

Week 12 IP

Ted Askoye Samuel

1. Business Understanding

1 a.) Defining the Question

A Kenyan entrepreneur has created an online cryptography course and would want to advertise it on her blog. She currently targets audiences originating from various countries. In the past, she ran ads to advertise a related course on the same blog and collected data in the process. She has employed the services of Skoko Limited, a Data Science Consultancy to help her identify which individuals are most likely to click on her ads.

2. Defining the Metrics of Success

The success of this analysis will occur when the target audience is known as per the adverts.

3. Context

Advertising is everywhere online, but we've gotten pretty good at ignoring it. To win back our attention, advertisers have adapted to our digital viewing habits by remembering what we read and buy online, then using this information to sell us things they think we might like. Part of this strategy is Targeted advertising. Targeted Advertising is a form of online advertising that focuses on the specific traits, interests, and preferences of a consumer. Advertisers discover this information by tracking your activity on the Internet.

4. Experimental Design

We will define the question, the metric of success, context and experimental design taken. This will be followed by reading and exploring the dataset and its appropriateness of the available data to answer the given question. This will be followed by cleaning the data off outliers, anomalies and null values from missing data, perform an exploratory data analysis after which we will record our observations and provide a conclusion and recommendation.

5. Data Relevance

Our data is very relevant to our research question. The more you know about your audience, the better you'll be able to sell advertising to them. The dataset provided has relevant information about the blog's audience.

6. Loading relevant Libraries and Reading the Data

```
# Importing the required packages  
  
library("data.table")  
library("dplyr")
```

```
##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:data.table':
##
##   between, first, last

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

library("tidyverse")

## -- Attaching packages ----- tidyverse 1.3.0 --

## v ggplot2 3.3.2      v purrr 0.3.4
## v tibble 3.0.3       v stringr 1.4.0
## v tidyr 1.1.2        v forcats 0.5.0
## v readr 1.3.1

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::between()   masks data.table::between()
## x dplyr::filter()    masks stats::filter()
## x dplyr::first()     masks data.table::first()
## x dplyr::lag()       masks stats::lag()
## x dplyr::last()      masks data.table::last()
## x purrr::transpose() masks data.table::transpose()

library("tidyr")
library("lubridate")

##
## Attaching package: 'lubridate'

## The following objects are masked from 'package:data.table':
##
##   hour, isoweek, mday, minute, month, quarter, second, wday, week,
##   yday, year

## The following objects are masked from 'package:base':
##
##   date, intersect, setdiff, union

library("ggcorrplot")
library("ggplot2")
library("corrplot")

## corrplot 0.84 loaded
```

```
library("moments")

# Loading the Dataset

ad_df <- read.csv(url("http://bit.ly/IPAdvertisingData"))
```

Previewing the data

```
# Previewing The First Seven records in the Dataset
```

```
head(ad_df, n=7)
```

```
##   Daily.Time.Spent.on.Site Age Area.Income Daily.Internet.Usage
## 1                68.95  35    61833.90                256.09
## 2                80.23  31    68441.85                193.77
## 3                69.47  26    59785.94                236.50
## 4                74.15  29    54806.18                245.89
## 5                68.37  35    73889.99                225.58
## 6                59.99  23    59761.56                226.74
## 7                88.91  33    53852.85                208.36
##               Ad.Topic.Line           City Male   Country
## 1   Cloned 5thgeneration orchestration Wrightburgh    0   Tunisia
## 2   Monitored national standardization   West Jodi    1     Nauru
## 3   Organic bottom-line service-desk     Davidton    0 San Marino
## 4 Triple-buffered reciprocal time-frame West Terrifurt    1     Italy
## 5   Robust logistical utilization        South Manuel    0   Iceland
## 6   Sharable client-driven software      Jamieberg    1     Norway
## 7   Enhanced dedicated support          Brandonstad    0   Myanmar
##   Timestamp Clicked.on.Ad
## 1 2016-03-27 00:53:11      0
## 2 2016-04-04 01:39:02      0
## 3 2016-03-13 20:35:42      0
## 4 2016-01-10 02:31:19      0
## 5 2016-06-03 03:36:18      0
## 6 2016-05-19 14:30:17      0
## 7 2016-01-28 20:59:32      0
```

```
# Previewing The Last Seven records in the Dataset
```

```
tail(ad_df, n=7)
```

```
##   Daily.Time.Spent.on.Site Age Area.Income Daily.Internet.Usage
## 994                64.20  27    66200.96                227.63
## 995                43.70  28    63126.96                173.01
## 996                72.97  30    71384.57                208.58
## 997                51.30  45    67782.17                134.42
## 998                51.63  51    42415.72                120.37
## 999                55.55  19    41920.79                187.95
## 1000               45.01  26    29875.80                178.35
```

```
##               Ad.Topic.Line           City Male
## 994      Phased zero tolerance extranet  Edwandsmouth  1
## 995      Front-line bifurcated ability  Nicholasland  0
## 996      Fundamental modular algorithm   Duffystad  1
## 997      Grass-roots cohesive monitoring  New Darlene  1
## 998      Expanded intangible solution  South Jessica  1
## 999  Proactive bandwidth-monitored policy  West Steven  0
## 1000     Virtual 5thgeneration emulation  Ronniemouth  0
##               Country           Timestamp Clicked.on.Ad
## 994      Isle of Man 2016-02-11 23:45:01           0
## 995      Mayotte 2016-04-04 03:57:48           1
## 996      Lebanon 2016-02-11 21:49:00           1
## 997  Bosnia and Herzegovina 2016-04-22 02:07:01           1
## 998      Mongolia 2016-02-01 17:24:57           1
## 999      Guatemala 2016-03-24 02:35:54           0
## 1000     Brazil 2016-06-03 21:43:21           1
```

```
# Checking the Data Dimensions
```

```
dim(ad_df)
```

```
## [1] 1000  10
```

The dataset has 1000 records and 10 columns

```
# Checking the Structure of the Dataset
```

```
str(ad_df)
```

```
## 'data.frame':  1000 obs. of  10 variables:
## $ Daily.Time.Spent.on.Site: num  69 80.2 69.5 74.2 68.4 ...
## $ Age                     : int  35 31 26 29 35 23 33 48 30 20 ...
## $ Area.Income             : num  61834 68442 59786 54806 73890 ...
## $ Daily.Internet.Usage     : num  256 194 236 246 226 ...
## $ Ad.Topic.Line           : chr  "Cloned 5thgeneration orchestration" "Monitored national standardi
## $ City                     : chr  "Wrightburgh" "West Jodi" "Davidton" "West Terrifurt" ...
## $ Male                     : int  0 1 0 1 0 1 0 1 1 1 ...
## $ Country                  : chr  "Tunisia" "Nauru" "San Marino" "Italy" ...
## $ Timestamp                : chr  "2016-03-27 00:53:11" "2016-04-04 01:39:02" "2016-03-13 20:35:42"
## $ Clicked.on.Ad            : int  0 0 0 0 0 0 0 1 0 0 ...
```

```
# Checking The Data present in each column
```

```
glimpse(ad_df)
```

```
## Rows: 1,000
## Columns: 10
## $ Daily.Time.Spent.on.Site <dbl> 68.95, 80.23, 69.47, 74.15, 68.37, 59.99, ...
## $ Age                     <int> 35, 31, 26, 29, 35, 23, 33, 48, 30, 20, 49...
## $ Area.Income             <dbl> 61833.90, 68441.85, 59785.94, 54806.18, 73...
```

```
## $ Daily.Internet.Usage      <dbl> 256.09, 193.77, 236.50, 245.89, 225.58, 22...
## $ Ad.Topic.Line            <chr> "Cloned 5thgeneration orchestration", "Mon...
## $ City                     <chr> "Wrightburgh", "West Jodi", "Davidton", "W...
## $ Male                     <int> 0, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, ...
## $ Country                  <chr> "Tunisia", "Nauru", "San Marino", "Italy",...
## $ Timestamp                <chr> "2016-03-27 00:53:11", "2016-04-04 01:39:0...
## $ Clicked.on.Ad            <int> 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, ...
```

7. Data Preparation

Uniformity

```
# Check column names
```

```
colnames(ad_df)
```

```
## [1] "Daily.Time.Spent.on.Site" "Age"
## [3] "Area.Income"             "Daily.Internet.Usage"
## [5] "Ad.Topic.Line"           "City"
## [7] "Male"                    "Country"
## [9] "Timestamp"               "Clicked.on.Ad"
```

```
# Renaming column names
```

```
names(ad_df)[1] <- "daily_time_spent_on_site"
names(ad_df)[2] <- "age"
names(ad_df)[3] <- "area_income"
names(ad_df)[4] <- "daily_internet_usage"
names(ad_df)[5] <- "ad_topic_line"
names(ad_df)[6] <- "city"
names(ad_df)[7] <- "male"
names(ad_df)[8] <- "country"
names(ad_df)[9] <- "timestamp"
names(ad_df)[10] <- "clicked_on_ad"
```

```
# Checking whether the column names have been changed
```

```
colnames(ad_df)
```

We'll rename the column names for Uniformity purposes

```
## [1] "daily_time_spent_on_site" "age"
## [3] "area_income"             "daily_internet_usage"
## [5] "ad_topic_line"           "city"
## [7] "male"                    "country"
## [9] "timestamp"               "clicked_on_ad"
```

```
# Checking for the length of unique values in each column
```

```
lapply(ad_df, function (x) {length(unique(x))})
```

```
## $daily_time_spent_on_site
## [1] 900
##
## $age
## [1] 43
##
## $area_income
## [1] 1000
##
## $daily_internet_usage
## [1] 966
##
## $ad_topic_line
## [1] 1000
##
## $city
## [1] 969
##
## $male
## [1] 2
##
## $country
## [1] 237
##
## $timestamp
## [1] 1000
##
## $clicked_on_ad
## [1] 2
```

We can observe that the 'Male' and 'Clicked_on_ad' columns are categorical since they only have 2 factor variables

Appropriateness

```
# Converting timestamp column to datetime datatype
```

```
ad_df[["timestamp"]] <- as.POSIXct(ad_df$timestamp, tz=Sys.timezone())
str(ad_df)
```

```
## 'data.frame': 1000 obs. of 10 variables:
## $ daily_time_spent_on_site: num 69 80.2 69.5 74.2 68.4 ...
## $ age : int 35 31 26 29 35 23 33 48 30 20 ...
## $ area_income : num 61834 68442 59786 54806 73890 ...
## $ daily_internet_usage : num 256 194 236 246 226 ...
## $ ad_topic_line : chr "Cloned 5thgeneration orchestration" "Monitored national standardi
```

```
## $ city : chr "Wrightburgh" "West Jodi" "Davidton" "West Terrifurt" ...
## $ male : int 0 1 0 1 0 1 0 1 1 1 ...
## $ country : chr "Tunisia" "Nauru" "San Marino" "Italy" ...
## $ timestamp : POSIXct, format: "2016-03-27 00:53:11" "2016-04-04 01:39:02" ...
## $ clicked_on_ad : int 0 0 0 0 0 0 0 1 0 0 ...
```

```
glimpse(ad_df)
```

```
## Rows: 1,000
## Columns: 10
## $ daily_time_spent_on_site <dbl> 68.95, 80.23, 69.47, 74.15, 68.37, 59.99, ...
## $ age <int> 35, 31, 26, 29, 35, 23, 33, 48, 30, 20, 49...
## $ area_income <dbl> 61833.90, 68441.85, 59785.94, 54806.18, 73...
## $ daily_internet_usage <dbl> 256.09, 193.77, 236.50, 245.89, 225.58, 22...
## $ ad_topic_line <chr> "Cloned 5thgeneration orchestration", "Mon...
## $ city <chr> "Wrightburgh", "West Jodi", "Davidton", "W...
## $ male <int> 0, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, ...
## $ country <chr> "Tunisia", "Nauru", "San Marino", "Italy",...
## $ timestamp <dtm> 2016-03-27 00:53:11, 2016-04-04 01:39:02,...
## $ clicked_on_ad <int> 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, ...
```

