

```
# Import necessary libraries
import pandas as pd
import seaborn as sns
import numpy as np
import matplotlib.pyplot as plt
from matplotlib import dates
from datetime import datetime
```

```
# Load dataset
data = pd.read_csv('/content/Walmart_Store_sales.csv')
data
```

	Store	Date	Weekly_Sales	Holiday_Flag	Temperature	Fuel_Price	CPI	Une
0	1	05-02-2010	1643690.90	0	42.31	2.572	211.096358	
1	1	12-02-2010	1641957.44	1	38.51	2.548	211.242170	
2	1	19-02-2010	1611968.17	0	39.93	2.514	211.289143	
3	1	26-02-2010	1409727.59	0	46.63	2.561	211.319643	
4	1	05-03-2010	1554806.68	0	46.50	2.625	211.350143	
...
		28-						

DATA PREPARATION

```
# Convert date to datetime format and show dataset information
data['Date'] = pd.to_datetime(data['Date'])
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
```

Double-click (or enter) to edit

```
data.shape
```

```
(6435, 8)
```

```
data.size
```

```
51480
```

```
# checking for missing values
data.isnull().sum()
```

```
Store          0  
Date           0  
Weekly_Sales   0  
Holiday_Flag   0  
Temperature    0  
Fuel_Price     0  
CPI            0  
Unemployment   0  
dtype: int64
```

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

	Store	Date	Weekly_Sales	Holiday_Flag	Temperature	Fuel_Price	CPI	Unemployment
0	1	2010-05-02	1643690.90	0	42.31	2.572	211.096358	2.007
1	1	2010-12-02	1641957.44	1	38.51	2.548	211.242170	2.007
2	1	2010-02-19	1611968.17	0	39.93	2.514	211.289143	2.007
3	1	2010-02-26	1409727.59	0	46.63	2.561	211.319643	2.007
4	1	2010-05-03	1554806.68	0	46.50	2.625	211.350143	2.007

DATA PREPARATION

6430 45 2012-05-05 712472.05 0 64.89 3.985 192.170412

Double-click (or enter) to edit

6431 45 2012-05-05 733455.07 0 64.89 3.985 192.170412

```
# Convert date to datetime format and show dataset information
```

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

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6435 entries, 0 to 6434
Data columns (total 11 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 
 8   Day              6435 non-null    int64  
 9   Month            6435 non-null    int64  
 10  Year             6435 non-null    int64  
dtypes: datetime64[ns](1), float64(5), int64(5)
memory usage: 553.1 KB
```

```
# checking for missing values
data.isnull().sum()
```

Store	0
Date	0
Weekly_Sales	0
Holiday_Flag	0
Temperature	0

```
Fuel_Price      0
CPI            0
Unemployment   0
Day            0
Month          0
Year           0
dtype: int64
```

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

	Store	Date	Weekly_Sales	Holiday_Flag	Temperature	Fuel_Price	CPI	Unemployment
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	

Q1: Which store has minimum and maximum sales?

```
plt.figure(figsize=(15,7))

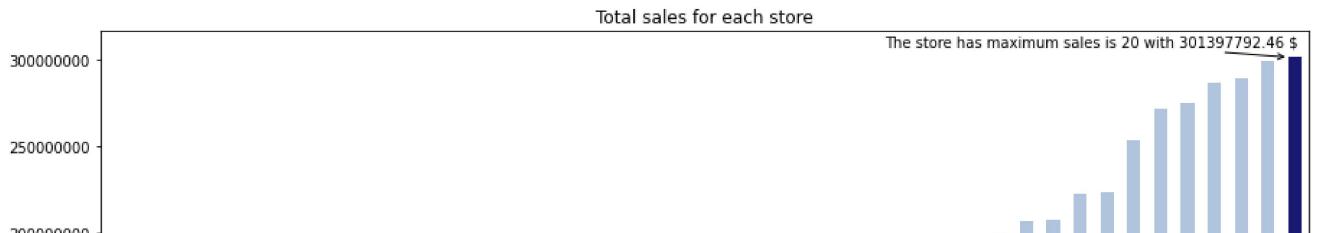
# Sum Weekly_Sales for each store, then sorted by total sales
total_sales_for_each_store = data.groupby('Store')['Weekly_Sales'].sum().sort_values()
total_sales_for_each_store_array = np.array(total_sales_for_each_store)

