

Define the Research Question

Identifying individuals most likely to click her ads.

The Metric of Success

Being able to identify individuals who are most likely to click her ads from our analysis.

The Context

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 would now like to employ your services as a Data Science Consultant to help her identify which individuals are most likely to click on her ads.

Experimental Design Taken

For us to meet our objective, the following experimental design was taken:

- Defining the research question.
- Setting the metric of success.
- Checking the appropriateness of the available data to answer the given question.
- Reading the dataset.
- Cleaning the dataset.
- Finding and dealing with outliers, anomalies and missing data within the dataset.
- Performing univariate and bivariate analysis.
- From the insights provide recommendations and a conclusion.

Appropriateness of the available data to answer the given question.

The data provided for analysis is very appropriate since it contains different variables which will help in answering our research question. The data also contains 1000 entries with no missing data or duplicates hence it is enough to conduct our analysis.

▼ Reading the dataset

```
#Reading the dataset
#Previewing the first six rows of the dataset.
library("data.table")
data = fread('http://bit.ly/IPAdvertisingData')
head(data)
```

A data.table: 6 × 10

Daily Time Spent on Site	Age	Area Income	Daily Internet Usage	Ad Topic Line	City	Male	Country	Timestamp
<dbl>	<int>	<dbl>	<dbl>	<chr>	<chr>	<int>	<chr>	<dtm>
68.95	35	61833.90	256.09	Cloned 5thgeneration orchestration	Wrightburgh	0	Tunisia	2016-03- 27 00:53:11
80.23	31	68441.85	193.77	Monitored national standardization	West Jodi	1	Nauru	2016-04- 04 01:39:02
69.47	26	59785.94	236.50	Organic bottom-line service-desk	Davidton	0	San Marino	2016-03- 13 20:35:42
74.15	29	54806.18	245.89	Triple-buffered reciprocal time-frame	West Terrifurt	1	Italy	2016-01- 10 02:31:19

```
#Previewing the last 10 rows of the dataset
tail(data, n=10)
```

A data.table: 10 × 10

Daily Time Spent on Site	Age	Area Income	Daily Internet Usage	Ad Topic Line	City	Male	Country	Timestamp
<dbl>	<int>	<dbl>	<dbl>	<chr>	<chr>	<int>	<chr>	<dt>
35.79	44	33813.08	165.62	Enterprise-wide tangible model	North Katie	1	Tonga	2016-01-01 13:36:00
38.96	38	36497.22	140.67	Versatile mission-critical application	Mauricefurt	1	Comoros	2016-01-01 16:02:00
69.17	40	66193.81	123.62	Extended leadingedge solution	New Patrick	0	Montenegro	2016-01-01 11:36:00
64.20	27	66200.96	227.63	Phased zero tolerance extranet	Edwardsmouth	1	Isle of Man	2016-01-01 23:45:00
				Front-line				2016-01-01 00:00:00

▼ Checking the dataset

```
12.91 30 11304.31 206.30 modular Dunystau 1 Lebanon
```

```
#Checking the number of columns in the dataset
ncol(data)
```

```
10
```

Our dataset has 10 columns.

```
#Checking the number of rows in the dataset
nrow(data)
```

```
1000
```

Our dataset has 1000 rows.

```
#Checking the dimensions of the dataset.
dim(data)
```

```
1000 10
```

```
#Checking the length of the dataset
length(data)
```

```
10
```

```
#Checking the structure of our dataset.
str(data)
```

```
Classes 'data.table' and '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 natio
 $ 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 01:39:0
 $ Clicked on Ad : int 0 0 0 0 0 0 0 1 0 0 ...
 - attr(*, ".internal.selfref")=<externalptr>
```

We can see the data types of the various variables in the dataset. Male and Clicked on Ad have the wrong data type and hence need to be corrected.

```
#Splitting the timestamp dataset to year, month, day and hour so that we can get as much info
data$Year <- format(as.POSIXct(data$Timestamp, format="%Y-%m-%d %H:%M:%S"), "%Y")
data$Month <- format(as.POSIXct(data$Timestamp, format="%Y-%m-%d %H:%M:%S"), "%m")
data$Day <- format(as.POSIXct(data$Timestamp, format="%Y-%m-%d %H:%M:%S"), "%d")
data$Hour <- format(as.POSIXct(data$Timestamp, format="%Y-%m-%d %H:%M:%S"), "%H")
head(data)
```

A data.table: 6 × 14

Daily Time Spent on Site	Age	Area Income	Daily Internet Usage	Ad Topic Line	City	Male	Country	Timestamp
<dbl>	<int>	<dbl>	<dbl>	<chr>	<chr>	<int>	<chr>	<dtm>
68.95	35	61833.90	256.09	Cloned 5thgeneration orchestration	Wrightburgh	0	Tunisia	2016-03-27 00:53:11
				Monitored				2016-04-

```
# drop the timestamp column since it is no longer useful
data$Timestamp <- NULL
colnames(data)
```

```
'Daily Time Spent on Site' · 'Age' · 'Area Income' · 'Daily Internet Usage' · 'Ad Topic Line' · 'City' · 'Male' ·
'Country' · 'Clicked on Ad' · 'Year' · 'Month' · 'Day' · 'Hour'
```

```
1 68.95 35 61833.90 256.09 Cloned 5thgeneration orchestration Wrightburgh 0 Tunisia 2016-03-27 00:53:11
2 68.95 35 61833.90 256.09 Monitored Wrightburgh 0 Tunisia 2016-04-01 00:01:10
```

```
#checking the data types of the new columns.
str(data)
```

```
Classes 'data.table' and 'data.frame': 1000 obs. of 13 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 natio
 $ 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" ...
 $ Clicked on Ad : int 0 0 0 0 0 0 0 1 0 0 ...
 $ Year : chr "2016" "2016" "2016" "2016" ...
 $ Month : chr "03" "04" "03" "01" ...
 $ Day : chr "27" "04" "13" "10" ...
 $ Hour : chr "00" "01" "20" "02" ...
 - attr(*, ".internal.selfref")=<externalptr>
```

The year, day, month and hour column have the wrong data type and so we have to change the datatype to factor.

```
# Correcting the data types of the year, month, day, hour and male columns.
data$Year <- as.factor(data$Year)
data$Month <- as.factor(data$Month)
data$Day <- as.factor(data$Day)
data$Hour <- as.factor(data$Hour)
data$Male <- as.factor(data$Male)
```

```
#checking to see if the columns have been assigned the right data types.
str(data)
```

```
Classes 'data.table' and 'data.frame': 1000 obs. of 13 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 natio
 $ City : chr "Wrightburgh" "West Jodi" "Davidton" "West Terrifurt"
 $ Male : Factor w/ 2 levels "0","1": 1 2 1 2 1 2 1 2 2 2 ...
 $ Country : chr "Tunisia" "Nauru" "San Marino" "Italy" ...
 $ Clicked on Ad : int 0 0 0 0 0 0 0 1 0 0 ...
 $ Year : Factor w/ 1 level "2016": 1 1 1 1 1 1 1 1 1 1 ...
 $ Month : Factor w/ 7 levels "01","02","03",...: 3 4 3 1 6 5 1 3 4 7
 $ Day : Factor w/ 31 levels "01","02","03",...: 27 4 13 10 3 19 28
 $ Hour : Factor w/ 24 levels "00","01","02",...: 1 2 21 3 4 15 21 2
 - attr(*, ".internal.selfref")=<externalptr>
```

```
#Removing white spaces from the column names
names(data)<-make.names(names(data),unique = TRUE)
```

```
#changing the data type of the column Clicked on Ad.
data$Clicked.on.Ad <- as.factor(data$Clicked.on.Ad)
```

```
#confirming if the column has been assigned the right data type.
str(data)
```

```
Classes 'data.table' and 'data.frame': 1000 obs. of 13 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 natio
 $ City : chr "Wrightburgh" "West Jodi" "Davidton" "West Terrifurt"
 $ Male : Factor w/ 2 levels "0","1": 1 2 1 2 1 2 1 2 2 2 ...
 $ Country : chr "Tunisia" "Nauru" "San Marino" "Italy" ...
 $ Clicked.on.Ad : Factor w/ 2 levels "0","1": 1 1 1 1 1 1 1 2 1 1 ...
 $ Year : Factor w/ 1 level "2016": 1 1 1 1 1 1 1 1 1 1 ...
 $ Month : Factor w/ 7 levels "01","02","03",...: 3 4 3 1 6 5 1 3 4 7
 $ Day : Factor w/ 31 levels "01","02","03",...: 27 4 13 10 3 19 28
 $ Hour : Factor w/ 24 levels "00","01","02",...: 1 2 21 3 4 15 21 2
 - attr(*, ".internal.selfref")=<externalptr>
```

```
#Previewing the dataset to see if the whitespaces have been removed
head(data)
```

A data.table: 6 × 13

Daily.Time.Spent.on.Site	Age	Area.Income	Daily.Internet.Usage	Ad.Topic.Line	
<dbl>	<int>	<dbl>	<dbl>	<chr>	
68.95	35	61833.90	256.09	Cloned 5thgeneration orchestration	Wrig
80.23	31	68441.85	193.77	Monitored national standardization	W
69.47	26	59785.94	236.50	Organic bottom- line service- desk	L
74.15	29	54806.18	245.89	Triple-buffered reciprocal time- frame	
68.37	35	73889.99	225.58	Robust logistical utilization	
59.99	23	59761.56	226.74	Sharable client- driven software	Ja

```
#Listing variables in our dataset.
names(data)
```

```
'Daily.Time.Spent.on.Site' · 'Age' · 'Area.Income' · 'Daily.Internet.Usage' · 'Ad.Topic.Line' · 'City' · 'Male' ·  
'Country' · 'Clicked.on.Ad' · 'Year' · 'Month' · 'Day' · 'Hour'
```

```
#Checking the class of Age column in the dataset.
class(data$Age)
```

```
'integer'
```