We can observe that the change has taken shape successfully. We now want to split the column to date and time

```
# Splitting datetime into date and time
```

```
Time <- format(as.POSIXct(strptime(ad_df$timestamp,"%Y-%m-%d %H:%M:%S",tz="")),format = "%H:%M:%S")
head(Time)
```

```
## [1] "00:53:11" "01:39:02" "20:35:42" "02:31:19" "03:36:18" "14:30:17"
```

```
Dates <- format(as.POSIXct(strptime(ad_df$timestamp,"%Y-%m-%d %H:%M:%S",tz="")),format = "%Y-%m-%d")
head(Dates)
```

```
## [1] "2016-03-27" "2016-04-04" "2016-03-13" "2016-01-10" "2016-06-03"
```

```
## [6] "2016-05-19"
```

```
ad_df$Dates <- Dates
ad_df$Time <- Time
```

```
str(ad_df)
```

```
## 'data.frame': 1000 obs. of 12 variables:
## $ daily_time_spent_on_site: num 69 80.2 69.5 74.2 68.4 ...
## $ age : int 35 31 26 29 35 23 33 48 30 20 ...
## $ area_income : num 61834 68442 59786 54806 73890 ...
## $ daily_internet_usage : num 256 194 236 246 226 ...
## $ ad_topic_line : chr "Cloned 5thgeneration orchestration" "Monitored national standardi
## $ city : chr "Wrightburgh" "West Jodi" "Davidton" "West Terrifurt" ...
```

```
## $ male : int 0 1 0 1 0 1 0 1 1 1 ...
## $ country : chr "Tunisia" "Nauru" "San Marino" "Italy" ...
## $ timestamp : POSIXct, format: "2016-03-27 00:53:11" "2016-04-04 01:39:02" ...
## $ clicked_on_ad : int 0 0 0 0 0 0 0 1 0 0 ...
## $ Dates : chr "2016-03-27" "2016-04-04" "2016-03-13" "2016-01-10" ...
## $ Time : chr "00:53:11" "01:39:02" "20:35:42" "02:31:19" ...
```

```
# Separating dates to hours minutes and days and dropping the timestamp column
```

```
ad_df <- separate(ad_df, "Dates", c("year", "month", "day"), sep = "-")
ad_df <- separate(ad_df, "Time", c("hour", "minutes", "seconds"), sep = ":")

colnames(ad_df)
```

```
## [1] "daily_time_spent_on_site" "age"
## [3] "area_income" "daily_internet_usage"
## [5] "ad_topic_line" "city"
## [7] "male" "country"
## [9] "timestamp" "clicked_on_ad"
## [11] "year" "month"
## [13] "day" "hour"
## [15] "minutes" "seconds"
```

```
# Changing the new derived columns to factors for ease of analysis
```

```
ad_df$Male = factor(ad_df$male)
ad_df$Year = factor(ad_df$year)
ad_df$Month = factor(ad_df$month)
ad_df$Day = factor(ad_df$day)
ad_df$Hour = factor(ad_df$hour)
ad_df$Minutes = factor(ad_df$minutes)
ad_df$Seconds = factor(ad_df$seconds)
ad_df$clicked_on_ad = factor(ad_df$clicked_on_ad)
```

We can see that the date and time have their respective columns

```
#——- ## Completeness
```

```
# Checking for missing values
```

```
colSums(is.na(ad_df))
```

```
## daily_time_spent_on_site      age      area_income
##                0                0                0
##      daily_internet_usage      ad_topic_line      city
##                0                0                0
##                male      country      timestamp
##                0                0                0
##      clicked_on_ad      year      month
```



```
##           0           0           0
##           day           hour           minutes
##           0           0           0
##           seconds           Male           Year
##           0           0           0
##           Month           Day           Hour
##           0           0           0
##           Minutes           Seconds
##           0           0
```

Our data is complete hence no missing values

```
#——- ## Consistency
```

```
# Checking for duplicate values
```

```
duplicates <- ad_df[duplicated(ad_df),]
duplicates
```

```
## [1] daily_time_spent_on_site age area_income
## [4] daily_internet_usage ad_topic_line city
## [7] male country timestamp
## [10] clicked_on_ad year month
## [13] day hour minutes
## [16] seconds Male Year
## [19] Month Day Hour
## [22] Minutes Seconds
## <0 rows> (or 0-length row.names)
```

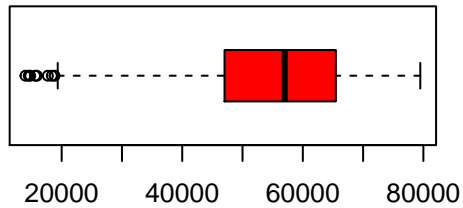
Our data is consistent due to no duplicate values present #——- ### Anomaly Detection ####
 # Checking for anomalies in our numerical variables i.e daily_time_spent_on_site, area income, age, and daily_internet usage

Boxplots

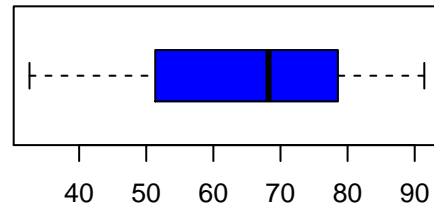
```
# Plotting boxplots for all the numerical variables
```

```
par(mfrow=c(2,2))
boxplot((ad_df$`area_income`), horizontal = TRUE, col = 'red', main = "boxplot of area income")
boxplot((ad_df$`daily_time_spent_on_site`), horizontal = TRUE, col = 'blue', main = "boxplot of daily t")
boxplot((ad_df$`age`), horizontal = TRUE, col = 'yellow', main = "boxplot of age")
boxplot((ad_df$`daily_internet_usage`), horizontal = TRUE, col = 'green', main = "boxplot of daily inte")
```

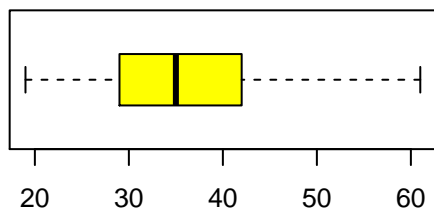
boxplot of area income



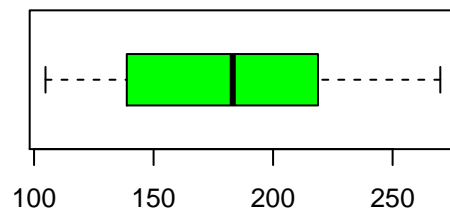
boxplot of daily time spent on site



boxplot of age



boxplot of daily internet usage



1. Area income variable has values ranging from below 0 to 80,000. We have a few values below 20,000 which are outliers but we'll keep them because they represent crucial data for analysis
2. Daily time spent on site has values from around 20 to 90 with the mode between 50 to 80
3. Age variable has observations from the age of 20 to 60 with the mode between 30 to 40
4. Daily internet usage has values from 100 to slightly above 250 with the mode between 150 to 200

8. Exploratory Data Analysis

Univariate Analysis

```
# Checking the statistical summary of the data
```

```
summary(ad_df)
```

```
##  daily_time_spent_on_site    age    area_income    daily_internet_usage
##  Min.   :32.60             Min.   :19.00   Min.   :13996   Min.   :104.8
##  1st Qu.:51.36             1st Qu.:29.00   1st Qu.:47032   1st Qu.:138.8
```

```

## Median :68.22          Median :35.00   Median :57012   Median :183.1
## Mean    :65.00          Mean    :36.01   Mean    :55000   Mean    :180.0
## 3rd Qu. :78.55          3rd Qu. :42.00   3rd Qu. :65471   3rd Qu. :218.8
## Max.    :91.43          Max.    :61.00   Max.    :79485   Max.    :270.0
##
## ad_topic_line          city              male              country
## Length:1000           Length:1000       Min.    :0.000     Length:1000
## Class :character       Class :character   1st Qu.:0.000     Class :character
## Mode  :character       Mode  :character   Median :0.000     Mode  :character
##                               Mean    :0.481
##                               3rd Qu.:1.000
##                               Max.    :1.000
##
## timestamp              clicked_on_ad   year
## Min.    :2016-01-01 02:52:10  0:500      Length:1000
## 1st Qu. :2016-02-18 02:55:42  1:500      Class :character
## Median  :2016-04-07 17:27:29      Mode  :character
## Mean    :2016-04-10 10:34:06
## 3rd Qu. :2016-05-31 03:18:14
## Max.    :2016-07-24 00:22:16
##
## month                  day              hour              minutes
## Length:1000           Length:1000       Length:1000       Length:1000
## Class :character       Class :character   Class :character   Class :character
## Mode  :character       Mode  :character   Mode  :character   Mode  :character
##
##
##
## seconds                Male      Year      Month      Day      Hour
## Length:1000            0:519    2016:1000  01:147    03      : 46    07      : 54
## Class :character       1:481
## Mode  :character
##                               02:160    17      : 42    20      : 50
##                               03:156    15      : 41    09      : 49
##                               04:147    10      : 37    21      : 48
##                               05:147    04      : 36    00      : 45
##                               06:142    26      : 36    05      : 44
##                               07:101    (Other):762  (Other):710
##
## Minutes                Seconds
## 02      : 26    22      : 28
## 07      : 24    10      : 27
## 13      : 24    35      : 27
## 10      : 22    37      : 27
## 21      : 21    38      : 24
## 33      : 21    15      : 23
## (Other):862    (Other):844

```

The timestamp has a conflicting datatype compared to what its normal date/time format as well as gender and and clicked on ad datatypes which should be categorical instead of integers

The daily time spent on the site seems to be in minutes and seconds ranging from 32.60 to 91.43. The values are likely to be close to normally distributed as the median is 68.22 and the mean is 65.

The area income are not likely to be close to normally distributed due to a large difference in ranges i.e from 13996 to 79485, with a median of 57012 and a mean of 55000.

The daily internet usage ranges from 104.8 to 270.0, with a median of 183.1 and a mean of 180.0. The values are likely to be close to normally distributed.

The ad topic line, City, male, Country are categorical features, with a different value for each record.

The feature male is categorical (binary) with a mean of 0.481, which means there are more records from individuals that are female.

The clicked on ad variable is categorical (binary) with a mean of 0.5, which means that the variable of interest is balanced in this dataset.