# CONVERT TO ARRAY

# Assigning a specific color for the stores have the lowest and highest sales
```

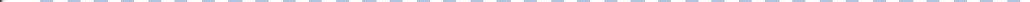
```
clrs = ['lightsteelblue' if ((x < max(total_sales_for_each_store_array)) and  
                           (x > min(total_sales_for_each_store_array))) else 'midnightblue'  
  
ax = total_sales_for_each_store.plot(kind='bar',color=clrs);  
  
# store have minimum sales  
p = ax.patches[0]  
print(type(p.get_height()))  
ax.annotate("The store has minimum sales is 33 with {0:.2f} $".format((p.get_height())),  
xy=(p.get_x(), p.get_height()), xycoords='data',  
xytext=(0.17, 0.32), textcoords='axes fraction',  
arrowprops=dict(arrowstyle="->", connectionstyle="arc3"),  
horizontalalignment='center', verticalalignment='center')  
  
# store have maximum sales  
p = ax.patches[44]  
ax.annotate("The store has maximum sales is 20 with {0:.2f} $".format((p.get_height())),  
xy=(p.get_x(), p.get_height()), xycoords='data',  
xytext=(0.82, 0.98), textcoords='axes fraction',  
arrowprops=dict(arrowstyle="->", connectionstyle="arc3"),  
horizontalalignment='center', verticalalignment='center')  
  
# plot properties  
plt.xticks(rotation=0)  
plt.ticklabel_format(useOffset=False, style='plain', axis='y')  
plt.title('Total sales for each store')  
plt.xlabel('Store')  
plt.ylabel('Total Sales');
```

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



Capital 150000000

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

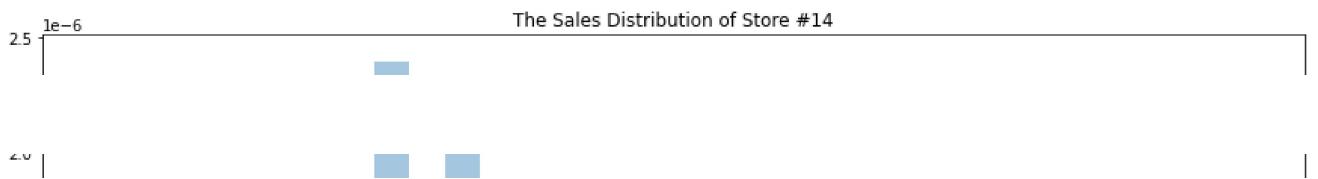
50000000 | 

```
# Which store has maximum standard deviation  
data_std = pd.DataFrame(data.groupby('Store')['Weekly_Sales'].std().sort_values(ascending=False))  
print("The store has maximum standard deviation is "+str(data_std.head(1).index[0])+"  
      with{0:.0f} $".format(data_std.head(1).Weekly_Sales[data_std.head(1).index[0]]))
```

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

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

