

# 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 promotional markdown events throughout the year. These markdowns precede prominent holidays, the four largest of all, which are the Super Bowl, Labour 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 complete/historical data. Historical sales data for 45 Walmart stores located in different regions are available.

## Dataset Description

This is the historical data that covers sales from 2010-02-05 to 2012-11-01, in the file Walmart\_Store\_sales. Within this file you will find the following fields:

- 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

## 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

## 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 store 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 – Utilize variables like date and restructure 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.

```
In [1]: #Importing necessary libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from matplotlib import style
import matplotlib inline
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn import metrics
from statsmodels.formula.api import ols
import datetime as dt
import warnings

In [2]: #read the dataset using pandas dataframe
walmart_df = pd.read_csv('Walmart_Store_sales.csv')
warnings.filterwarnings('ignore')

In [3]: #to display first 5 rows
walmart_df.head()

Out[3]:
   Store  Date  Weekly_Sales  Holiday_Flag  Temperature  Fuel_Price  CPI  Unemployment
0      1  05-02-2010  1643690.90          0          42.31      2.572  211.096358      8.106
1      1  12-02-2010  1641957.44          1          38.51      2.548  211.242170      8.108
2      1  19-02-2010  1611968.17          0          39.93      2.514  211.289143      8.106
3      1  26-02-2010  1409727.59          0          46.63      2.561  211.319643      8.106
4      1  05-03-2010  1554806.68          0          46.50      2.625  211.350143      8.108

In [4]: #to display last 5 rows
walmart_df.tail()

Out[4]:
   Store  Date  Weekly_Sales  Holiday_Flag  Temperature  Fuel_Price  CPI  Unemployment
6430  45  28-08-2012  713173.95          0          64.88      3.997  192.013558      8.684
6431  45  05-10-2012  733455.07          0          64.89      3.985  192.170412      8.667
6432  45  12-10-2012  734464.36          0          54.47      4.000  192.327265      8.667
6433  45  19-10-2012  718125.53          0          56.47      3.969  192.330854      8.667
6434  45  26-10-2012  760281.43          0          58.85      3.882  192.308899      8.667

In [5]: #info() is used to complete description of data
walmart_df.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        object
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: float64(5), int64(2), object(1)
memory usage: 402.3+ KB

In [6]: #to get column names
walmart_df.columns

Out[6]:
Index(['Store', 'Date', 'Weekly_Sales', 'Holiday_Flag', 'Temperature',
       'Fuel_Price', 'CPI', 'Unemployment'],
      dtype='object')

In [7]: #generalizing column names by lowercasing
walmart_df.columns = walmart_df.columns.str.lower()
walmart_df.columns

Out[7]:
Index(['store', 'date', 'weekly_sales', 'holiday_flag', 'temperature',
       'fuel_price', 'cpi', 'unemployment'],
      dtype='object')

In [8]: #looking for null values in dataset
walmart_df.isnull().sum()

Out[8]:
store          0
date           0
weekly_sales   0
holiday_flag   0
temperature    0
fuel_price     0
cpi            0
unemployment   0
dtype: int64

In [ ]:
```

## Maximum sales

```
In [9]: # groupby stores and get total sales
store_totalweeklysales = walmart_df.groupby('store')['weekly_sales'].sum()
store_totalweeklysales.to_frame()

Out[9]:
   weekly_sales
store
1  2.224029e+08
2  2.753824e+08
3  3.586714e+07
4  2.995440e+08
5  4.545409e+07
6  2.237561e+08
7  8.159828e+07
8  1.299512e+08
9  7.778922e+07
10 2.716177e+08
11 1.939628e+08
12 1.442872e+08
13 2.865177e+08
14 2.899990e+08
15 8.913688e+07
16 7.425243e+07
17 1.277821e+08
18 1.551147e+08
19 2.068349e+08
20 3.013978e+08
21 1.081790e+08
22 1.470756e+08
23 1.987508e+08
24 1.940160e+08
25 1.010612e+08
26 1.434164e+08
27 2.538550e+08
28 1.802037e+08
29 7.714155e+07
30 6.271689e+07
31 1.996139e+08
32 1.668192e+08
33 3.716022e+07
34 1.362498e+08
35 1.513207e+08
36 5.341211e+07
37 7.420274e+07
38 5.515963e+07
39 2.074455e+08
40 1.378703e+08
41 1.813418e+08
42 7.956575e+07
43 9.056544e+07
44 4.323009e+07
45 1.123930e+08

In [10]: print("{:.2f}".format(store_totalweeklysales.max()))
# using argmax to get the max sales store index
print(store_totalweeklysales.index[store_totalweeklysales.argmax()])

301397792.45
20

Store number : 20 has the maximum sales of 301397792.46
```