#——

Measures of Central Tendancy and Dispersion - Summary

Central Tendancy - Mode, Mean and Median

```
# First, a function for mode will be created since R does not have a built in function.
```

```
getmode <- function(v) {  
  uniqv <- unique(v)  
  uniqv[which.max(tabulate(match(v, uniqv)))]  
}
```

```
# City  
# This column represents the city where the most users are from  
mode.city <- getmode(ad_df$city)  
mode.city
```

```
## [1] "Lisamouth"
```

```
# Country  
# This column represents the country where the most users are from  
mode.country <- getmode(ad_df$country)  
mode.country
```

```
## [1] "Czech Republic"
```

```
# Age
# This column represents the Age That most users are, its mean and median
mode.age <- getmode(ad_df$age)
mode.age
```

```
## [1] 31
```

```
mean(ad_df$age)
```

```
## [1] 36.009
```

```
median(ad_df$age)
```

```
## [1] 35
```

```
# Daily Internet Usage
# This column represents the daily internet usage for most users, its mean and median
mode.usage <- getmode(ad_df$daily_internet_usage)
mode.usage
```

```
## [1] 167.22
```

```
mean(ad_df$daily_internet_usage)
```

```
## [1] 180.0001
```

```
median(ad_df$daily_internet_usage)
```

```
## [1] 183.13
```

```
# Area Income
# This column represents most of the Area Income
mode.income <- getmode(ad_df$area_income)
mode.income
```

```
## [1] 61833.9
```

```
mean(ad_df$area_income)
```

```
## [1] 55000
```

```
median(ad_df$area_income)
```

```
## [1] 57012.3
```

```
# Male
# This column represents gender with the most users
mode.male <- getmode(ad_df$male)
mode.male
```

```
## [1] 0
```

```
# Ad_Topic_line
# This column represents most advertisement topic line

mode.adline<-getmode(ad_df$ad_topic_line)
mode.adline
```

```
## [1] "Cloned 5thgeneration orchestration"
```

```
# Daily_Time_Spent
# This column represents most frequent daily time spent on site

mode.time <- getmode(ad_df$daily_time_spent_on_site)
mode.time
```

```
## [1] 62.26
```

```
mean(ad_df$daily_time_spent_on_site)
```

```
## [1] 65.0002
```

```
median(ad_df$daily_time_spent_on_site)
```

```
## [1] 68.215
```

```
# Month
# This column represents most frequent months during usage

mode.month <- getmode(ad_df$month)
mode.month
```

```
## [1] "02"
```

```
# Day
# This column represents most frequent day during usage

mode.day <- getmode(ad_df$day)
mode.day
```

```
## [1] "03"
```

```
# Hour  
# This column represents most frequent hour during usage
```

```
mode.hour <- getmode(ad_df$hour)  
mode.hour
```

```
## [1] "07"
```

```
# Minute  
# This column represents most frequent Minutes during usage
```

```
mode.minutes <- getmode(ad_df$minutes)  
mode.minutes
```

```
## [1] "02"
```

```
# Seconds  
# This column represents most frequent months during usage
```

```
mode.seconds <- getmode(ad_df$seconds)  
mode.seconds
```

```
## [1] "22"
```

```
# Age
```

```
sd.age <- sd(ad_df$age)  
sd.age
```

Measure of Dispersion - Standard Deviation, Variance, Skewness, Kurtosis and Range

```
## [1] 8.785562
```

```
var.age <- var(ad_df$age)  
var.age
```

```
## [1] 77.18611
```

```
range.age <- range(ad_df$age)  
range.age
```

```
## [1] 19 61
```

```
skew.age <- skewness(ad_df$age)  
skew.age
```

```
## [1] 0.4784227
```

```
kurt.age <- kurtosis(ad_df$age)
kurt.age
```

```
## [1] 2.595482
```

```
# Daily Internet Usage
```

```
sd.daily_internet_usage <- sd(ad_df$daily_internet_usage)
sd.daily_internet_usage
```

```
## [1] 43.90234
```

```
var.daily_internet_usage <- var(ad_df$daily_internet_usage)
var.daily_internet_usage
```

```
## [1] 1927.415
```

```
range.daily_internet_usage <- range(ad_df$daily_internet_usage)
range.daily_internet_usage
```

```
## [1] 104.78 269.96
```

```
skew.daily_internet_usage <- skewness(ad_df$daily_internet_usage)
skew.daily_internet_usage
```

```
## [1] -0.03348703
```

```
kurt.daily_internet_usage <- kurtosis(ad_df$daily_internet_usage)
kurt.daily_internet_usage
```

```
## [1] 1.727701
```

```
# Daily time spent on site
```

```
sd.daily_time_spent_on_site <- sd(ad_df$daily_time_spent_on_site)
sd.daily_time_spent_on_site
```

```
## [1] 15.85361
```

```
var.daily_time_spent_on_site <- var(ad_df$daily_time_spent_on_site)
var.daily_time_spent_on_site
```

```
## [1] 251.3371
```

```
range.daily_time_spent_on_site <- range(ad_df$daily_time_spent_on_site)
range.daily_time_spent_on_site
```

```
## [1] 32.60 91.43
```



```
skew.daily_time_spent_on_site <- skewness(ad_df$daily_time_spent_on_site)
skew.daily_time_spent_on_site
```

```
## [1] -0.3712026
```

```
kurt.daily_time_spent_on_site <- kurtosis(ad_df$daily_time_spent_on_site)
kurt.daily_time_spent_on_site
```

```
## [1] 1.903942
```

```
# Area Income
```

```
sd.area_income <- sd(ad_df$area_income)
sd.area_income
```

```
## [1] 13414.63
```

```
var.area_income <- var(ad_df$area_income)
var.area_income
```

```
## [1] 179952406
```

```
range.area_income <- range(ad_df$area_income)
range.area_income
```

```
## [1] 13996.5 79484.8
```

```
skew.area_income <- skewness(ad_df$area_income)
skew.area_income
```

```
## [1] -0.6493967
```

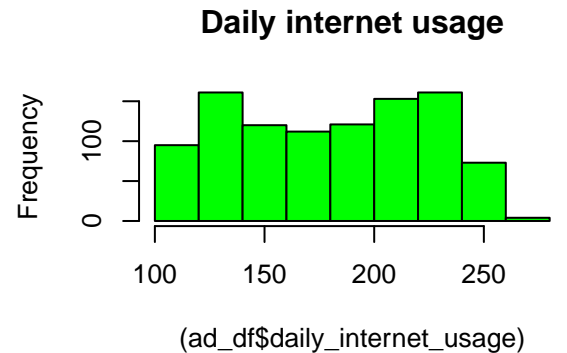
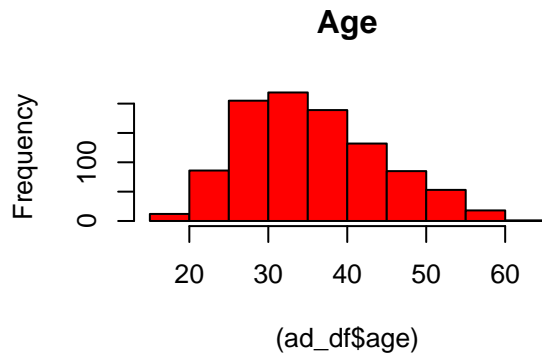
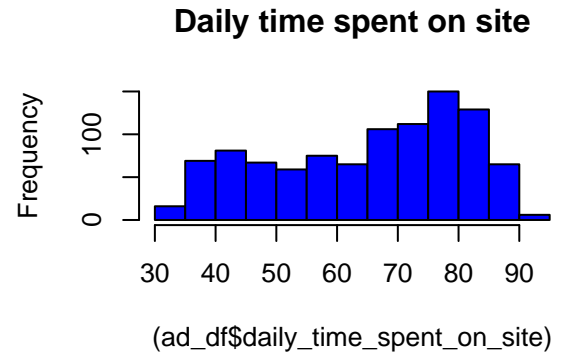
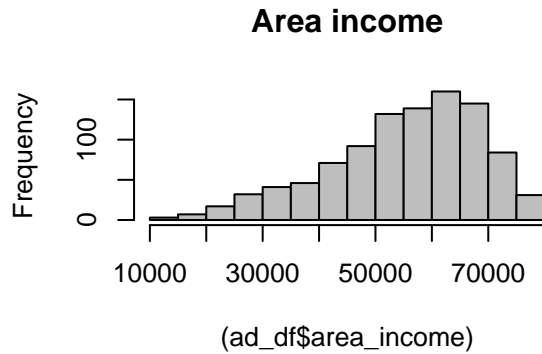
```
kurt.area_income <- kurtosis(ad_df$area_income)
kurt.area_income
```

```
## [1] 2.894694
```

```
# Plotting multiple histograms for Area income, Age, Daily time spent on site and Daily Internet Usage
```

```
par(mfrow=c(2,2))
```

```
hist((ad_df$`area_income`), col = 'grey', main = "Area income")
hist((ad_df$`daily_time_spent_on_site`), col = 'blue', main = "Daily time spent on site")
hist((ad_df$`age`), col = 'red', main = "Age")
hist((ad_df$`daily_internet_usage`), col = 'green', main = "Daily internet usage")
```



Histograms

Observations: ##### 1. Area income variable is negatively skewed as most of the observations recorded are lower compared to the high area income ##### 2. Age variable is positively skewed as most of the ages recorded are younger ##### 3. Daily internet usage and daily time spent on site are bimodal as they have an almost normal distribution

Bivariate Analysis

```
# Correlation Matrix
# Calling all the numerical data present

age<- ad_df$age
income<-ad_df$area_income
time<-ad_df$daily_time_spent_on_site
usage<-ad_df$daily_internet_usage

# Creating a new dataframe num with numerical data variables

num_data <- data.frame(age, income, time, usage)
head(num_data)
```

Correlation

```
##   age   income   time   usage
## 1  35 61833.90  68.95 256.09
## 2  31 68441.85  80.23 193.77
## 3  26 59785.94  69.47 236.50
```

```
## 4 29 54806.18 74.15 245.89
## 5 35 73889.99 68.37 225.58
## 6 23 59761.56 59.99 226.74
```

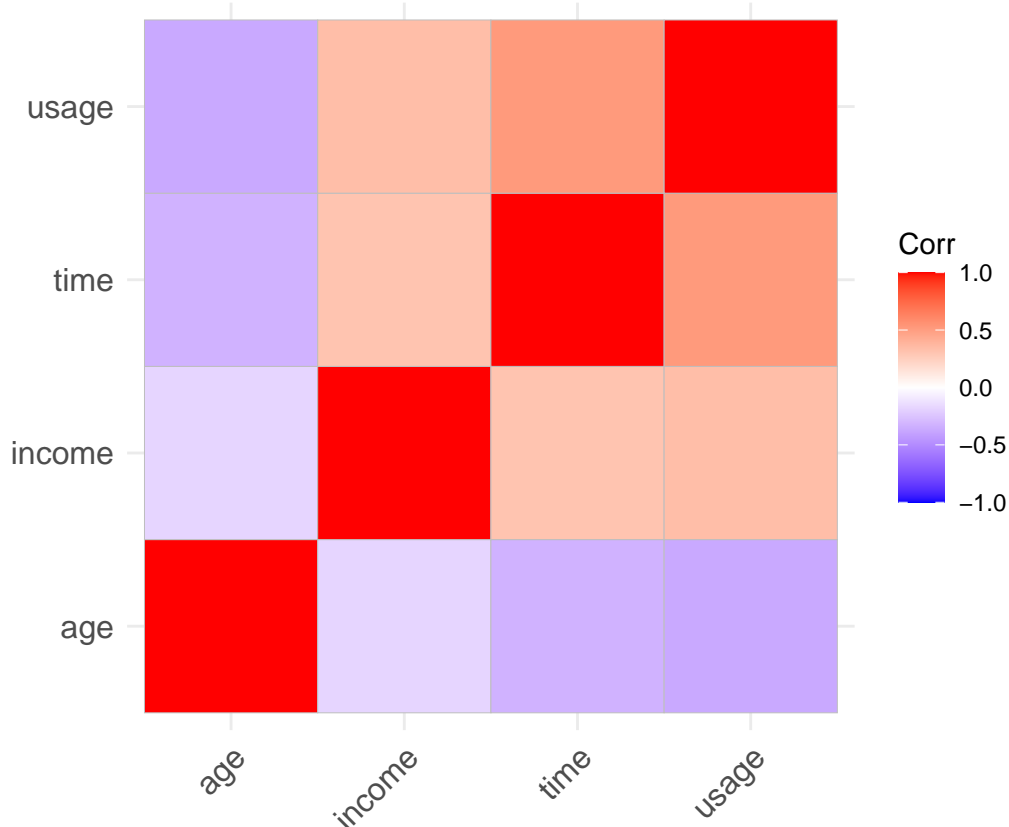
```
# Calculating the correlation matrix
```

```
corr <- cor(num_data)
head(corr)
```

```
##           age      income      time      usage
## age      1.0000000 -0.1826050 -0.3315133 -0.3672086
## income  -0.1826050  1.0000000  0.3109544  0.3374955
## time    -0.3315133  0.3109544  1.0000000  0.5186585
## usage   -0.3672086  0.3374955  0.5186585  1.0000000
```

```
# Plotting the correlation matrix
```

