

# Predicting Weekly Sales at Walmart Stores

*August 31, 2015*

- 1. Executive Summary
- 2. Introduction
  - 2.1 About the Solution Environment
  - 2.2 About the Data
    - 2.2.1 The Challenge
  - 2.3 Getting the Data
    - 2.3.1 The Data Files
    - 2.3.2 Ingesting the Data
  - 2.4 R Libraries Used
- 3. Stage 1: Data Exploration and Preparation
  - 3.1 Summary Statistics
    - 3.1.1 The Training Dataset (train)
    - 3.1.2 The Stores Dataset (stores)
    - 3.1.3 The Features Dataset (features)
    - 3.1.4 The Test Dataset (test)
  - 3.2 Data Preparation - Merging the Datasets
    - 3.2.1 Merging Train and Stores Datasets
    - 3.2.2 Merging Train, Stores and Features Datasets
    - 3.2.3 Merging Test, Stores and Features Datasets
  - 3.3 Data Exploration
    - 3.3.1 Total Sales Per Department in each Store
    - 3.3.2 Store Total Sales Vs. Size
    - 3.3.3 Total Sales Per Week - Time Series
    - 3.3.4 Store-Department-wise Sales per Week - Time Series
- 4. Stage 2: Formal Statistical Inferences
  - 4.1 On Average, Do Holiday Weeks Spike Sales Up?
    - 4.1.1 The Hypotheses
    - 4.1.2 The Data
  - 4.2 Do Bigger Stores contribute to Higher Sales Figures?
- 5. Stage 3: Linear Regression: Predicting Weekly\_Sales

## 1. Executive Summary

Retail stores need to be able to predict sales forecasts for the future and study the effect how strategic offers affect sales, especially during holiday season. Since the number of days in holidays are limited, it becomes more challenging to be able to accurately predict how different aspects affect sales.

This report will focus on a Walmart Data set that has Department-wise Weekly Sales of 45 Walmart stores. It will attempt to create a predictive model and also discuss the extent to which different factors affect the sales.

## 2.1 About the Solution Environment

## 2.2 About the Data

The training data set has more than 400K records. The testing data set has over 100K records.

## 2.3 Getting the Data

The files available are the following:

### 2.3.1 The Data Files

### 2.3.1.1 stores.csv

### 2.3.1.2 train.csv

- Store: store number
- Dept: the department number
- Date: week date
- Weekly\_Sales: sales for the given department in the given store
- IsHoliday: whether the week is a special holiday week

### 2.3.1.3 test.csv

data:text/html;charset=utf-8,%3Cdiv%20id%3D%22header%22%20style%3D%22box-sizing%3A%20border-box%3B%20color%3A%20rgb(85%2C%2085%... 2/21

model with unseen data and can be evaluated by uploading the data set to Kaggle.

### 2.3.1.4 features.csv

This data file contains additional relevant information relating to the physical and business environment around the store. The fields are as follows:

- Store: store number
- Date: the week date
- Temperature: the average temperature in the region
- Fuel\_Price: cost of fuel in the region
- Markdown1-5: data related to the markdowns that Walmart is running. Markdown data is only available after November 2011 and is not available for all stores all the time. Any missing value is marked with an NA.
- CPI - the Consumer Price Index
- Unemployment - the unemployment rate
- IsHoliday - whether the week is a special holiday week

The four holidays fall in the following weeks in the data set:

- Super Bowl: 12-Feb-10, 11-Feb-11, 10-Feb-12, 8-Feb-13
- Labor Day: 10-Sep-10, 9-Sep-11, 7-Sep-12, 6-Sep-13
- Thanksgiving: 26-Nov-10, 25-Nov-11, 23-Nov-12, 29-Nov-13
- Christmas: 31-Dec-10, 30-Dec-11, 28-Dec-12, 27-Dec-13

## 2.3.2 Ingesting the Data

```
## Ingesting the data from the Data folder: Training Dataset  
train <- read.csv("Data/train.csv")
```

```
## Ingesting the data from the Data folder: Stores Dataset  
stores <- read.csv("Data/stores.csv")
```

```
## Ingesting the data from the Data folder: features Dataset  
features <- read.csv("Data/features.csv")
```

```
## Ingesting the data from the Data folder: testing Dataset  
test <- read.csv("Data/test.csv")
```

## 2.4 R Libraries Used

The following libraries are used in this report:

```
# Grammar of Graphics Plotting Library  
library(ggplot2)  
# To use 'melt'  
library(reshape2)  
# to enable commas in graphs  
library(scales)  
# to get the month number from date variable
```

```
library(lubridate)
```

## 3. Stage 1: Data Exploration and Preparation

### 3.1 Summary Statistics

#### 3.1.1 The Training Dataset (train)

```
## Structure of Train Dataset
str(train)
```

```
## 'data.frame':    421570 obs. of  5 variables:
##  $ Store      : int   1 1 1 1 1 1 1 1 1 1 ...
##  $ Dept       : int   1 1 1 1 1 1 1 1 1 1 ...
##  $ Date       : Factor w/ 143 levels "2010-02-05","2010-02-12",...: 1 2 3 4 5
## 6 7 8 9 10 ...
##  $ Weekly_Sales: num   24924 46039 41596 19404 21828 ...
##  $ IsHoliday   : logi  FALSE TRUE  FALSE FALSE FALSE FALSE ...
```

`Date` is ingested as factor (as opposed to being ingested as date type). There are 143 dates in total.

```
## Changing the Date from "Format" type to "Date" Type
train$Date <- as.Date(train$Date)
## Getting the summary of the Data
summary(train)
```

```
##      Store      Dept      Date      Weekly_Sales
##  Min.   : 1.0    Min.   : 1.00   Min.   :2010-02-05   Min.   : -4989
##  1st Qu.:11.0    1st Qu.:18.00   1st Qu.:2010-10-08   1st Qu.:  2080
##  Median :22.0    Median :37.00   Median :2011-06-17   Median :   7612
##  Mean   :22.2    Mean   :44.26   Mean   :2011-06-18   Mean   : 15981
##  3rd Qu.:33.0    3rd Qu.:74.00   3rd Qu.:2012-02-24   3rd Qu.: 20206
##  Max.   :45.0    Max.   :99.00   Max.   :2012-10-26   Max.   :693099
##  IsHoliday
##  Mode :logical
##  FALSE:391909
##  TRUE :29661
##  NA's :0
##
##
```

There is no missing data in the data set.

