

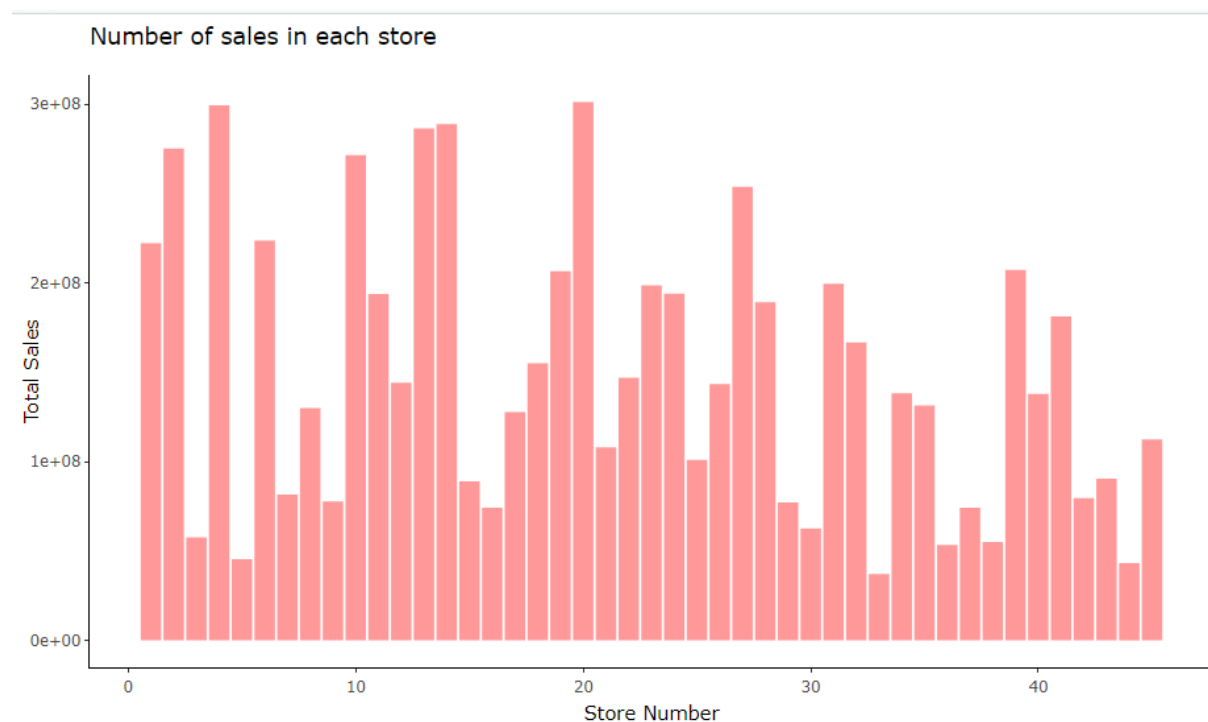
Analysis Tasks

Basic Statistics tasks

- Which store has maximum sales

Solution

```
1 library(ggplot2)
2 library(plotly)
3 library(data.table)
4 library(stats)
5 library(usdm)
6 library(caTools)
7 walmart<-read.csv('D:\\Data_Analysis_Simplilearn_materials\\Data_Science_with_R\\walmart_store_sales.csv')
8 #view(walmart)
9
10 # Analysis Task
11
12 # Basic Statistics tasks
13
14 # 1. which store has maximum sales
15
16 store_sales<-summarise(group_by(walmart,Store),sum(weekly_sales),mean(weekly_sales),sd(weekly_sales))
17 names(store_sales)<-c('Store','Total_sales','avg_weekly_sales','sales_sd')
18
19 ggplot(store_sales,aes(Store>Total_sales))+geom_col(fill='red',alpha=0.4)+
20   labs(x='Store Number',y='Total Sales',title='Number of sales in each store')+theme_classic()
21 ggplotly()
22 cat("Store number ",head(arrange(store_sales,-Total_sales),1)$Store,"has maximum sales =",
23     head(arrange(store_sales,-Total_sales),1)$Total_sales)
24
```



```
> cat("Store number ",head(arrange(store_sales,-Total_sales),1)$Store,"has maximum sales =",
+     head(arrange(store_sales,-Total_sales),1)$Total_sales)
Store number 20 has maximum sales = 301397792
```

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

```

25 # 2. Which store has maximum standard deviation i.e., the sales vary a lot.
26 # Also, find out the coefficient of mean to standard deviation
27
28 ggplot(store_sales,aes(store,sales_sd))+geom_col(fill='red',alpha=0.4)+
29   labs(x='Store Number',y='Standard Deviation of sales',title='Standard Deviation of sales in each store')+
30   theme_classic()
31 ggplotly()
32 cat("Store number ",head(arrange(store_sales,-sales_sd),1)$store,"has maximum standard deviation of sales =",
33     head(arrange(store_sales,-sales_sd),1)$sales_sd)
34
35 store_sales<-mutate(store_sales,sales_mean_to_sd=avg_weekly_sales/sales_sd)
36 ggplot(store_sales,aes(store,sales_mean_to_sd))+geom_col(fill='red',alpha=0.4)+
37   labs(x='Store Number',y='Mean/Std Dev',title='Mean/Standard Deviation of sales in each store')+
38   theme_classic()
39 ggplotly()
40 cat("Store number ",head(arrange(store_sales,-sales_mean_to_sd),1)$store,
41     "has maximum mean to standard deviation ratio of sales =",
42     head(arrange(store_sales,-sales_mean_to_sd),1)$sales_mean_to_sd)
43
44