```
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `d
warnings.warn(msg, FutureWarning)
```



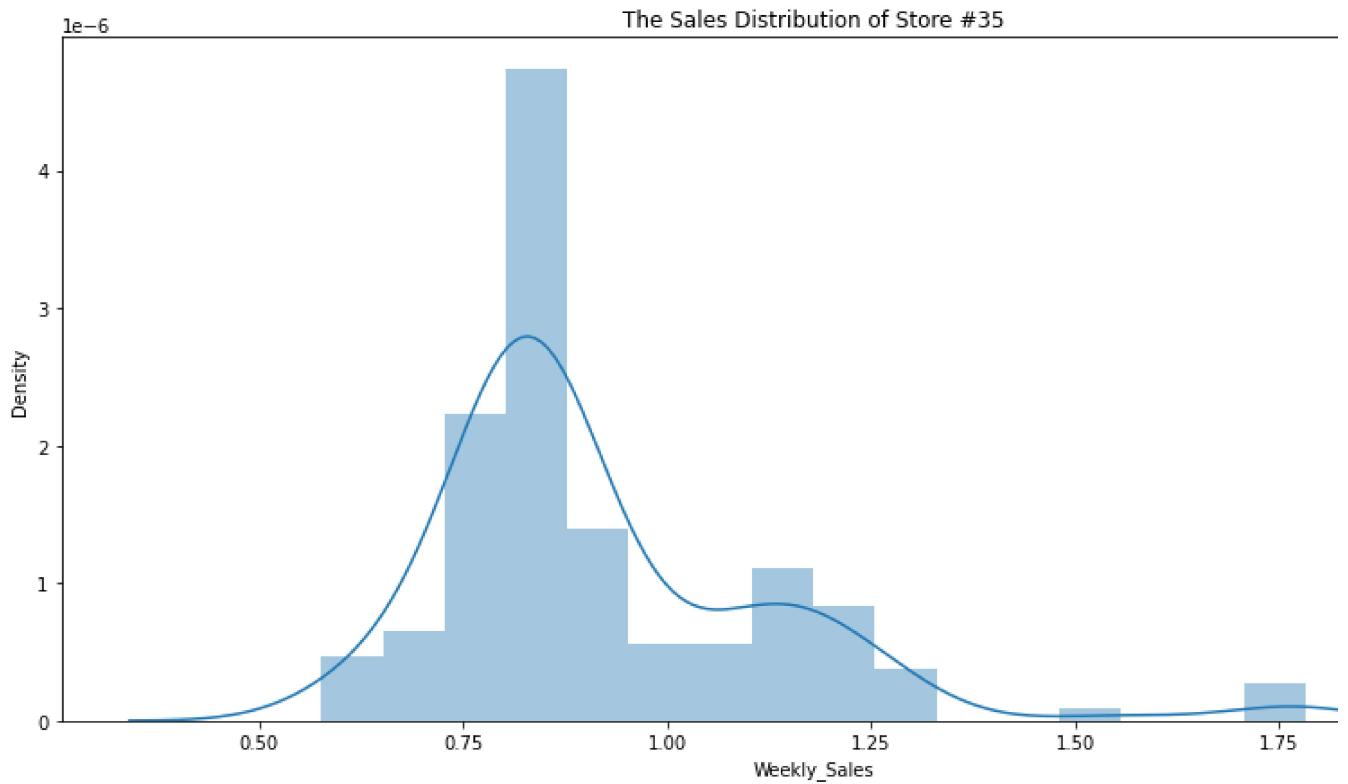
```
# Coefficient of mean to standard deviation
coef_mean_std = pd.DataFrame(data.groupby('Store')['Weekly_Sales'].std() / data.groupby('Stor
coef_mean_std = coef_mean_std.rename(columns={'Weekly_Sales':'Coefficient of mean to standarc
coef_mean_std
```

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	0.170721

```
from logging import warning
# Distribution of store has maximum coefficient of mean to standard deviation
coef_mean_std_max = coef_mean_std.sort_values(by='Coefficient of mean to standard deviation')
plt.figure(figsize=(15,7))
sns.distplot(data[data['Store'] == coef_mean_std_max.tail(1).index[0]]['Weekly_Sales'])
plt.title('The Sales Distribution of Store #' + str(coef_mean_std_max.tail(1).index[0]));
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `d
warnings.warn(msg, FutureWarning)
```



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

```
plt.figure(figsize=(15,7))
```

```
<Figure size 1080x504 with 0 Axes>
<Figure size 1080x504 with 0 Axes>
```

Sales for third quarterly in 2012

```
Q3 = data[(data['Date'] > '2012-07-01') & (data['Date'] < '2012-09-30')].groupby('Store')[['We
```

Double-click (or enter) to edit

Sales for second quarterly in 2012

```
Q2 = data[(data['Date'] > '2012-04-01') & (data['Date'] < '2012-06-30')].groupby('Store')[['We
```

Double-click (or enter) to edit

Plotting the difference between sales for second and third quarterly

```
Q2.plot(ax=Q3.plot('bar',legend=True),kind='bar',color='r',alpha=0.2,legend=True); plt.legend(["Q3'2012", "Q2'2012"]);
```

Double-click (or enter) to edit

```
# store/s has good quarterly growth rate in Q3'2012 - .sort_values(by='Weekly_Sales')
print('Store have good quarterly growth rate in Q3'2012 is Store '+str(Q3.idxmax())+
' With'+ str(Q3.max( ))+' $')
```

Store have good quarterly growth rate in Q3'2012 is Store 4 With25652119.35 \$

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

```
def plot_line(df,holiday_dates,holiday_label):
    fig, ax = plt.subplots(figsize = (15,5))
    ax.plot(df['Date'],df['Weekly_Sales'],label=holiday_label)