The starting date for training data set is 2010-02-05 . It starts on a Friday . The last date recorded in the data set is 2012-10-26 , which is also a Friday . There are 994 days between them - so the data consists of a total of 143 weeks of data.

It is interesting to note that for some departments the `Weekly_Sales` are **negative**. Returns and special offers cause these negative sales figures.

The **mean** is 15981.2581235 and **median** is 7612.03 . The mean and the median are very far apart, indicating that the data is skewed - in this case, extremely **right-skewed**. The following histogram depicts this relationship - where you can clearly observe the long tail towards the right making it extremely right-skewed.

data:text/html;charset=utf-8,%3Cdiv%20id%3D%22header%22%20style%3D%22box-sizing%3A%20border-box%3B%20color%3A%20rgb(85%2C%2085%... 5/21

```
## Structure of Stores Dataset
str(stores)
```

```
## 'data.frame':    45 obs. of  3 variables:
## $ Store: int  1 2 3 4 5 6 7 8 9 10 ...
## $ Type : Factor w/ 3 levels "A","B","C": 1 1 2 1 2 1 2 1 2 2 ...
## $ Size : int  151315 202307 37392 205863 34875 202505 70713 155078 125833 126
512 ...
```

```
## summary Statistics of Stores dataset
summary(stores)
```

```
##      Store      Type      Size
## Min.   : 1    A:22   Min.   : 34875
## 1st Qu.:12   B:17   1st Qu.: 70713
## Median :23   C: 6   Median :126512
## Mean   :23           Mean   :130288
## 3rd Qu.:34           3rd Qu.:202307
## Max.   :45           Max.   :219622
```

No missing data.

### 3.1.3 The Features Dataset (features)

```
## Structure of features dataset
str(features)
```

```
## 'data.frame':    8190 obs. of  12 variables:
## $ Store      : int  1 1 1 1 1 1 1 1 1 1 ...
## $ Date       : Factor w/ 182 levels "2010-02-05","2010-02-12",...: 1 2 3 4 5
6 7 8 9 10 ...
## $ Temperature : num  42.3 38.5 39.9 46.6 46.5 ...
## $ Fuel_Price  : num  2.57 2.55 2.51 2.56 2.62 ...
## $ Markdown1   : num  NA NA NA NA NA NA NA NA NA NA ...
## $ Markdown2   : num  NA NA NA NA NA NA NA NA NA NA ...
## $ Markdown3   : num  NA NA NA NA NA NA NA NA NA NA ...
## $ Markdown4   : num  NA NA NA NA NA NA NA NA NA NA ...
## $ Markdown5   : num  NA NA NA NA NA NA NA NA NA NA ...
## $ CPI         : num  211 211 211 211 211 ...
## $ Unemployment: num  8.11 8.11 8.11 8.11 8.11 ...
## $ IsHoliday   : logi  FALSE TRUE FALSE FALSE FALSE FALSE ...
```

`Date` is ingested as factor (as opposed to being ingested as date type). There are 182 dates in total. This data set is relevant for both the `train` and the `test` data set.

```
## Changing the Date from "Format" type to "Date" Type
features$Date <- as.Date(features$Date)
## Summary Statistics of Features Dataset
summary(features)
```



```
## Mean      :22.24    Mean      :44.34    Mean      :2013-03-14    NA's :0
## 3rd Qu.   :33.00    3rd Qu.   :74.00    3rd Qu.   :2013-05-24
## Max.      :45.00    Max.       :99.00    Max.      :2013-07-26
```

## 3.2 Data Preparation - Merging the Datasets

### 3.2.1 Merging Train and Stores Datasets

Since the `Type` & `Size` variables may influence the Weekly Sales, we are merging the `train` & data sets. We merge the data by `Store`.

```
## Merging train and stores by Store
trainStoresMerge <- merge(train , stores , by = "Store")
```

### 3.2.2 Merging Train, Stores and Features Datasets

Since `Markdown1-5` and other variables could play an important role at predicting `Weekly_Sales`, this should be merged with the `trainStoresMerge` data set. We merge the data by `Store` & `Date`.

```
## Merging trainStoresMerge and features datasets
trainStoresFeaturesMerge <-
  merge( trainStoresMerge , features , by = c( "Store" , "Date" ) )
## Clearing memory - removing intermediate datasets
rm( trainStoresMerge )
## Fixing the name of the Column
colnames(trainStoresFeaturesMerge)[5] <- "IsHoliday"
trainStoresFeaturesMerge$IsHoliday.y <- NULL
```

### 3.2.3 Merging Test, Stores and Features Datasets

We similarly merge the `test`, `stores` & `features` to create the `testStoresFeaturesMerge` data set.

```
## Merging test and stores by Store
testStoresMerge <- merge(test , stores , by = "Store")
## Merging testStoresMerge and features datasets
testStoresFeaturesMerge <-
  merge( testStoresMerge , features , by = c( "Store" , "Date" ) )
## Clearing Memory - removing intermediate Datasets
rm( testStoresMerge )
## Fixing the name of the Column
colnames(testStoresFeaturesMerge)[5] <- "IsHoliday"
testStoresFeaturesMerge$IsHoliday.y <- NULL
```

## 3.3 Data Exploration

### 3.3.1 Total Sales Per Department in each Store

The final goal of this report is to be able to predict the weekly sales for each department in a store. First



