**Research Project**
**On**
# Stock Market Analysis
**By**
**Bhakti Sangoi**
**Sagar Chettiyar**

**Point of Contact**
**Bhakti Sangoi**
Email:  sangoi.b@husky.neu.edu
Phone:  +1 (857) 389-5956

**Sagar Chettiyar**
Email: chettiyar.s@husky.neu.edu
Phone: +1 (617) 407-9130

Submitted on: April 22, 2018

**Instructor: Prof.  Kathleen Durant**

## Table of Contents

# 1. Introduction

Stocks:

A stock is a type of ownership in a corporation. It represents the claim on the company's earnings and assets. More you buy stocks of a company, more ownership stake in the company becomes greater.

Stock Market:

Stock Market is a place where bonds, securities and shares/stocks of public listed companies are traded. The main function of the stock market is to process the transaction between buyers and sellers. Buyers/Sellers can track the price change of stocks through the stock market.

One of the many things that people believe is that there is only one stock market as the name suggests but that's not true, because there are multiple other stock markets.

How do Stock Market works?

The stock market is consisting of different exchanges such as New York Stock Exchange(NYSE), National Association of Securities Dealers Automated Quotation System (Nasdaq) and American Stock Exchange (AMEX). Stocks/Shares are entered on one of these exchanges. Buyers and Sellers come together to the market to exchange stocks. Individuals typically buy and sell between one another and there is an auction occurring; i.e. the highest bidding price is matched with the lowest asking price.

NYSE and Nasdaq:

New York Stock Exchange (NYSE) is the world's largest and the oldest stock exchange and is located on Wall Street in New York City. It was founded in 1792. 2800 companies are listed under NYSE with trading 1.46 billion shares each day. Nasdaq is 46 years young. There are 3800 listings with the market capitalization of $11 trillion. NYSE and Nasdaq do the maximum business in the stock market.

Trading starts at 9:30 am to 4:00 pm EST from Monday to Friday for both the exchanges. The only difference between NYSE and Nasdaq is the way stocks are exchanged between buyers and sellers. Nasdaq is Dealer's market whereas NYSE is auction market.

# 2. Aim

Seemingly endless data is available for Stock market. It is difficult to collect a large dataset of stock prices which is structured, cleaned and has high granularity. Analysing stock prices helps an intelligent trader to invest wisely in the stock market. Thereby we provide a dataset which is trying to cover the 6V's of big data: Volume, Volatility, Veracity, Variety, Velocity, and Validity.

## 2.1 Objective

The Objective of this project is:

- **Collect** and clean the high-volume dataset for the year 2013-2016 on both daily and annual basis.
- Once the data is cleaned, create a schema with a variety of data and store it in the database. We are using SQL database to **store** the data.
- After the data is in the structured format, **retrieve** data using SQL queries.
- Perform Data Analysis using R programming

- Additionally, perform Machine Learning Algorithm on highly volatile data
- Create a Web App using R Shiny

## 2.2 Motivation

We both come from an IT domain with experience and a strong interest in Finance industry and Investment Management. Also, we plan to deal in the stock market as we all know that stock market has been on a bit of a tear. Is it a great time to buy stocks? Let's have a look.

# 3. Data Collection and Data Cleaning

We are collecting data from two sources. We followed the extraction with the manipulation phase on both the sources.

**Source 1: S&P Dow Jones Indices**
- S&P Dow Jones Indices is maintained by S&P Global company and we are using this data because it covers almost 80% of American equity market by capitalization.
- This data source proved to be a wise choice as it helped in exhibiting interpretation from raw data and extract the required data for changing and to improve on.

**Source 2: Kaggle NYSE and NASDAQ data**
- This dataset has raw, as-in daily prices of data spans from 2010 to the end of 2016.
- This dataset consists of high level of granularity which is enough to derive some other fundamental indicators.

**Source 3: Yahoo! Finance**
- Data from Yahoo Finance is used to generate Dynamic Time Series analysis. Yahoo Finance provides live and raw data to perform analysis and further explain interpretations.

Step 1: Scraping the S&P Dow Jones Indices website
   i) We used scraping techniques for Dow Jones Index data. Using Import.io data scrapping tool, we collected the necessary data for 500 companies required for this project.

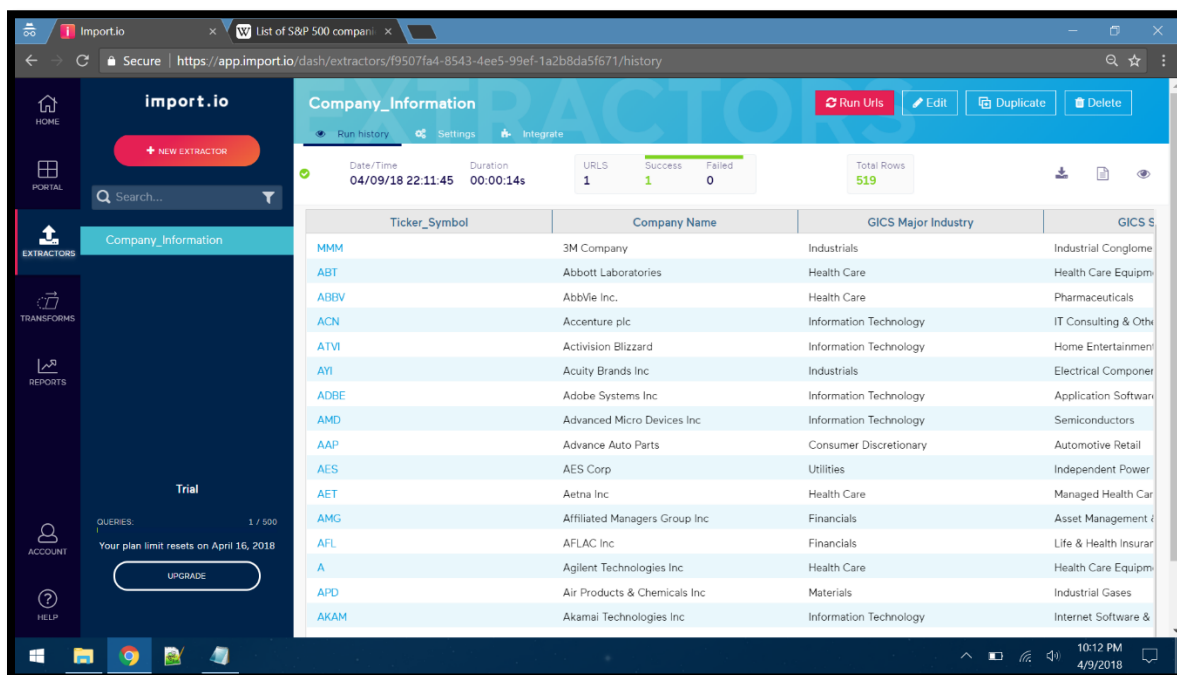ii) Selected the required columns needed for storing in the schema. It helps to reduce the process time for loading the data into the database system.

iii) Saving the data in the .csv format.



iv) The final output of the Scrapping. This dataset has 519 rows and 7 columns.

Step 2: Collecting Stocks data
We have collected the stocks data for the listed 500 companies from source 1 for both daily basis and annual basis. We planned to work with stocks coming from both the exchanges i.e. NYSE and NASDAQ because it will help in performing analysis on a huge dataset. Also, this will help us give an accurate result for forecasting and prediction algorithms.
The daily stocks dataset file contains 851013 rows and annual stocks file contains 1781.

```
#--------------------------DATA COLLECTION---------------------------------
Company_Info<-read_csv("Comapny_Information.csv")
Annual_Balance_Info<-read_csv("Annual_Balance_Sheet.csv")
Daily_Stock_Info<-read_csv("Daily_Stock_Prices.csv")
```

Step 3: Data Cleaning
A lot of time and efforts was put on data cleaning.

i) Splitting the address into City and State
        Splitting the City and State from Headquarters address for better analysis using the cSplit function

```
28 ▾ #------------------------------DATA CLEANING----------------------------------
29
30 ▾ #--------------------------For Company Info--------------------------------
31
32   #Splitting the address into City and State
33   Company_Info<-cSplit(Company_Info, "Address of Headquarters", ",", fixed = FALSE)
34
```

ii) Removing unwanted columns
        Removing all the unwanted columns such as SEC filings.

```
35   #Dropping unwanted columns
36   Company_Info$`SEC filings`<-NULL
37   Company_Info$`Address of Headquarters_3`<-NULL
38
```

iii) Renaming the column names
        Renaming the columns names as per database naming schema for easy understanding.

```
38
39   #Renaming the dataset for better understanding
40   names(Company_Info)<-c("Ticker_Symbol",
41                          "Company_Name",
42                          "GICS_Sector",
43                          "GICS_Major_Industry",
44                          "Date_Added",
45                          "CIK",
46                          "City",
47                          "State")
```

```
56 ▾ #--------------------------For Daily_Stock_Info----------------------------------
57  #Renaming the dataset for better understanding
58  names(Daily_Stock_Info)<-c("Date",
59                             "Ticker_Symbol",
60                             "Open_Price",
61                             "Close_Price",
62                             "Lowest_Price",
63                             "Highest_Price",
64                             "Volume")
```

```
99 ▾ #--------------------------For Annual_Balance_Info----------------------------------
100  #Renaming the dataset for better understanding
101  names(Annual_Balance_Info)<-c("Financial_Year",
102                        "Ticker_Symbol",
103                        "Cash_Cash_Equivalents",
104                        "Total_Current_Assets",
105                        "Total_Current_Liabilities",
106                        "Short_Term_Investments",
107                        "Total_Revenue",
108                        "Treasury_Stock",
109                        "Retained_Earnings",
110                        "Operation_Margin",
111                        "Quick_Ratio")
```

iv) Filling Missing Values in date column by defaulting a date to all the NA's

There were some missing dates of company information being added to the database. We thought of giving a default date value to NA's in the Date_Added field.

```
48
49  #Filling Missing values
50  #Replacing the NA's in date column with a default date:"1999-01-01"
51  Company_Info$Date_Added[is.na(Company_Info$Date_Added)] <- "1999-01-01"
```

v) Following one standard format for Major Sectors

There are rows which have Major Sector as IT and some have Information Technology. Following one format and converting all of them to Information Technology to get an accurate analysis.

```
53  #Replacing IT with Information Technology in GICS Sector
54  Company_Info$`GICS_Sector`[Company_Info$`GICS_Sector`=="IT"] <- "Information Technology"
55
```

vi) Cleaning the date column, formatting it and factorizing it into different columns for analysis purpose.

```
66  #Subsetting date column to keep only dates
67  #Factorizing Dates
68  Daily_Stock_Info$Date<-as.Date(Daily_Stock_Info$Date, "%m/%d/%Y") #Reformatting Date
69  Daily_Stock_Info$Day<-factor(day(as.POSIXlt(Daily_Stock_Info$Date, format="%m/%d/%Y"))) #Adding a Day column representing the
70  Daily_Stock_Info$Month<-factor(month(as.POSIXlt(Daily_Stock_Info$Date, format="%m/%d/%Y"))) #Adding a Month column representi
71  Daily_Stock_Info$Year<-factor(year(as.POSIXlt(Daily_Stock_Info$Date, format="%m/%d/%Y"))) #Sepereting the Year from the date
72  Daily_Stock_Info$Weekday<-factor(wday(as.POSIXlt(Daily_Stock_Info$Date, format="%m/%d/%Y"))) #Adding a weekday column represe
73
```

vii) Creating new data frames volume of stocks traded by month, Year, weekdays and each day of the month respectively.

```
74  #Creating new dataframes for future analyses
75  Volume_YearMonth<-Daily_Stock_Info %>%
76          group_by(Year, Month) %>%
77          summarise(Total = n())
78  #This dataframe contains volume of stocks traded by Month and Year
79
80  Volume_Month<-Daily_Stock_Info %>%
81          group_by(Month) %>%
82          summarise(Total = n())
83  Months<-c("Jan","Feb","Mar","Apr","May","Jun","Jul","Aug","Sep","Oct","Nov","Dec") #Assigning month names instead of month nu
84  Volume_Month$Month<-Months[Volume_M$Month]
85  #This dataframe contains volume of stocks traded by Month
86
87  Volume_Weekday<-Daily_Stock_Info %>%
88          group_by(Weekday) %>%
89          summarise(Total = n())
90  Weekday<-c("Monday","Tuesday","Wednesday","Thursday","Friday")  #Assigning Weekday names instead of month numbers
91  #Note: Saturday and Sunday are not included because the stock market doesnot remain functional on these days
92  Volume_Weekday$Weekday<-Weekday[Volume_W$Weekday]
93  #This dataframe contains volume of stocks traded by weekdays
94
95  Volume_Day<-Daily_Stock_Info %>%
96          group_by(Day) %>%
97          summarise(Total = n())
98  #This dataframe contains volume of stocks traded by each day of the month
```

## 4. Data Storage

Once the data was collected and cleaned, next step was to create the data structure which follows 3rd Normal Form(3NF). 3NF is used in normalizing the database design to reduce the duplication of data and ensuring referential integrity.

We created 3 tables by establishing the connection with the RSQLite package.