## Maximum Sales Standard Deviation

```
In [11]: # groupby stores and get std.dev
store_sales_std = walmart_df.groupby('store')['weekly_sales'].std()
store_sales_std.to_frame()

Out[11]:
   weekly_sales
store
1  155980.767761
2  237683.694682
3  46319.631557
4  286201.442297
5  37737.965745
6  212525.855882
7  125985.689220
8  106280.829881
9  69028.666585
10 302262.052504
11 165833.887863
12 139166.871880
13 265006.995776
14 317569.649476
15 120638.652043
16 85769.680133
17 112162.630087
18 176641.510839
19 1971722.638730
20 275900.562742
21 128752.812853
22 161251.350031
23 249788.039068
24 167745.677567
25 112976.788600
26 110431.288141
27 239630.135688
28 181758.967539
29 199120.136590
30 22809.665996
31 125855.942933
32 138017.252087
33 24132.927322
34 104630.164676
35 211243.457791
36 60725.173579
37 21837.461190
38 42768.169450
39 217466.454833
40 119002.112858
41 187907.162766
42 50282.925530
43 40598.413260
44 24782.832015
45 130168.526935

In [12]: print("{:.2f}".format(store_sales_std.max()))
# using argmax to get the max std.dev store index
print(store_sales_std.index[store_sales_std.argmax()])

317569.65
14

Store number : 14 has the maximum sales standard sevation 317569.95
```

## Coefficient of Mean to Standard Deviation

```
In [13]: # groupby stores and get mean values
store_sales_mean = walmart_df.groupby('store')['weekly_sales'].mean()
store_sales_mean.to_frame()

Out[13]:
   weekly_sales
store
1  1.555284e+08
2  1.925751e+08
3  4.022044e+05
4  2.094713e+06
5  3.180118e+06
6  1.564713e+06
7  5.706173e+05
8  9.087495e+05
9  5.343806e+05
10 1.899425e+06
11 1.356383e+06
12 1.009020e+06
13 2.003620e+06
14 2.029078e+06
15 6.233125e+05
16 5.192477e+05
17 8.930814e+05
18 1.084718e+06
19 1.444999e+06
20 2.107677e+06
21 7.560691e+05
22 1.028501e+06
23 1.389864e+06
24 1.356755e+06
25 7.067215e+05
26 1.002912e+06
27 1.752161e+06
28 5.394514e+05
29 4.385796e+05
30 1.395901e+06
31 1.166568e+06
32 2.596817e+05
33 9.667816e+05
34 9.197250e+05
35 3.735120e+05
36 1.589003e+05
37 5.189003e+05
38 3.857317e+05
39 1.450669e+06
40 9.641280e+05
41 5.604125e+06
42 5.560439e+05
43 6.332474e+05
44 3.027489e+05
45 7.859814e+05

In [14]: #cv=(mean/std.dev)*100
(covariance_std_mean = (store_sales_std / store_sales_mean) * 100

Covariance Values (Cv)
- Cv less than 10 = Very Good
- Cv between 10 - 20 = Good
- Cv between 21 - 30 = Acceptable
- Cv greater than 30 = Not Acceptable

In [15]: covariance_std_mean.to_frame()

Out[15]:
   weekly_sales
store
1  10.029212
2  12.342388
3  11.502141
4  12.708254
5  11.866844
6  13.552989
7  19.730469
8  11.695283
9  12.689547
10 15.913349
11 12.226183
12 13.762532
13 13.251363
14 15.713674
15 19.338399
16 16.518905
17 12.552067
18 16.284550
19 13.268912
20 17.002629
21 15.678288
22 17.972115
23 12.363738
24 15.986004
25 11.011066
26 13.515544
27 13.732974
28 18.379427
29 15.200804
30 9.016105
31 11.831040
32 9.286935
33 10.822524
34 22.968111
35 16.257891
36 4.208412
37 11.087545
38 14.990779
39 14.817711
40 9.033533
41 6.410363
42 8.179331
43 16.561273

Store numbers : 30,31,33,37,42,43,44 has very good co-efficients of mean to standard deviation
```