```

Standard Deviation of sales in each store

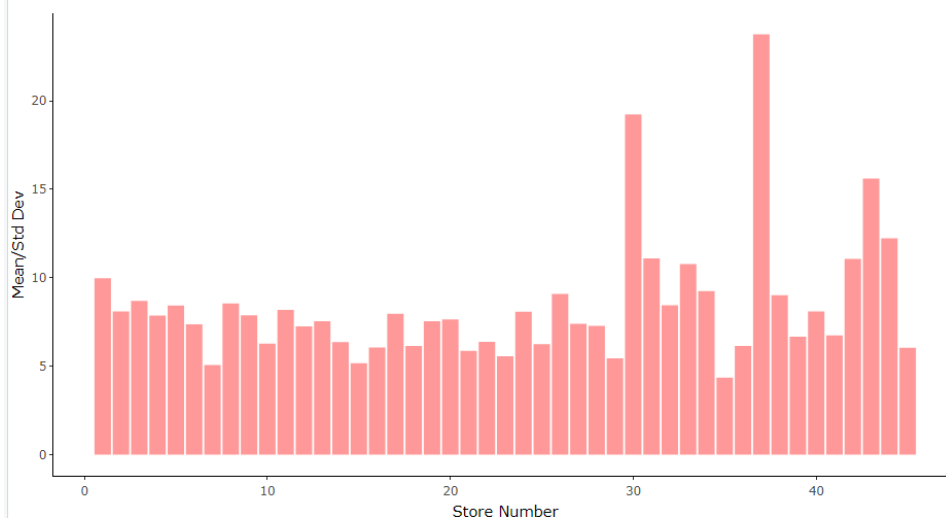


```

> cat("Store number ",head(arrange(store_sales,-sales_sd),1)$store,"has maximum standard deviation of sales =",
+     head(arrange(store_sales,-sales_sd),1)$sales_sd)
Store number 14 has maximum standard deviation of sales = 317569.9
>

```

Mean/Standard Deviation of sales in each store



```

> cat("Store number ",head(arrange(store_sales,-sales_mean_to_sd),1)$store,
+     "has maximum mean to standard deviation ratio of sales =",
+     head(arrange(store_sales,-sales_mean_to_sd),1)$sales_mean_to_sd)
Store number 37 has maximum mean to standard deviation ratio of sales = 23.76193
>

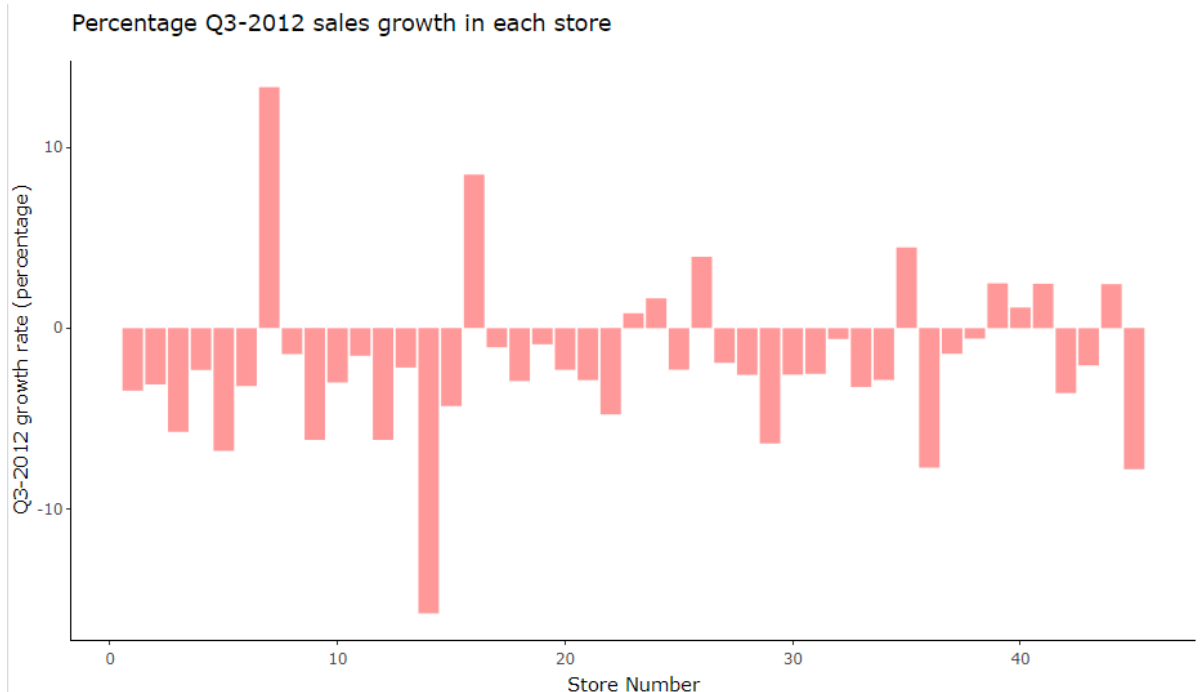
```

