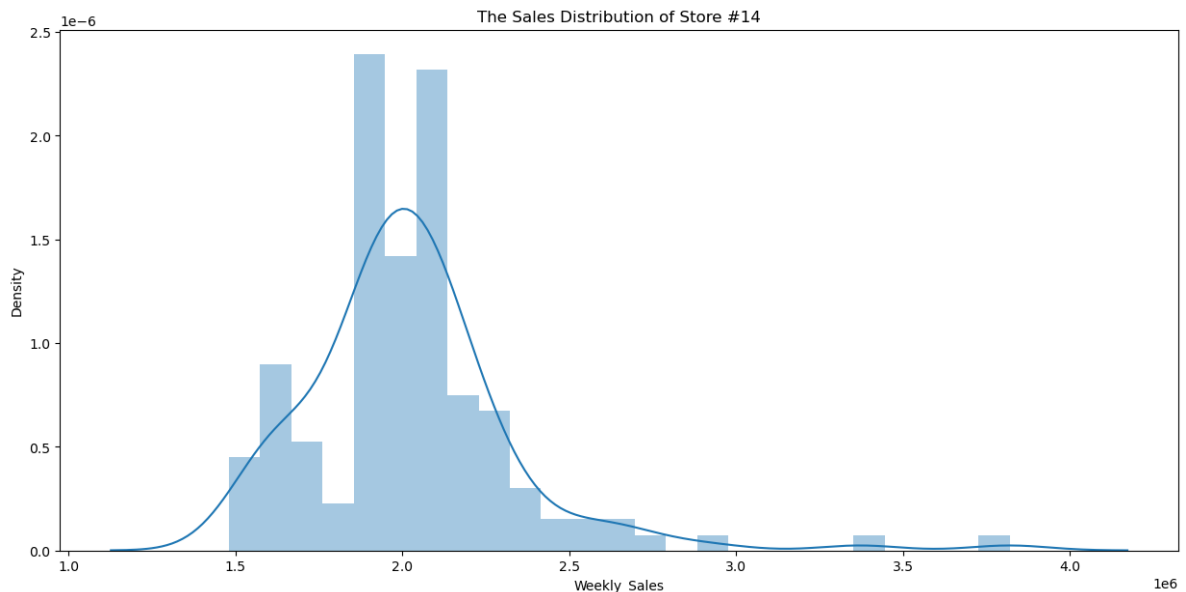
C:\Users\Parag\AppData\Local\Temp\ipykernel_10960\3470610508.py:3: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751> (<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>)

```
sns.distplot(data[data['Store'] == data_std.head(1).index[0]]['Weekly_Sales'])
```



In [10]:

```
1 # Coefficient of mean to standard deviation
2 coef_mean_std = pd.DataFrame(data.groupby('Store')['Weekly_Sales'].std()
3 coef_mean_std = coef_mean_std.rename(columns={'Weekly_Sales': 'Coefficient
4 coef_mean_std
```


Out[10]:

Coefficient of mean to standard deviation

Store	
1	0.100292
2	0.123424
3	0.115021
4	0.127083
5	0.118668
6	0.135823
7	0.197305
8	0.116953
9	0.126895
10	0.159133
11	0.122262
12	0.137925
13	0.132514
14	0.157137
15	0.193384
16	0.165181
17	0.125521
18	0.162845
19	0.132680
20	0.130903
21	0.170292
22	0.156783
23	0.179721
24	0.123637
25	0.159860
26	0.110111
27	0.135155
28	0.137330
29	0.183742
30	0.052008
31	0.090161
32	0.118310
33	0.092868
34	0.108225
35	0.229681

Coefficient of mean to standard deviation

Store	
36	0.162579
37	0.042084
38	0.110875
39	0.149908
40	0.123430
41	0.148177
42	0.090335
43	0.064104
44	0.081793
45	0.165613

```
In [11]: 1 # Distribution of store has maximum coefficient of mean to standard deviation
2 coef_mean_std_max = coef_mean_std.sort_values(by='Coefficient of mean to s
3 plt.figure(figsize=(15,7))
4 sns.distplot(data[data['Store'] == coef_mean_std_max.tail(1).index[0]]['Wee
5 plt.title('The Sales Distribution of Store #' + str(coef_mean_std_max.tail(1
```

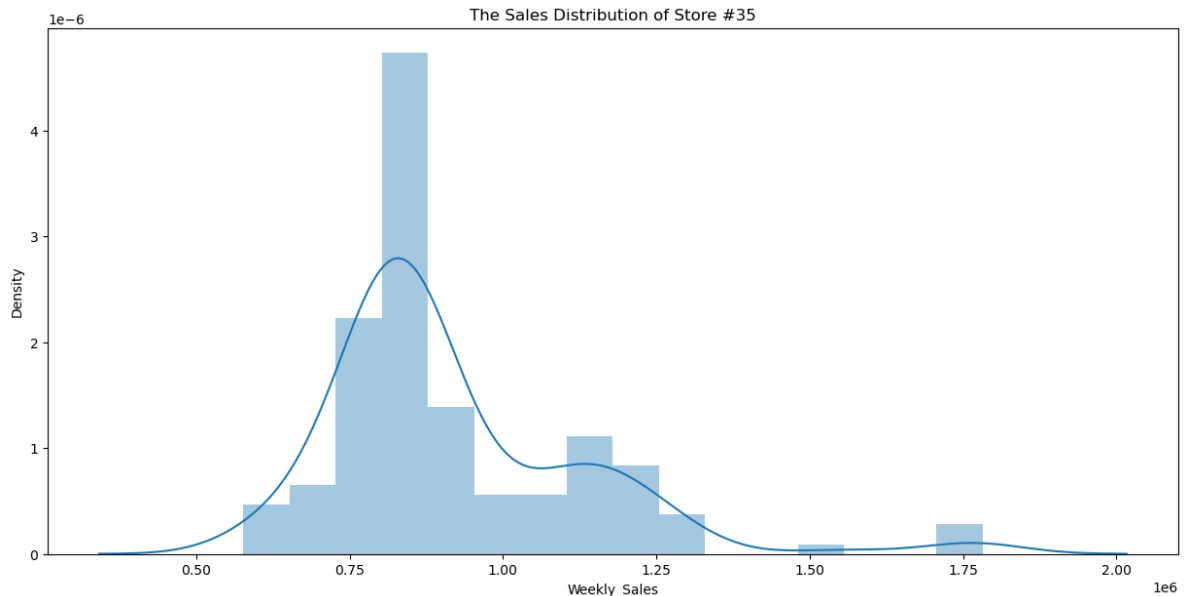
C:\Users\Parag\AppData\Local\Temp\ipykernel_10960\1932089423.py:4: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

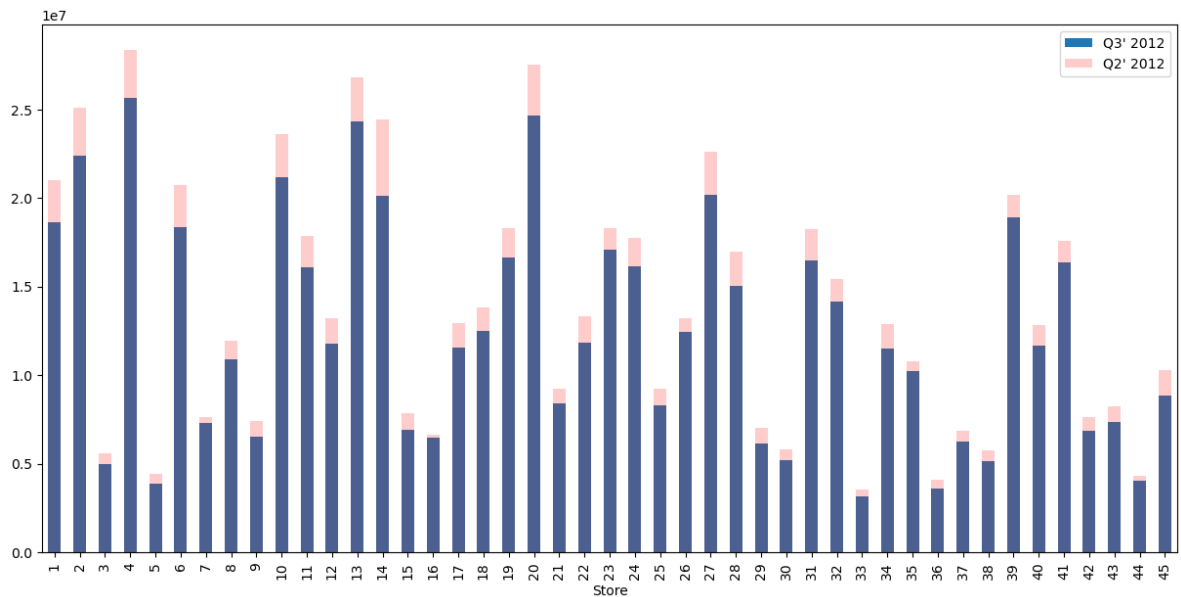
For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751> (<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>)

```
sns.distplot(data[data['Store'] == coef_mean_std_max.tail(1).index[0]]['Weekly_Sales'])
```