► Cleaning the dataset

[] ↳ 23 cells hidden

▼ Univariate Analysis

Measures of central tendency

```
#Checking the mean, median and mode of the Daily time spent on site column
dt.mean <- mean(data$Daily.Time.Spent.on.Site)
dt.mean
dt.median <- median(data$Daily.Time.Spent.on.Site)
dt.median
getmode <- function(v) {
  uniqv <- unique(v)
  uniqv[which.max(tabulate(match(v, uniqv)))]
}
dt.mode <- getmode(data$Daily.Time.Spent.on.Site)
dt.mode
```

```
65.0002
68.215
62.26
```

Daily Time Spent on Site Measures of central tendency

- Mean - 65.002
- Median - 68.215
- Mode - 62.26

```
#Checking the mean, median and mode of the age column
age.mean <- mean(data$Age)
age.mean
age.median <- median(data$Age)
age.median
getmode <- function(v) {
  uniqv <- unique(v)
  uniqv[which.max(tabulate(match(v, uniqv)))]
}
age.mode <- getmode(data$Age)
age.mode
```

```
36.009
35
31
```

Age Measures of central tendency

- Mean - 36.009
- Median - 35

- Mode - 31

```
#Checking the mean, median and mode of the area income column
ai.mean <- mean(data$Area.Income)
ai.mean
ai.median <- median(data$Area.Income)
ai.median
getmode <- function(v) {
  uniqv <- unique(v)
  uniqv[which.max(tabulate(match(v, uniqv)))]
}
ai.mode <- getmode(data$Area.Income)
ai.mode
```

```
55000.00008
57012.3
61833.9
```

Area income Measures of central tendency

- Mean - 55000.00008
- Median - 57012.3
- Mode -61833.9

```
#Checking the mean, median and mode of the daily internet usage column
diu.mean <- mean(data$Daily.Internet.Usage)
diu.mean
diu.median <- median(data$Daily.Internet.Usage)
diu.median
getmode <- function(v) {
  uniqv <- unique(v)
  uniqv[which.max(tabulate(match(v, uniqv)))]
}
diu.mode <- getmode(data$Daily.Internet.Usage)
diu.mode
```

```
180.0001
183.13
167.22
```

Daily Internet Usage Measures of central tendency

- Mean - 180.0001
- Median - 183.13
- Mode -167.22

```
#Checking the mode of the country column
getmode <- function(v) {
```

```

    uniqv <- unique(v)
    uniqv[which.max(tabulate(match(v, uniqv)))]
  }
country.mode <- getmode(data$Country)
country.mode

'Czech Republic'

```

Czech Republic is the most frequent country.

```

#Checking the mode of the city column
getmode <- function(v) {
  uniqv <- unique(v)
  uniqv[which.max(tabulate(match(v, uniqv)))]
}
city.mode <- getmode(data$City)
city.mode

'Lisamouth'

```

Lisamouth is the most frequent city.

```

#Checking the mode of the sex column
getmode <- function(v) {
  uniqv <- unique(v)
  uniqv[which.max(tabulate(match(v, uniqv)))]
}
sex.mode <- getmode(data$Sex)
sex.mode

0
► Levels:

```

Most people from the data collected are female.

```

#Checking the mode of the Daily.Internet.Usage column
getmode <- function(v) {
  uniqv <- unique(v)
  uniqv[which.max(tabulate(match(v, uniqv)))]
}
diu.mode <- getmode(data$Daily.Internet.Usage)
diu.mode

167.22

```

Most people had a daily internet usage of 167.22

```
#Checking the mode of the Year column
getmode <- function(v) {
  uniqv <- unique(v)
  uniqv[which.max(tabulate(match(v, uniqv)))]
}
yr.mode <- getmode(data$Year)
yr.mode
```

2016

► **Levels:**

2016 is the only year represented in the dataset.

```
#Checking the mode of the Year column
getmode <- function(v) {
  uniqv <- unique(v)
  uniqv[which.max(tabulate(match(v, uniqv)))]
}
mth.mode <- getmode(data$Month)
mth.mode
```

02

► **Levels:**

February is the month that appears multiple times.

```
#Checking the mode of the Year column
getmode <- function(v) {
  uniqv <- unique(v)
  uniqv[which.max(tabulate(match(v, uniqv)))]
}
dy.mode <- getmode(data$Day)
dy.mode
```

03

► **Levels:**

Day 3 of the month appears the most.

```
#Checking the mode of the Year column
getmode <- function(v) {
```

```

    uniqv <- unique(v)
    uniqv[which.max(tabulate(match(v, uniqv)))]
  }
hr.mode <- getmode(data$Hour)
hr.mode

```

07

► **Levels:**

7:00am appears the most in our dataset.

Measures of dispersion

```

#Checking the min, max, range, quantile, variance, standard deviation of Daily time spent or
dt.min <- min(data$Daily.Time.Spent.on.Site)
dt.min
dt.max <- max(data$Daily.Time.Spent.on.Site)
dt.max
dt.range <- range(data$Daily.Time.Spent.on.Site)
dt.range
dt.quantile <- quantile(data$Daily.Time.Spent.on.Site)
dt.quantile
dt.var <- var(data$Daily.Time.Spent.on.Site)
dt.var
dt.sd <- sd(data$Daily.Time.Spent.on.Site)
dt.sd

```

32.6

91.43

32.6 · 91.43

0%:	32.6	25%:	51.36	50%:	68.215	75%:	78.5475	100%:	91.43
-----	------	------	-------	------	--------	------	---------	-------	-------

251.337094854855

15.8536145675002

```

#Checking for the skewness of daily time spent on site.
install.packages("moments")
libraryskewness(data$Daily.Time.Spent.on.Site)y(moments)

```

Installing package into ‘/usr/local/lib/R/site-library’
(as ‘lib’ is unspecified)

-0.371202614867441

Since the skewness is negative, it means we have a left skewed distribution

```
#Checking for Kurtosis
kurtosis(data$Daily.Time.Spent.on.Site)

1.90394215401081
```

Has a leptokurtic distribution since the kurtosis is > 0 .

Daily Time Spent on Site Measures of dispersion

- Min - 32.6
- Max - 91.43
- Range - 32.691.43
- Quantile - 0% - 32.6 25% - 51.36 50% - 68.215 75% - 78.5475 100% - 91.43
- Variance - 251.337094854855
- Standard deviation - 15.8536145675002

```
#Checking the min, max, range, quantile, variance, standard deviation of the age column
age.min <- min(data$Age)
age.min
age.max <- max(data$Age)
age.max
age.range <- range(data$Age)
age.range
age.quantile <- quantile(data$Age)
age.quantile
age.var <- var(data$Age)
age.var
age.sd <- sd(data$Age)
age.sd

19
61
19 - 61
0%:      19 25%:      29 50%:      35 75%:      42 100%:      61
77.1861051051051
8.78556231012592
```

```
#checking for skewness of the age column
skewness(data$Age)

0.478422676206608
```

Since the skewness is positive, it means we have a right skewed distribution

```
#Checking for Kurtosis
```

```
kurtosis(data$Age)
```

```
2.59548176807726
```

Has a leptokurtic distribution since the kurtosis is > 0.

Age Measures of dispersion

- Min - 19
- Max - 61
- Range - 19-61
- Quantile - 0%: 19 25%: 29 50%: 35 75%: 42 100%: 61
- Variance - 77.1861051051051
- Standard deviation - 18.78556231012592

#Checking the min, max, range, quantile, variance, standard deviation of the area income column

```
ai.min <- min(data$Area.Income)
```

```
ai.min
```

```
ai.max <- max(data$Area.Income)
```

```
ai.max
```

```
ai.range <- range(data$Area.Income)
```

```
ai.range
```

```
ai.quantile <- quantile(data$Area.Income)
```

```
ai.quantile
```

```
ai.var <- var(data$Area.Income)
```

```
ai.var
```

```
ai.sd <- sd(data$Area.Income)
```

```
ai.sd
```

```
13996.5
```

```
79484.8
```

```
13996.5 79484.8
```

```
0%: 13996.5 25%: 47031.8025 50%: 57012.3 75%: 65470.635 100%:
```

```
79484.8
```

```
179952405.951775
```

```
13414 6340722824
```

#Checking for the skewness of the Area income column

```
skewness(data$Area.Income)
```

```
-0.649396701694076
```

Since the skewness is negative, it means we have a left skewed distribution.

#Checking for Kurtosis

```
kurtosis(data$Area.Income)
```

2 89469406161926

Has a leptokurtic distribution since the kurtosis is > 0 .

Area Income Measures of dispersion

- Min - 13996.5
- Max - 79484.8
- Range - 13996.579484.8
- Quantile - 0% - 13996.5 25% - 47031.8025 50% - 57012.3 75% - 65470.635 100% - 79484.8
- Variance - 179952405.951775
- Standard deviation - 13414.6340222824

```
#Checking the min, max, range, quantile, variance, standard deviation of the daily internet u
diu.min <- min(data$Daily.Internet.Usage)
diu.min
diu.max <- max(data$Daily.Internet.Usage)
diu.max
diu.range <- range(data$Daily.Internet.Usage)
diu.range
diu.quantile <- quantile(data$Daily.Internet.Usage)
diu.quantile
diu.var <- var(data$Daily.Internet.Usage)
diu.var
diu.sd <- sd(data$Daily.Internet.Usage)
diu.sd

104.78
269.96
104.78 · 269.96
0%:      104.78 25%:      138.83 50%:      183.13 75%:      218.7925 100%:      269.96
1927.41539618619
43.9023393019801
```

```
#Checking for the skewness of the daily internet usage column.
skewness(data$Daily.Internet.Usage)

-0.0334870316434409
```

Since the skewness is negative, it means we have a left skewed distribution.

```
#Checking for Kurtosis
kurtosis(data$Daily.Internet.Usage)
```

```
1 72770118001810
```

Has a leptokurtic distribution since the kurtosis is > 0 .