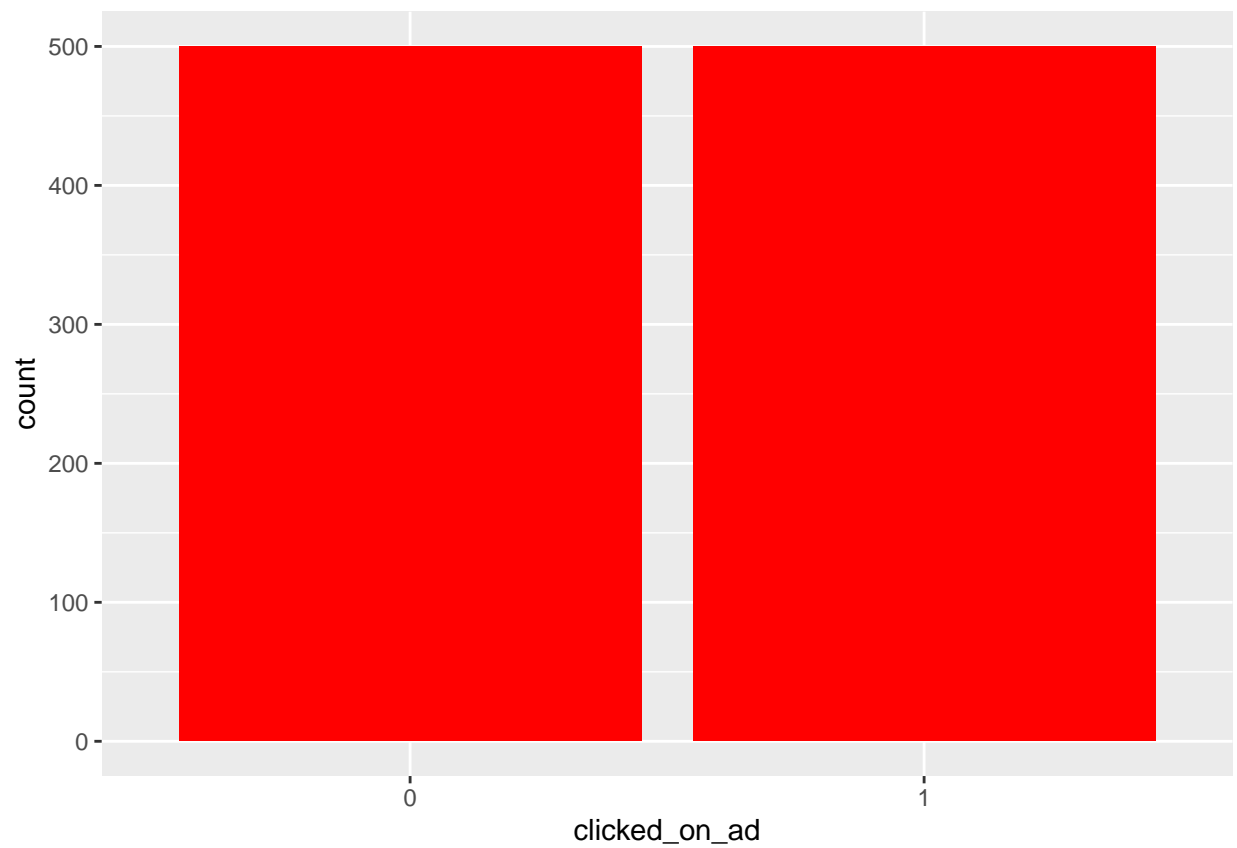
```
ggcorrplot(corr, hc.order = TRUE)
```



```
##### Observations #####
1. Daily_internet_usage and Daily_time_spent_on_site seem to have a moderate positive correlation #####
2. Daily_internet_usage and Age seem to have a negative correlation #####
3. Area Income and Age are weakly correlated
```

```
# Finding out and previewing the Number of clicked and no clicked ads
```

```
ggplot(ad_df, aes(clicked_on_ad)) + geom_bar(fill = "red")
```

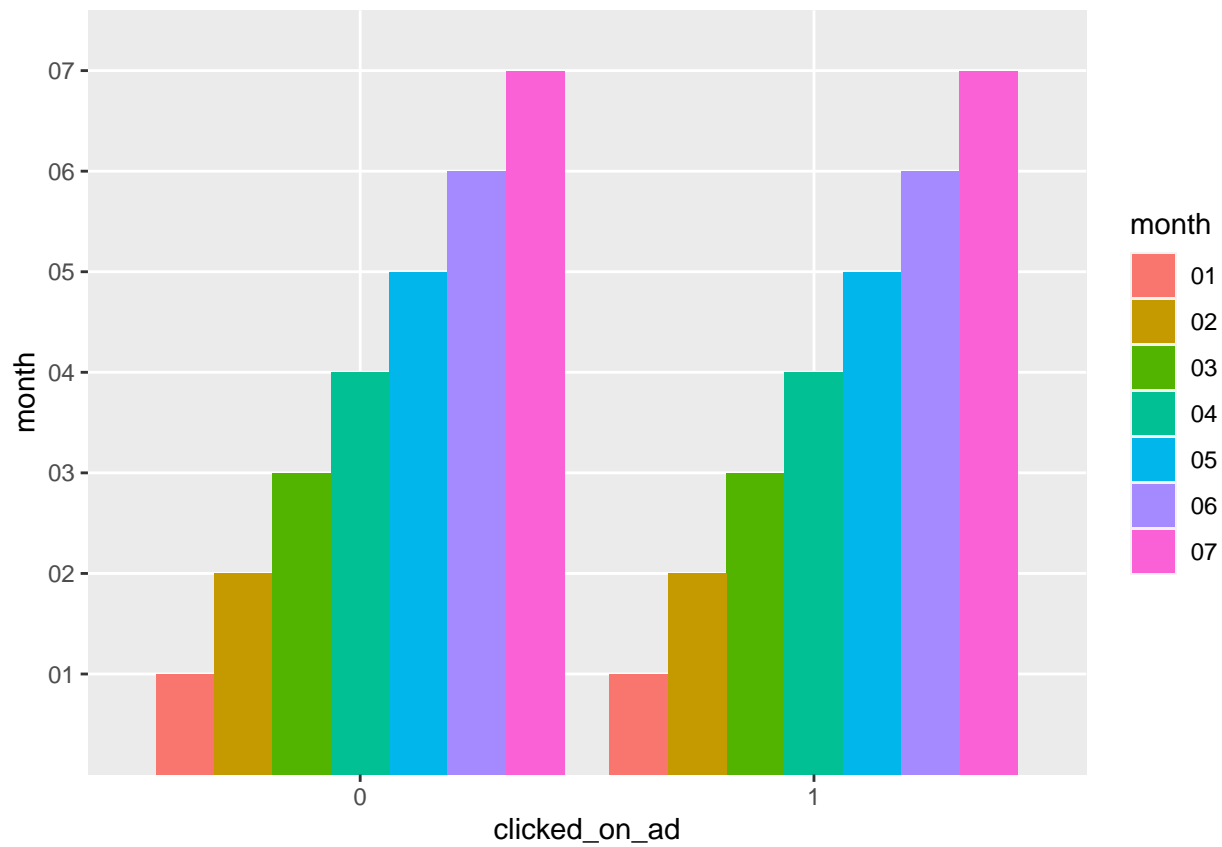


Barplots

The clicked ads and no clicked ads in our dataset were equal

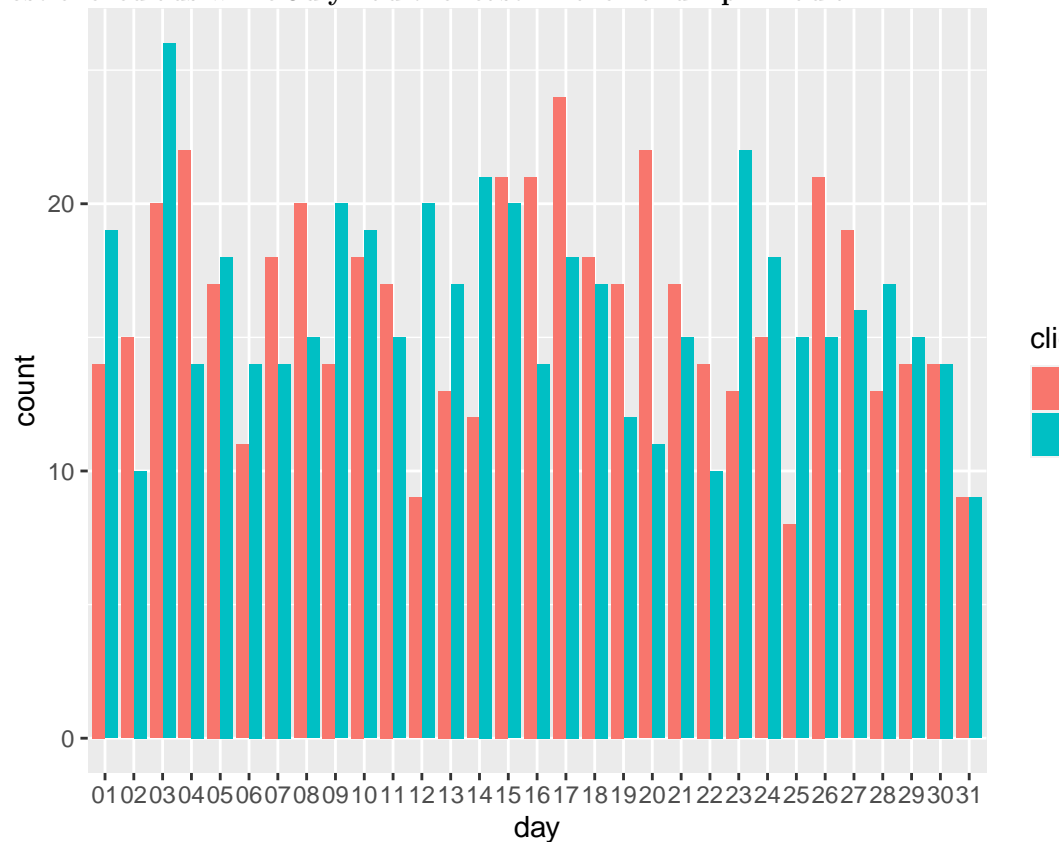
Finding out and previewing the month with the most clicked ads

```
ggplot(ad_df, aes(x = 'clicked_on_ad', y = 'month')) + geom_col(aes(fill = 'month'), position = "dodge")
```



```
# Finding out and previewing the day with the most clicked ads  
ggplot(data = ad_df) +  
  geom_bar(mapping = aes(x = day, fill = clicked_on_ad), position = "dodge")
```

February and May had the most clicked ads while July had the least. March and April had an



equal number of clicked ads.