Q3: Which store/s has good quarterly growth rate in Q3'2012

```
In [13]: 1 plt.figure(figsize=(15, 7))
2
3 # Sales for third quarterly in 2012
4 Q3 = data[(data['Date'] > '2012-07-01') & (data['Date'] < '2012-09-30')].groupby('Store').sum()
5
6 # Sales for second quarterly in 2012
7 Q2 = data[(data['Date'] > '2012-04-01') & (data['Date'] < '2012-06-30')].groupby('Store').sum()
8
9 # Plotting the difference between sales for second and third quarterly
10 Q2.plot(ax=Q3.plot(kind='bar', legend=True), kind='bar', color='r', alpha=0.5)
11 plt.legend(["Q3' 2012", "Q2' 2012"])
12 plt.show()
```



```
In [14]: 1 # store/s has good quarterly growth rate in Q3'2012 - .sort_values(by='Sales')
2 print('Store have good quarterly growth rate in Q3'2012 is Store ' + str(Q3.sort_values(by='Sales').index[0]))
```

Store have good quarterly growth rate in Q3'2012 is Store 4 With 25652119.35 \$

Q4: Some holidays have a negative impact on sales. Find out holidays which have higher sales than the mean sales in non-holiday season for all stores together

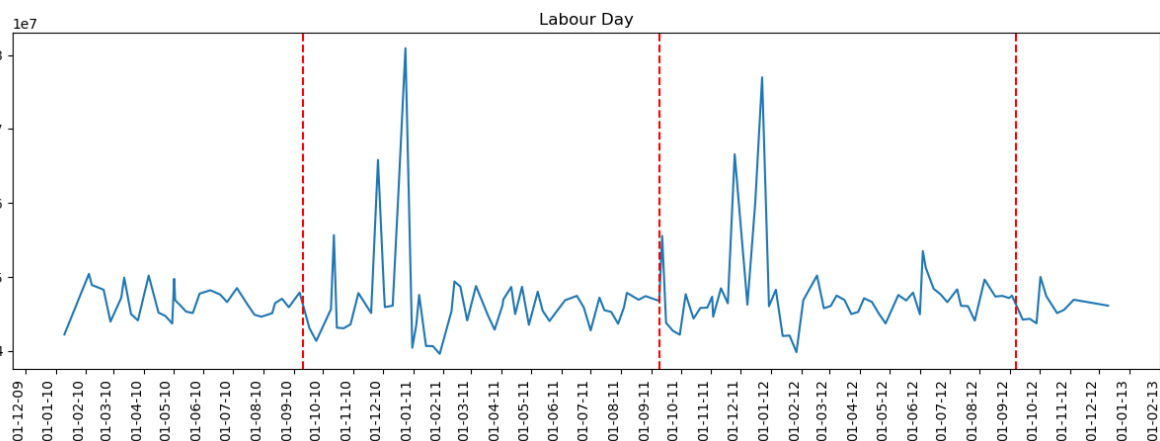
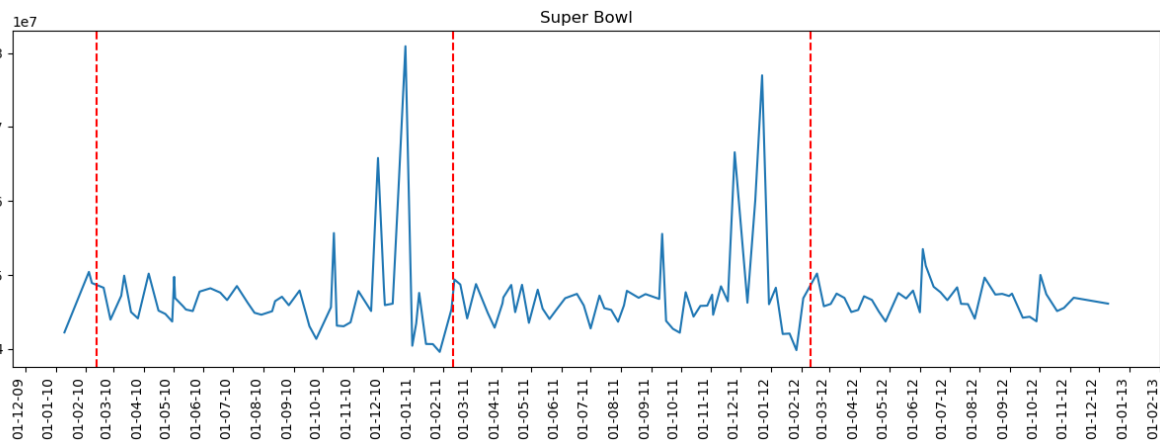
Holiday Events:

- Super Bowl: 12-Feb-10, 11-Feb-11, 10-Feb-12, 8-Feb-13
- Labour Day: 10-Sep-10, 9-Sep-11, 7-Sep-12, 6-Sep-13
- Thanksgiving: 26-Nov-10, 25-Nov-11, 23-Nov-12, 29-Nov-13
- Christmas: 31-Dec-10, 30-Dec-11, 28-Dec-12, 27-Dec-13

```

In [15]: 1 def plot_line(df,holiday_dates,holiday_label):
2         fig, ax = plt.subplots(figsize = (15,5))
3         ax.plot(df['Date'],df['Weekly_Sales'],label=holiday_label)
4
5         for day in holiday_dates:
6             day = datetime.strptime(day, '%d-%m-%Y')
7             plt.axvline(x=day, linestyle='--', c='r')
8
9
10        plt.title(holiday_label)
11        x_dates = df['Date'].dt.strftime('%Y-%m-%d').sort_values().unique()
12        xfmt = dates.DateFormatter('%d-%m-%y')
13        ax.xaxis.set_major_formatter(xfmt)
14        ax.xaxis.set_major_locator(dates.DayLocator(1))
15        plt.gcf().autofmt_xdate(rotation=90)
16        plt.show()
17
18
19        total_sales = data.groupby('Date')['Weekly_Sales'].sum().reset_index()
20        Super_Bowl = ['12-2-2010', '11-2-2011', '10-2-2012']
21        Labour_Day = ['10-9-2010', '9-9-2011', '7-9-2012']
22        Thanksgiving = ['26-11-2010', '25-11-2011', '23-11-2012']
23        Christmas = ['31-12-2010', '30-12-2011', '28-12-2012']
24
25        plot_line(total_sales,Super_Bowl,'Super Bowl')
26        plot_line(total_sales,Labour_Day,'Labour Day')
27        plot_line(total_sales,Thanksgiving,'Thanksgiving')
28        plot_line(total_sales,Christmas,'Christmas')

```





The sales increased during thanksgiving. And the sales decreased during christmas.



```
In [16]: 1 data.loc[data.Date.isin(Super_Bowl)]
```