Daily Internet Used Measures of dispersion

- Min - 104.78
- Max - 269.96
- Range - 104.78269.96
- Quantile - 0% - 104.78 25% - 138.83 50% - 183.13 75% - 218.7925 100% - 269.96
- Variance - 1927.41539618619
- Standard deviation - 43.9023393019801

Univariate graphs

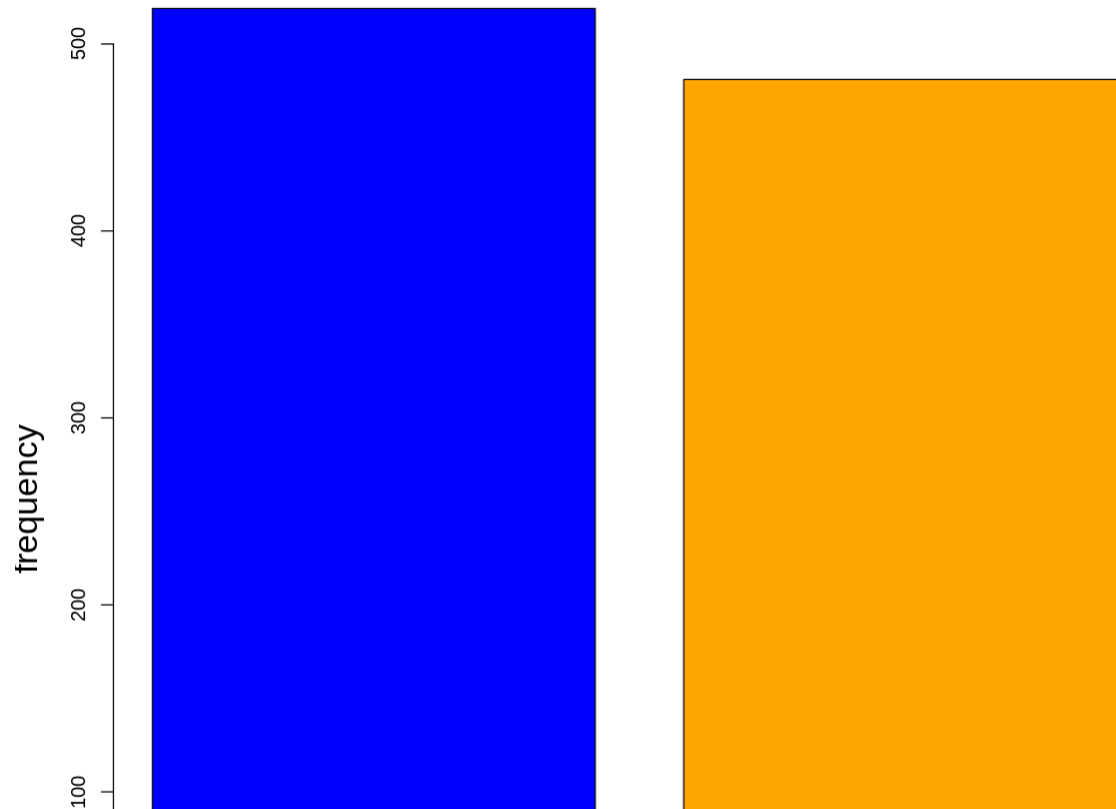
```
# Fetching the Sex column
# Computing the frequency of both male and female respondents.

sex <- data$Sex
sex_frequency <- table(sex)
sex_frequency

      sex
      0   1
519 481

#Bar graph representing Sex frequency.
options(repr.plot.width = 10, repr.plot.height = 10)
barplot(c(sex_frequency), main="A barplot representing the Sex column.",
        xlab="Sex",
        ylab="frequency",
        sub="From the graph we can see that females(0) are more than males(1)",
        cex.main=2, cex.lab=1.7, cex.sub=1.2,
        width=c(30,30),
        col=c("blue", "orange"))
```


A barplot representing the Sex column.



We can observe from the frequency table and from the graph that most respondents are female.

- 519 - Female
- 481 - Male

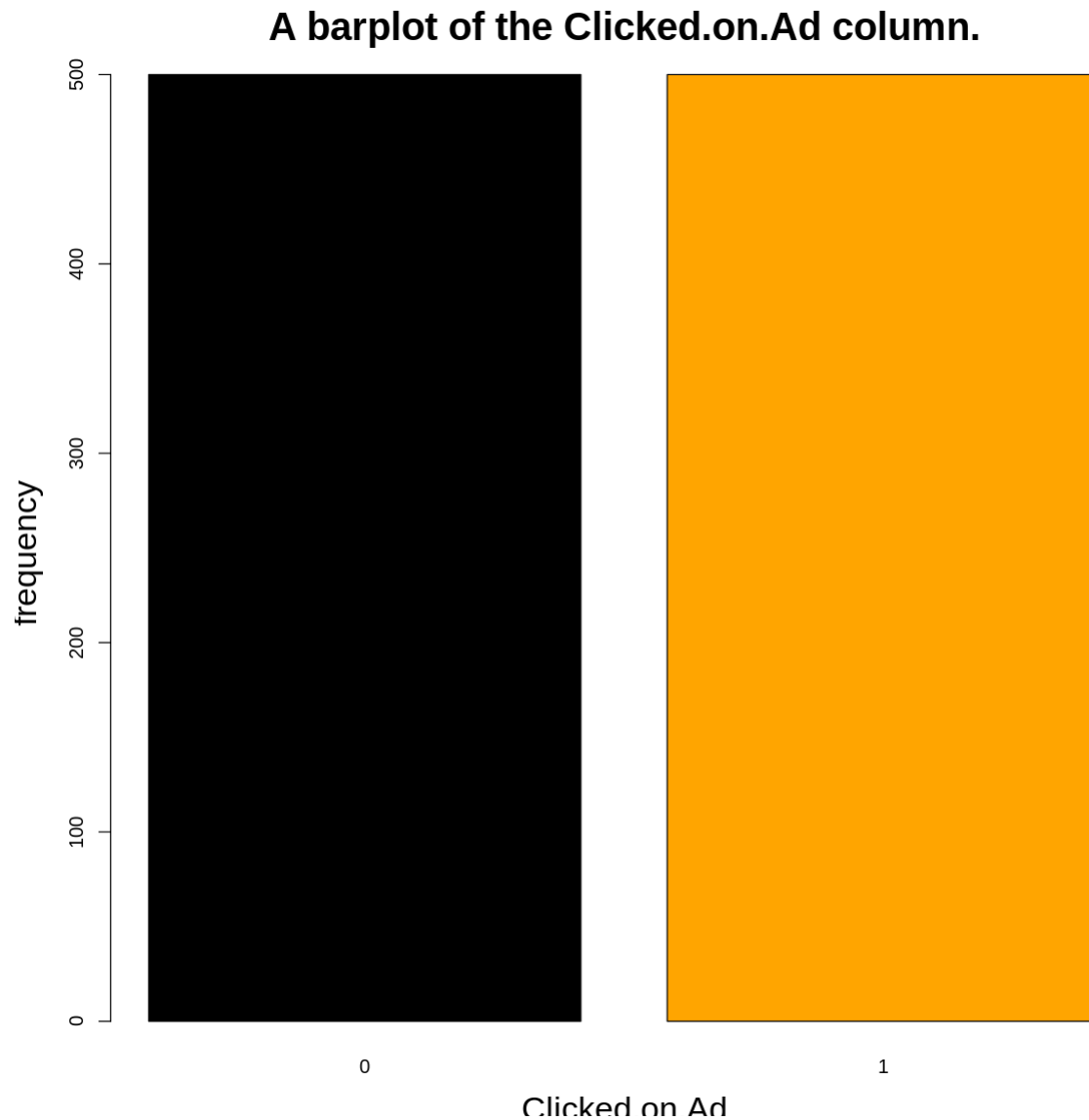
0

1

```
# Fetching the Clicked on ad column
# Computing the frequency of respondents who clicked on the ad and those who did not..
Clicked.on.Ad <- data$Clicked.on.Ad
Clicked.on.Ad_frequency <- table(Clicked.on.Ad)
Clicked.on.Ad_frequency

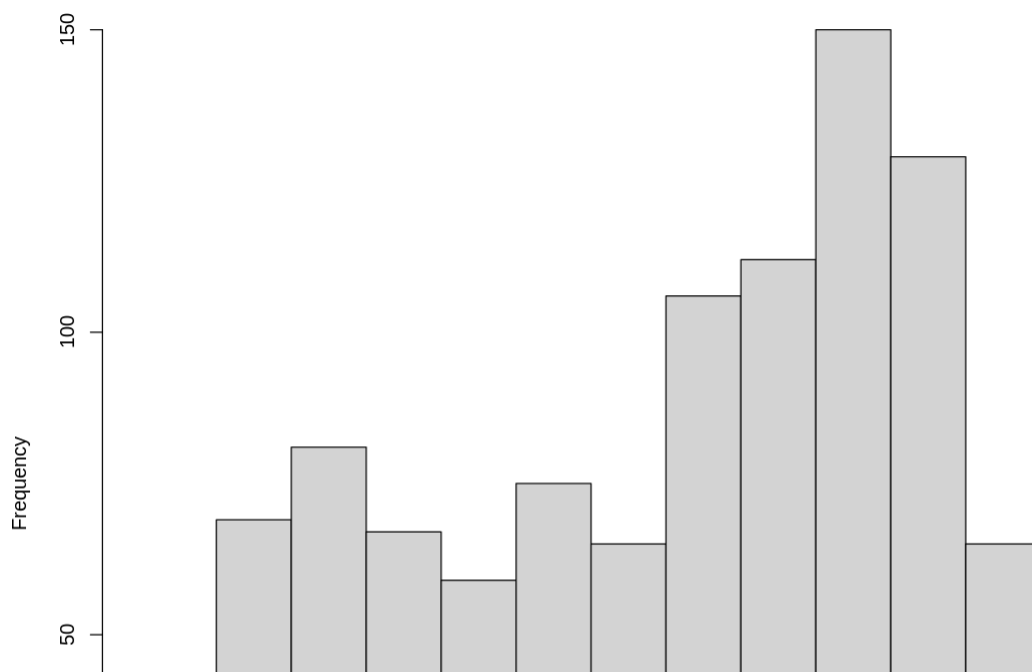
Clicked.on.Ad
  0    1
500 500

#Bar graph representing clicked on ad frequency
options(repr.plot.width = 10, repr.plot.height = 10)
barplot(c(Clicked.on.Ad_frequency), main="A barplot of the Clicked.on.Ad column.",
        xlab="Clicked.on.Ad",
        ylab="frequency",
        sub="The proportion of people who clicked on ad and those who did not is equal.",
        cex.main=2, cex.lab=1.7,cex.sub=1.2,
        col=c("black","orange"))
```



From the frequency table and the bar graph, we observe that there is a balance of those who clicked on the ad and those who did not.