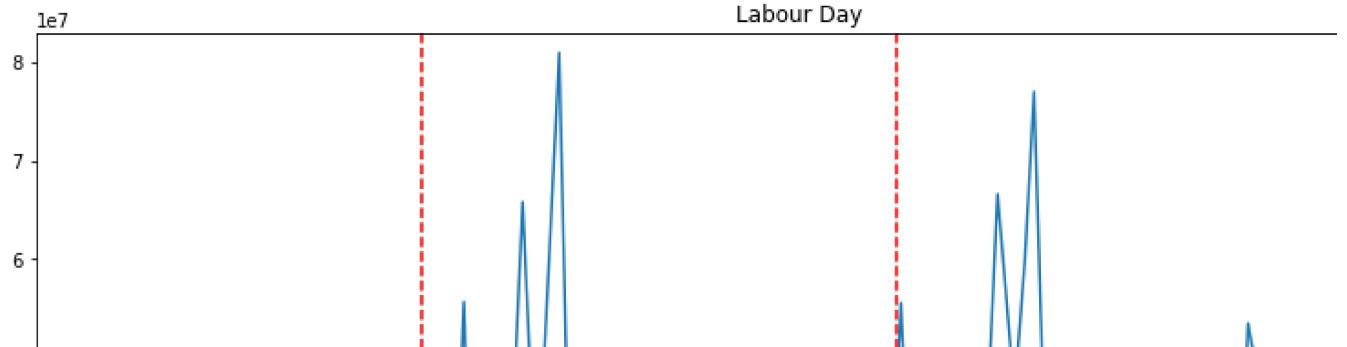
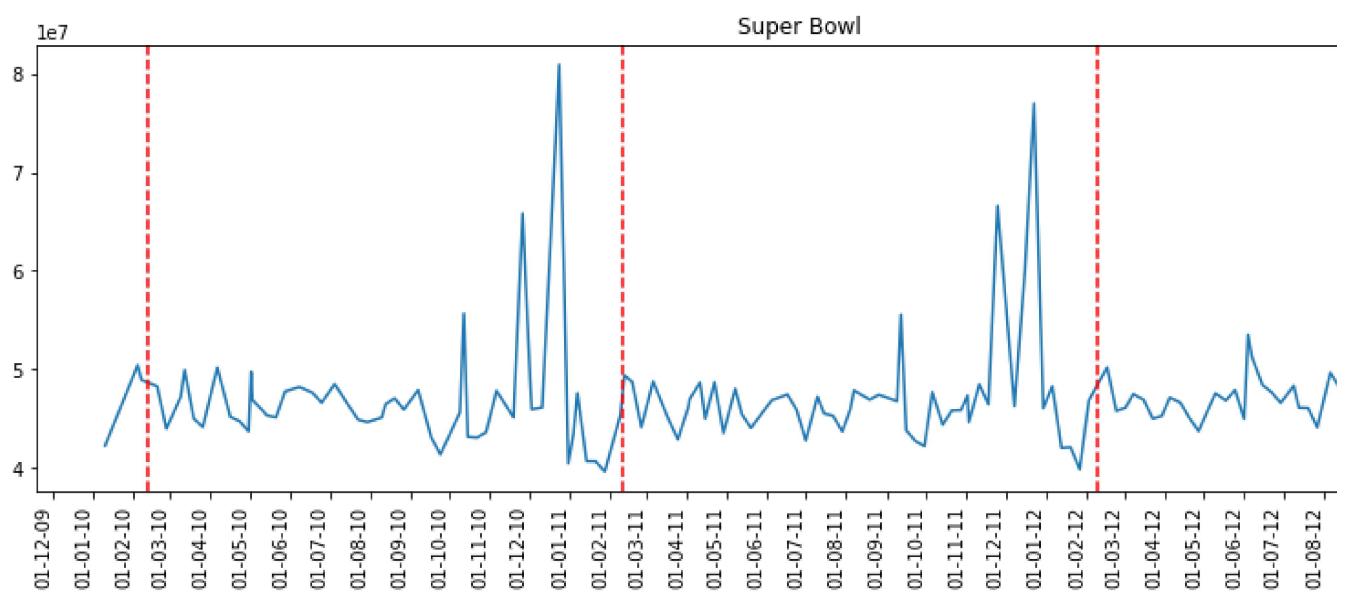
    for day in holiday_dates:
        day = datetime.strptime(day, '%d-%m-%Y')
        plt.axvline(x=day, linestyle='--', c='r')

    plt.title(holiday_label)
    x_dates = df['Date'].dt.strftime('%Y-%m-%d').sort_values().unique()
    xfmt = dates.DateFormatter('%d-%m-%y')
    ax.xaxis.set_major_formatter(xfmt)
    ax.xaxis.set_major_locator(dates.DayLocator(1))
    plt.gcf().autofmt_xdate(rotation=90)
    plt.show()

total_sales = data.groupby('Date')['Weekly_Sales'].sum().reset_index()
Super_Bowl =['12-2-2010', '11-2-2011', '10-2-2012']
Labour_Day = ['10-9-2010', '9-9-2011', '7-9-2012']
Thanksgiving = ['26-11-2010', '25-11-2011', '23-11-2012']
```