Out[16]:

	Store	Date	Weekly_Sales	Holiday_Flag	Temperature	Fuel_Price	CPI	Unemploy
1	1	2010-12-02	1641957.44	1	38.51	2.548	211.242170	8
53	1	2011-11-02	1649614.93	1	36.39	3.022	212.936705	7
105	1	2012-10-02	1802477.43	1	48.02	3.409	220.265178	7
144	2	2010-12-02	2137809.50	1	38.49	2.548	210.897994	8
196	2	2011-11-02	2168041.61	1	33.19	3.022	212.592862	8
...
6202	44	2011-11-02	307486.73	1	30.83	3.034	127.859129	7
6254	44	2012-10-02	325377.97	1	33.73	3.116	130.384903	8
6293	45	2010-12-02	656988.64	1	27.73	2.773	181.982317	8
6345	45	2011-11-02	766456.00	1	30.30	3.239	183.701613	8
6397	45	2012-10-02	803657.12	1	37.00	3.640	189.707605	8

135 rows × 11 columns

In [18]:

```

1 # Yearly Sales in holidays
2 Super_Bowl_df = pd.DataFrame(data.loc[data.Date.isin(Super_Bowl)].groupby('Year')['Weekly_Sales'].sum())
3 Thanksgiving_df = pd.DataFrame(data.loc[data.Date.isin(Thanksgiving)].groupby('Year')['Weekly_Sales'].sum())
4 Labour_Day_df = pd.DataFrame(data.loc[data.Date.isin(Labour_Day)].groupby('Year')['Weekly_Sales'].sum())
5 Christmas_df = pd.DataFrame(data.loc[data.Date.isin(Christmas)].groupby('Year')['Weekly_Sales'].sum())
6
7 Super_Bowl_df.plot(kind='bar', legend=False, title='Yearly Sales in Super Bowl holiday')
8 Thanksgiving_df.plot(kind='bar', legend=False, title='Yearly Sales in Thanksgiving holiday')
9 Labour_Day_df.plot(kind='bar', legend=False, title='Yearly Sales in Labour Day holiday')
10 Christmas_df.plot(kind='bar', legend=False, title='Yearly Sales in Christmas holiday')

```

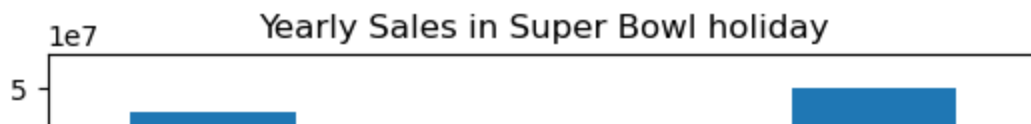
C:\Users\Parag\AppData\Local\Temp\ipykernel_10960\28647724.py:3: UserWarning: Parsing dates in DD/MM/YYYY format when dayfirst=False (the default) was specified. This may lead to inconsistently parsed dates! Specify a format to ensure consistent parsing.

```
Thanksgiving_df = pd.DataFrame(data.loc[data.Date.isin(Thanksgiving)].groupby('Year')['Weekly_Sales'].sum())
```

C:\Users\Parag\AppData\Local\Temp\ipykernel_10960\28647724.py:5: UserWarning: Parsing dates in DD/MM/YYYY format when dayfirst=False (the default) was specified. This may lead to inconsistently parsed dates! Specify a format to ensure consistent parsing.

```
Christmas_df = pd.DataFrame(data.loc[data.Date.isin(Christmas)].groupby('Year')['Weekly_Sales'].sum())
```

Out[18]: <Axes: title={'center': 'Yearly Sales in Christmas holiday'}, xlabel='Year'>



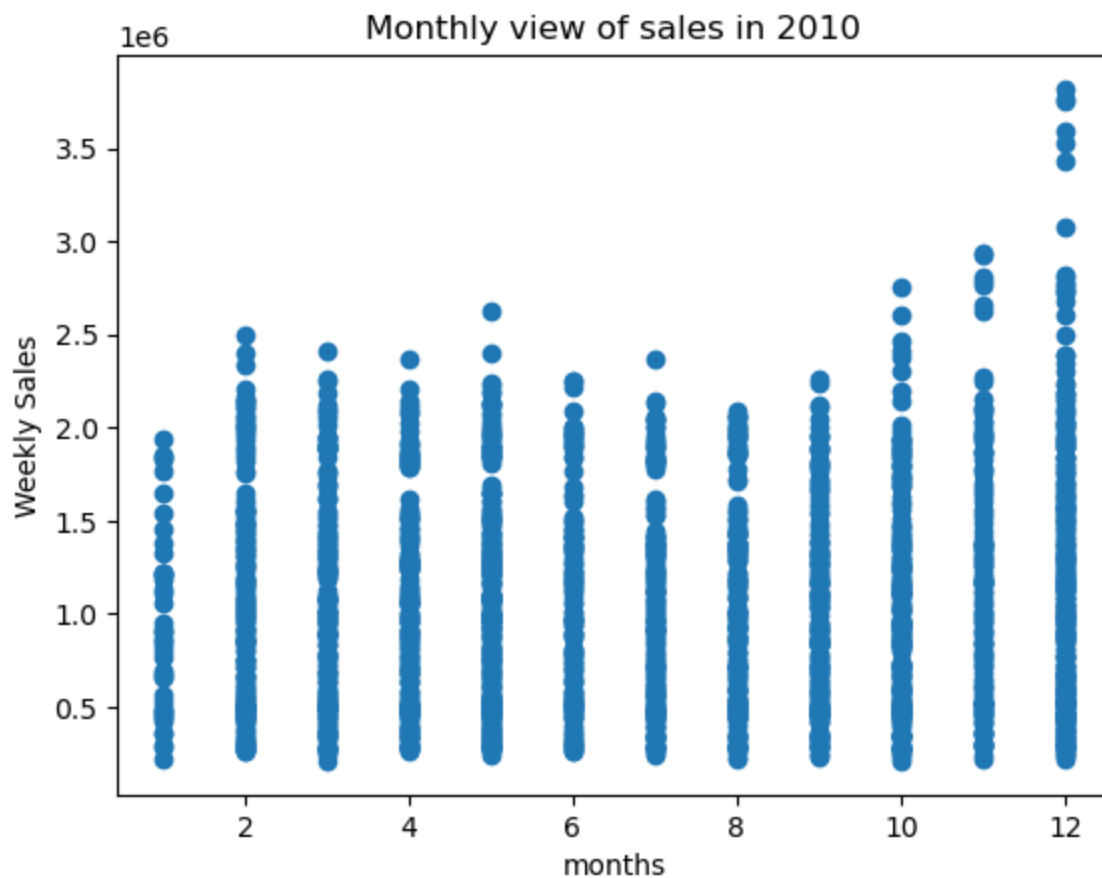
Q5: Provide a monthly and semester view of sales in units and give insights

