# Introduction

In this project, exploratory data analysis (EDA) was conducted on large, structured and semi-structured datasets. With the help of EDA, important analysis about data and its relation with various parameters was studied. The different types of datasets analyses were twitter data set, New York Times ad click dataset, Real Direct dataset and live twitter data.

# Problem 1: Data Acquisition

For this problem, the data source used was Twitter data using Twitter Search API, collected over a period of ten days. Twitter Streaming API is used in Problem 5. TwitteR and ROAuth were the packages that were used to carry out twitter data collection. TwitteR provided functions like setup_twitter_oauth which creates a connection to Twitter's Search API. The resulting list of retrieved tweets is stored in a data frame by using the function twListToDF(). Over a period of ten days, I created different data frames that stored that day's tweets. At the end, I combined all these data frames into one unit and stored this collective data as both cvs file and Rdata by using write.cv and save commands respectively.
tf is the local data frame containing collected tweets. These are the column variables it has.

```
> colnames(tf)
 [1] "text"          "favorited"     "favoriteCount" "replyToSN"     "created"       "truncated"     "replyToSID"
 [8] "id"            "replyToUID"    "statusSource"  "screenName"    "retweetCount"  "isRetweet"     "retweeted"
[15] "longitude"     "latitude"
>
```

This is the cvs file with tweets, stored on disk.

Falguni Bharadwaj

# Problem 2: Simple EDA

In this problem, we were provided with the New York Times dataset which had the information about the ad clicking behavior of readers with attributes such as age, gender, impressions, number of clicks, signed in or not.

First part of the problem is to categorize the age variable into intervals of age.

```
data1$agecat <-cut(data1$Age,c(-Inf,0,18,24,34,44,54,64,Inf))
```

Next, plotting impressions and click-through-rates for this interval of age. Here, upon noticing a high value for(Inf, 0], I realized that these users were not logged in. So setting the signed in value for these as NA and generated the bar plot again.



As is seen from the figure, the bar plot shows the number of impressions for the age intervals which is highest for age group (34,44].



This figure shows the number of clicks for each age group. It can be inferred that people from age group (54,64] appeared to click on the ads the most, followed by (64,Inf] (not counting the unsigned users). This shows that elderly people are more interested in (or easily distracted by) online advertisements. Similarly age group (18,24] seemed to click on ads the least which is probably because most of the people in this age group are students or young working professionals who might not spend much time clicking on ads.

Following shows the density distribution of CTR with respect to age. First figure is for Clicks>0 and second figure is for Impressions>0.



Following are a histogram and box plot that show the relation between CTR and age group.



Falguni Bharadwaj

```
ggplot(data=data1, aes(x=agecat, y=Clicks, fill=agecat)) + geom_bar(stat="identity") + theme_bw()
data1$agecat[data1$Signed_In == 0] = NA
data1$Gender[data1$Signed_In == 0] = NA
summary(data1)
ggplot(data=data1, aes(x=agecat, y=Clicks, fill=agecat)) + geom_bar(stat="identity") + theme_bw()
data2 = na.omit(subset(data1, Impressions>0)) %>% group_by(agecat) %>% summarise(Impressions = sum(Impressions), Clicks = sum(Clicks))
ggplot(data=data2, aes(x=agecat, y=Clicks/Impressions, fill=agecat)) + geom_bar(stat="identity") + theme_bw()
ggplot(subset(data1, Impressions>0), aes(x=Clicks/Impressions, colour=agecat)) + geom_density()
ggplot(subset(data1, Clicks>0), aes(x=Clicks/Impressions, colour=agecat)) + geom_density()
#First graph histogram impressions with age
ggplot(data1, aes(x=Impressions, fill=agecat)) +geom_histogram(binwidth=1)
#Second graph boxplot
ggplot(data1, aes(x=agecat, y=Impressions, fill=agecat)) +geom_boxplot()
```

Next, a new variable "clickcat" was defined to differentiate between users based on their click behavior which would hold three values: No Impressions, Only Impressions(No Click) and Clicks.

```
data1$clickcat[data1$Impressions==0] <- "NoImps"
data1$clickcat[data1$Impressions >0] <- "Imps"
data1$clickcat[data1$Clicks >0] <- "Clicks"
```

Next part of the problem was to explore the data and make visual and quantitative comparisons across user segments. Here are some observations. This following graph shows a normal density distribution graph for click through rate patterns of males and females of all age groups



The following graph compares the Click Through Rate patterns of less than 18 year old male & females.