```
Christmas = ['31-12-2010', '30-12-2011', '28-12-2012']
```

```
plot_line(total_sales,Super_Bowl,'Super Bowl')
plot_line(total_sales,Labour_Day,'Labour Day')
plot_line(total_sales,Thanksgiving,'Thanksgiving')
plot_line(total_sales,Christmas,'Christmas')
```



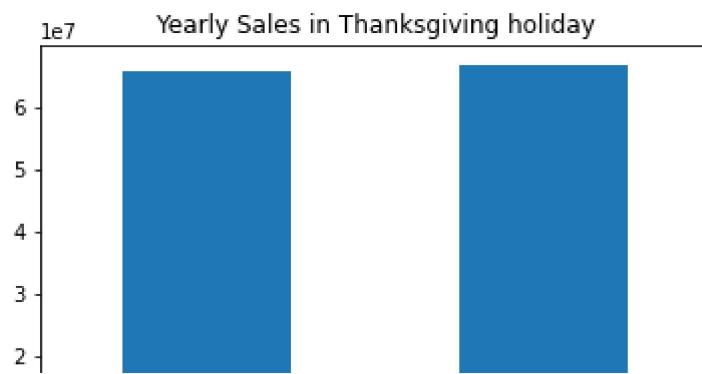
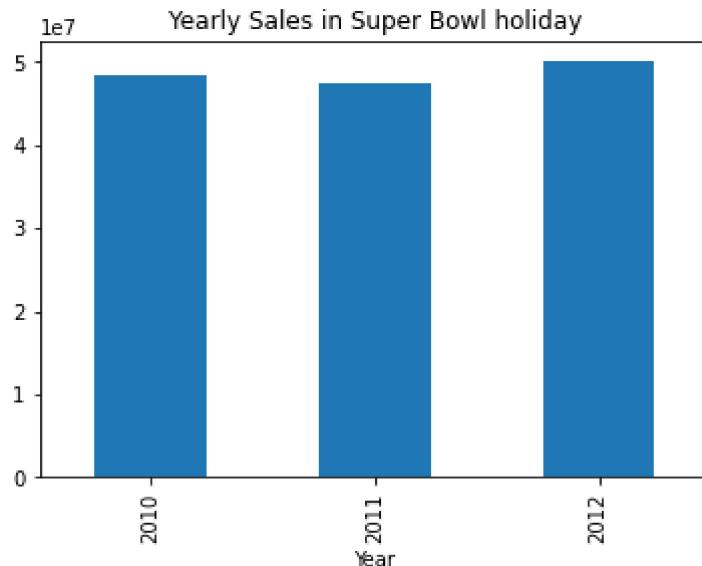
The sales increased during thanksgiving. And the sales decreased during christmas.

4 | Page

