# Predicting Weeekly Sales from Store Features
## Data Science: Capstone Project (HarvardX PH125.9x)

Aditya Prasad Dash

## Introduction

Understanding the structure and patterns in the data, developing data structures for effectively storing and visualizing it, and using insights from it to predict the target variable of interest from values of other variables is the realm of data science. In this project, I have used the Walmart dataset form kaggle [1] to predict the weekly sales of a store based on the various features of the store and the environment affecting the sales. In the first step, as the Weekly_Sales has numerical values, a linear regression model was trained to predict the Weekly_Sales from the other features. Then, a k-Nearest neighbors model was trained over the dataset for this regression task and the best value of k was found after evaluating the model performance for different k. Then as sometimes it is more useful to understand if a given set of conditions would produce a high or low weekly_sales, the Weekly_Sales column was categorized into high, medium and low weekly sales and a knn model was trained to predict the class from the feature variables.

## Methods

### Reading, visualizing and preprocessing the data

The first step is to load the required libraries for the analysis.

```r
if(!require(tidyverse)) install.packages("tidyverse")
library(tidyverse)
if(!require(tidyr)) install.packages("tidyr")
library(tidyr)
if(!require(caret)) install.packages("caret")
library(caret)
if(!require(stringr)) install.packages("stringr")
library(stringr)
if(!require(ggplot2)) install.packages("ggplot2")
library(ggplot2)
if(!require(lubridate)) install.packages("lubridate")
library(lubridate)
if(!require(corrplot)) install.packages("corrplot")
library(corrplot)
#library(neuralnet)
```

I downloaded the dataset into my computer, therefore, I used the read.csv function in R to read the Walmart dataset and store it into a dataframe called walmart_dataset. The function str(walmart_dataset) displays information about the number of rows and columns of the dataset, the data type in each column and the first few entries of each column.

```
walmart_dataset <- read.csv("/Users/aditya/Documents/Coursework/Online_Courses/Machine_Learning_and_Deep
str(walmart_dataset)
```

```
## 'data.frame':    6435 obs. of  8 variables:
##  $ Store       : int  1 1 1 1 1 1 1 1 1 1 ...
##  $ Date        : chr  "05-02-2010" "12-02-2010" "19-02-2010" "26-02-2010" ...
##  $ Weekly_Sales: num  1643691 1641957 1611968 1409728 1554807 ...
##  $ Holiday_Flag: int  0 1 0 0 0 0 0 0 0 0 ...
##  $ 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 ...
##  $ CPI         : num  211 211 211 211 211 ...
##  $ Unemployment: num  8.11 8.11 8.11 8.11 8.11 ...
```

We notice that the dataset contains 6435 rows and 8 columns. As we want to predict Weekly_Sales from the other variables, I will refer to Weekly_Sales as the target variable and all other variables as feature variables. The data type of "Store" and "Holiday_Flag" is integer, that for "Date" is chr and for the other variables, and that for "Weekly_Sales","Holiday_Flag","Temperature","Fuel_Price","CPI" and "Unemployment" is num.

The function head(dataset,n) displays the first n rows of the dataset, we can use it to visualize walmart_dataset shows the entries of the first 5 rows.

```
head(walmart_dataset,5)
```

```
##   Store       Date Weekly_Sales Holiday_Flag Temperature Fuel_Price      CPI
## 1     1 05-02-2010      1643691            0       42.31      2.572 211.0964
## 2     1 12-02-2010      1641957            1       38.51      2.548 211.2422
## 3     1 19-02-2010      1611968            0       39.93      2.514 211.2891
## 4     1 26-02-2010      1409728            0       46.63      2.561 211.3196
## 5     1 05-03-2010      1554807            0       46.50      2.625 211.3501
##   Unemployment
## 1        8.106
## 2        8.106
## 3        8.106
## 4        8.106
## 5        8.106
```

We first omit all the rows with nan with the function na.omit(dataset) function. Then we check the range of our target variable Weekly_Sales and as both minimum and maximum Weekly sales are positive, it seems reasonable. Next we look at the store column which has integer entries. To find the number of stores included in the dataset we can use the unique() function with argument walmart_dataset$Store to find the unique store numbers.

```
walmart_dataset<-na.omit(walmart_dataset)
range(walmart_dataset$Weekly_Sales)
```
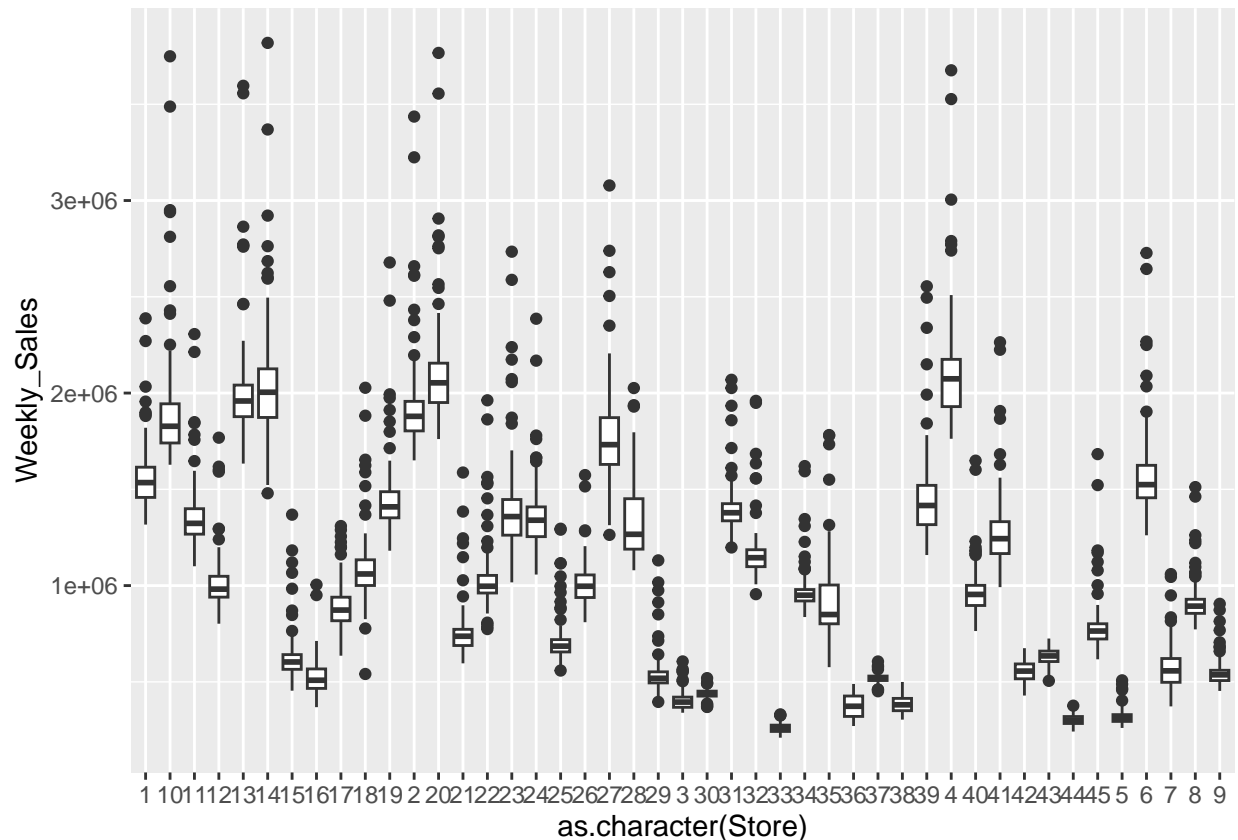
```
## [1]  209986.2 3818686.5
```

```
unique(walmart_dataset$Store)
```

```
##  [1]  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
## [26] 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45
```
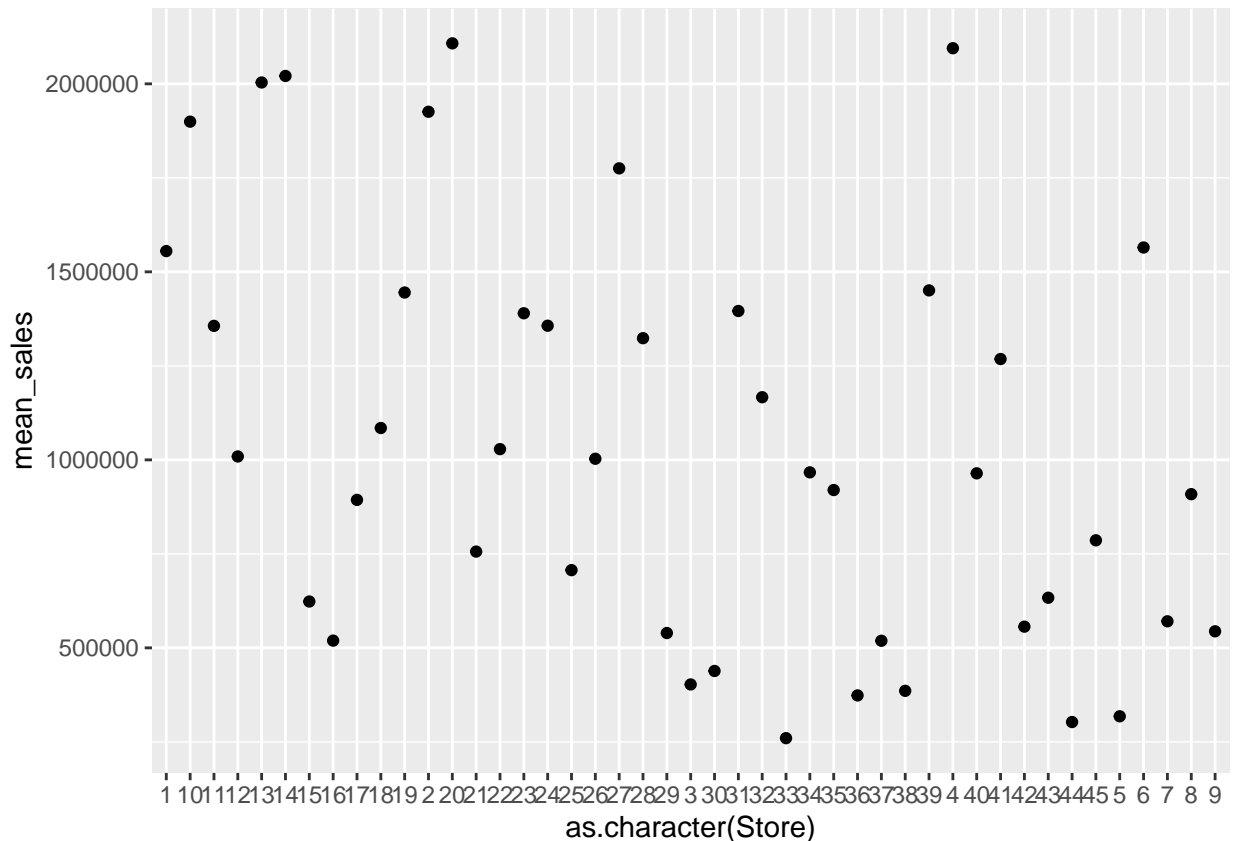
We can notice that there are 45 stores in the dataset with the store numbers ranging from 1 to 45. Some stores can have more sales than others due to location, total volume of products and other factors. Therefore, its useful to visualize the distribution of Weekly sales for different stores. For that, we can use a boxplot which displays the distribution of a numerical variable for each instance of a categorical variable.

```
walmart_dataset%>% ggplot(aes(as.character(Store), Weekly_Sales)) + geom_boxplot()
```