Falguni Bharadwaj

In the above graphs, figure 1 shows the density distribution for logged in users and figure 2 shows density distribution for not logged in users.



Here, there is a normal distribution of age with males usually twice the number of females except for >65 where difference starts to reduce.

Now looking at some statistics.

```
> s <- function(x){c(length = length(x),min = min(x),mean = mean(x), max = max(x), median = median(x))}
> summaryBy(Age~agecat, data =data1, FUN=s)
     agecat Age.length Age.min Age.mean Age.max Age.median
1  (-Inf,0]     137106       0  0.00000       0          0
2    (0,18]      19252       7 16.03350      18         16
3   (18,24]      35270      19 21.26904      24         21
4   (24,34]      58174      25 29.50335      34         30
5   (34,44]      70860      35 39.49468      44         39
6   (44,54]      64288      45 49.49258      54         49
7   (54,64]      44738      55 59.49819      64         60
8  (64, Inf]     28753      65 72.98870     108         72

> summaryBy(Gender+Signed_In+Impressions+Clicks~agecat,data =data1)
     agecat Gender.mean Signed_In.mean Impressions.mean Clicks.mean
1  (-Inf,0]   0.0000000              0         4.999657  0.14207985
2    (0,18]   0.6421151              1         4.998961  0.13105132
3   (18,24]   0.5338531              1         5.006635  0.04845478
4   (24,34]   0.5321621              1         4.993829  0.05048647
5   (34,44]   0.5316963              1         5.021507  0.05167937
6   (44,54]   0.5289790              1         5.010406  0.05027377
7   (54,64]   0.5361885              1         5.022308  0.10183736
8  (64, Inf]   0.3632664              1         5.012347  0.15128856
```
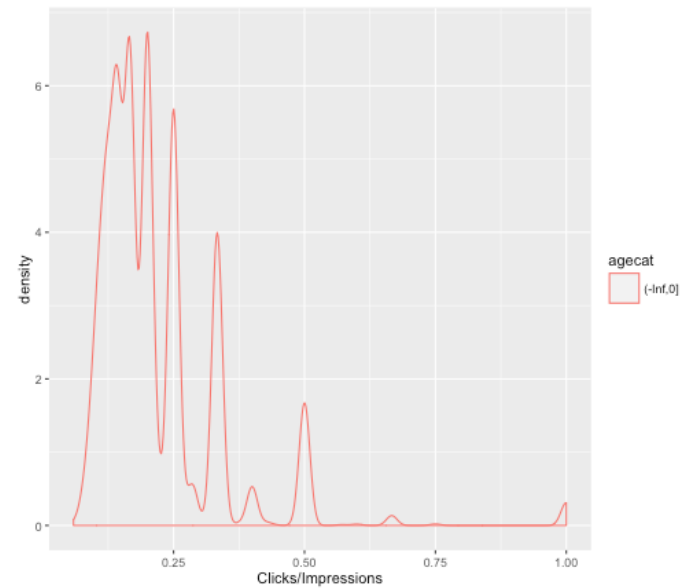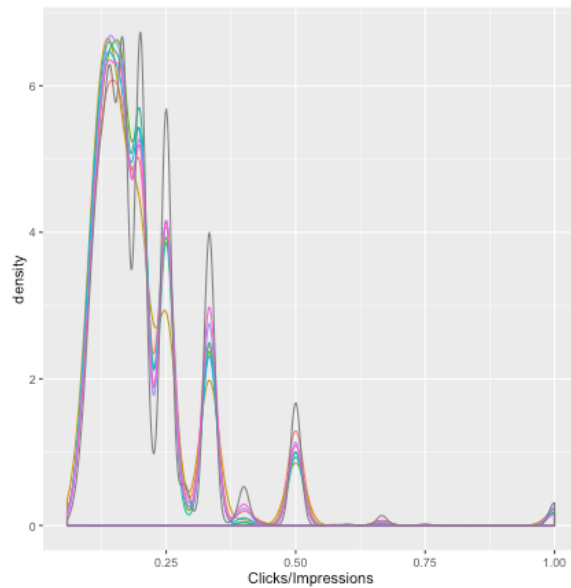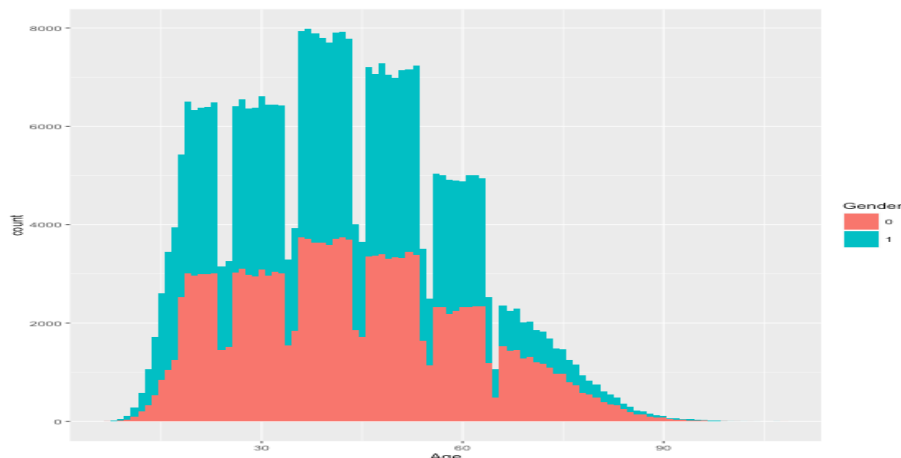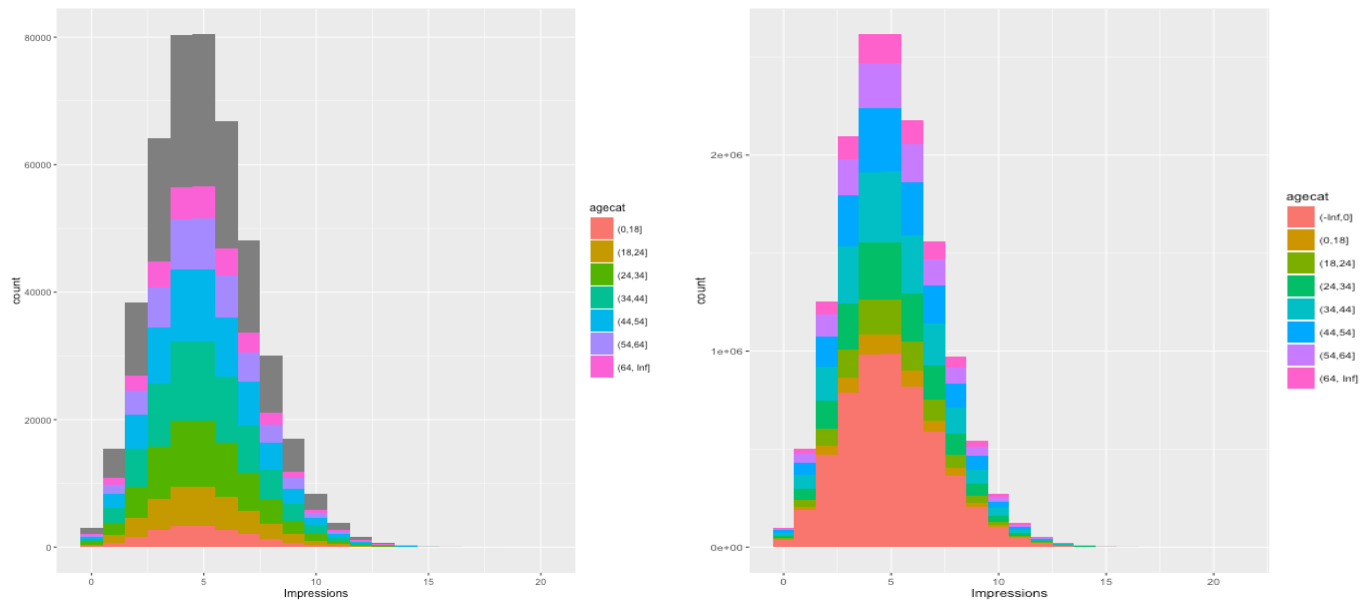
Last part of the problem, was to extend the findings and calculation analysis from one day to a month. As we can see from the graphs below, the data is normally distributed over the month too.





The following graph shows the density distribution over 30 days. It gives a click through rate pattern of users with respect to age over a period of 30 days.



Falguni Bharadwaj