Step 1: Creating schema design

1. Company Information
   This table has company information

   | Ticker Symbol (PK) | Company Name | City | State | GICS Sector | GICS Major Industry | CIK | Date Added |
   |---|---|---|---|---|---|---|---|
   | | | | | | | | |

   - o  Ticker_Symbol
     - Ticker_Symbol is the primary key for Company Information table
     - It contains companies Stock's Name
     - Eg: Apple – AAPL, Microsoft – MSFT
   - o  Company Name
     - It shows the Name of the company a stock belongs to.
   - o  City
     - It provides the City in which the Headquarter of the company is located.
   - o  State
     - Contains state in which headquarter of the company is located
   - o  GICS Sector
     - Provides the main domain of a company.
     - Eg: Healthcare Industry, Information Technology
   - o  GICS Major Industry
     - It gives a more detailed information about the GICS Major industry

- Eg: Application Software, Automotive Retail, Health Care Equipment or Pharmaceuticals
  - o Date_Added
    - Date at which the company information was added to the stock market

2. Daily Stock Information
   Step 1: Creating Schema
   This schema has stock information on a daily basis

| Daily Stock ID (PK) | Date | Ticker_Symbol | Open Price | Close Price | Highest Price | Lowest Price | Volume |
|---|---|---|---|---|---|---|---|

  - o Daily_Stock_ID
    - It is a primary key to the Daily Stock Information table
  - o Date
    - The daily date for which prices are displayed.
  - o Ticker_Symbol
    - It is a foreign key from the Company Information table (schema 1)
  - o Open Price
    - The price at which the stock market opens for a day
    - It need not be the same as previous day's closing price
  - o Close Price
    - The price at which the stock market closes for the day
  - o Highest Price
    - The highest price of a stock for one day
  - o Lowest Price
    - The lowest price of a stock for one day
  - o Volume
    - Total number of stocks traded in one day

3. Balance Sheet
   This table has stock information on an annual basis

| Annual Stock ID (PK) | Financial Year | Ticker_ Symbol | Cash and Cash Equivalents | Total Current Assets | Total Current Liabilities | Short Term Investments | Total Revenue | Treasury Stock | Retained Earnings | Operating Margin | Quick Ratio |
|---|---|---|---|---|---|---|---|---|---|---|---|

  - o Annual_Stock_ID
    - It is the primary key for Balance Sheet
  - o Financial Year
    - Since the balance sheet is calculated annually, this specifies the year for which it was calculated
  - o Ticker_Symbol:
    - It is a foreign key from Company Information Table (schema 1)
  - o Cash and Cash Equivalents
    - Value of a company's assets that are cash or can be converted into cash immediately.
    - These include bank accounts, marketable securities, commercial paper, Treasury bills and short-term government bonds with a maturity date of three months or less
  - o Total Current Assets

- The total sum of all the available assets of the company for a particular financial year
- Total Assets and Total Current Assets are different
- Total Current Liabilities
    - The total sum of all the available liabilities of the company for a particular financial year
    - Total liabilities and Total Current liabilities are different
- Short-Term Investments
    - Debt incurred by a company that is due within one year
    - The value of the short-term debt account is very important when determining a company's performance.
    - If the account is larger than the company's cash and cash equivalents, this suggests that the company may be in poor financial health and does not have enough cash to pay off its short-term debts.
- Total Revenue
    - It is the top line or gross income figure from which costs are subtracted to determine net income.
    - Revenue is calculated by multiplying the price at which goods or services are sold by the number of units or amount sold.
- Treasury Stock
    - A portion of shares that a company keeps in its own treasury.
    - Treasury stock may have come from a repurchase or buyback from shareholders, or it may have never been issued to the public in the first place.
    - These shares don't pay dividends, have no voting rights and should not be included in shares outstanding calculations.
- Retained Earnings
    - Retained earnings refer to the percentage of net earnings not paid out as dividends, but retained by the company to be reinvested in its core business, or to pay the debt
    - They can show a positive earnings accumulation or can turn negative and have a deficit if a current period's net loss exceeds the period's beginning retained earnings.
- Operating Margin
    - Operating margin is a measurement of what proportion of a company's revenue is left over after paying for variable costs of production such as wages, raw materials
    - Operating margin gives analysts an idea of how much a company makes on each dollar of sales.
- Quick Ratio
    - The quick ratio is an indicator of a company's short-term liquidity and measures a company's ability to meet its short-term obligations with its most liquid assets.
    - Because we're only concerned with the most liquid assets, the ratio excludes inventories from current assets.
    - Quick assets are current assets that can be converted to cash within 90 days or in the short-term
    - If greater than 1, liquid assets can cover for short-term investments
    - If less than 1, the company may not be able to pay off their debts

**Step 2**: Setting up Connection

Initializing Database connection using SQLite.

```
113 ▾ #------------------------DATA STORAGE----------------------------------
114   #Creating Database Connection
115   db_conn <- dbConnect(SQLite(), dbname="Stock_Market.sqlite")|
116
```

**Step 3**: Creating Tables in the database

```
117   #Table Creation
118 ▾ #------------------------For Company_Info----------------------------
119   dbGetQuery(db_conn, #Database Connector name
120              "create table Company_Info
121              (
122              Ticker_Symbol Text Primary Key,
123              Company_Name Text,
124              GICS_Sector Text,
125              GICS_Major_Industry Text,
126              Date_Added Date,
127              CIK Text,
128              City Text,
129              State Text
130              )") #Create Table Script
131   dbWriteTable(conn = db_conn, #Database Connector name
132                name = "Company_Info", #Table Name
133                value = Company_Info, #Load data from the newly created dataframe
134                append=TRUE,row.names = FALSE,header = TRUE)#Since the CSV contains headers in the columns
135   dbListFields(db_conn, "Company_Info") #Listing the field of the table
136
```

```
137 ▾ #------------------------For Daily_Stock_Info------------------------
138   dbGetQuery(db_conn, #Database Connector name
139              "create table Daily_Stock_Info
140              (
141              Daily_Stock_ID Integer Primary Key Autoincrement,
142              Date Date,
143              Ticker_Symbol Text,
144              Open_Price Numeric,
145              Close_Price Numeric,
146              Lowest_Price Numeric,
147              Highest_Price Numeric,
148              Volume Real
149              )") #Create Table Script                                    .
150   dbWriteTable(conn = db_conn, #Database Connector name
151                name = "Daily_Stock_Info", #Table Name
152                value = Daily_Stock_Info, #Load data from the newly created dataframe
153                append=TRUE,row.names = FALSE,header = TRUE)#Since the CSV contains headers in the columns
154   dbListFields(db_conn, "Daily_Stock_Info") #Listing the field of the table
155
```

```
156 ▾ #----------------------------For Annual_Balance_Info---------------------------------|
157  dbGetQuery(db_conn, #Database Connector name
158            "create table Annual_Balance_Info
159            (
160            Annual_Stock_ID Integer Primary Key Autoincrement,
161            Financial_Year Integer,
162            Ticker_Symbol Text,
163            Cash_Cash_Equivalents Numeric,
164            Total_Current_Assets Numeric,
165            Total_Current_Liabilities Numeric,
166            Short_Term_Investments Numeric,
167            Total_Revenue Numeric,
168            Treasury_Stock Numeric,
169            Retained_Earnings Numeric,
170            Operation_Margin Numeric,
171            Quick_Ratio Numeric
172            )") #Create Table Script
173  dbWriteTable(conn = db_conn, #Database Connector name
174            name = "Annual_Balance_Info", #Table Name
175            value = Annual_Balance_Info, #Load data from the newly created dataframe
176            append=TRUE,row.names = FALSE,header = TRUE)#Since the CSV contains headers in the columns
177  dbListFields(db_conn, "Annual_Balance_Info") #Listing the field of the table
178
```

# 5. Calculation

We have done few calculations on Current Ratio, Cash Ratio, %Change in volume/week and %Change in price/week was not present in our dataset. This calculation helps in knowing fluctuation of stocks in the stock market per week and the liability of a company. This additional calculation makes the analysis even more interesting and helps an intelligent trader invest wisely.

- o Current Ratio
  - It is a ratio of Total Current Assets to Total Current Liabilities
  - If greater than 1, company has more assets than liabilities
  - If less than 1, company has more liabilities which is not good
  - If equal to 1, company has zero assets left

$$Current\ Ratio = \frac{Total\ Current\ Assets}{Total\ Current\ Liabilities}$$

- o Cash Ratio
  - The ratio of company's cash and cash equivalents to its current liabilities
  - If the company is forced to pay all current liabilities immediately, this metric shows the company's ability to do so without having to sell or liquidate other assets
  - If a company's cash ratio is equal to 1, the company has the same amount of current liabilities as it does cash and cash equivalents to pay off those debts.
  - If a company's cash ratio is less than 1, there are more current liabilities than cash and cash equivalents. In this situation, there is insufficient cash on hand to pay off short-term debt.
  - If a company's cash ratio is greater than 1, the company has more cash and cash equivalents than current liabilities. In this situation, the company can cover all short-term debt and still have cash remaining

$$Cash\ Ratio\ = \frac{Cash\ \&\ Cash\ Equivalent}{Total\ Current\ Liabilities}$$

```
288  #Calculations
289  Annual_Info_CR_Ratio<-Annual_Balance_Info %>%
290               select(Financial_Year,Ticker_Symbol,Total_Current_Liabilities,Total_Current_Assets)
291  Annual_Info_CR_Ratio$Current_Ratio<-Annual_Info_CR_Ratio$Total_Current_Assets/Annual_Info_CR_Ratio$Total_Current_Liabilitie
     s #Calculating the Current Ratio for all the stocks
292  Annual_Profit<-filter(Annual_Info_CR_Ratio, Current_Ratio>=1) #This dataframe contains all the companies that have higher
     assets than liabilities
293  Annual_Loss<-filter(Annual_Info_CR_Ratio, Current_Ratio<1) #This dataframe contains all the companies that have more
     liabilities than assets
294
295  Annual_Balance_Info$Cash_Ratio<-Annual_Balance_Info$Cash_And_Cash_Equivalents/Annual_Balance_Info$Total_Current_Liabilities
     #Calculating the Cash Ratio for all the stocks
```

# 6. Data Retrieving

## 6.1 Analysis using SQL queries

To transform structured data into meaningful information. We performed analysis by retrieving data using SQL queries.

**1) Display the Companies having a current ratio greater than one.**

```
query<-"select a.Ticker_Symbol,c.Company_Name, (a.Total_Current_Assets/a.Total_Current_Liabilities) as Current_Ratio
        from Annual_Balance_Info a
        join Company_Info c on a.Ticker_Symbol=c.Ticker_Symbol
        group by a.Ticker_Symbol,c.Company_Name
        having (a.Total_Current_Assets/a.Total_Current_Liabilities)>1
        order by (a.Total_Current_Assets/a.Total_Current_Liabilities) desc
        limit 10"
dbGetQuery(db_conn,query) #Calling dbquery to run the query on db_conn and show output
```

```
##     Ticker_Symbol          Company_Name Current_Ratio
## 1              FB              Facebook            11
## 2            LLTC Linear Technology Corp.            9
## 3            SWKS    Skyworks Solutions            9
## 4            MCHP   Microchip Technology            8
## 5             ADI   Analog Devices, Inc.            6
## 6            FAST            Fastenal Co            6
## 7            MNST       Monster Beverage            6
## 8             FTR Frontier Communications            5
## 9            ISRG Intuitive Surgical Inc.            5
## 10            WAT     Waters Corporation            5
```

Result:
- Current ratio of Facebook is 11. Facebook has the highest current ratio followed by Linear Technology Corp and Skyworks Solutions with a current ratio of 8.
- Microchip Technology, Analog Devices, Inc., Fastenal Co, Monster Beverage, Frontier Communications, Intuitive Surgical Inc. and Waters Corporation having the current ratio of **greater than one**. It means these 10 companies have **more assets than liabilities**.

**2) Which state has the highest number of headquarters?**

```
query<-"select State,count(Ticker_Symbol) as No_Of_Headquarters
        from Company_Info
        group by State
        order by count(Ticker_Symbol) desc
        limit 5"
dbGetQuery(db_conn,query) #Calling dbquery to run the query on db_conn and show output
```

```
##        State No_Of_Headquarters
## 1 California                 67
## 2   New York                 60
## 3      Texas                 39
## 4   Illinois                 30
## 5       Ohio                 23
```

Result:

- **California** has the highest number of Headquarters i.e. 67, which means **maximum money flow** from there.
- **New York** is the second highest state with 60 no. of headquarters followed by Texas, Illinois, and Ohio.

**3) Display all the GICS Sector along with their Total revenue.**

```
query<-"select c.GICS_Sector, sum(a.Total_Revenue) as Total_Revenue
        from Annual_Balance_Info a
        join Company_Info c on a.Ticker_Symbol=c.Ticker_Symbol
        group by c.GICS_Sector
        order by sum(a.Total_Revenue) desc"
dbGetQuery(db_conn,query) #Calling dbquery to run the query on db_conn and show output
```

```
##                       GICS_Sector Total_Revenue
## 1             Consumer Staples 5730167251000
## 2        Consumer Discretionary 5706202107000
## 3                        Energy 4833442768000
## 4                   Health Care 4489612571000
## 5                   Industrials 4103821019000
## 6                    Financials 3962733623000
## 7        Information Technology 3702212820000
## 8    Telecommunications Services 1151218853000
## 9                     Materials 1124278980000
## 10                    Utilities 1069025171000
## 11                  Real Estate  267223066000
```

Result:

-Among GICS_Sector, **Consumer Staples** has the **highest** total revenue.
- Consumer Discretionary is very near to Consumer Discretionary. Energy Sector has the third highest Revenue. **Real Estate** has the **lowest** revenue.

**4) When was the highest volume of stocks traded? For which company? How much?**

```
query<-"select d.Date, c.Company_Name, max(d.Volume) as Max_Volume
        from Daily_Stock_Info d
        join Company_Info c on d.Ticker_Symbol=c.Ticker_Symbol"
dbGetQuery(db_conn,query) #Calling dbquery to run the query on db_conn and show output
```

```
##        Date        Company_Name Max_Volume
## 1 8/25/2011 Bank of America Corp  859643400
```

Result:
- Bank of America Corp has the **highest volume** of stocks traded on **08-25-2011** with the volume of 859643400. It means there was a huge demand for the Bank of America stock on 08-25-2011.

**5) Which Industry has the highest number of stocks traded?**

```
query<-"select c.GICS_Sector, count(d.Volume) as No_Of_Stocks_Traded
        from Daily_Stock_Info d
        join Company_Info c on d.Ticker_Symbol=c.Ticker_Symbol
        group by c.GICS_Sector
        order by count(d.Volume) desc"
dbGetQuery(db_conn,query) #Calling dbquery to run the query on db_conn and show output
```

```
##                        GICS_Sector No_Of_Stocks_Traded
## 1        Consumer Discretionary              143783
## 2                   Industrials              115950
## 3        Information Technology              113484
## 4                    Financials              103708
## 5                   Health Care               99898
## 6               Consumer Staples               61182
## 7                        Energy               61170
## 8                   Real Estate               51098
## 9                     Utilities               49336
## 10                    Materials               42594
## 11 Telecommunications Services                8810
```

Result:
- **Consumer Discretionary** sector has the highest number of stock traded. Consumer Discretionary sector has 143783 stocks traded from 2010 to 2016. Consumer Discretionary stocks are in **demand**.
- Industrials sector has 115950 stocks traded and Information Technology has 113484 stocks traded.

