

# Retail Analysis with Walmart Data

DESCRIPTION :

One of the leading retail stores in the US, Walmart, would like to predict the sales and demand accurately. There are certain events and holidays which impact sales on each day. There are sales data available for 45 stores of Walmart. The business is facing a challenge due to unforeseen demands and runs out of stock some times, due to the inappropriate machine learning algorithm. An ideal ML algorithm will predict demand accurately and ingest factors like economic conditions including CPI, Unemployment Index, etc.

Walmart runs several prominent markdown events throughout the year. These markdowns precede prominent holidays, the four largest of all, which are the Super Bowl, Labor Day, Thanksgiving, and Christmas. The weeks including these holidays are weighted five times higher in the evaluation than non-holiday weeks. Part of the challenge presented by this competition is modeling the effects of markdowns on these holiday weeks in the absence of comprehensive historical data. Historical sales data for 45 Walmart stores located in different regions are available.

**Holiday Events :**

1. **Super Bowl** : 12-Feb-10, 11-Feb-11, 10-Feb-12, 8-Feb-13
2. **Labour Day** : 10-Sep-10, 9-Sep-11, 7-Sep-12, 6-Sep-13
3. **Thanksgiving** : 26-Nov-10, 25-Nov-11, 23-Nov-12, 25-Nov-13
4. **Christmas** : 13-Dec-10, 30-Dec-11, 26-Dec-12, 27-Dec-13

**Analysis Tasks:**

Basic Statistics Tasks :

- Which store has maximum sales.
- Which store has maximum standard deviation i.e., the sales vary a lot. Also, find out the coefficient of mean to standard deviation
- Which stores has good quarterly growth rate in Q3/2012
- 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
- Provide a monthly and semester view of sales in units and give insights.

**Statistical Model :**

For Store 1 – Build prediction models to forecast demand

- Linear Regression – (Like variables like date and restructuring dates as 1 for 5 Feb 2010 (starting from the earliest date in order). Hypothesize if CPI, Unemployment, and fuel price have any impact on sales.
- Change dates into days by creating new variable.
- Select the model which gives best accuracy.

```
In [11]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import matplotlib

import scipy.stats as stats
import sklearn as sk
import sklearn.feature as fe

from datetime import datetime as dt
from datetime import timedelta as date

import warnings
warnings.filterwarnings('ignore')
```

```
In [12]: store = pd.read_csv('walmart_store_sales.csv')
```

```
Out[12]:
```

	Store	Date	Weekly_Sales	Holiday_Flag	Temperature	Fuel_Price	CPI	Unemployment
0	1	2010-02-09	164969.00	0	42.31	2.572	211.060560	8.106
1	1	2010-02-10	1641967.44	1	38.51	2.548	211.242170	8.106
2	1	18-02-2010	1611868.17	0	39.85	2.534	211.280143	8.106
3	1	26-02-2010	146077.59	0	46.63	2.563	211.318643	8.106
4	1	05-03-2010	155400.00	0	46.50	2.629	211.380143	8.106

**EDA**

```
In [13]: data.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 6425 entries, 0 to 6424
Data columns (total 8 columns):
0      Store      Non-Holiday_Count      dtype: int64
1      Date      Date      dtype: object
2      Weekly_Sales      dtype: float64
3      Holiday_Flag      dtype: int64
4      Temperature      dtype: float64
5      Fuel_Price      dtype: float64
6      CPI      dtype: float64
7      Unemployment      dtype: float64
dtypes: float64(6), int64(1), object(1)
memory usage: 492.3+ KB

• Charge date to datetime format
• Weekly_Sales is target variable here
• No missing values.
• Shape is (6425, 8)
```

```
In [14]: # Date Formatting :
data['Date'] = pd.to_datetime(data['Date'])
data.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 6425 entries, 0 to 6424
Data columns (total 8 columns):
0      Store      Non-Holiday_Count      dtype: int64
1      Date      Date      dtype: object
2      Weekly_Sales      dtype: float64
3      Holiday_Flag      dtype: int64
4      Temperature      dtype: float64
5      Fuel_Price      dtype: float64
6      CPI      dtype: float64
7      Unemployment      dtype: float64
dtypes: datetime(1), float64(5), int64(2)
memory usage: 492.3+ KB
```

```
In [15]: data['Holiday_Flag'].value_counts()

Out[15]:
0      5985
485    Holiday_Flag      dtype: int64
```

```
In [16]: data['DateType'] = [dt.strftime(date, '%Y-%m-%d') for date in data['Date'].astype(str).values.tolist()]
```