## Good Quarterly3 growth rate in year 2010

```
In [16]: # Converting date column from object dtype to datetime dtype
walmart_df['date'] = walmart_df['date'].apply(lambda x:dt.datetime.strptime(x,"%d-%m-%Y"))

In [17]: # Using datetimeindex to get year/month/day from dataset
walmart_df['year'] = pd.DatetimeIndex(walmart_df['date']).year
walmart_df['month'] = pd.DatetimeIndex(walmart_df['date']).month
walmart_df['day'] = pd.DatetimeIndex(walmart_df['date']).day

In [18]: walmart_df.head()

Out[18]:
   store  date  weekly_sales  holiday_flag  temperature  fuel_price  cpi  unemployment  year  month  day
0      1  05-02-2010  1643690.90          0          42.31      2.572  211.096358      8.106  2010    Feb    5
1      1  12-02-2010  1641957.44          1          38.51      2.548  211.242170      8.106  2010    Feb    12
2      1  19-02-2010  1611968.17          0          39.93      2.514  211.289143      8.106  2010    Feb    19
3      1  26-02-2010  1409727.59          0          46.63      2.561  211.319643      8.106  2010    Feb    26
4      1  05-03-2010  1554806.68          0          46.50      2.625  211.350143      8.106  2010    Mar    5

In [19]: # Converting date back to string with our desired format
walmart_df['date'] = walmart_df['date'].apply(lambda x:dt.datetime.strftime(x,"%d-%b-%Y"))

In [20]: month_dict = {'1':'Jan',2:'Feb',3:'Mar',4:'Apr',5:'May',6:'Jun',7:'Jul',8:'Aug',9:'Sep',10:'Oct',11:'Nov',
                    '12':'Dec'}
# Converting month int values to corresponding month names using apply function and dictionary
walmart_df['month'] = walmart_df['month'].apply(lambda x: month_dict[x])

In [21]: walmart_df.head()

Out[21]:
   store  date  weekly_sales  holiday_flag  temperature  fuel_price  cpi  unemployment  year  month  day
0      1  05-Jan-10  1643690.90          0          42.31      2.572  211.096358      8.106  2010    Jan    5
1      1  12-Feb-10  1641957.44          1          38.51      2.548  211.242170      8.106  2010    Jan    13
2      1  19-Feb-10  1611968.17          0          39.93      2.514  211.289143      8.106  2010    Feb    19
3      1  26-Mar-10  1409727.59          0          46.63      2.561  211.319643      8.106  2010    Feb    26
4      1  03-Feb-12  1636339.65          0          56.55      3.360  220.172015      7.348  2012    Feb    3

In [22]: # grouping data by years
yearly_walmartdata = walmart_df.groupby('year')
yearly_walmartdata.ngroups

Out[22]:
3

In [23]: # getting groups of different years(2010,2011,2012)
stores_2012_sales = yearly_walmartdata.get_group(2012)
stores_2011_sales = yearly_walmartdata.get_group(2011)
stores_2010_sales = yearly_walmartdata.get_group(2010)

In [24]: stores_2012_sales.head()

Out[24]:
   store  date  weekly_sales  holiday_flag  temperature  fuel_price  cpi  unemployment  year  month  day
100     1  06-Jan-12  1550369.92          0          49.01      3.157  219.714258      7.348  2012    Jan    6
101     1  13-Jan-12  1459601.17          0          48.53      3.261  219.892526      7.348  2012    Jan    13
102     1  20-Jan-12  1394393.84          0          54.11      3.268  219.985689      7.348  2012    Jan    20
103     1  27-Jan-12  1319325.59          0          54.26      3.290  220.078852      7.348  2012    Jan    27
104     1  03-Feb-12  1636339.65          0          56.55      3.360  220.172015      7.348  2012    Feb    3

In [25]: # adding new column to find and store quartile values
stores_2012_sales['quartile'] = 0

In [26]: stores_2012_sales.head()

Out[26]:
   store  date  weekly_sales  holiday_flag  temperature  fuel_price  cpi  unemployment  year  month  day  quartile
100     1  06-Jan-12  1550369.92          0          49.01      3.157  219.714258      7.348  2012    Jan    6    0
101     1  13-Jan-12  1459601.17          0          48.53      3.261  219.892526      7.348  2012    Jan    13    0
102     1  20-Jan-12  1394393.84          0          54.11      3.268  219.985689      7.348  2012    Jan    20    0
103     1  27-Jan-12  1319325.59          0          54.26      3.290  220.078852      7.348  2012    Jan    27    0
104     1  03-Feb-12  1636339.65          0          56.55      3.360  220.172015      7.348  2012    Feb    3    0

In [27]: stores_2012_sales['month'].unique()

Out[27]:
array(['Jan', 'Feb', 'Mar', 'Apr', 'May', 'Jun', 'Jul', 'Aug', 'Sep',
       'Oct'], dtype=object)

In [28]: # using python code to split months into quartiles
for i in stores_2012_sales['month']:
    if i in ['Jan','Feb','Mar']:
        stores_2012_sales['quartile'][stores_2012_sales[stores_2012_sales['month'] == i].index] = "Q1"
    elif i in ['Apr','May','Jun']:
        stores_2012_sales['quartile'][stores_2012_sales[stores_2012_sales['month'] == i].index] = "Q2"
    elif i in ['Jul','Aug','Sep']:
        stores_2012_sales['quartile'][stores_2012_sales[stores_2012_sales['month'] == i].index] = "Q3"
    elif i in ['Oct','Nov','Dec']:
        stores_2012_sales['quartile'][stores_2012_sales[stores_2012_sales['month'] == i].index] = "Q4"

In [29]: stores_2012_sales

Out[29]:
   store  date  weekly_sales  holiday_flag  temperature  fuel_price  cpi  unemployment  year  month  day  quartile
100     1  06-Jan-12  1550369.92          0          49.01      3.157  219.714258      7.348  2012    Jan    6    Q1
101     1  13-Jan-12  1459601.17          0          48.53      3.261  219.892526      7.348  2012    Jan    13    0
102     1  20-Jan-12  1394393.84          0          54.11      3.268  219.985689      7.348  2012    Jan    20    0
103     1  27-Jan-12  1319325.59          0          54.26      3.290  220.078852      7.348  2012    Jan    27    0
104     1  03-Feb-12  1636339.65          0          56.55      3.360  220.172015      7.348  2012    Feb    3    Q1
6430  45  28-Sep-12  713173.95          0          64.88      3.997  192.013558      8.684  2012    Sep    28    Q3
6431  45  05-Oct-12  733455.07          0          64.89      3.985  192.170412      8.667  2012    Oct    5    Q4
6432  45  12-Oct-12  734464.36          0          54.47      4.000  192.327265      8.667  2012    Oct    12    Q4
6433  45  19-Oct-12  718125.53          0          56.47      3.969  192.330854      8.667  2012    Oct    19    Q4
6434  45  26-Oct-12  760281.43          0          58.85      3.882  192.308899      8.667  2012    Oct    26    Q4

1935 rows x 12 columns

In [30]: # grouping data by quartiles
quartiles_2012sales = stores_2012_sales.groupby('quartile')
quartiles_2012sales.ngroups

Out[30]:
4

In [31]: # get individual quartile groups(Q1,Q2,Q3,Q4)
q1_data = quartiles_2012sales.get_group('Q1')
q2_data = quartiles_2012sales.get_group('Q2')
q3_data = quartiles_2012sales.get_group('Q3')
q4_data = quartiles_2012sales.get_group('Q4')

In [32]: # getting total sales of all stores in each quarters
q1_sales = q1_data.groupby('store')['weekly_sales'].sum()
q2_sales = q2_data.groupby('store')['weekly_sales'].sum()
q3_sales = q3_data.groupby('store')['weekly_sales'].sum()
q4_sales = q4_data.groupby('store')['weekly_sales'].sum()

In [33]: q1_sales.head(3)

Out[33]:
store
1  20978760.12
2  25083604.88
3  5620316.49
Name: weekly_sales, dtype: float64

In [34]: q3_sales.head(3)

Out[34]:
store
1  20253947.78
2  24303354.86
3  5298005.47
Name: weekly_sales, dtype: float64

In [35]: # finding the Q3 growth rate for each store
q3_salesgrowth = q3_sales - q2_sales
q3_salesgrowth.to_frame()

Out[35]:
   weekly_sales
store
1  -724812.34
2  -780260.02
3  -322311.02
4  -657571.21
5  -302572.70
6  -666597.68
7  971928.12
8  -10678.25
9  -462785.55
10 -713110.41
11 -271290.51
12 -820084.21
13 -987974.84
14 -343162.04
15 557205.66
16 -132947.88
17 -406429.38
18 -163745.39
19 -632670.34
20 -2696907.03
21 -642754.35
22 152680.33
23 292158.81
24 -21330.25
25 520356.34
26 -436301.34
27 -454073.36
28 -147612.43
29 -480524.05
30 -92742.10
31 -115380.03
32 -367622.08
33 484108.12
34 -320299.94
35 -96481.13
36 -32436.44
37 500867.77
38 145457.84
39 -271479.93
40 -108264.19
41 10845.38
42 -188284.39
43 -809499.45

Store number : 7 has good quarterly 3 sales growth than other stores
```