**6) Display the Retained earnings for Industries under Telecommunications Sector?**

```
query<-"select c.GICS_Major_Industry, sum(a.Retained_Earnings) as Total_Retained_Earnings
        from Annual_Balance_Info a
        join Company_Info c on a.Ticker_Symbol=c.Ticker_Symbol
        group by c.GICS_Major_Industry
        having c.GICS_Sector='Consumer Discretionary'
        order by sum(a.Retained_Earnings) desc"
dbGetQuery(db_conn,query) #Calling dbquery to run the query on db_conn and show output
```

```
##                        GICS_Major_Industry Total_Retained_Earnings
## 1              Broadcasting & Cable TV            283610552000
## 2                          Restaurants            208440752000
## 3              Home Improvement Retail            142949000000
## 4               Automobile Manufacturers            137702000000
## 5          Hotels, Resorts & Cruise Lines            125837853000
## 6                      Specialty Stores            110537149000
## 7              General Merchandise Stores            105199275000
## 8   Apparel, Accessories & Luxury Goods             87349309000
## 9                       Apparel Retail             58390048000
## 10              Household Appliances             54392300000
## 11                        Advertising             38617200000
## 12   Internet & Direct Marketing Retail             36499653000
## 13              Motorcycle Manufacturers             32580178000
## 14                          Publishing             31899259000
## 15                        Homebuilding             31098539000
## 16                    Leisure Products             29344493000
## 17                   Cable & Satellite             28242000000
## 18                   Department Stores             24180000000
## 19               Auto Parts & Equipment             21920800000
## 20                      Tires & Rubber             16667000000
## 21                 Consumer Electronics             15787289000
## 22          Computer & Electronics Retail             15051000000
```

```
## 23              Home Furnishings       13148618000
## 24    Housewares & Specialties          8739100000
## 25             Specialty Retail          7147124000
## 26                 Distributors          6161206000
## 27            Automotive Retail          5390430000
## 28             Casinos & Gaming           330724000
```

Result:

- Under Telecommunication Sector, there are **28 GICS Major Industry**.
- Broadcasting & Cable TV has retained earnings of 283610552000. Broadcasting & Cable TV has the highest net earning which is not paid out as dividends but retained by the company to be **reinvested** in its core business whereas Casinos and Gaming have the lowest retained earnings in the Telecommunication Sector.

## 6.2. Data Visualization using R libraries

To perform advanced analysis, which gives more clarity in understanding the data. We have performed Data Visualization using R Libraries.

**1) What are the Total Current Liabilities and Assets of American Airlines Group?**

```
231 ▾ #------------------------DATA ANALYSIS---------------------------------
232   #1)What are the Total Current Liabilities and Assets of American Airlines Group?
233   Anual_Info_CR_Ratio<-Annual_Balance_Info %>%
234                 select(Financial_Year,Ticker_Symbol,Total_Current_Liabilities,Total_Current_Assets)
235   Anual_Info_CR_Ratio$Current_Ratio<-Anual_Info_CR_Ratio$Total_Current_Assets/Anual_Info_CR_Ratio$Total_Current_Liabilities #
236   Annual_Profit<-filter(Anual_Info_CR_Ratio, Current_Ratio>=1) #This dataframe contains all the companies that have higher as
237   Annual_Loss<-filter(Anual_Info_CR_Ratio, Current_Ratio<1) #This dataframe contains all the companies that have more liabil
238
239   AAL_Info<-Anual_Info_CR_Ratio %>%
240       select(Financial_Year,Ticker_Symbol,Total_Current_Liabilities,Total_Current_Assets,Current_Ratio) %>%
241       filter(Ticker_Symbol=='AAL') #Filtering for American Airlines Group
242   AAL_Info
243
244   highchart() %>%
245     hc_xAxis(categories = AAL_Info$Financial_Year) %>%
246     hc_title(text = "Total_Current_Liabilites vs Total_Currents_Assets for American Airlines Group from 2012 to 2015") %>% #
247     hc_add_series(
248       name='Total_Current_Liabilities', #Legend
249       color='red',#Color of the line graph
250       data=AAL_Info$Total_Current_Liabilities #Column to Display
251     ) %>%
252     hc_add_series(
253       name='Total_Current_Assets', #Legend
254       color='green', #Color of the line graph
255       data=AAL_Info$Total_Current_Assets #Column to Display
256     )
```
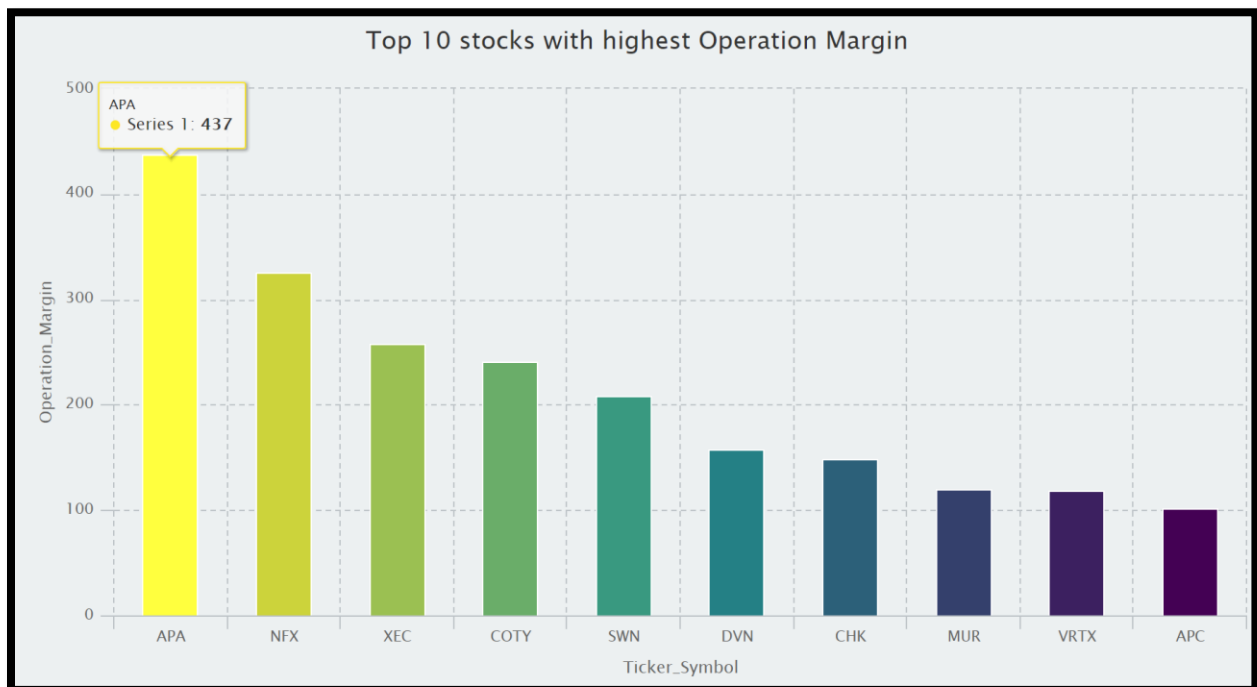


Total_Current_Liabilites vs Total_Currents_Assets for American Airlines Group from 2012 to 2015

Result:
- Only 2013 has a Current Ratio **greater than one** which is a **good sign** otherwise current ratio is less than one in the year of 2012, 2014 and 2015 which means total liabilities are greater than total assets.

**2) Plot a graph of how much a company makes on each dollar of sales [Operation Margin]**

```
258  #2)Plot a graph of how much a company make's on each dollar of sales [Operation Margin].
259  Opt_Margin<-Annual_Balance_Info%>%
260    select(Ticker_Symbol,Operation_Margin)%>%
261    arrange(desc(Operation_Margin))%>% #Arranging in descending order
262    head(n=10)%>% #Selecting top 10 Highest Operation Margin
263    distinct() #Selecting only distinct values as the Ticker_Symbl may repeat because of mutlple financial year
264
265  hchart(Opt_Margin,'column', hcaes(x = Ticker_Symbol, y = Operation_Margin, color = Operation_Margin)) %>%  #Plot typ
266    hc_add_theme(hc_theme_flat()) %>%  #Using flat theme for appealing visualization
267    hc_title(text="Top 10 stocks with highest Operation Margin") #Graph Title
268
```



Top 10 stocks with highest Operation Margin

Result:
- Operation Margin measures **profitability**. APA stocks have the highest Operation Margin which says they earn the highest profit on each dollar of revenue.
- APA stocks belong to Apache Corporation. Investing in APA stocks seems to be a safe bet for investors. The bar graph displays top 10 stocks having highest Operation Margin.

**3) Display the top 5 states having a higher number of headquarters.**

```
269  #3)Display the top 5 states having higher number of headquarters.
270  State_HQ<-Company_Info%>%
271          select(State)%>%
272          group_by(State)%>%
273          summarise(No_of_HQ = n())%>% #Counting the number of HeadQuarters
274          arrange(desc(No_of_HQ)) #Arranging in descending order [Max No. of HQ on top]
275
276  hchart(State_HQ[1:5,],'pie', hcaes(x = State, y = No_of_HQ, color = No_of_HQ)) %>%  #Plot type pie a
277     hc_add_theme(hc_theme_google()) %>%  #Theming for better visualization
278     hc_title(text="Top 5 states having higher number of headquarters") #Graph Title
279
```



Top 5 states having higher number of headquarters

Result:
- California has the 67 no. of headquarters from top 500 companies. Followed by New York, Texas, Illinois and Ohio. These are the top 5 states having highest no of headquarters.

**4) The volume of Stocks traded by Year and Month from 2010 to 2016.**

Volume of Stocks traded by Year and Month[2010-2016]

Result:

- Through this analysis, we can infer that, from 2010 to 2016, the highest volume of stocks was traded in **August 2016** and lowest volume of stocks were traded in **September 2012**.
- It is also observed that more than 10K volume of stocks is traded in the month of August from 2010 to 2016. August month has never seen a downfall in the volume of stocks traded. Investors usually **prefer investing** in the month of **August** because it's the middle of the year and maximum product launches at that time and hence the market is always in the boom.
- Also, there is another interesting insight that **less than 10K volume** of stocks were traded in the month of **February** always. This could be because it is the start of a year and budgets are usually not allocated thus it is always a risk.

## 5) The volume of Stocks traded by Month

```
288   #5)Volume of Stocks traded by Month
289   hchart(Volume_Month,'line', hcaes(x = Month, y = Total, color = Total)) %>% #Line plot
290     hc_title(text = "Volume of Stocks traded by month of the year")  #Title
```
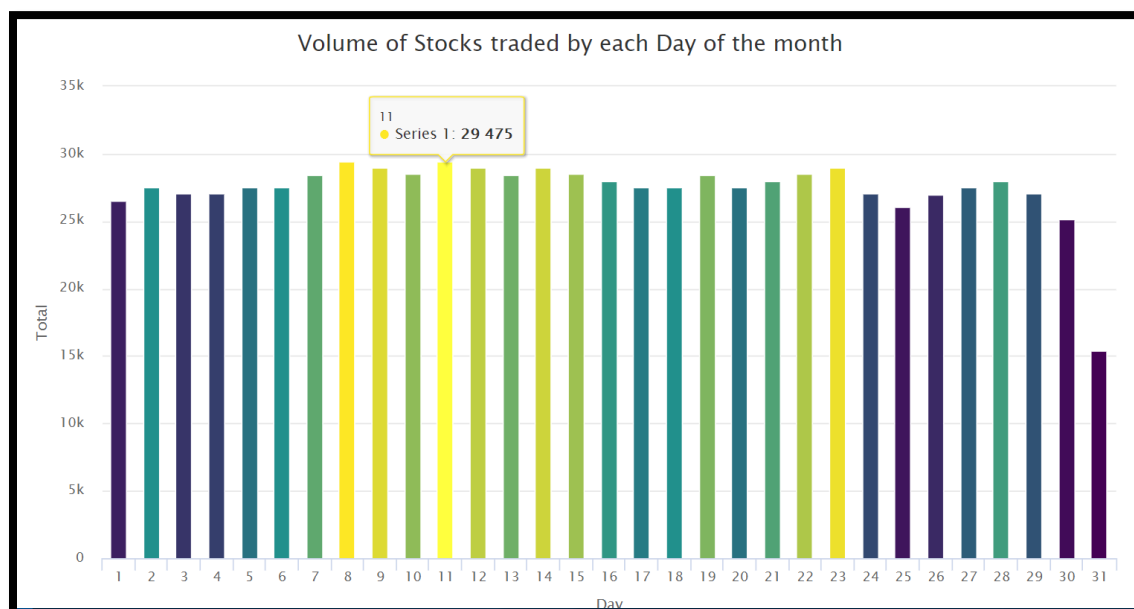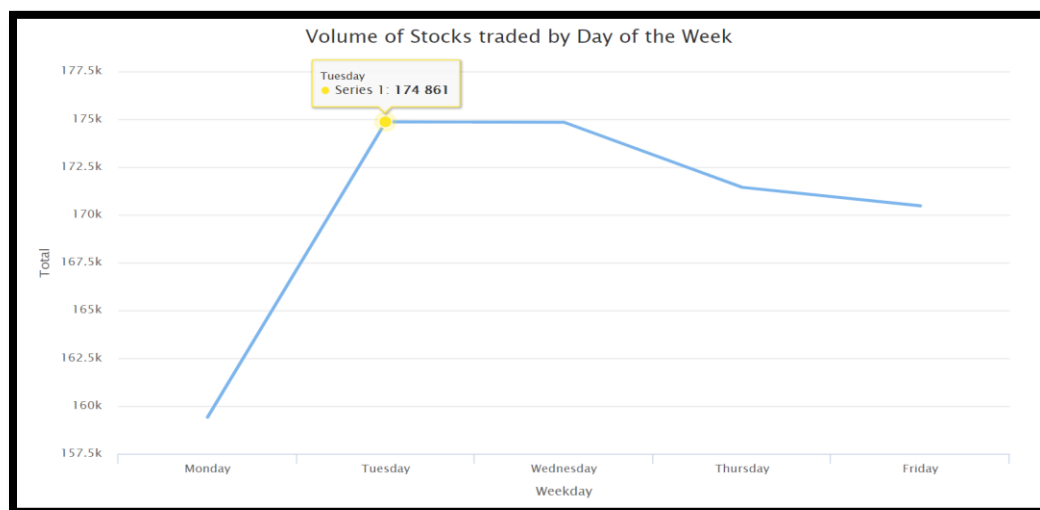
Volume of Stocks traded by month of the year

Result:
- Line graph clearly shows that **August has the highest** volume of stocks i.e. 74914 whereas February has the lowest number of stocks i.e. 65095.
- **January, February doesn't** seem to be a good time at Stock Market being the start of the year.
- The market seems to be very **uncertain** from **September to December**.