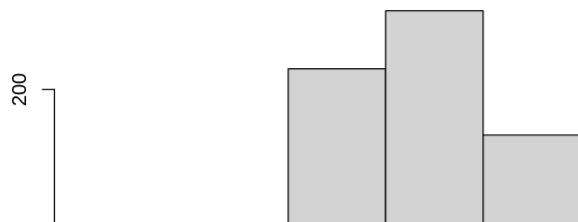
```
#A histogram of the daily time spent on site.  
options(repr.plot.width = 10, repr.plot.height = 10)  
hist(data$Daily.Time.Spent.on.Site)
```

Histogram of data\$Daily.Time.Spent.on.Site

The histogram appears to be relatively uniform.

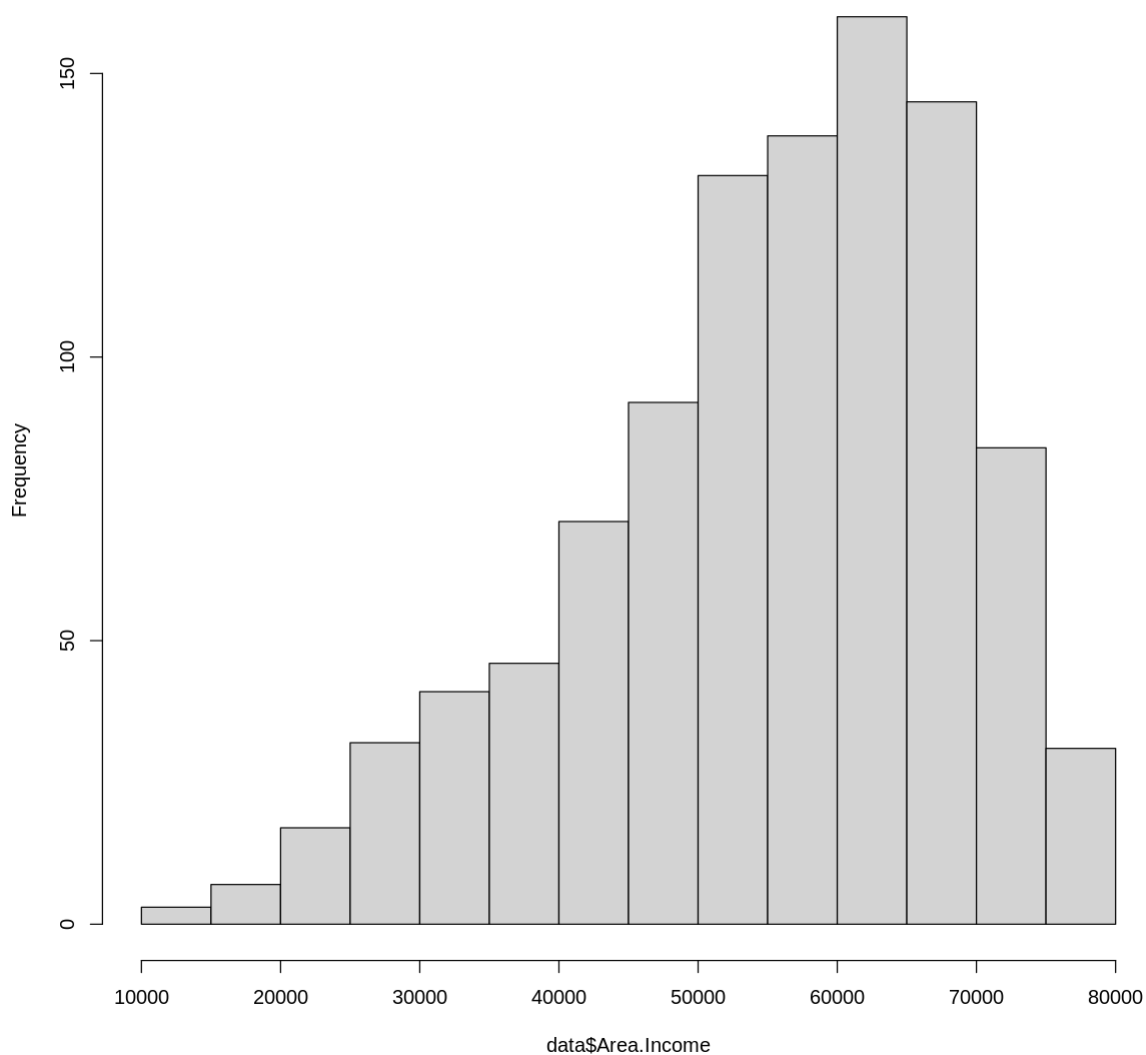


```
#A histogram of age  
options(repr.plot.width = 10, repr.plot.height = 10)  
hist(data$Age)
```

Histogram of data\$Age

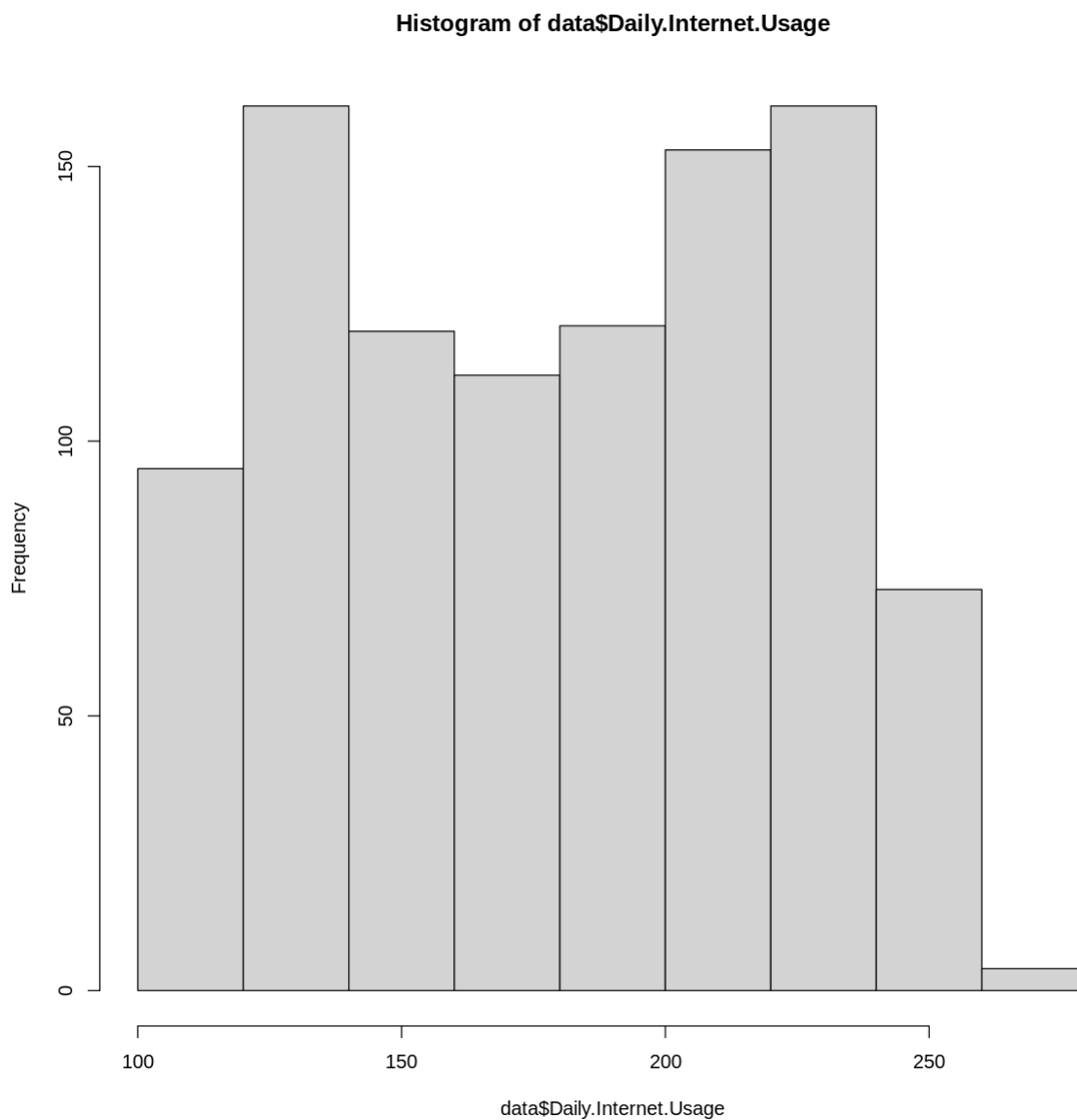
The histogram is skewed to the right.

```
#A histogram of Area Income
options(repr.plot.width = 10, repr.plot.height = 10)
hist(data$Area.Income)
```

Histogram of data\$Area.Income

The above histogram is skewed to the left.

```
#A histogram of daily internet usage
options(repr.plot.width = 10, repr.plot.height = 10)
hist(data$Daily.Internet.Usage)
```



The histogram appears to be relatively normal.

```
#Checking the most frequent countries.
sort(table(data$Country), decreasing=TRUE)[1:10]
```

