

WALMART - STORE SALES FORECAST




```
In [1]: 1  ## import all the necessary libraries
2  import warnings
3  warnings.filterwarnings("ignore")
4  import pandas as pd
5  import sqlite3
6  import csv
7  import matplotlib.pyplot as plt
8  import seaborn as sns
9  import numpy as np
10 from wordcloud import WordCloud
11 import re
12 import chart_studio.plotly as py
13 import plotly.express as px
14
15 import os
16 import plotly.graph_objects as go
17 import numpy as np
18 import seaborn as sns
19 from scipy import stats
20 from plotly.subplots import make_subplots
21 from mlxtend.feature_selection import SequentialFeatureSelector as SFS
22 from sklearn.ensemble import RandomForestRegressor
23 import statsmodels.api as sm
24 from statsmodels.formula.api import ols
25 from sklearn.metrics import mean_absolute_error
26 from sqlalchemy import create_engine # database connection
27 import datetime as dt
28 from sklearn import metrics
29 import pickle as pkl
30 from sklearn.metrics import f1_score, precision_score, recall_score
31 from datetime import datetime
32 pd.set_option('max_rows', 80000, 'max_columns', 200)
```

1. Business Problem

Description

Walmart is leading multinational Retail corporation (headquartered in the USA) operates a chain of hypermarkets.

The sales prediction problem was launched by Walmart via Kaggle as a part of its recruitment process to predict their weekly store sales using the historical sales data of 45 stores.

Why is sales forecasting important?

- sales forecasting allows companies to estimate their cost and revenue effectively accurately based on which they are able to predict their short-term and long-term performance
- In short we can say forecasts help in financial planning, Sales Planning, Inventory Control, Supply chain management, cash-flow management.

How is machine learning useful in this scenario ?

There is not clear cut pattern humans have set for buying and selling, However backing up the datasets spanning years into the past it's highly possible to draw patterns in sales and consumption. Machine learning becomes prominent here because of its ability to mine through years of data to detect patterns and repetitive behaviour, which can then be leveraged to forecast sales and demand. Conventional methods like moving averages and exponential moving averages we are assuming future sales values might be some way nearer to the recent past values, thereby cannot capture the external driving factors that determine the abrupt spike in sales values in some weeks. For example in this case study we have factors like if the given week is a holiday or not. If the week is a holiday week we capture the insight that there will be spike in sales due to this factor and if there are any promotional events conducted during the week by the retailer which again will account for an impact. Using machine learning we are not just able to capture them but also these algorithms are self-correcting. (i.e.) as newer data gets added in the future, the model accustoms itself to it and learns many new uncertainties in data.

Problem Statement

Predict weekly sales for each store and department based on the past years sales behaviour.

Source: <https://www.kaggle.com/c/walmart-recruiting-store-sales-forecasting>

1.2 Source / useful links

Data Source : [https://www.kaggle.com/c/walmart-recruiting-store-sales-forecasting_\(https://www.kaggle.com/c/walmart-recruiting-store-sales-forecasting\)](https://www.kaggle.com/c/walmart-recruiting-store-sales-forecasting_(https://www.kaggle.com/c/walmart-recruiting-store-sales-forecasting))

Blogs : [https://medium.com/analytics-vidhya/walmart-sales-forecasting-d6bd537e4904_\(https://medium.com/analytics-vidhya/walmart-sales-forecasting-d6bd537e4904\)](https://medium.com/analytics-vidhya/walmart-sales-forecasting-d6bd537e4904_(https://medium.com/analytics-vidhya/walmart-sales-forecasting-d6bd537e4904))

Research paper : "[http://cs229.stanford.edu/proj2017/final-reports/5244336.pdf_\(http://cs229.stanford.edu/proj2017/final-reports/5244336.pdf%E2%80%9D\)](http://cs229.stanford.edu/proj2017/final-reports/5244336.pdf_(http://cs229.stanford.edu/proj2017/final-reports/5244336.pdf%E2%80%9D))

Research paper : "[https://www.researchgate.net/publication/339362837_Learnings_from_Kaggle's_Forecasting_Competitions_\(https://www.researchgate.net/publication/339362837_Learnings_from_Kaggle's_Forecasting_Competitions\)](https://www.researchgate.net/publication/339362837_Learnings_from_Kaggle's_Forecasting_Competitions_(https://www.researchgate.net/publication/339362837_Learnings_from_Kaggle's_Forecasting_Competitions))"

(https://www.researchgate.net/publication/339362837_Learnings_from_Kaggle's_Forecasting_Competitions)

Blogs : " [https://machinelearningmastery.com/multivariate-time-series-forecasting-lstms-keras/\(https://machinelearningmastery.com/multivariate-time-series-forecasting-lstms-keras/\)](https://machinelearningmastery.com/multivariate-time-series-forecasting-lstms-keras/(https://machinelearningmastery.com/multivariate-time-series-forecasting-lstms-keras/))

1.3 Real World / Business Objectives and Constraints

1. **Robustness** *Ensuring using Cross validation*
2. **Interpretability** As we forecast sales and hand it over to businesses, it is important to answer Why's ? and How's
3. Accurate prediction matters as it might impact taking business decisions
4. **No latency** requirements
5. Models should be **Evolutionary** as consumer demands/supplies can change over time.

2. Machine Learning problem

2.1 Data

2.1.1 Data Overview

Data was collected from year 2010 - 2012 for 45 Walmart stores. We are tasked with predicting department wise sales for each store.

We are provided with datasets : Stores.csv, train.csv, test.csv, fetaures.csv

Dataset Size

- Trainset : 421570 rows (12.2MB)
- Testset : 115064 rows (2.47MB)
- Stores.csv : 4KB
- features.csv: 580 KB

Column Descriptions

Stores

- Type: Type of the store namely "A", "B", "C"
- Size: Size of the store .Size refers to the number of products inside the store .Range 34000 -210,000

Sales

- Store - The store number. Range 1- 45.
- Dept - One of 1–99 that shows the department
- Date - the week
- Weekly_Sales - sales for the given department in the given store
- IsHoliday - whether the week is a special holiday week.

Features

- Temperature - Average temperature in the region during that week.
- Fuel_Price - cost of fuel in the region during that week
- Markdown1-5 - anonymized data related to promotional markdowns that Walmart is running. Markdown data is only available after Nov 2011, and is not available for all stores all the time. Any missing value is marked with an NA.
represents what quantity was available during that week. Type is from 1 - 5.
- CPI - the consumer price index during that week.
- Unemployment - the unemployment rate during that week.

2.2 Mapping the real-world problem to a Machine Learning Problem

2.2.1 Type of Machine Learning Problem

It's a regression problem where we are given time series(from 05-02-2010 to 26-10-2012) data with multiple categorical and numerical features.We are tasked to predict department wise sales of each store from 02-11-2012 to 26-07-2013 . In addition, Walmart runs several promotional markdown events throughout the year. These markdowns precede prominent holidays, the four largest of 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 is modeling the effects of markdowns on these holiday weeks in the absence of complete/ideal historical data

2.2.2 Performance metric

1. Weighted mean absolute error

$$WMAE = \frac{1}{\sum w_i} \sum_{i=1}^n w_i |y_i - \hat{y}_i|$$

where n is the number of rows.

\hat{y}_i is the predicted sales

y_i is the actual sales

w_i are weights. $w = 5$ if the week is a holiday week, 1 otherwise

2. MAE

3. R2-(R-squared)

Why not MSE/RMSE/MAPE ??

- MAE seems to be a good choice for performance metric as it is robust towards outlier points than compared to RMSE,MSE.We are going to train multiple models and we want to check how each model explains variance of the outcome.Hence we also use R-2 (R-squared) for comparing realtive model performance.Why we dont want to go forth choosing MAPE for this study . MAPE produces undefined outcome when actual values are 0 .((1/n * sum(0-prediction / 0) ~ undefined).We also have more than 100% MAPE which is tough to interpret. MAPE treats underforecast and overforecast differently.For underforecasts we have MAPE < 100% and over forecasts have MAPE > 100%.

3. Exploratory Data Analysis

3.1 Data Loading

```
In [2]: 1 ### Loading our trainset and test set
        2 trainset = pd.read_csv("train.csv")
        3 testset = pd.read_csv("test.csv")
        4 featureset = pd.read_csv("features.csv")
        5 storeset = pd.read_csv("stores.csv")
```

```
In [3]: 1 ## Check number of rows and columns
        2 print("Shape of trainset: ",trainset.shape)
        3 print("Shape of testset: ",testset.shape)
```

Shape of trainset: (421570, 5)

Shape of testset: (115064, 4)

3.2 Identifying Variables and data types

```
In [4]: 1 ### Check field names of all the columns
        2 print("Field name in trainset ",trainset.columns)
        3 print("\n")
        4 print("*****100")
        5 print("Field name in featureset ",featureset.columns)
        6 print("\n")
        7 print("*****100")
        8 print("Field name in storeset ",storeset.columns)
        9 print("\n")
        10 print("*****100")
```

Field name in trainset Index(['Store', 'Dept', 'Date', 'Weekly_Sales', 'IsHoliday'], dtype='object')

```
*****
Field name in featureset Index(['Store', 'Date', 'Temperature', 'Fuel_Price', 'MarkDown1', 'MarkDown2',
                               'MarkDown3', 'MarkDown4', 'MarkDown5', 'CPI', 'Unemployment',
                               'IsHoliday'],
                               dtype='object')
```

```
*****
Field name in storeset Index(['Store', 'Type', 'Size'], dtype='object')
```

```
*****
```

```
In [5]: 1 ### Check numerical and categorical variables
        2 print("In the above dataset the numerical variables are :\tStore,Dept,Weekly_Sales,Markdown1:5,Temperature,Fuel_Price,CPI,Unemployment,size")
        3 print("Categorical variables are :      IsHoliday,Type")
```

In the above dataset the numerical variables are : Store,Dept,Weekly_Sales,Markdown1:5,Temperature,Fuel_Price,CPI,Unemployment,size
Categorical variables are : IsHoliday,Type

```
In [6]: 1 ## Let's check the dtypes of the columns
        2 trainset.dtypes.to_frame().rename(columns={0:"Datatype"})
```

Out[6]:

| Datatype | |
|--------------|---------|
| Store | int64 |
| Dept | int64 |
| Date | object |
| Weekly_Sales | float64 |
| IsHoliday | bool |

- we can see that the date field is in string format lets convert to datetime using pandas to_datetime feature.

3.3 Treating missing values and duplicates

How null values can affect our data ?

“ Missing values are common occurrences in data. Unfortunately, most predictive modeling techniques cannot handle any missing values. Therefore, this problem must be addressed prior to modeling.” Missing/null values in dataset arises problems of "Value Error" while modeling

```

In [7]: 1  ## Lets check missing values in our Store set
        2  print('Number of unique values in Store = ',storeset.nunique())
        3  print('*'*50)
        4  print('Null values check in Store :')
        5  print(storeset.apply(lambda x: sum(x.isnull()),axis=0))
        6  ## Lets check missing values in our feature set
        7  print('*'*50)
        8  print('Unique values in feature:',featureset.nunique())
        9  print('*'*50)
       10  print('Null values check in Feature :')
       11  print(featureset.apply(lambda x: sum(x.isnull()),axis=0))
       12  ## Lets check missing values in our trainset
       13  print('*'*50)
       14  print('Unique values in train:',trainset.nunique())
       15  print('*'*50)
       16  print('Null values check in train :')
       17  print(trainset.apply(lambda x: sum(x.isnull()),axis=0))
       18  ## Lets check missing values in our testset
       19  print('*'*50)
       20  print('Unique values in test:',testset.nunique())
       21  print('*'*50)
       22  print('Null values check in test :')
       23  print(testset.apply(lambda x: sum(x.isnull()),axis=0))

```

Number of unique values in Store = Store 45

Type 3

Size 40

dtype: int64

Null values check in Store :

Store 0

Type 0

Size 0

dtype: int64

Unique values in feature: Store 45

Date 182

Temperature 4178

Fuel_Price 1011

Markdown1 4023

Markdown2 2715

Markdown3 2885

Markdown4 3405

Markdown5 4045

CPI 2505

Unemployment 404

IsHoliday 2

dtype: int64

Null values check in Feature :

Store 0

Date 0

Temperature 0

Fuel_Price 0

Markdown1 4158

Markdown2 5269

Markdown3 4577

Markdown4 4726

```

Markdown5      4140
CPI             585
Unemployment    585
IsHoliday       0
dtype: int64
*****
Unique values in train: Store      45
Dept           81
Date           143
Weekly_Sales   359464
IsHoliday       2
dtype: int64
*****
Null values check in train :
Store          0
Dept           0
Date           0
Weekly_Sales   0
IsHoliday       0
dtype: int64
*****
Unique values in test: Store      45
Dept           81
Date           39
IsHoliday       2
dtype: int64
*****
Null values check in test :
Store          0
Dept           0
Date           0
IsHoliday       0
dtype: int64

```

As we can see CPI,Unemployment,Markdown Contains null values we can replace them with rolling mean values with period 20

```

In [8]: 1 ## Number of unique stores and number unique departments in trainset
        2 ustore = sorted(trainset.Store.unique()) ## unique stores list
        3 udept = sorted(trainset.Dept.unique()) ## unique department list
        4 udates = sorted(trainset.Date.unique()) ## unique date list
        5 print('Total number of unique store and department and date set is ',len(ustore)*len(udept)*len(udates))

```

Total number of unique store and department and date set is 521235

```

In [9]: 1 print('Our trainset has store department combination of size :',len(trainset.groupby(['Store','Dept','Date']).count()))

```

Our trainset has store department combination of size : 421570

We see that few stores have missing departments and dates and their reported weekly sales !! We will mark their weekly sales as 0 and 'IsHoliday' column as False

We can merge all the different fields from stores,features,train,test to one detailed dataset which contains all the field before performing the analysis,so that it's becomes handy in analysis.As we can observe that "Store" field is common in stores and features set we can perform a merge operation on "store" field and then merge it to the train set using three common fields (i.e) "Stores,Date,IsHoliday"

As we can see other than markdown cell no field has missing values in them..Let's calculate the % of null values in them


```
In [10]: 1 print("% of missing value in Markdown1: ",len(featureset[featureset.MarkDown1.isnull() == True])/len(featureset) * 100,"%")
2 print("% of missing value in Markdown2: ",len(featureset[featureset.MarkDown2.isnull() == True])/len(featureset) * 100,"%")
3 print("% of missing value in Markdown3: ",len(featureset[featureset.MarkDown3.isnull() == True])/len(featureset) * 100,"%")
4 print("% of missing value in Markdown4: ",len(featureset[featureset.MarkDown4.isnull() == True])/len(featureset) * 100,"%")
5 print("% of missing value in Markdown5: ",len(featureset[featureset.MarkDown5.isnull() == True])/len(featureset) * 100,"%")
```

```
% of missing value in Markdown1:  50.76923076923077 %
% of missing value in Markdown2:  64.33455433455434 %
% of missing value in Markdown3:  55.885225885225886 %
% of missing value in Markdown4:  57.7045177045177 %
% of missing value in Markdown5:  50.54945054945055 %
```

- We are going to treat missing values by two ways :
 - If there are more than 75% of null values we will drop the column as we can't even retain 75% information from them.
 - If there are less than 75% of null values we will perform imputation such as replacing with mean value per store

Data preparation for Analysis

- We are going to clean our featureset,train and test set
- Impute missing values in columns CPI,Umenployment with rolling averages and mean for markdown columns
- Generate new rows for missing store ,dept combos with weekly sales as 0
- merge all the dataset and prepare final clean datset

```

In [11]: 1 def clean_feature_(feature):
2         ## Lets fill the nan values in CPI and unemployment with rolling mean of window size= 20 for every store and dept.
3         print('Filling Null values in CPI and Unemployment with rolling averages .....')
4         print('\n')
5         for store in range(1,46) :
6             values = feature[feature.Store == store].set_index('Date')
7             feature.loc[(feature.Store == store),['CPI']] = np.array(values.CPI.fillna(values.CPI.rolling(20,min_periods=1).mean()))
8             feature.loc[(feature.Store == store),['Unemployment']] = np.array(values.Unemployment.fillna(values.Unemployment.rolling(20,min_periods=1).mean()))
9
10        ## Let's fill markdown nan values with mean value in each store
11        print('Filling Null values in Markdown with mean values in each store .....')
12        print('\n')
13        grp = feature.groupby('Store')
14        for i in range(1,46):
15            feature.loc[(feature.Store==i)&(feature.MarkDown1.isnull()==True),['MarkDown1']] = grp.MarkDown1.mean()[i]
16            feature.loc[(feature.Store==i)&(feature.MarkDown2.isnull()==True),['MarkDown2']] = grp.MarkDown2.mean()[i]
17            feature.loc[(feature.Store==i)&(feature.MarkDown3.isnull()==True),['MarkDown3']] = grp.MarkDown3.mean()[i]
18            feature.loc[(feature.Store==i)&(feature.MarkDown4.isnull()==True),['MarkDown4']] = grp.MarkDown4.mean()[i]
19            feature.loc[(feature.Store==i)&(feature.MarkDown5.isnull()==True),['MarkDown5']] = grp.MarkDown5.mean()[i]
20
21        print('\n Null values computed for feature set .....')
22        ### Recheck null values in feature
23        print('Null values recheck in Feature :')
24        print(feature.apply(lambda x: sum(x.isnull()),axis=0))
25        return feature
26
27    def generate_rows(data,s,dep,d,index,type_):
28        '''Function generates new row for missing values of store,department and date combination,with weekly sales as 0
29        and IsHoliday as False'''
30        if type_ == 'train':
31            data.loc[index] = [s,dep,d,0,False]
32        else:
33            data.loc[index] = [s,dep,d,False]
34
35    def trte_clean(data,type_):
36        updates = data.Date.unique()
37        index=len(data)
38        for s in ustore:
39            for dep in udept:
40                for d in updates:
41                    ## check if s,dep,d combo is present in test set
42                    ispresent = bool(len(data.loc[(data.Store == s)&(data.Dept == dep)& (data.Date == d)]))
43                    ## if not present generate new row with default values
44                    if not ispresent:
45                        generate_rows(data,s,dep,d,index,type_)
46                        index+=1
47                    else:pass
48        data.to_csv('Clean_{}.csv'.format(type_),index=False) ## contains missing values imputed csv
49
50
51    def data_cleaning(train,test,feature,store) :
52        start_ = datetime.now()
53        print('Converting datetime type from object to Datetime[ns].....')
54        print('\n')
55
56        ### Chaning datatype to relevent type
57        train['Date'] = pd.to_datetime(train['Date'])
58        test['Date'] = pd.to_datetime(test['Date'])

```

```

59 feature['Date'] = pd.to_datetime(feature['Date'])
60 print('Conversion done in all sets.....')
61 print('\n')
62 print("Type of trainset datetime after conversion:", trainset['Date'].dtype)
63 print('\n')
64 print("Type of testset datetime after conversion:", testset['Date'].dtype)
65 print('\n')
66
67 ### Missing value imputation in feature dataset
68 clean_feature = clean_feature_(feature)
69 print('Computing missing values in train and test.....')
70 print('\n')
71 ### Missing values in train and test
72 index = len(train)
73 if (not os.path.isfile('Clean_train.csv')) and (not os.path.isfile('Clean_test.csv')):
74     cleaned_tr = trte_clean(train, 'train')
75     start = datetime.now()
76     print('Computing done in train.....', datetime.now() - start)
77     print('\n')
78     start = datetime.now()
79     cleaned_te = trte_clean(test, 'test')
80     print('Computing done in test.....', datetime.now() - start)
81     print('\n')
82 else:
83     print('Missing Values imputed train/test is present in disk!!')
84     cleaned_tr = pd.read_csv("Clean_train.csv")
85     cleaned_te = pd.read_csv("Clean_test.csv")
86
87 ### Merging the data to perform univariate analysis / multivariate analysis
88 print('Final Step !! Merging columns!!.....')
89 print('\n')
90 ### feature, train/test, store
91 feat_store = pd.merge(clean_feature, store, how='inner', on=["Store"])
92
93 feat_store['Date'] = pd.to_datetime(feat_store['Date'])
94 cleaned_tr['Date'] = pd.to_datetime(cleaned_tr['Date'])
95 cleaned_te['Date'] = pd.to_datetime(cleaned_te['Date'])
96
97 trainset_merged = pd.merge(cleaned_tr, feat_store, how='inner', on=['Store', 'Date']).sort_values(by=["Store", "Dept", 'Date']).reset_index(drop=True)
98 testset_merged = pd.merge(cleaned_te, feat_store, how='inner', on=['Store', "Date"]).sort_values(by=["Store", "Dept", "Date"]).reset_index(drop=True)
99
100 trainset_merged['set_type'] = 'Train'
101 testset_merged['set_type'] = 'Test'
102
103 print('Datasets cleaned and merged..... ')
104 print('\n')
105 print('Shape of trainset after merge operation is :', trainset_merged.shape)
106 print('\n')
107 print('Shape of testset after merge operation is :', testset_merged.shape)
108 print('\n')
109 print('Let\'s concat train and test data together!! ')
110 print('\n')
111 trte_set = pd.concat([trainset_merged, testset_merged])
112 print('Sorting final set in ascending based on Store, Dept, Date combination!! ')
113 trte_set.sort_values(by=['Store', 'Dept', 'Date'], inplace=True)
114 print('Done-----', datetime.now() - start_)
115 return clean_feature, trainset_merged, testset_merged, trte_set

```



```
In [12]: 1 ### preparing sets
        2 featureset,trainset_merged,testset_merged,trte_set = data_cleaning(trainset,testset,featureset,storeset)
```

Converting datetime type from object to Datetime[ns].....

Conversion done in all sets.....

Type of trainset datetime after conversion: datetime64[ns]

Type of testset datetime after conversion: datetime64[ns]

Filling Null values in CPI and Unemployment with rolling averages

Filling Null values in Markdown with mean values in each store

Null values computed for feature set

Null values recheck in Feature :

Store 0

Date 0

Temperature 0

Fuel_Price 0

MarkDown1 0

MarkDown2 0

MarkDown3 0

MarkDown4 0

MarkDown5 0

CPI 0

Unemployment 0

IsHoliday 0

dtype: int64

Computing missing values in train and test.....

Missing Values imputed train/test is present in disk!!

Final Step !! Merging columns!!.....

Datasets cleaned and merged.....

Shape of trainset after merge operation is : (521235, 18)

Shape of testset after merge operation is : (142155, 17)

Let's concat train and test data together!!

Sorting final set in ascending based onnStore,Dept,Date combination!!