**Holidays :**

```
In [17]: data['SuperBowl'] = np.where((data['DateType'] == dt(2010, 2, 12).date()) | (data['DateType'] == dt(2011, 2, 13).date()))
data['LabourDay'] = np.where((data['DateType'] == dt(2010, 9, 30).date()) | (data['DateType'] == dt(2011, 9, 0).date()) | (data['DateType'] == dt(2012, 9, 0).date()) | (data['DateType'] == dt(2013, 9, 0).date()))
data['Christmas'] = np.where((data['DateType'] == dt(2010, 12, 21).date()) | (data['DateType'] == dt(2011, 12, 20).date()) | (data['DateType'] == dt(2012, 12, 27).date()) | (data['DateType'] == dt(2013, 12, 26).date()))
data['Thanksgiving'] = np.where((data['DateType'] == dt(2010, 11, 26).date()) | (data['DateType'] == dt(2011, 11, 25).date()) | (data['DateType'] == dt(2012, 11, 23).date()) | (data['DateType'] == dt(2013, 11, 28).date()))
```

```
In [18]: print(data.SuperBowl.value_counts())
print(data.LabourDay.value_counts())
print(data.Christmas.value_counts())
print(data.Thanksgiving.value_counts())

0      6425
Name: SuperBowl, dtype: int64
0      6395
1      435
Name: LabourDay, dtype: int64
0      6345
1      8345
Name: Thanksgiving, dtype: int64
0      6345
1      49
Name: Christmas, dtype: int64
```

**1. Which store has maximum sales**

```
In [19]: store_sales = data.groupby(['Store'])['Weekly_Sales'].sum().sort_values(ascending = False)
round(store_sales, 1).head()
```

```
Out[19]:
Store
0      301297792.3
4      29643903.4
14     28666611.3
13     28637760.8
27     27138641.8
Name: Weekly_Sales, dtype: float64
```

**So, the store 0 has 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**

```
In [110]: store_std = data.groupby(['Store'])['Weekly_Sales'].std().sort_values(ascending = False)
round(store_std, 2).head()
```

```
Out[110]:
Store
14      317669.95
10     305282.46
2      292051.44
6      265051.44
5      262051.44
Name: Weekly_Sales, dtype: float64
```

```
In [111]: store_mean = data.groupby(['Store'])['Weekly_Sales'].mean().sort_values(ascending = False)
coeff_variance = round(store_std / store_mean, 2)
coeff_variance.sort_values(ascending = False).head()
```

```
Out[111]:
Store
0      0.23
35     0.23
15     0.19
23     0.18
29     0.18
Name: Weekly_Sales, dtype: float64
```

**Insights :**

- **coeff\_variance** gives the coefficient of variance w.r.t particular stores of the dataset, as shown above.
- **Store 14** has the maximum standard deviation.