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

```

44
45 # 3. Which store/s has good quarterly growth rate in Q3'2012
46
47 walmart$Date<-as.Date(walmart$Date,format='%d-%m-%Y')
48 storewise_sales_Q3_2012<-walmart%>%filter(Date>='2012-07-01',Date<='2012-09-30')%>%
49   group_by(Store)%>%summarise(sum(weekly_sales))
50 names(storewise_sales_Q3_2012)<-c('Store','sales_q3_2012')
51 storewise_sales_Q2_2012<-walmart%>%filter(Date>='2012-04-01',Date<='2012-06-30')%>%
52   group_by(Store)%>%summarise(sum(weekly_sales))
53 names(storewise_sales_Q2_2012)<-c('Store','sales_q2_2012')
54 storewise_sales_Q2_Q3_2012<-merge(storewise_sales_Q2_2012,storewise_sales_Q3_2012)
55 storewise_sales_Q2_Q3_2012$q3_growth_rate<-
56   ((storewise_sales_Q2_Q3_2012$sales_q3_2012/storewise_sales_Q2_Q3_2012$sales_q2_2012)-1)*100;
57
58 ggplot(storewise_sales_Q2_Q3_2012,aes(Store,q3_growth_rate))+
59   geom_col(fill='red',alpha=0.4)+labs(x='Store Number',y='Q3-2012 growth rate (percentage)',
60   title='Percentage Q3-2012 sales growth in each store')+theme_classic()
61 ggplotly()
62 print('Top 4 stores with good Q3 growth rate')
63 print(head(arrange(storewise_sales_Q2_Q3_2012,-q3_growth_rate),4))
64

```



```

> print('Top 4 stores with good Q3 growth rate')
[1] "Top 4 stores with good Q3 growth rate"
> print(head(arrange(storewise_sales_Q2_Q3_2012,-q3_growth_rate),4))
  Store sales_q2_2012 sales_q3_2012 q3_growth_rate
1     7       7290859       8262787      13.330776
2    16       6564336       7121542       8.488378
3    35      10838313      11322421       4.466637
4    26      13155336      13675692       3.955478
>

```

- 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

```

64
65 # 4. Some holidays have a negative impact on sales. Find out holidays which have higher
66 # sales than the mean sales in non-holiday season for all stores together
67 holidayType<-function(x){
68   if(x=='February'){return('Super Bowl')}
69   else if(x=='September'){return('Labour Day')}
70   else if(x=='November'){return('Thanksgiving')}
71   else if(x=='December'){return('Christmas')}
72   else{return('Not Holiday Month')}};
73
74
75 daywise_sales<-summarise(group_by(walmart,Date,Holiday_Flag),sum(weekly_sales));
76 names(daywise_sales)<-c("Date","Holiday_Flag","Total_weekly_Sales");
77 daywise_sales<-as.data.frame(daywise_sales)
78 working_week_sales<-summarise(filter(daywise_sales,Holiday_Flag==0),mean(Total_weekly_sales))[[1]]
79 print("Holiday weeks when sales are more than mean sales")
80 print(mutate(filter(daywise_sales,Holiday_Flag==1,Totale_weekly_sales>working_week_sales),
81   festival_season=apply(months(Date),HolidayType)))
82
83 holiday_sales<-summarise(group_by(filter(daywise_sales,Holiday_Flag==1),months(Date)),
84   mean(Total_weekly_sales))
85 names(holiday_sales)<-c('HolidayMonth','MeanWeeklySales')
86 holiday_sales<-as.data.frame(holiday_sales)
87 holiday_sales<-mutate(holiday_sales,Holiday_Festival=apply(HolidayMonth,HolidayType))
88 print("Festive Seasons when mean sales is higher than mean sales in non-holiday season")
89 print(filter(holiday_sales,MeanWeeklySales>working_week_sales))
90
91
92