## Holiday and Non-Holiday Sales

```
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, 29-Nov-13
4) Christmas: 31-Dec-10, 30-Dec-11, 28-Dec-12, 27-Dec-13

In [37]: # Finding mean sales on superbowl for all stores each year
superbowl_sales = walmart_df[walmart_df['date'] == '12-Feb-10']|walmart_df['date'] == '11-Feb-11']|(walmart_df['date'] == '10-Feb-12'])
print("{:.2f}".format(superbowl_sales['weekly_sales'].mean()))

1079127.99

In [38]: # finding mean sales on labourday for all stores each year
labourday_sales = walmart_df[walmart_df['date'] == '10-Sep-10']|walmart_df['date'] == '09-Sep-11']|(walmart_df['date'] == '07-Sep-12'])
print("{:.2f}".format(labourday_sales['weekly_sales'].mean()))

1042427.29

In [39]: # finding mean sales on thanksgiving for all stores each year
thanksgiving_sales = walmart_df[walmart_df['date'] == '12-Nov-10']|walmart_df['date'] == '25-Nov-11']|(walmart_df['date'] == '23-Nov-12'])
print("{:.2f}".format(thanksgiving_sales['weekly_sales'].mean()))

1471273.43

In [40]: # finding mean sales on christmas for all stores each year
christmas_sales = walmart_df[walmart_df['date'] == '31-Dec-10']|walmart_df['date'] == '30-Dec-11']|(walmart_df['date'] == '28-Dec-12'])
print("{:.2f}".format(christmas_sales['weekly_sales'].mean()))

960833.11

Thanksgiving Holiday have Higher mean sales than any other Holiday sales
```