We see that for each store the weekly sales are distributed over a wide range. Moreover, the average (mean) weekly sales for a given store, shown by the central line in the boxplot is different accross different stores. To visualize the mean sales better, we can group the dataset by store numbers and then extract the mean sales for each store and store it in a tibble called mean_sales_store. Next, we can use ggplot and the geom_point function to make a scatterplot of mean scales as a function of the store number.

```
mean_sales_store<-walmart_dataset%>%group_by(Store)%>%summarise(mean_sales=mean(Weekly_Sales))
mean_sales_store %>% ggplot(aes(as.character(Store), mean_sales)) + geom_point()
```

We notice that the weekly sales are spread in each store and they have different averages for different stores, therefore it is useful to use this information in predicting weekly sales. As the store number should not have a heirarchy, we should use one-hot-encoding (https://stackoverflow.com/questions/74706268/how-to-create-dummy-variables-in-r-based-on-multiple-values-within-each-cell-in) to make separate columns for each store in which the entry for a row (observation) is 1 if the Store corresponds to that store and 0 otherwise. To do this we can use the pivot wider function in R to separate the Store numbers into different columns and use length as a funcion to put 1 if that observation is from that Store and 0 (values_fill is set to 0) otherwise and update the walmart_dataset.

```
walmart_dataset<-walmart_dataset%>% pivot_wider(names_from = Store, names_prefix="Store",values_from =
```