```

```

> print("Holiday weeks when sales are more than mean sales")
[1] "Holiday weeks when sales are more than mean sales"
> print(mutate(filter(daywise_sales,Holiday_Flag==1,Totale_weekly_sales>working_week_sales),
+   festival_season=apply(months(Date),HolidayType)))
  Date Holiday_Flag Total_weekly_Sales festival_season
1 2010-02-12         1         48336678      Super Bowl
2 2010-11-26         1         65821003    Thanksgiving
3 2011-02-11         1         47336193      Super Bowl
4 2011-11-25         1         66593605    Thanksgiving
5 2012-02-10         1         50009408      Super Bowl
6 2012-09-07         1         48330059    Labour Day
> holiday_sales<-summarise(group_by(filter(daywise_sales,Holiday_Flag==1),months(Date)),
+   mean(Total_weekly_sales))
> names(holiday_sales)<-c('HolidayMonth','MeanWeeklySales')
> holiday_sales<-as.data.frame(holiday_sales)
> holiday_sales<-mutate(holiday_sales,Holiday_Festival=apply(HolidayMonth,HolidayType))
> print("Festive Seasons when mean sales is higher than mean sales in non-holiday season")
[1] "Festive Seasons when mean sales is higher than mean sales in non-holiday season"
> print(filter(holiday_sales,MeanWeeklySales>working_week_sales))
  HolidayMonth MeanWeeklySales Holiday_Festival
1    February         48560759      Super Bowl
2    November         66207304    Thanksgiving
3    September         46909228    Labour Day
>

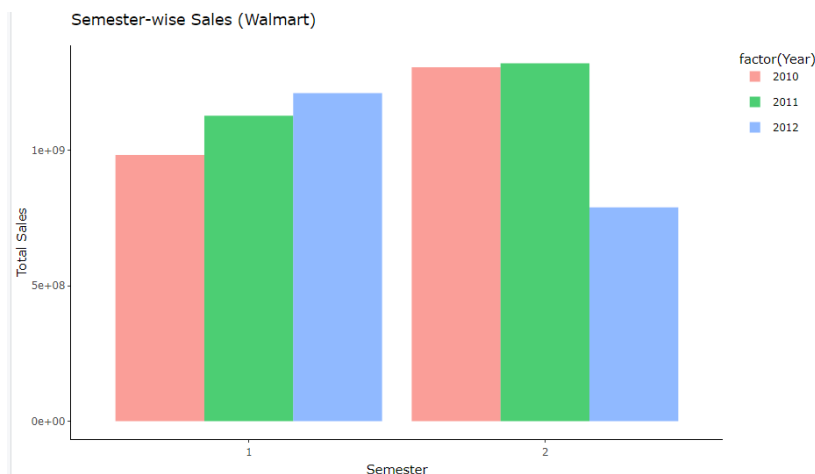
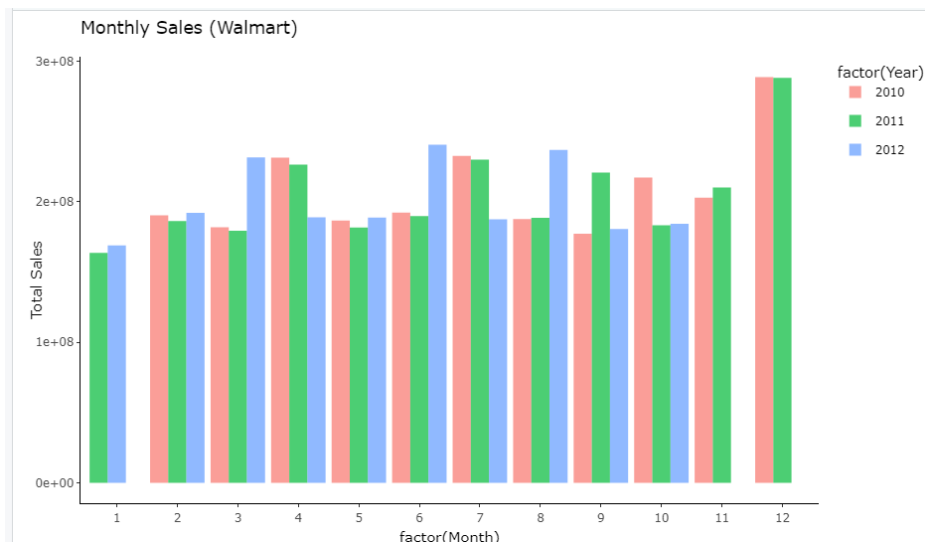
```

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

```

91
92
93 # 5. Provide a monthly and semester view of sales in units and give insights
94 monthly_sales<-summarise(group_by(walmart,month(as.IDate(Date))),year(as.IDate(Date))),sum(weekly_sales))
95 names(monthly_sales)<-c("Month","Year","TotalSales")
96 arrange(monthly_sales,Month,Year)
97
98 view(monthly_sales)
99
100 ggplot(monthly_sales,aes(x=factor(Month),y=TotalSales))+
101   geom_col(aes(fill=factor(Year)),alpha=0.7,position = position_dodge(preserve = "single"))+
102   labs(y='Total Sales',title='Monthly Sales (walmart)')+theme_classic()+
103   scale_x_discrete(name='Month',labels=month_name)
104   ggplotly()
105
106 monthly_sales=as.data.frame(monthly_sales)
107 monthly_sales$Semester=ifelse(monthly_sales$Month<=6,1,2)
108 view(summarise(group_by(monthly_sales,Semester,Year),TotalSales=sum(TotalSales)))
109 ggplot(summarise(group_by(monthly_sales,Semester,Year),TotalSales=sum(TotalSales)),
110   aes(x=factor(Semester),y=TotalSales))+
111   geom_col(aes(fill=factor(Year)),alpha=0.7,position = position_dodge(preserve = "single"))+
112   labs(x='Semester',y='Total Sales',title='Semester-wise Sales (walmart)')+theme_classic()
113   ggplotly()
114

```



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.

```
114
115 # Statistical Model
116
117 # For Store 1 – Build prediction models to forecast demand
118
119 # 1. Linear Regression – Utilize variables like date and restructure dates as 1 for 5 Feb 2010
120 # (starting from the earliest date in order). Hypothesize if CPI, unemployment,
121 # and fuel price have any impact on sales.
122
123 # 2. Change dates into days by creating new variable.
124
125 view(walmart)
126 walmart_store1<-filter(walmart,store==1);
127
128 view(walmart_store1)
129
130
131
132
133 walmart_store1$Day_number<-as.numeric(walmart_store1$Date-as.Date('2010-02-04'))
134 walmart_store1$month=month(walmart_store1$Date);
135 walmart_store1$year=year(walmart_store1$Date);
136 walmart_store1$Super_Bowl=as.numeric(walmart_store1$Holiday_Flag & walmart_store1$month==2)
137 walmart_store1$Labour_Day=as.numeric(walmart_store1$Holiday_Flag & walmart_store1$month==9)
138 walmart_store1$Thanksgiving=as.numeric(walmart_store1$Holiday_Flag & walmart_store1$month==11)
139 walmart_store1$Christmas=as.numeric(walmart_store1$Holiday_Flag & walmart_store1$month==12)
140 view(walmart_store1)
141
```