```
data.loc[data.Date.isin(Super_Bowl)]
```

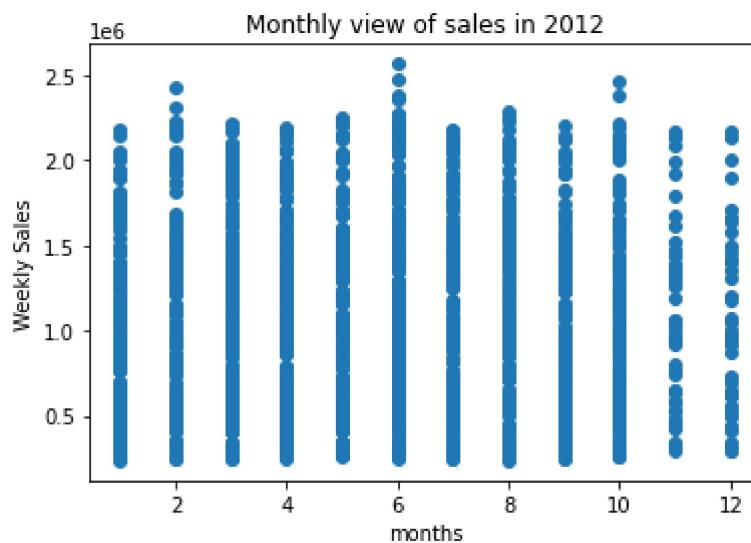
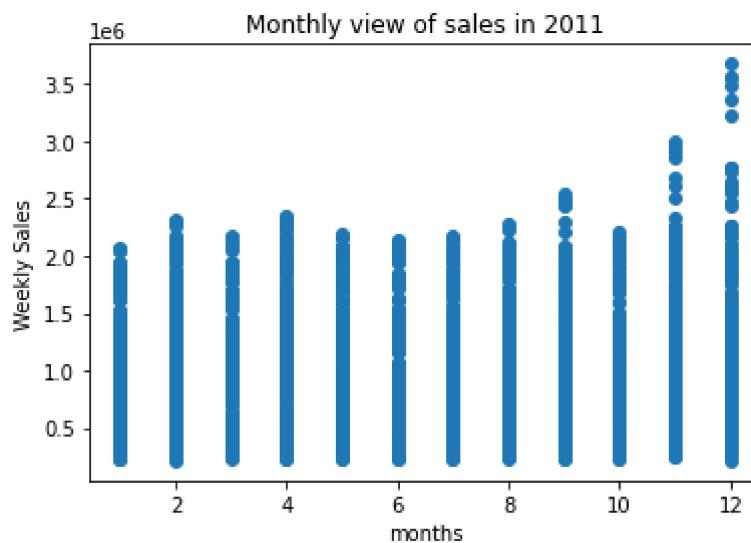
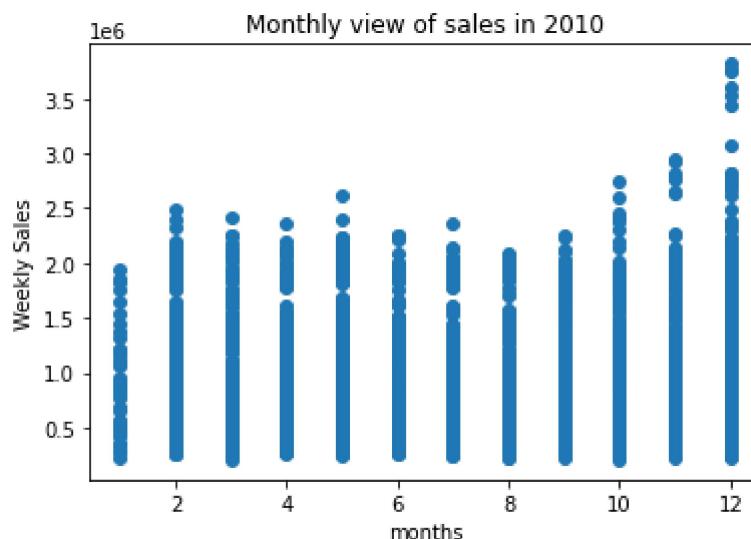
```
Store Date Weekly_Sales Holiday_Flag Temperature Fuel_Price          CPI  Unemployment_Rate  
# Yearly Sales in holidays  
Super_Bowl_df = pd.DataFrame(data.loc[data.Date.isin(Super_Bowl)].groupby('Year')['Weekly_Sales'])  
Thanksgiving_df = pd.DataFrame(data.loc[data.Date.isin(Thanksgiving)].groupby('Year')['Weekly_Sales'])  
Labour_Day_df = pd.DataFrame(data.loc[data.Date.isin(Labour_Day)].groupby('Year')['Weekly_Sales'])  
Christmas_df = pd.DataFrame(data.loc[data.Date.isin(Christmas)].groupby('Year')['Weekly_Sales'])  
  