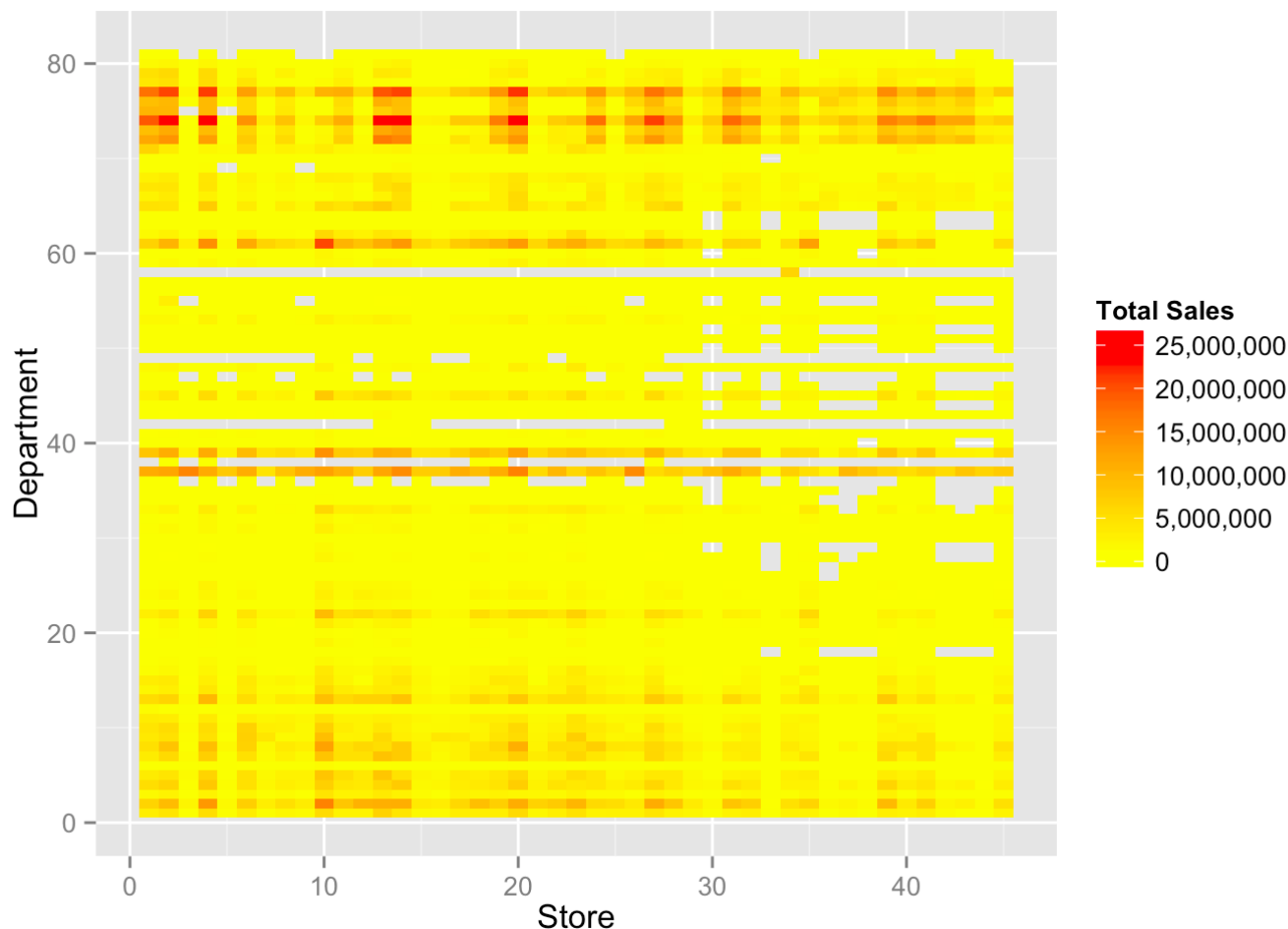
we would like to understand which departments are present in the 45 different stores and their total sales.

```
## running the sum function for each store & department
storeDeptTotalSales <- tapply(
  trainStoresFeaturesMerge$Weekly_Sales ,
  trainStoresFeaturesMerge[, c("Store", "Dept")] ,
  FUN = sum )
## Converting the matrix to a dataframe
storeDeptTotalSalesDataFrame <- as.data.frame( storeDeptTotalSales )
## Setting the Store Number into the table so we can analyze it further
storeDeptTotalSalesDataFrame$Store <-
  as.integer( rownames( storeDeptTotalSalesDataFrame ) )
## Move Store to the 1st column in the dataframe
storeDeptTotalSalesDataFrame <-
  storeDeptTotalSalesDataFrame[ , c( ncol(storeDeptTotalSalesDataFrame) , 1:ncol(
    storeDeptTotalSalesDataFrame)-1 )]
## Melting the columns into rows to enable analysis
storeDeptTotalSalesDataFrame <-
  melt(storeDeptTotalSalesDataFrame , id="Store" )
## removing the NA variables - where the department does not exist in a store
storeDeptTotalSalesDataFrame <- storeDeptTotalSalesDataFrame[ complete.cases(storeDeptTotalSalesDataFrame),]
## Renaming the Columns in the Dataframe
colnames( storeDeptTotalSalesDataFrame )[2:3] <- c("Dept" , "TotalSales" )
## Changing the Dept Type from String to Numeric
storeDeptTotalSalesDataFrame$Dept <-
  as.integer(storeDeptTotalSalesDataFrame$Dept)
## Freeing Memory - Removing the intermediate Matrix
rm( storeDeptTotalSales )
## printing out summary statistics
summary( storeDeptTotalSalesDataFrame)
```

```
##      Store      Dept      TotalSales
##  Min.   : 1.0    Min.   : 1.00    Min.   : -3567
##  1st Qu.:11.0    1st Qu.:19.00    1st Qu.: 137763
##  Median :22.0    Median :40.00    Median : 880317
##  Mean   :22.5    Mean   :40.49    Mean   : 2022582
##  3rd Qu.:33.0    3rd Qu.:62.00    3rd Qu.: 2609819
##  Max.   :45.0    Max.   :81.00    Max.   :26101498
```

### 3.3.1.1 Heatmap - Store & Department Total Sales

```
## Generating a Heatmap of the Department's Total Sales in each of 45 stores
ggplot( storeDeptTotalSalesDataFrame , aes(x = Store, y = Dept)) +
  geom_tile(aes(fill = TotalSales)) +
  scale_fill_gradient(
    low="yellow", high="red" , labels = comma , name="Total Sales") +
  scale_y_continuous( name="Department" )
```



From the heat map we can draw the following broad conclusions:

- The departments between 70-80 account for more sales than other departments
- Some departments are missing in some stores

```
## removing dataframe to free up memory
rm( storeDeptTotalSalesDataFrame )
```

