# First Independent Project

STEPHEN ODHIAMBO OGAJA

2022-05-28

## CRYPTOGRAPHY ADVERTISING

## ## 1a). Defining the Question

###–» Which individuals are more likely to click on adverts on cryptography?

### b). Defining the Metric of Success

###–» The project will be considered a success when we can identify which individuals will click on the advert.

### c). Understanding 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.

### d). Recording the Experimental Design

###–» (i) Find and deal with outliers, anomalies, and missing data within the dataset. (ii) Perform uni variate and bivariate analysis. (iii) From your insights provide a conclusion and recommendation

### e). Data Relevance

###–» The data is valid and has been provided by the entrepreneur, it was collected from the previous adverts.

# # 2. Reading the data

**let's import the dataset**

```
adverts <- read.csv("advertising.csv")
```

# # 3. Checking the Data

let's preview the top 6 records of the dataset

```
head(adverts)
```

```
##   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
##                            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
##               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
```

let's check the last 6 records of the dataset

```
tail(adverts)
```

```
##      Daily.Time.Spent.on.Site Age Area.Income Daily.Internet.Usage
## 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
## 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
## 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
```

**let's see the shape of our dataset**

```
dim(adverts)
```

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

###–» The dataframe has 1000 observations and 10 variables

**let's see the data types of the variables**

```
str(adverts)
```

```
## '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 ...
```

###–» R stores the dataframe and views the variables as lists so to see the the various data types of this list we use the str function.

**let's check for duplicates in the dataframe**

```
duplicates <- adverts[duplicated(adverts), ]
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
## <0 rows> (or 0-length row.names)
```

3

**−» The dataframe does not contain duplicate values.**

**let's check for missing data in each column**

```
colSums(is.na(adverts))
```

```
## 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
##                        0
```

###−» The dataset's columns does not have missing data.

**let's check for outliers in the dataset**

**selecting only numeric columns**

```
num_cols <- adverts[,unlist(lapply(adverts, is.numeric))]
head(num_cols)
```

```
##   Daily.Time.Spent.on.Site Age Area.Income Daily.Internet.Usage Male
## 1                    68.95  35    61833.90               256.09    0
## 2                    80.23  31    68441.85               193.77    1
## 3                    69.47  26    59785.94               236.50    0
## 4                    74.15  29    54806.18               245.89    1
## 5                    68.37  35    73889.99               225.58    0
## 6                    59.99  23    59761.56               226.74    1
##   Clicked.on.Ad
## 1             0
## 2             0
## 3             0
## 4             0
## 5             0
## 6             0
```
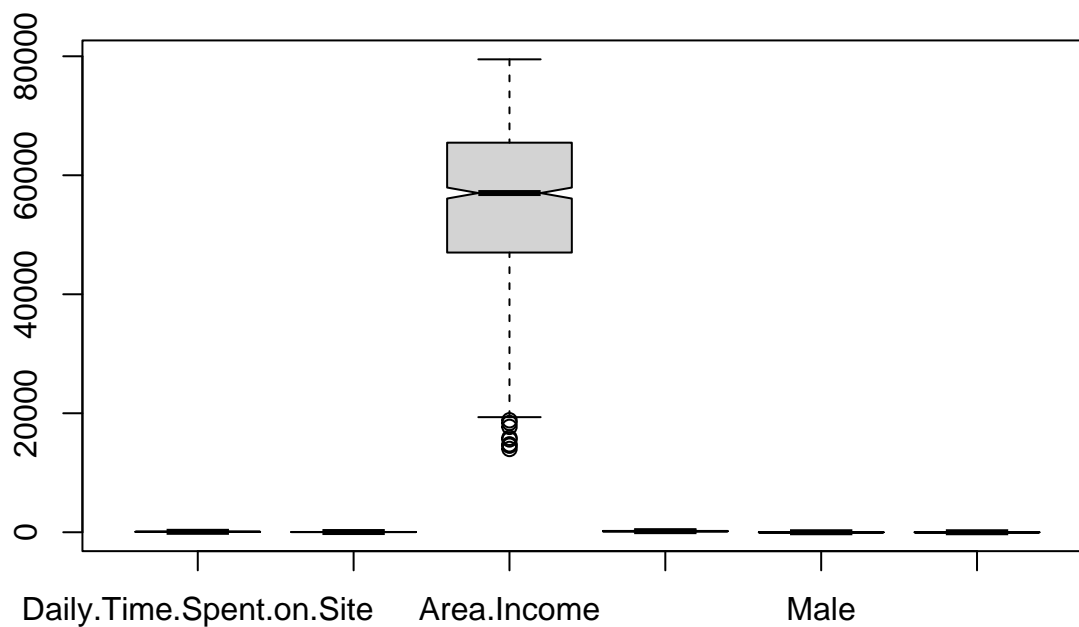
###−» 6 columns are numerical in nature