Super_Bowl_df.plot(kind='bar',legend=False,title='Yearly Sales in Super Bowl holiday')  
Thanksgiving_df.plot(kind='bar',legend=False,title='Yearly Sales in Thanksgiving holiday')  
Labour_Day_df.plot(kind='bar',legend=False,title='Yearly Sales in Labour_Day holiday')  
Christmas_df.plot(kind='bar',legend=False,title='Yearly Sales in Christmas holiday')
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f02090f8f10>
```



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

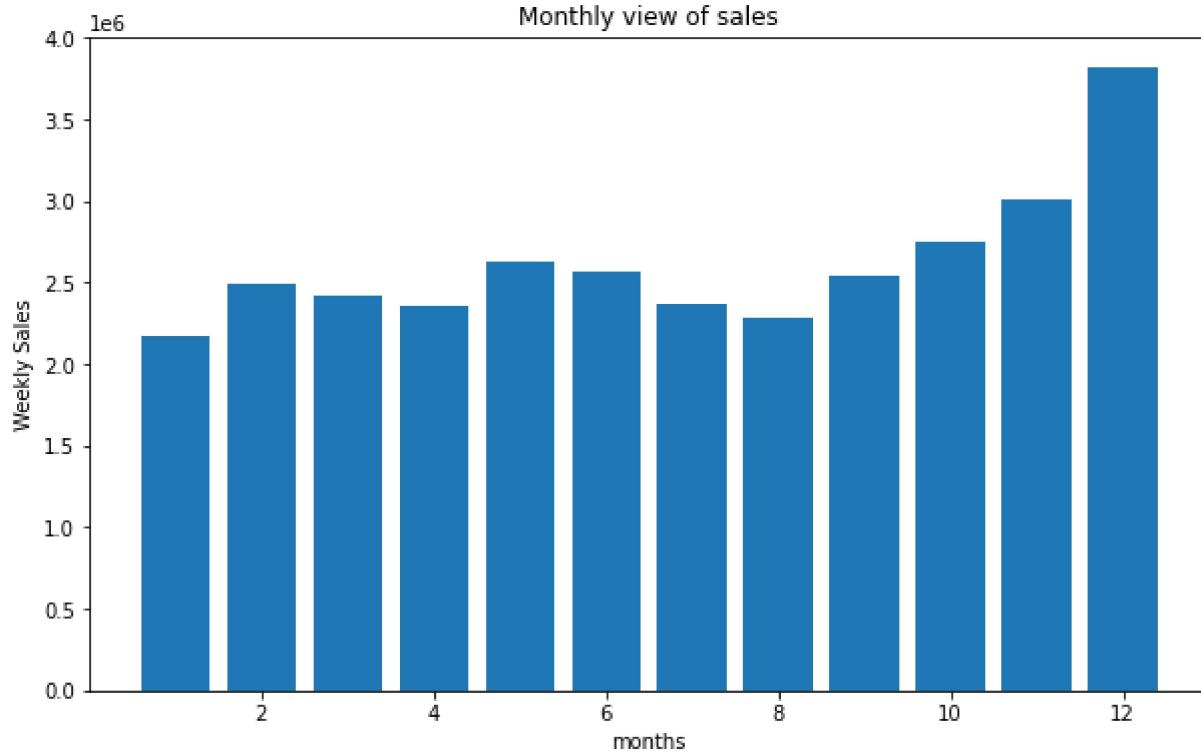
```
# Monthly view of sales for each years  
plt.scatter(data[data.Year==2010]["Month"],data[data.Year==2010]["Weekly_Sales"])  
plt.xlabel("months")  
plt.ylabel("Weekly Sales")  
plt.title("Monthly view of sales in 2010")  
plt.show()  
  
plt.scatter(data[data.Year==2011]["Month"],data[data.Year==2011]["Weekly_Sales"])  
plt.xlabel("months")  
plt.ylabel("Weekly Sales")  
plt.title("Monthly view of sales in 2011")  
plt.show()  
  
plt.scatter(data[data.Year==2012]["Month"],data[data.Year==2012]["Weekly_Sales"])  
plt.xlabel("months")  
plt.ylabel("Weekly Sales")  
plt.title("Monthly view of sales in 2012")  
plt.show()
```