Done----- 0:00:29.282953

In [13]: ▶

```
1 # Lets drop IsHoiday_y as it is a duplicate
2 trainset_merged.drop(columns=['IsHoliday_y'],inplace=True)
3 trainset_merged=trainset_merged.rename(columns={'IsHoliday_x':'IsHoliday'})
4 testset_merged.drop(columns=['IsHoliday_y'],inplace=True)
5 testset_merged=testset_merged.rename(columns={'IsHoliday_x':'IsHoliday'})
6 trte_set.drop(columns=['IsHoliday_y'],inplace=True)
7 trte_set=trte_set.rename(columns={'IsHoliday_x':'IsHoliday'})
```

In [14]: ▶

```
1 trte_set.head(3)
```

Out[14]:

| | Store | Dept | Date | Weekly_Sales | IsHoliday | Temperature | Fuel_Price | MarkDown1 | MarkDown2 | MarkDown3 | MarkDown4 | MarkDown5 | CPI | Unemployment | Type | Size | set_type |
|---|-------|------|------------|--------------|-----------|-------------|------------|-------------|-------------|-------------|-------------|-------------|------------|--------------|------|--------|----------|
| 0 | 1 | 1 | 2010-02-05 | 24924.50 | False | 42.31 | 2.572 | 8536.592778 | 3346.401918 | 1670.797978 | 3653.631444 | 4428.307667 | 211.096358 | 8.106000 | A | 151315 | Train |
| 1 | 1 | 1 | 2010-02-12 | 46039.49 | True | 38.51 | 2.548 | 8536.592778 | 3346.401918 | 1670.797978 | 3653.631444 | 4428.307667 | 211.242170 | 8.106000 | A | 151315 | Train |
| 2 | 1 | 1 | 2010-02-19 | 41595.55 | False | 39.93 | 2.514 | 8536.592778 | 3346.401918 | 1670.797978 | 3653.631444 | 4428.307667 | 211.289143 | 8.106000 | A | 151315 | Train |
| 3 | 1 | 1 | 2010-02-26 | 19403.54 | False | 46.63 | 2.561 | 8536.592778 | 3346.401918 | 1670.797978 | 3653.631444 | 4428.307667 | 211.319643 | 8.106000 | A | 151315 | Train |
| 4 | 1 | 1 | 2010-03-05 | 21827.90 | False | 46.50 | 2.625 | 8536.592778 | 3346.401918 | 1670.797978 | 3653.631444 | 4428.307667 | 211.350143 | 8.106000 | A | 151315 | Train |
| 5 | 1 | 1 | 2010-03-12 | 21043.39 | False | 57.79 | 2.667 | 8536.592778 | 3346.401918 | 1670.797978 | 3653.631444 | 4428.307667 | 211.380643 | 8.106000 | A | 151315 | Train |
| 6 | 1 | 1 | 2010-03-19 | 22136.64 | False | 54.58 | 2.720 | 8536.592778 | 3346.401918 | 1670.797978 | 3653.631444 | 4428.307667 | 211.215635 | 8.106000 | A | 151315 | Train |
| | | | 2010-03-26 | | | | | | | | | | | | | | |

We have total of 45 (store) * 81 (dept) = 3645 different time series in our dataset!!

```
In [15]: 1 ## Let's check the dtypes of the columns
        2 trte_set.dtypes.to_frame().rename(columns={0: 'Type'})
```

Out[15]:

| | Type |
|--------------|----------------|
| Store | int64 |
| Dept | int64 |
| Date | datetime64[ns] |
| Weekly_Sales | float64 |
| IsHoliday | bool |
| Temperature | float64 |
| Fuel_Price | float64 |
| Markdown1 | float64 |
| Markdown2 | float64 |
| Markdown3 | float64 |
| Markdown4 | float64 |
| Markdown5 | float64 |
| CPI | float64 |
| Unemployment | float64 |
| Type | object |
| Size | int64 |
| set_type | object |

```
In [16]: 1 ## check null in final merged
        2 print(trte_set.apply(lambda x: sum(x.isnull()),axis=0))
```

| | |
|--------------|--------|
| Store | 0 |
| Dept | 0 |
| Date | 0 |
| Weekly_Sales | 142155 |
| IsHoliday | 0 |
| Temperature | 0 |
| Fuel_Price | 0 |
| Markdown1 | 0 |
| Markdown2 | 0 |
| Markdown3 | 0 |
| Markdown4 | 0 |
| Markdown5 | 0 |
| CPI | 0 |
| Unemployment | 0 |
| Type | 0 |
| Size | 0 |
| set_type | 0 |
| dtype: | int64 |

Let's find duplicate value in the data based on store,department and date column

```
In [17]: 1 duplicates = len(trte_set[trte_set.duplicated(subset=['Store','Dept','Date'])])
2 print("Does the dataset contain duplicate values : ",bool(duplicates))
```

Does the dataset contain duplicate values : False

- let's perform some validation on dataset if each store has all the unique departments !! and each store and department combination has 182 unique rows of dates !

```
In [18]: 1 print('Unique Counts of departments for each store!!',trte_set.groupby(['Store','Date']).count()['Dept'].unique())
2 print('Number of dates for each store and departments',trte_set.groupby(['Store','Dept']).count()['Date'].unique())
```

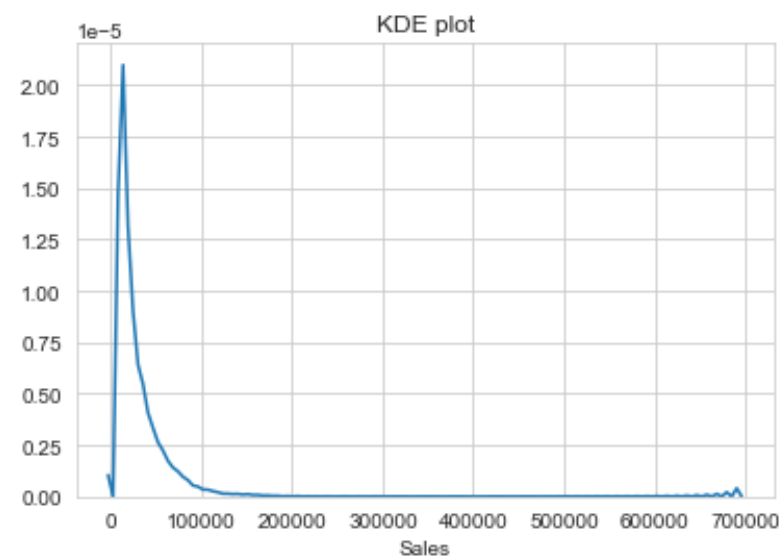
Unique Counts of departments for each store!! [81]

Number of dates for each store and departments [182]

- We can see there are 81 departments for each store and with 182 unique rows of weekly sales.

Analysis on the target column 'Weekly Sales' !!

```
In [19]: 1 ### Let's check the distribution of the weekly sales
2 x=trte_set[trte_set['Weekly_Sales']>0]['Weekly_Sales'] #x~N(0,1)
3 sns.set_style('whitegrid')
4 sns.kdeplot(np.array(x))
5 plt.title('KDE plot')
6 plt.xlabel('Sales')
7 plt.show()
```

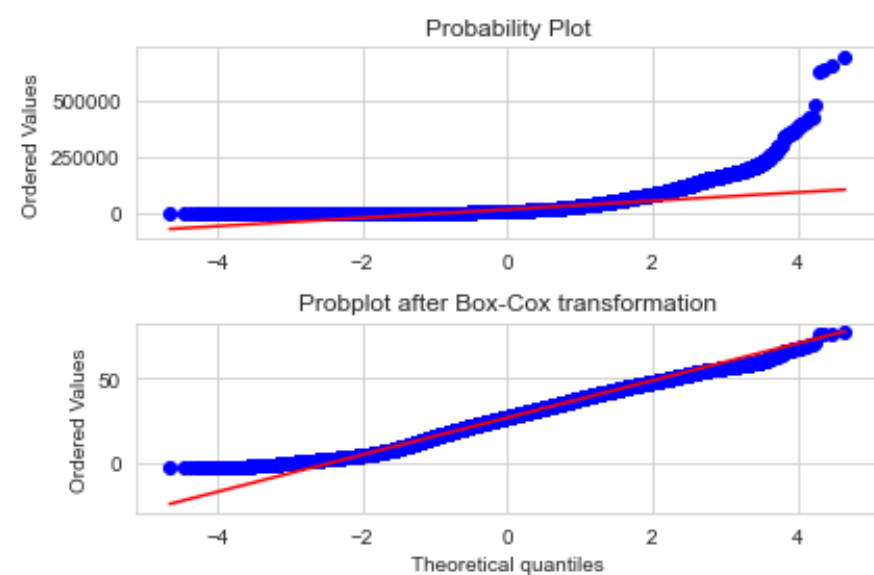


The distribution is **Powerlaw**


```

In [20]: 1  ### performing box-cox transformation to see if we can tranform to gaussian
2  from scipy import stats
3  import matplotlib.pyplot as plt
4  fig=plt.figure()
5  ax1=fig.add_subplot(211)
6  ## plot before box-cox
7  x = trte_set[trte_set['Weekly_Sales']>0]['Weekly_Sales']
8  prob=stats.probplot(x,dist=stats.norm,plot=ax1)
9  ax1.set_xlabel('')
10 ## plot after box-cox
11 ax2 = fig.add_subplot(212)
12 xt, _ = stats.boxcox(x)
13 prob = stats.probplot(xt, dist=stats.norm, plot=ax2)
14 ax2.set_title('Probplot after Box-Cox transformation')
15 fig.tight_layout()
16 plt.show()

```



- As we can observe clearly that the weekly_Sales data follows Power law distribution.
- After performing box-cox transformations we can observe the distribution is directed towards normal-distribution to a good extent. However we are not completely able to make it gaussian.

```
In [21]: 1 #calculating total fare amount values at each percentile 0,10,20,30,40,50,60,70,80,90,100
2 for i in range(0,100,10):
3     var = trte_set["Weekly_Sales"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6 print("100 percentile value is ",var[-1])
```

```
0 percentile value is -4988.94
10 percentile value is 0.0
20 percentile value is 161.09
30 percentile value is 1857.1
40 percentile value is 4732.28
50 percentile value is 9287.7
60 percentile value is 17123.77
70 percentile value is 35148.19
80 percentile value is nan
90 percentile value is nan
100 percentile value is nan
```

```
In [22]: 1 for i in range(0,9,1):
2     var = trte_set["Weekly_Sales"].values
3     var = np.sort(var,axis = None)
4     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
5 print("100 percentile value is ",var[-1])
```

```
0 percentile value is -4988.94
1 percentile value is 0.0
2 percentile value is 0.0
3 percentile value is 0.0
4 percentile value is 0.0
5 percentile value is 0.0
6 percentile value is 0.0
7 percentile value is 0.0
8 percentile value is 0.0
100 percentile value is nan
```

- We see negative weekly sales reported values in the dataset!! Weekly Sales cannot be negative!

```
In [23]: 1 print('Number of rows which has negative sales reported == ',len(trte_set.loc[trte_set['Weekly_Sales']<0]))
```

```
Number of rows which has negative sales reported == 1285
```

- Transform negatively reported sales to 0 .

```
In [24]: 1 trte_set.loc[trte_set['Weekly_Sales']<0,'Weekly_Sales'] = 0
```

```
In [25]: 1 print('Number of negative sales now in set!!',len(trte_set.loc[trte_set['Weekly_Sales']<0]))
```

```
Number of negative sales now in set!! 0
```

- Let's check time series (Store and Dept combos) which has zero recorded sales for all the dates , if yes then delete those series

```
In [26]: 1  ### Check time series with all 0 weekly sales values!!
2  grp = trte_set.groupby(['Store', 'Dept'])['Weekly_Sales'].sum()
3  print('Number of time series with all zero sales == ', len(grp[grp==0]))
4  store_dept_with_all_zeros_sales = list(grp[grp==0].index)
```

Number of time series with all zero sales == 322

```
In [27]: 1  ## reset the index
2  trte_set = trte_set.reset_index(drop=True)
3  for i,j in store_dept_with_all_zeros_sales:
4      trte_set.drop(trte_set[(trte_set.Store == i) & (trte_set.Dept == j)].index ,inplace=True)
```

```
In [28]: 1  print('Total number of rows deleted == ', 182 * len(store_dept_with_all_zeros_sales))
2  print('Length of train and test after deletion ==', len(trte_set))
```

Total number of rows deleted == 58604
Length of train and test after deletion == 604786

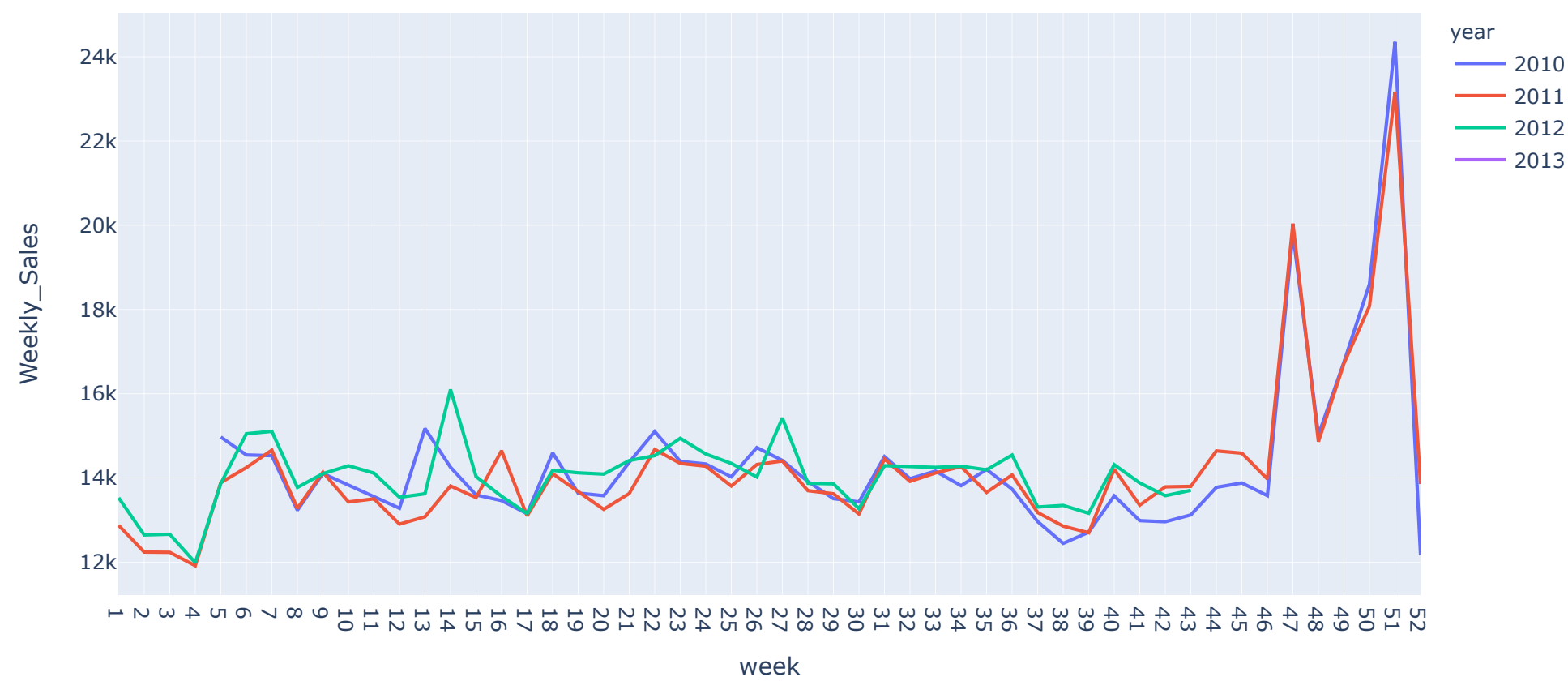
3.4 Univarite analysis

Weekly Sales x Date

```
In [47]: 1  ## Let's craete two date features here week,year for analysis
2  trte_set['week'] = trte_set['Date'].dt.week
3  trte_set['year'] = trte_set['Date'].dt.year
4  sales_per_week_year = trte_set.groupby(['Date', 'week', 'year'], as_index=False)['Weekly_Sales'].mean()
```

```
In [48]: 1  ## Plotting mean sales over years(2010-2012)
2  ## we are plotting using plotly library to get more readability
3  fig = px.line(sales_per_week_year, y="Weekly_Sales",x='week',
4               title='Mean Sales - per week over the years',color='year',line_group='year',)
5  fig.update_layout(
6      xaxis = dict(
7          tickmode = 'array',
8          tickvals = np.arange(1,53) ,
9          ticktext = np.arange(1,53)
10     )
11 )
12 fig.show()
13
```

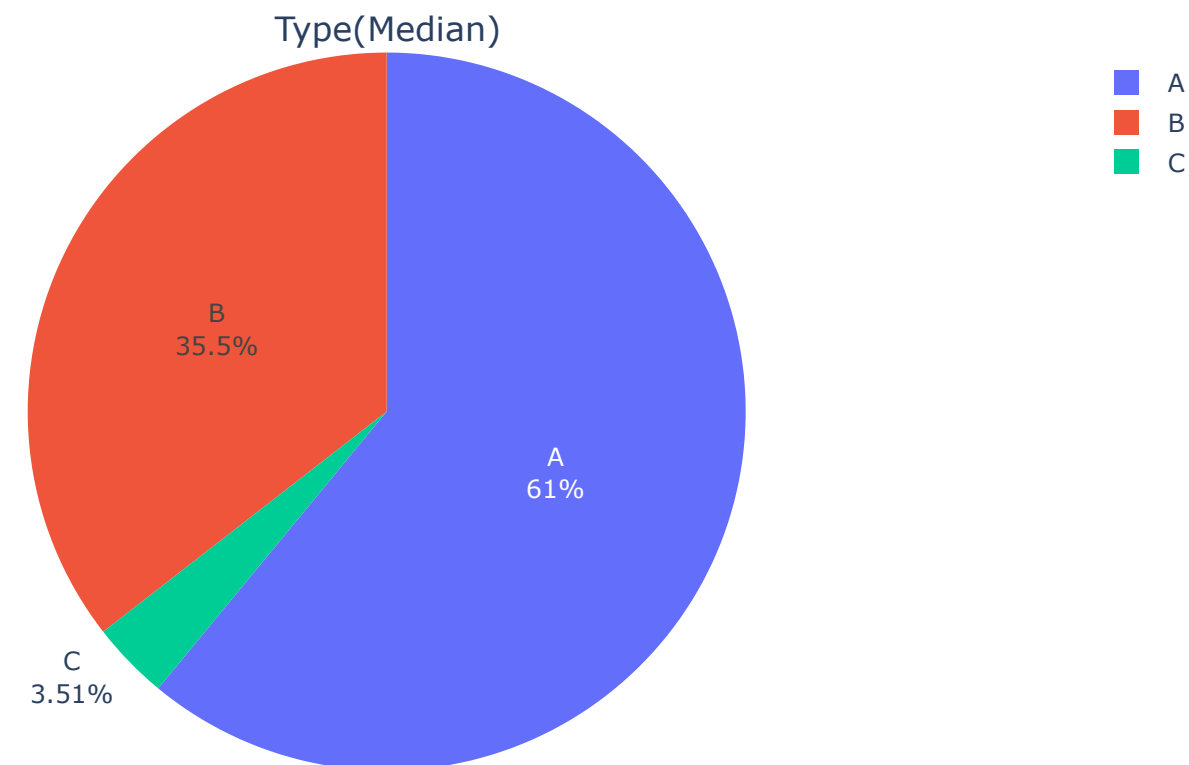
Mean Sales - per week over the years



- We can observe that during thanksgiving week(Nov 26) and christmas week(December 24) soaring peaks of sales.
- Apart from that we can observe peaks during easter week i.e April 2,2010 , April 22 ,2011 and April 6,2012 .However,these are not included as holiday weeks in our datasets.So we will include them as holiday week.
- Holidays like superbowl occurs in first sunday of feb month,labour day on first monday of september !! It will be useful to include features like week number, day, month's week number
- Over the years 2010-2012 the time-series data followed a similar pattern.

Type X Sales


```
In [49]: 1 ## Let's plot sales per store
2 ## Using Plotly we are creating pie chart !! Let's plot the mean sales for every week
3 sales_per_type = trte_set.groupby(['Type'],as_index=False)['Weekly_Sales'].median()
4 fig = make_subplots(rows=1, cols=1,subplot_titles=['Type(Median)', 'Store(Median)'],
5               specs=[[{"type": "pie"}]])
6 fig.add_trace(go.Pie(labels=sales_per_type['Type'], values=sales_per_type['Weekly_Sales'],text=sales_per_type['Type']),row=1,col=1)
7 fig.show()
```



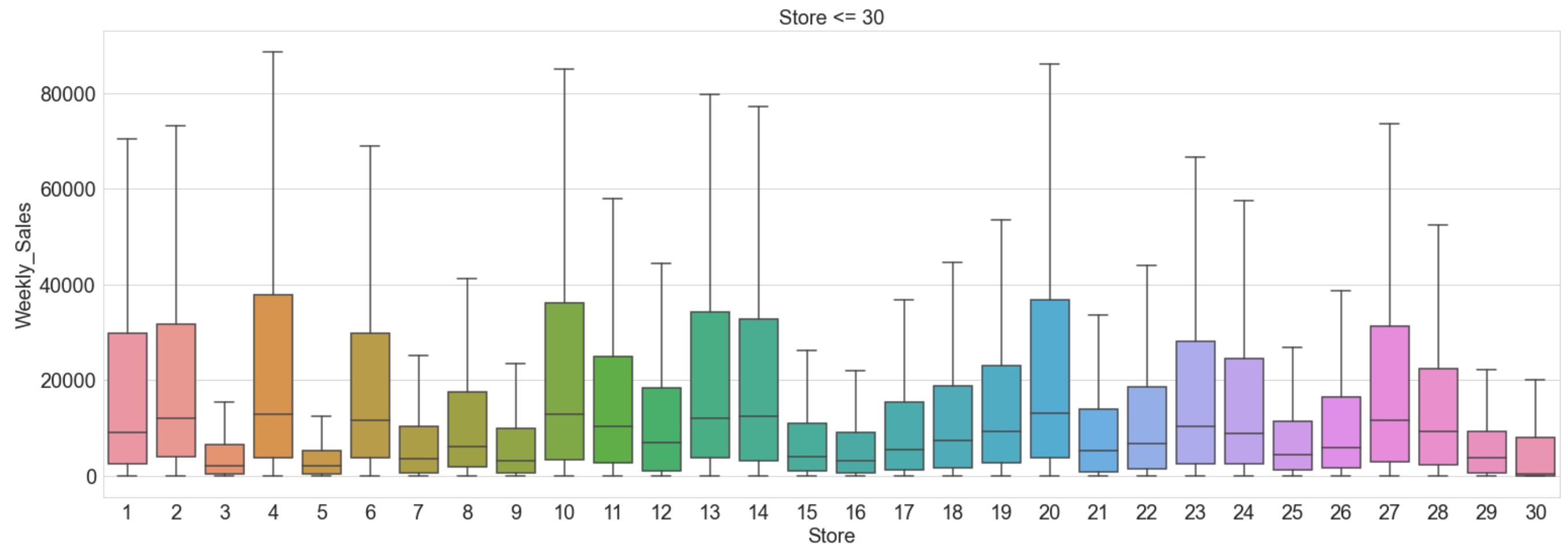
As observed from the pie-cart

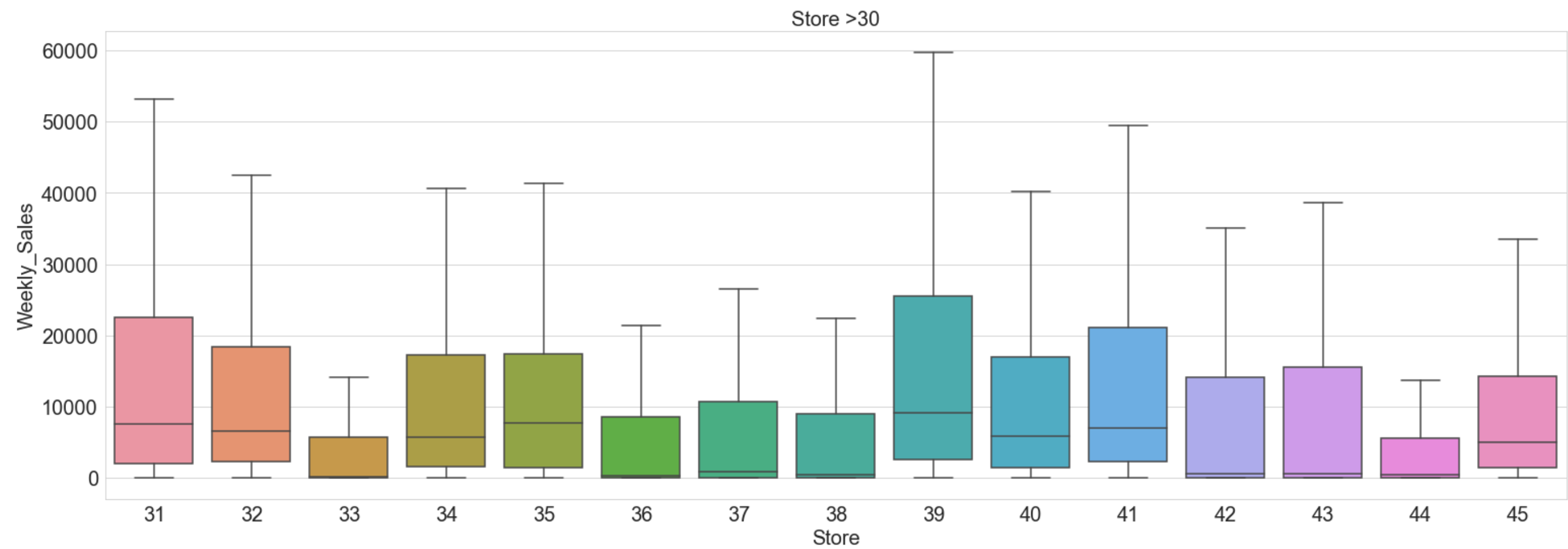
- From the plot we can see that type "A" > Type "B" > Type "C" in producing sales.