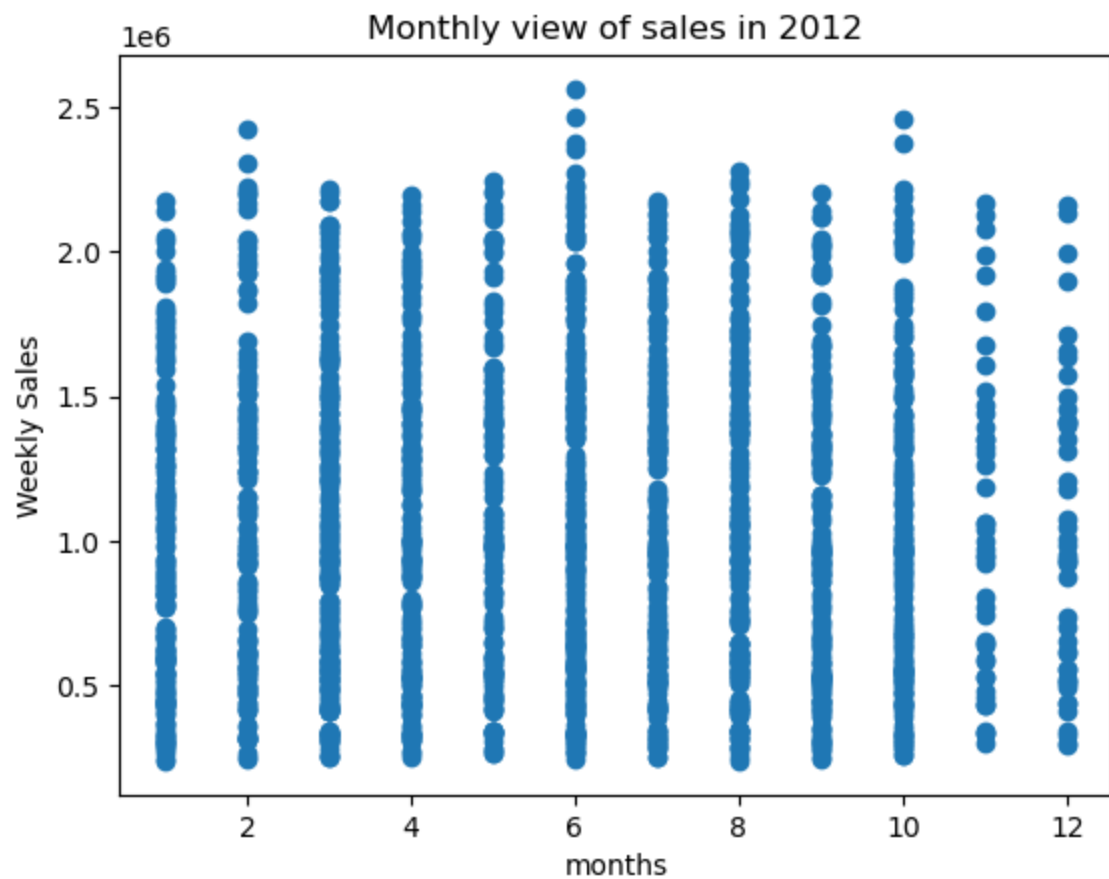
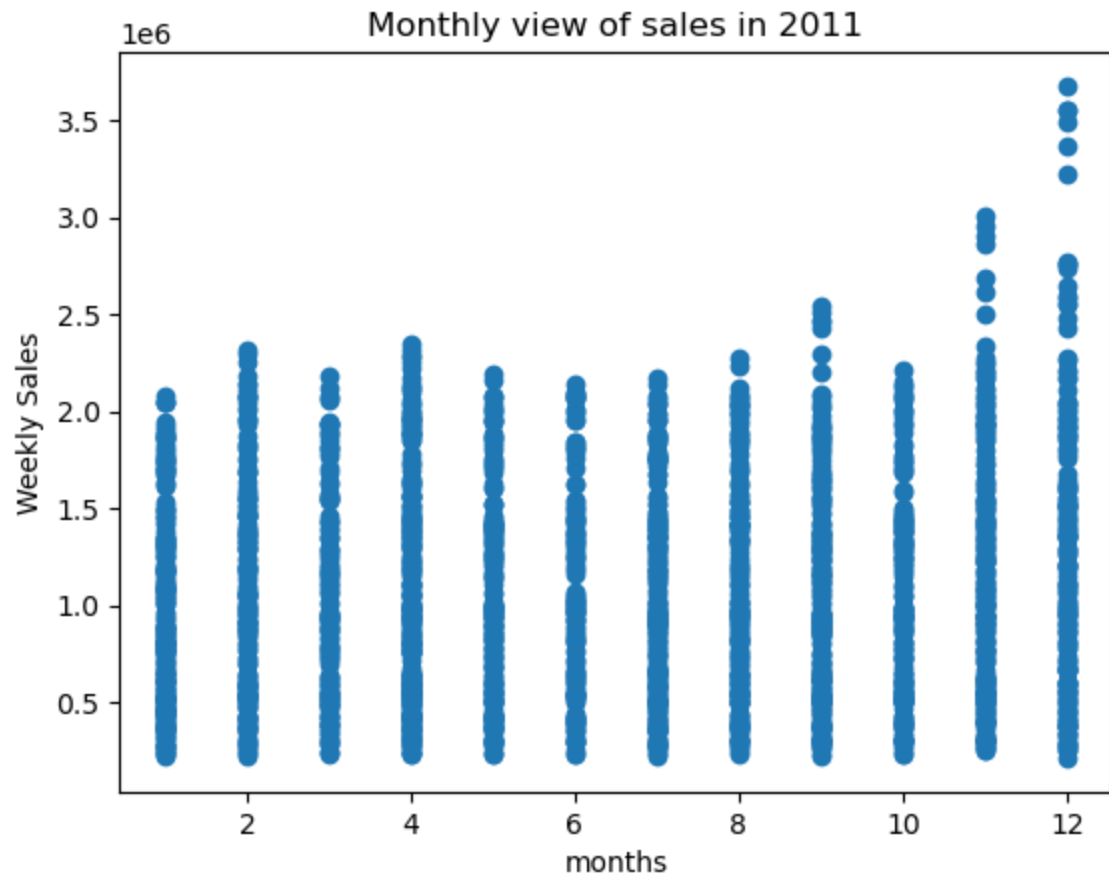
In [19]:

```

1 # Monthly view of sales for each years
2 plt.scatter(data[data.Year==2010]["Month"],data[data.Year==2010]["Weekly_Sales"])
3 plt.xlabel("months")
4 plt.ylabel("Weekly Sales")
5 plt.title("Monthly view of sales in 2010")
6 plt.show()
7
8 plt.scatter(data[data.Year==2011]["Month"],data[data.Year==2011]["Weekly_Sales"])
9 plt.xlabel("months")
10 plt.ylabel("Weekly Sales")
11 plt.title("Monthly view of sales in 2011")
12 plt.show()
13
14 plt.scatter(data[data.Year==2012]["Month"],data[data.Year==2012]["Weekly_Sales"])
15 plt.xlabel("months")
16 plt.ylabel("Weekly Sales")
17 plt.title("Monthly view of sales in 2012")
18 plt.show()

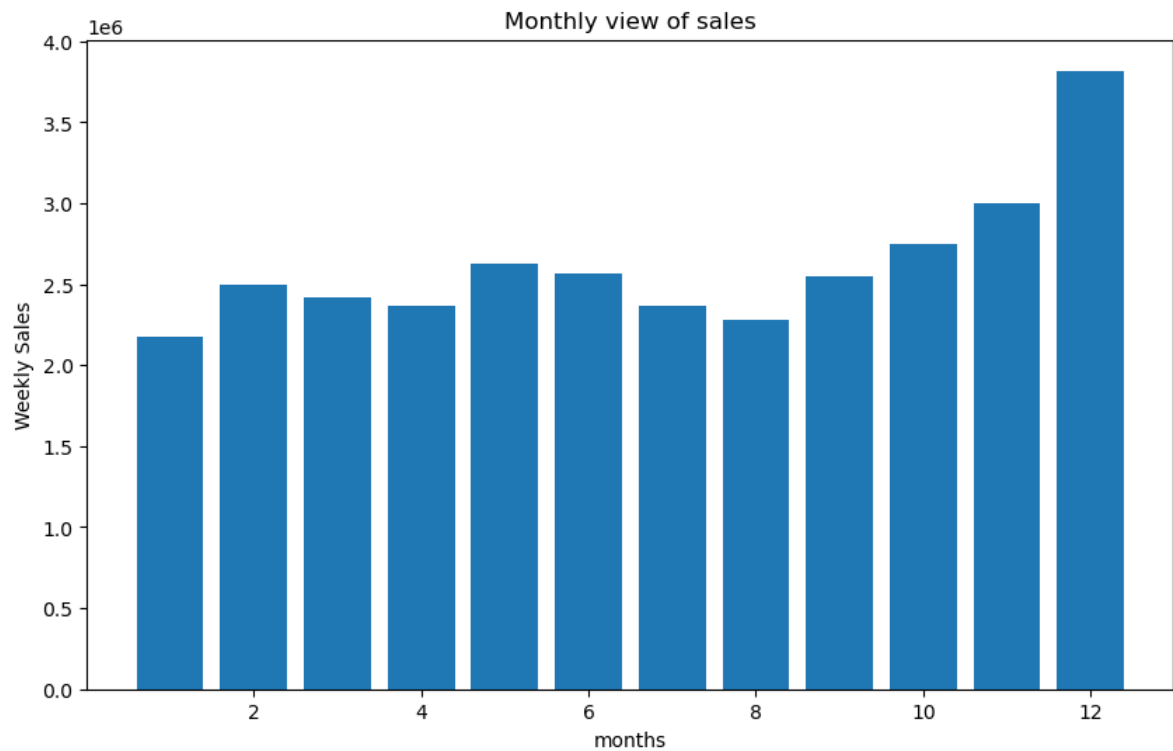
```