**6) The volume of Stocks traded by each day of the month**

```
292   #6)Volume of Stocks traded by each day of the month
293   hchart(Volume_Day,'column', hcaes(x = Day, y = Total, color = Total)) %>%  #Column Plot
294     hc_title(text = "Volume of Stocks traded by each Day of the month") #Title
```


Volume of Stocks traded by each Day of the month

Result:
- There is the **highest peak in the middle** of each month. There has been a record that on 11th of every month, the highest volume of stocks is traded.
- Bars with yellow to green color says high volume of stocks has been traded, which is seen in the middle of every month. Whereas dark green to purple and blue says the **low volume** of stocks have been traded, which is seen in **start and end** of every month.
- People get their salary at the start of every month and the first week usually they spend on clearing their dues which helps them give a clear idea on how much to invest. Thus, middle of the month is the time when people seem to invest in stock market.

**7) The volume of Stocks traded by Weekday**

```
296  #7)Volume of Stocks traded by Weekday
297  hchart(Volume_Weekday,'line', hcaes(x = Weekday, y = Total, color = Total)) %>%  #Line Plot
298    hc_title(text = "Volume of Stocks traded by Day of the Week") #Title
```
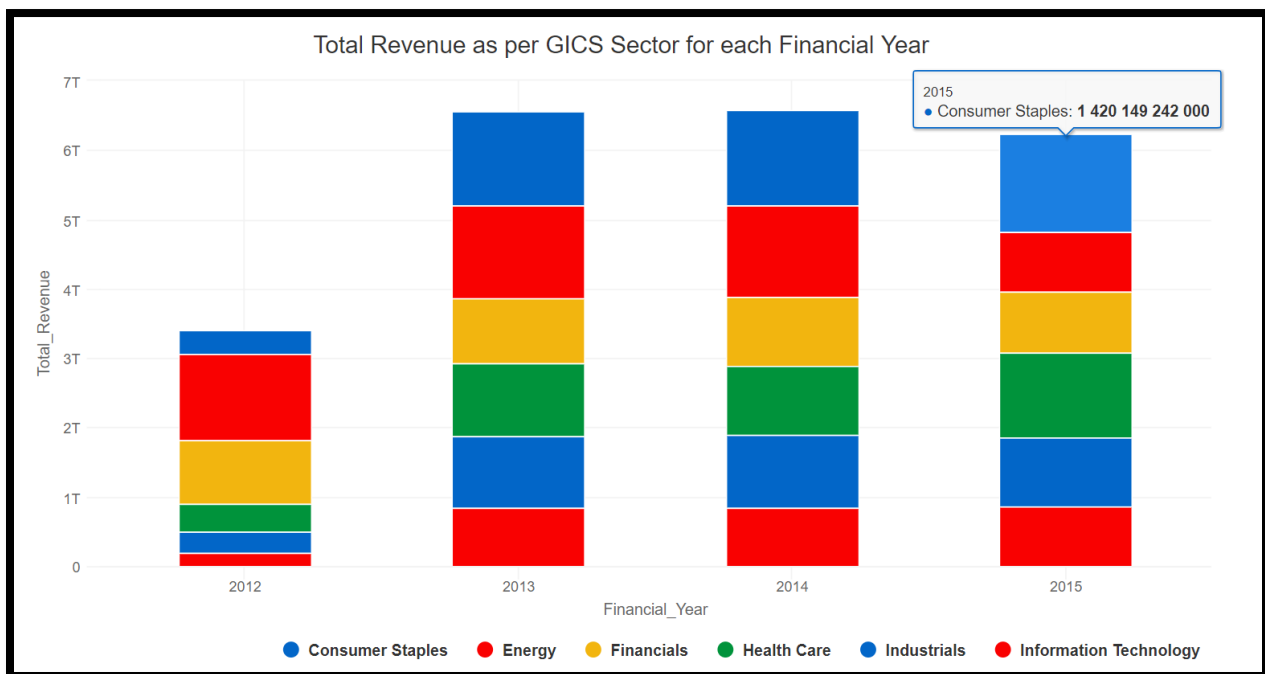


Result:
- **Lowest volume** of stocks are traded on **Monday**. As Monday is the first day of the week and it is also the busiest day of the week as well as Monday being the first day of the week for the stock market, the market is highly volatile thus making it unpredictable for investments.
- **The maximum volume** of stocks is traded on **Tuesdays and Wednesdays**. Tuesdays and Wednesdays are comparatively more predictable. Hence it is the preferred choice for the investors.

**8) Display the Total Revenue for Energy, Financials, Information Technology, Health Care, Consumer Staples and Industrials from 2012 to 2015**

```
300  #8)Display the Total Revenue for Energy,Financials,Information Technology,Health Care,Consumer Staples and Industrials fr
301  query<-"select a.Financial_Year,c.GICS_Sector, sum(a.Total_Revenue) as Total_Revenue
302          from Annual_Balance_Info a
303  join Company_Info c on a.Ticker_Symbol=c.Ticker_Symbol
304  group by a.Financial_Year,c.GICS_Sector
305  order by a.Financial_Year,c.GICS_Sector"
306
307  GICS<-dbGetQuery(db_conn,query) #Calling dbquery to run the query on db_conn and show output and storing it for future re
308  GICS<-na.omit(GICS) #Removing Nulls
309  GICS<-GICS%>% filter(Financial_Year %in% c("2012","2013","2014","2015")) #Filtering out for the period 2012- 2015
310  GICS<-GICS%>% filter(GICS_Sector %in% c("Energy",
311                                          "Financials",
312                                          "Information Technology",
313                                          "Health Care",
314                                          "Consumer Staples",
315                                          "Industrials"))
316  #Choosing only Top GICS sector for clear visualization
317
318  hchart(GICS,'column', hcaes(x = Financial_Year, y = Total_Revenue, group = GICS_Sector)) %>%
319    hc_add_theme(hc_theme_google()) %>%
320    hc_plotOptions(column = list(stacking = 'normal')) %>%
321    hc_legend(align = 'right', float = T)%>%
322    hc_title(text = "Total Revenue as per GICS Sector for each Financial Year")
```
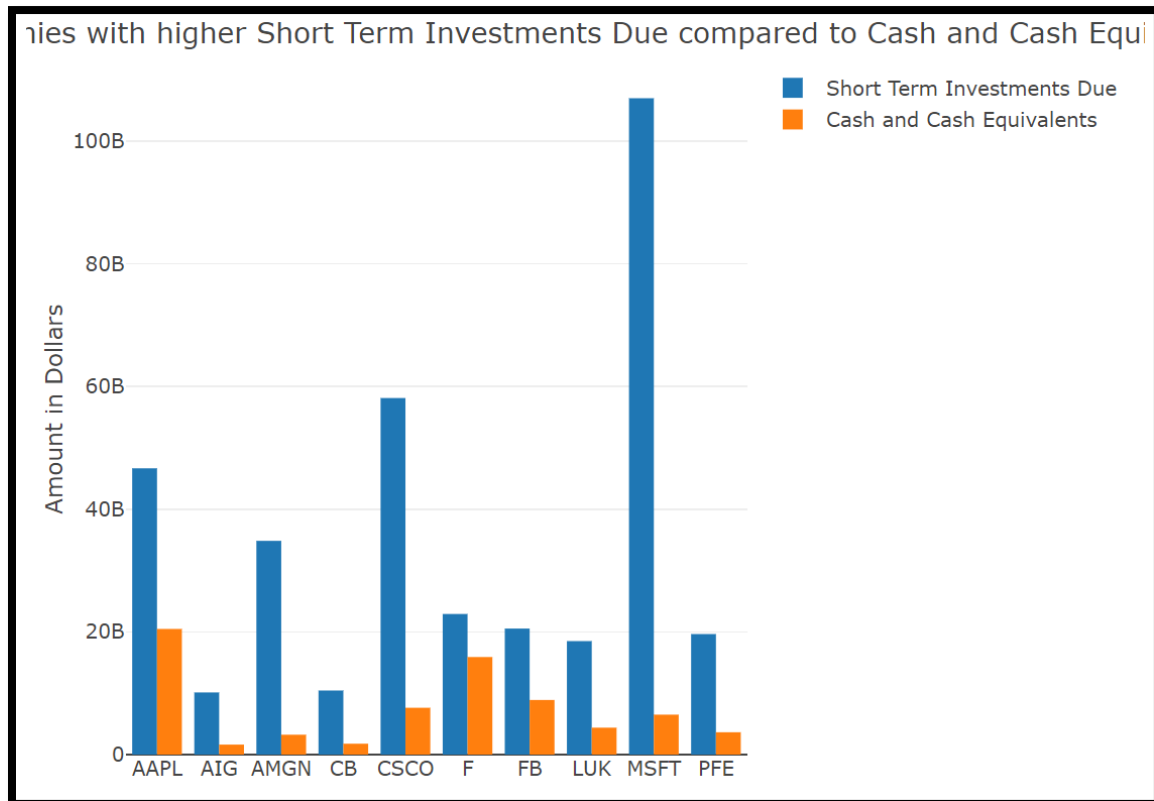


Result:
- From Data Retrieval Question 5 we found that Consumer Staples, Energy, Financials, Health Care, Industrials and Information Technology are top 5 GICS Sector having the highest number of stocks traded.
- To perform further analysis, we created a stacked graph to look for total revenue of top 5 GICS Sectors.
- Above graph shows that Total Revenue for Consumer Staples has increased every year from 2012 to 2015. **Consumer Staples** are basic necessities and what so ever be the individuals' financial conditions, consumption cannot be compromised. Hence it has an **ever-growing market**.
- Total Revenue is decreased in 2015 as compared to 2014 and 2013 because of a decrease in revenue from Energy, Financials, Industrials and Information Technology sectors.

**9) Which companies have more Short-Term Investments Due compared to Cash and Cash Equivalent.**

```
#9) Which companies have more Short Term Investments Due compared to Cash and Cash Equivalents
query<-"select Ticker_Symbol,Short_Term_Investments,Cash_Cash_Equivalents, Short_Term_Investments - Cash_Cash_Equivalents
        from Annual_Balance_Info
        group by Ticker_Symbol
        having Short_Term_Investments > Cash_Cash_Equivalents and Short_Term_Investments > 0 and Cash_Cash_Equivalents > 0
        order by (Short_Term_Investments - Cash_Cash_Equivalents) desc
        limit 10"
Loss_Stock<-dbGetQuery(db_conn,query) #Calling dbquery to run the query on db_conn and show output


plot_ly(Loss_Stock, x = Loss_Stock$Ticker_Symbol, y = Loss_Stock$Short_Term_Investments, type = 'bar', name = 'Short Term I
  add_trace(y = Loss_Stock$Cash_Cash_Equivalents, name = 'Cash and Cash Equivalents') %>%
  layout(yaxis = list(title = 'Amount in Dollars'), barmode = 'group',title='Comapnies with higher Short Term Investments D
```



Result:
- Short-Term Investments Due are the **debts** incurred by the company which is due within one year. This value helps in determining a company's performance.
- Cash and Cash Equivalents are company's assets that are **cash or can be converted** into cash immediately. Companies with a high amount of cash and cash equivalents are better and show higher liquidity.
- Thus, having less short-term investments due and higher cash & cash equivalents makes company stable and reliable to buy stocks.
-  The above bar graph shows **top 10 companies** which have more short-term investments due and less cash & cash equivalent, which makes them **less reliable**. This analysis can make investors think twice before buying their stocks.
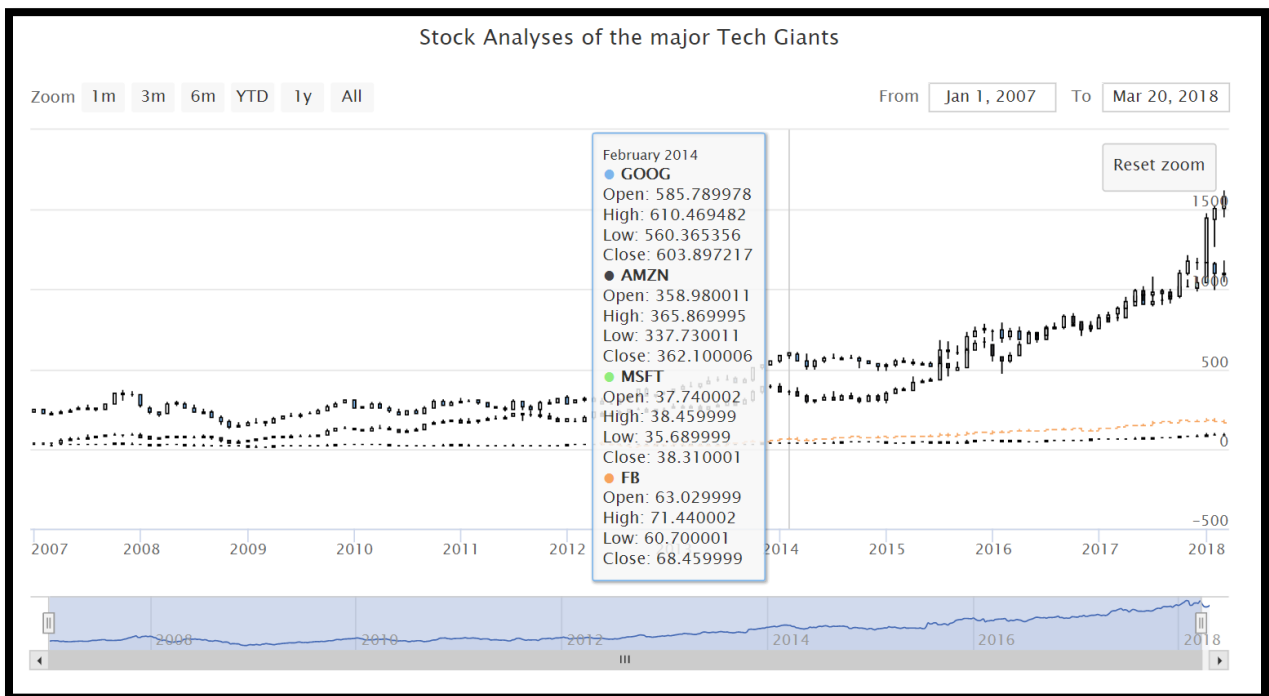
# 7. Additional