Czech Republic	France	Afghanistan	Australia	Cyprus
9	9	8	8	8
Greece	Liberia	Micronesia	Peru	Senegal
8	8	8	8	8

Czech Republic, France and Afghanistan are among the first 10 countries with the highest frequency.

```
#Checking the most frequent cities.
sort(table(data$City), decreasing=TRUE)[1:10]
```

Lisamouth	Williamsport	Benjaminchester	East John	East Timothy
3	3	2	2	2
Johnstad	Joneston	Lake David	Lake James	Lake Jose
2	2	2	2	2

Lisamouth, Williamsport Benjaminchester, East John and East Timothy are among the first 10 cities with the highest frequency.

```
#Checking the most frequent ad topic line
sort(table(data$Ad.Topic.Line), decreasing=TRUE)[1:10]
```

Adaptive 24hour Graphic Interface	Adaptive asynchronous attitude
1	1
Adaptive context-sensitive application	Adaptive contextually-based methodology
1	1
Adaptive demand-driven knowledgebase	Adaptive uniform capability
1	1
Advanced 24/7 productivity	Advanced 5thgeneration capability
1	1
Advanced didactic conglomeration	Advanced disintermediate data-warehouse
1	1

All entries in this column occurred once.

```
#Checking the most frequent hour.
sort(table(data$Hour), decreasing=TRUE)[1:10]
```

07 20 09 21 00 05 23 08 14 22
54 50 49 48 45 44 44 43 43 43

7:00am and 9:00pm are the most frequent hours.

```
#Checking the most frequent month.
sort(table(data$Month), decreasing=TRUE)[1:7]
```

February, March and January are the most frequent months.

#Checking the most frequent day.

```
sort(table(data$Day), decreasing=TRUE)[1:10]
```

```
03 17 15 10 04 26 05 08 16 18
46 42 41 37 36 36 35 35 35 35
```

Day 3, 17 and 15 are the most frequent days.

▼ Bivariate Analysis.

Covariance

#Covariance of daily time spent on site and age.

```
Daily.Time.Spent.on.Site <- data$Daily.Time.Spent.on.Site
```

```
Age <- data$Age
```

```
cov(Daily.Time.Spent.on.Site, Age)
```

```
-46.1741459459459
```

The covariance is negative showing that greater values of one variable corresponds to smaller values of the other.

#Covariance of daily time spent on site and Area.Income.

```
Daily.Time.Spent.on.Site <- data$Daily.Time.Spent.on.Site
```

```
Area.Income <- data$Area.Income
```

```
cov(Daily.Time.Spent.on.Site, Area.Income)
```

```
66130.8109081922
```

The covariance is positive showing that greater values of one variable correspond to greater values of the other.

#Covariance of daily time spent on site and Daily.Internet.Usage.

```
Daily.Internet.Usage <- data$Daily.Internet.Usage
```

```
Daily.Time.Spent.on.Site <- data$Daily.Time.Spent.on.Site
```

```
cov(Daily.Time.Spent.on.Site, Daily.Internet.Usage)
```

360.991882662663

The covariance is positive showing that greater values of one variable correspond to greater values of the other.

```
#Covariance of Age and Area.Income.  
Age<- data$Age  
Area.Income <- data$Area.Income  
cov(Age, Area.Income)
```

-21520.9257965165

The covariance is negative showing that greater values of one variable corresponds to smaller values of the other.

```
#covariance of age and daily internet usage.  
Age<- data$Age  
Daily.Internet.Usage <- data$Daily.Internet.Usage  
cov(Age, Daily.Internet.Usage <- data$Daily.Internet.Usage)
```

-141.634815715716

The covariance is negative showing that greater values of one variable corresponds to smaller values of the other.

```
#covariance of area income and daily internet usage.  
Area.Income <- data$Area.Income  
Daily.Internet.Usage <- data$Daily.Internet.Usage  
cov(Area.Income, Daily.Internet.Usage <- data$Daily.Internet.Usage)
```

198762.531532925

The covariance is positive showing that greater values of one variable correspond to greater values of the other.

Correlation

```
#Correlation of daily time spent on site and age.  
Daily.Time.Spent.on.Site <- data$Daily.Time.Spent.on.Site  
Age <- data$Age  
cor(Daily.Time.Spent.on.Site, Age)
```

-0.331513342786584

The two variables have a weak negative correlation.

```
#Correlation of daily time spent on site and Area.Income.  
Daily.Time.Spent.on.Site <- data$Daily.Time.Spent.on.Site  
Area.Income <- data$Area.Income  
cor(Daily.Time.Spent.on.Site, Area.Income)
```

0.310954412522883

The two variables have a weak positive correlation.

```
#Correlation of daily time spent on site and Daily.Internet.Usage.  
Daily.Internet.Usage <- data$Daily.Internet.Usage  
Daily.Time.Spent.on.Site <- data$Daily.Time.Spent.on.Site  
cor(Daily.Time.Spent.on.Site, Daily.Internet.Usage)
```

0.518658475337186

The two variables have a positive correlation.

```
#Correlation of Age and Area.Income.  
Age<- data$Age  
Area.Income <- data$Area.Income  
cor(Age, Area.Income)
```

-0.182604955032622

The two variables have a weak negative correlation.

```
#Correlation of age and daily internet usage.  
Age<- data$Age  
Daily.Internet.Usage <- data$Daily.Internet.Usage  
cor(Age, Daily.Internet.Usage <- data$Daily.Internet.Usage)
```

-0.367208560147359

The two variables have a weak negative correlation.

```
#Correlation of area income and daily internet usage.  
Area.Income <- data$Area.Income  
Daily.Internet.Usage <- data$Daily.Internet.Usage  
cor(Area.Income, Daily.Internet.Usage <- data$Daily.Internet.Usage)
```

0.337495532865276

The two variables have a weak positive correlation.

▼ Selecting data that consists of people who clicked on ad.

▼ Univariate Analysis of the people who clicked on the ad

```
#Selecting the data with click on ad as 1
clicked <- data[data$Clicked.on.Ad ==1,]
head(clicked)
```

A data.table: 6 × 5

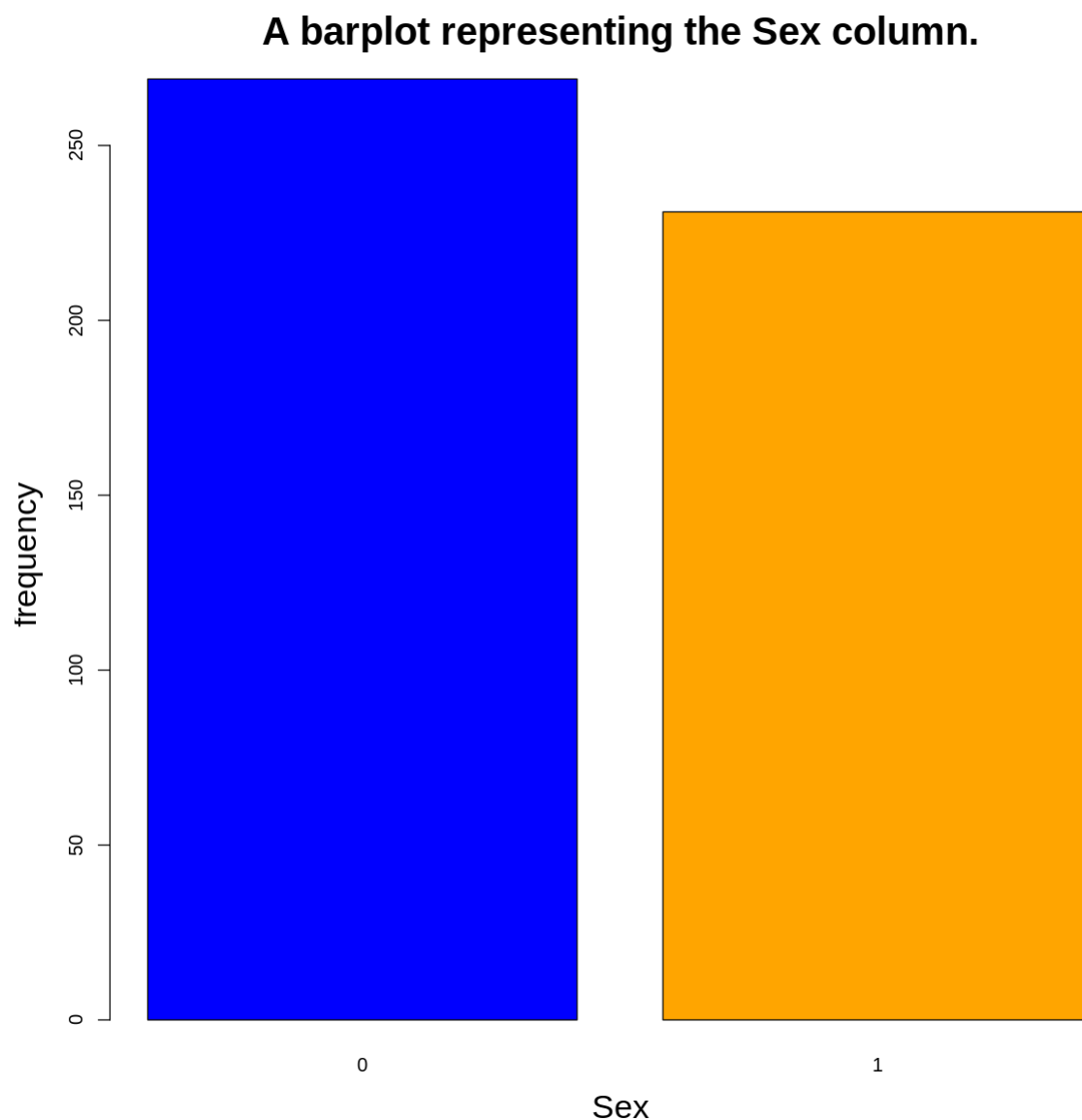
Daily.Time.Spent.on.Site	Age	Area.Income	Daily.Internet.Usage	Ad.Topic.Line
<dbl>	<int>	<dbl>	<dbl>	<chr>
66.00	48	24593.33	131.76	Reactive local challenge
47.64	49	45632.51	122.02	Centralized neutral neural-net
69.57	48	51636.92	113.12	Centralized content-based focus group
42.95	33	30976.00	143.56	Grass-roots coherent extranet
63.45	23	52182.23	140.64	Persistent demand-driven interface
55.39	37	23936.86	129.41	Customizable multi-tasking website

```
# Fetching the Sex column
sex <- clicked$Sex
sex_frequency <- table(sex)
sex_frequency
```

```

options(repr.plot.width = 10, repr.plot.height = 10)
barplot(c(sex_frequency), main="A barplot representing the Sex column.",
        xlab="Sex",
        ylab="frequency",
        sub="From the graph we can see that females(0) are more than males(1)",
        cex.main=2, cex.lab=1.7, cex.sub=1.2,
        width=c(30,30),
        col=c("blue","orange"))