**3. Which store/s has good quarterly growth rate in Q3/2012**

```
In [122]: data['Year'] = data['Date'].dt.year
data['Month'] = data['Date'].dt.month
data['Quarter'] = data['Date'].dt.quarter
data.drop(columns = ['DateType'], inplace = True)
data.head(3)
```

```
Out[122]:
```

	Store	Date	Weekly_Sales	Holiday_Flag	Temperature	Fuel_Price	CPI	Unemployment	SuperBowl	Labour_Day	Christmas	Thanksgiving	Year	Month	Quarter
0	1	2010-02-09	164969.00	0	42.31	2.572	211.060560	8.106	0	0	0	0	2010	2	1
1	1	2010-02-10	1641967.44	1	38.51	2.548	211.242170	8.106	0	0	0	0	2010	12	4
2	1	2010-02-19	1611868.17	0	39.85	2.534	211.280143	8.106	0	0	0	0	2010	12	4

```
In [123]: q3 = data[(data['Quarter'] == 3) & (data['Year'] == 2012)].groupby(['Store'])['Weekly_Sales'].sum().sort_values(ascending = False)
q3.head(3)
```

```
Out[123]:
Store
20     2565219.35
20     2466598.11
13     2451894.36
Name: Weekly_Sales, dtype: float64
```

**Insight :**

- Store 4 has good quarterly growth rate in Q3/2012 with maximum profit of 2565219.35 compared with other stores.

**Growth Rate :**

- Growth rate formula is defined as the ratio of difference in present value to past value by past value whole multiplied with 100 (since it is in percentage)
- (Present value – Past value) / Past value \* 100

```
In [143]: Q2_date = pd.Timestamp(date(2012, 7, 1))
Q3_date_to = pd.Timestamp(date(2012, 9, 30))
Q2_date_from = pd.Timestamp(date(2012, 4, 1))
Q3_date_to = pd.Timestamp(date(2012, 6, 30))

#Collecting the data of Q2 and Q3 from original dataset:
Q2data = data[(data['Date'] > Q2_date_from & (data['Date'] <= Q2_date_to)]
Q3data = data[(data['Date'] > Q3_date_from & (data['Date'] <= Q3_date_to)]

#Finding the sum weekly sales of each store in Q2
Q2 = pd.DataFrame(Q2data.groupby('Store')['Weekly_Sales'].sum())
Q2.reset_index(inplace=True)
Q2.columns = ['Weekly_Sales', 'Q2.Weekly_Sales', 'inplace = True]

#Finding the sum weekly sales of each store in Q3
Q3 = pd.DataFrame(Q3data.groupby('Store')['Weekly_Sales'].sum())
Q3.reset_index(inplace=True)
Q3.columns = ['Weekly_Sales', 'Q3.Weekly_Sales', 'inplace = True]

#Merging Q2 and Q3 data as Store as a common column
Q3_Growth_Q2_merged[Q3_how = 'inner', on = 'Store']

#Calculating growth rate of each store and collecting it into a dataframe
Q3_Growth['Growth_Rate'] = (Q3_Growth['Q3.Weekly_Sales'] - Q3_Growth['Q2.Weekly_Sales']) / Q3_Growth['Q2.Weekly_Sales']
Q3_Growth['Growth_Rate'] = round(Q3_Growth['Growth_Rate'], 2)
Q3_Growth.sort_values('Growth_Rate', ascending = False).head(5)
```

```
Out[143]:
```

	Store	Q2_Weekly_Sales	Q3_Weekly_Sales	Growth_Rate
15	15	6026133.44	6441311.11	-0.03
7	7	7612003.02	7222933.02	-0.04
35	35	1070870.87	1050212.69	-0.05

```
In [145]: Q3_Growth.sort_values('Growth_Rate', ascending = False).tail(5)
```

```
Out[145]:
```

	Store	Q2_Weekly_Sales	Q3_Weekly_Sales	Growth_Rate
35	35	4900278.50	3576123.50	-0.23
45	45	1027000.16	683242.22	-0.34
14	14	14477799.05	2004369.40	-0.38

**Insights :**

- From above information we can say that Q3 growth rate is in loss.
- Store 10 has the least loss of 0% compared with other stores.
- Store 14 has the highest loss rate of 38%.