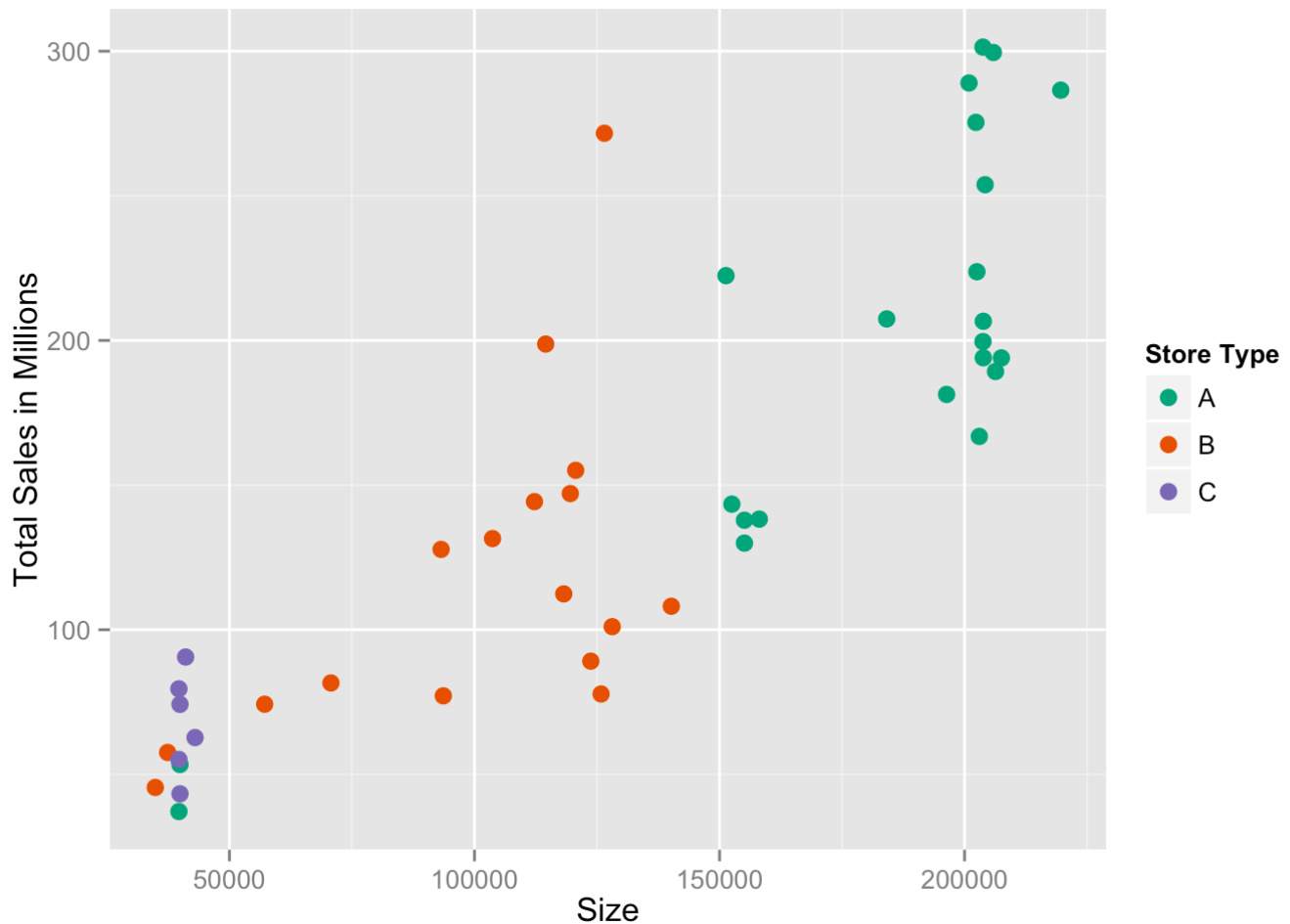
### 3.3.2 Store Total Sales Vs. Size

Plotting the total sales of a store vs. Store Size. We first calculate the total sales per Store and plot it as a response (y-axis) to the Store size (x-axis) to understand the relationship between them.

```
## Total Sales vs. Store Size - plotting the relationship
## calculating the sum of all the store sales
StoreTotalSales <-
  tapply(
    trainStoresFeaturesMerge$Weekly_Sales,
    trainStoresFeaturesMerge$Store,
    FUN = sum)
## converting the table to a DataFrame
stores$TotalSales <- StoreTotalSales
stores$TotalSalesInMillion <- stores$TotalSales/1000000
rm( StoreTotalSales )
```

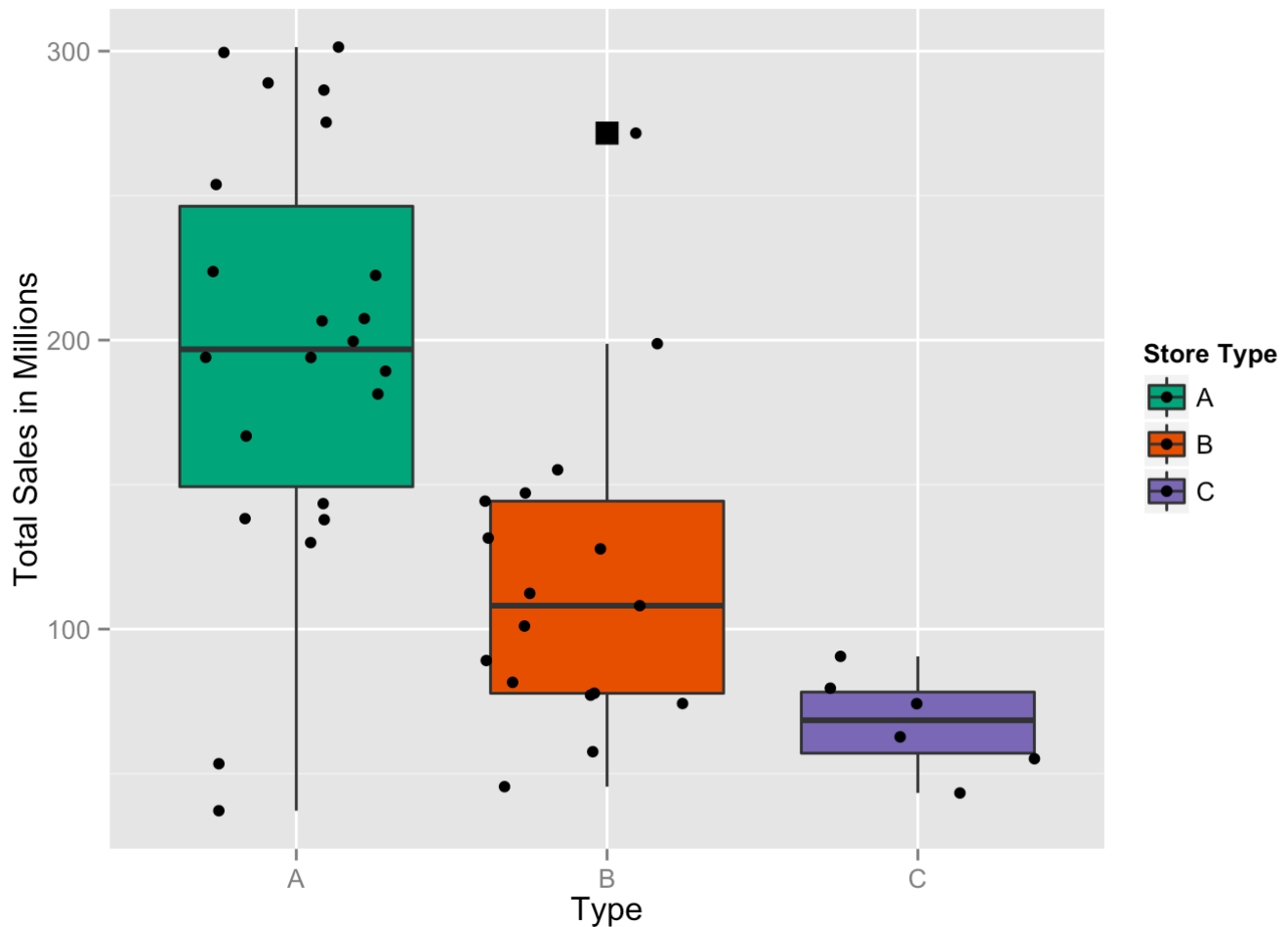
```
## Plotting the Total Sales vs. Store Size
```

```
ggplot( stores , aes(x=Size , y=TotalSalesInMillion , color = Type ) ) +
  geom_point( size=3 ) +
  scale_y_continuous(name="Total Sales in Millions" ) +
  scale_color_brewer(palette = "Dark2", name="Store Type" )
```