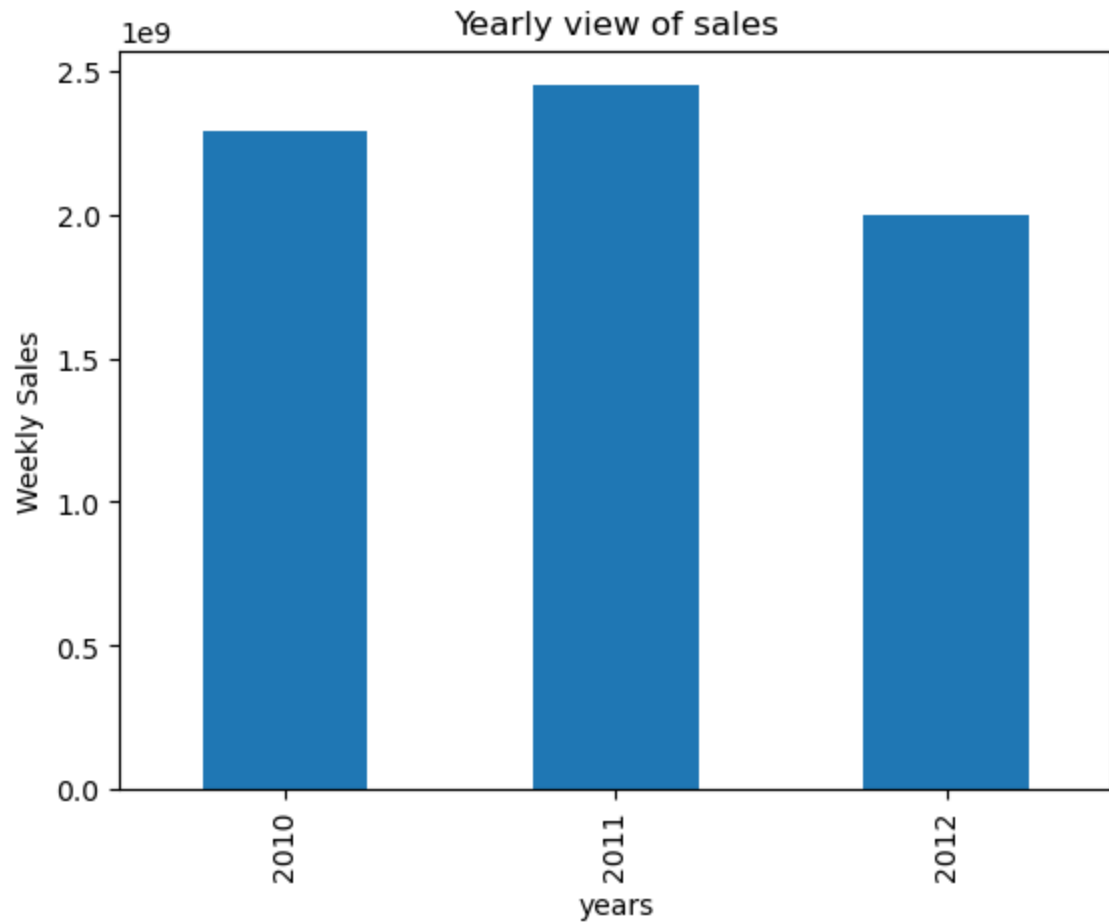
```
In [20]: 1 # Monthly view of sales for all years
2 plt.figure(figsize=(10,6))
3 plt.bar(data["Month"],data["Weekly_Sales"])
4 plt.xlabel("months")
5 plt.ylabel("Weekly Sales")
6 plt.title("Monthly view of sales")
```

Out[20]: Text(0.5, 1.0, 'Monthly view of sales')



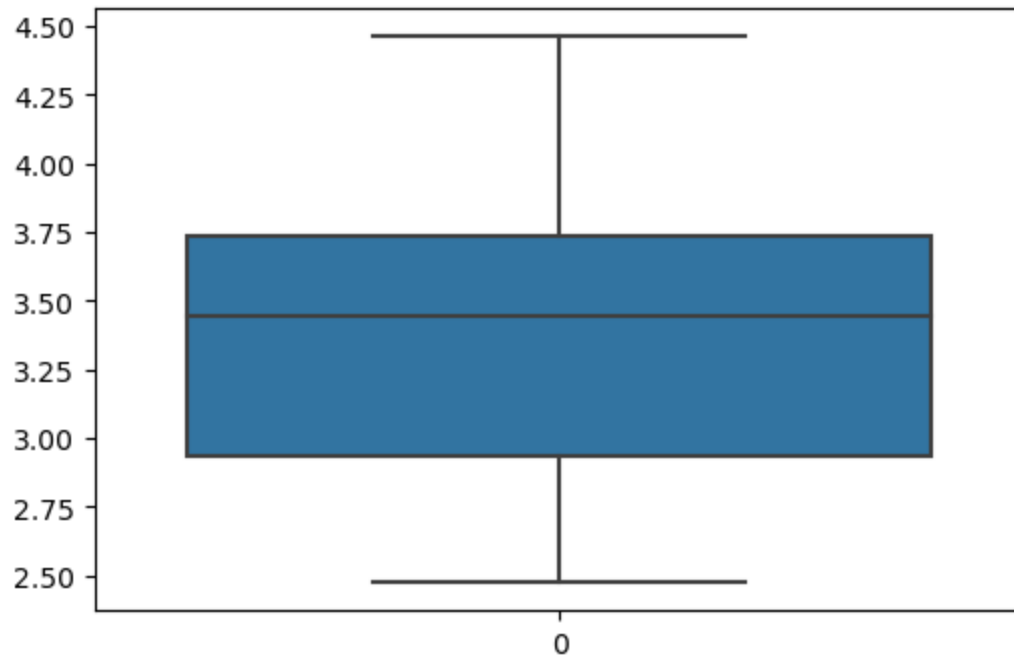
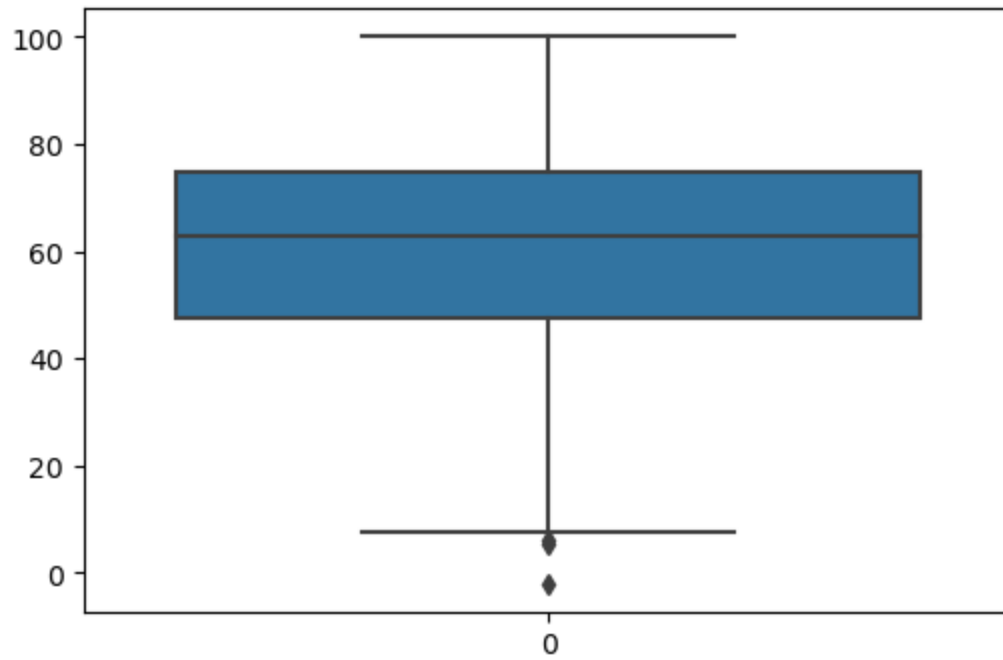
```
In [21]: 1 # Yearly view of sales
2 plt.figure(figsize=(10,6))
3 data.groupby("Year")[["Weekly_Sales"]].sum().plot(kind='bar',legend=False)
4 plt.xlabel("years")
5 plt.ylabel("Weekly Sales")
6 plt.title("Yearly view of sales");
```