**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**

```
In [146]: round(data.groupby(['Holiday_Flag'])['Weekly_Sales'].sum(), 1)
```

```
Out[146]:
Holiday_Flag
0      6.251935e+09
1      6.052998e+08
Name: Weekly_Sales, dtype: float64
```

```
In [147]: print('0', 6.251935e+09)
print('1', 6.052998e+08)

6.251935e+09
1: 605299800.0
```

```
In [148]: Spr_sales = data.groupby(['SuperBowl'])['Weekly_Sales'].mean()
Ld_sales = data.groupby(['Labour_Day'])['Weekly_Sales'].mean()
Thksg_sales = data.groupby(['Thanksgiving'])['Weekly_Sales'].mean()
Christmas_sales = data.groupby(['Christmas'])['Weekly_Sales'].mean()
Non_Holi_Sales = data[(data['Holiday_Flag'] == 0)].groupby('Holiday_Flag')['Weekly_Sales'].mean()
```

```
print(round(Spr_sales, 2))
print(round(Ld_sales, 2))
print(round(Thksg_sales, 2))
print(round(Christmas_sales, 2))
print(round(Non_Holi_Sales, 2))
```

```
SuperBowl
Mean: 145464.88
Name: Weekly_Sales, dtype: float64
Labour_Day
Mean: 145479.48
1: 1031242.83
Name: Weekly_Sales, dtype: float64
Thanksgiving
Mean: 147179.48
Name: Weekly_Sales, dtype: float64
Christmas
Mean: 184186.49
1: 908653.11
Name: Weekly_Sales, dtype: float64
Holiday_Flag
Mean: 145126.38
Name: Weekly_Sales, dtype: float64
```

**Visualizing**

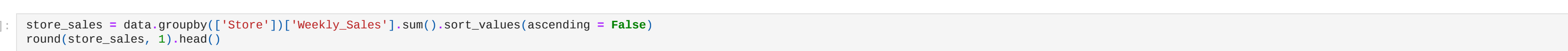
**1. Sales in Super Bowl holiday**

```
In [149]: plt.style.use('seaborn-whitegrid')
plt.figure(figsize=(10, 6))
Spr_sales.plot(kind = 'bar', legend = False, title = 'Sales in Super Bowl holiday', color = 'b18e76')
plt.show()
```



**2. Sales in Labour Day holiday**

```
In [150]: plt.figure(figsize=(10, 6))
Ld_sales.plot(kind = 'bar', legend = False, title = 'Sales in Labour Day holiday', color = 'plum', lavender))
plt.show()
```



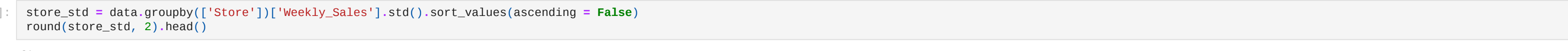
**3. Sales in Thanksgiving holiday**

```
In [151]: plt.figure(figsize=(10, 6))
Thksg_sales.plot(kind = 'bar', legend = False, title = 'Sales in Thanksgiving holiday', color = 'palevioletred', skyblue))
plt.show()
```



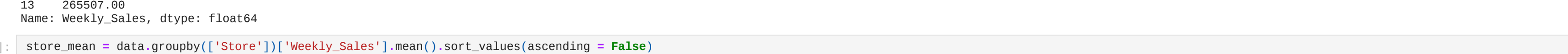
**4. Sales in Christmas holiday**

```
In [152]: plt.figure(figsize=(10, 6))
Christmas_sales.plot(kind = 'bar', legend = False, title = 'Sales in Christmas holiday', color = 'plum', lavender))
plt.show()
```



**5. Non Holiday Sales.**

```
In [153]: plt.figure(figsize=(10, 6))
Non_Holi_Sales.plot(kind = 'bar', legend = False, title = 'Sales in Labour Day holiday', color = 'lavender')
plt.show()
```