## Holiday and Non-Holiday data

```
In [41]: # grouping dataset by holidayflag
holidays = walmart_df.groupby('holiday_flag')

Out[41]:
holiday_flag
0      1  05-02-2010  1643690.90          0          42.31      2.572  211.096358      8.106
1      1  12-02-2010  1641957.44          1          38.51      2.548  211.242170      8.108
2      1  19-02-2010  1611968.17          0          39.93      2.514  211.289143      8.106
3      1  26-02-2010  1409727.59          0          46.63      2.561  211.319643      8.106
4      1  05-03-2010  1554806.68          0          46.50      2.625  211.350143      8.108
...
```



```
In [42]: # getting only holiday data
holiday_sales = holidaysdf.get_group(1)
print("{:.2f}".format(holiday_sales['weekly_sales'].mean()))
1122887.89
```

```
In [43]: # getting only non-holiday data
nonholiday_sales = holidaysdf.get_group(0)
print("{:.2f}".format(nonholiday_sales['weekly_sales'].mean()))
1041256.38
```

## Holidays have higher sales mean than non-holiday sales mean

```
In [44]: # Each holiday vs Non-holiday sales
```

```
In [45]: print("{:.2f}".format(superbowlsales['weekly_sales'].mean()))
print("{:.2f}".format(labourdaysales['weekly_sales'].mean()))
print("{:.2f}".format(thanksgiving_sales['weekly_sales'].mean()))
print("{:.2f}".format(christmas_sales['weekly_sales'].mean()))
1071972.99
1042427.59
1471273.43
960833.11
```

```
In [46]: print("{:.2f}".format(nonholiday_sales['weekly_sales'].mean()))
1041256.38
```

## Christmas holiday sales has a negative impact on sales than other holidays when compared to Non-holiday sales mean

```
In [ ]:
```