**let's check for outliers in the numerical columns using BOXPLOT**

```
boxplot(num_cols, notch = TRUE)
```

```
## Warning in (function (z, notch = FALSE, width = NULL, varwidth = FALSE, : some
## notches went outside hinges ('box'): maybe set notch=FALSE
```

###–» The Area.Income variable has outliers which will be imputed.

**let's see the values which are outliers in the Area.Income variable**

```
boxplot.stats(adverts$Area.Income)$out
```

```
## [1] 17709.98 18819.34 15598.29 15879.10 14548.06 13996.50 14775.50 18368.57
```

**let's check for outliers using Z-SCORES**

**The z-score indicates the number of standard deviations a given value deviates from the mean.**

```
z_scores <- as.data.frame(sapply(num_cols, function(num_cols) (abs(num_cols-mean(num_cols))/sd(num_cols)
head(z_scores)
```

```
##   Daily.Time.Spent.on.Site       Age Area.Income Daily.Internet.Usage      Male
## 1                0.2491419 0.1148475  0.50943618            1.7331628 0.9622138
## 2                0.9606516 0.5701399  1.00202882            0.3136484 1.0382307
## 3                0.2819420 1.1392555  0.35677007            1.2869451 0.9622138
## 4                0.5771428 0.7977862  0.01444841            1.5008289 1.0382307
## 5                0.2125572 0.1148475  1.40816290            1.0382112 0.9622138
## 6                0.3160289 1.4807248  0.35495265            1.0646335 1.0382307
```

5

```
##   Clicked.on.Ad
## 1    0.9994999
## 2    0.9994999
## 3    0.9994999
## 4    0.9994999
## 5    0.9994999
## 6    0.9994999
```

###–» We will drop values with a Z-Score of more than 3 or -3. They are the outliers

**Removing the outliers**

```
no_outliers <- z_scores[!rowSums(z_scores>3), ]
head(no_outliers)
```

```
##   Daily.Time.Spent.on.Site       Age Area.Income Daily.Internet.Usage      Male
## 1                0.2491419 0.1148475  0.50943618            1.7331628 0.9622138
## 2                0.9606516 0.5701399  1.00202882            0.3136484 1.0382307
## 3                0.2819420 1.1392555  0.35677007            1.2869451 0.9622138
## 4                0.5771428 0.7977862  0.01444841            1.5008289 1.0382307
## 5                0.2125572 0.1148475  1.40816290            1.0382112 0.9622138
## 6                0.3160289 1.4807248  0.35495265            1.0646335 1.0382307
##   Clicked.on.Ad
## 1    0.9994999
## 2    0.9994999
## 3    0.9994999
## 4    0.9994999
## 5    0.9994999
## 6    0.9994999
```