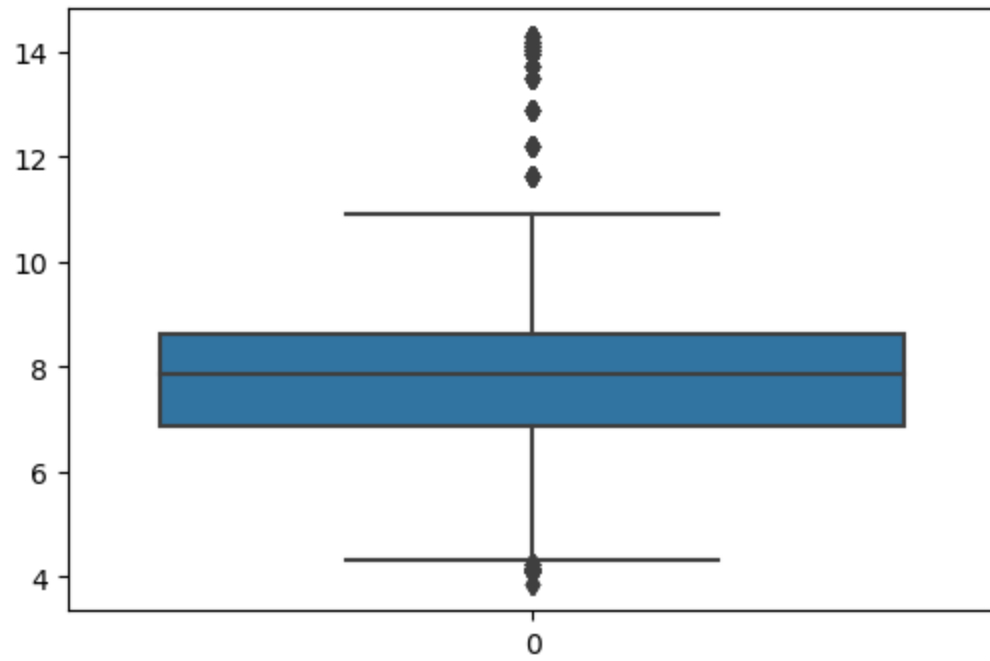
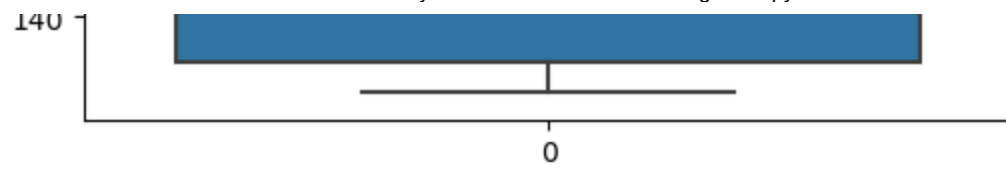
<Figure size 1000x600 with 0 Axes>



Build prediction models to forecast demand (Modeling)

```
In [22]: 1 # find outliers
2 fig, axs = plt.subplots(4,figsize=(6,18))
3 X = data[['Temperature','Fuel_Price','CPI','Unemployment']]
4 for i,column in enumerate(X):
5     sns.boxplot(data[column], ax=axs[i])
6
```



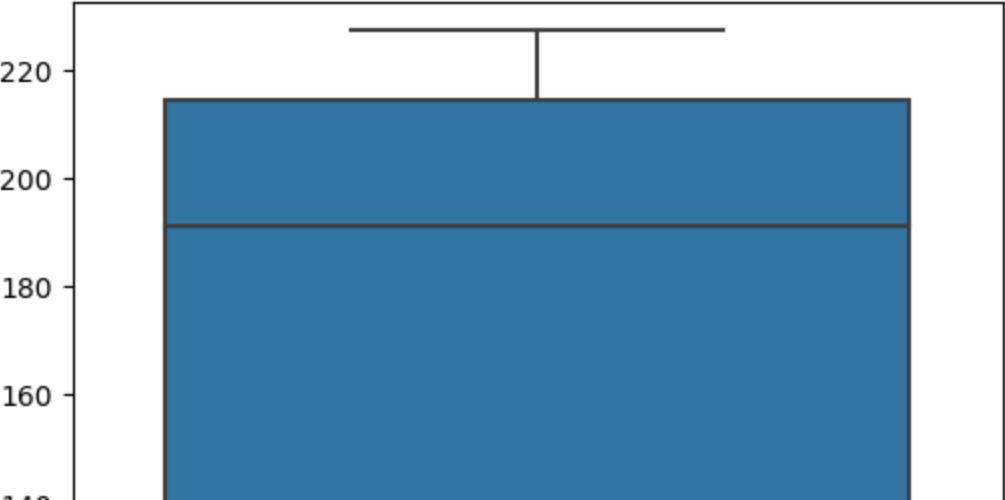
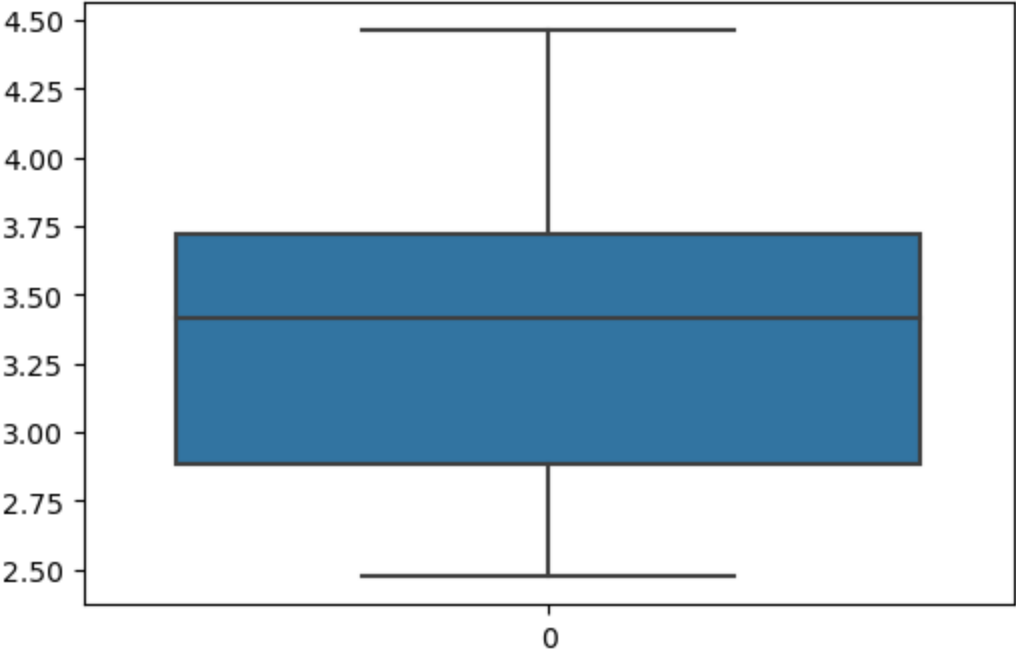
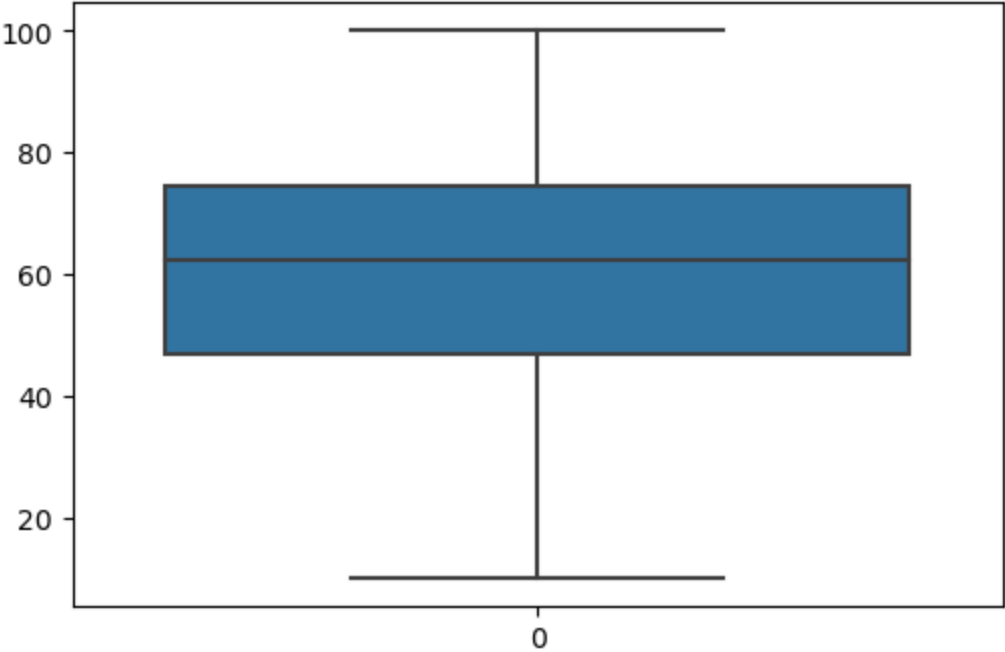

```
In [23]: 1 # drop the outliers
        2 data_new = data[(data['Unemployment']<10) & (data['Unemployment']>4.5) & (
        3 data_new
```