This plot indicates that there is a positive relationship between the size of the store and total sales. Also Type 'A' Stores are mostly larger stores with bigger sales and Type 'C' Stores are small with lower sales.

```
## box plot to show the summary statistics of the Type of Stores
ggplot(data=stores,
  aes(x=Type, y=TotalSalesInMillion, fill=Type) ) +
  geom_boxplot(outlier.shape = 15, outlier.size = 4) +
  ## to show how the individual store sales are distributed
  geom_jitter() +
  scale_y_continuous(name="Total Sales in Millions" ) +
  scale_fill_brewer(name = "Store Type" , palette = "Dark2")
```



```
## calculating the Summary Statistics for each Type
tableTypeWiseSummaryStatistics <-
  tapply(stores$TotalSalesInMillion , stores$Type , summary)

## Changing the Labels of the tabled Summary Statistics for printing
attributes(tableTypeWiseSummaryStatistics)$dimnames[[1]] <-
  c(
    "Type A Store Summary Statistics" ,
    "Type B Store Summary Statistics" ,
    "Type C Store Summary Statistics" )

## Printing Summary Statistics for each Type of Store (based on Total Sales)
tableTypeWiseSummaryStatistics
```

```
## $`Type A Store Summary Statistics`
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
##    37.16 149.30  196.80  196.90  246.30  301.40
##
## $`Type B Store Summary Statistics`
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
##    45.48   77.79  108.10  117.70  144.30  271.60
##
## $`Type C Store Summary Statistics`
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
##    43.29   57.05   68.46   67.58   78.22   90.57
```

```
## Removing table - not needed for further calculations
rm( tabledTypeWiseSummaryStatistics )
```

From the graph, we can notice a Type B Outlier. A value is considered an outlier when it's more than 3 Standard Deviations from the mean.

From this we can hypothesize that the `Type` of store could be an important predictor of `Weekly_Sales`.