## 7.1 Data Mining – Time Series Analysis

Quantmod Package is a financial modeling R package which has a lot of financial trading options such as adjust OHLC (Open, High, Low, Close) stock price. Getsymbols is a wrapper to load data from different sources and it defaults to yahoo source. It helps in creating interesting and intuitive financial charts.
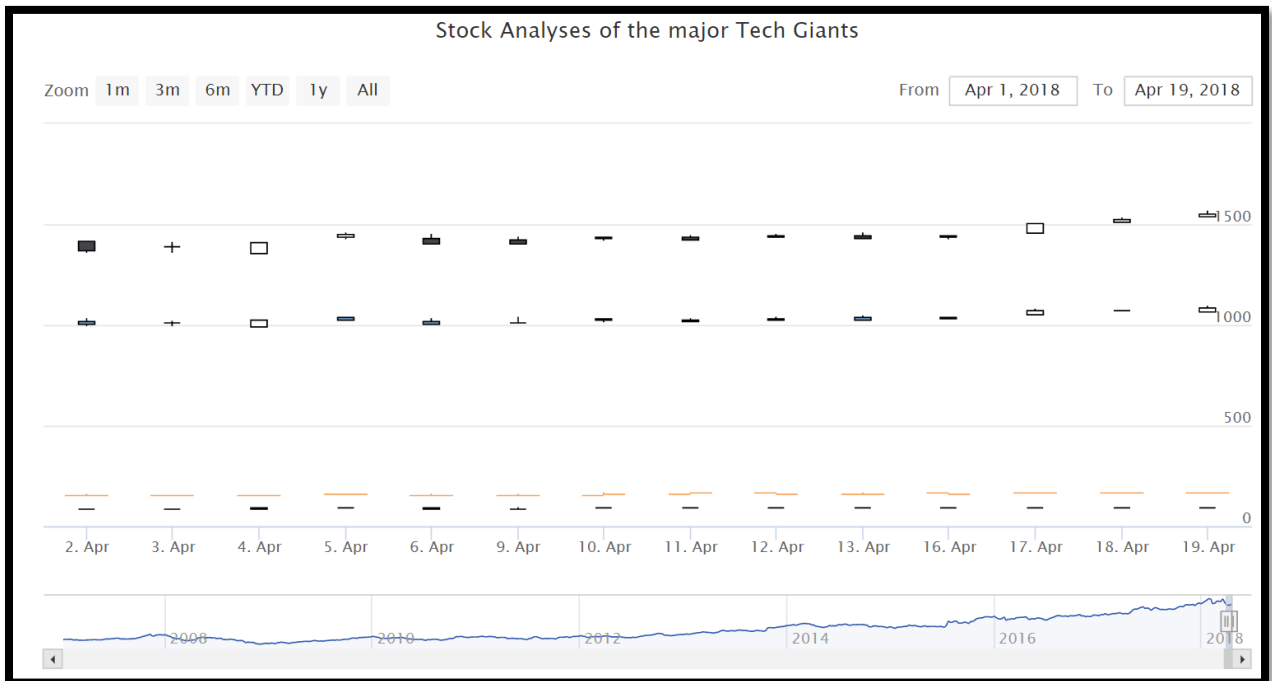
To generate Time Series Analysis, we are using Yahoo! Finance to get livestock market data from January 1, 2007, till now. Google, Amazon, Microsoft, and Facebook run the Information Technology sector, hence we decided to perform Time Series Analysis on top 4 Information Technology companies.

```
324 #---------------------------TIME SERIES ANALYSES---------------------------------
325 #Using quantmod package gathering live data for the following companies
326 Google<-getSymbols(Symbols="GOOG", src = "yahoo", auto.assign = FALSE) #Pulls data from Yahoo Finance for Google Stocks
327 Amazon<-getSymbols(Symbols="AMZN", src = "yahoo", auto.assign = FALSE) #Pulls data from Yahoo Finance for Amazon Stocks
328 Microsoft<-getSymbols(Symbols="MSFT", src = "yahoo", auto.assign = FALSE) #Pulls data from Yahoo Finance for Microsoft Stock
329 Facebook<-getSymbols(Symbols="FB", src = "yahoo", auto.assign = FALSE) #Pulls data from Yahoo Finance for Facebook Stocks
330
331 highchart(type="stock") %>%
332   hc_title(text="Stock Analyses of the major Tech Giants") %>%
333   hc_add_series(Google) %>%
334   hc_add_series(Amazon) %>%
335   hc_add_series(Microsoft) %>%
336   hc_add_series(Facebook,type="ohlc")
337
```
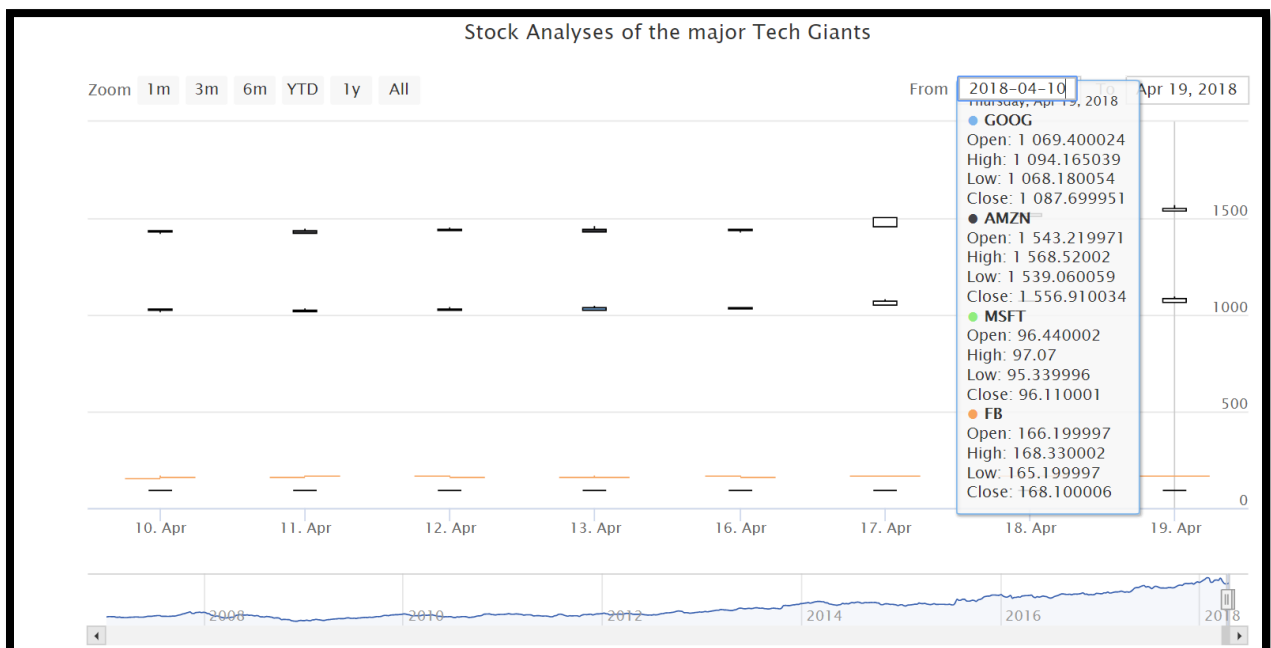


- Above Time series analysis takes the **live data** containing Open, Close, Highest and Lowest price of Google, Amazon, Microsoft and Facebook stocks from **Yahoo Finance source**. It helps you in comparing the price of 4 stocks. This analysis is very useful for investors who are keen on investing in these stocks and looking for the right time to gain maximum profit.
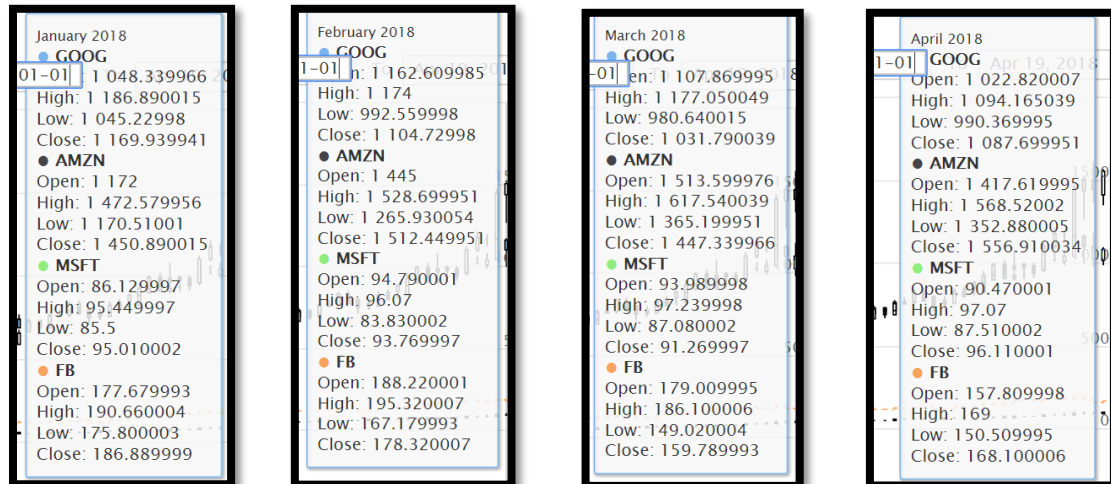
- Refer the following Link for **Dynamic Time Series**, as above graph is understood better when seen through the HTML format: (Download the HTML to view it)
- https://drive.google.com/open?id=140j85Tjm6BuBsgcVJMVwtXCJotdakrR2



- Dynamic Time Series graph helps you change and look for the desired Start and End date of the time series graph.
- In the above screenshot, we are looking for the latest stock prices of Amazon, Google, Facebook and Microsoft from April 1, 2018, to April 19, 2018.

- On April 21, Google stocks open at 1069 which is the highest while comparing with Amazon, Microsoft, and Facebook.
- High and Low stock price of Microsoft and Facebook is not fluctuating. It is very stable and not risky to invest.



- We have taken a screenshot from the Dynamic Time Series graph to compare the stock prices for Amazon, Google, Microsoft, and Facebook for the month of January, February, March and April (until now) which will help to analyse the stocks before investing in them.
- As seen above, Google's stock prices are highly fluctuating since February. Hence, due to such a volatile nature of the stock, an investor should consider other parameters too before investing.

## 7.2 Forecasting/Prediction Algorithm

Using the historical data, we are forecasting the volume of stocks traded for the future years. We are predicting the data using two packages. One is Prophet package, and another is forecast and tries to build the ARIMA model.

```
#Forecasting Packages
library("tseries")
library("StanHeaders")
library("rstan")
library("prophet")
library("forecast")
```

### 7.2.1 Prophet – Facebook's Forecasting R Package

Prophet, Facebook's R package is a procedure to forecast time series data. It is an open source software released by Facebook's Data Science core team. It is used for non-linear trends and best used for daily periodic values having historical data.
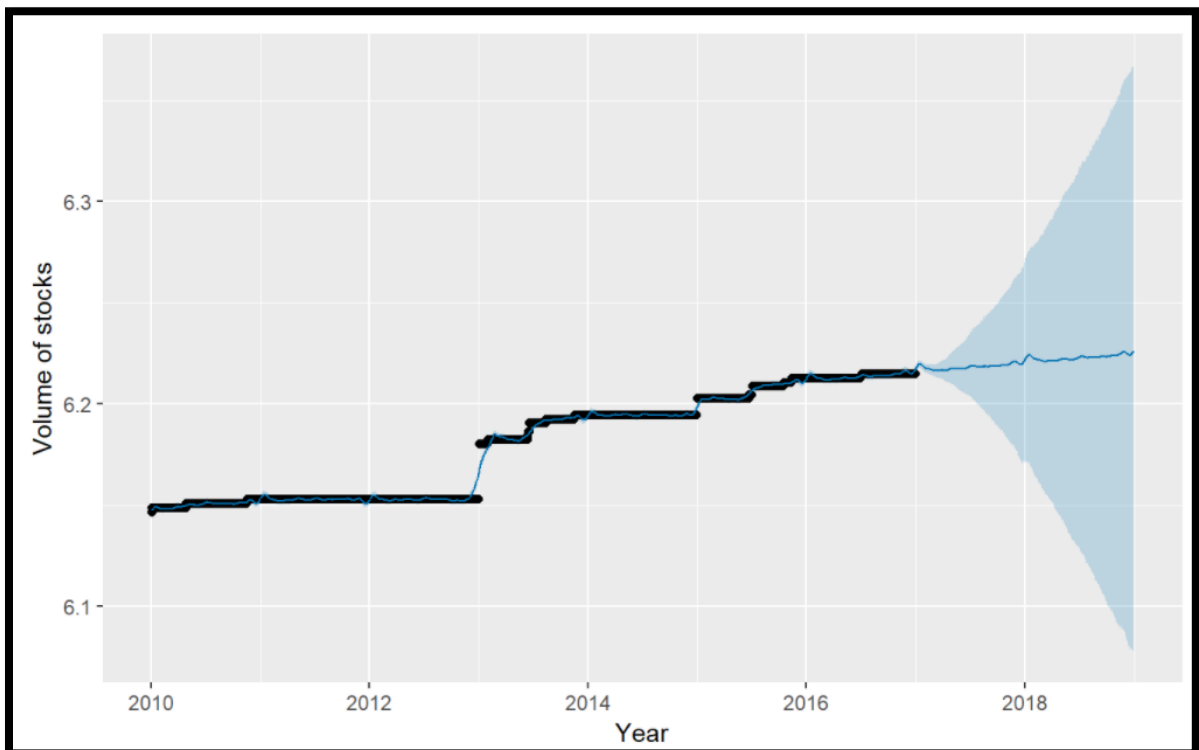
Using the historical data from 2010 and assuming to have all unpredictable factors (such as geo political factors like government changing policies, civil instabilities or