**let's check the number of observations after removing outliers**
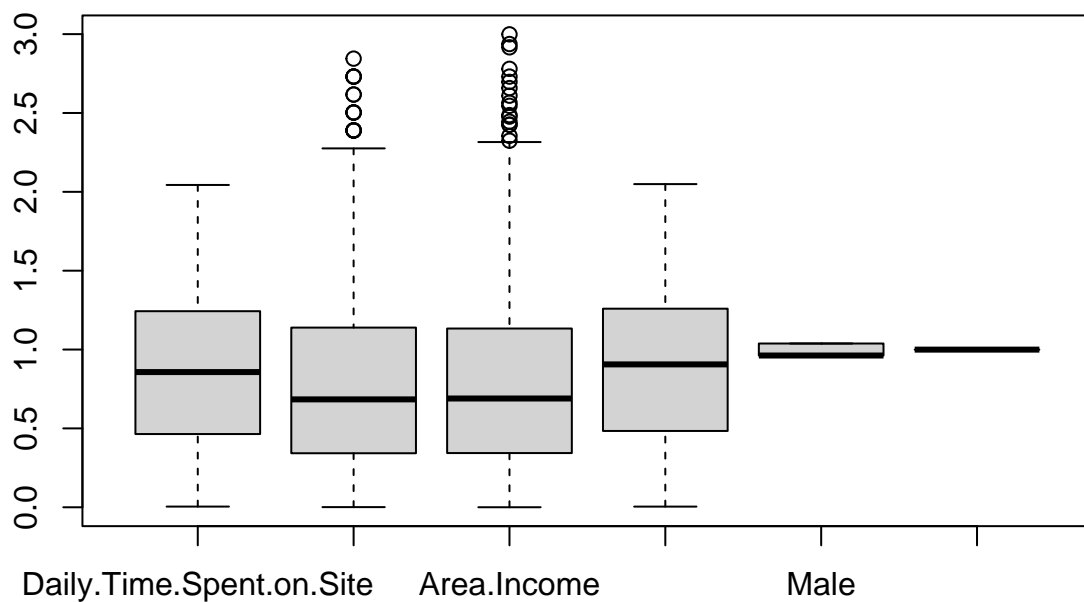
```
dim(num_cols)
```

```
## [1] 1000    6
```

```
dim(no_outliers)
```

```
## [1] 998    6
```

###–» We removed 2 observations.

**let's check for outliers in the new dataframe after removing them**

```
boxplot(no_outliers)
```

6

###-» There are still outliers so we will use interquantile range method to remove outliers

**checking and removing outliers using IQR**

**The Area.Income column had outliers so we focus on it**

```
income.IQR <- 65471-47032
income.IQR <-IQR(adverts$`Area.Income`)
income.IQR
```

```
## [1] 18438.83
```

**let's save the dataframe without outliers into a new dataframe by assigning it to a variable**

```
adverts_2 <- subset(adverts, adverts$`Area.Income`> (47032 - 1.5*income.IQR) & adverts$`Area.Income`<(65
```

**let's see the shape of the new dataframe**

```
dim(adverts_2)
```

```
## [1] 991  10
```

###–» We have lost 9 observations that included the outliers. We proceed with analysis.

# ### 4. {UNIVARIATE ANALYSIS}

**let's get the mean of the numerical columns**

```
summary(num_cols)
```

```
##  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
##       Male           Clicked.on.Ad
##  Min.   :0.000   Min.   :0.0
##  1st Qu.:0.000   1st Qu.:0.0
##  Median :0.000   Median :0.5
##  Mean   :0.481   Mean   :0.5
##  3rd Qu.:1.000   3rd Qu.:1.0
##  Max.   :1.000   Max.   :1.0
```

###–» The summary shows: ###–» 1. The minimum value for each numerical variable. ###–» 2. The first quantile for each numerical variable ###–» 3. The median value for all numeric variables across the dataframe. ###–» 4. The mean value for all numeric variables. ###–» 5. The third quantile. ###–» 6. The maximum value for all numerical columns.

**let's get the variance for the numeric variables**

```
variance <- var(num_cols)
variance
```

```
##                          Daily.Time.Spent.on.Site          Age     Area.Income
## Daily.Time.Spent.on.Site                251.3370949 -4.617415e+01   6.613081e+04
## Age                                     -46.1741459  7.718611e+01  -2.152093e+04
## Area.Income                           66130.8109082 -2.152093e+04   1.799524e+08
## Daily.Internet.Usage                    360.9918827 -1.416348e+02   1.987625e+05
## Male                                     -0.1501864 -9.242142e-02   8.867509e+00
## Clicked.on.Ad                            -5.9331431  2.164665e+00  -3.195989e+03
##                          Daily.Internet.Usage        Male Clicked.on.Ad
## Daily.Time.Spent.on.Site         3.609919e+02 -0.15018639  -5.933143e+00
## Age                             -1.416348e+02 -0.09242142   2.164665e+00
## Area.Income                      1.987625e+05  8.86750903  -3.195989e+03
## Daily.Internet.Usage             1.927415e+03  0.61476667  -1.727409e+01
## Male                             6.147667e-01  0.24988889  -9.509510e-03
## Clicked.on.Ad                   -1.727409e+01 -0.00950951   2.502503e-01
```