Weekly Sales x Store

```
In [50]: 1 a=trte_set.query('Store <= 30' )
2 b=trte_set.query('Store >30 and Store <=60')
3 fig,axes = plt.subplots(2,1,figsize=(25,20))
4 fig.suptitle('Store Wise Sales',fontsize=25)
5
6 fig.subplots_adjust(wspace=0.5,hspace=0.5)
7 fig.subplots_adjust(left=0.1,right=0.9,bottom=0.1,top=0.9)
8 ### subplot for Store<30
9 sns.boxplot(x=a['Store'],y=a['Weekly_Sales'],showfliers=False,ax=axes[0])
10 axes[0].set_title('Store <= 30',fontsize=20)
11 axes[0].set_xlabel('Store',fontsize=20)
12 axes[0].set_ylabel('Weekly_Sales',fontsize=20)
13 axes[0].tick_params(labelsize=20)
14
15 ### Store >30
16 sns.boxplot(x=b['Store'],y=b['Weekly_Sales'],showfliers=False,ax=axes[1])
17 axes[1].set_title('Store >30',fontsize=20)
18 axes[1].set_xlabel('Store',fontsize=20)
19 axes[1].set_ylabel('Weekly_Sales',fontsize=20)
20 axes[1].tick_params(labelsize=20)
21
22
23
```

Store Wise Sales





- Observations :
 - Few stores show higher mean sales and higher interquartile range depicting that few stores are bigger than the others .
 - Some stores have their 25th and 50th percentile values very close and 75th percentile is far away showing that there is not much sales values difference between point below 25% and point below 50% .There was sudden increase in sales after 50 % of points.(might be flier points showing sudden increase in sales due to some external factor)
 - Theory : We can understand that the stores with the higher weekly sales(1,2,4,10...) can be categorized as type 'A' stores with medium sales as type 'B'(11,2,19,20...) and the rest as type 'C'(3,5,29,30...)

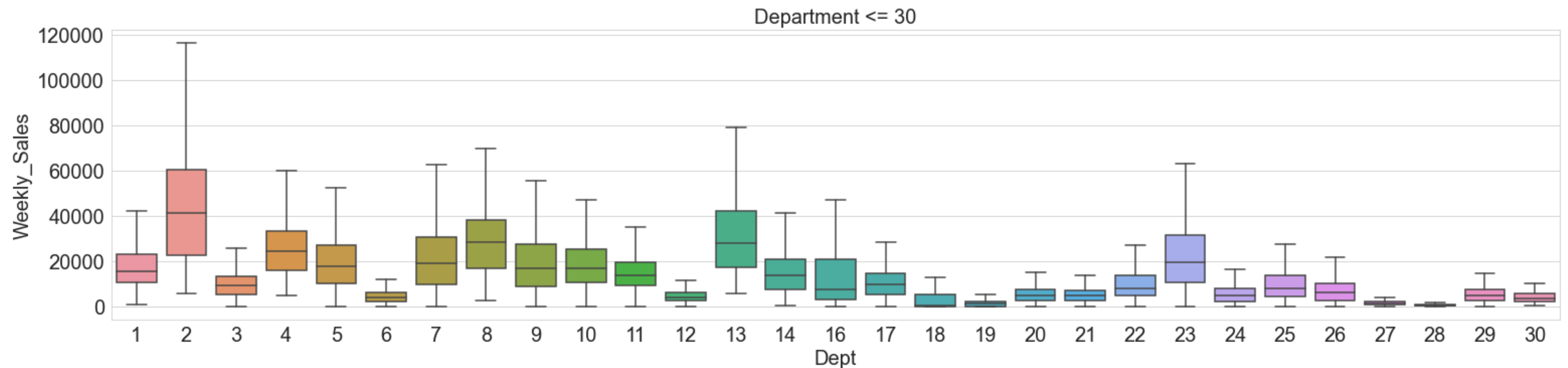
Weekly Sales x Dept

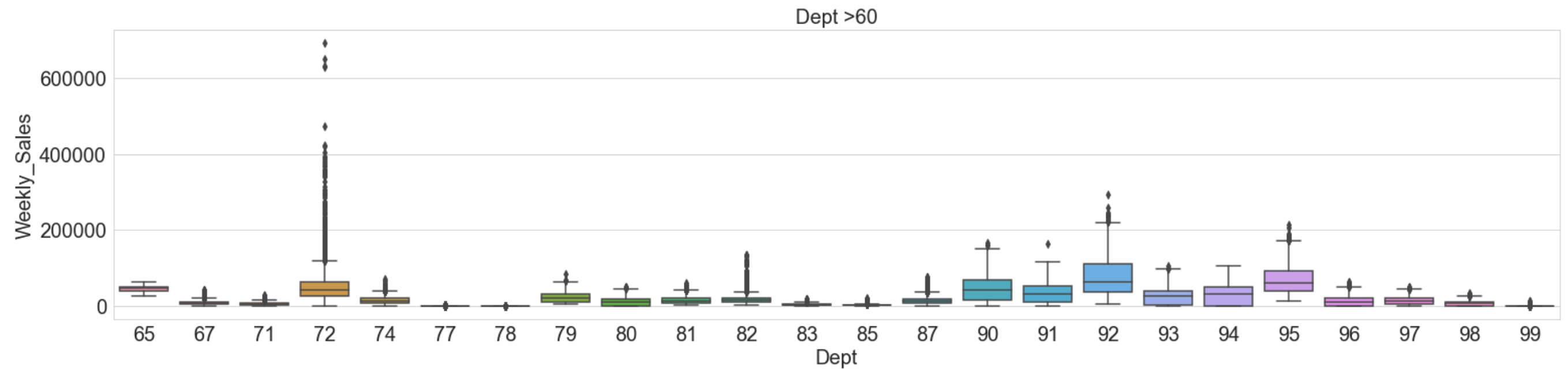
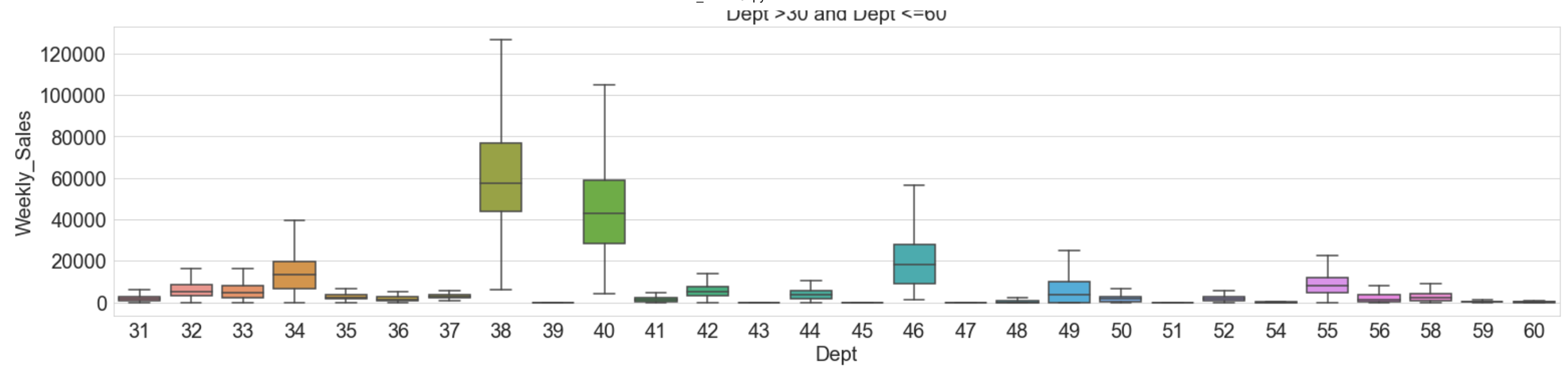
```

In [51]: 1 a=trte_set.query(' Dept <= 30' )
2 b=trte_set.query('Dept >30 and Dept <=60')
3 c=trte_set.query('Dept >60 ')
4 fig,axes = plt.subplots(3,1,figsize=(25,20))
5 fig.suptitle('Department Wise Sales',fontsize=25)
6
7 fig.subplots_adjust(wspace=0.5,hspace=0.5)
8 fig.subplots_adjust(left=0.1,right=0.9,bottom=0.1,top=0.9)
9 ### subplot for dept<30
10 sns.boxplot(x=a['Dept'],y=a['Weekly_Sales'],showfliers=False,ax=axes[0])
11 axes[0].set_title('Department <= 30',fontsize=20)
12 axes[0].set_xlabel('Dept',fontsize=20)
13 axes[0].set_ylabel('Weekly_Sales',fontsize=20)
14 axes[0].tick_params(labelsize=20)
15
16 ### Dept >30 and Dept <=60
17 sns.boxplot(x=b['Dept'],y=b['Weekly_Sales'],showfliers=False,ax=axes[1])
18 axes[1].set_title('Dept >30 and Dept <=60',fontsize=20)
19 axes[1].set_xlabel('Dept',fontsize=20)
20 axes[1].set_ylabel('Weekly_Sales',fontsize=20)
21 axes[1].tick_params(labelsize=20)
22
23 # Dept >60
24 sns.boxplot(x=c['Dept'],y=c['Weekly_Sales'],ax=axes[2])
25 axes[2].set_title('Dept >60',fontsize=20)
26 axes[2].set_xlabel('Dept',fontsize=20)
27 axes[2].set_ylabel('Weekly_Sales',fontsize=20)
28 axes[2].tick_params(labelsize=20)
29
30

```

Department Wise Sales

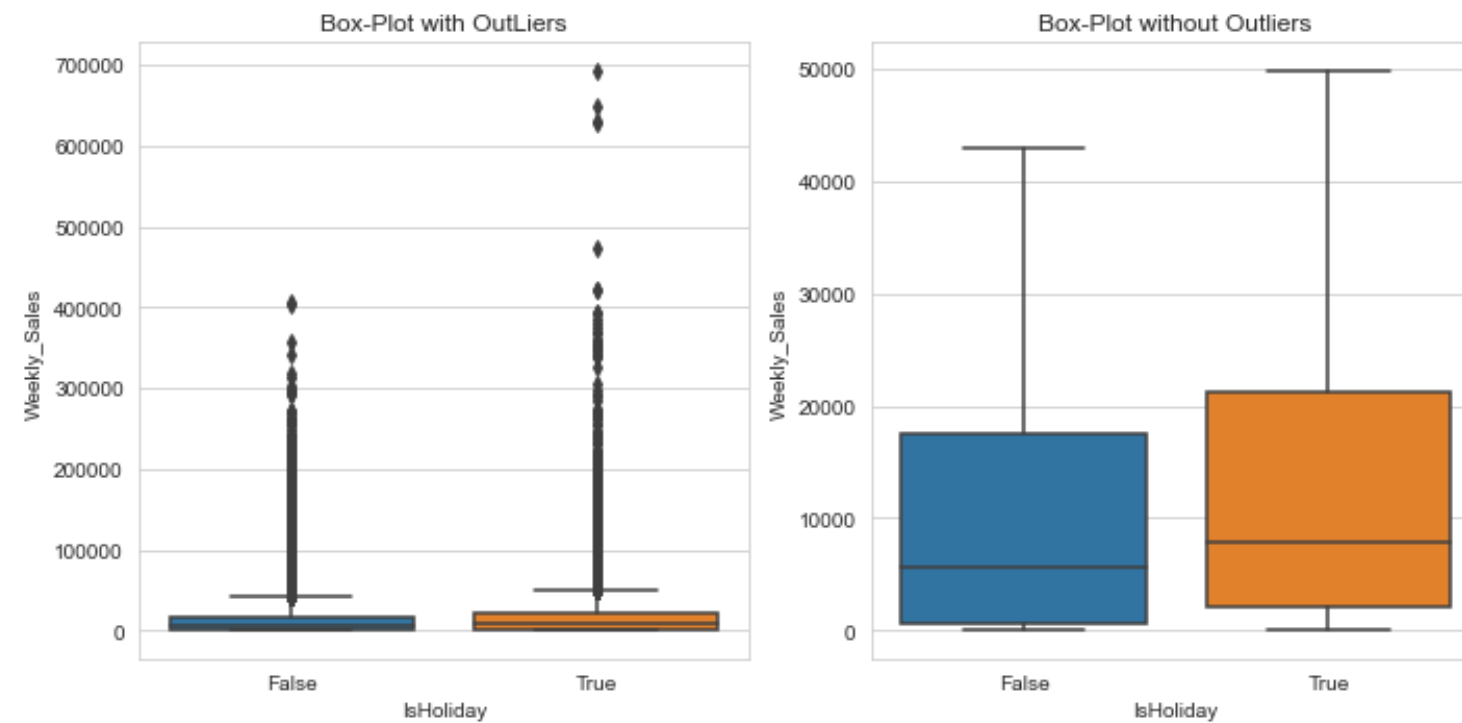




- We have calculated the mean sales per department and plotted it separately for readability.
- From this bar plot we can see department 2,38,40 has greater mean sales compared to other departments which shows that these department have more useful products for the customer. People are more likely to buy products from these departments.
- Department 72,7,5 we can observe high sales during few weeks which is farther away from their mean sales. There is a possibility that these departments supply customers with products which are needed during Holiday Weeks.

Weekly Sales on Holidays

```
In [52]: 1 fig = plt.figure(figsize=(10,5))
2 plt.subplot(1,2,1)
3 sns.boxplot(data=trte_set,x='IsHoliday',y='Weekly_Sales',)
4 plt.title("Box-Plot with OutLiers")
5 plt.subplot(1,2,2)
6 sns.boxplot(data=trte_set,x='IsHoliday',y='Weekly_Sales',showfliers=False)
7 plt.title('Box-Plot without Outliers')
8 plt.tight_layout()
```



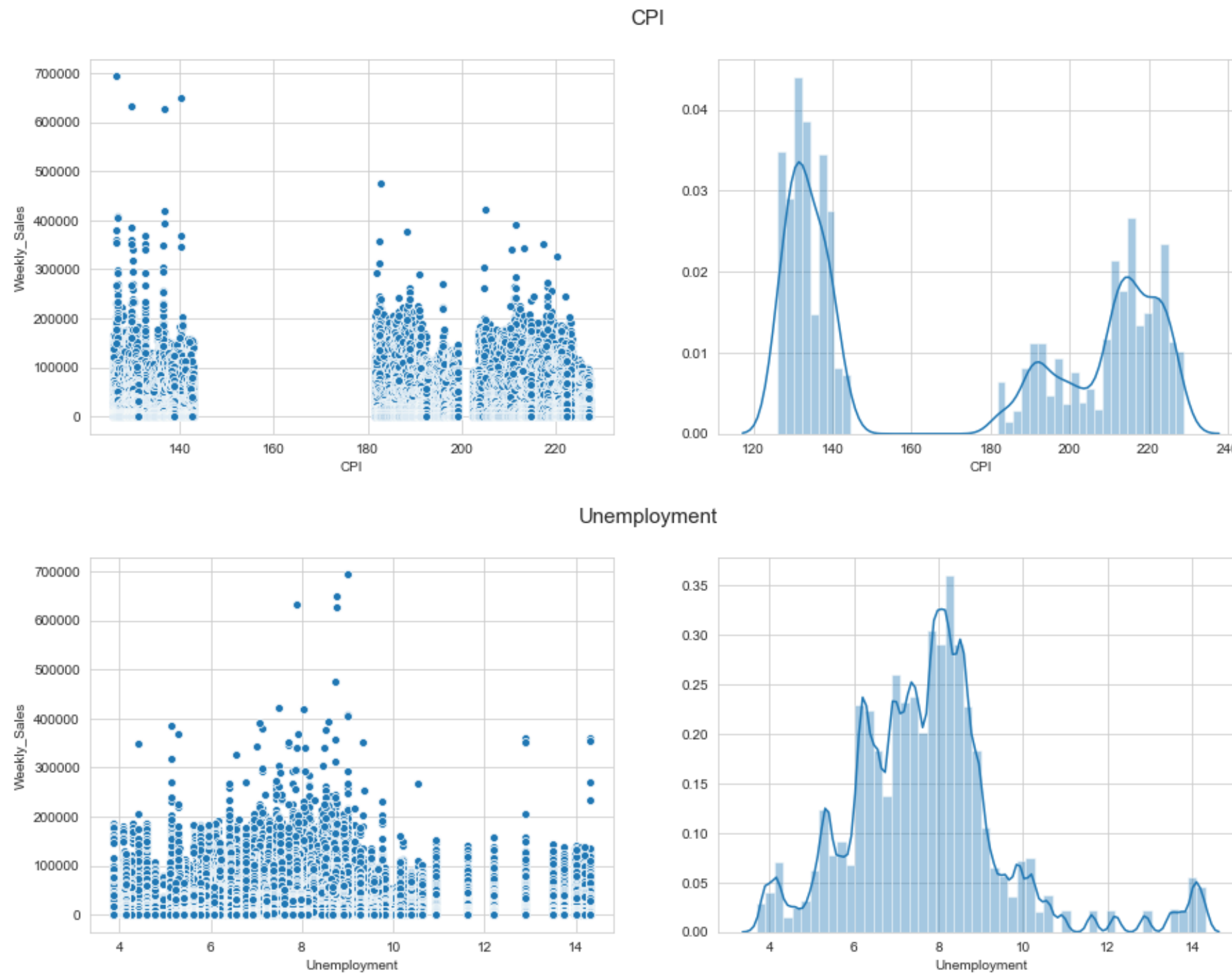
Observation 1 : As we can see outlier points on Holidays are greater than non-Holidays. The mean sales/inter-quartile range values for both Holiday and nonholiday weeks are overlapping, which doesn't make 'IsHoliday' a good predictor. Which indicates are only a few Holiday weeks which clearly increases sales values than non-Holiday WEEKs.