Next we try to understand the date column. The format of the date in the dataset is day-month-year as a character. However, to extract meaningful insights from we should separate the month day and year into separate columns. To do this we can the dmy function in lubridate package to convert Date into the appropriate format and then mutate the dataset with the Day, Month and Year columns after extracting it from the Date variable. Displaying the first 5 entries using the head function shows that indeed we have created new columns corresponding to the day, month and year.
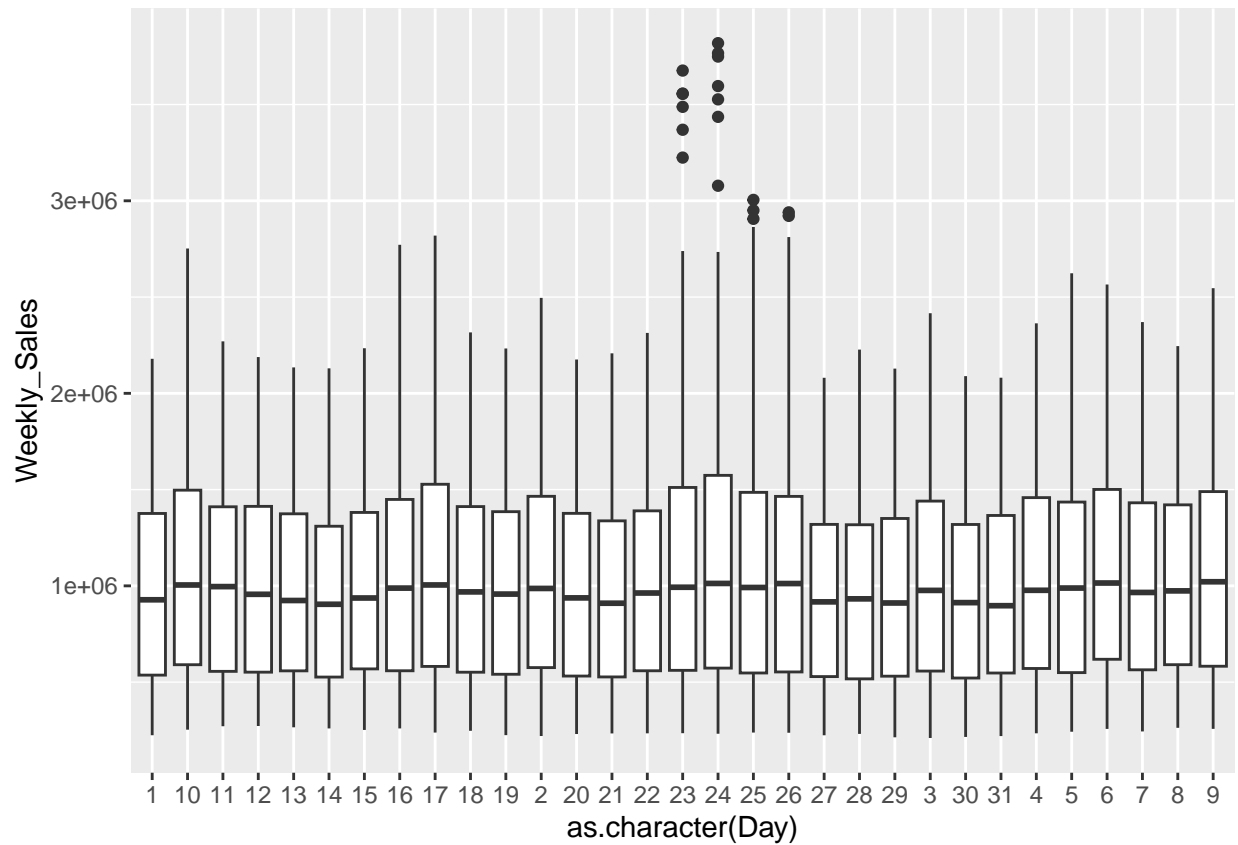
```
walmart_dataset<-walmart_dataset%>%mutate(Day=day(dmy(Date)),Month=month(dmy(Date)), Year=year(dmy(Date
head(walmart_dataset%>%select(c("Day","Month","Year")),5)
```

```
## # A tibble: 5 x 3
##     Day Month  Year
##   <int> <dbl> <dbl>
## 1     5     2  2010
```
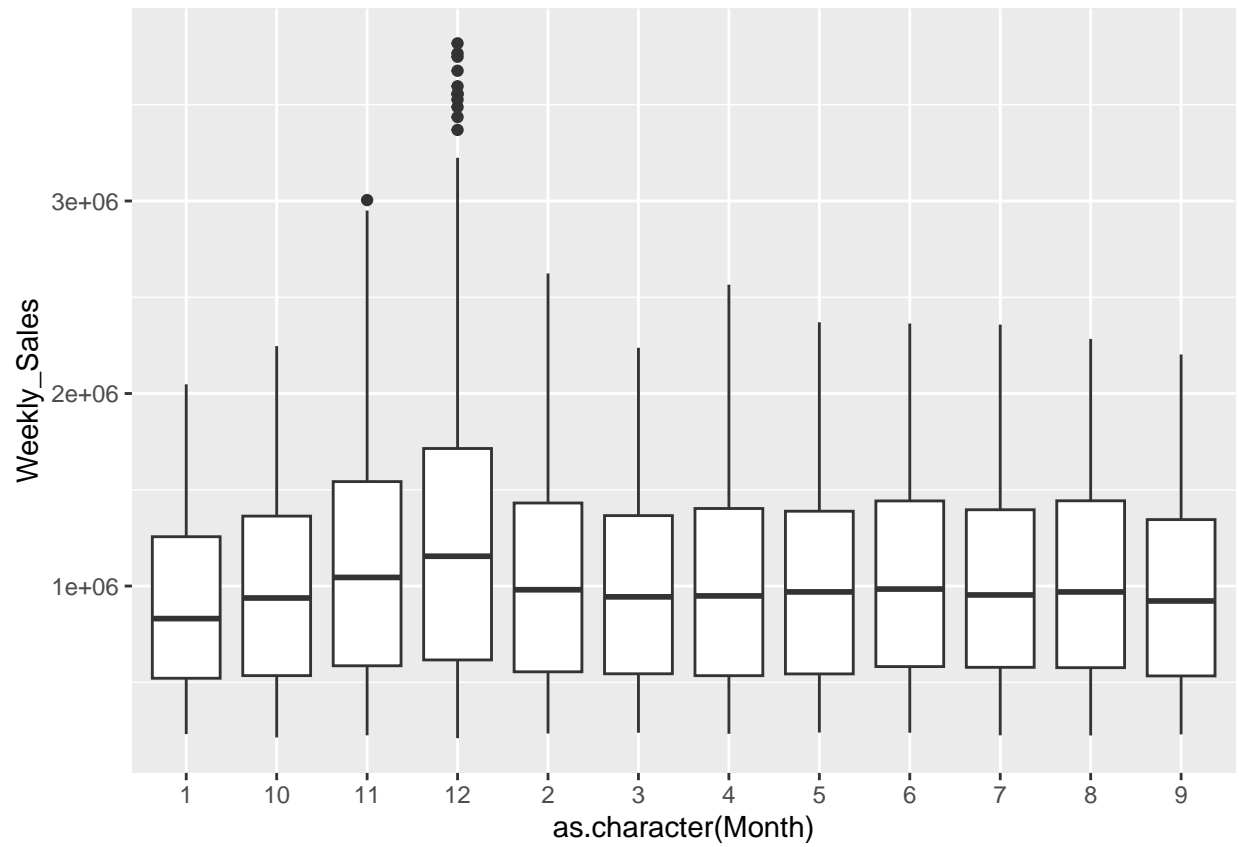
4

```
## 2     12      2   2010
## 3     19      2   2010
## 4     26      2   2010
## 5      5      3   2010
```