```



From the graph we can see that females(0) are more than males(1)

From those who clicked on the ad, 269 were female and 231 were male.

```

#Checking the most frequent Daily.Time.Spent.on.Site.
sort(table(clicked$Daily.Time.Spent.on.Site), decreasing=TRUE)[1:10]

```

```
75 55 22 6 25 19 25 66 25 98 28 25 29 86 29 96 11 19 11 72
```

```
#Checking the most frequent Age.
```

```
sort(table(clicked$Age), decreasing=TRUE)[1:10]
```

```
45 36 38 41 42 40 43 50 39 49
27 25 25 22 20 19 19 19 17 17
```

We can see the most frequent age of respondents who clicked on the ad as 40's, 30's and 50's.

```
#Checking the most frequent Daily.Internet.Usage.
```

```
sort(table(clicked$Daily.Internet.Usage), decreasing=TRUE)[1:10]
```

```
113.53 115.91 117.3 119.3 120.06 125.45 132.38 135.24 136.18 138.35
      2      2      2      2      2      2      2      2      2      2
```

```
#Checking the most frequent cities.
```

```
sort(table(clicked$City), decreasing=TRUE)[1:10]
```

```
Lake David  Lake James  Lisamouth Michelleside  Millerbury  Robertfurt
      2      2      2      2      2
South Lisa  West Amanda West Shannon Williamsport
      2      2      2      2
```

the cities that are more frequent are Lake David, Lake James and so on.

```
#Checking the most frequent Country.
```

```
sort(table(clicked$Country), decreasing=TRUE)[1:10]
```

```
Australia  Ethiopia  Turkey  Liberia Liechtenstein
      7      7      7      6      6
South Africa  Afghanistan  France  Hungary  Mayotte
      6      5      5      5      5
```

The most frequent country is Australia.

```
#Checking the most frequent Month.
```

```
sort(table(clicked$Month), decreasing=TRUE)[1:10]
```

```
02 05 03 04 06 01 07 <NA> <NA> <NA>
83 79 74 74 71 69 50
```

The most frequent month is February.

```
#Checking the most frequent Day.
```

```
sort(table(clicked$Day), decreasing=TRUE)[1:10]
```

```
03 23 14 09 12 15 01 10 05 17
26 22 21 20 20 20 19 19 18 18
```

The most frequent day of the month is day 3.

```
#Checking the most frequent Hour.
```

```
sort(table(clicked$Hour), decreasing=TRUE)[1:10]
```

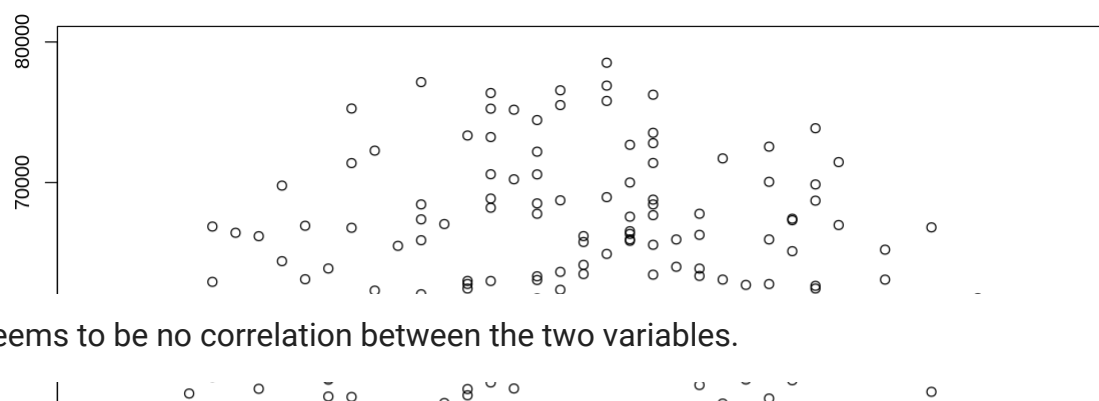
```
09 00 07 18 11 20 03 06 17 04
28 26 26 25 24 24 23 23 23 21
```

The most frequent hour is 9 am.

Scatter Plots

```
#A scatter plot of Age vs Area Income.
```

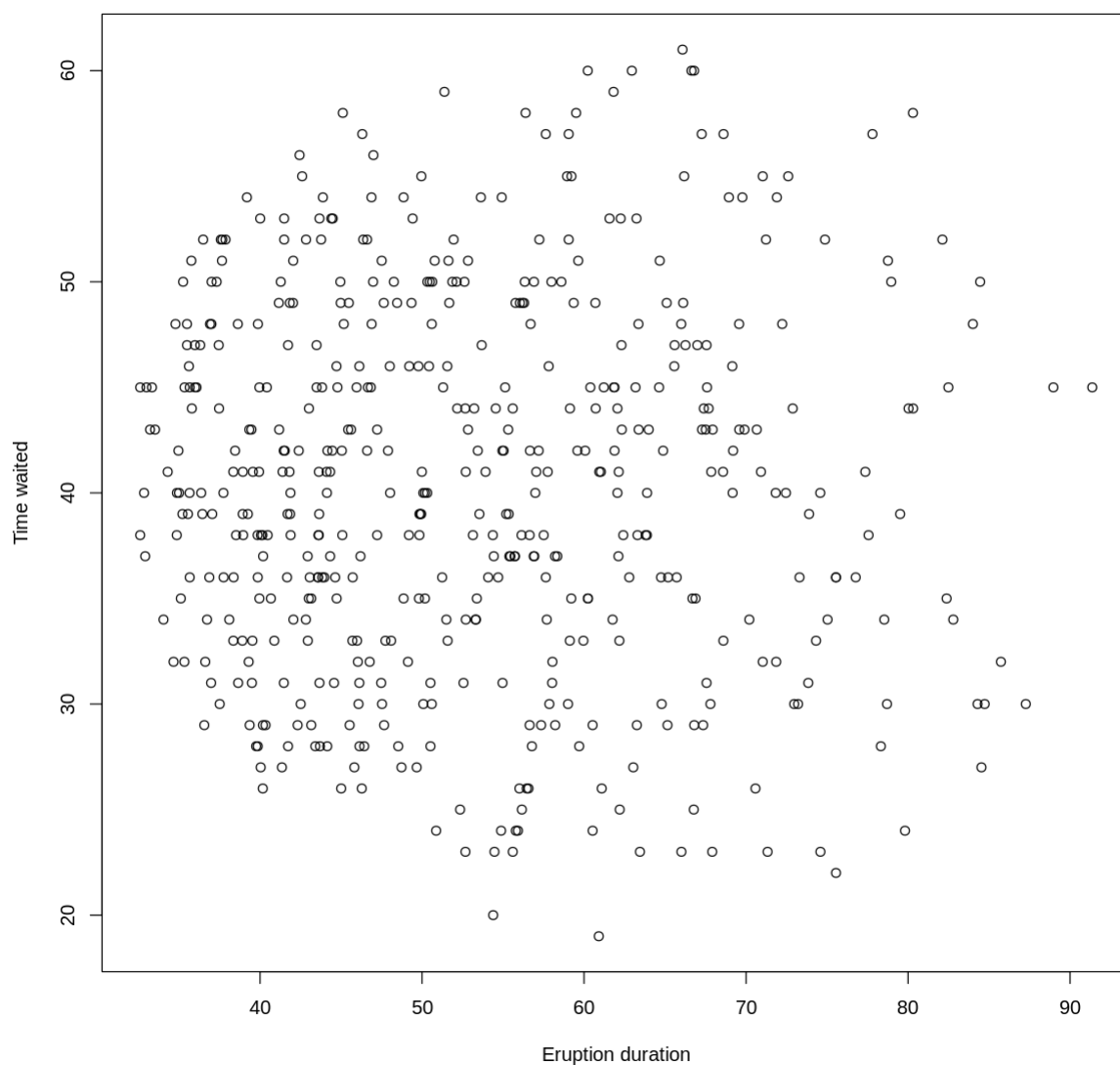
```
plot(clicked$Age, clicked$Area.Income, xlab="Age", ylab="Area.Income")
```



There seems to be no correlation between the two variables.

#A scatter plot of Age vs Daily Time Spent on Site.