###–» variance is a measure of how far the set of data points per column is spread out from their mean eg. those of the area income seem to be far spread out from their mean when compared to that of the age column.

**let's get the standard deviation of the numeric variables**

**let's create a function to get the standard deviations**

```
sd.function <- function(column) {
  standard.deviations <- sd(column)
  print(standard.deviations)
}
```

**standard deviation for daily time spent on site**

```
sd.function(adverts_2$Daily.Time.Spent.on.Site)
```

```
## [1] 15.9005
```

**standard deviation for Age**

```
sd.function(adverts_2$Age)
```

```
## [1] 8.804716
```

**standard deviation for Area.Income**

```
sd.function(adverts_2$Area.Income)
```

```
## [1] 12961.5
```

**standard deviation for Daily.Internet.Usage**

```
sd.function(adverts_2$Daily.Internet.Usage)
```

```
## [1] 44.05386
```

###–» Where a low standard deviation indicates that values are closer to the mean a high one indicates the standard deviation is far from the mean e.g the age column standard deviation of 8.8 displays that its values are closer to their mean than that of the Area income column whose value is 12961

**let's get the skewness of the numerical column**

```
library(moments)
skewness(num_cols)
```

```
## Daily.Time.Spent.on.Site                    Age              Area.Income
##              -0.37120261             0.47842268             -0.64939670
##        Daily.Internet.Usage                 Male             Clicked.on.Ad
##              -0.03348703             0.07605493              0.00000000
```

###–» The skewness of the Age variable being positive indicates that its distribution has a longer right tail than left tail while the rest of the columns' left tails.

# ### 5. {BIVARIATE ANALYSIS}

**let's get the covariance of the numeric variables**

```
cov(num_cols)
```

```
##                          Daily.Time.Spent.on.Site          Age    Area.Income
## Daily.Time.Spent.on.Site                251.3370949 -4.617415e+01  6.613081e+04
## Age                                     -46.1741459  7.718611e+01 -2.152093e+04
## Area.Income                           66130.8109082 -2.152093e+04  1.799524e+08
## Daily.Internet.Usage                    360.9918827 -1.416348e+02  1.987625e+05
## Male                                     -0.1501864 -9.242142e-02  8.867509e+00
## Clicked.on.Ad                            -5.9331431  2.164665e+00 -3.195989e+03
##                          Daily.Internet.Usage        Male Clicked.on.Ad
## Daily.Time.Spent.on.Site        3.609919e+02 -0.15018639 -5.933143e+00
## Age                            -1.416348e+02 -0.09242142  2.164665e+00
## Area.Income                     1.987625e+05  8.86750903 -3.195989e+03
## Daily.Internet.Usage            1.927415e+03  0.61476667 -1.727409e+01
## Male                            6.147667e-01  0.24988889 -9.509510e-03
## Clicked.on.Ad                  -1.727409e+01 -0.00950951  2.502503e-01
```

###–» The age variable is the only column with a positive covariance with the ad click variable, the rest have negative covariances.