Now, we try to visualize the distribution of weekly sales accross day, month and year by creating 3 different boxplots correspondin to each of them.
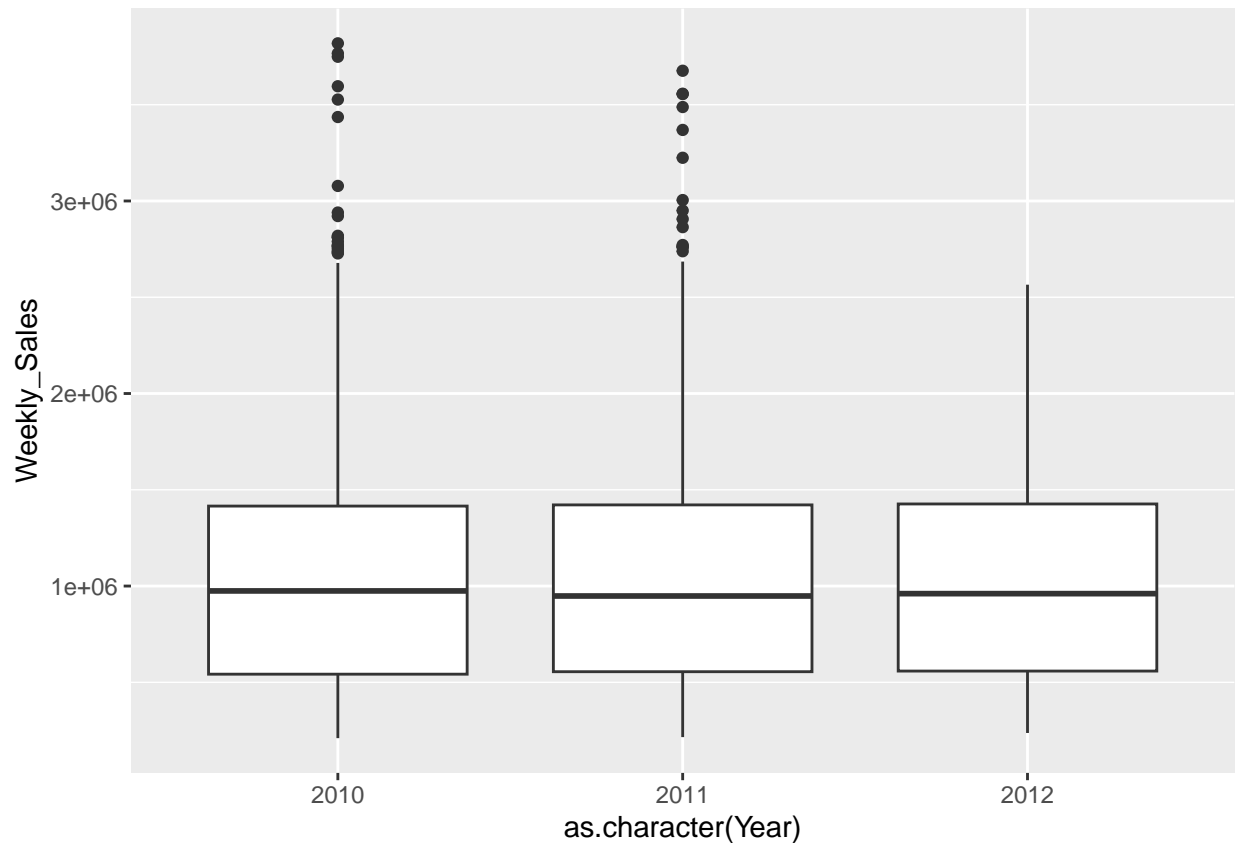
```
walmart_dataset%>% ggplot(aes(as.character(Day), Weekly_Sales)) + geom_boxplot()
```



```
walmart_dataset%>% ggplot(aes(as.character(Month), Weekly_Sales)) + geom_boxplot()
```
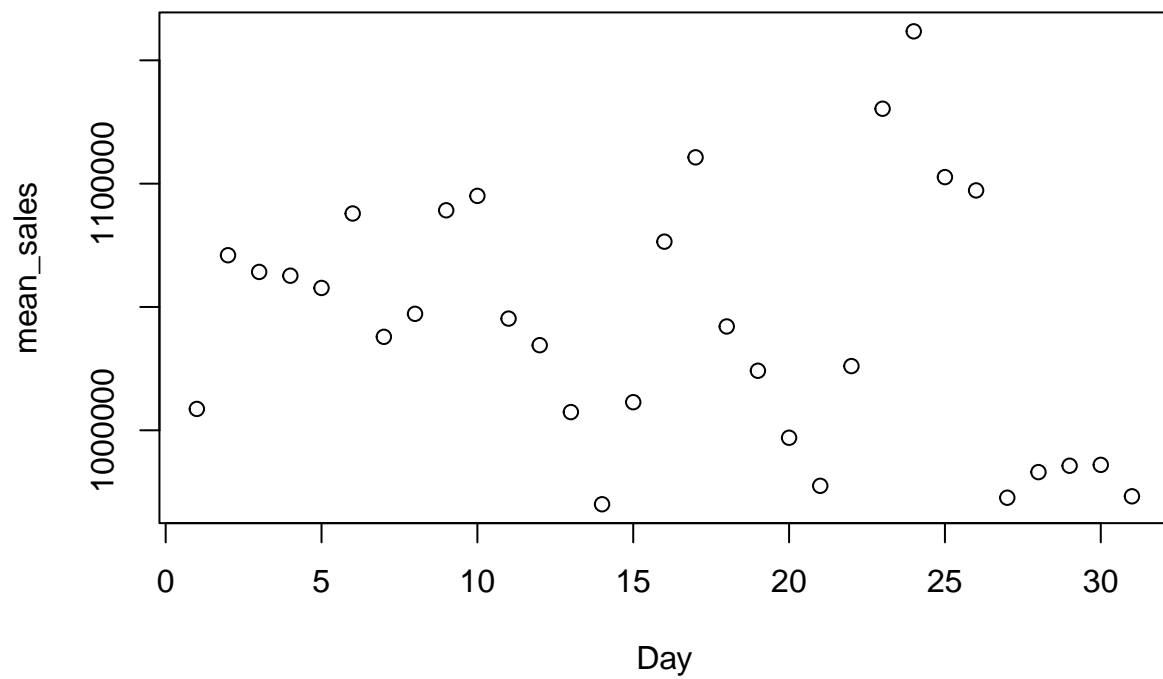
```
walmart_dataset%>% ggplot(aes(as.character(Year), Weekly_Sales)) + geom_boxplot()
```

We see that the distribution of sales shows slight variation with day, month and year. Especially, the average sales is higher around december and 2010 had higher mean weekly sales than 2011 which was higher than that in 2012. To visualize the distribution of mean weekly sales, we can plot the scatter plot of the mean of the weekly sales with day, month and year which agrees with our findings from the boxplots.

```
plot(walmart_dataset%>%group_by(Day)%>%summarise(mean_sales=mean(Weekly_Sales)))
```

```
plot(walmart_dataset%>%group_by(Month)%>%summarise(mean_sales=mean(Weekly_Sales)))
```

```r
plot(walmart_dataset%>%group_by(Year)%>%summarise(mean_sales=mean(Weekly_Sales)))
```
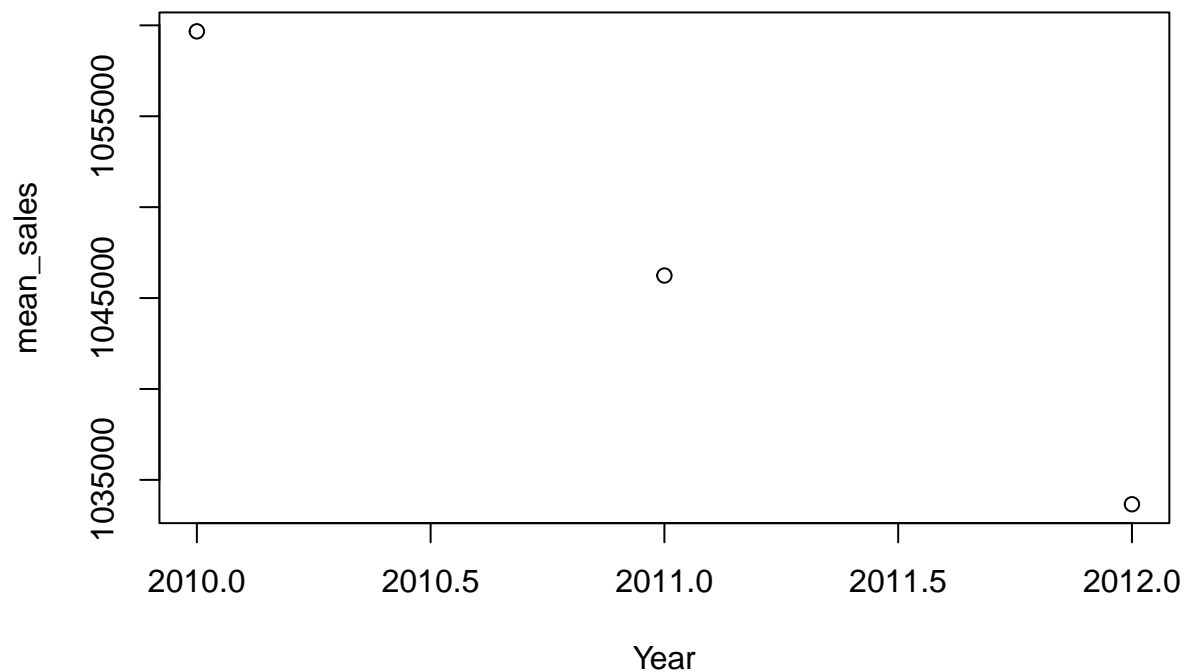
To understand the dependence of the mean weekly sales on Day, Month and year we can group the dataset with the correspoding variable and arrange the categories with descending values of mean sales. We find that on average, the sales were highest on day and minimum on day 14, highest in December and lowest in January and highest in 2010 and loweset in 2012.

```
walmart_dataset%>%group_by(Day)%>%summarise(n=n(),mean_sales=mean(Weekly_Sales))%>%arrange(desc(mean_sa
```

```
## # A tibble: 31 x 3
##        Day     n mean_sales
##      <int> <int>      <dbl>
## 1     24   225   1161746.
## 2     23   225   1130371.
## 3     17   225   1110669.
## 4     25   225   1102654.
## 5     26   225   1097265.
## 6     10   225   1095039.
## 7      9   225   1089198.
## 8      6   225   1087901.
## 9     16   225   1076483.
## 10     2   225   1070976.
## # i 21 more rows
```

```
walmart_dataset%>%group_by(Month)%>%summarise(n=n(),mean_sales=mean(Weekly_Sales))%>%arrange(desc(mean_s
```

```
## # A tibble: 12 x 3
```

```
##      Month     n mean_sales
##      <dbl> <int>      <dbl>
##  1      12   450   1281864.
##  2      11   360   1147266.
##  3       6   585   1064325.
##  4       2   540   1053200.
##  5       8   585   1048017.
##  6       7   630   1031748.
##  7       5   540   1031714.
##  8       4   630   1026762.
##  9       3   585   1013309.
## 10      10   585    999632.
## 11       9   585    989335.
## 12       1   360    923885.
```

```r
walmart_dataset%>%group_by(Year)%>%summarise(n=n(),mean_sales=mean(Weekly_Sales))%>%arrange(desc(mean_sa
```

```
## # A tibble: 3 x 3
##     Year     n mean_sales
##    <dbl> <int>      <dbl>
## 1   2010  2160   1059670.
## 2   2011  2340   1046239.
## 3   2012  1935   1033660.
```