The most activity recorded is in the first 3 months, from both who clicked the ads and those who didn't. ##### January (1), March (3) and July (7) had more activity from those who did not click on the ads as compared to those who clicked on the ads. ##### Months February (2), April (4) and May (5) had more people who clicked on the ads as compared to those who did not click on the ads ##### June (6) had an equal number of people who clicked on the ads and those who did Not ##### We observe that at around mid month we had more people who were not clicking on the ads as compared to the beginning and the end of the month

Finding out and previewing the hours with the most clicked ads

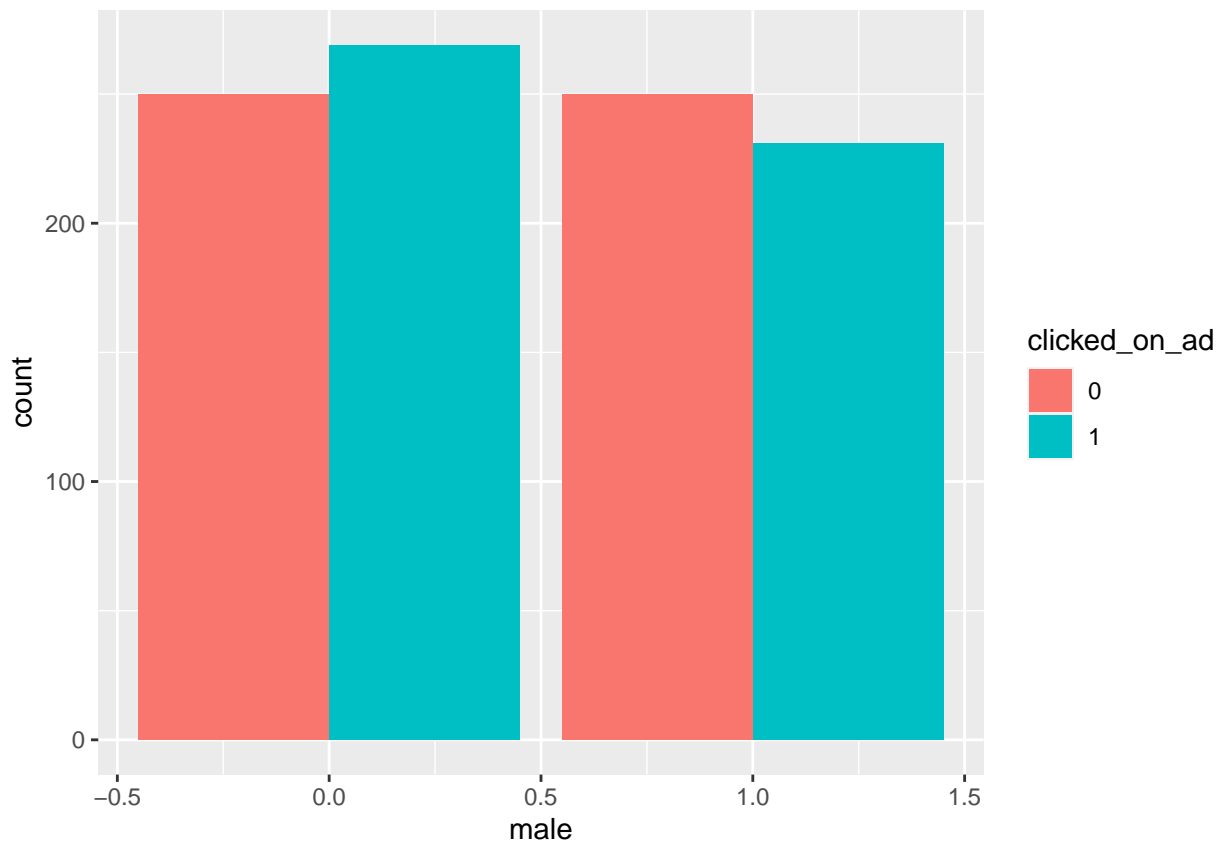
```
ggplot(data = ad_df) +
  geom_bar(mapping = aes(x = hour, fill = clicked_on_ad), position = "dodge")
```



Observations - From around 8 pm to 11 pm, we have more people not clicking on ads as compared to those who clicked on the ads before 8 pm and a little after Midnight. 3, 6, 9 and 11 am are the morning hours with the most clicked ads while 3,5 and 6 pm are the hours with the most clicks on the ads in the evening.

Finding out and previewing the gender with the most clicked ads

```
ggplot(data = ad_df) +
  geom_bar(mapping = aes(x = male, fill = clicked_on_ad), position = "dodge")
```

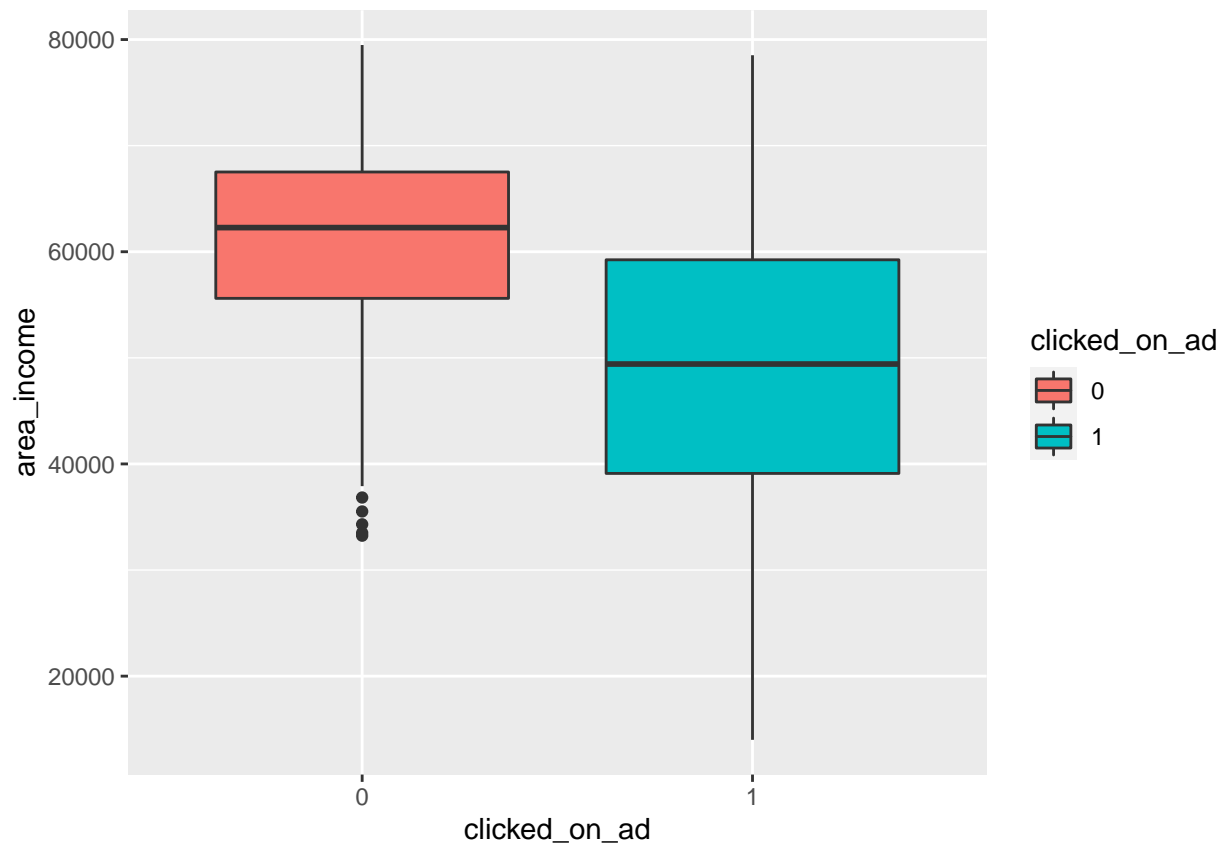


Observations - We have more number of females who clicked on the ads as compared to those who did not. Most males did not click on the ads.

Boxplots

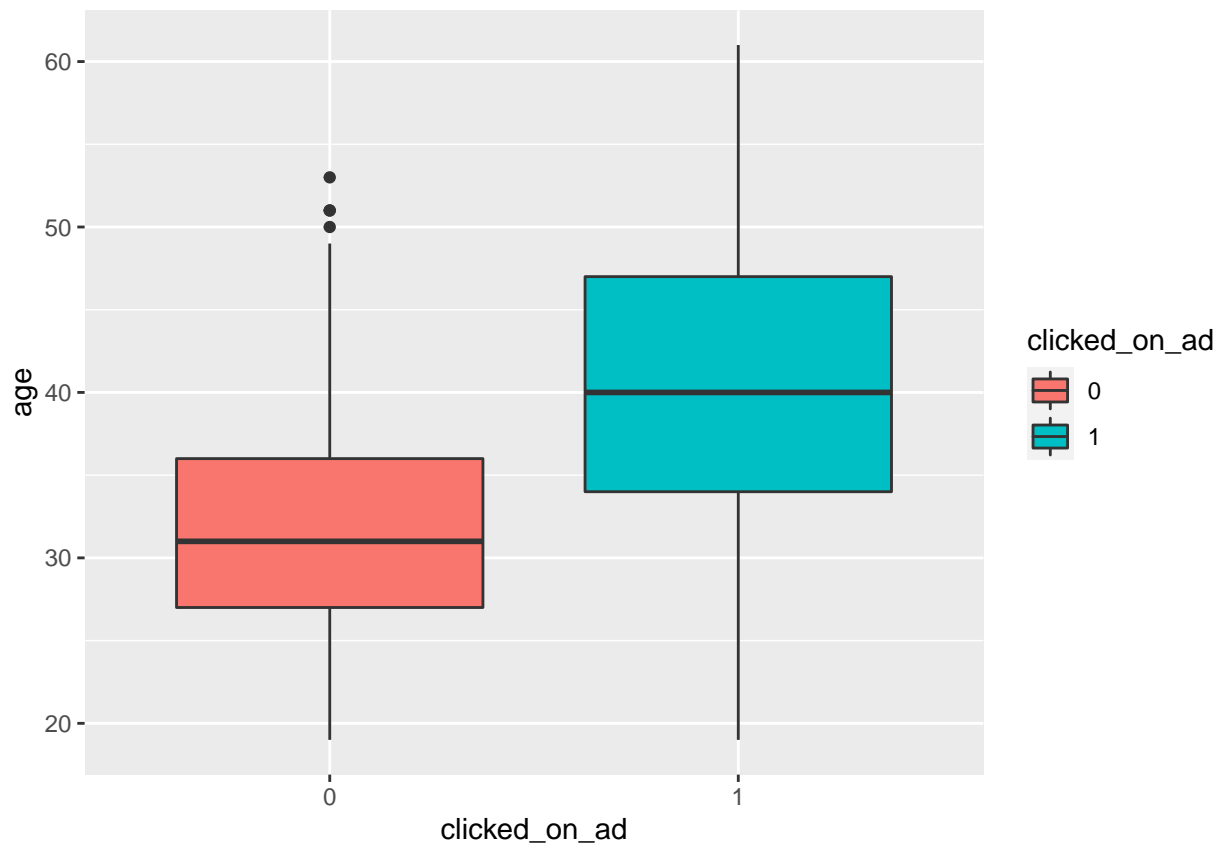
```
# Area Income vs Number of ad clicks
# Finding out and previewing boxplots to show how the area income relates with the number of clicks

ggplot(data = ad_df, mapping = aes( x = area_income, y = clicked_on_ad, fill = clicked_on_ad)) +
  geom_boxplot() +
  coord_flip()
```