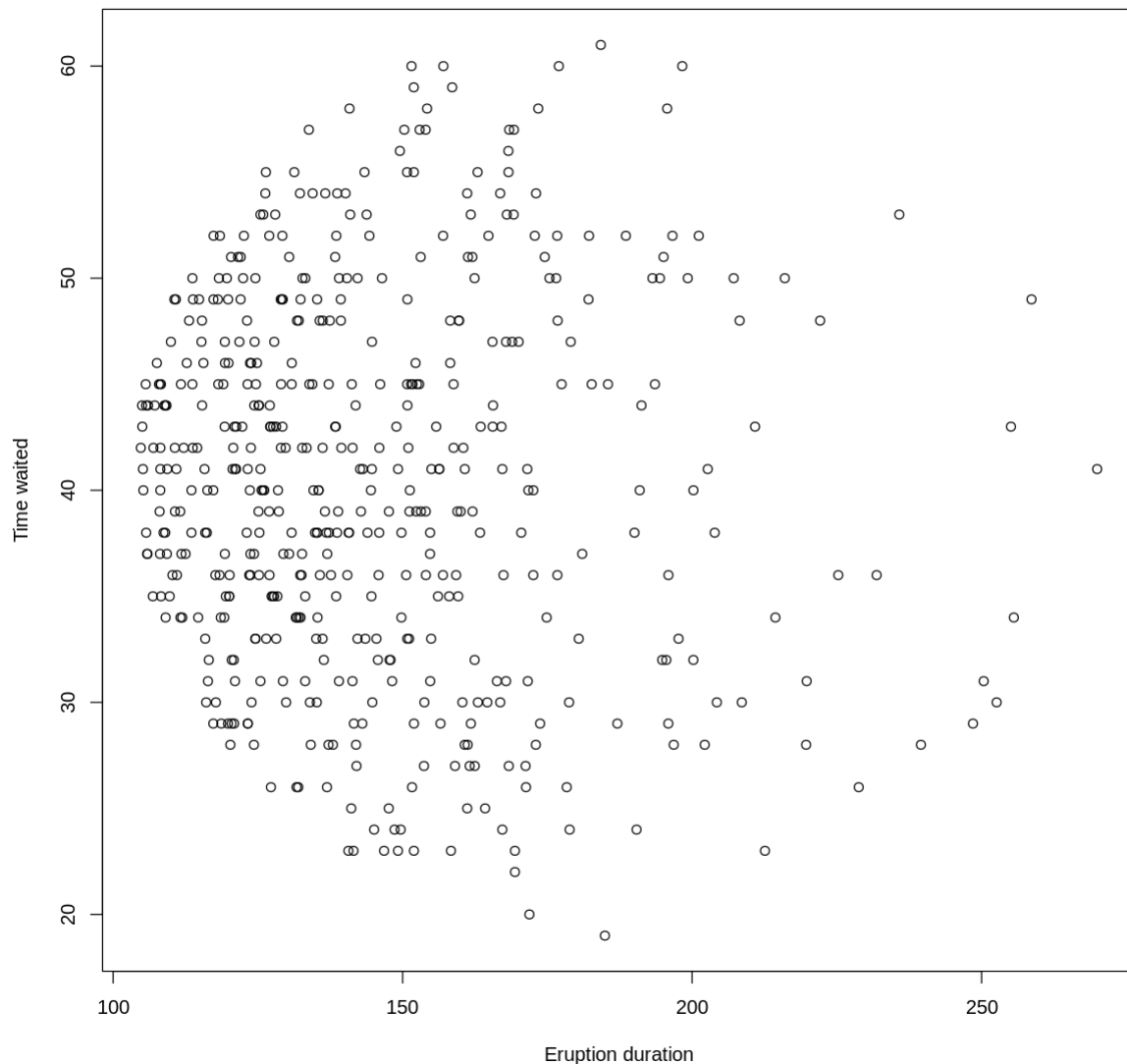
```
plot(clicked$Daily.Time.Spent.on.Site, clicked$Age, xlab="Daily.Time.Spent.on.Site", ylab="Age")
```



There seems to be no correlation between the two variables.

#A scatter plot of Age vs Daily.Internet.Usage.

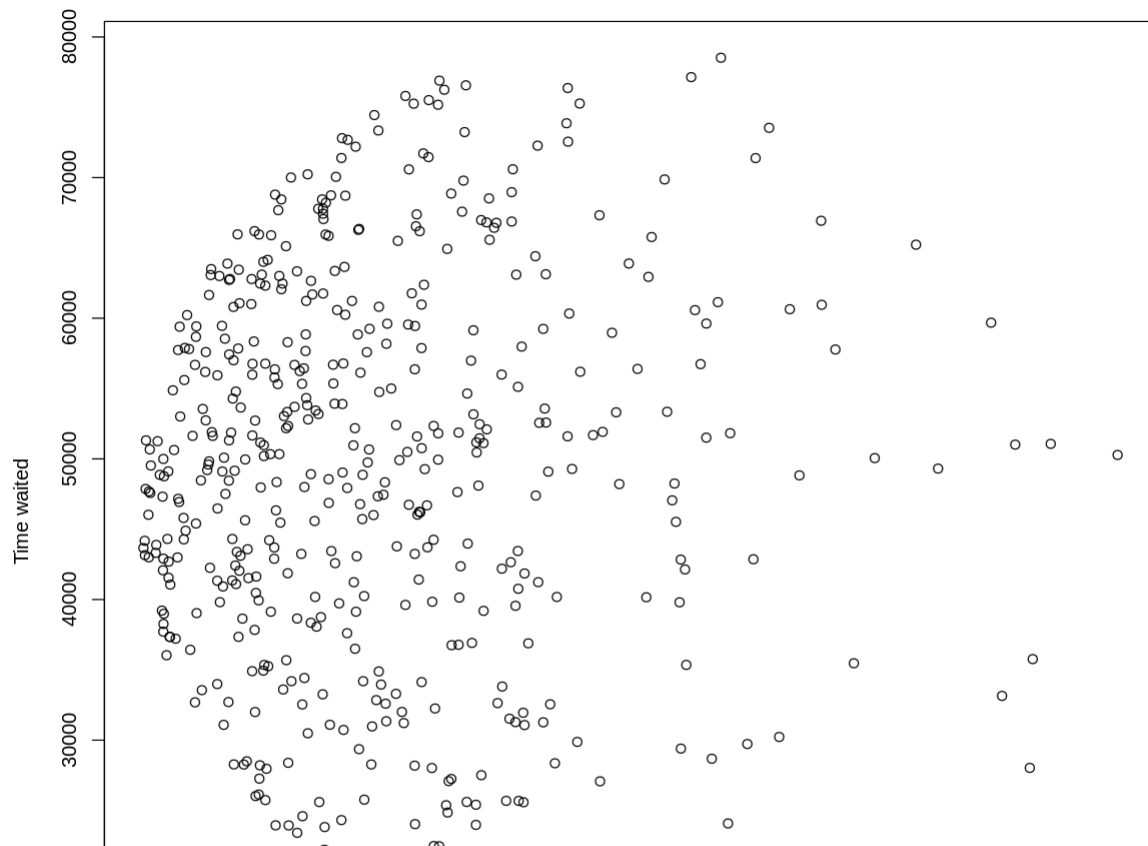
```
plot(clicked$Daily.Internet.Usage, clicked$Age, xlab="Daily.Internet.Usage", ylab="Age")
```



There seems to be no correlation between the two variables.

#A scatter plot of Area.Income vs Daily.Internet.Usage.

```
plot(clicked$Daily.Internet.Usage, clicked$Area.Income, xlab="Daily.Internet.Usage", ylab="Ar
```



There seems to be no correlation between the two variables.

```
#A scatter plot of Daily.Internet.Usage vs Daily Time Spent on Site.  
plot(clicked$Daily.Internet.Usage, clicked$Daily.Time.Spent.on.Site, xlab="Daily.Internet.Usa
```



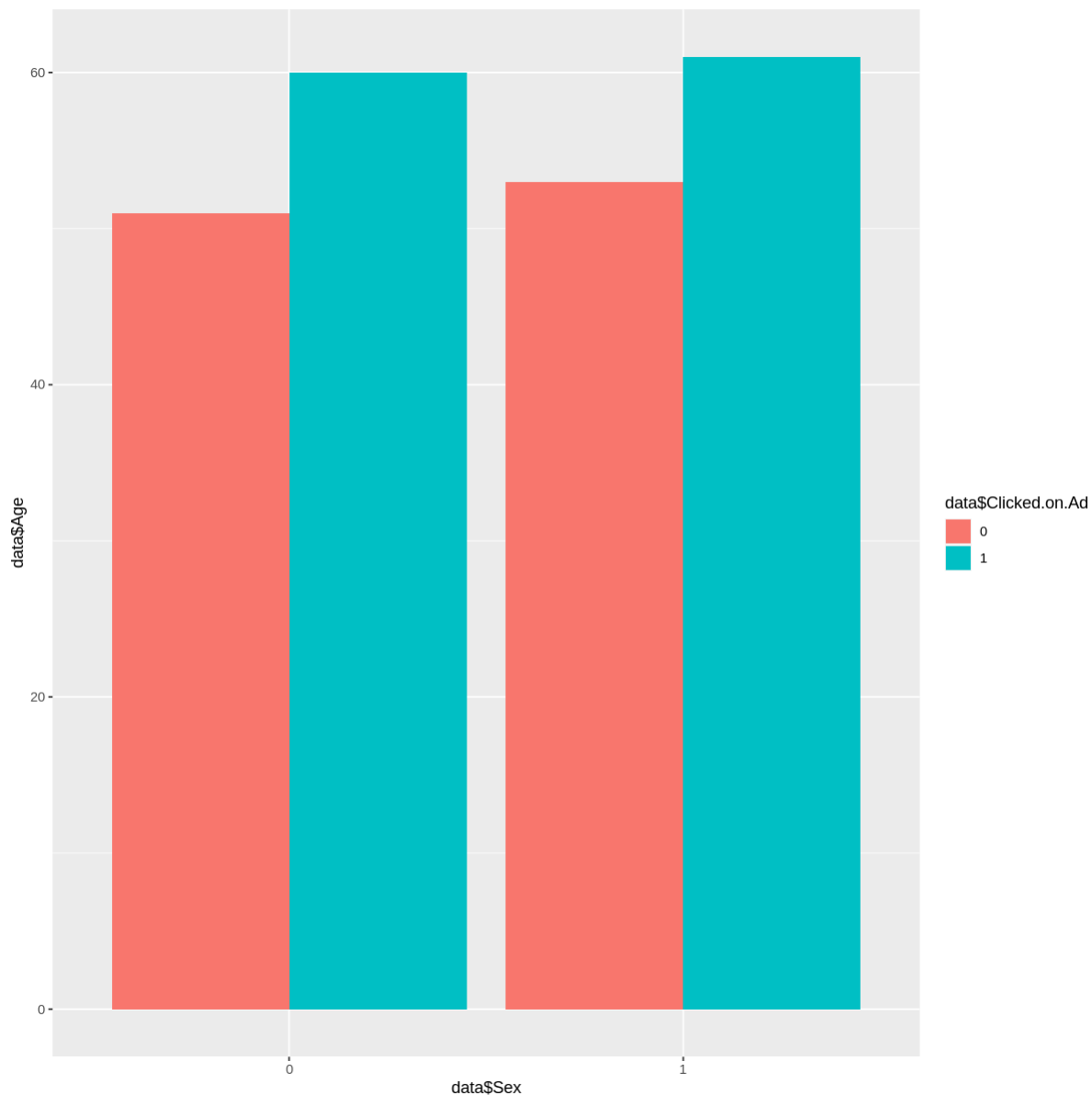

```
plot(clicked$Daily.Time.Spent.on.Site, clicked$Area.Income, xlab="Area.Income", ylab="Time wa
```