As we see a dependene of the weekly sales on day, month and year, we can separate them into different columns using one-hot encoding similar to that done for Stores. However, as we have 31 days, including each day as a feature can overcomplicate the model, therefore we only consider the effect of Month and Year. We can use pivot wider to create separate columns for each month and Year which takes the value 1 if the observation corresponds to that month or year and 0 otherwise.

```r
walmart_dataset<-walmart_dataset%>% pivot_wider(names_from = Month, names_prefix="Month",values_from = 
walmart_dataset<-walmart_dataset%>% pivot_wider(names_from = Year, names_prefix="Year",values_from = Yea
colnames(walmart_dataset)
```

```
##  [1] "Date"         "Weekly_Sales" "Holiday_Flag" "Temperature"  "Fuel_Price"
##  [6] "CPI"          "Unemployment" "Store1"       "Store2"       "Store3"
## [11] "Store4"       "Store5"       "Store6"       "Store7"       "Store8"
## [16] "Store9"       "Store10"      "Store11"      "Store12"      "Store13"
## [21] "Store14"      "Store15"      "Store16"      "Store17"      "Store18"
## [26] "Store19"      "Store20"      "Store21"      "Store22"      "Store23"
## [31] "Store24"      "Store25"      "Store26"      "Store27"      "Store28"
## [36] "Store29"      "Store30"      "Store31"      "Store32"      "Store33"
## [41] "Store34"      "Store35"      "Store36"      "Store37"      "Store38"
## [46] "Store39"      "Store40"      "Store41"      "Store42"      "Store43"
## [51] "Store44"      "Store45"      "Day"          "Month2"       "Month3"
## [56] "Month4"       "Month5"       "Month6"       "Month7"       "Month8"
## [61] "Month9"       "Month10"      "Month11"      "Month12"      "Month1"
## [66] "Year2010"     "Year2011"     "Year2012"
```

Also, we remove the Date column from the dataset as its information is already contained in columns corresponding to different Month and Year.

```r
walmart_dataset<-walmart_dataset%>%select(-c("Date"))
walmart_dataset<-walmart_dataset%>%select(-c("Day"))
```

**Scaling the dataset**

As the numerical values of each of the feautre variables can have a different magnitude, giving unequal importance to the training parameters corresponding to them during the training process. Therefore, we should scale them to have similar magnitudes, two common methods of scaling are the standard scaler which centers the values of the column at the column mean and then divides each entry with the standard devaition. However, as the outliers are also included in calculating the mean standard deviation, it can bias the centering of the column. Therefore, we can use minmax scaling (https://www.geeksforgeeks.org/how-to-normalize-and-standardize-data-in-r/) in which we subtract the minimum value of the column from each entry and then divide each entry of the column by the range (max-min) of that column.

```r
#Scaling the dataset
fMinMaxSclaer <- function(x) ( #https://www.geeksforgeeks.org/how-to-normalize-and-standardize-data-in-
  (x - min(x)) / (max(x) - min(x))
)

walmart_dataset[-1]<-as.data.frame(lapply(walmart_dataset[-1],fMinMaxSclaer))
head(walmart_dataset)
```
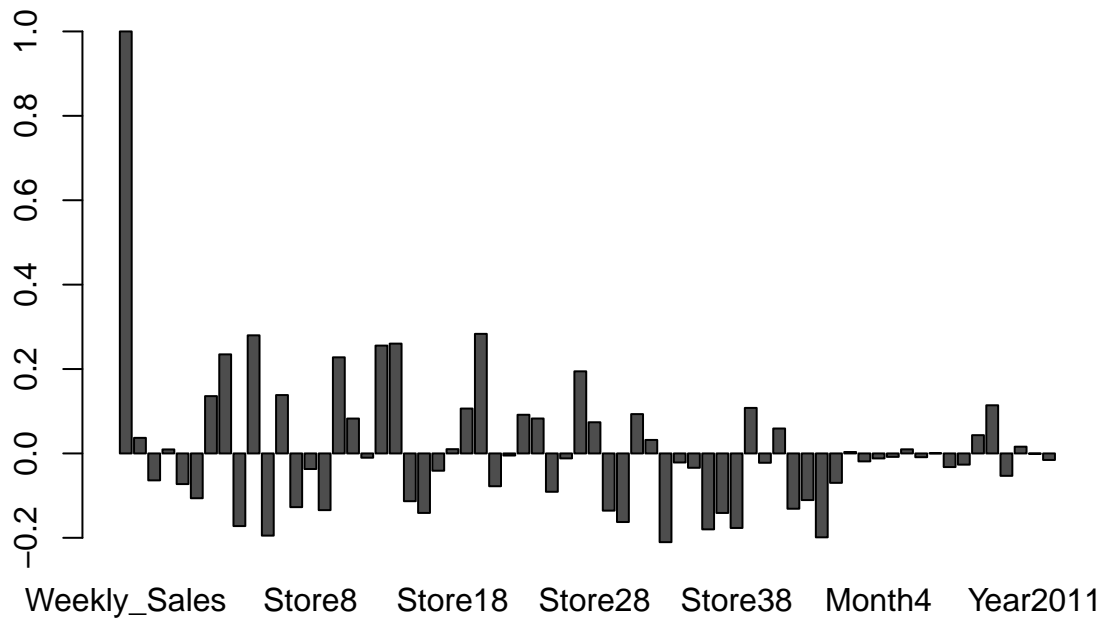
```
## # A tibble: 6 x 66
##   Weekly_Sales Holiday_Flag Temperature Fuel_Price   CPI Unemployment Store1
##          <dbl>        <dbl>       <dbl>      <dbl> <dbl>        <dbl>  <dbl>
## 1    1643691.            0       0.434     0.0501 0.840        0.405      1
## 2    1641957.            1       0.397     0.0381 0.842        0.405      1
## 3    1611968.            0       0.411     0.0210 0.842        0.405      1
## 4    1409728.            0       0.476     0.0446 0.843        0.405      1
## 5    1554807.            0       0.475     0.0767 0.843        0.405      1
## 6    1439542.            0       0.586     0.0977 0.843        0.405      1
## # i 59 more variables: Store2 <dbl>, Store3 <dbl>, Store4 <dbl>, Store5 <dbl>,
## #   Store6 <dbl>, Store7 <dbl>, Store8 <dbl>, Store9 <dbl>, Store10 <dbl>,
## #   Store11 <dbl>, Store12 <dbl>, Store13 <dbl>, Store14 <dbl>, Store15 <dbl>,
## #   Store16 <dbl>, Store17 <dbl>, Store18 <dbl>, Store19 <dbl>, Store20 <dbl>,
## #   Store21 <dbl>, Store22 <dbl>, Store23 <dbl>, Store24 <dbl>, Store25 <dbl>,
## #   Store26 <dbl>, Store27 <dbl>, Store28 <dbl>, Store29 <dbl>, Store30 <dbl>,
## #   Store31 <dbl>, Store32 <dbl>, Store33 <dbl>, Store34 <dbl>, ...
```