Most people who clicked on the ads have a lower income as compared to those who did Not click on the ads

```
# Age vs Number of ad clicks  
# Finding out and previewing boxplots to show how the age relates with the number of clicks  
  
ggplot(data = ad_df, mapping = aes( x = age, y = clicked_on_ad, fill = clicked_on_ad)) +  
  geom_boxplot() +  
  coord_flip()
```

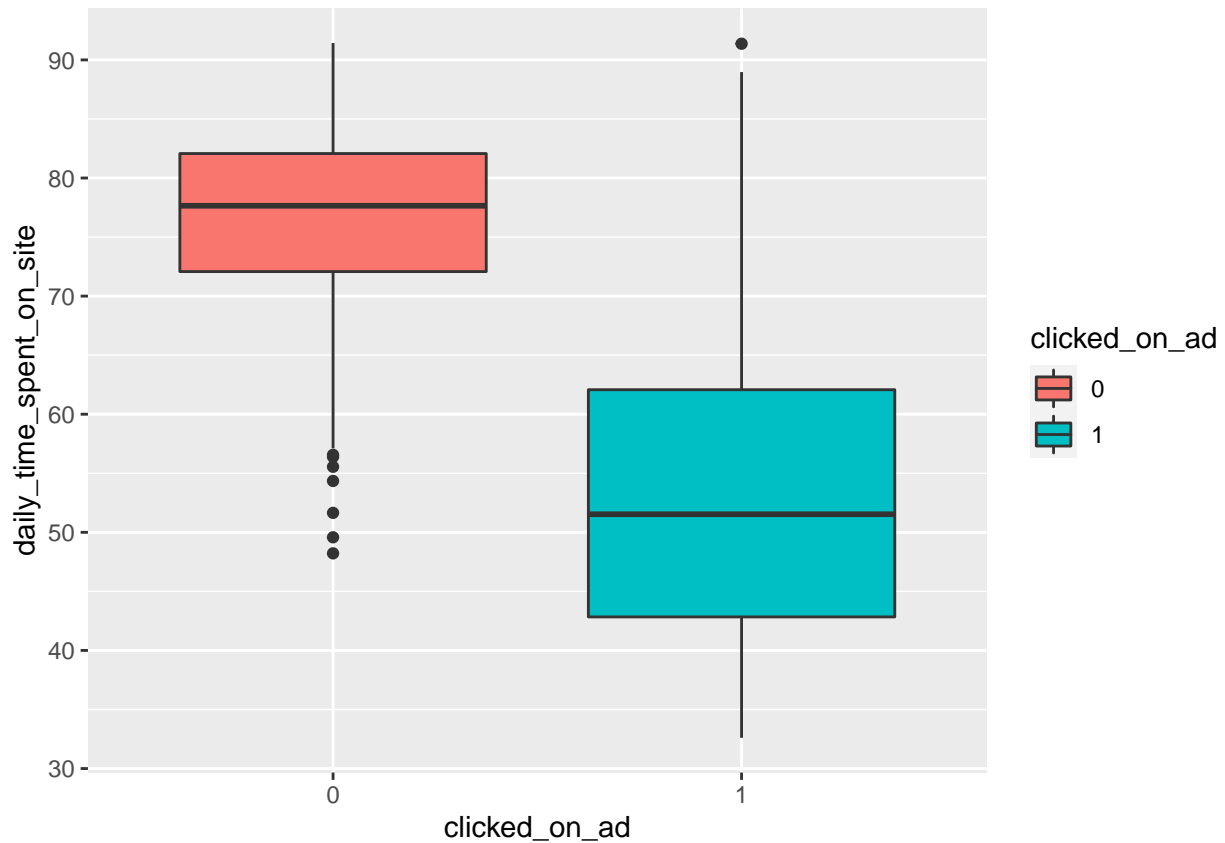


Most people who clicked on the ads were older than those who did NOT click on the ads

Daily Time spent on site vs Number of ad clicks

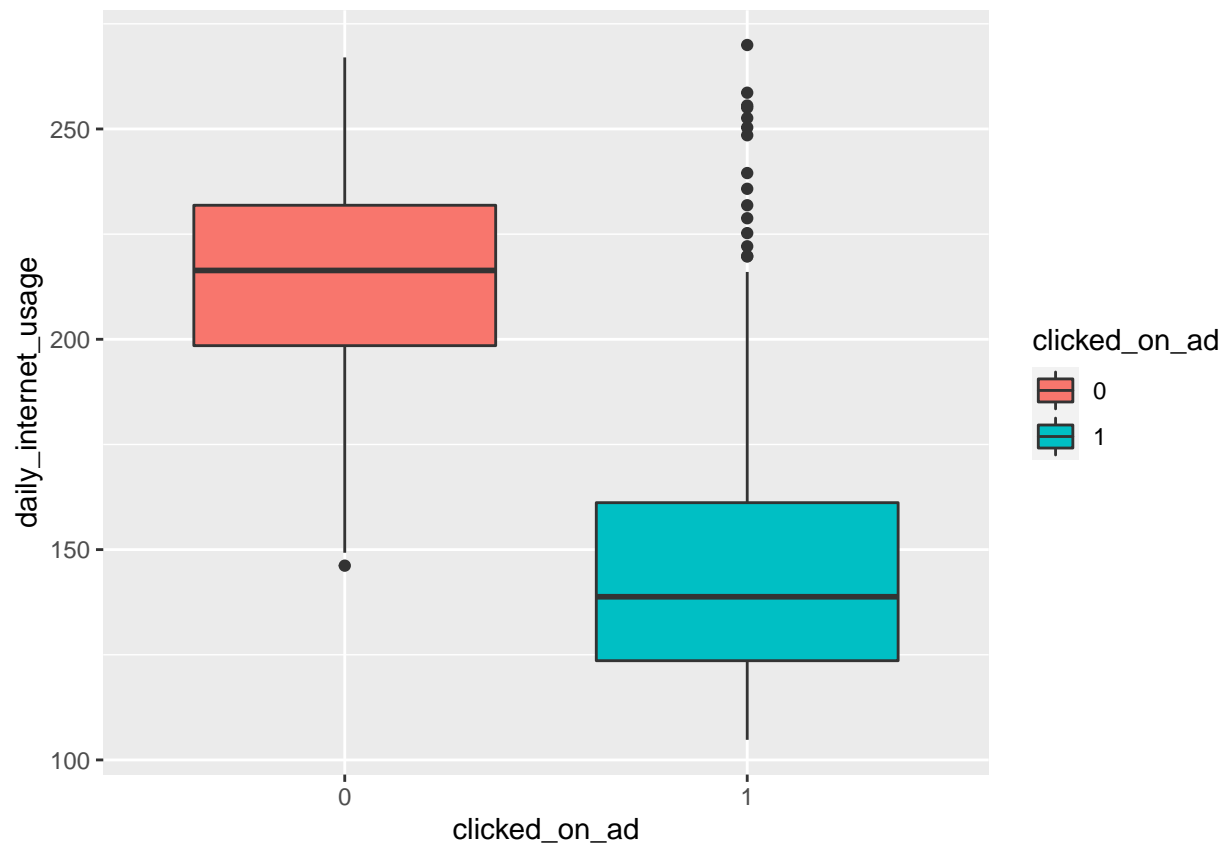
Finding out and previewing boxplots to show how the daily time spent on site relates with the number

```
ggplot(data = ad_df, mapping = aes( x = daily_time_spent_on_site, y = clicked_on_ad, fill = clicked_on_ad)) +
  geom_boxplot() +
  coord_flip()
```



Most people who clicked on the ads spent way less time on the site as compared to thos who did not click on the ads

```
# Daily internet usage vs Number of ad clicks  
# Finding out and previewing boxplots to show how the daily internet usage relates with the number of c  
  
ggplot(data = ad_df, mapping = aes( x = daily_internet_usage, y = clicked_on_ad, fill = clicked_on_ad))  
  geom_boxplot() +  
  coord_flip()
```

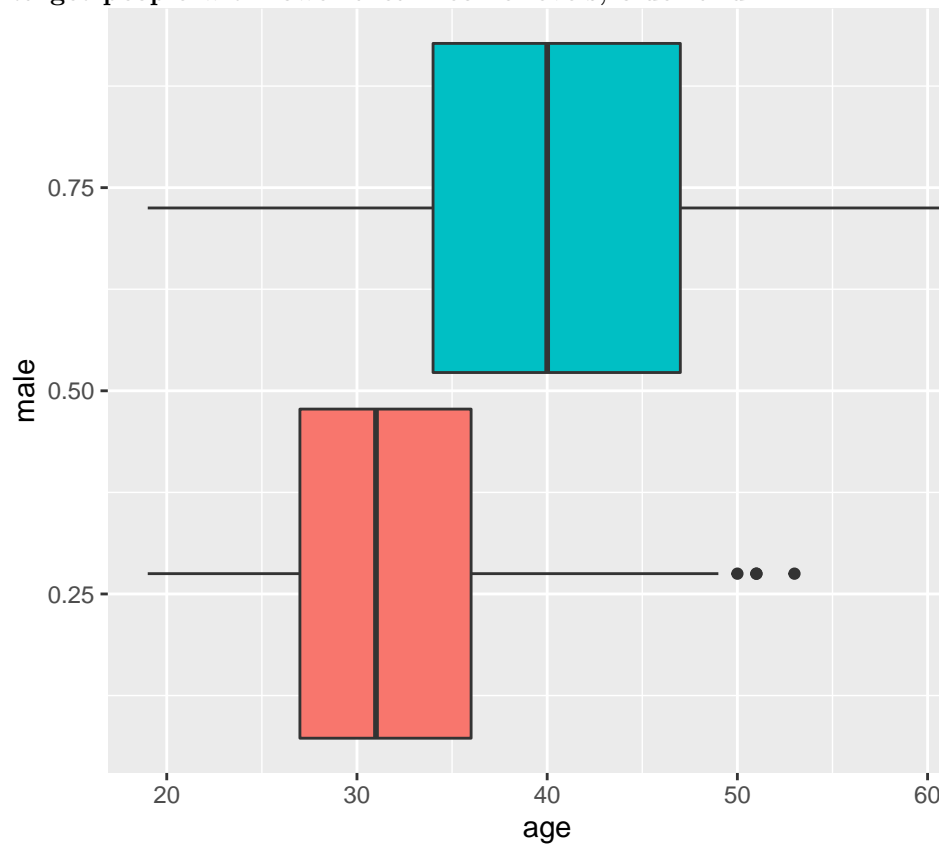


The daily internet usage of most people who clicked on the ads is way less than those who did NOT click on the ads

```
# Age vs Gender
# Finding out and previewing boxplots to show how the Age relates with the gender

ggplot(data = ad_df, mapping = aes( x = male, y = age, fill = clicked_on_ad)) +
  geom_boxplot() +
  coord_flip()
```

Conclusion - The entrepreneur should target people with lower area income levels, older and