3.5 Skewness Measure and Outlier Detection

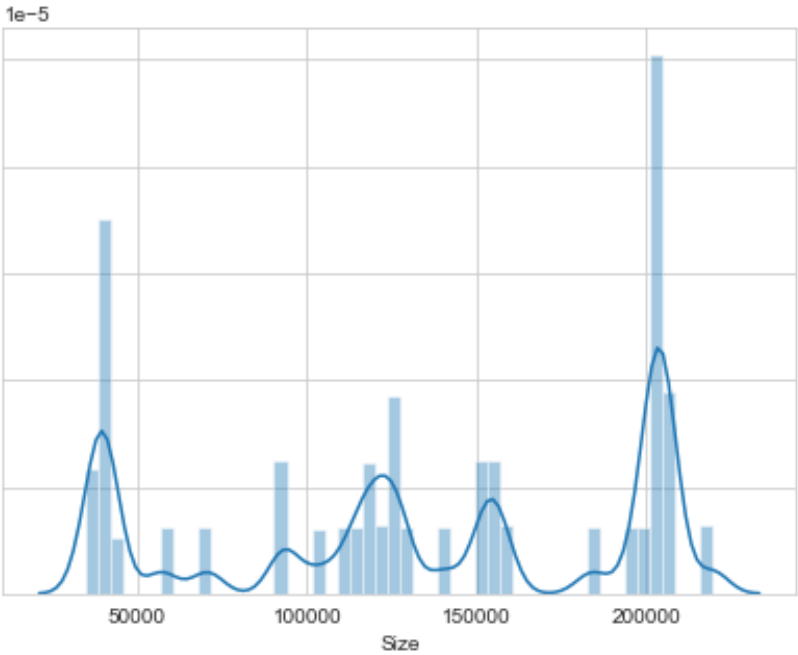
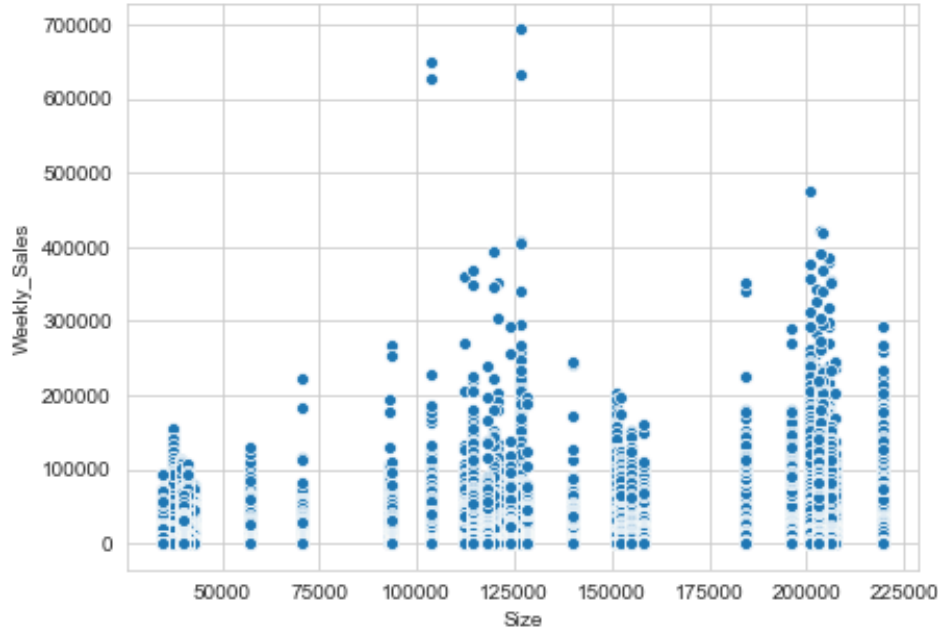
- Linear Regression Models makes assumptions such as the data is normally distributed, without noise, has no-collinearity which is why Skewness in data reduces the predictive power of our models,
- Let's Analyse some skewness in data.

```
In [53]: ▶ 1 ### for the numerical variable CPI,employment,fuel_price Lets plot kde and scatter plot
2 def kde_scatter(data,x,y,title):
3     '''Funcation plots KDE plot and Scatter plot as subplots '''
4     fig, axes = plt.subplots(1, 2, figsize=(15, 5))
5     fig.suptitle(title,fontsize=15)
6     sns.scatterplot(x=data[x],y=data[y],ax=axes[0])
7     sns.distplot(data[x], kde=True,ax=axes[1])
8     plt.show()
```

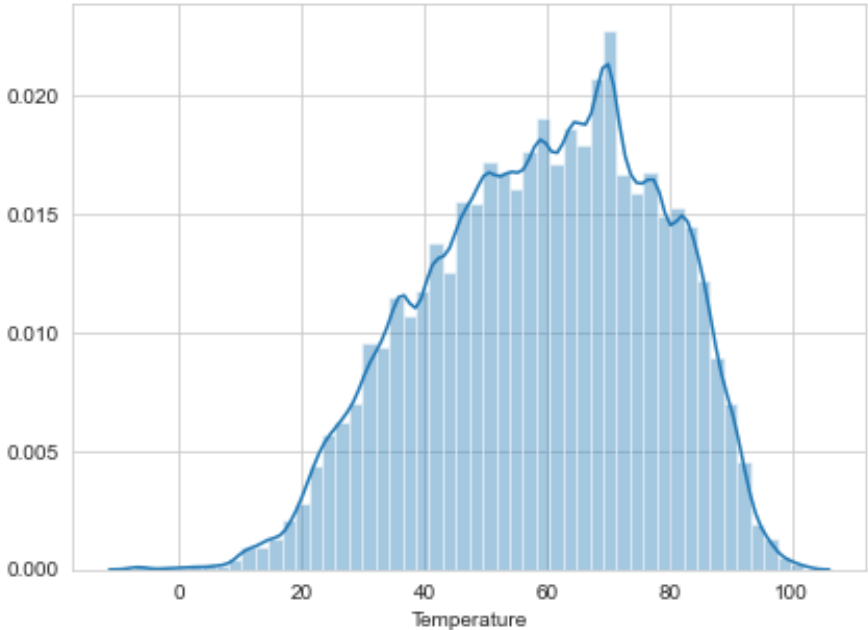
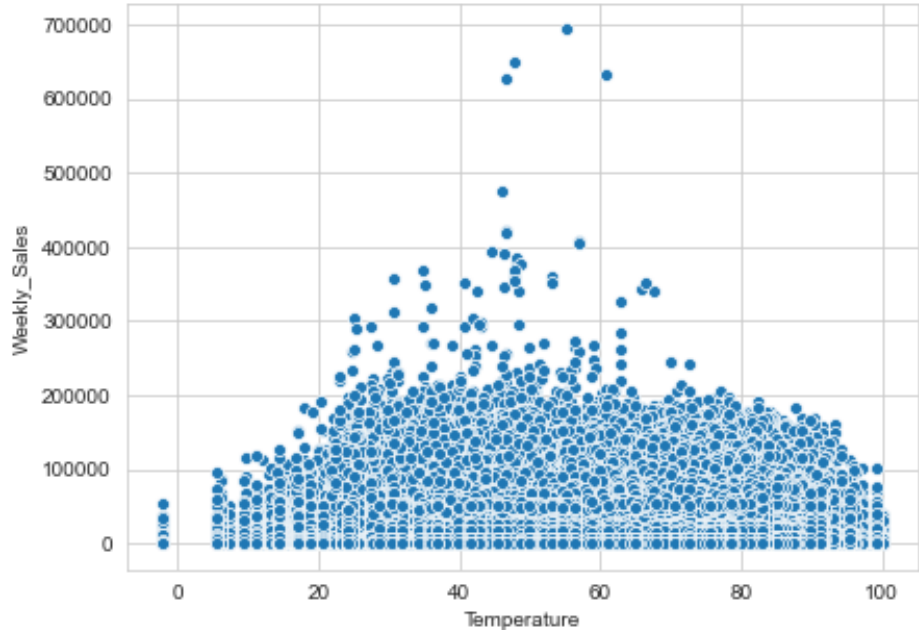
```
In [54]: 1 ## plotting KDE and scatter for Numerical Variables
2 kde_scatter(trte_set, 'CPI', 'Weekly_Sales', "CPI")
3 kde_scatter(trte_set, 'Unemployment', 'Weekly_Sales', "Unemployment")
4 kde_scatter(trte_set, 'Size', 'Weekly_Sales', 'Size')
5 kde_scatter(trte_set, 'Temperature', 'Weekly_Sales', 'Temperature')
6 kde_scatter(trte_set, 'Markdown1', 'Weekly_Sales', 'Markdown1')
7 kde_scatter(trte_set, 'Markdown2', 'Weekly_Sales', 'Markdown2')
8 kde_scatter(trte_set, 'Markdown3', 'Weekly_Sales', 'Markdown3')
9 kde_scatter(trte_set, 'Markdown5', 'Weekly_Sales', 'Markdown5')
```



Size



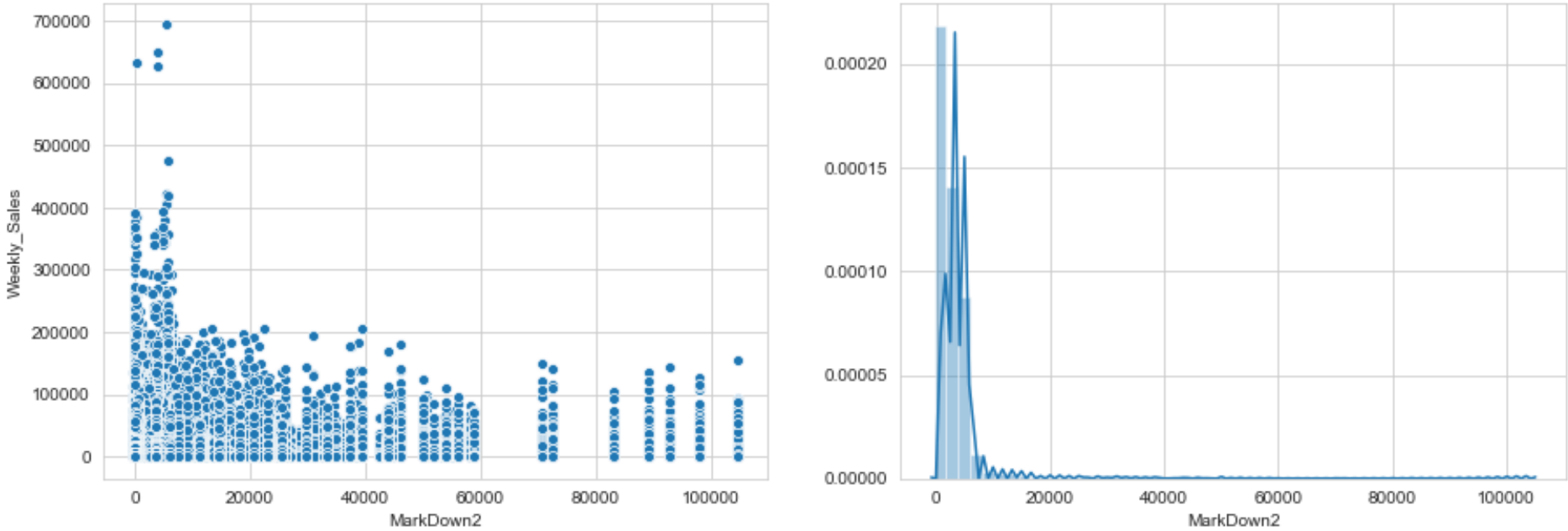
Temperature



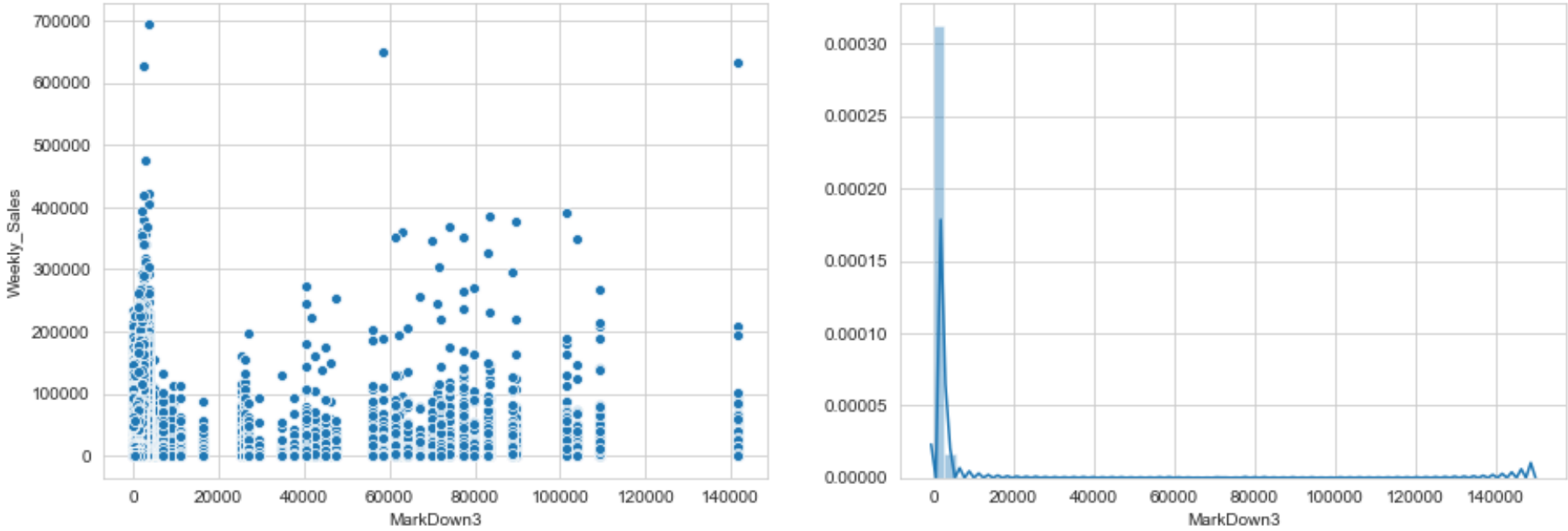
Markdown1

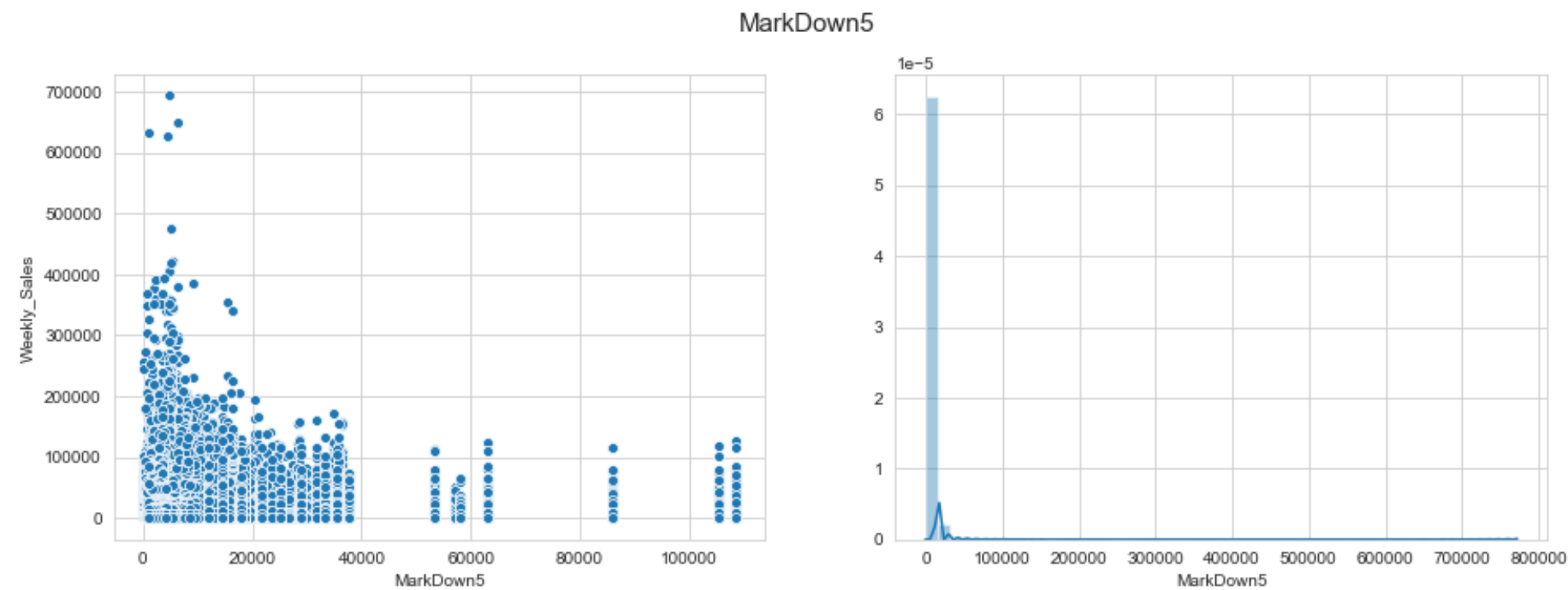


Markdown2



Markdown3





Observation: We can see all the numerical features are Right-skewed

Theory: We don't have to take in account of normalizing features as we are not implementing any linear models (assumes disribution to be gaussian).

Let's try to reduce skewness !! P.S This is just for observation purpose. We need not apply any normalizing techniques.

```
In [55]: 1 from scipy.stats import skew
2 numerical_vars = ['CPI','Unemployment','Temperature','Size','Markdown1','Markdown2','Markdown3','Markdown4','Markdown5']
3 ## printing skewness for each variable
4 for var in numerical_vars:
5     print("Column name:",var,'\t',"Skew:",skew(trte_set[var]),' ',' neg-values',
6           bool(len(trte_set.query('{} < 0'.format(var)))))
7
```

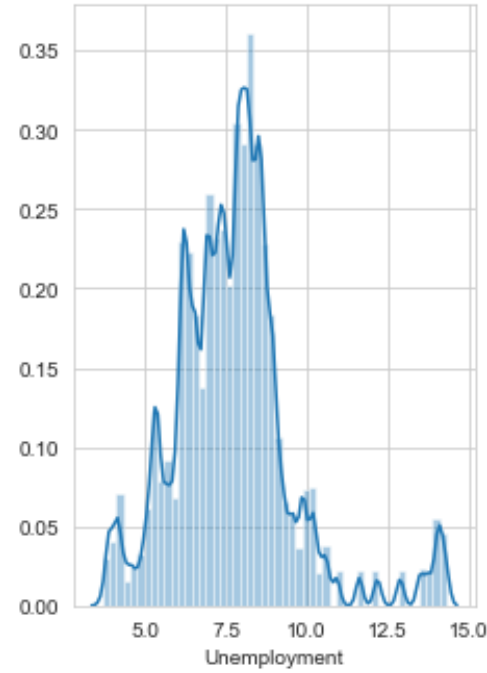
```
Column name: CPI          Skew: 0.08296710665638728    neg-values False
Column name: Unemployment  Skew: 1.0150771848243956    neg-values False
Column name: Temperature   Skew: -0.2741277765467456    neg-values True
Column name: Size          Skew: -0.2713788675521774    neg-values False
Column name: Markdown1     Skew: 4.4791479745936424    neg-values True
Column name: Markdown2     Skew: 7.46910934822387    neg-values True
Column name: Markdown3     Skew: 11.883689905525937    neg-values True
Column name: Markdown4     Skew: 6.499750652818528    neg-values False
Column name: Markdown5     Skew: 65.3566601396212    neg-values True
```

- Since Markdown 2,3 and Temperature contain negative values we cannot perform log tranform
- let's consider the columns without negative vals.

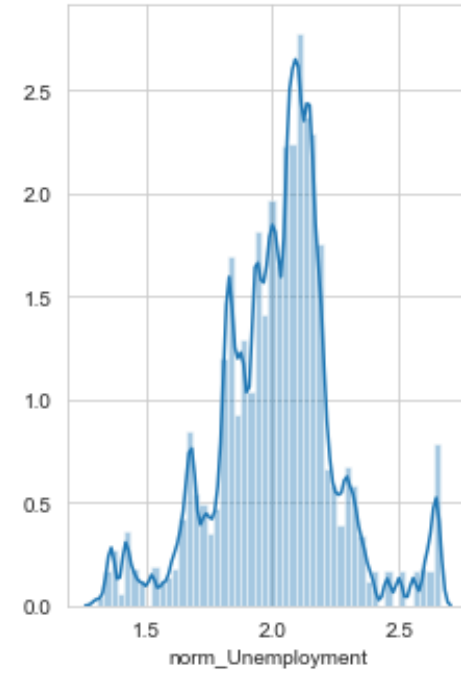
```
In [56]: ▶ 1 def skew_corr_plot(num_var):
2           '''Function to visualize reduction of skewness after performing log-transform'''
3           for var in num_var:
4               fig = plt.figure(figsize=(15,5))
5               c = np.log(trte_set[var]) ## perform Log-transform
6               trte_set['norm_{}'.format(var)] = c ## create a new norm feature
7               plt.subplot(141)
8               sns.distplot(trte_set[var],kde=True) ### plot kde
9               ## title to represent skewness of var and correlation of var
10              plt.title('Skew: {} Corr: {}'.format(round(skew(trte_set.query('{}>0'.format(var))[var]),4),
11                                                    round(trte_set.corr()['Weekly_Sales'][var],4)),fontdict={'fontsize':20})
12
13              plt.subplot(142)
14              sns.distplot(trte_set['norm_{}'.format(var)],kde=True)
15              ## title to represent skewness of var and correlation of var
16              plt.title('Skew: {} Corr: {}'.format(round(skew(trte_set.query('norm_{}>0'.format(var))['norm_{}'.format(var)]),4),
17                                                    round(trte_set.corr()['Weekly_Sales']['norm_{}'.format(var)],4)),fontdict={'fontsize':20})
18
19              ax=plt.subplot(143)
20              stats.probplot(trte_set[var], dist=stats.norm,plot=ax) ## probability plot before transform
21              ax=plt.subplot(144)
22              stats.probplot(trte_set['norm_{}'.format(var)], dist=stats.norm,plot=ax ) ## probability plot after transform
23              plt.tight_layout()
24              plt.show()
```

```
In [57]: 1 skew_corr_plot(['Unemployment', 'Size', 'CPI', 'Markdown1', 'Markdown4', 'Markdown5'])
```