### 3.3.3 Total Sales Per Week - Time Series

We discuss here the effect holidays have on Total Sales of 45 Stores.

```
## Running tapply with sum to find the total sales per week
totalSalesPerWeek <-
  tapply(
    trainStoresFeaturesMerge$Weekly_Sales ,
    trainStoresFeaturesMerge$Date ,
    FUN = sum )
## Converting table to Data Frame
totalSalesPerWeekDataFrame <- as.data.frame( totalSalesPerWeek )
## Converting date from String-Factor to Date Type
totalSalesPerWeekDataFrame$Date <-
  as.Date( rownames(totalSalesPerWeekDataFrame) ) )
## Renaming the Column to "TotalSales"
colnames(totalSalesPerWeekDataFrame)[1] <- "TotalSales"
## Calculating the Total sales in Millions
totalSalesPerWeekDataFrame$TotalSalesInMillion =
  totalSalesPerWeekDataFrame$TotalSales/1000000
```

```
## function to handle lag
## Since the in-built function in R to handle LAG is not working
## x - vector that needs to be lagged
## k - is the no of lags that need to be returned as a vector
## - may be positive or negative
## - it returns 0 instead of NA for the missing values
## returns a vector contained the lagged data with padded 0s for missing data
lagpad <- function(x, k) {
  if( k > 0 ) {
    # It should actually be NA in the rep function
    c(rep(0, k), x)[1 : length(x)]
  } else {
    # It should actually be NA in the rep function
    c(x[ (abs(k)+1) : length(x)] , rep(0, abs(k) ) )
  }
}
```

```
## Getting the holiday List
## Extracting the Holiday List
holidayDateTable <-
  table(trainStoresFeaturesMerge$Date , trainStoresFeaturesMerge$IsHoliday)
## Converting from Table to Data Frame
holidayDateTableDataFrame <- as.data.frame( holidayDateTable)
```

```

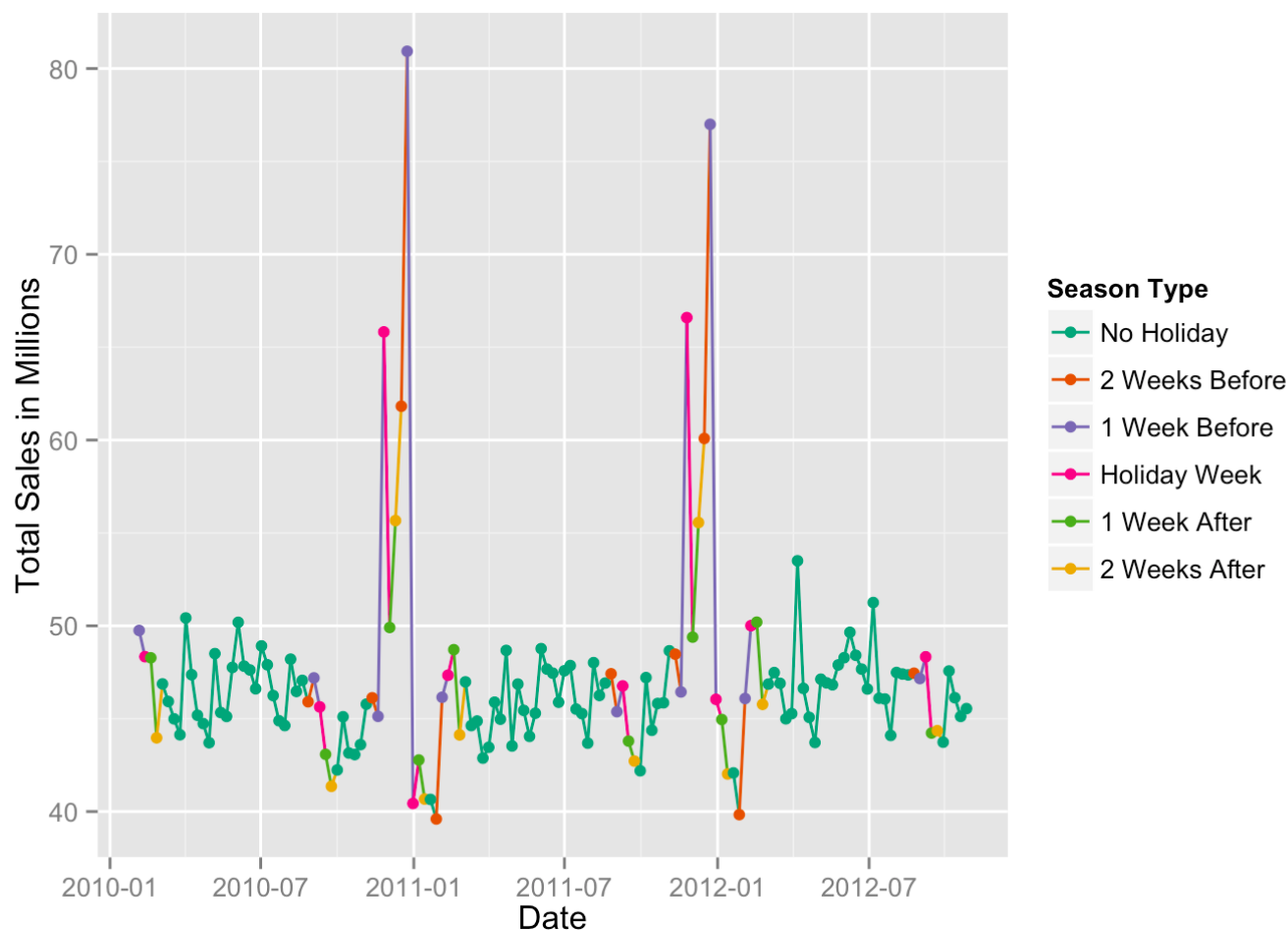
## Extracting the Holidays
holidayDataTableDataFrame <-
  subset( holidayDataTableDataFrame, holidayDataTableDataFrame$Var2==T)
## Marking the Holidays in the dataset
holidayDataTableDataFrame$IsHoliday <-
  ifelse( holidayDataTableDataFrame$Freq >0 , 1 , 0 )
## Converting Date from String-Factor to Date Type
holidayDataTableDataFrame$Date <- as.Date( holidayDataTableDataFrame$Var1 )
# creating the holiday season - before and after the holiday week - 2 weeks
holidayDataTableDataFrame$Week1BeforeHoliday <-
  lagpad(holidayDataTableDataFrame$IsHoliday , -1)
holidayDataTableDataFrame$Week2BeforeHoliday <-
  lagpad(holidayDataTableDataFrame$IsHoliday , -2)
holidayDataTableDataFrame$Week1AfterHoliday <-
  lagpad(holidayDataTableDataFrame$IsHoliday , 1)
holidayDataTableDataFrame$Week2AfterHoliday <-
  lagpad(holidayDataTableDataFrame$IsHoliday , 2)
## Creating an ordered Holiday Season Type
holidayDataTableDataFrame$HolidaySeasonType =
  ifelse( holidayDataTableDataFrame$Week1BeforeHoliday == 1 ,
    "1 Week Before" ,
    ifelse(
      holidayDataTableDataFrame$Week2BeforeHoliday == 1 ,
        "2 Weeks Before" ,
        ifelse(
          holidayDataTableDataFrame$Week1AfterHoliday == 1 ,
            "1 Week After" ,
            ifelse(
              holidayDataTableDataFrame$Week2AfterHoliday == 1 ,
                "2 Weeks After" ,
                ifelse(
                  holidayDataTableDataFrame$IsHoliday == 1 ,
                    "Holiday Week" , "No Holiday" )))))
holidayDataTableDataFrame$HolidaySeasonType = factor(
  holidayDataTableDataFrame$HolidaySeasonType ,
  ordered=TRUE ,
  levels=c(
    "No Holiday" ,
    "2 Weeks Before" ,
    "1 Week Before" ,
    "Holiday Week" ,
    "1 Week After" ,
    "2 Weeks After"
  )
)

## Creating a variable to mark if it is a Holiday season -
## that is 2 weeks before + holiday week + 2 weeks after
holidayDataTableDataFrame$IsHolidaySeason <-
  holidayDataTableDataFrame$IsHoliday +
  holidayDataTableDataFrame$Week1BeforeHoliday * .6 +
  holidayDataTableDataFrame$Week2BeforeHoliday *.2 +
  holidayDataTableDataFrame$Week1AfterHoliday * .6 +
  holidayDataTableDataFrame$Week2AfterHoliday *.2

```

```
## Mering the sales per week and holiday list
totalSalesPerWeekDataFrame <-
  merge( totalSalesPerWeekDataFrame, holidayDateTableDataFrame , by=2)
```

```
# Plotting Sales Per Week
ggplot( totalSalesPerWeekDataFrame ,
        aes(x=Date , y=TotalSalesInMillion , color = HolidaySeasonType ) ) +
  geom_line( aes(group=1) ) +
  geom_point(size = 2) +
  scale_y_continuous(name="Total Sales in Millions" ) +
  scale_color_brewer(palette="Dark2" , name = "Season Type")
```



- Sales go up a week before the holiday week
- It is a repeating pattern - towards the end of the year, there is more sales - perhaps because of

Thanksgiving and Christmas.