For training data

```
141
142 set.seed(101);
143 train_sample<-sample.split(walmart_store1$weekly_Sales,splitRatio = 0.8)
144 train_data<-subset(walmart_store1,train_sample==T);
145 test_data<-subset(walmart_store1,train_sample==F);
146
147 # Impact of CPI on sales
148 ggplot(train_data,aes(CPI,weekly_Sales))+geom_point()
149 cor.test(train_data$CPI,train_data$weekly_Sales,method='spearman')
150
```

```
> cor.test(train_data$CPI,train_data$weekly_Sales,method='spearman')

Spearman's rank correlation rho

data:  train_data$CPI and train_data$weekly_sales
s = 174142, p-value = 0.001521
alternative hypothesis: true rho is not equal to 0
sample estimates:
rho
0.2947004
```

```

150
151 # Impact of Unemployment on sales
152 ggplot(train_data,aes(Unemployment,weekly_sales))+geom_point()
153 cor.test(train_data$Unemployment,train_data$weekly_sales,method='spearman')
154

```

```

> ggplot(train_data,aes(Unemployment,weekly_sales))+geom_point()
> cor.test(train_data$Unemployment,train_data$weekly_sales,method='spearman')

spearman's rank correlation rho

data: train_data$Unemployment and train_data$weekly_sales
S = 292195, p-value = 0.05076
alternative hypothesis: true rho is not equal to 0
sample estimates:
rho
-0.1834317

```

```

154
155 # Impact of Fuel Price on sales
156 ggplot(train_data,aes(Fuel_Price,weekly_sales))+geom_point()
157 cor.test(train_data$Fuel_Price,train_data$weekly_sales,method='spearman')
158

```

```

> cor.test(train_data$Fuel_Price,train_data$weekly_sales,method='spearman')

spearman's rank correlation rho

data: train_data$Fuel_Price and train_data$weekly_sales
S = 190905, p-value = 0.01524
alternative hypothesis: true rho is not equal to 0
sample estimates:
rho
0.2268093

```

```

159 # Impact of Temperature on sales
160 ggplot(train_data,aes(Temperature,weekly_sales))+geom_point()
161 cor.test(train_data$Temperature,train_data$weekly_sales,method='spearman')
162

```

```

> ggplot(train_data,aes(Temperature,weekly_sales))+geom_point()
> cor.test(train_data$Temperature,train_data$weekly_sales,method='spearman')

spearman's rank correlation rho

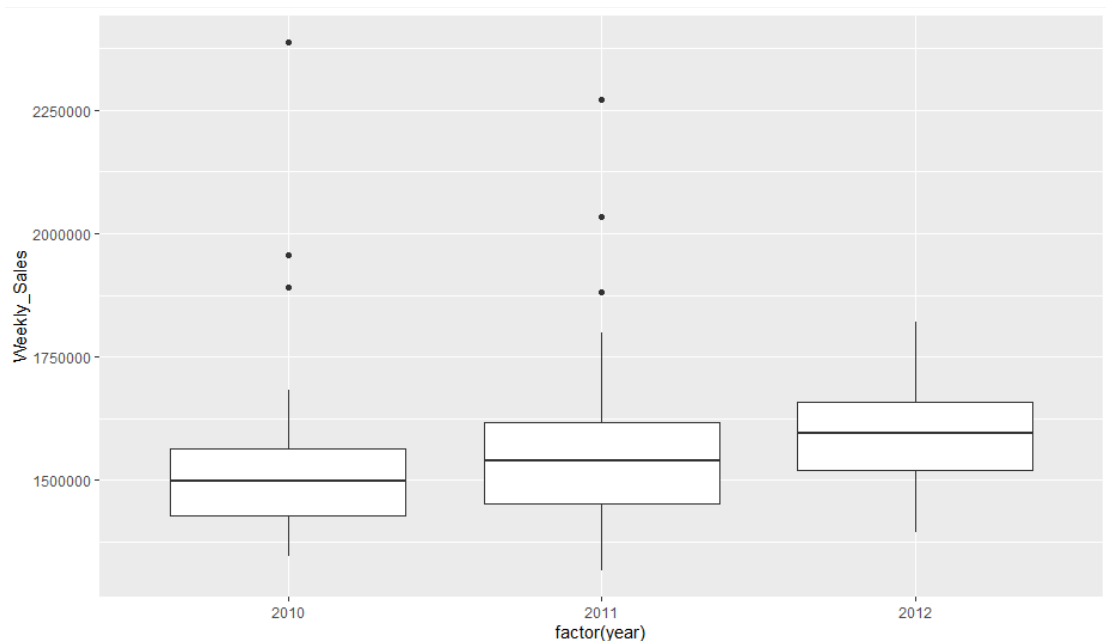
data: train_data$Temperature and train_data$weekly_sales
S = 293284, p-value = 0.04548
alternative hypothesis: true rho is not equal to 0
sample estimates:
rho
-0.1878415