## Monthly and Semester view of Sales (Year-wise)

```
In [47]: # grouping by month for all years
totalmonthly_sales = walmart_df.groupby('month')
totalmonthly_sales.agggroups
```

```
In [48]: totalmonthly_sales['weekly_sales'].describe()
Out[48]:
```

	count	mean	std	min	25%	50%	75%	max
month								
Apr	6300	1.02878e+06	543864.824192	23769.09	534983.5075	948789.575	1402726.700	2565259.92
Aug	5850	1.04197e+06	542663.050946	224031.19	575997.7800	990387.480	1428191.280	2283540.30
Dec	4500	1.28188e+06	774037.720767	209986.25	616295.8450	1154880.920	1714442.515	3818686.45
Feb	5400	1.05320e+06	564207.057354	234218.03	554628.6350	980765.180	1431376.635	2623469.95
Jan	3800	9.23884e+05	542616.460339	231155.90	521051.0125	830944.935	1256202.130	2047766.07
Jul	6300	1.03174e+06	531141.778886	224806.96	577830.8550	953770.830	1396269.865	2358055.30
Jun	5850	1.06432e+06	548693.953606	238172.66	581745.7200	984360.040	1442920.080	2325444.75
Mar	5850	1.03130e+06	529865.743801	238084.08	541408.1400	943951.650	1368524.970	2237544.75
May	5400	1.03714e+06	536589.412470	239206.28	543588.2225	969562.080	1388630.735	2307166.52
Nov	3600	1.14726e+06	648795.327036	224639.76	558963.6500	959567.800	1437824e+06	2207141.81
Oct	5850	9.96321e+05	517186.653614	213538.32	534738.4300	837956.890	1363592e+06	2246411.89
Sep	5850	9.89335e+05	510532.949375	229731.98	533761.6400	822404.600	1345167.610	2202742.90

```
In [ ]:
```

```
In [49]: # grouping by month for year 2010
monthlysales_2010 = stores_2010_sales.groupby('month')
monthlysales_2010['weekly_sales'].describe()
```

```
Out[49]:
```

	count	mean	std	min	25%	50%	75%	max
month								
Apr	2250	1.028499e+06	545024.247786	257361.30	527019.7800	969894.470	1.391259e+06	2495630.51
Aug	1900	1.041441e+06	549322.149980	224031.19	573235.4675	961237.730	1.497055e+06	22719813.50
Dec	2250	1.283380e+06	793498.247836	209986.25	599730.0700	1149612.040	1.702798e+06	3818686.45
Feb	1800	1.057400e+06	57029.268079	267956.30	534468.2225	997074.875	1.470289e+06	2623469.95
Jul	2250	1.036386e+06	543430.296029	242047.03	570231.2100	959229.090	1.371987e+06	2334788.42
Jan	1800	1.03130e+06	556849.016817	239419.91	551924.8925	1002244.605	1.423940e+06	2336301.47
Jun	1800	1.010665e+06	511192.520079	262893.76	519202.3825	964763.600	1.340757e+06	2237544.75
Mar	1800	1.03726e+06	542865.021379	267085.35	530713.0675	951724.460	1.368593e+06	2370116.52
Nov	1800	1.128963e+06	648795.522239	224639.76	558963.6500	959567.800	1.437824e+06	2039946.38
Oct	2250	9.961637e+05	546358.362753	213538.32	534738.4300	837956.890	1.363592e+06	2246411.89
Sep	1800	9.848216e+05	508207.747513	231978.64	530738.3500	810702.500	1.304980e+06	2019167.76

```
In [50]: semester_2010 = stores_2010_sales[stores_2010_sales['month'].isin(['Jan','Feb','Mar','Apr','May','Jun'])]
semester_2010.head()
```

```
Out[50]:
```

	store	date	weekly_sales	holiday_flag	temperature	fuel_price	cpi	unemployment	year	month	day
0	1	05-Feb-10	1643690.90	0	42.31	2.572	211.096358	8.106	2010	Feb	5
1	1	12-Feb-10	1641957.44	1	38.51	2.548	211.242170	8.106	2010	Feb	12
2	1	19-Feb-10	1611968.17	0	39.93	2.514	211.289143	8.106	2010	Feb	19
3	1	26-Feb-10	1409727.59	0	46.63	2.561	211.319643	8.106	2010	Feb	26
4	1	05-Mar-10	1554806.68	0	46.50	2.625	211.350143	8.106	2010	Mar	5

```
In [51]: semester_2010 = stores_2010_sales[stores_2010_sales['month'].isin(['Jul','Aug','Sep','Oct','Nov','Dec'])]
semester_2010.head()
```

```
Out[51]:
```