- Average Total Sales per week is between 45-50 million Dollars (for 45 stores in the dataset)
- Perhaps the Data given incorrectly marks the Holiday Week for Christmas - it is marked as the dates 27, 28, 30 & 31 across different years. But looking at the peak, the highest sales for Christmas is the week before - which is perhaps the more accurate holiday week
- To be able to make better predictors, perhaps it would help to actually create separate factors for each of the holidays
- It will be interesting to study the period between Aug 27, 2010 to Feb 25, 2011. This section has all the holidays - starting with Labor Day, Thanksgiving, Christmas & Super Bowl. Since the dataset will be smaller, we will perhaps understand the data better.

```
#####
# Creating separate holiday dummy variables
holidayDataTableDataFrame$IsHolidayDefined <-
  ifelse(
    holidayDataTableDataFrame$IsHoliday == 1 &
      month( holidayDataTableDataFrame$Date ) == 2 ,
      2, ifelse(
        holidayDataTableDataFrame$IsHoliday == 1 &
          month( holidayDataTableDataFrame$Date ) == 9 ,
          9 , ifelse(
            holidayDataTableDataFrame$IsHoliday == 1 &
              month( holidayDataTableDataFrame$Date ) == 11 ,
              11 , ifelse(
                holidayDataTableDataFrame$IsHoliday == 1 &
                  month(
                    holidayDataTableDataFrame$Date ) == 12 ,
                    12 , 0 ) ) ) )
# creating the holiday season lag - before and after the holiday week - 2 weeks
holidayDataTableDataFrame$Week1BeforeHoliday <-
  lagpad(holidayDataTableDataFrame$IsHolidayDefined , -1)
holidayDataTableDataFrame$Week2BeforeHoliday <-
  lagpad(holidayDataTableDataFrame$IsHolidayDefined , -2)
holidayDataTableDataFrame$Week1AfterHoliday <-
  lagpad(holidayDataTableDataFrame$IsHolidayDefined , 1)
holidayDataTableDataFrame$Week2AfterHoliday <-
  lagpad(holidayDataTableDataFrame$IsHolidayDefined , 2)

## Creating an ordered Holiday Season Type for Super Bowl
holidayDataTableDataFrame$HolidaySeasonType <-
  ifelse(
    holidayDataTableDataFrame$Week1BeforeHoliday == 2 ,
    "1 Week Before Super Bowl" ,
    ifelse(
      holidayDataTableDataFrame$Week2BeforeHoliday == 2 ,
      "2 Weeks Before Super Bowl" ,
      ifelse(
        holidayDataTableDataFrame$Week1AfterHoliday == 2 ,
        "1 Week After Super Bowl" ,
        ifelse(
          holidayDataTableDataFrame$Week2AfterHoliday == 2 ,
          "2 Weeks After Super Bowl" ,
```



```

        ifelse(
          holidayDateTableDataFrame$IsHolidayDefined
            == 2 ,
          "Super Bowl" , "No Holiday" )))))

## Creating an ordered Holiday Season Type for Labor Day
holidayDateTableDataFrame$HolidaySeasonType <-
ifelse(
  holidayDateTableDataFrame$Week1BeforeHoliday == 9 ,
  "1 Week Before Labor Day" ,
  ifelse(
    holidayDateTableDataFrame$Week2BeforeHoliday == 9 ,
    "2 Weeks Before Labor Day" ,
    ifelse(
      holidayDateTableDataFrame$Week1AfterHoliday == 9 ,
      "1 Week After Labor Day" ,
      ifelse(
        holidayDateTableDataFrame$Week2AfterHoliday == 9 ,
        "2 Weeks After Labor Day" ,
        ifelse(
          holidayDateTableDataFrame$IsHolidayDefined
            == 9 , "Labor Day" ,
          holidayDateTableDataFrame$HolidaySeasonType
            ) ) ) ) )

## Creating an ordered Holiday Season Type for Thanksgiving
holidayDateTableDataFrame$HolidaySeasonType <-
ifelse(
  holidayDateTableDataFrame$Week1BeforeHoliday == 11 ,
  "1 Week Before Thanksgiving" ,
  ifelse(
    holidayDateTableDataFrame$Week2BeforeHoliday == 11 ,
    "2 Weeks Before Thanksgiving" ,
    ifelse(
      holidayDateTableDataFrame$Week1AfterHoliday == 11 ,
      "1 Week After Thanksgiving" ,
      ifelse(
        holidayDateTableDataFrame$Week2AfterHoliday == 11 ,
        "2 Weeks After Thanksgiving" ,
        ifelse(
          holidayDateTableDataFrame$IsHolidayDefined
            == 11 , "Thanksgiving" ,
          holidayDateTableDataFrame$HolidaySeasonType
            ) ) ) ) )

## Creating an ordered Holiday Season Type for Christmas
holidayDateTableDataFrame$HolidaySeasonType <-
ifelse(
  holidayDateTableDataFrame$Week1BeforeHoliday == 12 ,
  "1 Week Before Christmas" ,
  ifelse(
    holidayDateTableDataFrame$Week2BeforeHoliday == 12 ,
    "2 Weeks Before Christmas" ,
    ifelse(

```

```

holidayDataTableDataFrame$Week1AfterHoliday == 12 ,
  "1 Week After Christmas" ,
  ifelse(
    holidayDataTableDataFrame$Week2AfterHoliday == 12 ,
    "2 Weeks After Christmas" ,
    ifelse(
      holidayDataTableDataFrame$IsHolidayDefined
      == 12 , "Christmas" ,
      holidayDataTableDataFrame$HolidaySeasonType
    ) ) ) ) )

holidayDataTableDataFrame$HolidaySeasonType = factor(
  holidayDataTableDataFrame$HolidaySeasonType ,
  ordered=TRUE ,
  levels=c(
    "No Holiday" ,
    "2 Weeks Before Super Bowl" ,
    "1 Week Before Super Bowl" ,
    "Super Bowl" ,
    "1 Week After Super Bowl" ,
    "2 Weeks After Super Bowl" ,
    "2 Weeks Before Labor Day" ,
    "1 Week Before Labor Day" ,
    "Labor Day" ,
    "1 Week After Labor Day" ,
    "2 Weeks After Labor Day" ,
    "2 Weeks Before Thanksgiving" ,
    "1 Week Before Thanksgiving" ,
    "Thanksgiving" ,
    "1 Week After Thanksgiving" ,
    "2 Weeks After Thanksgiving" ,
    "2 Weeks Before Christmas" ,
    "1 Week Before Christmas" ,
    "Christmas" ,
    "1 Week After Christmas" ,
    "2 Weeks After Christmas"
  )
)