```
# Monthly view of sales for all years
plt.figure(figsize=(10,6))
plt.bar(data["Month"],data["Weekly_Sales"])
plt.xlabel("months")
```

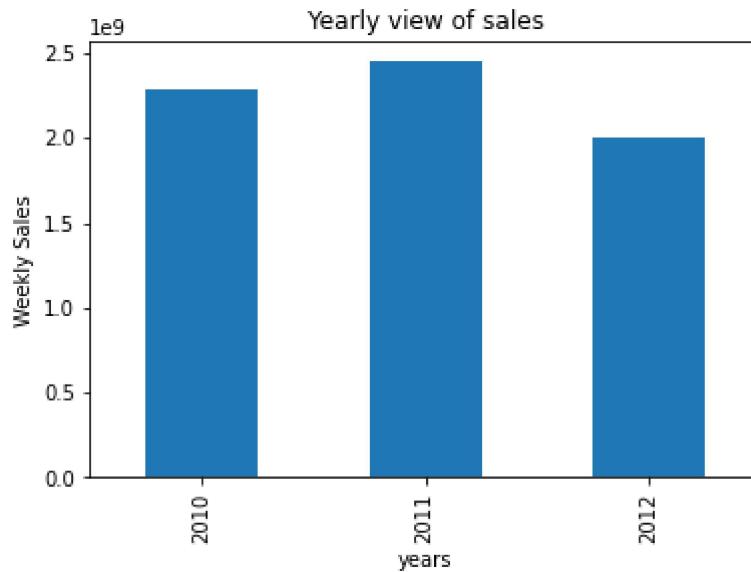
```
plt.ylabel("Weekly Sales")
plt.title("Monthly view of sales")
```

```
Text(0.5, 1.0, 'Monthly view of sales')
```



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

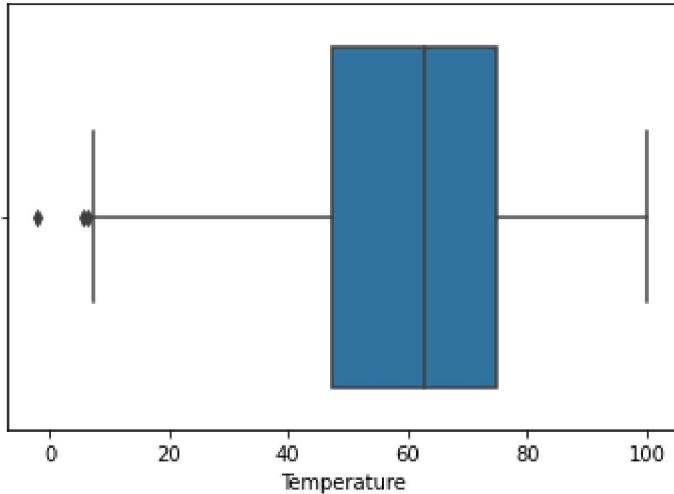
```
<Figure size 720x432 with 0 Axes>
```



Build prediction models to forecast demand (Modeling)

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

```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass t
  FutureWarning
```



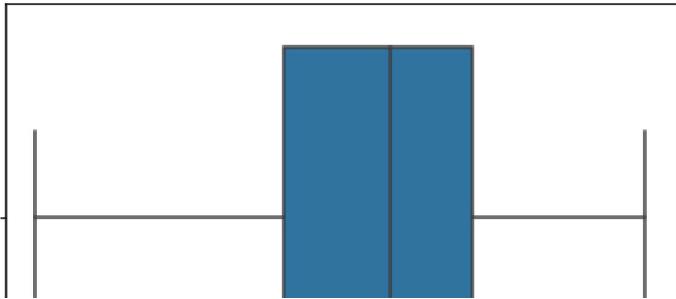
```
# drop the outliers
data_new = data[(data['Unemployment'] < 10) & (data['Unemployment'] > 4.5) & (data['Temperature'] >
data_new
```

	Store	Date	Weekly_Sales	Holiday_Flag	Temperature	Fuel_Price	CPI	Unemployment
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	

```
# check outliers
fig, axs = plt.subplots(4, figsize=(6,18))
```

```
X = data_new[['Temperature', 'Fuel_Price', 'CPI', 'Unemployment']]  
for i,column in enumerate(X):  
    sns.boxplot(data_new[column], ax=axs[i])
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass t
  FutureWarning
```



| | | |

Build Model

 | | | |

```
# Import sklearn
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import train_test_split
from sklearn import metrics
from sklearn.linear_model import LinearRegression
| |      | | | |
```

Select features and target
X = data_new[['Store','Fuel_Price','CPI','Unemployment','Day','Month','Year']]
y = data_new['Weekly_Sales']

Split data to train and test (0.80:0.20)
X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.2)
| |

Linear Regression model
print('Linear Regression:')
print()
reg = LinearRegression()
reg.fit(X_train, y_train)
y_pred = reg.predict(X_test)
print('Accuracy:',reg.score(X_train, y_train)*100)

print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, y_pred))
print('Mean Squared Error:', metrics.mean_squared_error(y_test, y_pred))
print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(y_test, y_pred)))

sns.scatterplot(y_pred, y_test);