```

```

163 # Impact of year on sales
164 ggplot(train_data,aes(factor(year),weekly_sales))+geom_boxplot()
165 cor.test(train_data$year,train_data$weekly_sales,method='spearman')

```



```
> cor.test(train_data$year,train_data$weekly_sales,method='spearman')
```

Spearman's rank correlation rho

```

data: train_data$year and train_data$weekly_sales
S = 177268, p-value = 0.002365
alternative hypothesis: true rho is not equal to 0
sample estimates:
      rho
0.2820402

```

```

166
167 # Impact of day number on sales
168 ggplot(train_data,aes(Day_number,weekly_sales))+geom_point()
169 cor.test(train_data$Day_number,train_data$weekly_sales,method='spearman')
170

```

```
> cor.test(train_data$Day_number,train_data$weekly_sales,method='spearman')
```

Spearman's rank correlation rho

```

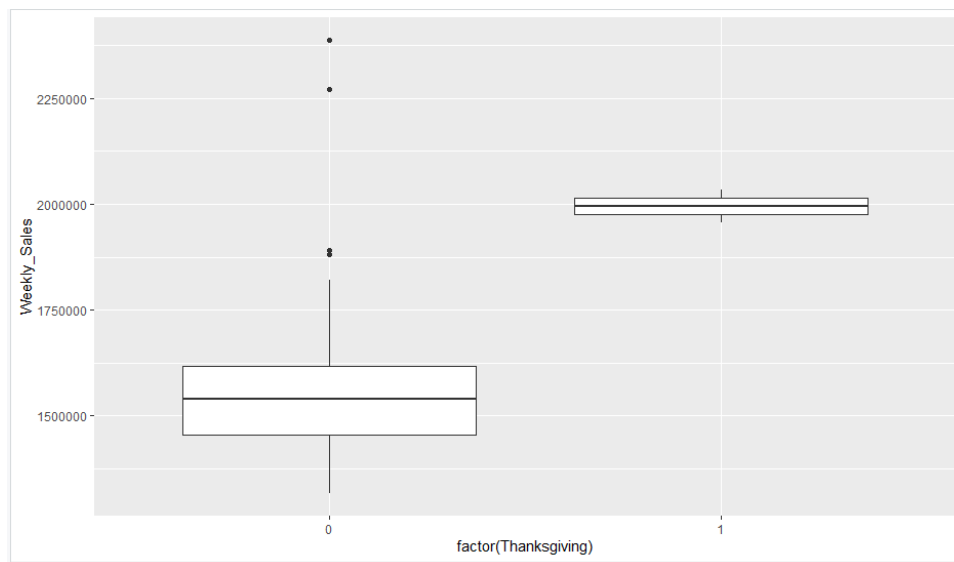
data: train_data$Day_number and train_data$weekly_sales
S = 173870, p-value = 0.001458
alternative hypothesis: true rho is not equal to 0
sample estimates:
      rho
0.295802

```

```

171
172 # Impact of Thanksgiving season on sales
173 ggplot(train_data,aes(factor(Thanksgiving),weekly_sales))+geom_boxplot()
174 cor.test(train_data$Thanksgiving,train_data$weekly_sales,method='spearman')

```

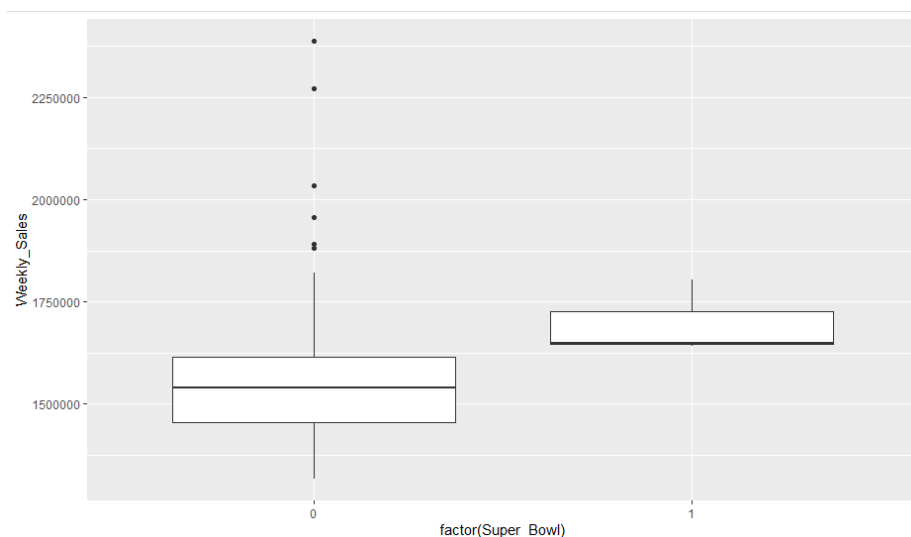



```
> ggplot(train_data,aes(factor(Thanksgiving),weekly_sales))+geom_boxplot()
> cor.test(train_data$Thanksgiving,train_data$weekly_sales,method='spearman')
```