**Correlation between features and Weekly_Sales**

To understand how each of the feature variable impacts the weekly sales, we can calculate the correlation between the data of that feature variable with Weeekly_Sales.

```r
cor_data<-cor(walmart_dataset$Weekly_Sales,walmart_dataset%>%select(where(is.numeric)))
barplot(cor_data)
```

We see that different feature variables have different correlations with Weekly sales ranging from -0.2102702149 for Store33 to 0.2833633 for Store20. ## Building the Machine Learning Model

**Machine Learning Model**

**Partitioning the dataset into training and testing data**  A machine learning model is build to get trained on some dataset and be used to predict a quantity of interest in some unknown dataset which may not even exist at the present. Therefore, it is necessary to test the performance of the model by training it over a subset of the dataset(training_set) and test its performance over the other subset of the dataset (testing_set). The performance is evaluated by predicting the output for the target variable over the testing set and compare to the actual values of the target variable in in

First we find out indices of 15% of rows smapled randomly from the walmart_dataset (using Weekly_Sales column) using the createDataPartition function in R. Then we use the test index to create the testing_set which will be used for testing the final model performance after fitting it over the training set which is the subset of the dataset excluding the testing set

```
test_index <- createDataPartition(y = walmart_dataset$Weekly_Sales, times = 1,
                                  p = 0.15, list = FALSE)
testing_set <- walmart_dataset[test_index,]
training_set <- walmart_dataset[-test_index,]
```

We can see the number of observations in each set using the nrow function and see the first 5 rows of the training set using the head(training_set,5) function.

```
nrow(training_set)
```

```
## [1] 5467
```

```
nrow(testing_set)
```

```
## [1] 968
```

```
head(training_set,5)
```

```
## # A tibble: 5 x 66
##   Weekly_Sales Holiday_Flag Temperature Fuel_Price   CPI Unemployment Store1
##          <dbl>        <dbl>       <dbl>      <dbl> <dbl>        <dbl>  <dbl>
## 1     1643691.            0       0.434     0.0501 0.840        0.405      1
## 2     1641957.            1       0.397     0.0381 0.842        0.405      1
## 3     1611968.            0       0.411     0.0210 0.842        0.405      1
## 4     1409728.            0       0.476     0.0446 0.843        0.405      1
## 5     1554807.            0       0.475     0.0767 0.843        0.405      1
## # i 59 more variables: Store2 <dbl>, Store3 <dbl>, Store4 <dbl>, Store5 <dbl>,
## #   Store6 <dbl>, Store7 <dbl>, Store8 <dbl>, Store9 <dbl>, Store10 <dbl>,
## #   Store11 <dbl>, Store12 <dbl>, Store13 <dbl>, Store14 <dbl>, Store15 <dbl>,
## #   Store16 <dbl>, Store17 <dbl>, Store18 <dbl>, Store19 <dbl>, Store20 <dbl>,
## #   Store21 <dbl>, Store22 <dbl>, Store23 <dbl>, Store24 <dbl>, Store25 <dbl>,
## #   Store26 <dbl>, Store27 <dbl>, Store28 <dbl>, Store29 <dbl>, Store30 <dbl>,
## #   Store31 <dbl>, Store32 <dbl>, Store33 <dbl>, Store34 <dbl>, ...
```

**Linear Regression Model**   In linear regression we fit the data to a function to the form $y = \Sigma c_i f_i$

where t is the target variable and f_i are the feature variables. The coefficients $c_i$ are the machine learning parameters which need to be optimized from the training data. I use the caret package to train a linear model to the training_data.

```
fit_lm<-train(Weekly_Sales~.,method="lm",data=training_set)
summary(fit_lm)
```

```
##
## Call:
## lm(formula = .outcome ~ ., data = dat)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -656516  -62093   -3393   44678 1621861
##
## Coefficients: (3 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    560524     107080   5.235 1.72e-07 ***
## Holiday_Flag    40052       8392   4.773 1.87e-06 ***
## Temperature     61677      37735   1.634 0.102221
## Fuel_Price     -19442      27645  -0.703 0.481919
## CPI            426786     170165   2.508 0.012168 *
## Unemployment  -391469      50759  -7.712 1.46e-14 ***
```

14

```
## Store1       594506     54283  10.952  < 2e-16 ***
## Store2       971930     53664  18.111  < 2e-16 ***
## Store3      -591421     60275  -9.812  < 2e-16 ***
## Store4      1452296     97644  14.873  < 2e-16 ***
## Store5      -695885     57121 -12.183  < 2e-16 ***
## Store6       561681     58191   9.652  < 2e-16 ***
## Store7      -240727     23045 -10.446  < 2e-16 ***
## Store8      -121378     61781  -1.965 0.049505 *
## Store9      -489625     62288  -7.861 4.58e-15 ***
## Store10     1338404     98304  13.615  < 2e-16 ***
## Store11      364298     60298   6.042 1.63e-09 ***
## Store12      620081    104295   5.945 2.93e-09 ***
## Store13     1387461     97779  14.190  < 2e-16 ***
## Store14     1222395     18382  66.498  < 2e-16 ***
## Store15       31259     87385   0.358 0.720567
## Store16     -374805     25908 -14.467  < 2e-16 ***
## Store17      278925     97716   2.854 0.004327 **
## Store18      516140     88347   5.842 5.45e-09 ***
## Store19      839968     87458   9.604  < 2e-16 ***
## Store20     1182172     43882  26.940  < 2e-16 ***
## Store21     -204215     53805  -3.795 0.000149 ***
## Store22      421181     81273   5.182 2.27e-07 ***
## Store23      673917     86943   7.751 1.08e-14 ***
## Store24      779325     87678   8.888  < 2e-16 ***
## Store25     -227257     43925  -5.174 2.38e-07 ***
## Store26      405722     87543   4.635 3.66e-06 ***
## Store27     1167423     80852  14.439  < 2e-16 ***
## Store28      936863    104303   8.982  < 2e-16 ***
## Store29       13308     89228   0.149 0.881445
## Store30     -524483     53829  -9.744  < 2e-16 ***
## Store31      432663     53903   8.027 1.22e-15 ***
## Store32      344257     22401  15.368  < 2e-16 ***
## Store33     -299958     98491  -3.046 0.002334 **
## Store34      472239    100316   4.708 2.57e-06 ***
## Store35      340390     81699   4.166 3.14e-05 ***
## Store36     -576999     52106 -11.073  < 2e-16 ***
## Store37     -427638     52127  -8.204 2.88e-16 ***
## Store38        1829    104321   0.018 0.986011
## Store39      515156     52256   9.858  < 2e-16 ***
## Store40      253844     86857   2.923 0.003486 **
## Store41      386602     24409  15.839  < 2e-16 ***
## Store42       -8775     98334  -0.089 0.928899
## Store43     -211933     40182  -5.274 1.38e-07 ***
## Store44     -317793     97778  -3.250 0.001161 **
## Store45          NA        NA      NA       NA
## Month2       119383     11420  10.454  < 2e-16 ***
## Month3        80612     13227   6.094 1.17e-09 ***
## Month4        86608     15454   5.604 2.19e-08 ***
## Month5        87302     17551   4.974 6.76e-07 ***
## Month6       110557     18749   5.897 3.94e-09 ***
## Month7        72358     19894   3.637 0.000278 ***
## Month8        89091     20121   4.428 9.71e-06 ***
## Month9        24197     18672   1.296 0.195059
## Month10       41274     16326   2.528 0.011498 *
```

```
## Month11          196369        15716  12.495  < 2e-16 ***
## Month12          338356        14583  23.202  < 2e-16 ***
## Month1               NA           NA     NA       NA
## Year2010          50043        21814   2.294 0.021822 *
## Year2011          24764        10378   2.386 0.017057 *
## Year2012             NA           NA     NA       NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 143400 on 5404 degrees of freedom
## Multiple R-squared:  0.9366, Adjusted R-squared:  0.9359
## F-statistic:  1288 on 62 and 5404 DF,  p-value: < 2.2e-16
```

The summary function displays the optimized fit coefficients after the training.with their standard errors and t values. The t-statistic is the ratio of the estimate of the parameter to its standard error (assuming the null-hypothesis is model parameter = 0). As the t-statistic folows a student-t distribution, we can estimate the probablity that such a parameter could occur due to statistical fluctuations and it is mentioned in the Pr(>|t|) column. As lower Pr(>|t|) means that the non-zero value of the parameter is less likely to have arisen due to statistical fluctuations.The residual standard error is an estimate of the errors made by asssuming the functional form of Weekly_Sales as a linear function of the feature variables evaluated on the training set.