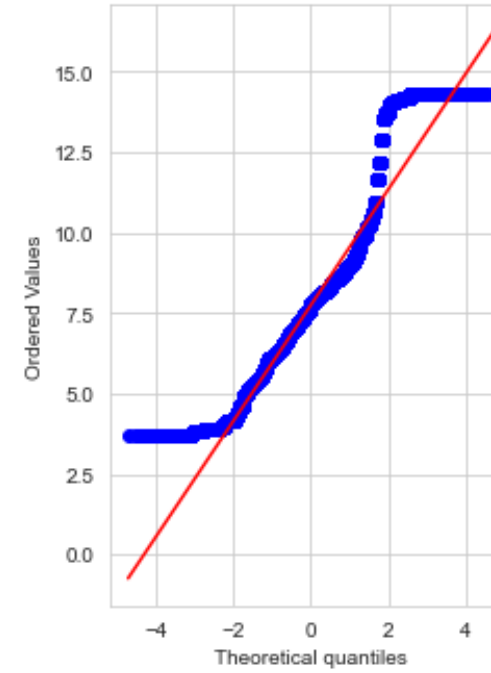
Skew: 1.0151 Corr: -0.0288



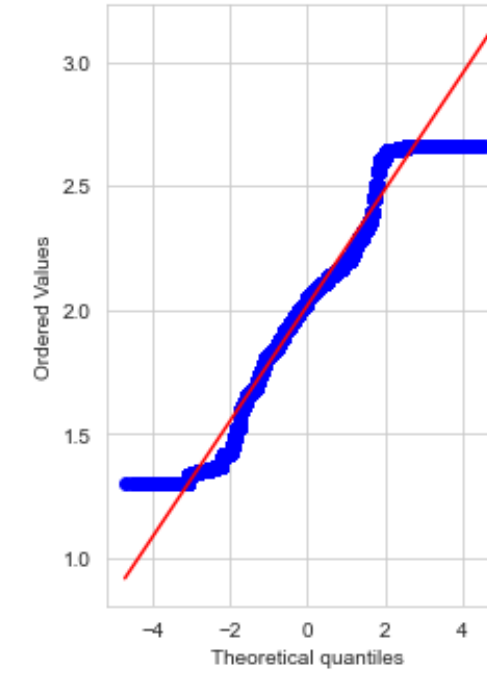
Skew: -0.0991 Corr: -0.0284



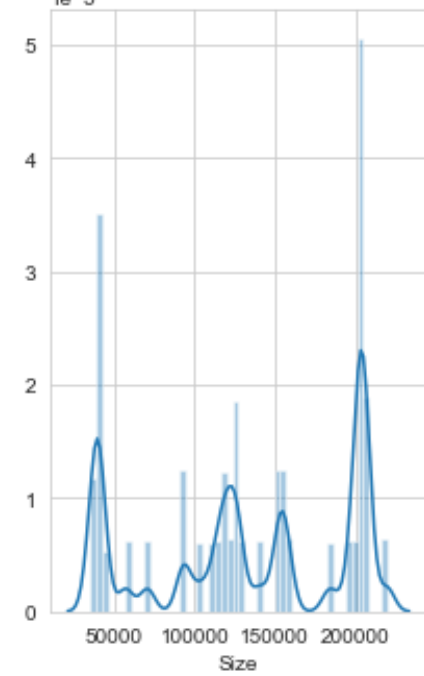
Probability Plot



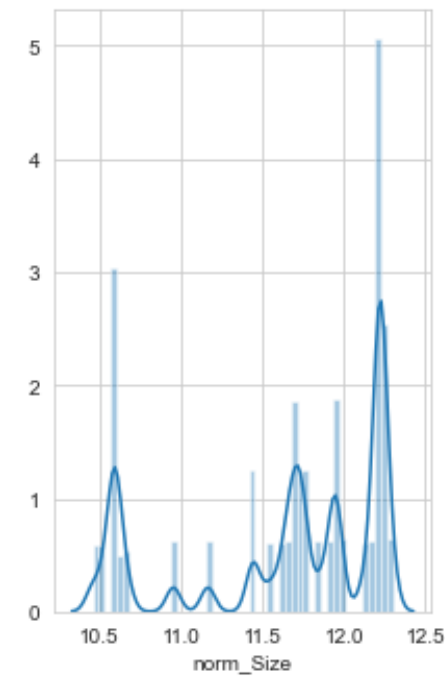
Probability Plot



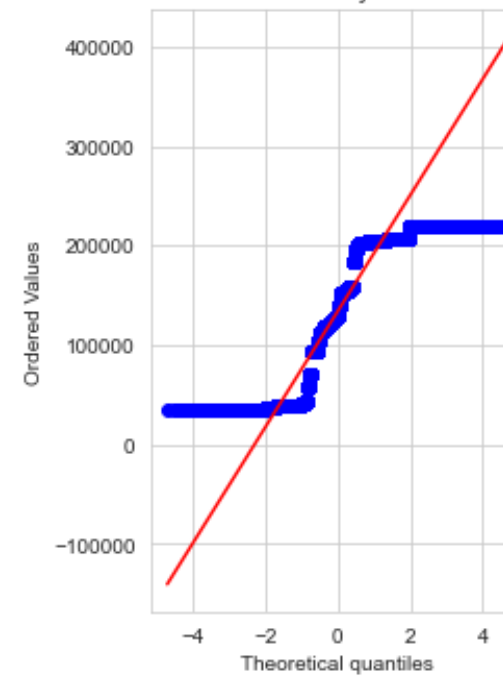
Skew: -0.2714 Corr: 0.2477



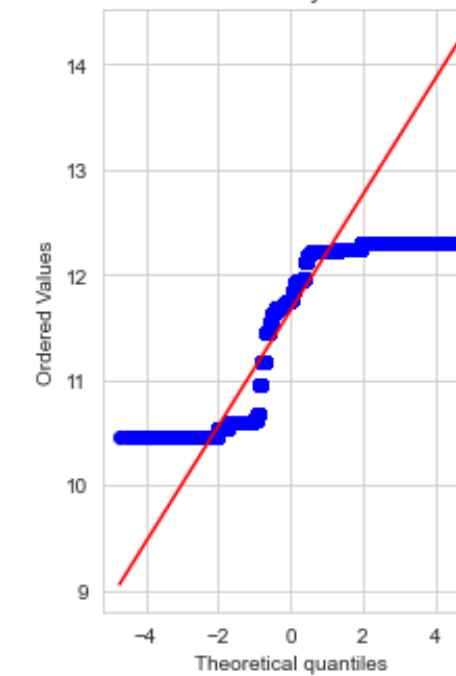
Skew: -0.8335 Corr: 0.2351



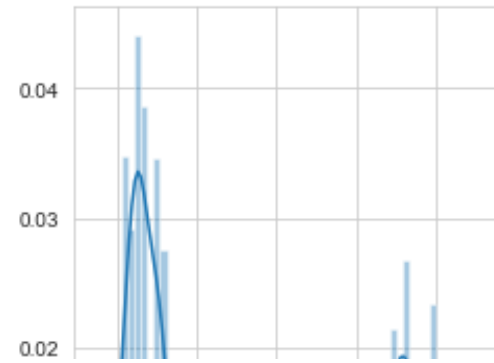
Probability Plot



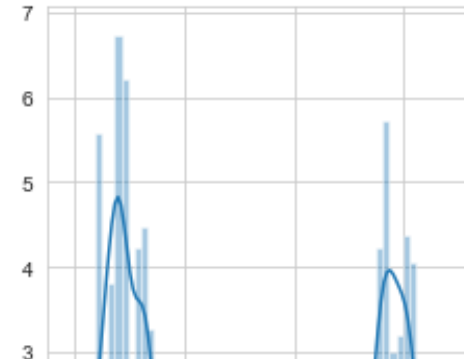
Probability Plot



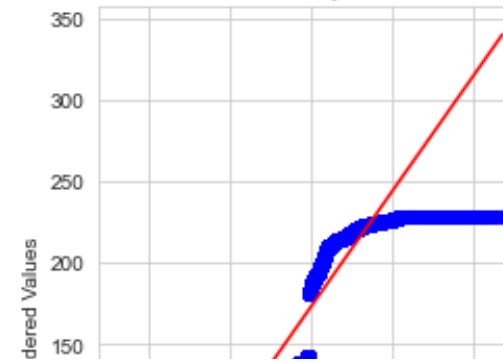
Skew: 0.083 Corr: -0.0214



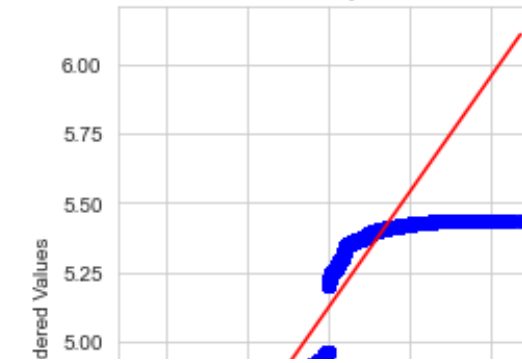
Skew: 0.029 Corr: -0.0208



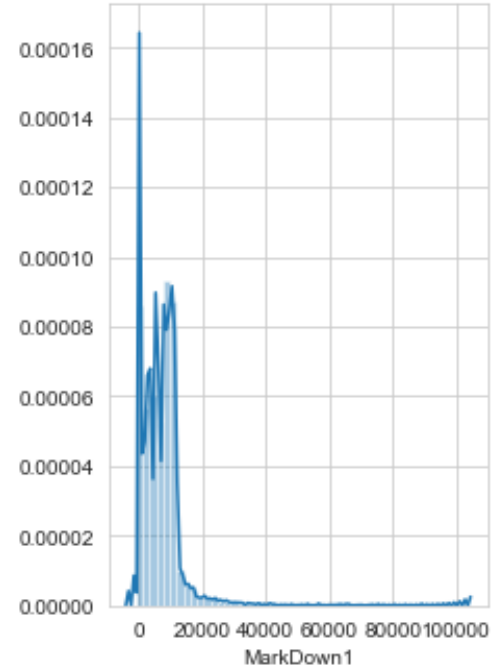
Probability Plot



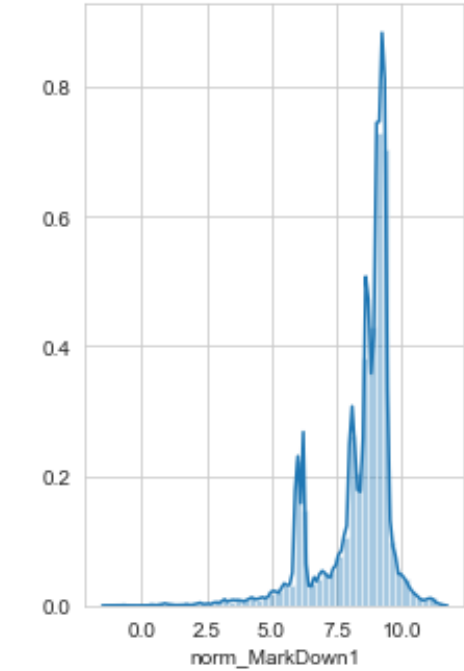
Probability Plot



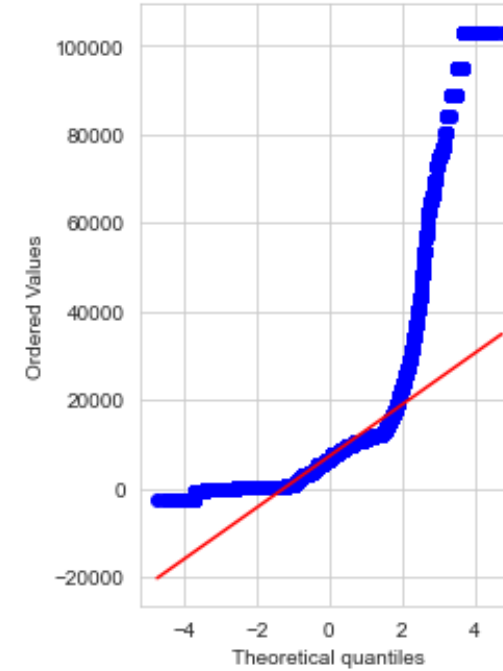
Skew: 4.4814 Corr: 0.1489



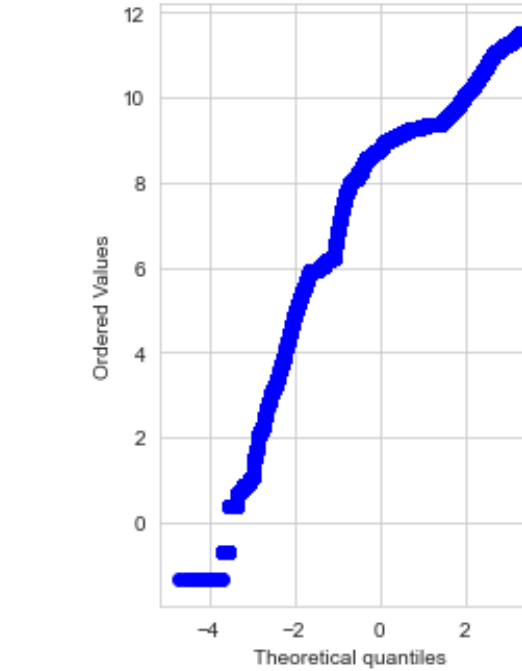
Skew: -1.5912 Corr: 0.1611



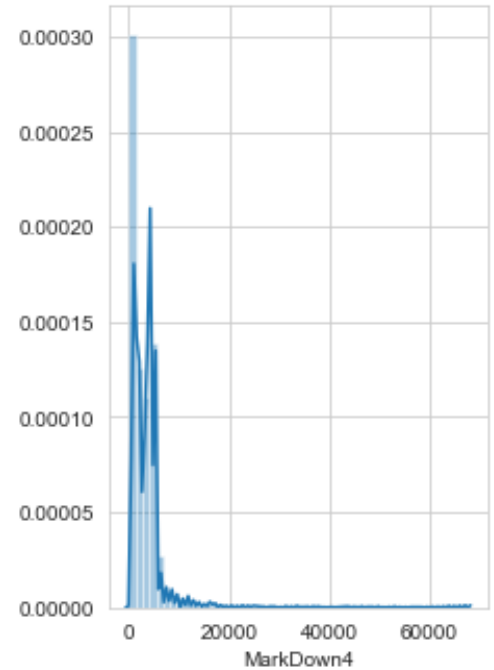
Probability Plot



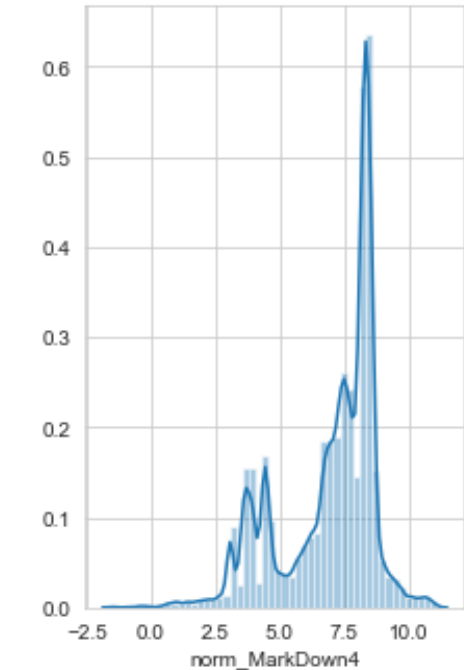
Probability Plot



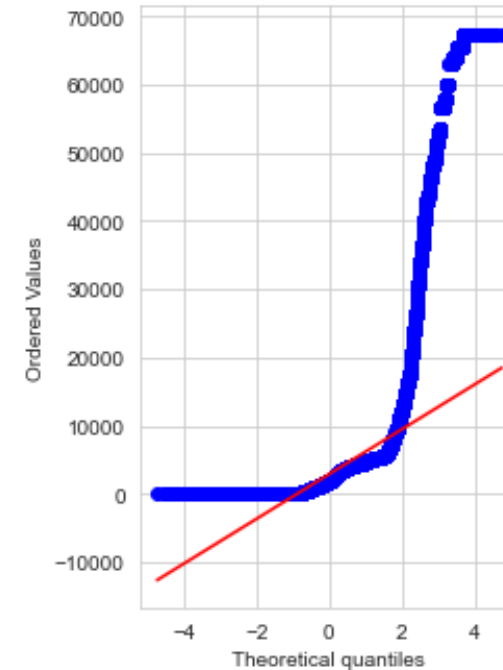
Skew: 6.4998 Corr: 0.1108



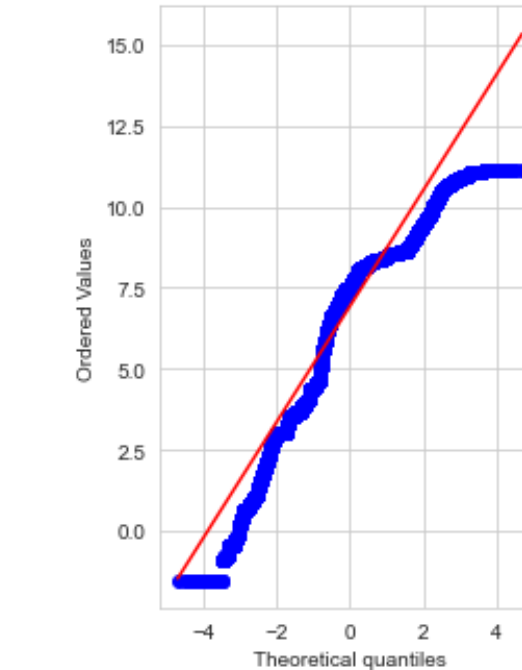
Skew: -0.9114 Corr: 0.1612



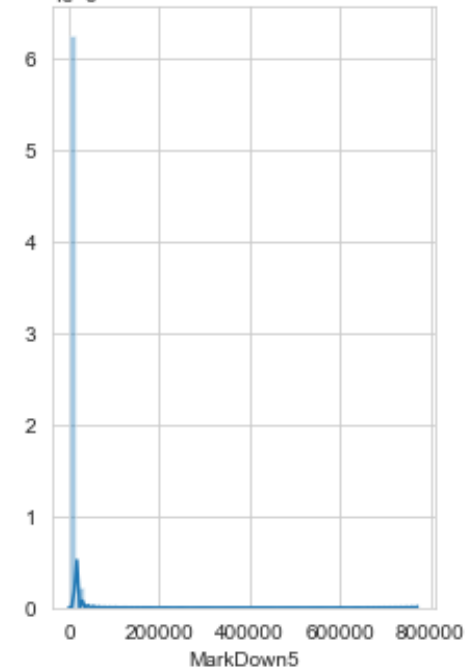
Probability Plot



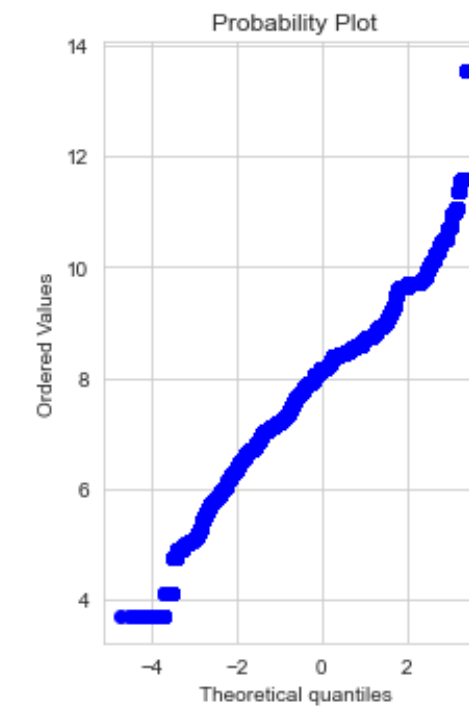
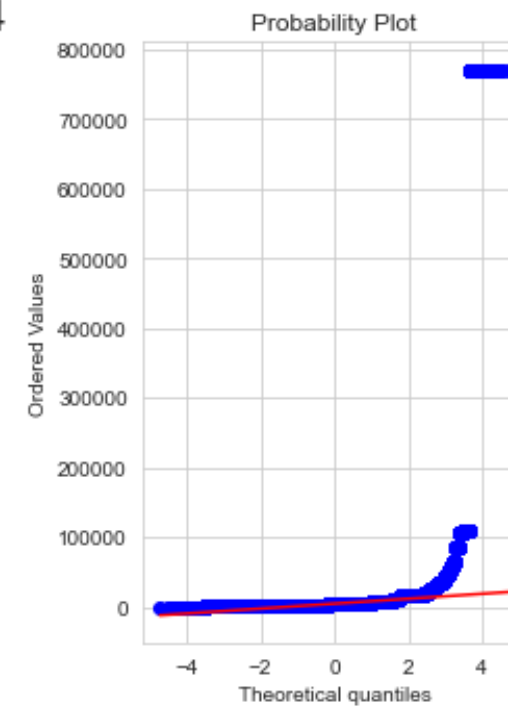
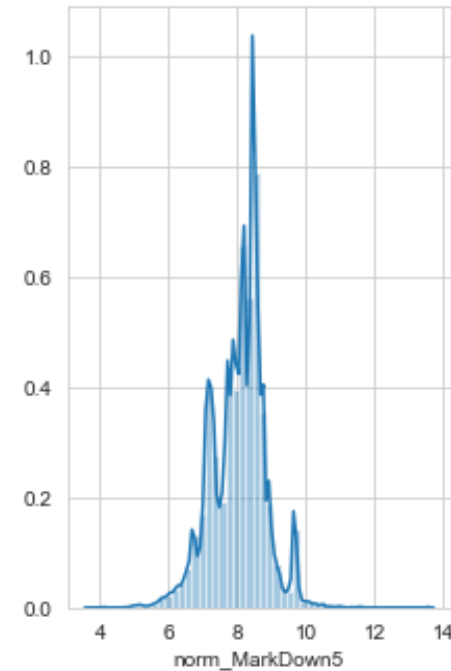
Probability Plot



Skew: 65.3532 Corr: 0.1206



Skew: -0.2347 Corr: 0.1904



Observations :

- The above graph depicts data distributions before and after performing log-transformation.
- We can see that performing log operations on the numerical distributions reduces skewness. However it doesn't completely become gaussian distributed correlations still remains the low with respect to the target variables.
- Since these markdown variables had lot of nan values and we imputed with mean values .Hence we see a very high spike at mean value.

```
In [58]: 1 ## As we dont need normalization we will drop them
          2 trte_set=trte_set[trte_set.columns.drop(list(trte_set.filter(regex='norm_')))]
```

3.6 Detecting outlier points in numerical Features

```
In [59]: 1 #calculating total fare amount values at each percentile 0,10,20,30,40,50,60,70,80,90,100
2 for i in range(0,100,10):
3     var = trte_set["CPI"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6 print("100 percentile value is ",var[-1])
```

```
0 percentile value is 126.064
10 percentile value is 129.1507742
20 percentile value is 131.686
30 percentile value is 134.17777420000002
40 percentile value is 138.3771935
50 percentile value is 182.60429219999997
60 percentile value is 197.780931
70 percentile value is 211.1886931
80 percentile value is 215.6488731
90 percentile value is 222.2174395
100 percentile value is 228.9764563
```

```
In [60]: 1 #calculating total fare amount values at each percentile 0,10,20,30,40,50,60,70,80,90,100
2 for i in range(0,100,10):
3     var = trte_set["Unemployment"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6 print("100 percentile value is ",var[-1])
```

```
0 percentile value is 3.683999999999997
10 percentile value is 5.539
20 percentile value is 6.29
30 percentile value is 6.875999999999979
40 percentile value is 7.287000000000001
50 percentile value is 7.742000000000001
60 percentile value is 8.058
70 percentile value is 8.36
80 percentile value is 8.743
90 percentile value is 9.652999999999999
100 percentile value is 14.312999999999999
```

```
In [61]: 1 #calculating total fare amount values at each percentile 0,10,20,30,40,50,60,70,80,90,100
2 for i in range(0,100,10):
3     var = trte_set["Temperature"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6 print("100 percentile value is ",var[-1])
```

```
0 percentile value is -7.29
10 percentile value is 33.0
20 percentile value is 41.72
30 percentile value is 48.71
40 percentile value is 54.62
50 percentile value is 60.32
60 percentile value is 65.79
70 percentile value is 70.71
80 percentile value is 76.58
90 percentile value is 83.04
100 percentile value is 101.95
```

```
In [62]: 1 #calculating total fare amount values at each percentile 0,10,20,30,40,50,60,70,80,90,100
2 for i in range(0,100,10):
3     var = trte_set["Size"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6 print("100 percentile value is ",var[-1])
```

```
0 percentile value is 34875
10 percentile value is 39910
20 percentile value is 57197
30 percentile value is 103681
40 percentile value is 120653
50 percentile value is 128107
60 percentile value is 155078
70 percentile value is 200898
80 percentile value is 203742
90 percentile value is 204184
100 percentile value is 219622
```

- We dont observe any outlier points in the dataset.