Spearman's rank correlation rho

```
data: train_data$Thanksgiving and train_data$weekly_sales
S = 192763, p-value = 0.01907
alternative hypothesis: true rho is not equal to 0
sample estimates:
      rho
0.2192816
```

```
175 # Impact of Super Bowl season on sales
176 ggplot(train_data,aes(factor(Super_Bowl),weekly_sales))+geom_boxplot()
177 cor.test(train_data$Super_Bowl,train_data$weekly_sales,method='spearman')
```



```
> cor.test(train_data$Super_Bowl,train_data$weekly_sales,method='spearman')
```

Spearman's rank correlation rho

```
data: train_data$Super_Bowl and train_data$weekly_sales
S = 196127, p-value = 0.02815
alternative hypothesis: true rho is not equal to 0
sample estimates:
      rho
0.2056588
```

Based on training dataset, Weekly Sales of Store 1 has correlation with factors like CPI, Fuel Price, year, Temperature, Thanksgiving Season, Super Bowl Season and Day Number. Since, probability of these features being uncorrelated with Weekly_Sales < 0.05, we can reject the Null Hypothesis which states that they are not correlated with Weekly_Sales.

However, Unemployment column has p-value > 0.05 and therefore we cannot reject Null Hypothesis for Unemployment column for Training Dataset and do not include it in Linear Regression Model

Now, Variance Inflation Factor (VIF) is calculated to check for Multi-collinearity.

Features with greatest VIF are removed one by one and VIF is recomputed till all remaining features have VIF < 10

```
179 # Sales of Store 1 depend on factors like CPI, Fuel Price, year, Temperature,
180 # Thanksgiving Season, Super Bowl season, Day Number
181
182 vif(train_data[c('CPI', 'Fuel_Price', 'year', 'Temperature', 'Thanksgiving', 'Super_Bowl', 'Day_number')])
183 vif(train_data[c('CPI', 'Fuel_Price', 'year', 'Temperature', 'Thanksgiving', 'Super_Bowl')])
184 vif(train_data[c('CPI', 'Fuel_Price', 'Temperature', 'Thanksgiving', 'Super_Bowl')])
185
```



```
cannot compute exact p-value with ties
> vif(train_data[c('CPI', 'Fuel_Price', 'year', 'Temperature', 'Thanksgiving', 'Super_Bowl', 'Day_number')])
Variables      VIF
1      CPI 22.122533
2 Fuel_Price 3.359369
3      year 13.707501
4 Temperature 1.296530
5 Thanksgiving 1.057822
6 Super_Bowl 1.122562
7 Day_number 23.204299
> vif(train_data[c('CPI', 'Fuel_Price', 'year', 'Temperature', 'Thanksgiving', 'Super_Bowl')])
Variables      VIF
1      CPI 10.012136
2 Fuel_Price 3.252955
3      year 12.759756
4 Temperature 1.255378
5 Thanksgiving 1.033406
6 Super_Bowl 1.112749
> vif(train_data[c('CPI', 'Fuel_Price', 'Temperature', 'Thanksgiving', 'Super_Bowl')])
Variables      VIF
1      CPI 2.427314
2 Fuel_Price 2.482621
3 Temperature 1.155696
4 Thanksgiving 1.008151
5 Super_Bowl 1.099281
>
```

After feature reduction using VIF, the remaining Features are CPI, Fuel_Price, Temperature, Thanksgiving, Super_Bowl

Now, we perform Linear Regression and remove features with $p\text{-val} > 0.05$, till essential features for Linear Regression model having $p\text{-value} < 0.05$ are remaining

```
185 Sales_estimate<-lm(weekly_Sales~CPI+Fuel_Price+Temperature+Thanksgiving+Super_Bowl,data=train_data)
186 Sales_estimate
187 summary(Sales_estimate)
```

```
> Sales_estimate<-lm(weekly_Sales~CPI+Fuel_Price+Temperature+Thanksgiving+Super_Bowl,data=train_data)
> Sales_estimate
```

```
Call:
lm(formula = weekly_Sales ~ CPI + Fuel_Price + Temperature +
    Thanksgiving + Super_Bowl, data = train_data)

Coefficients:
(Intercept)          CPI      Fuel_Price    Temperature  Thanksgiving      Super_Bowl
   -699213.0     11506.0    -15186.0    -2714.0      433227.0      82296.0

> summary(Sales_estimate)
```

```
Call:
lm(formula = weekly_Sales ~ CPI + Fuel_Price + Temperature +
    Thanksgiving + Super_Bowl, data = train_data)

Residuals:
    Min       1Q   Median       3Q      Max
-260692  -84196  -10910   63496  840658

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) -699213.4    932574.6  -0.750   0.45503
CPI           11505.7     4872.9    2.361   0.02001 *
Fuel_Price   -15185.5     51867.6  -0.293   0.77026
Temperature  -2714.1      989.6   -2.743   0.00714 **
Thanksgiving 433227.5    104163.5   4.159 6.43e-05 ***
Super_Bowl    82296.5     89209.1    0.923   0.35832
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 145400 on 108 degrees of freedom
Multiple R-squared:  0.2542,    Adjusted R-squared:  0.2197
F-statistic: 7.361 on 5 and 108 DF,  p-value: 5.676e-06
```

```
189
190 Sales_estimate<-lm(weekly_Sales~CPI+Temperature+Thanksgiving+Super_Bowl,data=train_data)
191 Sales_estimate
192 summary(Sales_estimate)
193
```

```
> Sales_estimate<-lm(weekly_Sales~CPI+Temperature+Thanksgiving+Super_Bowl,data=train_data)
> Sales_estimate
```

```
Call:
lm(formula = weekly_Sales ~ CPI + Temperature + Thanksgiving +
    Super_Bowl, data = train_data)

Coefficients:
(Intercept)          CPI      Temperature  Thanksgiving      Super_Bowl
   -511947.0     10426.0    -2750.0      435187.0      83053.0

> summary(Sales_estimate)
```

```
Call:
lm(formula = weekly_Sales ~ CPI + Temperature + Thanksgiving +
    Super_Bowl, data = train_data)

Residuals:
    Min       1Q   Median       3Q      Max
-262917  -84075  -9555   64438  839767

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) -511946.9    675804.1  -0.758   0.45036
CPI           10425.7     3170.7    3.288   0.00136 **
Temperature  -2750.3      977.7   -2.813   0.00582 **
Thanksgiving 435187.1    103511.4   4.204 5.39e-05 ***
Super_Bowl    83052.5     88796.9    0.935   0.35170
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 144800 on 109 degrees of freedom
Multiple R-squared:  0.2536,    Adjusted R-squared:  0.2262
F-statistic: 9.258 on 4 and 109 DF,  p-value: 1.771e-06
```

```

193 sales_estimate<-lm(weekly_Sales~CPI+Temperature+Thanksgiving,data=train_data)
194 sales_estimate
195 summary(sales_estimate)
196
197