**let's get the correlation coefficient**

```
cor(num_cols)
```

```
##                          Daily.Time.Spent.on.Site          Age   Area.Income
## Daily.Time.Spent.on.Site              1.00000000 -0.33151334  0.310954413
## Age                                  -0.33151334  1.00000000 -0.182604955
## Area.Income                           0.31095441 -0.18260496  1.000000000
## Daily.Internet.Usage                  0.51865848 -0.36720856  0.337495533
```

```
## Male                             -0.01895085 -0.02104406  0.001322359
## Clicked.on.Ad                     -0.74811656  0.49253127 -0.476254628
##                    Daily.Internet.Usage       Male Clicked.on.Ad
## Daily.Time.Spent.on.Site       0.51865848 -0.018950855   -0.74811656
## Age                           -0.36720856 -0.021044064    0.49253127
## Area.Income                    0.33749553  0.001322359   -0.47625463
## Daily.Internet.Usage          1.00000000  0.028012326   -0.78653918
## Male                           0.02801233  1.000000000   -0.03802747
## Clicked.on.Ad                 -0.78653918 -0.038027466    1.00000000
```

###–» The variables have a negative correlation with the target variable apart from the age variable which has a positive correlation. Let's see that in the correlogram below
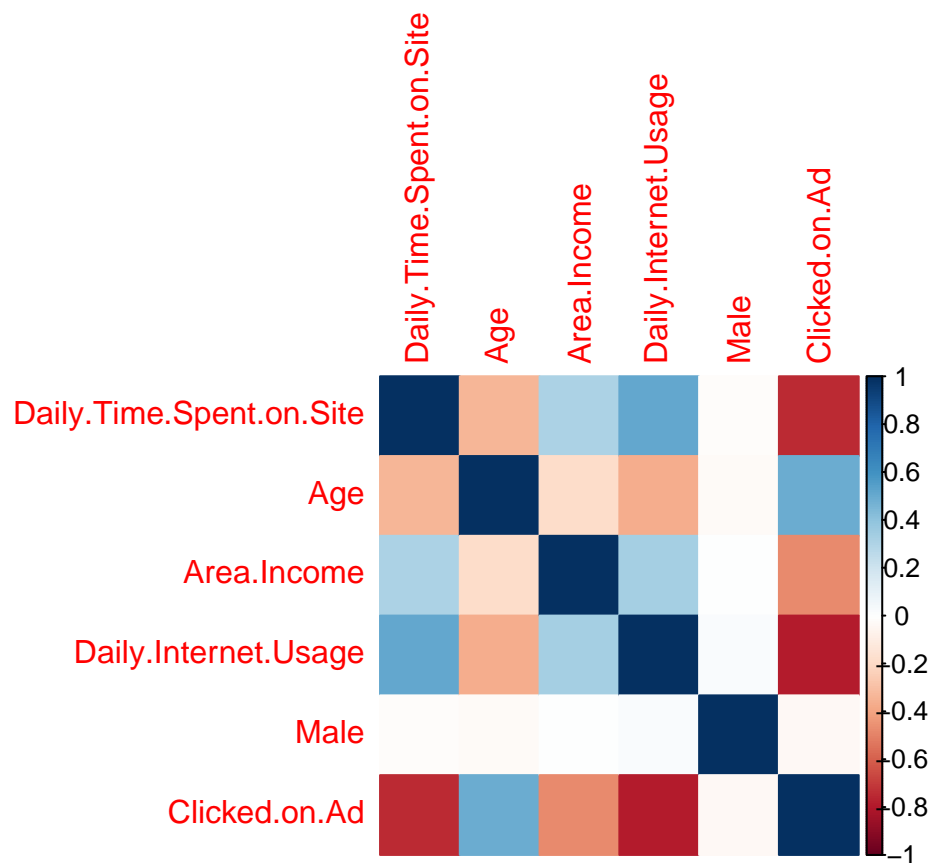
**let's see the corrplot of the numeric variables**

```
library(corrplot)
```

```
## corrplot 0.92 loaded
```

```
corr_ <- cor(num_cols)
```

```
corrplot(corr_, method = 'color')
```

### {RECOMMENDATIONS}

a. The entrepreneur should focus on the older population as the correlation between age and advert clicks is slightly positive indicating that as age increases the more likely the clicks are made.

b. The entreoreneur should focus on regions with bigger area coverage as those with a smaller area since the correlation between area income and advert clicks is negatively weak one indicating that as area income decreases the more likely the clicks are made and vice versa.

c. She should focus on the regions with low daily internet usage because the correlation between the daily internet usage and clicks on ads is negative indicating that as internet use decreases the more likely the clicks will be made.