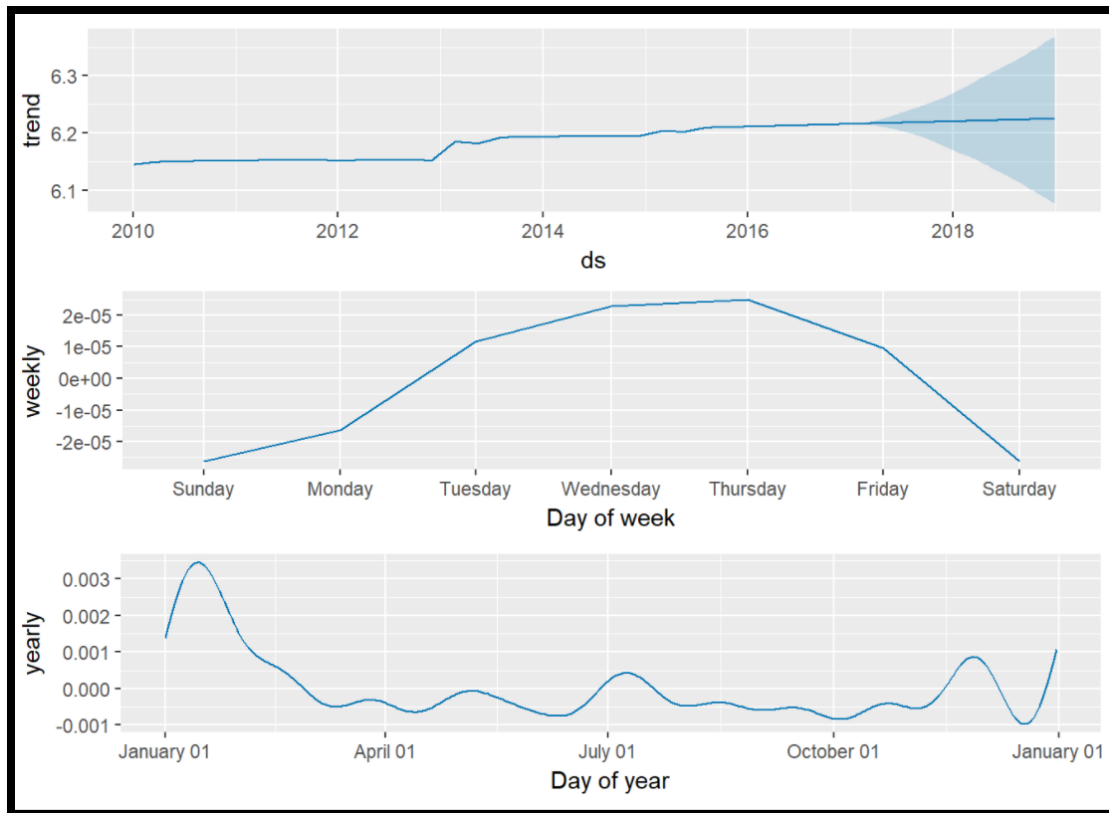natural calamities), we are predicting the volume of stocks traded for the year 2017 and 2018.

```r
#FORECASTING
##Prophet(Forecasting the Volume of Stocks)
```{r}
#Step1:Collect Data
Volume_Date<-Daily_Stock_Info %>%
            group_by(Date) %>%   #Grouping data by date
            summarise(Total = n()) #Volume of stocks

#Step2:Calculate the logarithmic parameter
names(Volume_Date)<-c("ds","y")
Volume_Date<-Volume_Date %>%
            mutate(y = log(y))
Volume_Prophet<-prophet(Volume_Date)

#Step3:Forecasting
Future<-make_future_dataframe(Volume_Prophet,periods = 730) #Period = 730 because we are forecastig it for two years
Forecast_Prophet<-predict(Volume_Prophet,Future)
plot(Volume_Prophet,Forecast_Prophet,ylab="Volume of stocks",xlab="Year") #Plotting a graph for the Volume of stocks with Year
prophet_plot_components(Volume_Prophet,Forecast_Prophet) #Plotting other forecasting components such as trends, weekly and yearly
```
```



- As per the historical data, it can be forecasted that there is a minor increase in the volume of stock traded. The volume of stocks increased to 6.23 in 2018 to 6.21 in 2016.

- There has been a change in trend which defines peoples trading pattern, the gradual shift from Tuesday and Wednesday which was seen in Data Visualization section to Wednesday and Thursday in the year 2017 and 2018.
- August always has been a crucial year for trading until 2016, but it has been forecasted that January will supersede August in the year 2018.

### 7.2.2 ARIMA Model

- ARIMA model is the Autoregressive integrated moving average. ARMIA model is applied in cases where data show evidence of non-stationarity, where an initial differencing step (corresponding to the "integrated" part of the model) can be applied one or more times to eliminate the non-stationarity.
- ARIMA model is used to forecasting model that utilize historical information to make predictions.
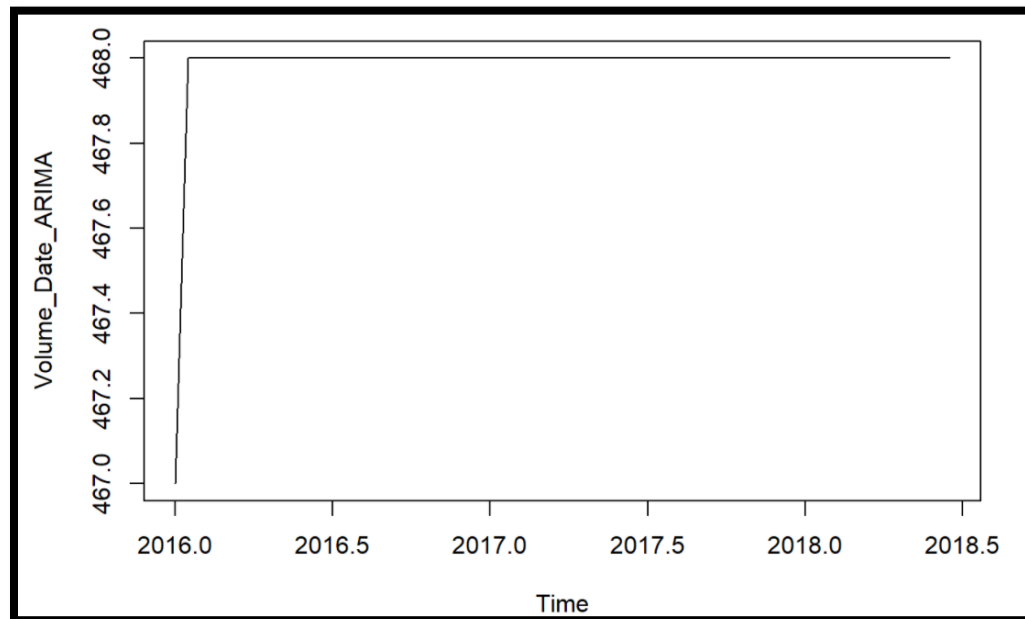
Step 1: Collecting data

Collecting data and limiting it to few companies for better understanding and clear analysis.

```
#Step1:Prepare Data
Volume_Date_ARIMA<-Daily_Stock_Info %>%
  group_by(Date) %>%  #Grouping data by date
  summarise(Total = n()) #Volume of stocks
```
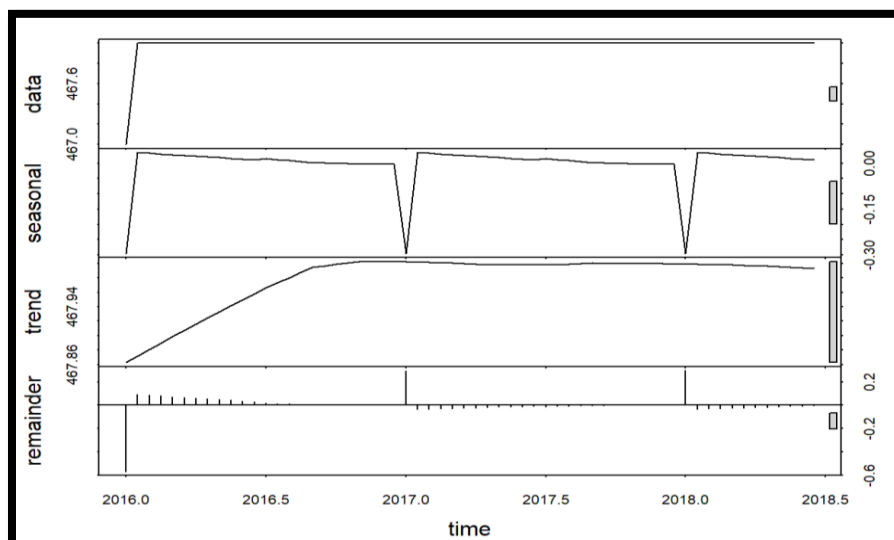
**Step 2:** Creating Time Series

```
#Step2:Creating Time Series
Volume_Date_ARIMA = ts(na.omit(Volume_Date_ARIMA$Total), start=c(2016,1), end=c(2018,12),frequency=24) #Preparing time
series
plot(Volume_Date_ARIMA)#Plot of the High Price of the stocks
```



**Step 3:** Decomposing data

Decomposing the data to remove seasonality trends and outliers.

```
#Step3:Decomposing the Data
Decomp = stl(Volume_Date_ARIMA, s.window="periodic") #STL is a flexible function for decomposing and forecasting the
series.
Deseasonal<-seasadj(Decomp) #Returns seasonally adjusted data constructed by removing the seasonal component.
plot(Decomp)
```