Now, that we have fitted the linear model on the training set we can use it to predict the Weekly_Sales from the the other variables in the validation set using the predict function

```
y_hat<-predict(fit_lm,testing_set)
```

We can find out the mean absolute error betweent the predicted and actual Weekly_Sales in the validation_set by using the MAE function (ref Introduction to Data Science by Rafael A. Irizarry)

```
MAE<-function(true_values,predicted_values){
  mean(abs(predicted_values - true_values))
}
```

```
mae_lm<-MAE(testing_set$Weekly_Sales,y_hat)
```

We see that the mean absolute error (MAE) is 94271.92. Similarly the root mean square deviation can be found by defining the RMSE function

```
RMSE <- function(true_values, predicted_values){
  sqrt(mean((true_values - predicted_values)^2))}
```

and using it to calulate the RMSE as

```
rmse_lm<-RMSE(testing_set$Weekly_Sales,y_hat)
```

which is calulated to be 155826.1. However, to get a sense of the model performance, we should calculate the ratio of the error to the mean of the Weekly_Sales, as scaling the Weekly_Sales would also scale up the errors in the prediction. This is known as the relative errors and can relative_mae and relative_rmse

```
rmse_lm<-RMSE(testing_set$Weekly_Sales,y_hat)
mean_weeklysales_val<-mean(testing_set$Weekly_Sales)
relative_mae_lm<-mae_lm/mean_weeklysales_val
relative_rmse_lm<-rmse_lm/mean_weeklysales_val

print(mae_lm)
```

```
## [1] 81893.94
```

```
print(relative_mae_lm)
```

```
## [1] 0.07844218
```

```
print(rmse_lm)
```

```
## [1] 133537.1
```

```
print(relative_rmse_lm)
```

```
## [1] 0.1279086
```

We find the relative_mae is 0.08956513 and the relative_rmse is 0.1480461 which are both around 10%. Thus, we find that using the linear model and just a handful of information about the store and environment, we can predict the weekly sales to an accuracy of around 10% relative error which is quite amazing.

**K-Nearest Neighbours Model** Usually the data points which are close in the feature space also have the same target value, for example the phones having similar features have similar price. We can utilize this information to make a prediction using the k-nearest neightbors (knn) model which uses knnreg for continuous variables (https://github.com/topepo/caret/blob/master/models/files/knn.R). In this model, for a given featueset, k observations in the training data are found which have the smallest distance to the given feature set. Then the predicted value is given by the average of the target value of these k-data points in the training set. However, we do not know which value of k should provide the best result, therefore I train the knn model using k values form 1-30 in steps of 2 and obtain the RMSE and print the RMSE vs k (#Neighbors) in a line plot.

```
fit_knnreg<-train(Weekly_Sales~.,method="knn",data=training_set,metric="RMSE",tuneGrid = data.frame(k =
summary(fit_knnreg)
```

```
##              Length Class      Mode
## learn         2     -none-     list
## k             1     -none-     numeric
## theDots       0     -none-     list
## xNames       65     -none-     character
## problemType   1     -none-     character
## tuneValue     1     data.frame list
## obsLevels     1     -none-     logical
## param         0     -none-     list
```

```
print(fit_knnreg)
```

```
## k-Nearest Neighbors
##
## 5467 samples
##   65 predictor
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
## Summary of sample sizes: 5467, 5467, 5467, 5467, 5467, 5467, ...
## Resampling results:
##
##   RMSE       Rsquared   MAE
##   388581.4   0.5449151   289449.1
##
## Tuning parameter 'k' was held constant at a value of 13
```

```
fit_knnreg$bestTune
```

```
##    k
## 1 13
```

The k-value giving the lowest RMSE can be found using fit_knn$bestTune and is found to be 15. The trained model contatins the best paramter (k=15) and can be used to predict the Weeekly_Sales in the validation dataset. Comparing to the true values of the Weekly_Sales in the dataset the mae,relative_mae,rmse and relative_rmse can be found.

```
y_hat<-predict(fit_knnreg,testing_set)
mae_knnreg<-MAE(testing_set$Weekly_Sales,y_hat)
relative_mae_knnreg<-mae_knnreg/mean_weeklysales_val
rmse_knnreg<-RMSE(testing_set$Weekly_Sales,y_hat)
relative_rmse_knnreg<-rmse_knnreg/mean_weeklysales_val
```

```
print(mae_knnreg)
```

```
## [1] 275444.1
```

```
print(relative_mae_knnreg)
```

```
## [1] 0.2638343
```

```
print(rmse_knnreg)
```

```
## [1] 353575.3
```

```
print(relative_rmse_knnreg)
```

```
## [1] 0.3386724
```

We find the relative mae is 0.07010366 and the relative rmse is 0.1372029 which is an improvement from the linear model. It suggests that a linear model is not able to capture the effect of the feature variables on weekly sales but the average of 15 nearest neighbors provides a better estimate of the Weekly_Sales in this dataset.

**Logistic Regression of the category of Weekly Sales**  Many a times, it is important to categorize the target variable into different classes ranging from its minimum to maximum value and use the Machine Learning model to output the cateogry into which the target variables should fall given its set of feature values. For example, the store might be interested in making its overall strategy based on weather a given set of conditions would tend to produce high, medium or low weekly_sales and optimize the conditions based on that. Therefore, in this section, I create 3 classes of the weekly sales using its maximum and minimum values and mutate it as the Category_Weekly_Sales column.

```
hist(training_set$Weekly_Sales)
```

### Histogram of training_set$Weekly_Sales



```
range(training_set$Weekly_Sales)
```

```
## [1]  209986.2 3818686.5
```
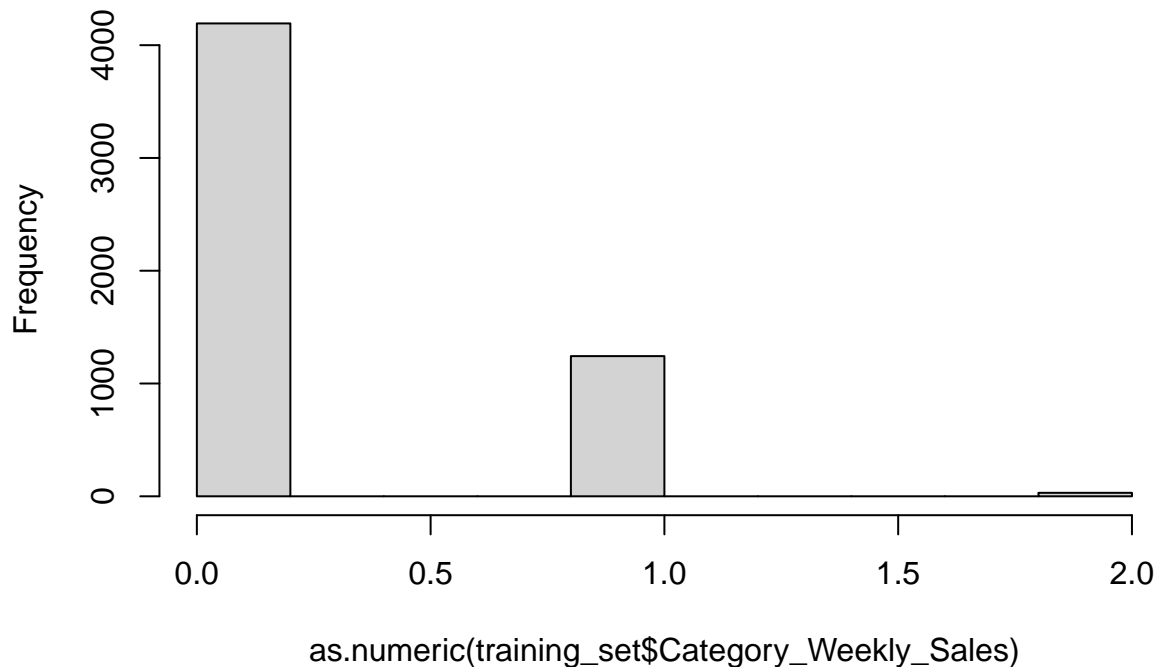
```
min_weekly_sales<-min(training_set$Weekly_Sales)
range_weekly_sales<-max(training_set$Weekly_Sales)-min(training_set$Weekly_Sales)

training_set<-training_set%>%mutate(Category_Weekly_Sales=floor((Weekly_Sales-min_weekly_sales)/range_we
testing_set<-testing_set%>%mutate(Category_Weekly_Sales=floor((Weekly_Sales-min_weekly_sales)/range_weel
```