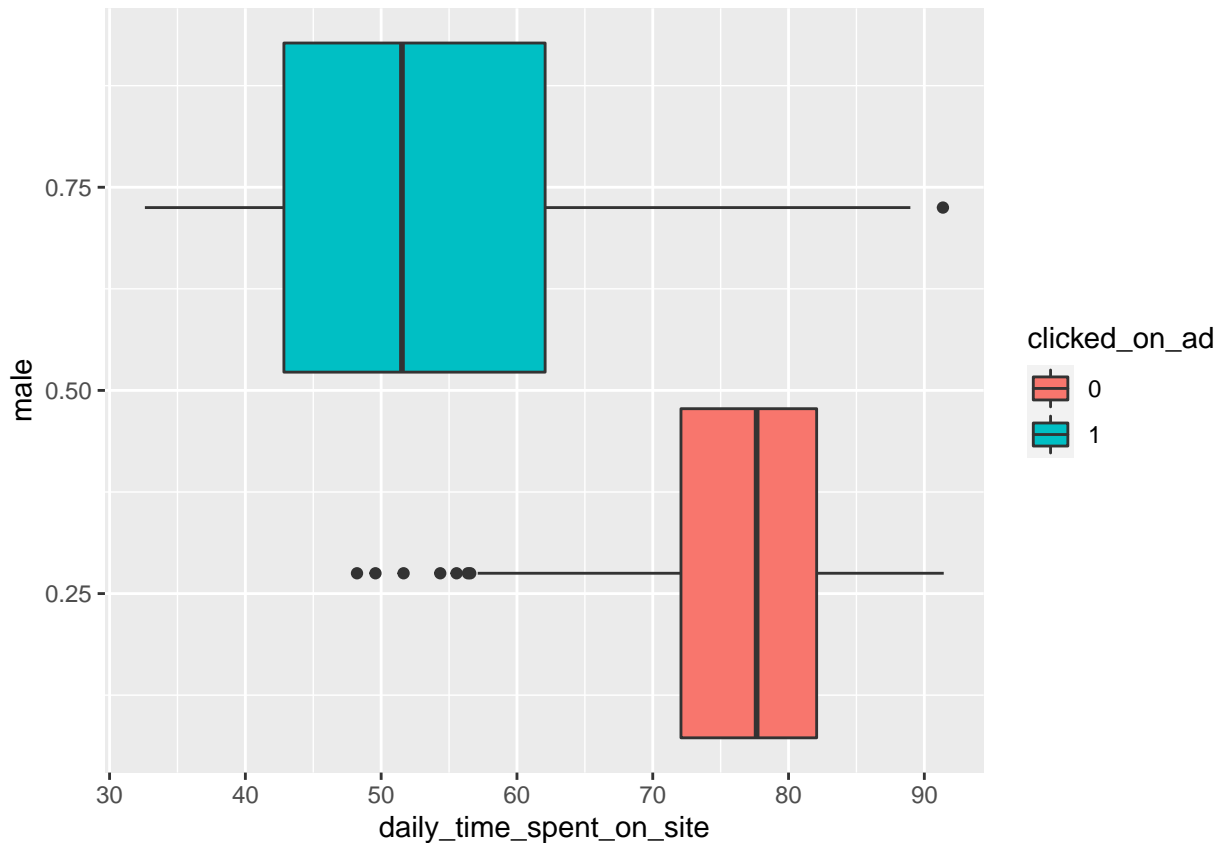


those who spend less time on the site.

Generally, those who clicked on the ads were older, but the males were slightly older than the females

```
# Daily time spent on site vs Gender
# Finding out and previewing boxplots to show how the Age relates with the gender

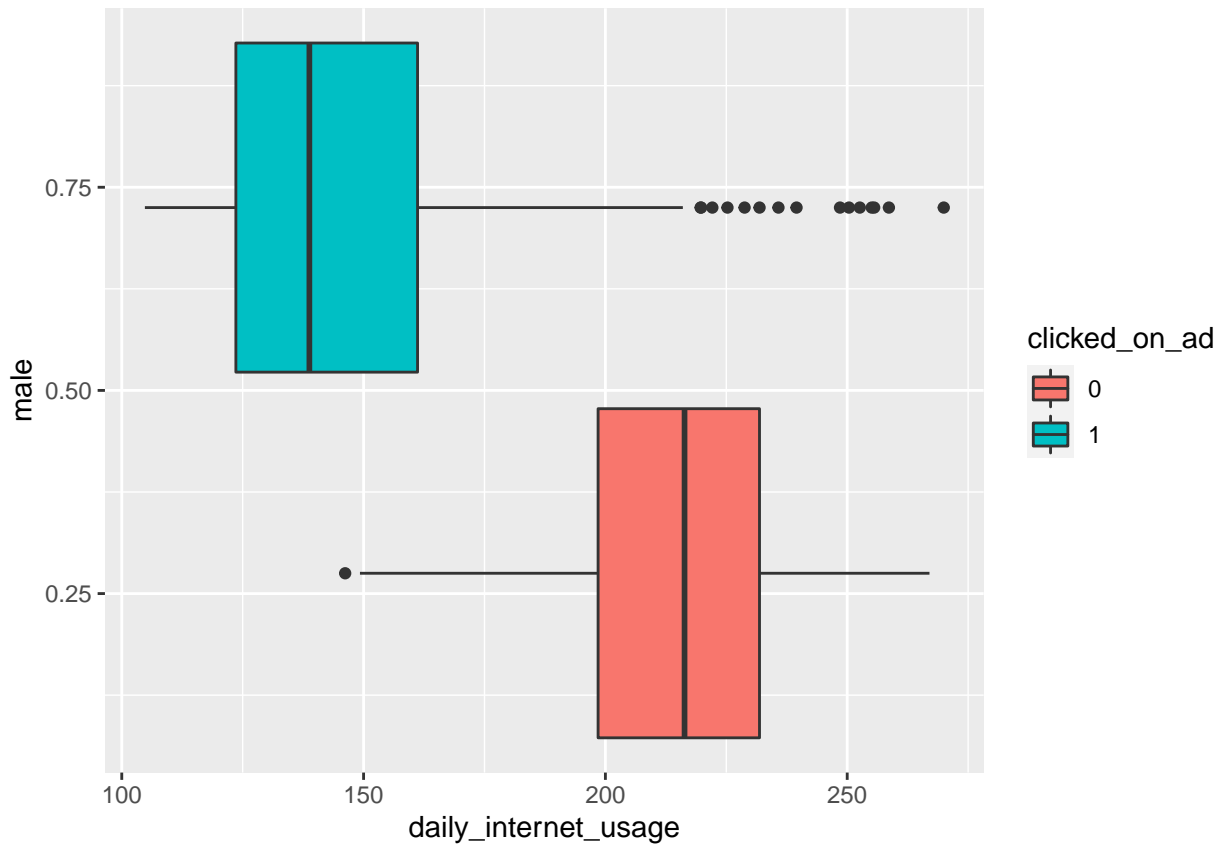
ggplot(data = ad_df, mapping = aes( x = male , y = daily_time_spent_on_site, fill = clicked_on_ad)) +
  geom_boxplot() +
  coord_flip()
```



More of those who click on the ads spend less time on the site. Of those who click on the ads, the females generally spend more time on the site as compared to the males

```
# Daily internet usage vs Gender
# Finding out and previewing boxplots to show how the Age relates with the gender

ggplot(data = ad_df, mapping = aes( x = male , y = daily_internet_usage, fill = clicked_on_ad)) +
  geom_boxplot() +
  coord_flip()
```

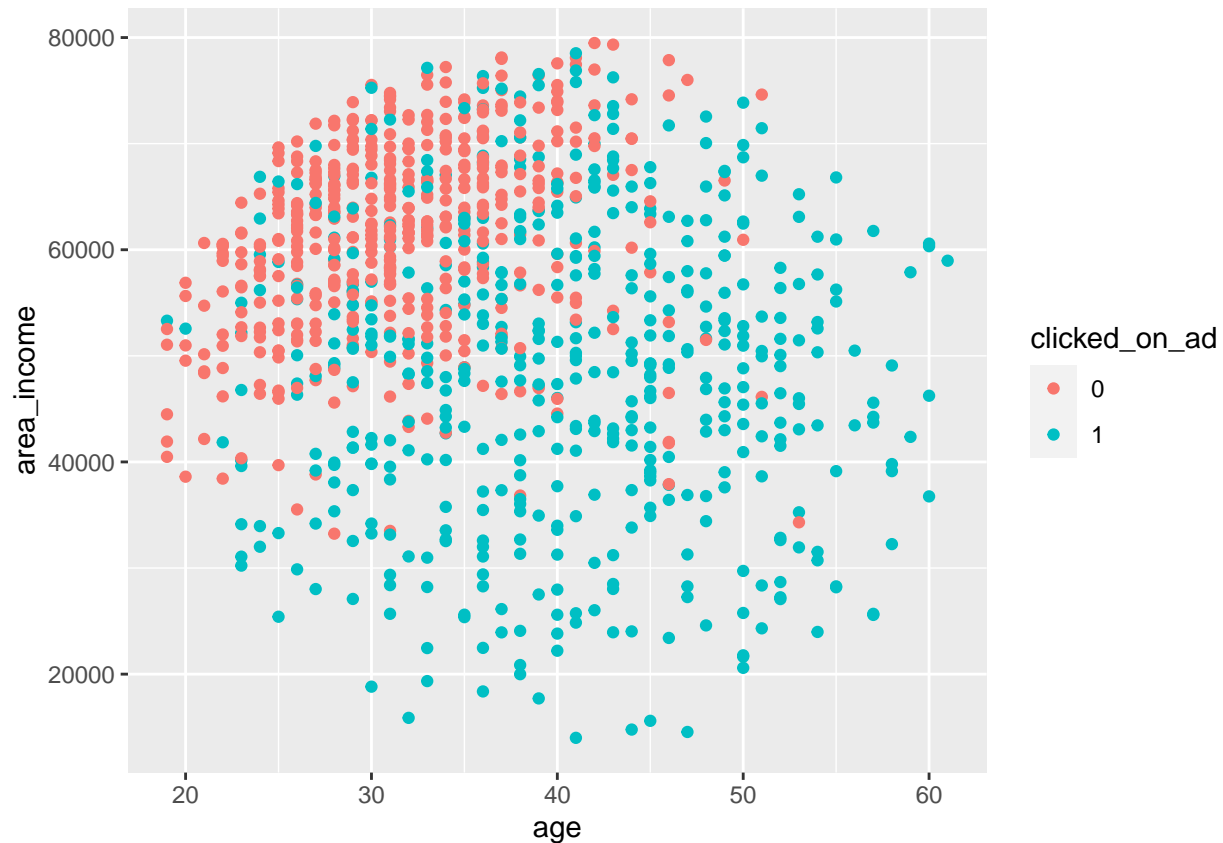


In general, those who click on the ads have a lower daily internet usage, with a few observations as outlier values with the males were slightly more than the females