	store	date	weekly_sales	holiday_flag	temperature	fuel_price	cpi	unemployment	year	month	day
21	1	02-Jul-10	1492418.14	0	80.91	2.689	211.223533	7.787	2010	Jul	2
22	1	09-Jul-10	1546074.18	0	80.48	2.642	211.108414	7.787	2010	Jul	9
23	1	16-Jul-10	1448938.92	0	83.15	2.623	211.100385	7.787	2010	Jul	16
24	1	23-Jul-10	1350655.20	0	83.36	2.608	211.235154	7.787	2010	Jul	23
25	1	30-Jul-10	1371986.60	0	81.84	2.640	211.369903	7.787	2010	Jul	30

```
In [ ]:
```

```
In [52]: # grouping by month for year 2011
monthlysales_2011 = stores_2011_sales.groupby('month')
monthlysales_2011['weekly_sales'].describe()
```

```
Out[52]:
```

	count	mean	std	min	25%	50%	75%	max
month								
Apr	2250	1.049651e+06	558163.253234	249798.75	548516.0000	950743.050	1.456415e+06	2565259.92
Aug	1800	1.067020e+06	571613.468771	245226.80	564122.8100	974532.6150	1.445838e+06	2273470.62
Dec	2250	1.260347e+06	758483.504954	215359.21	630327.2800	1158708.980	1.781529e+06	3676388.98
Feb	1800	1.036386e+06	542442.387986	234218.03	55183.5800	961012.850	1.420504e+06	2334788.42
Jan	1800	9.904665e+05	526122.257495	231155.90	51597.3125	812167.545	1.215770e+06	1889393.94
Jul	2250	1.021828e+06	548417.700020	224806.96	562381.9500	936001.980	1.396972e+06	2123787.79
Jun	1800	1.054297e+06	543819.984741	238172.66	568639.7750	974263.155	1.441112e+06	2182246.69
Mar	1800	9.96427e+05	52341.244321	238084.08	536784.3775	929977.755	1.352847e+06	2143424.61
May	1800	1.009156e+06	523438.331445	239206.28	543719.8250	943238.085	1.371667e+06	2095599.93
Nov	1800	1.167506e+06	656014.430247	236157.12	597806.1500	961735.140	1.559386e+06	3004702.33
Oct	1800	1.018118e+06	533292.121588	231319.96	537294.8075	843238.485	1.406233e+06	2202742.90
Sep	2250	9.815456e+05	511896.762925	229731.98	537124.7600	869834.750	1.347608e+06	2202742.90

```
In [53]: semester_2011 = stores_2011_sales[stores_2011_sales['month'].isin(['Jan','Feb','Mar','Apr','May','Jun'])]
semester_2011.head()
```

```
Out[53]:
```

	store	date	weekly_sales	holiday_flag	temperature	fuel_price	cpi	unemployment	year	month	day
48	1	07-Jan-11	1444732.28	0	48.27	2.976	211.404742	7.742	2011	Jan	7
49	1	14-Jan-11	1397013.96	0	35.40	2.963	211.457411	7.742	2011	Jan	14
50	1	21-Jan-11	1327405.42	0	44.04	3.016	211.827234	7.742	2011	Jan	21
51	1	28-Jan-11	1318899.31	0	43.83	3.010	212.197058	7.742	2011	Jan	28
52	1	04-Feb-11	1600262.58	0	42.27	2.989	212.566881	7.742	2011	Feb	4

```
In [54]: semester_2011 = stores_2011_sales[stores_2011_sales['month'].isin(['Jul','Aug','Sep','Oct','Nov','Dec'])]
semester_2011.head()
```

```
Out[54]:
```

	store	date	weekly_sales	holiday_flag	temperature	fuel_price	cpi	unemployment	year	month	day
73	1	01-Jul-11	1485838.09	0	85.55	3.524	215.194317	7.982	2011	Jul	1
74	1	08-Jul-11	1534846.64	0	85.83	3.480	215.277175	7.982	2011	Jul	8
75	1	15-Jul-11	1455116.97	0	88.54	3.575	215.381109	7.982	2011	Jul	15
76	1	22-Jul-11	1396026.82	0	85.77	3.651	215.422278	7.982	2011	Jul	22
77	1	29-Jul-11	1352218.79	0	86.63	3.682	215.483446	7.982	2011	Jul	29

```
In [ ]:
```

```
In [55]: # grouping by month for year 2012
monthlysales_2012 = stores_2012_sales.groupby('month')
monthlysales_2012['weekly_sales'].describe()
```

```
Out[55]:
```