4. Feature Engineering

4.1 Adding Basic Date,Holiday Features

- Let's create some useful date features like month,dayofweek , MajorHolidays
- We observed that in our trainset we are given with IsHoliday as true for week 52(supposedly our christmas week) but there is no spike in sales , week 51 has spike in sales which is our week before chirstmas .
- lets' create a categorical feature Major_Holidays which calculates the holiday week of Christmas,Thanksgiving and Labour day for every year.This will allow us to spot the sales spike effectively
- we will categorize them as '0' for 'Non-Major Holiday', '1' for Thanksgiving and '2' for Labour and '3' for Christmas and 4 for superbowl


```

In [63]: 1 #Ref::: ## Create a list which stores holiday weeknumbers as it values
2 ## all the holidays occurs in same week from 2010-2013
3 def return_holweek(data):
4     d = []
5     ## thanks giving (occurs in the last thursday of november)
6     thurs_weeknum_ = data.query('month == 11')[['week','day']] ## extracting weel an day of month 11
7     thanks_week_num = thurs_weeknum_.groupby('week',as_index=False).count()['week'][3] ## extracting the last week from month 11
8
9     ## Labour (occurs in the first monday of september)
10    mon_weeknum = data.query('month == 9')
11    if mon_weeknum['Date'].dt.day.iloc[0] < 4:
12        labour_week_num = mon_weeknum.week.iloc[1] ## if the monday occurs in 2nd week of month
13    else :
14        labour_week_num = mon_weeknum.week.iloc[0] ## if the monday occurs in 1nd week of month
15
16    ## christmas (extract week number just before december 25)
17    christmas_week_num = data[(data['month'] == 12 )& (data['Date'].dt.day <25) ][['week']].iloc[-1]
18
19    ### superbowl occurs in the 6th week from 2010-2013
20    d = [thanks_week_num,labour_week_num,christmas_week_num,6]
21    return d

```

```

In [64]: 1 def basic_features(data):
2     '''Returns new dataframe with basic features : month,dayofyear,quater,major sales'''
3
4     # Creating month field
5     data['month'] = data['Date'].dt.month
6
7     # craeting day field
8     data['day'] = data['Date'].dt.day
9
10    # creating day of the year
11    data['day_ofyear'] = data['Date'].dt.dayofyear
12
13    ## creating quatrer
14    data['quater'] = data['Date'].dt.quarter
15
16    return data

```

```

In [65]: 1 trte_basic_feat = basic_features(trte_set.copy())

```

```

In [66]: 1 # Creating Labeled feature Major Holidays
2 # References : 0:Not Holiday ,1:Thanksgiving ,2:Labour day, 3:Christmas, 4: Superbowl
3 d = return_holweek(trte_basic_feat)
4 trte_basic_feat['major_holiday'] = trte_basic_feat['week'].map({d[0]:1,d[1]:2,d[2]:3,d[3]:4})
5 trte_basic_feat['major_holiday'].fillna(0,inplace=True)

```

In [67]: ▶

1trte_basic_feat.head(2)

Out[67]:

| | Store | Dept | Date | Weekly_Sales | IsHoliday | Temperature | Fuel_Price | MarkDown1 | MarkDown2 | MarkDown3 | MarkDown4 | MarkDown5 | CPI | Unemployment | Type | Size | set_type | week | year | mc |
|---|-------|------|------------|--------------|-----------|-------------|------------|-------------|-------------|-------------|-------------|-------------|------------|--------------|------|--------|----------|------|------|----|
| 0 | 1 | 1 | 2010-02-05 | 24924.50 | False | 42.31 | 2.572 | 8536.592778 | 3346.401918 | 1670.797978 | 3653.631444 | 4428.307667 | 211.096358 | 8.106 | A | 151315 | Train | 5 | 2010 | |
| 1 | 1 | 1 | 2010-02-12 | 46039.49 | True | 38.51 | 2.548 | 8536.592778 | 3346.401918 | 1670.797978 | 3653.631444 | 4428.307667 | 211.242170 | 8.106 | A | 151315 | Train | 6 | 2010 | |

5. Multi-Variate Analysis

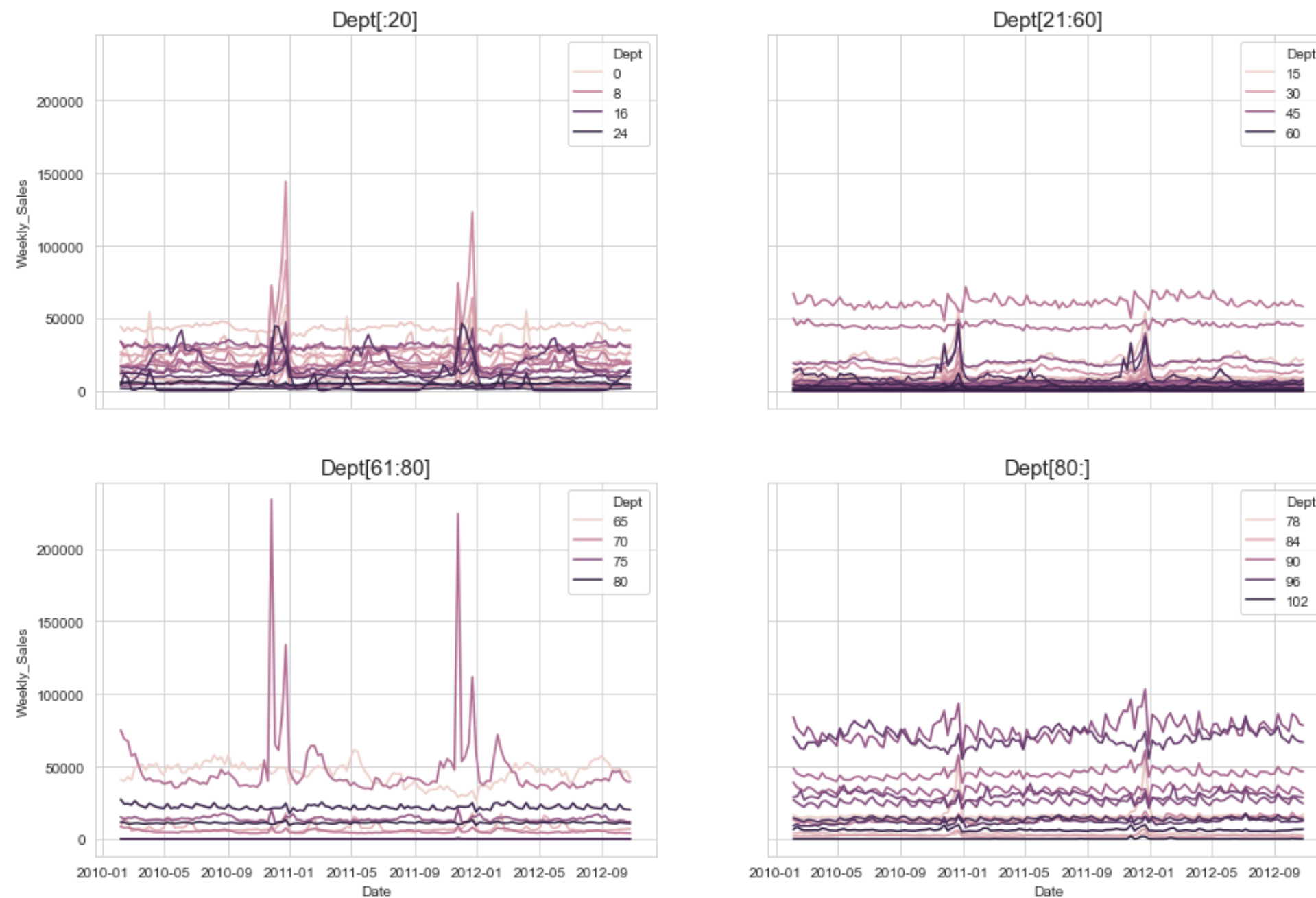
In [68]:

```

1  ## Date * Department - MEans sales lineplot
2  dept_mean = trte_basic_feat.groupby(['Dept', 'Date'],as_index=False).mean()
3  fig, axes = plt.subplots(2, 2, figsize=(15, 10),sharey=True,sharex=True)
4  fig.suptitle('Mean Department Sales',fontsize=17)
5  fig.subplots_adjust(wspace=0.2,hspace=0.2)
6  fig.subplots_adjust(left=0.1,right=0.9,bottom=0.1,top=0.9)
7
8  sns.lineplot(ax=axes[0,0],data=dept_mean.query('Dept <= 20'),x='Date',y='Weekly_Sales',hue='Dept')
9  sns.lineplot(ax=axes[0,1],data=dept_mean.query('Dept > 20 and Dept <=60 '),x='Date',y='Weekly_Sales',hue='Dept')
10 sns.lineplot(ax=axes[1,0],data=dept_mean.query('Dept > 60 and Dept <=80'),x='Date',y='Weekly_Sales',hue='Dept')
11 sns.lineplot(ax=axes[1,1],data=dept_mean.query('Dept > 80'),x='Date',y='Weekly_Sales',hue='Dept')
12
13 axes[0,0].set_title('Dept[:20]',fontsize=15)
14 axes[0,1].set_title('Dept[21:60]',fontsize=15)
15 axes[1,0].set_title('Dept[61:80]',fontsize=15)
16 axes[1,1].set_title('Dept[80:]',fontsize=15)
17 plt.show()

```

Mean Department Sales



Observations :

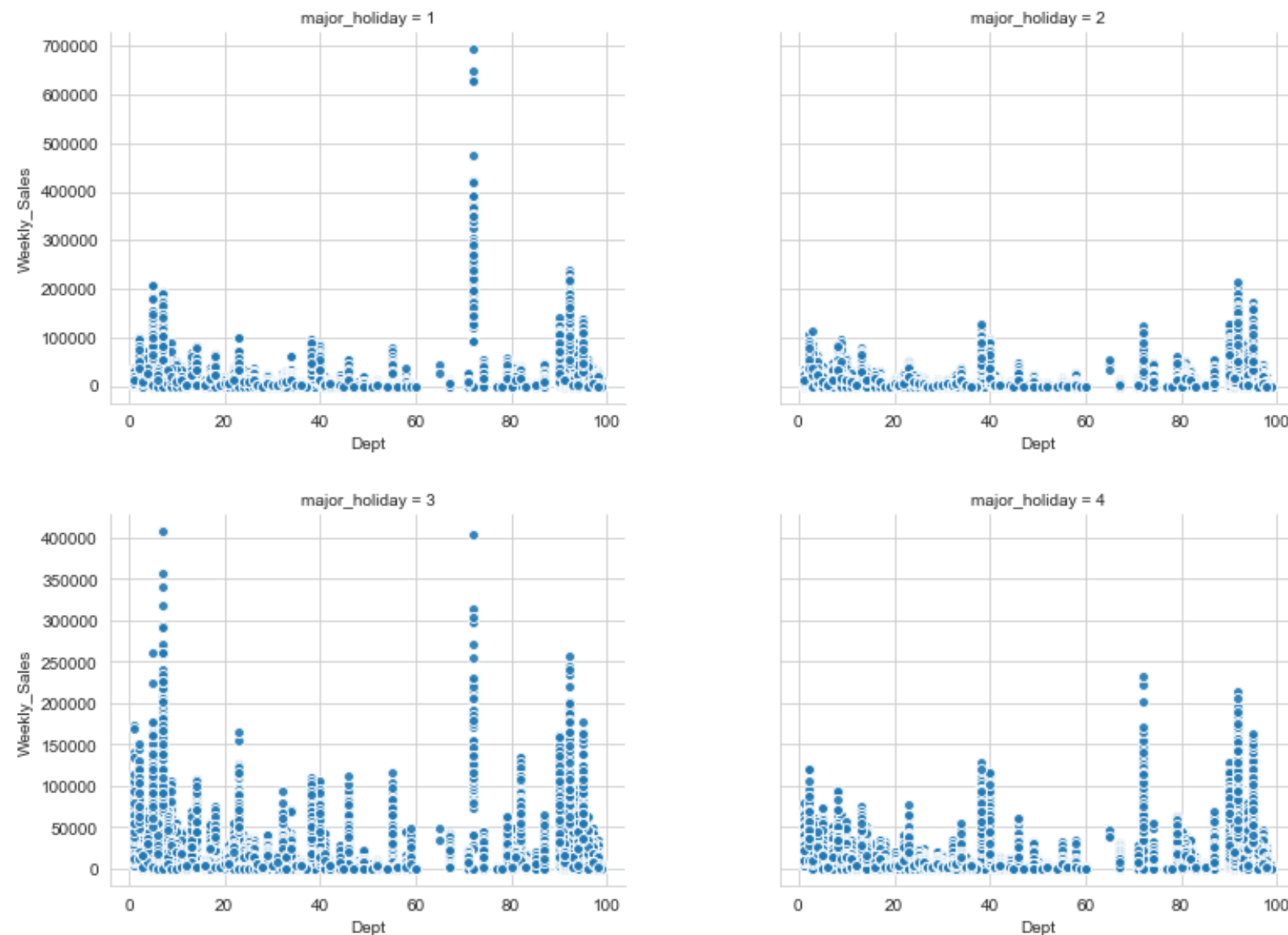
- From the above plot we can see department seems to be a good predictor for sales.
- From plot 3 it can be observed that few departments contribute to very high mean weekly sales
- **Theory 1** : These departments might contribute majorly for sales during holidays /Major Holidays

```

In [69]: 1 # Let's plot Effect of department on holidays
2 g = sns.FacetGrid(trte_basic_feat, col="major_holiday",
3                 col_order=[1,2],
4                 aspect=1.2, size=4.5, palette=['Yellow', 'Blue', "Green", 'Pink'])
5 g.map(plt.scatter, "Dept", "Weekly_Sales", alpha=0.9,
6       edgecolor='white')
7 fig = g.fig
8 fig.subplots_adjust(top=0.8, wspace=0.3)
9 fig.suptitle('Effect of Department on Holidays (1:Thanksgiving ,2:Labour day, 3:Christmas, 4: Superbowl)', fontsize=14)
10
11 g = sns.FacetGrid(trte_basic_feat, col="major_holiday",
12                 col_order=[3,4],
13                 aspect=1.2, size=4.5, palette=['Yellow', 'Blue', "Green", 'Pink'])
14 g.map(plt.scatter, "Dept", "Weekly_Sales", alpha=0.9,
15       edgecolor='white')
16 fig = g.fig
17 fig.subplots_adjust(top=0.8, wspace=0.3)
18

```

Effect of Department on Holidays (1:Thanksgiving ,2:Labour day, 3:Christmas, 4: Superbowl)



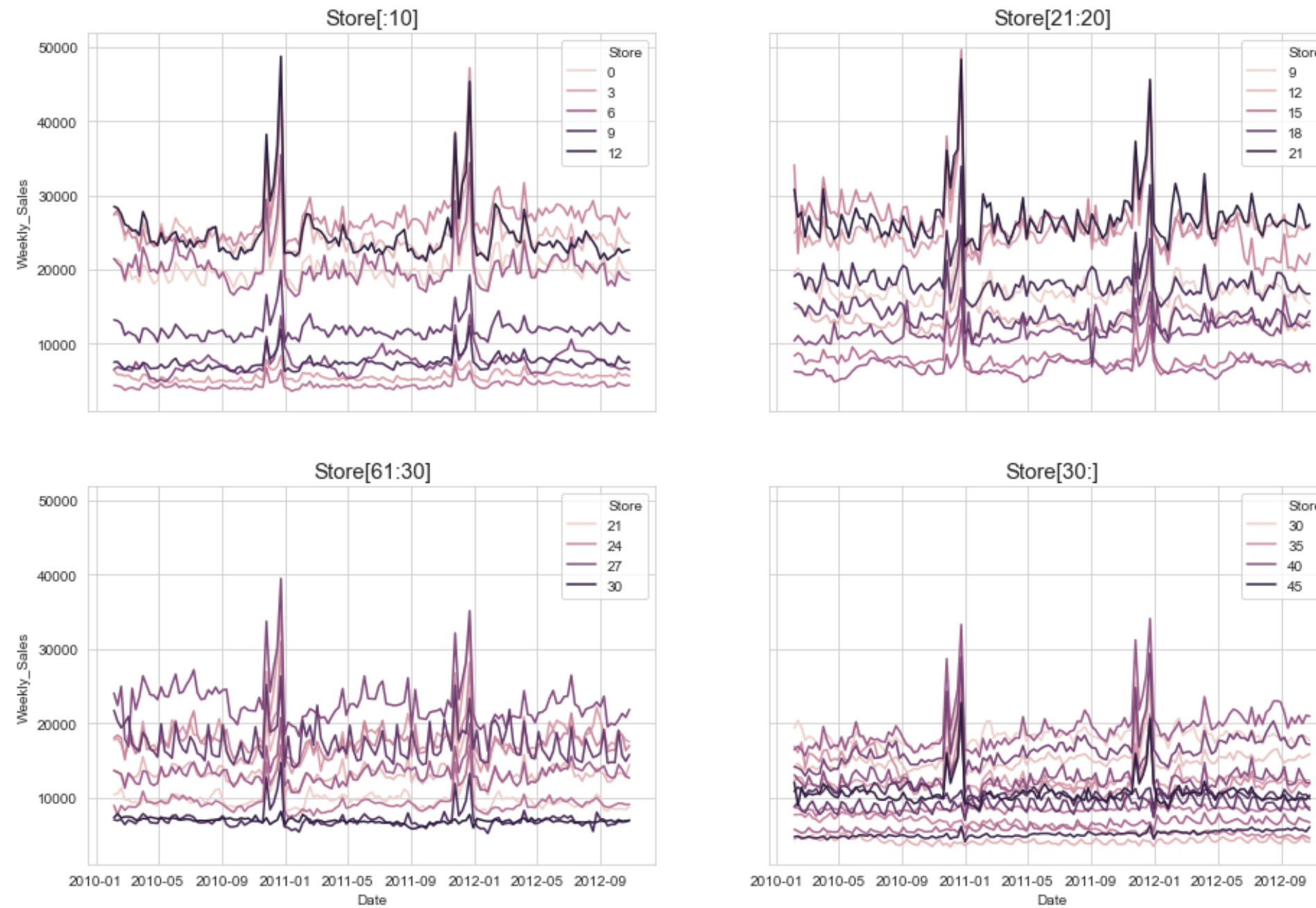
Observations :

- We can say that major_hols is a good predictor of sales based on major holidays
- There is a clear cut effect of certain department on Holidays.
- major_hols=3 indicates Christmas .We can see peaks in sales in frist 10-20 departments which was not observed in other holidays
- Sales of few departments also seems to have reduced during the Holiday week than the non-holiday weeks.
- We can See that although Type 'B' stores are smaller tha type 'A' Store during certain holidays Type-B store produced high sales in than A.

In [70]:

```
1  ## Date * Department - MEans sales lineplot
2  Store_mean = trte_basic_feat.groupby(['Store', 'Date'], as_index=False).mean()
3  fig, axes = plt.subplots(2, 2, figsize=(15, 10), sharey=True, sharex=True)
4  fig.suptitle('Mean Store Sales', fontsize=17)
5  fig.subplots_adjust(wspace=0.2, hspace=0.2)
6  fig.subplots_adjust(left=0.1, right=0.9, bottom=0.1, top=0.9)
7
8  ## we will divide our plots into 4 equal parts and plot Linegraph
9  sns.lineplot(ax=axes[0,0], data=Store_mean.query('Store <= 10'), x='Date', y='Weekly_Sales', hue='Store')
10 sns.lineplot(ax=axes[0,1], data=Store_mean.query('Store > 10 and Store <=20 '), x='Date', y='Weekly_Sales', hue='Store')
11 sns.lineplot(ax=axes[1,0], data=Store_mean.query('Store > 20 and Store <=30'), x='Date', y='Weekly_Sales', hue='Store')
12 sns.lineplot(ax=axes[1,1], data=Store_mean.query('Store > 30'), x='Date', y='Weekly_Sales', hue='Store')
13
14 axes[0,0].set_title('Store[:10]', fontsize=15)
15 axes[0,1].set_title('Store[21:20]', fontsize=15)
16 axes[1,0].set_title('Store[61:30]', fontsize=15)
17 axes[1,1].set_title('Store[30:]', fontsize=15)
18 plt.show()
```


Mean Store Sales



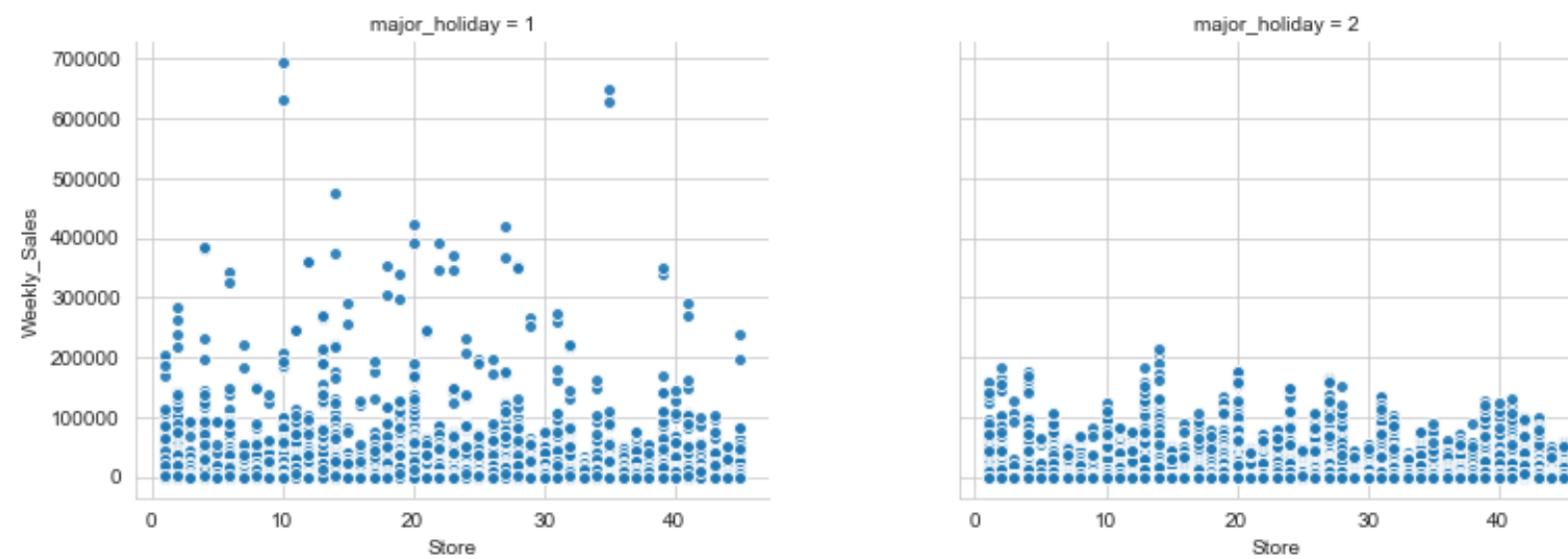
- Similar to deparment we can see few stores also has higher mean sales

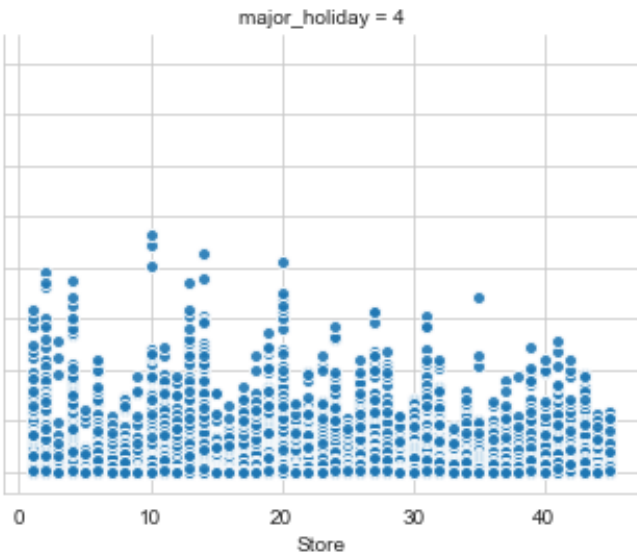
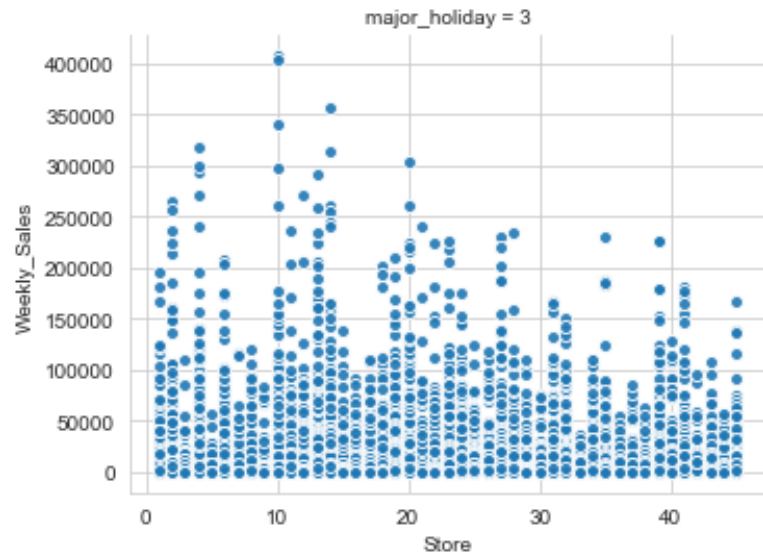
```

In [71]: 1  ## effect of holidays on Stores
2  g = sns.FacetGrid(trte_basic_feat, col="major_holiday",
3                  col_order=[1,2],
4                  aspect=1.2, size=4.5, palette=['Yellow', 'Blue', "Green", 'Pink'])
5  g.map(plt.scatter, "Store", "Weekly_Sales", alpha=0.9,
6        edgecolor='white')
7  fig = g.fig
8  fig.subplots_adjust(top=0.8, wspace=0.3)
9  fig.suptitle('Effect of Holidays on Store', fontsize=14)
10
11 g = sns.FacetGrid(trte_basic_feat, col="major_holiday",
12                  col_order=[3,4],
13                  aspect=1.2, size=4.5, palette=['Yellow', 'Blue', "Green", 'Pink'])
14 g.map(plt.scatter, "Store", "Weekly_Sales", alpha=0.9,
15       edgecolor='white')
16 fig = g.fig
17 fig.subplots_adjust(top=0.8, wspace=0.3)
18
19