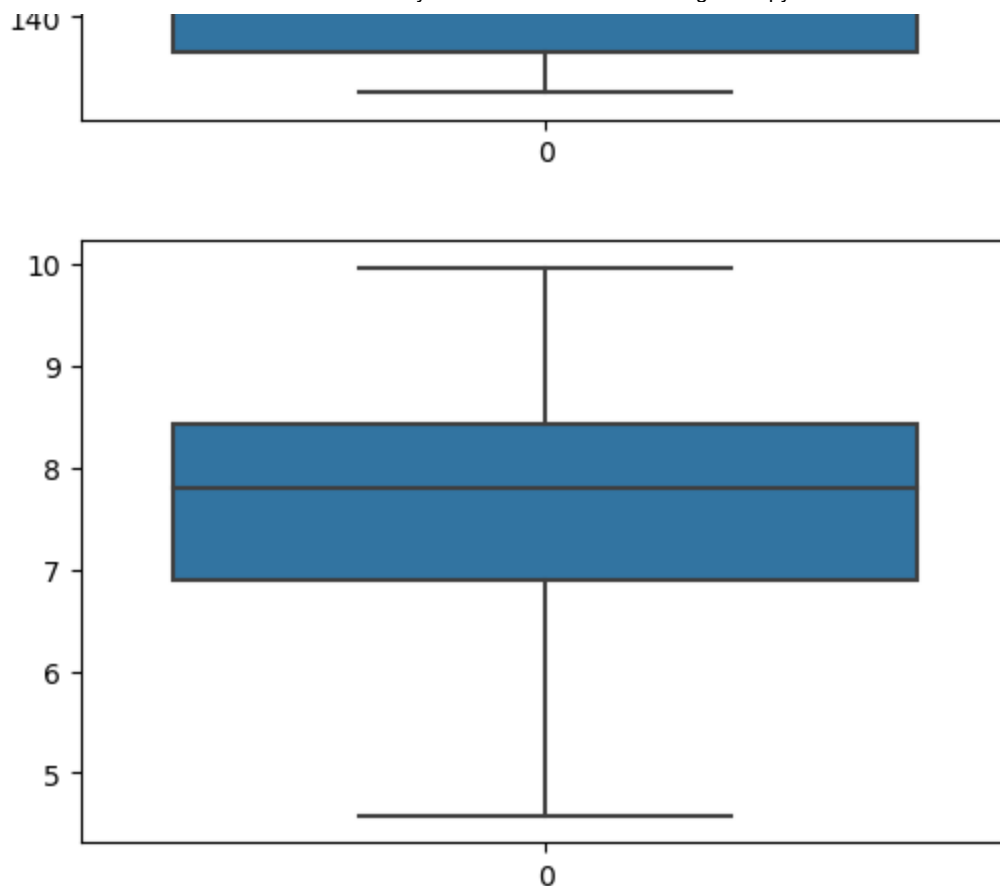
Out[23]:

	Store	Date	Weekly_Sales	Holiday_Flag	Temperature	Fuel_Price	CPI	Unemploy
0	1	2010-05-02	1643690.90	0	42.31	2.572	211.096358	{
1	1	2010-12-02	1641957.44	1	38.51	2.548	211.242170	{
2	1	2010-02-19	1611968.17	0	39.93	2.514	211.289143	{
3	1	2010-02-26	1409727.59	0	46.63	2.561	211.319643	{
4	1	2010-05-03	1554806.68	0	46.50	2.625	211.350143	{
...
6430	45	2012-09-28	713173.95	0	64.88	3.997	192.013558	{
6431	45	2012-05-10	733455.07	0	64.89	3.985	192.170412	{
6432	45	2012-12-10	734464.36	0	54.47	4.000	192.327265	{
6433	45	2012-10-19	718125.53	0	56.47	3.969	192.330854	{
6434	45	2012-10-26	760281.43	0	58.85	3.882	192.308899	{

5658 rows × 11 columns

```
In [24]: 1 # check outliers
2 fig, axs = plt.subplots(4,figsize=(6,18))
3 X = data_new[['Temperature','Fuel_Price','CPI','Unemployment']]
4 for i,column in enumerate(X):
5     sns.boxplot(data_new[column], ax=axs[i])
```



Build Model

In [25]:

```
1 # Import sklearn
2 from sklearn.ensemble import RandomForestRegressor
3 from sklearn.model_selection import train_test_split
4 from sklearn import metrics
5 from sklearn.linear_model import LinearRegression
```

In [26]:

```
1 # Select features and target
2 X = data_new[['Store', 'Fuel_Price', 'CPI', 'Unemployment', 'Day', 'Month', 'Year']]
3 y = data_new['Weekly_Sales']
4
5 # Split data to train and test (0.80:0.20)
6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