	count	mean	std	min	25%	50%	75%	max
month								
Apr	1800	1.049651e+06	558163.253234	249798.75	548516.0000	950743.050	1.456415e+06	2565259.92
Aug	2250	1.052670e+06	571613.468771	245226.80	564122.8100	974532.6150	1.445838e+06	2273470.62
Feb	1800	1.067020e+06	571613.468771	245226.80	564122.8100	974532.6150	1.445838e+06	2273470.62
Jan	1800	9.363026e+05	519887.953606	236920.49	540523.1525	855680.105	1.269834e+06	2047766.07
Jul	1800	1.04171e+06	535223.017499	249134.32	577204.6500	960529.395	1.456741e+06	238055.30
Jun	2250	1.069379e+06	548322.510590	244338.31	583646.5900	988764.840	1.451778e+06	2245287.18
Mar	2250	1.028932e+06	536489.204740	246997.07	557547.8000	952264.100	1.427781e+06	2214987.44
May	1800	1.048703e+06	545296.384619	261851.74	549055.7525	988833.255	1.428744e+06	2207214.81
Oct	1800	1.024232e+06	526157.386422	233731.13	548820.2100	962230.855	1.431769e+06	2246411.89
Sep	1800	1.003595e+06	513786.298200	242813.51	532561.7925	940096.910	1.371256e+06	2165796.31

```
In [56]: semester_2012 = stores_2012_sales[stores_2012_sales['month'].isin(['Jan','Feb','Mar','Apr','May','Jun'])]
semester_2012.head()
```

```
Out[56]:
```

	store	date	weekly_sales	holiday_flag	temperature	fuel_price	cpi	unemployment	year	month	day	quarter
100	1	06-Jan-12	1550089.92	0	49.01	3.157	219.714258	7.348	2012	Jan	6	Q1
101	1	13-Jan-12	1459601.17	0	45.53	3.261	219.892526	7.348	2012	Jan	13	Q1
102	1	20-Jan-12	1394393.84	0	54.11	3.268	219.985689	7.348	2012	Jan	20	Q1
103	1	27-Jan-12	1319325.59	0	54.26	3.290	220.078852	7.348	2012	Jan	27	Q1
104	1	03-Feb-12	156339.65	0	56.55	3.360	220.172015	7.348	2012	Feb	3	Q1

```
In [57]: semester_2012 = stores_2012_sales[stores_2012_sales['month'].isin(['Jul','Aug','Sep','Oct','Nov','Dec'])]
semester_2012.head()
```

```
Out[57]:
```

```

# Converting string object to datetime object
walmart_df['date'] = pd.to_datetime(walmart_df['date'])

# Changing dates to the day of the week
walmart_df['dayofweek'] = walmart_df['date'].dt.day_name()

```

	store	date	weekly_sales	holiday_flag	temperature	fuel_price	cpi	unemployment	year	month	day	dayofweek
0	1	2010-02-05	1643690.90	0	42.31	2.572	211.096358	8.106	2010	Feb	5	Friday
1	1	2010-02-12	1641957.44	1	38.51	2.548	211.242170	8.106	2010	Feb	12	Friday
2	1	2010-02-19	1611988.17	0	39.93	2.514	211.289143	8.106	2010	Feb	19	Friday
3	1	2010-02-26	1409727.59	0	46.83	2.561	211.319643	8.106	2010	Feb	26	Friday
4	1	2010-03-05	1554806.68	0	46.50	2.625	211.350143	8.106	2010	Mar	5	Friday
5	1	2010-03-12	1439541.59	0	57.79	2.667	211.380643	8.106	2010	Mar	12	Friday
6	1	2010-03-19	1472515.79	0	54.58	2.720	211.215635	8.106	2010	Mar	19	Friday
7	1	2010-03-26	1404429.92	0	51.45	2.719	210.818042	8.106	2010	Mar	26	Friday
8	1	2010-04-02	1494968.28	0	62.27	2.719	210.820450	7.808	2010	Apr	2	Friday
9	1	2010-04-09	145418.53	0	65.86	2.770	210.622857	7.808	2010	Apr	9	Friday