```

Effect of Holidays on Store

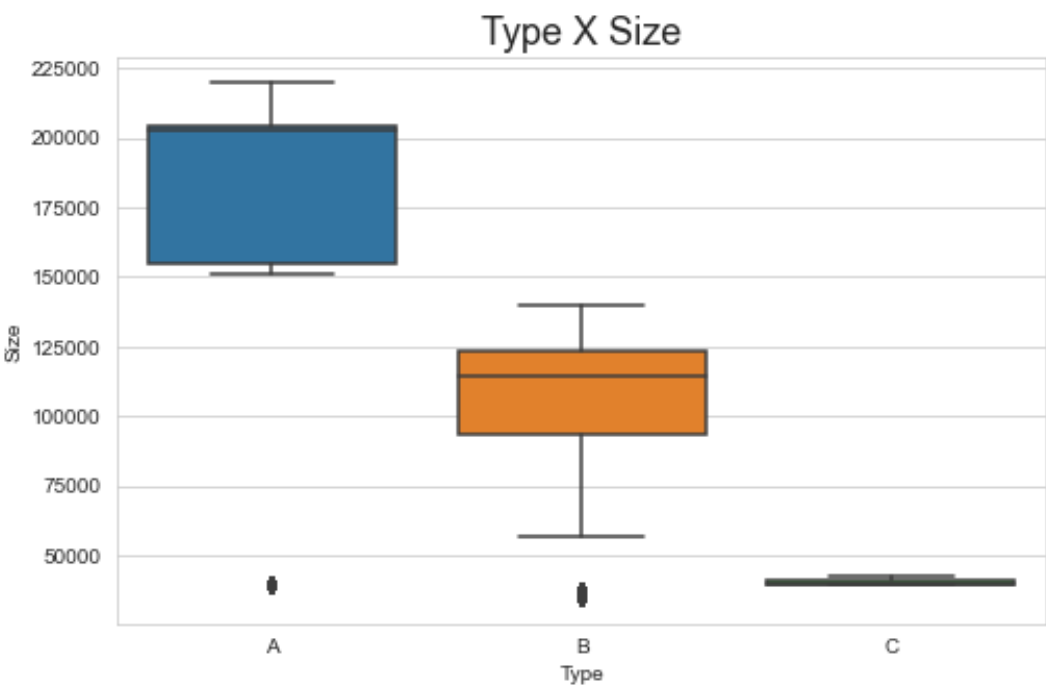




Observations :

- Similar to Department we can see that stores also have similar effect on sales during Holidays
- We can say that Christmas and Thanksgiving are the major holidays which uplifts the sales scale.

```
In [72]: 1 ## Size*Type -scatter bubble
2 fig = plt.figure(figsize=(8,5))
3 sns.boxplot(data=trte_basic_feat,x='Type',y='Size')
4 plt.title('Type X Size',fontsize=18)
5 plt.show()
```



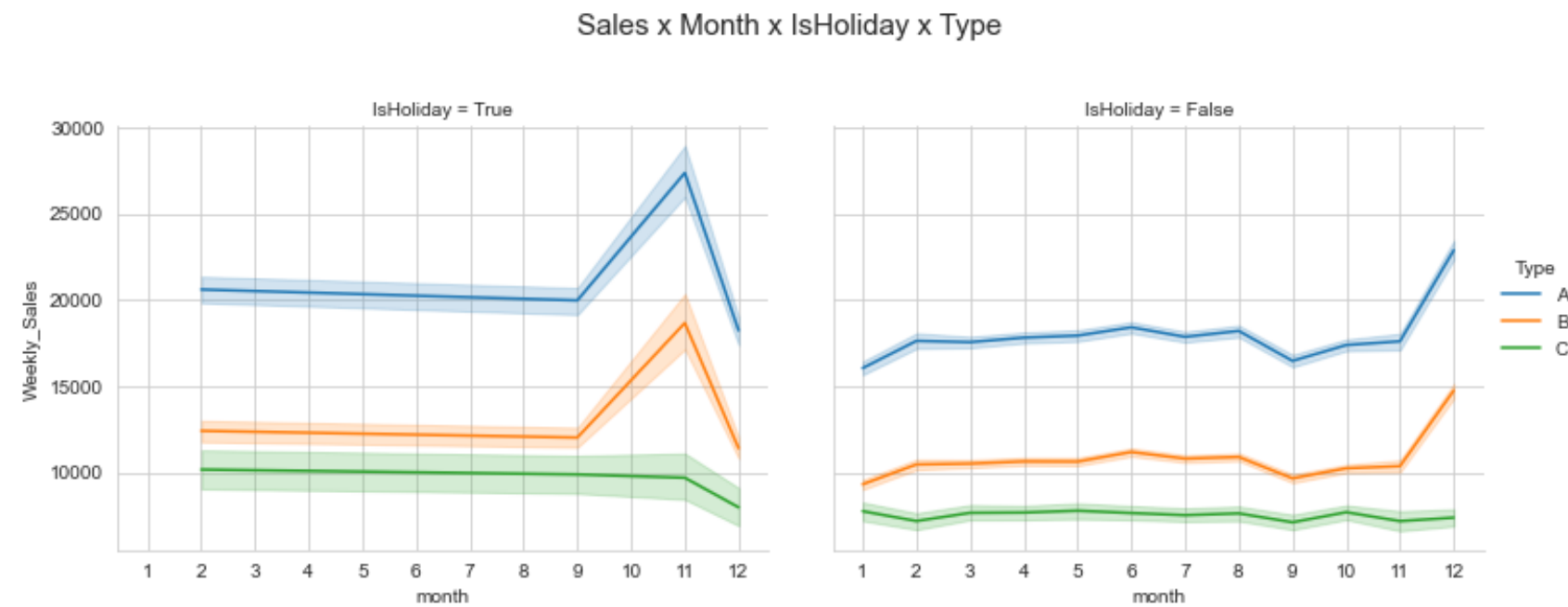
Observations :

- Here we can see how A,B and C stores are clearly differentiated by their size
- We already explored how Type A > B > C.Size here confirms bigger the store bigger the size.

```

In [73]: 1 ## Month*Holiday*Type - Box plot
2 g = sns.FacetGrid(trte_basic_feat,col='IsHoliday',col_order=[True,False],hue='Type',hue_order=['A','B','C'],
3                 aspect=1.2,size=4.5)
4 g.map(sns.lineplot,'month','Weekly_Sales')
5 g.set(xticks=np.arange(1,13))
6 fig = g.fig
7 fig.suptitle('Sales x Month x IsHoliday x Type',size=15)
8 fig.subplots_adjust(top=0.8,wspace=0.1)
9 l = g.add_legend(title='Type')

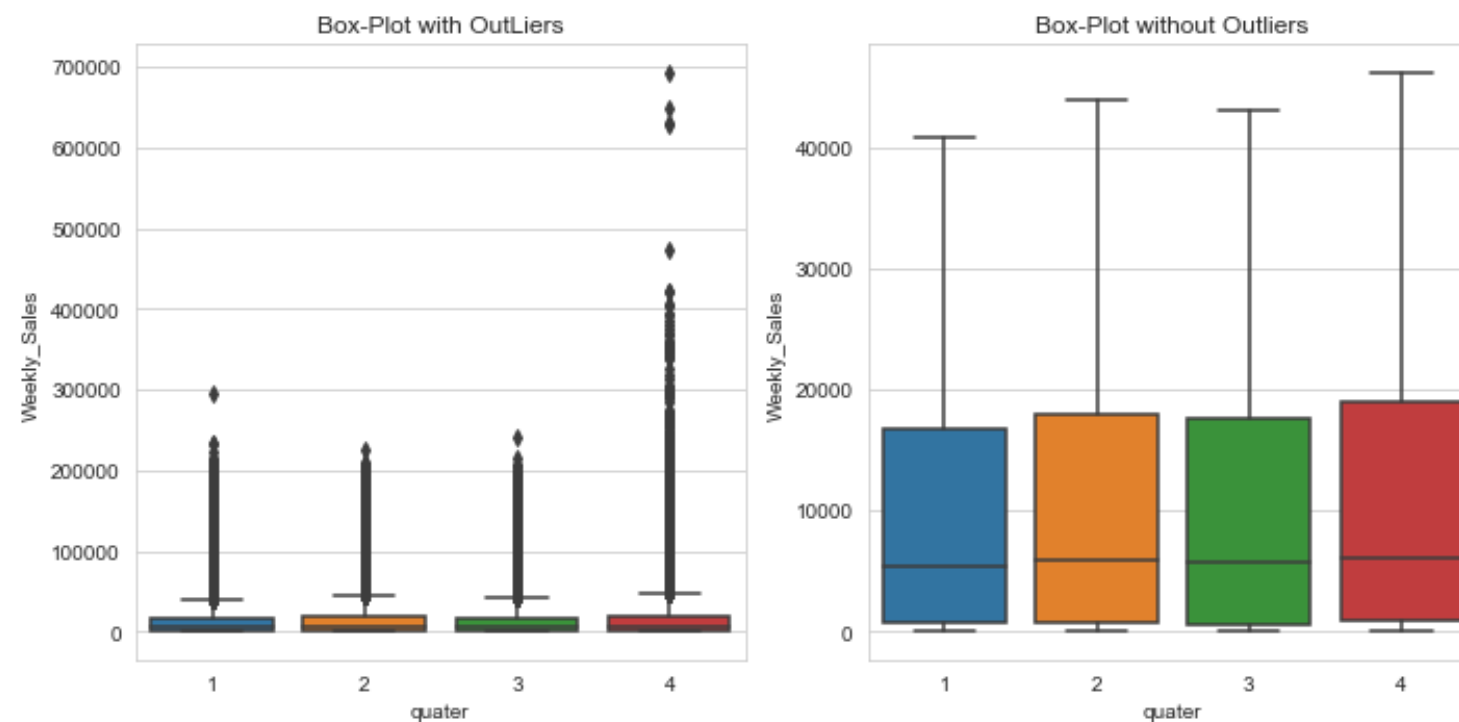
```



Observations:

- Month&Holiday is also a good predictor of sales as we can see during months 9-12 during holidays there is sales peak in stores A and B .

```
In [74]: 1 ### Let's check how quaters affect the sales
2 fig = plt.figure(figsize=(10,5))
3 plt.subplot(1,2,1)
4 sns.boxplot(data=trte_basic_feat,x='quater',y='Weekly_Sales',)
5 plt.title("Box-Plot with OutLiers")
6 plt.subplot(1,2,2)
7 sns.boxplot(data=trte_basic_feat,x='quater',y='Weekly_Sales',showfliers=False)
8 plt.title('Box-Plot without Outliers')
9 plt.tight_layout()
```



Observations:

- We see that mean sales are almost same across all the quaters.
- Quater 4 has has flier points depecting sudden peak in sales in some week.

5.1 Correlation

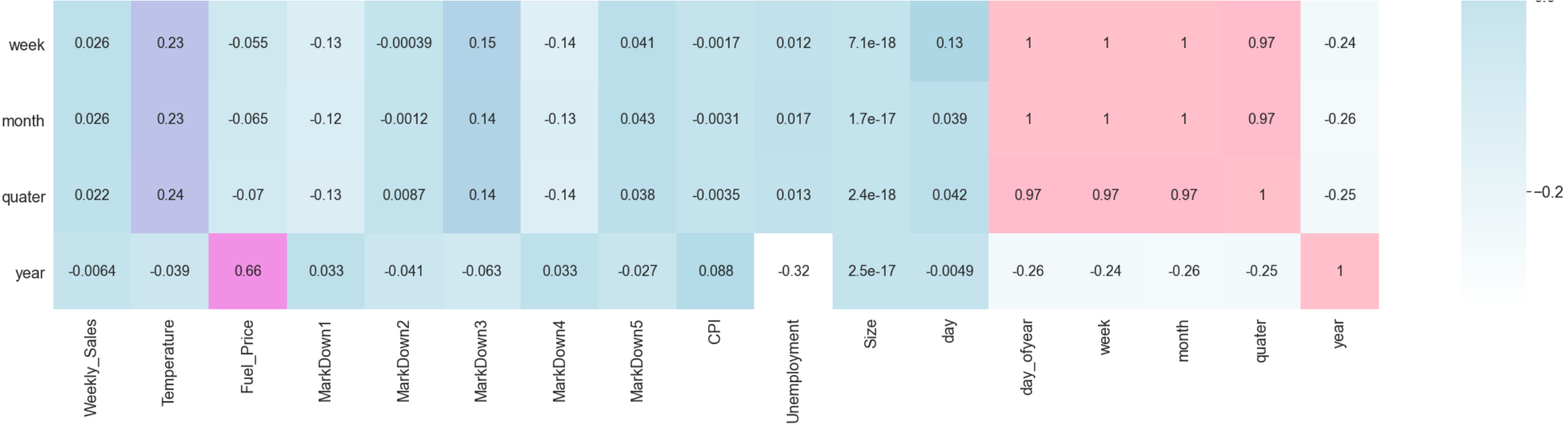
```
In [75]: 1 numerical_featues = ['Weekly_Sales', 'Temperature', 'Fuel_Price', 'MarkDown1', 'MarkDown2', 'MarkDown3', 'MarkDown4',
2                          'MarkDown5', 'CPI', 'Unemployment', 'Size', 'day', 'day_ofyear', 'week', 'month', 'quater', 'year']
3 categorical_features = ['Store', 'Dept', 'IsHoliday', 'Type', 'major_holiday']
```

Correlation between numerical features

In [76]: ▶

```
1 fig = plt.figure(figsize=(25,20))
2 sns.set(font_scale=1.8)
3 cmap = sns.blend_palette(colors=['White','lightblue','violet','pink'],n_colors=4,as_cmap=True)
4 ssn = sns.heatmap(trte_basic_feat[numerical_feats].corr(),annot=True,annot_kws={'fontsize':18},cmap=cmap)
5 plt.title("Correlation Map",fontdict={'size':30},pad=20)
6 fig = plt.gcf() # or by other means, like plt.subplots
7 figsize = fig.get_size_inches()
8 fig.set_size_inches(figsize * 1.5) # scale current size by 1.5
9 plt.show()
```





Observations:

- Fuel price and year are highly correlated,hence we can decide to drop Fuel_price.
- Markdown1 and Markdown4 are corelated ,lets drop markdown 4
- Size has the strongest corelation(0.2) with Weekly_Sales.
- Fuel_price and Temperature has very weak coorelation with sales
- Week is strongly co-related with month,dayofyear and quater we will drop them
- let's drop column dayofyear and month as it's strongly corelated to week.

```
In [77]: 1 ## DRopping Fuel_price and markdown4
        2 trte_basic_feat.drop(columns=['Fuel_Price', 'Markdown4', 'day_ofyear', 'quater'],inplace=True)
        3 print('Columns : Fuel_Price,Markdown4,day_ofyear,quater has been dropped')
```

Columns : Fuel_Price,Markdown4,day_ofyear,quater has been dropped

Correlation between numerical features and categorical features

Anova -Test : is a statistical test which can be used to determine the correlation between categorical and continuous variable.Anova test performs analysis of variance by calculating population means of each category.Annova test performed when there are atleast three categories.Here the null hypothesis(H0) states that the variables are not correlated. we can say if p_value < alpha(0.05) value we reject our null hypotheseis.If p_value > alpha value and F value > f critical we accept our null hypothesis

```

In [78]: 1 def anova_table(aov):
2         aov['mean_sq'] = aov[:]['sum_sq']/aov[:]['df']  ## calculating mean_sq
3         aov['eta_sq'] = aov[:,-1]['sum_sq']/sum(aov['sum_sq']) ### will tell the variance explained sseffect /sstotal
4         cols = ['sum_sq', 'df', 'mean_sq', 'F', 'PR(>F)', 'eta_sq']
5         aov = aov[cols]
6         return aov
7     def ols_(column):
8         model = ols('Sales ~ C({})'.format(column), data=data).fit()
9         aov_table = sm.stats.anova_lm(model, typ=2)
10        return aov_table
11
12    for i in categorical_features:
13        x = trte_basic_feat.loc[(trte_basic_feat.set_type=='Train')][i]
14        y = trte_basic_feat.loc[(trte_basic_feat.set_type=='Train')]['Weekly_Sales']
15        data = pd.DataFrame({i:x, 'Sales':y})
16        aov_table = ols_(i)
17        aov = anova_table(aov_table)
18        print('Stats Summary for feature ',i,'**20)
19        print(aov)
20        print('Variance explained (eta-squared) by the feature ',i,' :',aov_table['eta_sq'][-2])
21        print('\n\n')

```

```

Stats Summary for feature Store *****
      sum_sq      df      mean_sq      F  PR(>F)      eta_sq
C(Store)  2.100315e+13   44.0  4.773443e+11 1087.347135    0.0  0.091481
Residual  2.085878e+14 475144.0  4.389990e+08      NaN      NaN      NaN
Variance explained (eta-squared) by the feature Store : 0.09148076121863773

```

```

Stats Summary for feature Dept *****
      sum_sq      df      mean_sq      F  PR(>F)      eta_sq
C(Dept)   1.255263e+14   80.0  1.569078e+12 7163.639917    0.0  0.546739
Residual  1.040646e+14 475108.0  2.190337e+08      NaN      NaN      NaN
Variance explained (eta-squared) by the feature Dept : 0.546738823094827

```

```

Stats Summary for feature IsHoliday *****
      sum_sq      df      mean_sq      F      PR(>F)  \
C(IsHoliday)  2.583926e+11     1.0  2.583926e+11 535.400812  2.196219e-118
Residual      2.293325e+14 475187.0  4.826153e+08      NaN      NaN

      eta_sq
C(IsHoliday)  0.001125
Residual      NaN
Variance explained (eta-squared) by the feature IsHoliday : 0.001125447972901624

```

```

Stats Summary for feature Type *****
      sum_sq      df      mean_sq      F  PR(>F)      eta_sq
C(Type)   7.990171e+12     2.0  3.995085e+12 8566.79769    0.0  0.034802
Residual  2.216007e+14 475186.0  4.663453e+08      NaN      NaN      NaN
Variance explained (eta-squared) by the feature Type : 0.03480177417489362

```

```
Stats Summary for feature major_holiday *****
              sum_sq      df      mean_sq      F      PR(>F)  \
C(major_holiday)  8.571957e+11    4.0  2.142989e+11  445.1964    0.0
Residual          2.287337e+14  475184.0  4.813582e+08    NaN    NaN

              eta_sq
C(major_holiday)  0.003734
Residual          NaN
Variance explained (eta-squared) by the feature major_holiday : 0.0037335787735222523
```

P.S degree of correlation using (eta-squared)

- 0.001-0.009 - weak correlation
- 0.01 - 0.2 - moderate correlation
- 0.2 - 0.5 - good correlation

Observations :

- Using the anova test we checked for correlation between categorical(store,dept,isholiday etc) and numerical variable (sales).
- We see that the department feature is well correlated to weekly sales with 0.54 correlation.
- Store feature is moderately corelated (0.09) with weekly Sales
- We see that type,major holiday and isholiday possesses weak correlation

6. Time series Analysis

Why do we perform Time-Series Analysis ?

Using time-series we can get the insight of how time can add weightage in determining the value of an asset/commodity or variable of our concern. In our usual Machine Learning analysis we check how the independent variables affecting our dependent variable,order doesn't matter. In our time series analysis we check how the past observations of our target variables varies with respect to time .Here the data is ordered with respect to time which creates a dependency between past and future observations.

Time-Series Components !!!