Step 4: Stationary

ARIMA model requires series to be stationary. ADF (Augmented Dickey-Fuller) test is a statistical test for stationarity. The test checks whether the change in Y can be explained by lagged value and a linear trend.
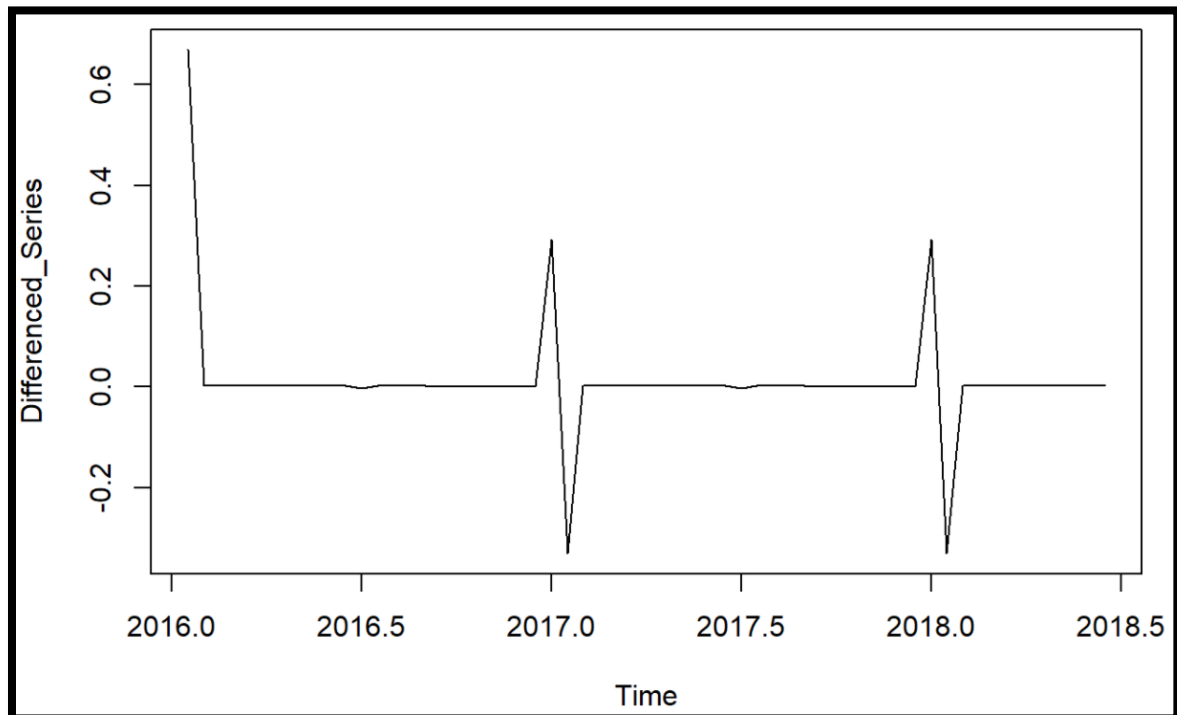
```
#Step4:Stationary
adf.test(Volume_Date_ARIMA,alternative="stationary") #ADF procedure tests whether the change in Y can be explained by
lagged value and a linear trend.
```
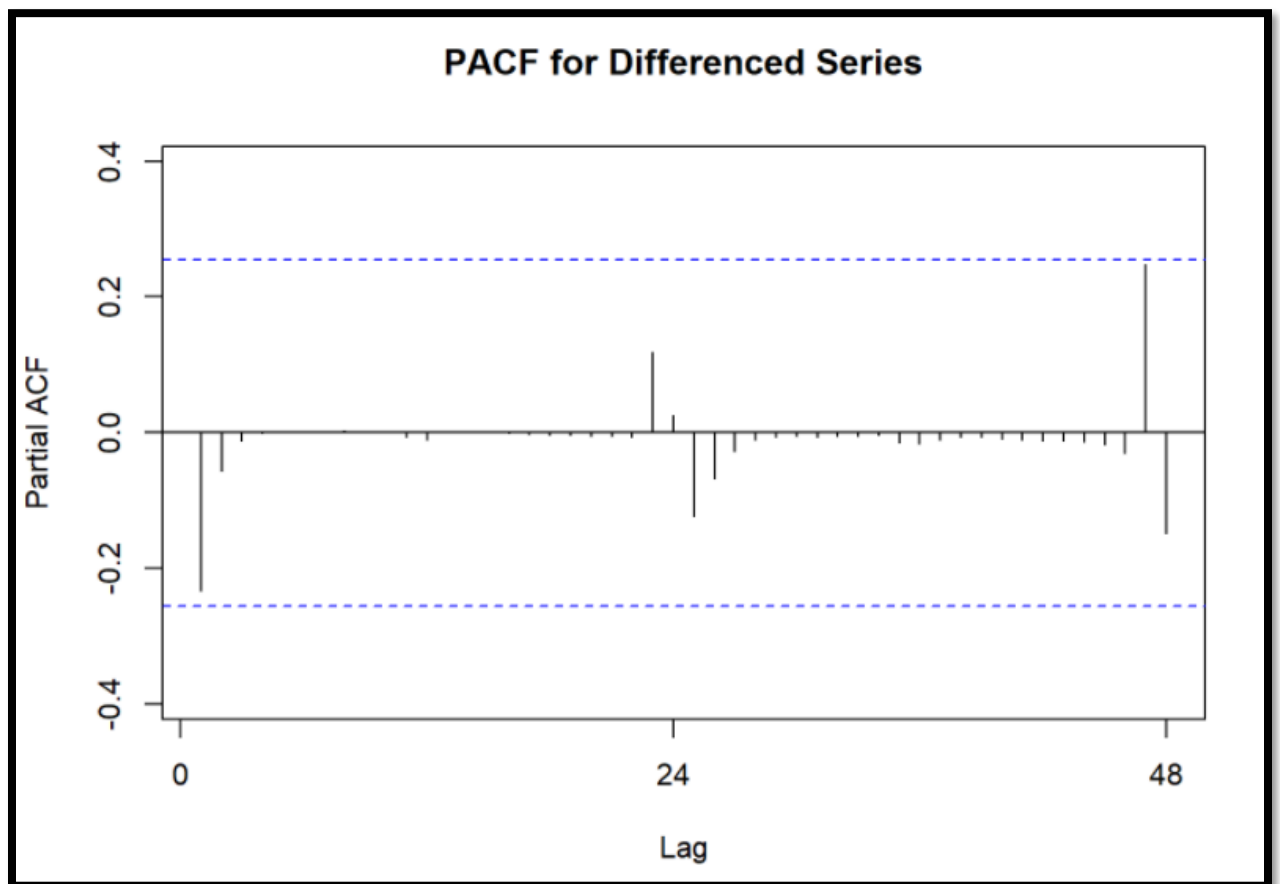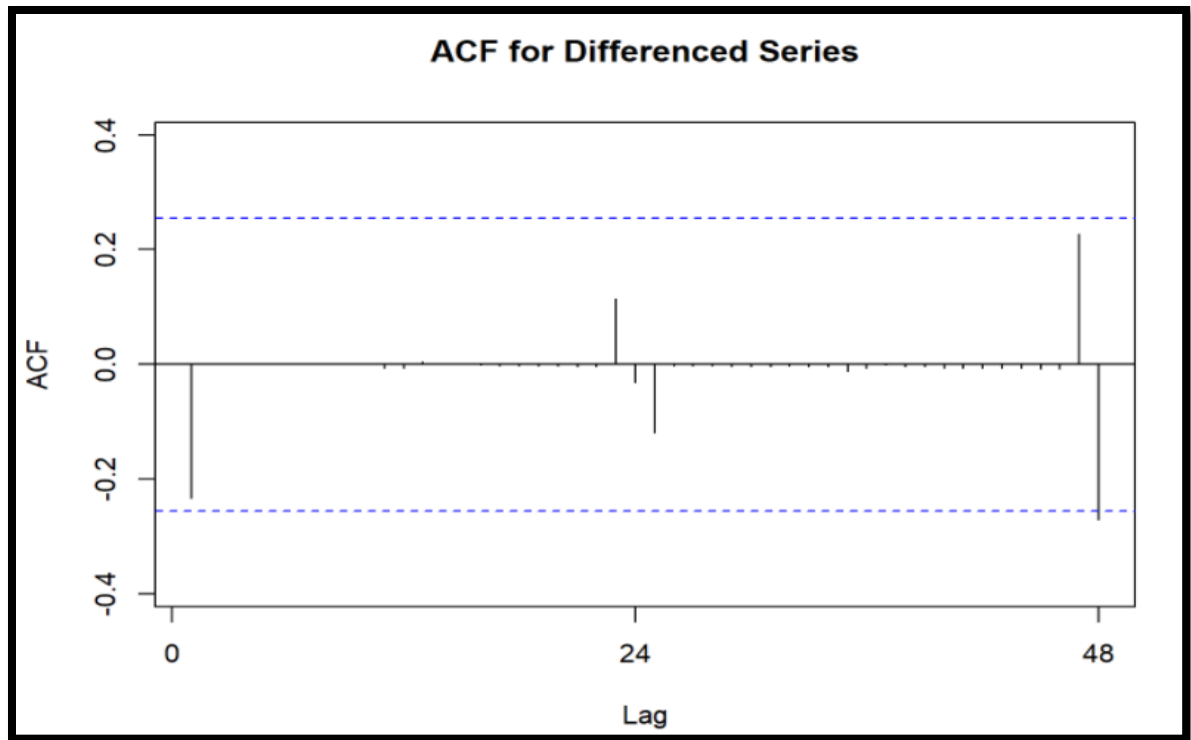
Step 5: Choosing Model Order

ACF() plots correlation between a series and its lags. PACF () at K lag is autocorrelation function to plot the correlation between all data points that are exactly k steps apart

```
#Step5:Autocorrelations and Choosing Model Order
Differenced_Series = diff(Deseasonal,differences = 1)
plot(Differenced_Series)
adf.test(Differenced_Series,alternative="stationary")

Acf(Differenced_Series,main='ACF for Differenced Series')
Pacf(Differenced_Series,main='PACF for Differenced Series')
```
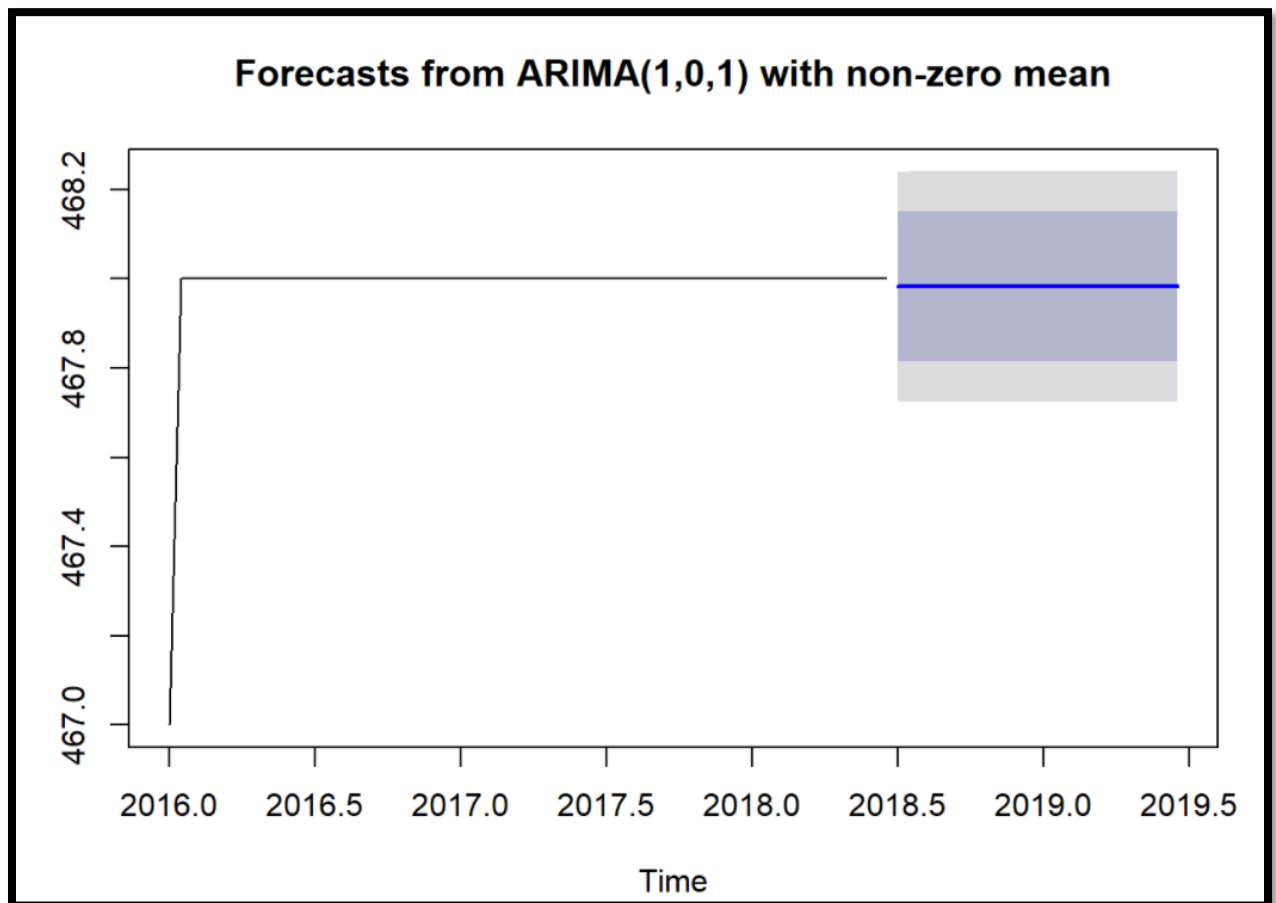
## ACF for Differenced Series



## PACF for Differenced Series



ACF and PACF residuals show no significant autocorrelations.

Step 6: Fitting ARIMA Model

An auto.arima function automatically selects an appropriate lag value using statistical tests and trains a linear regression model.

```
476  #Step6:Fitting ARIMA Model
477  Arima=auto.arima(Volume_Date_ARIMA,trace = TRUE,test="kpss",ic="aic") #Fits the model using approximation and
     no-approximation and picks the Best Model
478  plot.ts(Arima$residuals) #Plots the Residual component of the ARIMA
479
480  Arima_Forecast = forecast(Arima, h=24) #Forecasts the model for 24 months (2 Year36
481  Arima_Forecast
482  plot(Arima_Forecast, xlab="Time") #Plots the forecast graph for the best model selected above by Auto_ARIMA
```



Forecasts from ARIMA(1,0,1) with non-zero mean

- Using ARIMA model it is predicted that Volume of Stocks traded remains same in 2018 and 2019 compared to previous year
- One can find minimal difference stocks traded from 2016 to 2019.

### 7.2.3 Comparison between Prophet and ARIMA
- Prophet utilizes a Bayesian-based curve fitting method to forecast the time series data while
- ARIMA uses the traditional method for forecasting.
- The cool thing about Prophet is that it doesn't require much prior knowledge or experience of forecasting time series data since it automatically finds seasonal trends beneath the data and offers a set of 'easy to understand' parameters.

- Hence, it allows non-statisticians to start using it and get reasonably good results that are often equal or sometimes even better than the ones produced by the experts.
- But this is not always the case that Prophet is better ARIMA always.

## 7.3. Stock Market Web Application – RShiny

R lets us create easy and interactive web applications using the R Shiny Package. We can host standalone apps on a webpage or embed them in R Markdown documents or dashboards. To make the project more interactive, we decided to build few R Shiny applications.

All the data frames including input files are on the cloud which makes the R Shiny Application accessible from any device including mobile phones.

**Application 1: Volume of Stocks Traded**
Link: https://stockmarket.shinyapps.io/Volume_Of_Stocks_Traded/
Description:
- This app displays the volume of stocks traded between 2012 to 2016 by all the S&P 500 Companies.
- The user can choose the duration from a SliderInput for which he/she wants to view the volume of stocks. Eg: Monthly
- This will display a plot for all the volume of stocks for the period of 2012 to 2016 by month.
- The graph is displayed using plotly so it has additional functions of downloading the graph as .png or taking a snapshot, variable zoom options and reset axes.