We can visualize the distribution of Category_Weekly_Sales by plotting its histogram and also by using the group_by function to show the number of occurences and the mean weekly sales in each of these categories. We observe that the categories 0, 1 and 2 contain 4193,1245 and 29 entries while the mean weekly sales in them are 799830, 1833005 and 3087096 respectively.

19

```r
hist(as.numeric(training_set$Category_Weekly_Sales))
```

## Histogram of as.numeric(training_set$Category_Weekly_Sales)



```r
training_set%>%group_by(Category_Weekly_Sales)%>%summarise(n=n(), mean_weekly_sales=mean(Weekly_Sales))
```

```
## # A tibble: 3 x 3
##   Category_Weekly_Sales     n mean_weekly_sales
##                   <dbl> <int>            <dbl>
## 1                     0  4193          798924.
## 2                     1  1243         1835041.
## 3                     2    31         3089619.
```

Now, we can train a machine learning model to predict the category of Weekly_Sales from the other features (except Weekly_Sales as it was used to create the Category_Weekly_Sales so it should not be used as a feature variable). We can use knn, this time to be used for a classification task instead of the previous regression task. For that, we should first convert the Category_Weekly_Sales as factors an d then the knn model over it. As k=1 gave the lowest RMSE for regression, I have used k=1 although the model can be tested for different k values and the k value from the model with the lowest RMSE can be chosen for the final model.

```r
training_set<-training_set%>%mutate(Category_Weekly_Sales=as.factor(Category_Weekly_Sales))
testing_set<-testing_set%>%mutate(Category_Weekly_Sales=as.factor(Category_Weekly_Sales))

fit_knncl<-train(Category_Weekly_Sales~.,method="knn",data=training_set[,-1],tuneGrid = data.frame(k =
summary(fit_knncl)
```

```
##            Length Class      Mode
## learn         2    -none-     list
## k             1    -none-     numeric
## theDots       0    -none-     list
## xNames       65    -none-     character
## problemType   1    -none-     character
## tuneValue     1    data.frame list
## obsLevels     3    -none-     character
## param         0    -none-     list
```

```
fit_knncl$bestTune
```

```
##   k
## 1 1
```

Now, with the fitted knn model with the best tune, we can test its performance over the testing set.

```
y_hat_knncl<-predict(fit_knncl,testing_set[,-1],type="raw")
confusionMatrix(y_hat_knncl, testing_set$Category_Weekly_Sales)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   0   1   2
##          0 707  34   0
##          1  32 186   1
##          2   1   5   2
##
## Overall Statistics
##
##                Accuracy : 0.9246
##                  95% CI : (0.9061, 0.9404)
##     No Information Rate : 0.7645
##     P-Value [Acc > NIR] : <2e-16
##
##                   Kappa : 0.7918
##
##  Mcnemar's Test P-Value : 0.2925
##
## Statistics by Class:
##
##                      Class: 0 Class: 1 Class: 2
## Sensitivity            0.9554   0.8267 0.666667
## Specificity            0.8509   0.9556 0.993782
## Pos Pred Value         0.9541   0.8493 0.250000
## Neg Pred Value         0.8546   0.9479 0.998958
## Prevalence             0.7645   0.2324 0.003099
## Detection Rate         0.7304   0.1921 0.002066
## Detection Prevalence   0.7655   0.2262 0.008264
## Balanced Accuracy      0.9031   0.8911 0.830225
```

We find that the knn model gives an accuracy of 0.8874. The confusion matrix shows the comparision between the reference and prediction values which should be a diagonal matrix in case of perfect prediction. The

non-diagonal entries represent errors during the classification similar to the case when there are 2 classes in which false positive (when the prediction is positive and actual value is negative) and false_negatives (when the prediction is negative and actual value is positive) represent the errors made during the classification. Accurracy represents the ratio of number of instances which were correctly classified with the total number of instances.

## Summary of Results

In this project we have successfully built a machine learning model to predict the Weekly_Sales from the store number and other features in the dataset. In the first model, a linear model was used to regress the value of weekly sales from other features which was trained over the training set. To evaluate the model perfromance two different metrics, the mae (Mean absolute error) and rmse (Root Mean Squared Error) were used which was 83550.32 and 136923 respectively for the linear regression model. To compare this error with the magnitude of the variable we are predicting, the relative_mae and relative_rmse which are the ratio of the mae and rmse to the mean of the Weekly_Sales were calculated and were found to be 0.07956794 and 0.1303967 respectively. Next a k-Nearest Neighbors model with regression mode was trained with different values of k to predict the weekly sales from other features. Although the RMSE showed a local minima at k=15, the lowest RMSE was found to be at k=1 which was chosen for the knn model. The mae, rmse, relative_mae and relative_rmse for the knn model were found to be 84540.2, 165459.1, 0.08051064 and 0.1575726 respectively. Next, the Weekly_Sales column was categorized into high, medium and low weekly sales and a knn model with classification mode was trained over the trainining set and used to predict the category of the Weekly_Sales in the testing set. The performance was visualized using the confusion matrix and accuracy and the accuracy was found to be 89.26%.

## Conclusion

In this project, I have utilized the concepts of machine learning to build a model to predict the weekly sales of a store based on the store number and other conditions (Holiday, Temperature,Fuel_Price,CPI and Unemployment index). The dataset was first preprocessed, scaled and then partitioned into the training and testing sets. A linear model and a knn model with regression mode was then fitted over the training set and evaluated over the testing set. The linear model was found to have a better relatively mean absolute error and relative root mean squared error. The relative mean squared errors were different for different k which shows that hyperparameter tunning is important during training a machine learning model. Then, the Weekly_Sales was columns was categorized into high, medium and low weekly sales and a knn model was fit with classification to predict the cateogry of Weekly_Sales. It was found to have an accuracy of 89.26%. This model could be improved by using L1 and L2 type regularization and by using cross validation and decision trees for classification among other machine learning methods that can also be used and tested. Through this project, I have realized the importance of utilizing data to understand quantities of significance which can be used for optimizing the parameters at hand to achive the best possible value for a target variable and this model can be generalized to predict the weekly sales of other stores by preprocessing the features to the appropriate format.

## References

1. https://www.kaggle.com/datasets/yasserh/walmart-dataset/data

2. Data Science: Capstone, HarvardX PH125.9x provided via the edX platform

3. Introduction to Data Science by Rafael A. Irizarry

4. https://stackoverflow.com/questions/74706268/how-to-create-dummy-variables-in-r-based-on-multiple-values-within-each-cell-in

5. https://www.geeksforgeeks.org/how-to-normalize-and-standardize-data-in-r/

6. https://stackoverflow.com/questions/62679940/geom-bar-of-named-number-vector