Scatterplots

```
# Age vs Area Income
# Finding out and previewing scatterplots showing how the Age relates with the Area Income

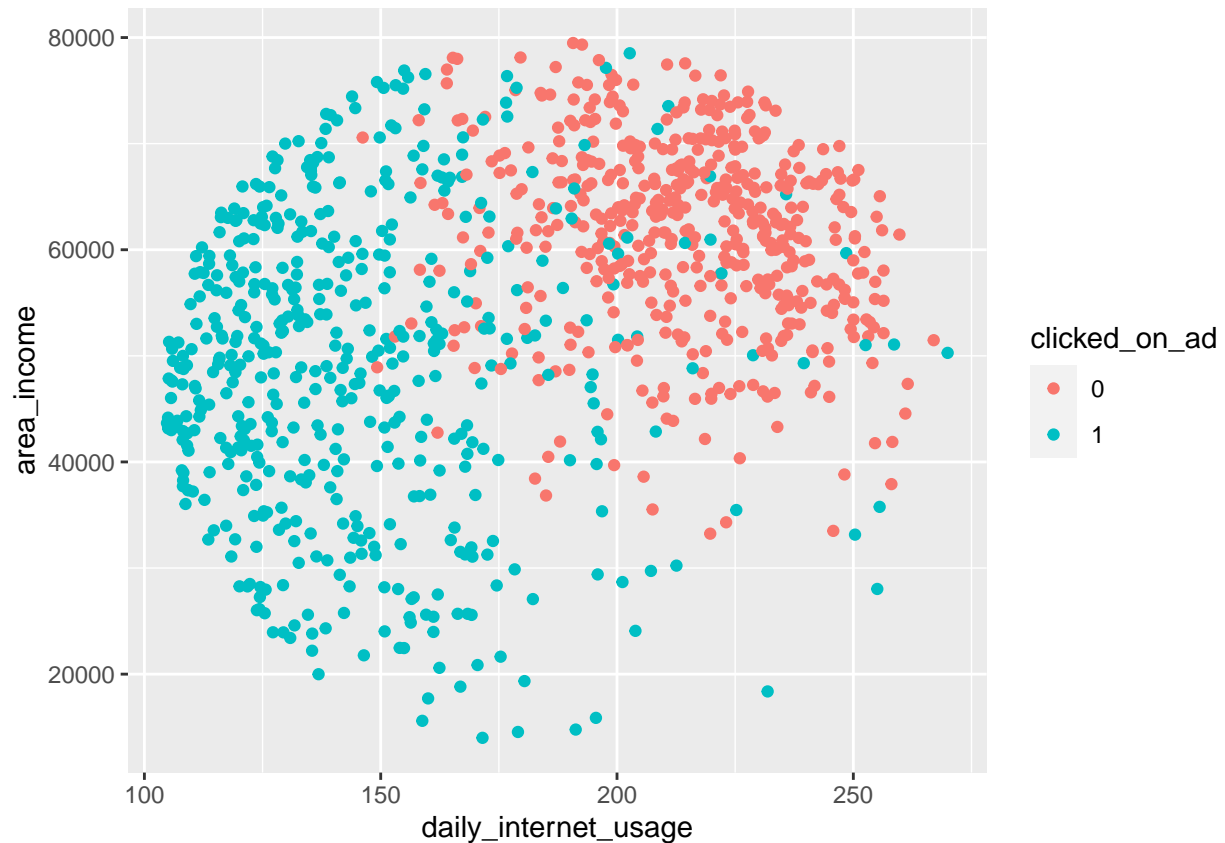
ggplot(data = ad_df) +
  geom_point(mapping = aes(x = age , y = area_income, color = clicked_on_ad))
```



We observe that the number of people who clicked on the ads are more evenly distributed while most of the people who did not click on the ads have a higher area income and a bit younger

```
# Daily Internet usage vs Area Income
# Finding out and previewing scatterplots showing how the Daily internet usage relates with the Area In

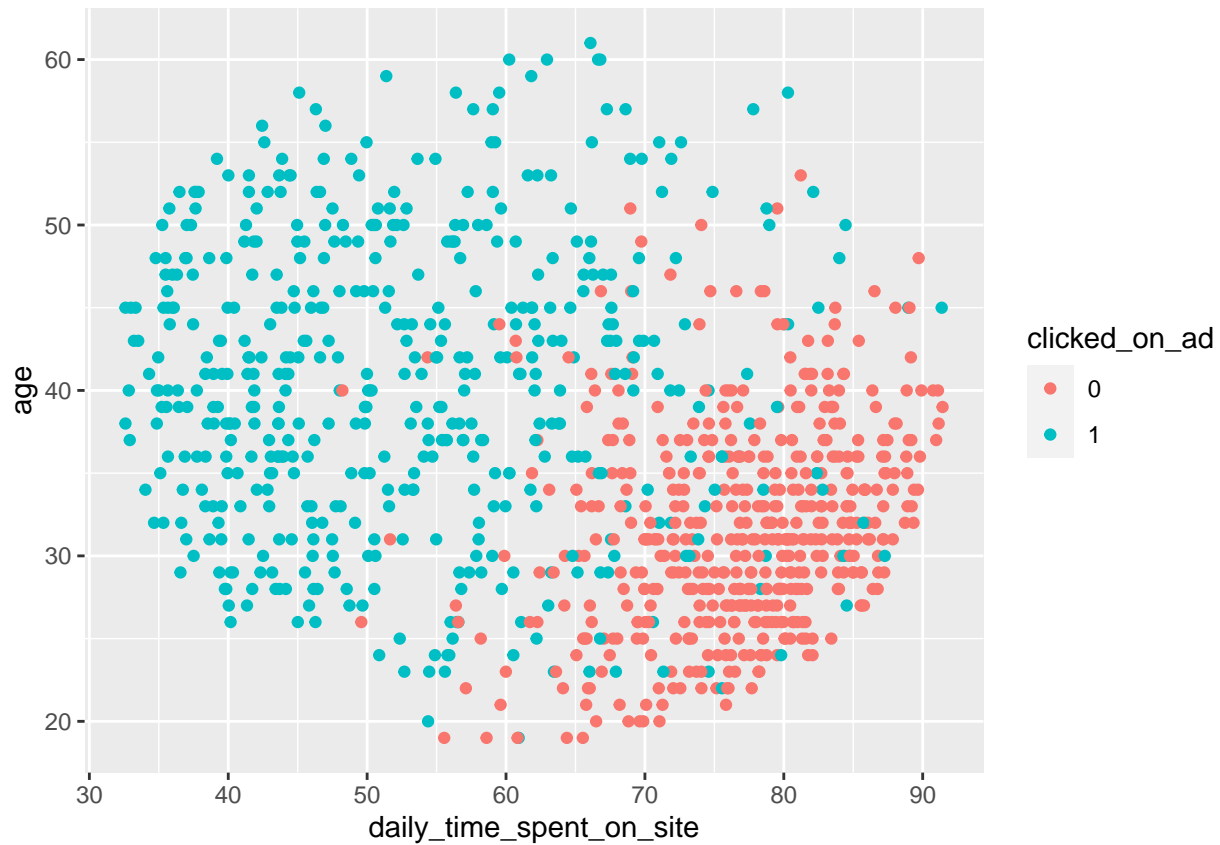
ggplot(data = ad_df) +
  geom_point(mapping = aes(x = daily_internet_usage , y = area_income, color = clicked_on_ad))
```

A great number of clicks comes from people who's daily internet usage is quite low and area income is also lower as compared to those who do Not click on the ads whose daily internet usage is significantly higher

```
# Age vs Daily time spent on site
# Finding out and previewing scatterplots showing how the Daily time spent on sites relates with the Ag

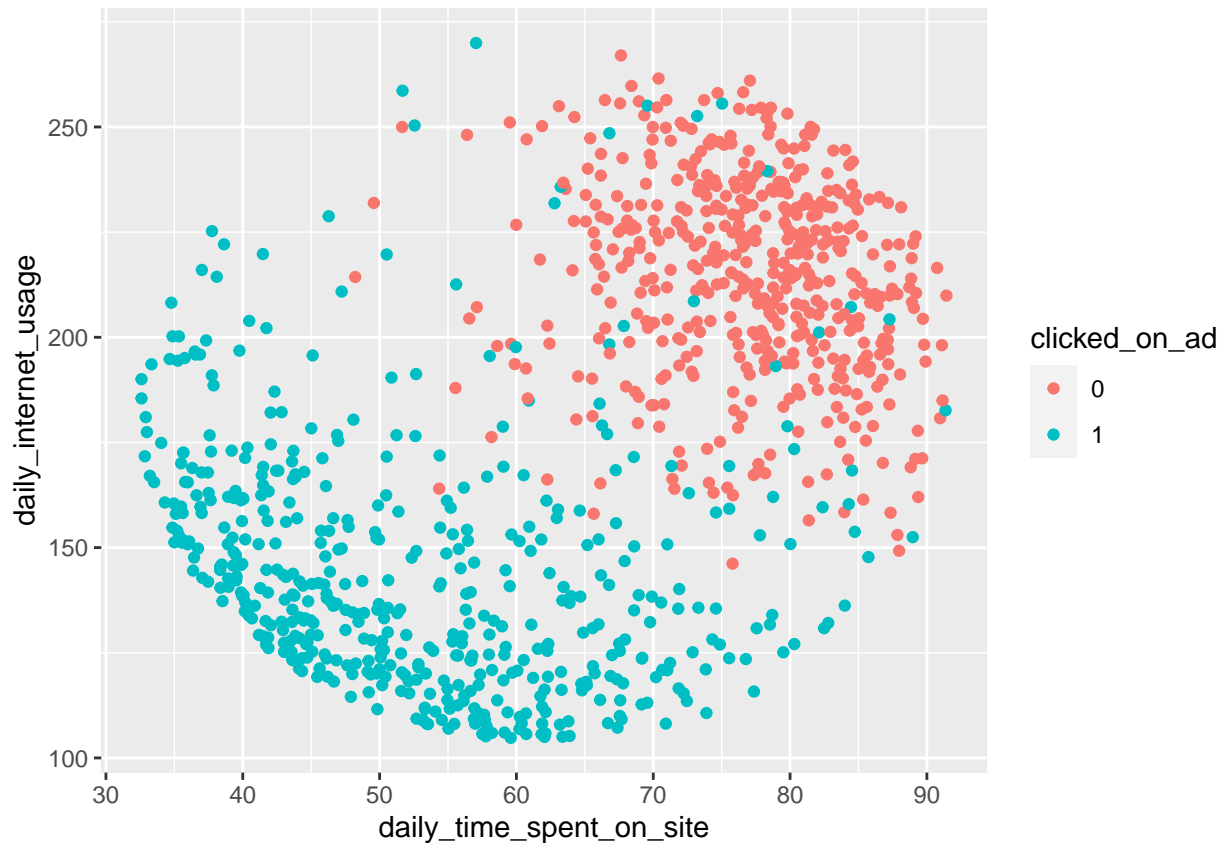
ggplot(data = ad_df) +
  geom_point(mapping = aes(x = daily_time_spent_on_site , y = age, color = clicked_on_ad))
```



A huge chunk of clicks come from people who spend significantly little time on the site as compared to those who spend more time on the site regardless of age

```
# Daily Internet Usage vs Daily time spent on site
# Finding out and previewing scatterplots showing how the Daily time spent on sites relates with the Da

ggplot(data = ad_df) +
  geom_point(mapping = aes(x = daily_time_spent_on_site , y = daily_internet_usage, color = clicked_on_a
```



Clearly, more clicks come from people who spend less time on the site and people whose daily internet usage is significantly lower as compared to those who spend more time on the site and have a high daily internet usage

The ads are getting more clicks from people who spend less time on the site and those whose daily internet usage is low.

9. Challenging the solution

Conclusion

Older people were more likely to be interested in cryptography than young users. The mean age of a person who clicked the ad was 40 years of age.

Females were more likely to click the cryptography ad than males however more analysis can be carried out in this particular area to determine the cause of this action.

The individuals from Lisamouth city were more likely to click the ad

People from the middle income areas clicked the ads more than the ones from a higher income area.

The lower daily internet usage users clicked the ads more than the ones who had a higher internet usage

Recommendations

We have observed that the users who were mostly interested in the ads were females who were older, had a lower area income and spent less time on the ads as they had less daily internet usage

To generate more intakes in the course, the company is better off increasing the number of ads towards the end and the beginning of the month and year as compared to the middle of the month and year.

Overall, we can say the study was successful based on the metrics of success.

Follow up Questions

Given that we had access to more data, would we be able to obtain better results?

END