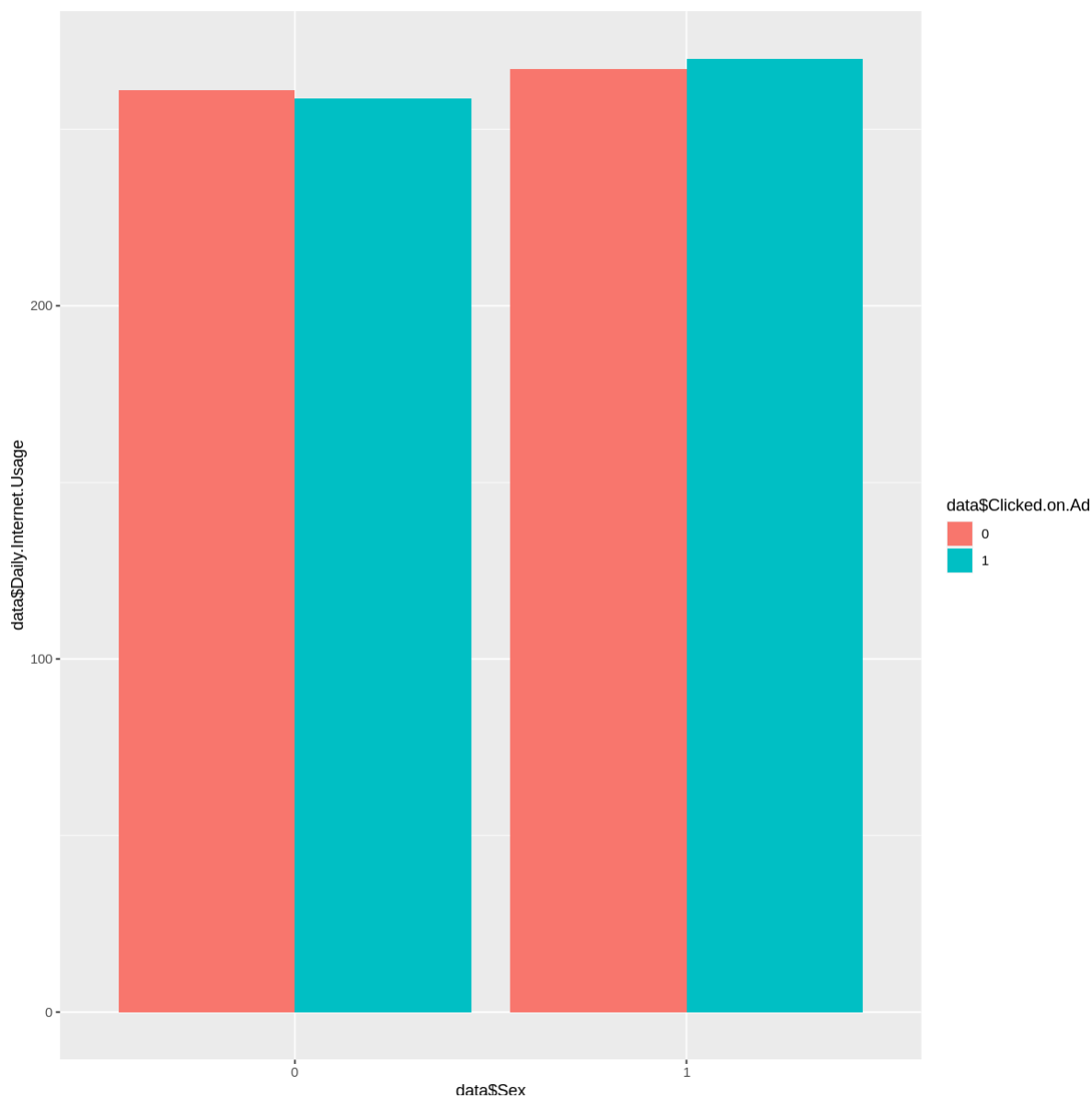
There seems to be no correlation between the two variables.

▼ Multivariate Analysis.

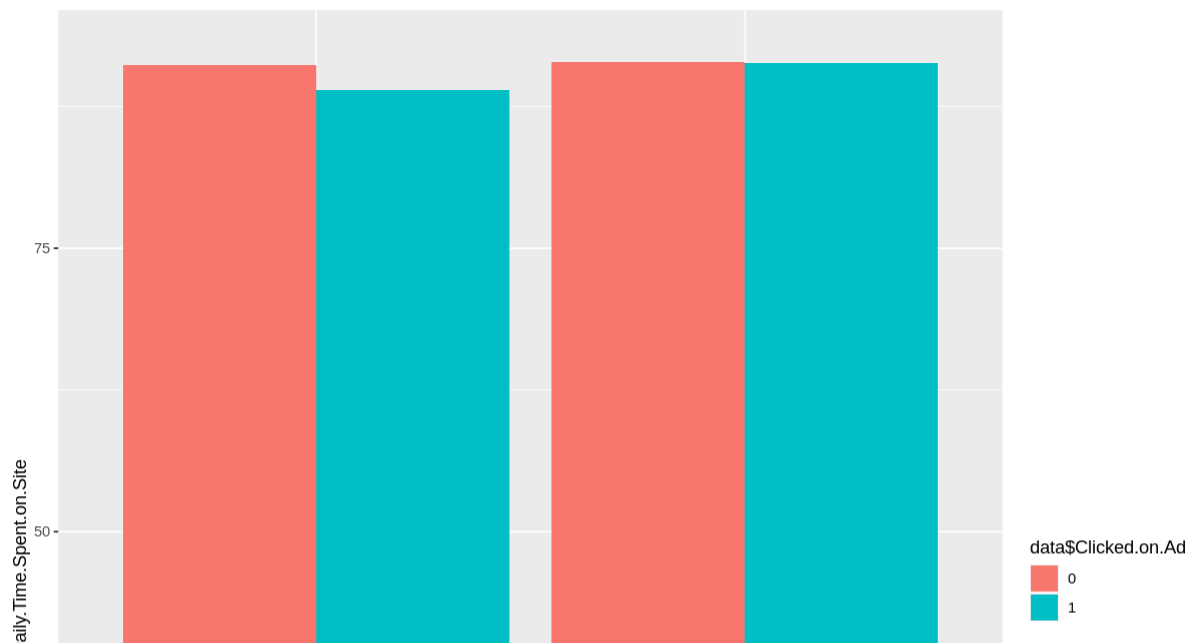
```
#A multivariate plot showing the relationship between Age, Sex and Clicked.on.Ad.  
library(ggplot2)  
ggplot(data, aes(fill=data$Clicked.on.Ad, y=data$Age, x=data$Sex)) +  
  geom_bar(position="dodge", stat="identity")
```



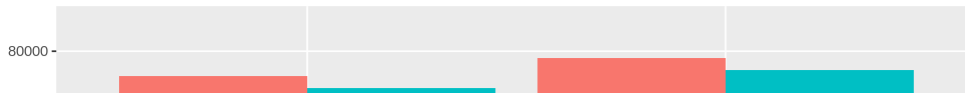
```
#A multivariate plot showing the relationship between Daily.Internet.Usage, Sex and Clicked.c  
ggplot(data, aes(fill=data$Clicked.on.Ad, y=data$Daily.Internet.Usage, x=data$Sex)) +  
  geom_bar(position="dodge", stat="identity")
```



#A multivariate plot showing the relationship between Daily.Time.Spent.on.Site, Sex and Click
ggplot(data, aes(fill=data\$Clicked.on.Ad, y=data\$Daily.Time.Spent.on.Site, x=data\$Sex)) +
geom_bar(position="dodge", stat="identity")



#A multivariate plot showing the relationship between Area.Income, Sex and Clicked.on.Ad.
ggplot(data, aes(fill=data\$Clicked.on.Ad, y=data\$Area.Income, x=data\$Sex)) +
geom_bar(position="dodge", stat="identity")



▼ Conclusions and Recommendations



Conclusions

- Most people who clicked on the ad are female.
- Most people who clicked on the ad are in their 30's, 40's and 50's.
- People from Lake David and Lake James are the most frequent in terms of clicking the ad.
- People from Australia are found to click the ad the most.
- The most frequent month in which the ad was clicked is February.
- The most frequent day of the month is day 3.
- The most frequent hour is 9:00 am.



Recommendations

I would recommend the Kenyan entrepreneur to Target the following market:

- Mostly females since they seem to click her ad the most.
- People in the age of 30's, 40's and 50's since they are the ones who are the most interested.
- People from cities like King David and King James.
- People from Australia and other leading countries in terms of clicking the ad.
- She should also consider the time, day and month that people are most likely to click her ad such as, February, day 3 of the month and also 9:00 am.