**5. Provide a monthly and semester view of sales in units and give insights.**

**I. Monthly sales view**

```
In [154]: monthly = data.groupby(pd.Grouper(key = 'Date', freq = 'MS')).sum()
monthly = monthly.reset_index()
fig, ax = plt.subplots(figsize=(13, 6))
ax = monthly['Date']
x = monthly['Weekly_Sales']
plt.plot(x, color = 'blue')
plt.title('Month Wise Sales')
plt.xlabel('monthly')
plt.ylabel('Weekly Sales')
plt.legend(['Sales'])
plt.show()
```



**II. Semester wise sales view**

```
In [155]: Semester = data.groupby(pd.Grouper(key = 'Date', freq = 'MS')).sum()
Semester = Semester.reset_index()
fig, ax = plt.subplots(figsize=(13, 6), dpi = 80)
x = Semester['Date']
y = Semester['Weekly_Sales']
plt.plot(x, y)
plt.title('Semester Wise Sales')
plt.xlabel('Semester')
plt.ylabel('Semester')
plt.ylabel('Weekly Sales')
plt.legend(['Sales'])
plt.show()
```



**Insights :**

- We can infer that there's a big spike in sales from **February 2010 to February 2011**. Exactly for one year we can say.
- Then spike goes **bit down** in February 2011 after that again there are low ups-downs in further.
- From **August-2012** sales goes **down**. We can acknowledge that there is less in sales.

**For Store 1 – Build prediction models to forecast demand**

```
In [156]: data.drop(columns = ['SuperBowl', 'Labour_Day', 'Christmas', 'Thanksgiving'], inplace = True)
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.preprocessing import LabelEncoder
```

```
Out[156]:
```

```
In [157]: data['Store']

Out[157]:
0      1
1      1
2      1
3      1
4      1
5      1
6      1
7      1
8      1
9      1
10     1
11     1
12     1
13     1
14     1
15     1
16     1
17     1
18     1
19     1
20     1
21     1
22     1
23     1
24     1
25     1
26     1
27     1
28     1
29     1
30     1
31     1
32     1
33     1
34     1
35     1
36     1
37     1
38     1
39     1
40     1
41     1
42     1
43     1
44     1
45     1
Name: Store, dtype: int64
```

```
In [158]: data['Store'] = data['Store'].astype(str)
data['Store'] = 'Store_' + data['Store'].astype(str)
```

```
Out[158]:
```

```
In [159]: Store

Out[159]:
0      Store_1
1      Store_1
2      Store_1
3      Store_1
4      Store_1
5      Store_1
6      Store_1
7      Store_1
8      Store_1
9      Store_1
10     Store_1
11     Store_1
12     Store_1
13     Store_1
14     Store_1
15     Store_1
16     Store_1
17     Store_1
18     Store_1
19     Store_1
20     Store_1
21     Store_1
22     Store_1
23     Store_1
24     Store_1
25     Store_1
26     Store_1
27     Store_1
28     Store_1
29     Store_1
30     Store_1
31     Store_1
32     Store_1
33     Store_1
34     Store_1
35     Store_1
36     Store_1
37     Store_1
38     Store_1
39     Store_1
40     Store_1
41     Store_1
42     Store_1
43     Store_1
44     Store_1
45     Store_1
Name: Store, dtype: object
```

```
In [160]: labelEncoder = LabelEncoder()
store_1 = data[data['Store'] == 'Store_1']
store_1['X'] = store_1.drop('Store', axis=1)
```

```
Out[160]:
```