```

In [28]: 1 # Linear Regression model
2 print('Linear Regression:')
3 print()
4 reg = LinearRegression()
5 reg.fit(X_train, y_train)
6 y_pred = reg.predict(X_test)
7 print('Accuracy:', reg.score(X_train, y_train) * 100)
8
9 print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, y_pred))
10 print('Mean Squared Error:', metrics.mean_squared_error(y_test, y_pred))
11 print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(y_test, y_pred)))
12
13 # Plotting predicted vs. actual values
14 sns.scatterplot(x=y_pred, y=y_test)
15 plt.xlabel('Predicted Values')
16 plt.ylabel('Actual Values')
17 plt.title('Predicted vs. Actual Values')
18 plt.show()

```

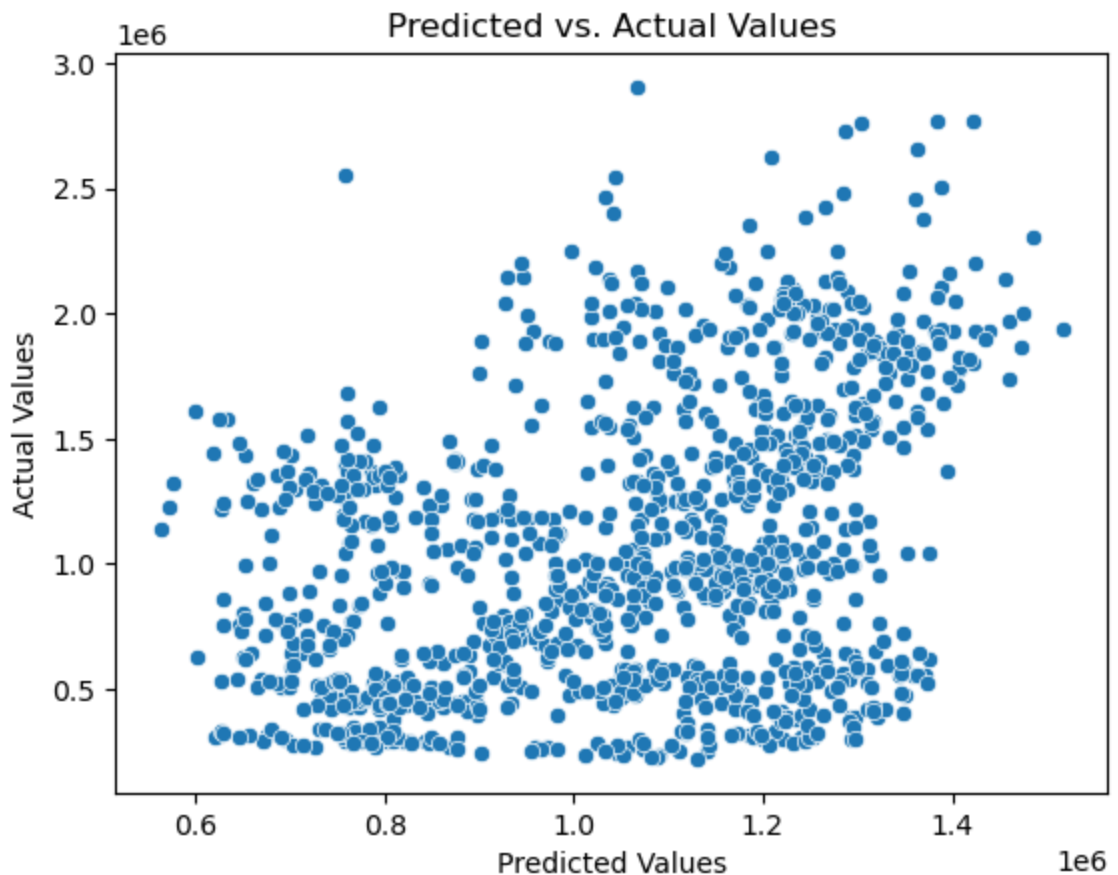
Linear Regression:

Accuracy: 12.984547620355913

Mean Absolute Error: 443498.0054904802

Mean Squared Error: 282327320550.53516

Root Mean Squared Error: 531344.8226439542



```

In [30]: 1 # Random Forest Regressor
2 print('Random Forest Regressor:')
3 print()
4 rfr = RandomForestRegressor(n_estimators=400, max_depth=15, n_jobs=5)
5 rfr.fit(X_train, y_train)
6 y_pred = rfr.predict(X_test)
7 print('Accuracy:', rfr.score(X_test, y_test) * 100)
8
9 print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, y_pred))
10 print('Mean Squared Error:', metrics.mean_squared_error(y_test, y_pred))
11 print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(y_test, y_pred)))
12
13 # Plotting predicted vs. actual values
14 sns.scatterplot(x=y_pred, y=y_test)
15 plt.xlabel('Predicted Values')
16 plt.ylabel('Actual Values')
17 plt.title('Predicted vs. Actual Values')
18 plt.show()

```

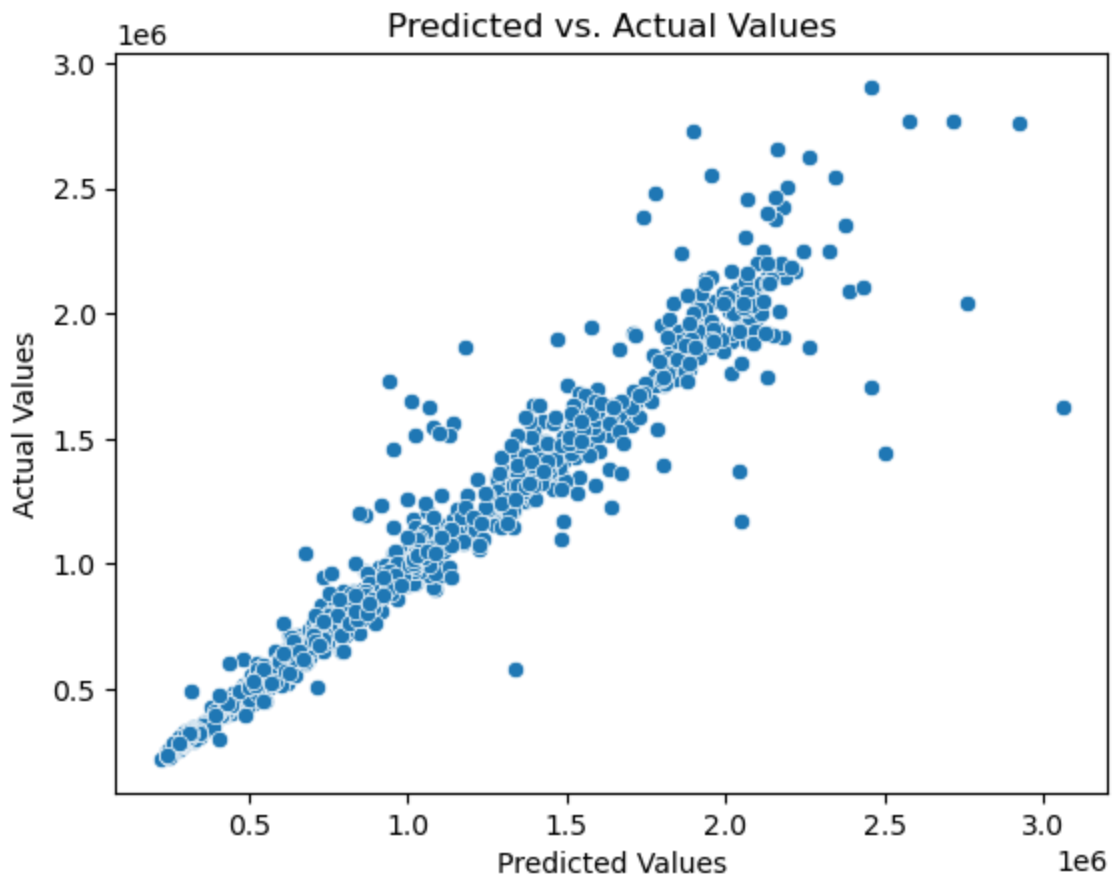
Random Forest Regressor:

Accuracy: 94.66082007163735

Mean Absolute Error: 65704.30705622295

Mean Squared Error: 17231608877.274364

Root Mean Squared Error: 131269.2228866857



In []:

1	
---	--