Code:

```r
2   #Packages Required for the Application's Components to work
3   library("tidyverse")
4   library("dplyr")
5   library("plotly")
6   library("shiny")
7
8   # Use a fluid Bootstrap layout
9   ui =
10    fluidPage(
11      titlePanel("Volume of Stock Traded from 2012-2016"), #Title of the page/application
12
13      sidebarLayout(   #Generates a row with SideBar input
14        sidebarPanel( # Define the sidebar with one input
15
16          selectInput("duration", #InputID
17                      "Duration:", #Input Title
18                      choices= c( #Input Value
19                        "Yearly"=1,
20                        "Monthly"=2,
21                        "Weekly"=3,
22                        "Daily"=4 )
23                      )
24        ),
25
26        mainPanel( #Create a spot for the Plot
27          plotlyOutput("stockPlot")
28                  )
29
30                )
31              )
32
```
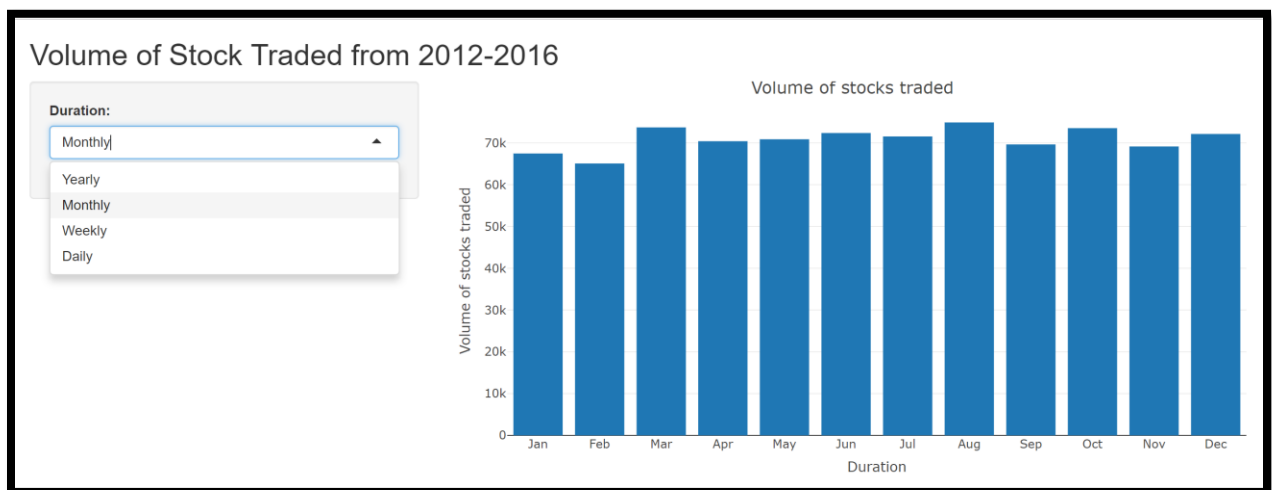
User Interface R Code for Volume of Stock Traded Analysis

```
33  # Define a server for the Shiny app
34 ▾ server<-function(input, output) { #Define a server for the Shiny app
35
36    names(Volume_Year)<-c("Duration","Total")
37    names(Volume_Month)<-c("Duration","Total")
38    names(Volume_Weekday)<-c("Duration","Total")
39    names(Volume_Day)<-c("Duration","Total")
40
41    #Creating a dataframe consisting of total Volume of stocks as per different duration
42    Volume_Stocks<-rbind(Volume_Year,Volume_Month,Volume_Weekday,Volume_Day)
43
44    #Creating reactive components which will assign dynamic values
45 ▾  Duration<-reactive({
46      switch(input$duration,
47             "1"=Volume_Stocks$Duration[1:5],
48             "2"=Volume_Stocks$Duration[6:17],
49             "3"=Volume_Stocks$Duration[18:22],
50             "4"=Volume_Stocks$Duration[23:53]
51            )
52    })
53
54 ▾  Total<-reactive({
55      switch(input$duration,
56             "1"=Volume_Stocks$Total[1:5],
57             "2"=Volume_Stocks$Total[6:17],
58             "3"=Volume_Stocks$Total[18:22],
59             "4"=Volume_Stocks$Total[23:53]
60            )
61    })
62
63 ▾  output$stockPlot <- renderPlotly({ #Fill in the spot we created for a plot
64      p<-plot_ly(x = Duration(), y = Total(), type = 'bar', name = 'Duration') #Render the Plot
65      layout(p, xaxis = list(title = "Duration"),yaxis = list(title = 'Volume of stocks traded'),title='Volume of stocks trad
66    })
67  }
68
69  #Connecting UI and Server
70  shinyApp(ui=ui, server=server)
71
```

Server R Code for RShiny Application 1.



Shiny Application 1 displaying Volume of Stocks Traded with Scroll Down options to select Duration

**Application 2: Annual Financial Stock Analysis**
Link: https://stockmarket.shinyapps.io/Annual_Financial_Stock_Analysis/
Description:
- This app displays the different financial parameters to choose for analysis.
- Using the Radio Button, the user can choose the Financial parameters such as Operation Margin, Current Ratio, Total Current Assets and Liabilities and so on.
- A Slider Input allows the user to pick the Stock Company in particular for which he/she wants to perform the analysis.

- The graph has a fixed X-Axis representing the financial year ranging from 2013 to 2016.
- This graph is also displayed using plot function of R.

Code:

```
2    #Packages Required for the Application's Components to work
3    library("tidyverse")
4    library("dplyr")
5    library("plotly")
6    library("shiny")
7
8    # Use a fluid Bootstrap layout
9    ui =
10     fluidPage(
11       titlePanel("Annual Financial Stock Analysis"), #Title of the page/application
12
13         sidebarLayout(   #Generates a row with SideBar input
14           sidebarPanel( # Define the sidebar with one input
15
16             selectInput("stock", #InputID
17                         "Choose the Stock:", #Input Title
18                         choices= c( #Input Value
19                           "Apple"=1,
20                           "Bristol-Myers Squibb"=2,
21                           "Costco Inc"=3,
22                           "Microsoft"=4,
23                           "TripAdvisor"=5,
24                           "Visa Inc."=6,
25                           "Yahoo Inc."=7
26                         )
27                       ),
28
29           #RadioButtons to choose the Financial Parameter
30           radioButtons("param", #InputID
31                         "Choose the Financial Parameter:", #Input Title
32                         choices= c( #Input Value
33                           "Cash and Cash Equivalents"=1,
34                           "Total Current Assets"=2,
35                           "Total Current Liabilities"=3,
36                           "Short Term Investments Due"=4,
37                           "Current Ratio"=5,
38                           "Operation Margin"=6,
39                           "Quick Ratio"=7
```

```
45
46             mainPanel( #Create a spot for the Plot
47               plotlyOutput("stockPlot")
48                 )
49
50                 )
51               )
52
```
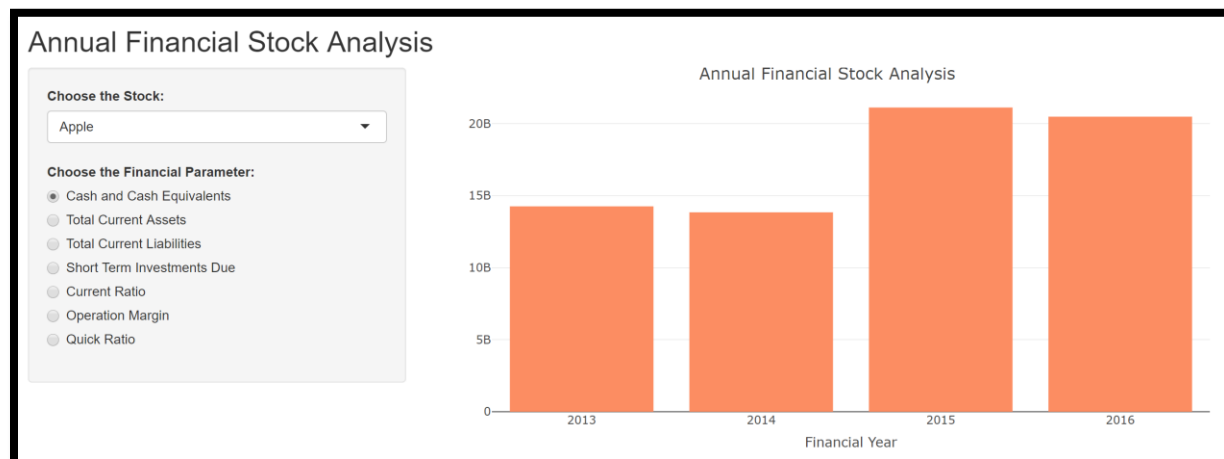
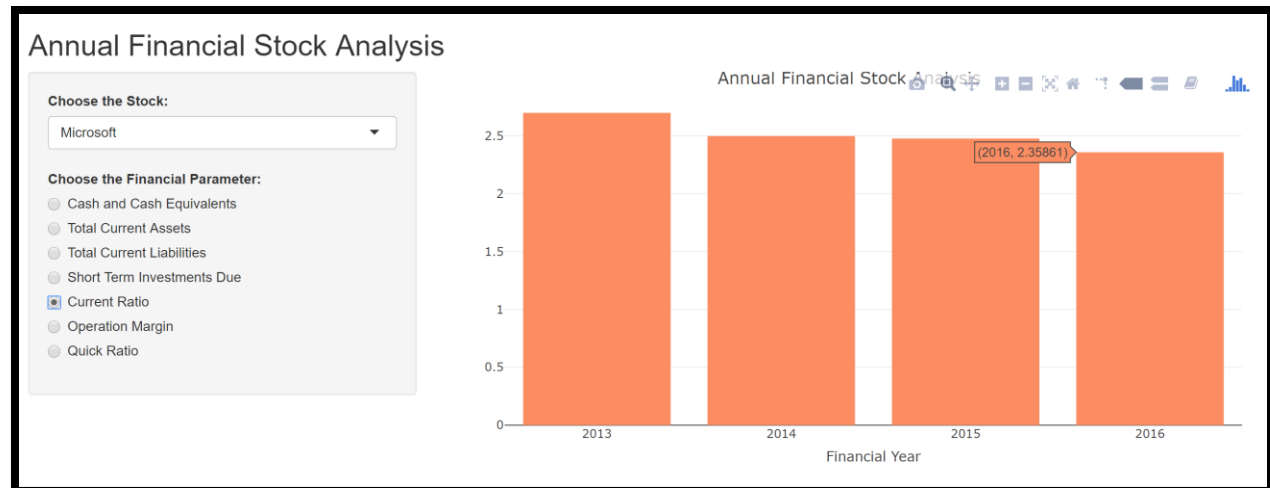User Interface Code for Annual Stock Financial Analysis

```
53    # Define a server for the Shiny app
54    server<-function(input, output) { #Define a server for the Shiny app
55
56
57        #Creating a dataframe consisting of all the financial parameters of stocks
58        Annual_Analysis<-Annual_Balance_Info%>% select("Financial_Year","Ticker_Symbol","Cash_And_Cash_Equivalents","Total_
59        Annual_Analysis<-na.omit(Annual_Analysis)
60        Annual_Analysis$Current_Ratio<-Annual_Analysis$Total_Current_Assets/Annual_Analysis$Total_Current_Liabilities
61
62        #Creating reactive components which will assign dynamic values
63        Stock<-reactive({
64          switch(input$stock,
65                "1"="AAPL",
66                "2"="BMY",
67                "3"="COST",
68                "4"="MSFT",
69                "5"="TRIP",
70                "6"="V",
71                "7"="YHOO"
72                )
73        })
74
75        Year<-c("2013","2014","2015","2016") #Data Frame consisting of the Financial Years to be displayed on X-Axes
76
77        output$stockPlot <- renderPlotly({ #Fill in the spot we created for a plot
78          Choose_Stock<-Annual_Analysis%>%
79            filter(Ticker_Symbol==Stock()) #Filter based on the Stock Name chosen
80
81          Fin_Param<-reactive({
82            switch(input$param,
83                "1"=Choose_Stock$Cash_And_Cash_Equivalents,
84                "2"=Choose_Stock$Total_Current_Assets,
85                "3"=Choose_Stock$Total_Current_Liabilities,
86                "4"=Choose_Stock$Short_Term_Investments,
87                "5"=Choose_Stock$Current_Ratio,
88                "6"=Choose_Stock$Operation_Margin,
89                "7"=Choose_Stock$Quick_Ratio
90
91            )
```

```
93          p<-plot_ly(x = ~Year, y = Fin_Param(), type = 'bar', color='Red') #Render the Plot
94          layout(p, xaxis = list(title = "Financial Year"),title='Annual Financial Stock Analysis') #Name the axes of Plot
95        })
96    }
97
98    #Connecting UI and Server
99    shinyApp(ui=ui, server=server)
100
```

Server Code for Annual Financial Stock Analysis



Shiny App 2 for Annual Financial Stock Analysis

Shiny Application 2: Selecting the Stocks to visualize different Annual Financial Stock Analysis.

# 8. Challenges Faced

1) Collecting data from heterogeneous sources and transforming them into a common format to work with.
2) Storing data in SQLite database. Since the data was collected from different sources each of them had different formats and types. So, while creating tables in the database it was critical to managing these different data types properly.
3) Deriving new columns by performing additional calculations. Since there was a high amount of missing and dirty data, the data first needed to be cleaned properly to carry out further calculations.
4) Working with Quantmod package to plot time series graph. Since, it was the first time we used the finance Quantmod package, to understand the components and plot it over a time series took a lot of effort.
5) Forecasting for additional 2 years using ARIMA model. Understanding how to evaluate the model to plot ACF and PACF.
6) It was extremely difficult building the R shiny app. Since this was the first time we were building the app from scratch. Tutorials from R Shiny website helped us a lot. The main challenge was to make the application more interactive and reactive.

# 9. Learning Outcomes

1) Scrapping, gathering, cleaning and consolidating techniques on the huge volume of data to fit 6 V's of big data is learned.
2) Creating meaningful and structured schema in its 3NF and then establishing an SQLite connection is skilled
3) Data retrieval is done using SQL queries and data visualization techniques are performed using ggplot2, plotly, ggthemes, highcharter packages
4) Data Mining Technique for Dynamic Time Series Analysis is performed.
5) Implemented forecasting technique using Prophet package and learned how to build ARIMA model to predict the data for future years.
6) Constructed a web application to perform parameter selection tool using RShiny.

## 10. Future Scope

1) Expanding the dataset by adding Stock Indexes, International Market Overview, Interest Rates etc. can increase the complexity in this project for which faster and efficient methods to collect, store and retrieve needs to be found.
2) Training the ARIMA model with more data to increase the accuracy. Also implementing other machine learning algorithms,
3) R Shiny app can be developed which more features and technicality. The user interface can be made more attractive.

## 11. References

- https://ntguardian.wordpress.com/2017/03/27/introduction-stock-market-data-r-1/
- http://www.visualcapitalist.com/difference-nyse-nasdaq/
- https://www.investopedia.com/terms/c/cashandcashequivalents.asp
- https://www.nerdwallet.com/blog/investing/stock-market-basics-everything-beginner-investors-know/
- http://www.visualcapitalist.com/difference-nyse-nasdaq/
- https://cran.r-project.org/web/packages/quantmod/quantmod.pdf
- https://dantonnoriega.github.io/ultinomics.org/post/2017-04-05-highcharter-explainer.html
- http://jkunst.com/highcharter/highstock.html
- https://www.datascience.com/blog/introduction-to-forecasting-with-arima-in-r-learn-data-science-tutorials
- http://ucanalytics.com/blogs/step-by-step-graphic-guide-to-forecasting-through-arima-modeling-in-r-manufacturing-case-study-example/
  https://shiny.rstudio.com/tutorial/