```

```

## Mering the sales per week and holiday list
totalSalesPerWeekDataFrame$HolidaySeasonType <-
  holidayDataTableDataFrame$HolidaySeasonType

```

```

## Subsetting only the holidays
totalSalesPerWeekDataFrameDuringHolidays <-
  subset( totalSalesPerWeekDataFrame ,
    totalSalesPerWeekDataFrame$Date >= '2010-08-27' &
    totalSalesPerWeekDataFrame$Date <= '2011-02-25' )

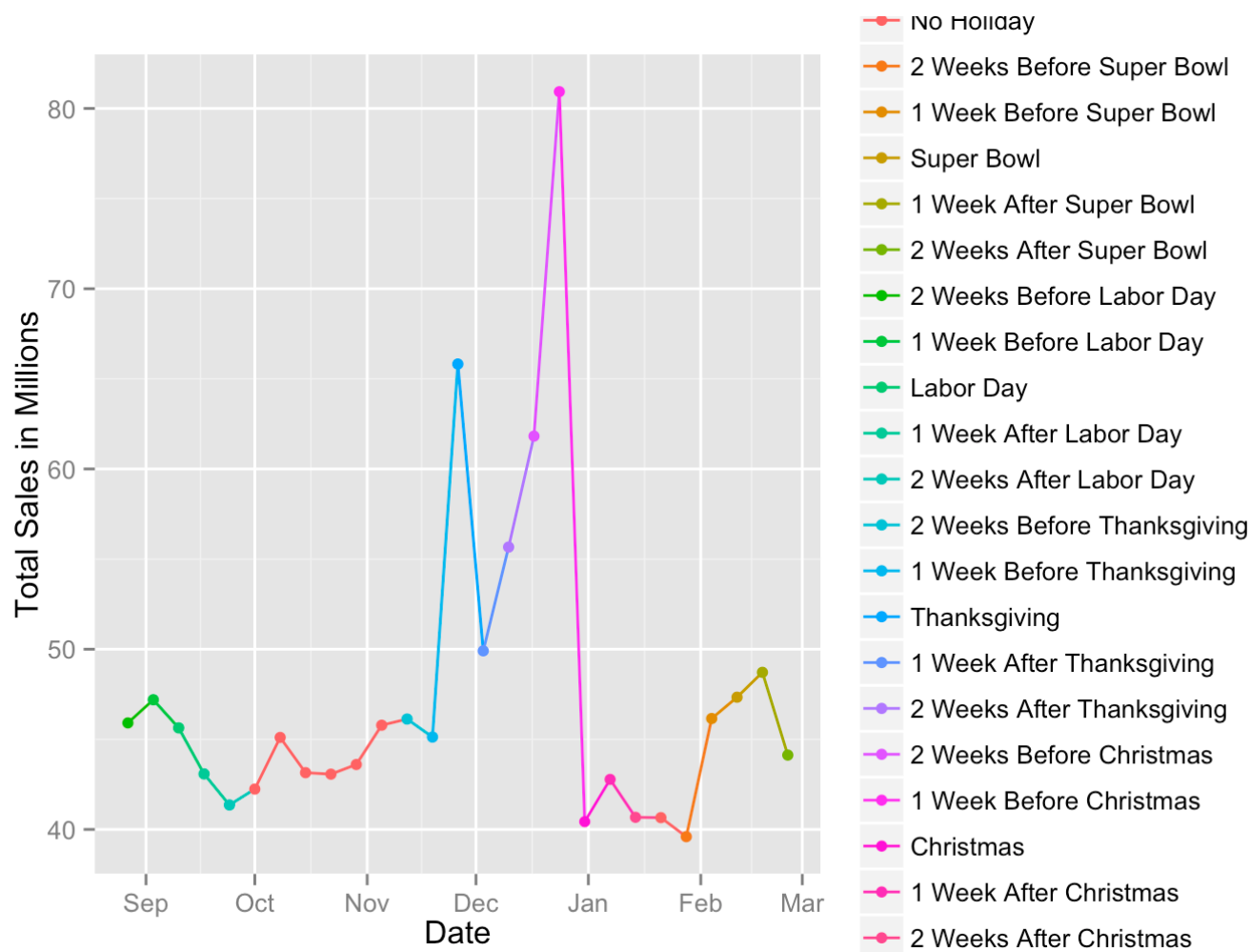
```

```

## Plotting the subset of totalSalesPerWeekDataFrame
ggplot( totalSalesPerWeekDataFrameDuringHolidays ,
  aes(x=Date , y=TotalSalesInMillion , color = HolidaySeasonType ) ) +

```

```
geom_line( aes(group=1) ) +
geom_point(size = 2) +
scale_y_continuous(name="Total Sales in Millions" )
```



```
## removing the following columns because it may cause
## multi-collinearity issues once merged with the main data and
## building a model with that
holidayDataTableDataFrame$Week1BeforeHoliday = NULL
holidayDataTableDataFrame$Week2BeforeHoliday = NULL
holidayDataTableDataFrame$Week1AfterHoliday = NULL
holidayDataTableDataFrame$Week2AfterHoliday = NULL
holidayDataTableDataFrame$IsHoliday = NULL
holidayDataTableDataFrame$IsHolidayDefined = NULL
```

### 3.3.4 Store-Department-wise Sales per Week - Time Series

To see a representation of the granularity of the data, we would like to plot all the data points of Weekly Sales vs Time (Week)

```
## plotting all the Weekly Sales figures - colored by Dept
ggplot(trainStoresFeaturesMerge ,
       aes(x=Date , y = Weekly_Sales , color = Dept ) ) +
geom_point() +
scale_y_continuous(name="Weekly Sales" )
```



- Mathematically, the hypotheses are expressed below:

- ## 5. Stage 3: Linear Regression: Predicting Weekly\_Sales