	Store	Date	Weekly_Sales	Holiday_Flag	Temperature	Fuel_Price	CPI	Unemployment	Year	Month	Quarter
0	Store_1	2010-02-09	164969.00	0	42.31	2.572	211.060560	8.106	2010	2	1
1	Store_1	2010-02-10	1641967.44	1	38.51	2.548	211.242170	8.106	2010	12	4
2	Store_1	2010-02-19	1611868.17	0	39.85	2.534	211.280143	8.106	2010	12	4
3	Store_1	2010-02-26	146077.59	0	46.63	2.563	211.318643	8.106	2010	2	1
4	Store_1	2010-03-05	155400.00	0	46.50	2.629	211.380143	8.106	2010	2	1

```
In [161]: store_1['Date'] = labelEncoder.fit_transform(store_1['Date'])
store_1['Store'] = store_1['Store'].astype('category').cat.rename_categories(['Year', 'Month', 'Quarter'])
store_1['X'] = store_1.drop('Store', axis=1)
```

```
Out[161]:
```

```
In [162]: import seaborn as sns
corr = store_1.corr()
plt.figure(figsize=(13, 6), dpi = 80)
corrmap = sns.heatmap(store_1.corr(), cmap = 'PuBuGn_r', annot = True)
corrmap
```



Hypothesize if CPI, unemployment, and fuel price have any impact on sales.

**Insights :**

- 1. As we can see **unemployment** is highly correlated with **days** and **sales** as correlation with **Weekly Sales** is quite low.
- 2. Also **temperature** and **unemployment** are negatively impacting the sales.
- 3. However **Fuel Price** and **CPI** are positively impacting the Sales.

**Model Building : Linear Regression**

```
In [163]: from sklearn.model_selection import train_test_split

# Store_1 'Date', 'Fuel_Price', 'CPI', 'Unemployment'
y = store_1['Weekly_Sales']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 123)
```

```
print('Shape of X_train :', X_train.shape)
print('Shape of y_train :', y_train.shape)
print('Shape of X_test :', X_test.shape)
print('Shape of y_test :', y_test.shape)
Shape of X_train : (114, 4)
Shape of y_train : (114,)
Shape of X_test : (29, 4)
Shape of y_test : (29,)
```

```
In [164]: from sklearn.linear_model import LinearRegression
linear_reg = LinearRegression()
linear_reg.fit(X_train, y_train)
```

```
Out[164]:
```

```
In [165]: y_pred = linear_reg.predict(X_test)

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)))
Mean Absolute Error: 284768.4857773899
Mean Squared Error: 236965885.154665
Root Mean Squared Error: 148118.41554546286
```

```
In [166]: coefficients = linear_reg.coef_()

Coefficients:
[ 712.88276817 -97286.78823978 19976.81518442 125161.27927183]
```

```
In [167]: print('Variance score: %2f' % linear_reg.score(X_test, y_test))
Variance score: 0.16
```

**Note:**

- From above I can infer that Linear Regression model performs very poorly on our dataset.
- As we can see the variance is in negative numbers, which indicates that our model is Poor.
- Minimum variance helps our model to be **ideal & accurate**. Therefore further down I'll implement **RandomForest** algorithm and check for accuracy.
- Accuracy increases, then I will pick the model.
- Another way we can check accuracy by showing comparison between actual values and predicted values.

**RandomForestRegressor**

```
In [168]: rf = RandomForestRegressor(n_estimators = 488, max_depth = 15)
rf.fit(X_train, y_train)
y_pred = rf.predict(X_test)
```

```
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)))
print('Variance score: %2f' % rf.score(X_test, y_test))
Mean Absolute Error: 284768.4857773899
Mean Squared Error: 236965885.154665
Root Mean Squared Error: 148118.41554546286
Variance score: 0.14
```

```
In [169]: Actual_vs_Pred = pd.DataFrame({'Actual Sales': y_test, 'Predicted Sales': y_pred})
Actual_vs_Pred.head()
```

```
Out[169]:
```

	Actual Sales	Predicted Sales
126	1517428.87	1.58190e+06
138	1497954.70	1.51820e+06
118	1021031.70	1.57220e+06
45	3851034.93	1.59150e+06
1	1507401.60	1.55280e+06

**Predicted Demand for Store 1**

```
In [141]: # checking errors
```