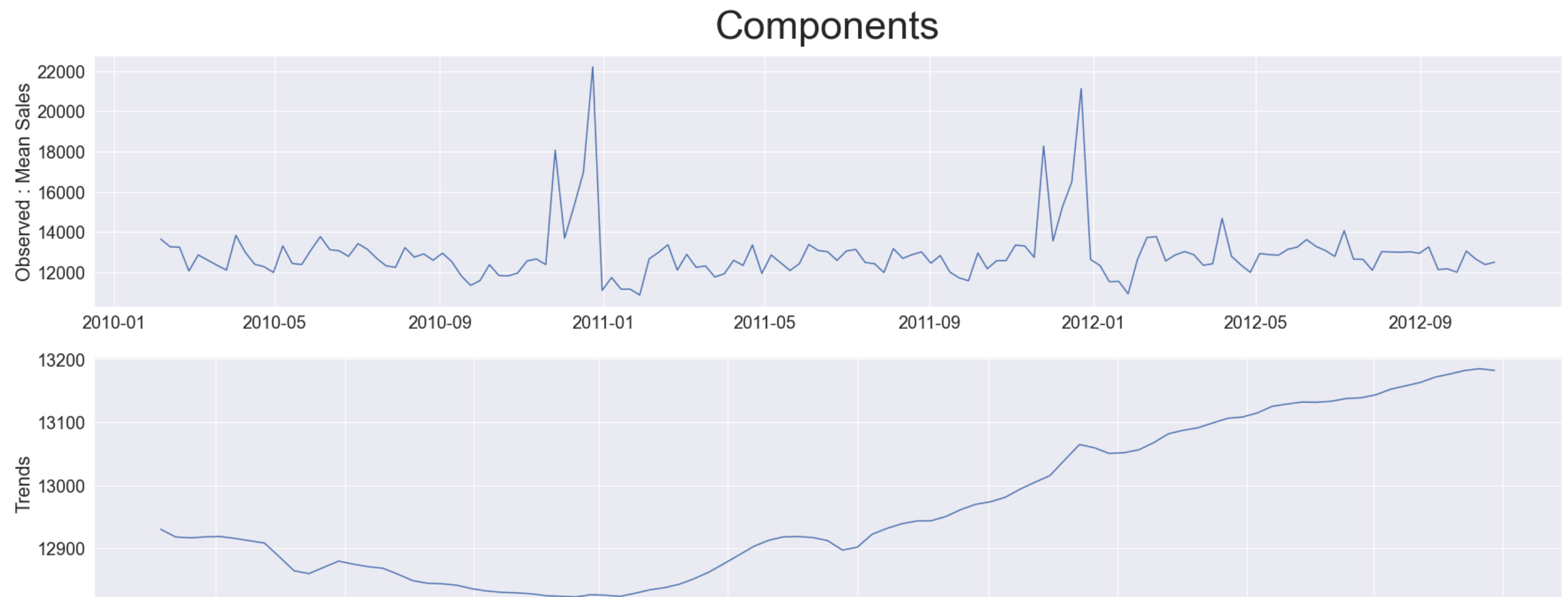
There are 4 major components of Time-series

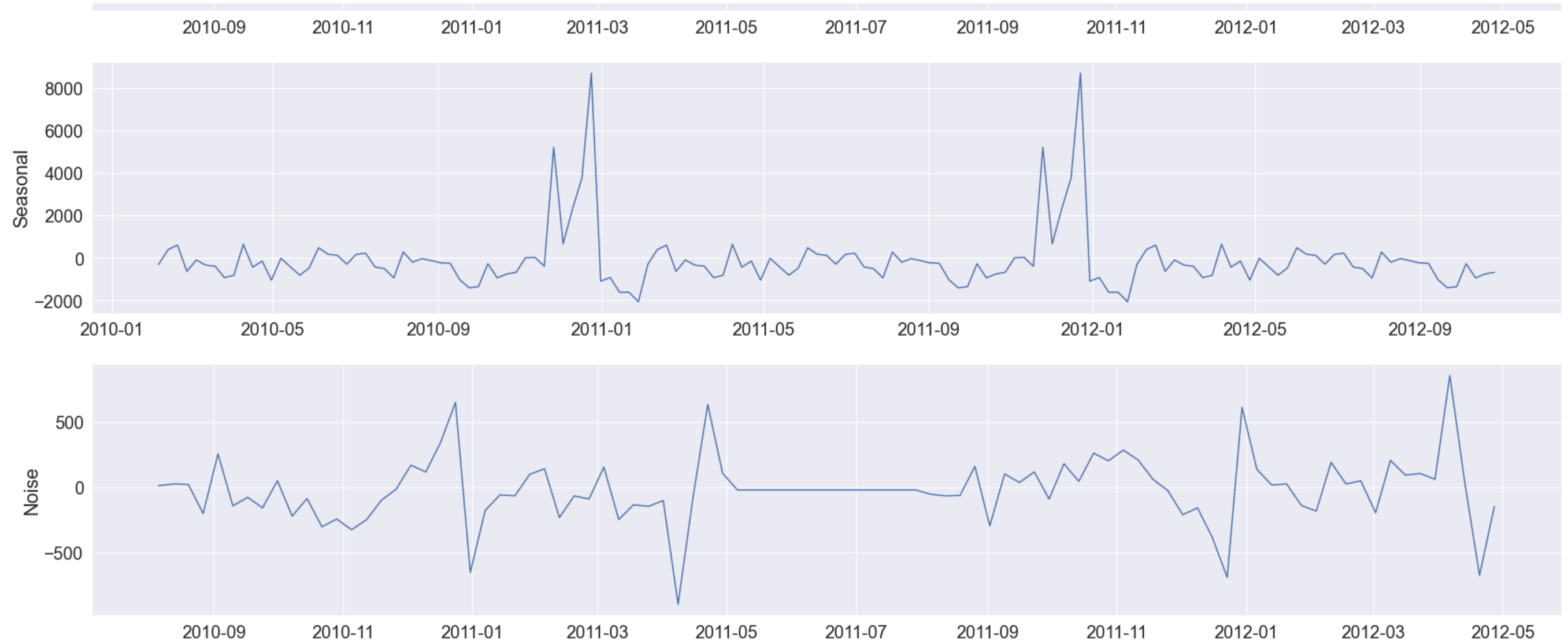
- Level: Average value of the series*
- Trend: Can be interpreted as a slope.Rate of change within two adjacent values.This is a optional component*
- Seasonality: Shows if there is any cyclic behaviour in the data.(i.e)Repeatedness after a time period of time*
- Noise: Uncertain variation which cannot be interpreted the model*

```
In [79]: 1  ## In time -series analysis we are just going to consider the
2  ## date feature and the weekly sales vale
3
4  ## Data Loading
5  ## Date has to be parsed to perform date related analysis
6  ## squeeze will convert data frame to series
7  time_series = pd.read_csv('Clean_train.csv',
8                             header=None,
9                             usecols=[2,3],skiprows=1,
10                             names=['Date','Weekly_Sales'],
11                             parse_dates=["Date"],
12                             index_col=0)
```

```
In [80]: 1  time_series_grouped = time_series.groupby('Date').mean()
```

```
In [81]: 1 ### Check for trend and seasonality noise in data
2 from statsmodels.tsa.seasonal import seasonal_decompose
3 decompose = seasonal_decompose(time_series_grouped)
4
5 trend = decompose.trend
6 seasonal = decompose.seasonal
7 noise = decompose.resid
8
9 fig = plt.figure(figsize=(25,20))
10 #plt.subplots_adjust(top=0.2, wspace=0.12)
11
12 plt.subplot(411)
13 plt.plot(time_series_grouped)
14 plt.title('Components', fontdict={'fontsize':45}, pad=20)
15 plt.ylabel('Observed : Mean Sales')
16
17 plt.subplot(412)
18 plt.plot(trend)
19 plt.ylabel('Trends')
20
21 plt.subplot(413)
22 plt.plot(seasonal)
23 plt.ylabel('Seasonal')
24
25 plt.subplot(414)
26 plt.plot(noise)
27 plt.ylabel('Noise')
28 plt.tight_layout()
```





- We can see clearly that data contains both upwards and downward trend.
- Additive Seasonality : We can see that there is 90% spike in sales during year end every year additively.
- Data its is also quite noisy, which depicts uncertainty in sales in some periods .Here is when we use smoothing technique.We are going to average out some past values which will eventually reduce noise in data.

6.1. Time-Series Smoothing Techniques

We are going to implement some Averaging techniques for each of our store and department columns as features.

Lag features : We are using shift() function from pandas library.Here we consider lags based on trend,the trend is yearly so lets take 52 lags

Rolling window : We are using rolling() function from pandas library.Here weConsider a constant window size(3) and take average of past values within that window.

Expanding window :We are using expand() function from pandas library.Here we Consider a expanding window size(2) and take average of past values within that window.

Exponential weighted moving average :We are using ewm() function from pandas library. Here recent values are given more weight and the weights reduce exponentially as we move further past values

Holt's Winters triple exponential smoothing: Consider smoothing by taking into account of level,trend and seasons

- Ref : <https://towardsdatascience.com/moving-averages-in-python-16170e20f6c> (<https://towardsdatascience.com/moving-averages-in-python-16170e20f6c>)
<https://machinelearningmastery.com/basic-feature-engineering-time-series-data-python/> (<https://machinelearningmastery.com/basic-feature-engineering-time-series-data-python/>)

```

In [82]: 1 def creating_smoothing_columns(df,cols,yhat,alphas=[0.3],lags=[52],esize=1,wsiz=3):
2         '''creates lag columns,rolling mean column,exponential weighted moving average column to the dataset'''
3
4         grouped = df.groupby(cols)
5
6         ## calculate rolling mean of specified window size
7         # we are shifting by 52 because there is a seasonal trend and we will be able to forecast better.
8         ## As we will be able to take mean of nearest points since we have 40 weeks to predict
9         df['rolling_mean'] = np.round(grouped[yhat].apply(lambda x: x.shift(52).rolling(window=wsiz,min_periods=1).mean()),3)
10
11        ## calculate expanding window of specified size
12        df['expanding_mean'] = np.round(grouped[yhat].apply(lambda x: x.shift(52).expanding(esize).mean()),3)
13
14        ## calculate ewm with different alpha combinations
15        for a in alphas:
16            df['EWM_{}'.format(a)] = np.round(grouped[yhat].apply(lambda x: x.shift(52).ewm(alpha=a,adjust=False).mean()),3)
17
18        ## create lag columns with specified lags
19        for l in lags :
20            df['lag_{}'.format(l)] = np.round(grouped[yhat].apply(lambda x: x.shift(l)),3)
21
22        return df
23
24
25    ##function call to create smoothed features
26
27    trte_featured = creating_smoothing_columns(trte_basic_feat.copy(),['Store','Dept'],'Weekly_Sales',[0.1,0.4,0.5,0.7],[52],1,3)
28
29

```

HOLT WINTER's - TRIPLE EXPONENT SMOOTHENING:

As our data contains additive seasonality.Holt winter's triple exponent smoothening is apt to capture forecast.The values are smoothened at 3 levels namely level,trend and seasonality.

$$\begin{aligned}
 \hat{y}_{t+h|t} &= \ell_t + hb_t + s_{t+h-m(k+1)} \\
 \ell_t &= \alpha(y_t - s_{t-m}) + (1 - \alpha)(\ell_{t-1} + b_{t-1}) \\
 b_t &= \beta^*(\ell_t - \ell_{t-1}) + (1 - \beta^*)b_{t-1} \\
 s_t &= \gamma(y_t - \ell_{t-1} - b_{t-1}) + (1 - \gamma)s_{t-m},
 \end{aligned}$$

The forecast $(\hat{y}_{t+h|t})$ (h is the number of predictions to be made from t) is calculated with the combination of level (ℓ_t) , trend (b_t) , and seasonal component (s_t) . The level, trend and seasonal factors are smoothened with $(\alpha, \beta^*, \gamma)$. In level we are taking the exponential weighted of non-seasonal forecast $(\ell_{t-1} + b_{t-1})$ and seasonal adjustment (how much difference from past season) $(y_t - s_{t-m})$. Seasonal frequency m can be quarterly (m=4), yearly (m=12) based on events. Trend can be intuitively understood as the slope. Trend tells the rate of change between adjacent observations t, (t-1). In trend we are taking exponential weighted average between past trend (b_{t-1}) and level adjustment $(\ell_t - \ell_{t-1})$ (how much difference from past level). Similarly in seasonal component we are considering two factors level and trend. Equation $(y_t - \ell_{t-1} - b_{t-1})$ is the trend and level adjustment from the current observation.

References - <https://grisha.org/blog/2016/02/17/triple-exponential-smoothing-forecasting-part-iii/> (<https://grisha.org/blog/2016/02/17/triple-exponential-smoothing-forecasting-part-iii/>).

```

In [83]: 1 def initial_trend(series, slen):
2         sum = 0.0
3         for i in range(slen):
4             sum += float(series[i+slen] - series[i]) / slen
5         return sum / slen

```

```

In [84]: 1 def initial_seasonal_components(series, slen):
2         seasonals = {}
3         season_averages = []
4         n_seasons = int(len(series)/slen)
5         # compute season averages
6         for j in range(n_seasons):
7             season_averages.append(sum(series[slen*j:slen*j+slen])/float(slen))
8         # compute initial values
9         for i in range(slen):
10            sum_of_vals_over_avg = 0.0
11            for j in range(n_seasons):
12                sum_of_vals_over_avg += series[slen*j+i]-season_averages[j]
13            seasonals[i] = sum_of_vals_over_avg/n_seasons
14        return seasonals

In [85]: 1 def triple_exponential_smoothing(series,slen, alpha, beta, gamma, n_preds):
2         result = []
3         # n_preds is the number of predictions to be made
4         abs_err = []
5         seasonals = initial_seasonal_components(series, slen)
6         for i in range(len(series)+n_preds):
7             if i == 0: # initial values
8                 smooth = series[0]
9                 abs_err.append(0)
10                trend = initial_trend(series, slen)
11                result.append(series[0])
12                continue
13            if i >= len(series): # we are forecasting
14                m = i - len(series) + 1
15                result.append((smooth + m*trend) + seasonals[i%slen])
16            else:
17                val = series[i]
18                last_smooth, smooth = smooth, alpha*(val-seasonals[i%slen]) + (1-alpha)*(smooth+trend)
19                trend = beta * (smooth-last_smooth) + (1-beta)*trend
20                seasonals[i%slen] = gamma*(val-smooth) + (1-gamma)*seasonals[i%slen]
21                final_val = smooth+trend+seasonals[i%slen]
22                abs_err.append(abs(series[i] - final_val))
23                result.append(final_val)
24
25        return result,abs_err

```

Hyperparameter Tuning Alpha,Beta and Gamma values


```

In [86]: 1 season_len = 52
2 mape = []
3 ## Lets manually fine tune alpha beta and gamma
4 alpha = [0.25,0.2,0.15,0.1]
5 beta = [0.1,0.15,0.20,0.25]
6 gamma = [0.1,0.15,0.20,0.08,0.05]
7 fine_tune = []
8 length = len(trainset_merged)
9 start = datetime.now()
10 if not os.path.isfile('holts.pkl'):
11     for a in alpha :
12         for b in beta:
13             for g in gamma:
14                 predict_list_2,act_sum,predict_sum=[],[],[]
15                 for s in ustore:
16                     for d in udept:
17                         ## for each store and dept pass the target column values
18                         data=trainset_merged.loc[(trainset_merged.Store ==s )&(trainset_merged.Dept==d)]['Weekly_Sales'].reset_index(drop=True)
19
20                         # perform smoothening (return predict values,absolute error)
21                         predict_values_2,abs_val = triple_exponential_smoothing(data, season_len, a, b, g, 0)
22
23                         ## Add to final predict list
24                         predict_list_2.extend(predict_values_2)
25                         act_sum.append(sum(data))
26                         predict_sum.append(sum(abs_val))
27
28                         ## calculate MAPE per combination
29                         mape_err = (sum(predict_sum)/length)/(sum(act_sum)/length)
30                         print('Mape% with alpha : {} ,beta:{}, gamma:{} is '.format(a,b,g),mape_err*100)
31                         fine_tune.append((a,b,g,mape_err*100))
32     pkl.dump(fine_tune,open('holts.pkl','wb'))
33 else:
34     fine_tune = pkl.load(open('holts.pkl','rb'))
35     print('Parameters already tuned for Triple smoothening !!!! ')
36 print('Time required to run this cell:',datetime.now()-start)

```

Parameters already tuned for Triple smoothening !!!!
Time required to run this cell: 0:00:00.001997

```

In [87]: 1 ## Lets get the parameters which gives minimum mape value
2 from operator import itemgetter
3 best_params = fine_tune[np.argmin(list(map(itemgetter(3),fine_tune)))]
4 print('For triple exponential smoothening the minimum mape value was found to be ',best_params[3])
5 print('The best alpha,beta and gamma parameters for triple exponential smoothening having minimum MAPE value are',best_params[:3])

```

For triple exponential smoothening the minimum mape value was found to be 4.187637418634715
The best alpha,beta and gamma parameters for triple exponential smoothening having minimum MAPE value are (0.25, 0.1, 0.2)

```
In [88]: 1 ### train with best_alpha ,alpha beta and gama
2 holt_winter_avg = []
3 a,b,g=best_params[:3]
4 predict_weeks = 39
5 for s in ustore:
6     for d in udept:
7         if (s,d) not in store_dept_with_all_zeros_sales:
8             data=trte_basic_feat.loc[(trte_basic_feat.Store ==s )&(trte_basic_feat.Dept==d)
9                                     &(trte_basic_feat.set_type=='Train')][ 'Weekly_Sales'].reset_index(drop=True)
10             predict_values_2,abs_val = triple_exponential_smoothing(data, season_len, a, b, g, predict_weeks)
11             holt_winter_avg.extend(predict_values_2)
```

```
In [89]: 1 trte_featured['holt_avg'] = np.round(holt_winter_avg,3)
```

```
In [90]: 1 ### Lets Look at our featured data
2 trte_featured.head(2)
```

Out[90]:

| | Store | Dept | Date | Weekly_Sales | IsHoliday | Temperature | Markdown1 | Markdown2 | Markdown3 | Markdown5 | CPI | Unemployment | Type | Size | set_type | week | year | month | day | major_holiday |
|---|-------|------|------------|--------------|-----------|-------------|-------------|-------------|-------------|-------------|------------|--------------|------|--------|----------|------|------|-------|-----|---------------|
| 0 | 1 | 1 | 2010-02-05 | 24924.50 | False | 42.31 | 8536.592778 | 3346.401918 | 1670.797978 | 4428.307667 | 211.096358 | 8.106 | A | 151315 | Train | 5 | 2010 | 2 | 5 | 0.0 |
| 1 | 1 | 1 | 2010-02-12 | 46039.49 | True | 38.51 | 8536.592778 | 3346.401918 | 1670.797978 | 4428.307667 | 211.242170 | 8.106 | A | 151315 | Train | 6 | 2010 | 2 | 12 | 4.0 |

Label encoding Categorical features

```
In [91]: 1 trte_featured['IsHoliday'] = trte_featured['IsHoliday'].map({True:1,False:0})
2 trte_featured['Type'] = trte_featured['Type'].map({'A':1,"B":2,"C":3})
3
4 trte_basic_feat['IsHoliday'] = trte_basic_feat['IsHoliday'].map({True:1,False:0})
5 trte_basic_feat['Type'] = trte_featured['Type']
6
7
8 trainset_merged['IsHoliday'] = trainset_merged['IsHoliday'].map({True:1,False:0})
9 trainset_merged['Type'] = trainset_merged['Type'].map({'A':1,"B":2,"C":3})
10
11
12 testset_merged['IsHoliday'] = testset_merged['IsHoliday'].map({True:1,False:0})
13 testset_merged['Type'] = testset_merged['Type'].map({'A':1,"B":2,"C":3})
14
15
16
17
```

```
In [92]: 1 ### Let's put the dataframe to csv
2 trte_featured.to_csv('trte_fetaured.csv',index=False,encoding='utf-8', date_format='%Y-%m-%d')
3 trainset_merged.to_csv('trainset_merged.csv',index=False,encoding='utf-8', date_format='%Y-%m-%d')
4 testset_merged.to_csv('testset_merged.csv',index=False,encoding='utf-8', date_format='%Y-%m-%d')
5 trte_basic_feat.to_csv('trte_basic_feat.csv',index=False,encoding='utf-8', date_format='%Y-%m-%d')
```

Feature Importance using Forward Feature Selection for basic features

```
In [75]: 1 input_ = trte_basic_feat.loc[trte_basic_feat.set_type=='Train']
2 input_['Date'] = input_.Date.map(datetime.toordinal)
3 input_ = input_.fillna(-1).reset_index(drop=True)
4 output = input_.Weekly_Sales
5 input_ = input_.drop(columns = ['Weekly_Sales', 'set_type'],axis=1)
```

```
In [76]: 1 seqfs = SFS(RandomForestRegressor(n_estimators=200,random_state=2),k_features='best',forward=True,
2           floating=False,verbose=2,scoring='neg_mean_absolute_error',cv=0)
```

```
In [103]: 1 seqfs.fit(input_,output,custom_feature_names=input_.columns)
```

```
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 16.2s remaining: 0.0s
[Parallel(n_jobs=1)]: Done 18 out of 18 | elapsed: 13.6min finished
```

```
[2021-07-20 13:49:48] Features: 1/18 -- score: -7998.225877432984[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 42.1s remaining: 0.0s
[Parallel(n_jobs=1)]: Done 17 out of 17 | elapsed: 20.3min finished
```

```
[2021-07-20 14:10:09] Features: 2/18 -- score: -2370.980427488272[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 2.7min remaining: 0.0s
[Parallel(n_jobs=1)]: Done 16 out of 16 | elapsed: 28.0min finished
```

```
[2021-07-20 14:38:10] Features: 3/18 -- score: -542.2262471597617[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 2.9min remaining: 0.0s
[Parallel(n_jobs=1)]: Done 15 out of 15 | elapsed: 51.0min finished
```

```
[2021-07-20 15:29:11] Features: 4/18 -- score: -468.4979528302423[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 3.7min remaining: 0.0s
[Parallel(n_jobs=1)]: Done 14 out of 14 | elapsed: 67.3min finished
```

```
In [140]: 1 df = pd.DataFrame.from_dict(seqfs.get_metric_dict()).T
2 sorted_df = df[['feature_names', 'avg_score']].sort_values(by='avg_score',ascending=False).reset_index(drop=True)
```

```
In [149]: 1 sorted_df
```

Out[149]:

| | feature_names | avg_score |
|----|---|-----------|
| 0 | (Store, Dept, Date, IsHoliday, Size, week, mon... | -384.817 |
| 1 | (Store, Dept, Date, IsHoliday, Type, Size, wee... | -384.818 |
| 2 | (Store, Dept, Date, Size, week, month, day, ma... | -384.846 |
| 3 | (Store, Dept, Date, IsHoliday, Type, Size, wee... | -384.95 |
| 4 | (Store, Dept, Date, Size, week, day, major_hol... | -385.134 |
| 5 | (Store, Dept, Date, IsHoliday, CPI, Type, Size... | -387.06 |
| 6 | (Store, Dept, Date, Size, week, day) | -387.828 |
| 7 | (Store, Dept, Date, IsHoliday, CPI, Unemployme... | -390.099 |
| 8 | (Store, Dept, Date, IsHoliday, Temperature, CP... | -394.64 |
| 9 | (Store, Dept, Date, IsHoliday, Temperature, Ma... | -409.507 |
| 10 | (Store, Dept, Date, week, day) | -419.531 |
| 11 | (Store, Dept, Date, IsHoliday, Temperature, Ma... | -419.913 |
| 12 | (Store, Dept, Date, IsHoliday, Temperature, Ma... | -425.013 |
| 13 | (Store, Dept, Date, IsHoliday, Temperature, Ma... | -429.66 |
| 14 | (Store, Dept, Date, week) | -468.498 |
| 15 | (Store, Dept, Date) | -542.226 |
| 16 | (Store, Dept) | -2370.98 |
| 17 | (Dept,) | -7998.23 |

```
In [155]: 1 print('Below feature combination gives the lowest mean absolute error on performing model:\n ',sorted_df.iloc[1]['feature_names'])
```

Below feature combination gives the lowest mean absolute error on performing model:
('Store', 'Dept', 'Date', 'IsHoliday', 'Type', 'Size', 'week', 'month', 'day', 'major_holiday')

- REfs: <https://plotly.com/python/figure-labels/> (<https://plotly.com/python/figure-labels/>).

<https://machinelearningmastery.com/decompose-time-series-data-trend-seasonality/#:~:text=Level%3A%20The%20average%20value%20in,random%20variation%20in%20the%20series>
(<https://machinelearningmastery.com/decompose-time-series-data-trend-seasonality/#:~:text=Level%3A%20The%20average%20value%20in,random%20variation%20in%20the%20series>).

<https://towardsdatascience.com/moving-averages-in-python-16170e20f6c> (<https://towardsdatascience.com/moving-averages-in-python-16170e20f6c>) <https://www.analyticsvidhya.com/blog/2019/12/6-powerful-feature-engineering-techniques-time-series/> (<https://www.analyticsvidhya.com/blog/2019/12/6-powerful-feature-engineering-techniques-time-series/>).

<https://www.r-craft.org/r-news/using-anova-to-get-correlation-between-categorical-and-continuous-variables/> (<https://www.r-craft.org/r-news/using-anova-to-get-correlation-between-categorical-and-continuous-variables/>).