```

```

> sales_estimate<-lm(weekly_Sales~CPI+Temperature+Thanksgiving,data=train_data)
> sales_estimate

Call:
lm(formula = weekly_Sales ~ CPI + Temperature + Thanksgiving,
    data = train_data)

Coefficients:
(Intercept)          CPI    Temperature    Thanksgiving
   -501424         10472         -3021          431527

> summary(sales_estimate)

Call:
lm(formula = weekly_Sales ~ CPI + Temperature + Thanksgiving,
    data = train_data)

Residuals:
    Min       1Q   Median       3Q      Max
-271450  -84654   -6415    60996   833573

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  -501423.9    675325.8   -0.742  0.45937
CPI           10472.3     3168.5     3.305  0.00128 **
Temperature   -3021.0     933.4    -3.237  0.00160 **
Thanksgiving  431526.8    103378.5     4.174  6e-05 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 144700 on 110 degrees of freedom
Multiple R-squared:  0.2476,    Adjusted R-squared:  0.2271
F-statistic: 12.07 on 3 and 110 DF,  p-value: 6.896e-07

```

Final Linear Regression Model

CPI → Prevailing consumer price index

Temperature → Temperature on the day of sale

Thanksgiving → This value is '1' if Thanksgiving holiday lies during the week of sale. Otherwise, it is '0'

Predicted Weekly Sales → Prediction of Weekly Sales

$$\text{Predicted Weekly Sales} = 10472.3 \times \text{CPI} - 3021 \times \text{Temperature} + 431526.8 \times \text{Thanksgiving}$$

Training Data has R^2 value of 0.248 and Adjusted R^2 value of 0.227

Performance on Test Dataset

```
197
198 sales_prediction<-predict(Sales_estimate,test_data[c('CPI','Temperature','Thanksgiving')])
199 testing_prediction<-cbind(sales_prediction,test_data$weekly_sales)
200 df.test.predict<-data.frame(testing_prediction)
201 colnames(df.test.predict)<-c('Prediction','True_value')
202 MSE_value<-mean((df.test.predict$Prediction-df.test.predict$True_value)^2)
203 cat("Mean Square Error is", MSE_value)
204 RMSE_value<-sqrt(mean((df.test.predict$Prediction-df.test.predict$True_value)^2))
205 cat("Root Mean Square Error is", RMSE_value)
206 SSR_value<-sum((df.test.predict$Prediction-mean(df.test.predict$True_value))^2);
207 SST_value<-sum((df.test.predict$True_value-mean(df.test.predict$True_value))^2);
208 test_R2_value<-SSR_value/SST_value;
209 cat("Coefficient of determination (R^2) on Test Data is equal to",test_R2_value )
```

```
> sales_prediction<-predict(Sales_estimate,test_data[c('CPI','Temperature','Thanksgiving')])
> testing_prediction<-cbind(sales_prediction,test_data$weekly_sales)
> df.test.predict<-data.frame(testing_prediction)
> colnames(df.test.predict)<-c('Prediction','True_value')
> MSE_value<-mean((df.test.predict$Prediction-df.test.predict$True_value)^2)
> cat("Mean Square Error is", MSE_value)
Mean Square Error is 13477465480
> RMSE_value<-sqrt(mean((df.test.predict$Prediction-df.test.predict$True_value)^2))
> cat("Root Mean Square Error is", RMSE_value)
Root Mean Square Error is 116092.5
> SSR_value<-sum((df.test.predict$Prediction-mean(df.test.predict$True_value))^2);
> SST_value<-sum((df.test.predict$True_value-mean(df.test.predict$True_value))^2);
> test_R2_value<-SSR_value/SST_value;
> cat("Coefficient of determination (R^2) on Test Data is equal to",test_R2_value )
Coefficient of determination (R^2) on Test Data is equal to 0.3594749
>
```

Test Data has R^2 value of 0.359