

Student Number: D20124995

Name: Mina Jamshidian

Programme Code: Probability And Statistical Inference Math 9102(TU059)

For this project the dataset of data mining assignment was used.(Bank Marketing)

Table of contents:	1
Section 1 - Research Question(s)	3
Research Question	3
Hypothesis	3
Section 2 - Dataset	3
Statistical Measurement	6
Duration variable:	6
Report of normality analysis for duration variable:	8
Transforming the duration variable:	8
Report of normality analysis for transformed duration variable:	9
Has the client subscribed to a term deposit?	10
Age of client:	11
Clients Education:	12
Section 3 - Result	13
Hypothesis 1:	13
Statistical Evidence	13
Check the difference in the last call duration for clients by different age groups.	13
Report of Difference Analysis:	16
Checking the difference in last duration call leads clients to subscribed a term deposit.	16
Report of Difference Analysis	16
Model 1	17
Report of Linear Modelling Analysis	21
Hypothesis 2:	21
Statistical Evidence	21
Check difference in the last call duration for clients with different education	22
Report of Difference Analysis:	24
Model 2	24
Report of Linear Modelling Analysis	28
Summary Comparison	28
Model Comparison Results	28
Section 4 – Discussion/Conclusion	29
Reference	29

Section 1 - Research Question(s)

Research Question

This analysis is about predicting whether a client with different individual and social characteristics will subscribe to term deposit or not based on last contact duration?

This analysis would help us to find which factors could be considered to be perfect for prediction of determining a subscription situation.

Furthermore this is a project that concerns whether clients' education, their age, last contact duration with two contact types can be influenced on subscription of term deposit by clients.

Hypothesis

- Hypothesis 1:
 - H0: There will be no significant predictor for last call duration that leads clients to subscribed term deposits by different age groups?
 - H1: There will be a significant predictor for last call duration that leads clients to subscribed term deposits by different age groups?
- Hypothesis 2:
 - H0: There will be no significant predictor for last call duration that leads clients to subscribed term deposits by different age and different education?
 - H1: There will be a significant predictor for last call duration that leads clients to subscribed term deposits by different age and different education?

Section 2 - Dataset

For this project the dataset of data mining assignment was used which is [bank marketing](#), this dataset is about the Portuguese banking direct marketing, the marketing was based on calling with phone. The result of this marketing shows whether the customers will subscribe to term deposit or not.

21 variables exist in a dataset that has 11 categorical variables and 9 numerical variables with 41188 rows. The y variable is the target variable which should be predicted. The datasets variables were divided to following parts:

1. Bank client data: 1- age, 2- job, 3- marital, 4- education, 5- default, 6- housing, 7- loan.
2. Depends on the last contact of the current campaign: 8- contact, 9- month, 10- day_of_week, 11- duration.

3. Other attributes: 12- campaign, 13- pdays, 14- previous, 15- poutcome.
4. Social and Economic Context Attributes: 16- emp.var.rate, 17- cons.price.idx, 18-cons.conf.idx, 19- euribor3m, 20- nr.employed.
5. Target Variable: 21- y

The dataset is csv format which is *bank-additional-full.csv*, it is imported to the variable with the name of the **bank**.

Loading Libraries:

```
library(lmSupport)
library(stargazer)
library(dplyr)
library(rcompanion)
library(pastecs)
library(lavaan)
library(pander)
library(psych)
library(ggplot2)
library(varhandle)
library(effsize)
library(Hmisc)
library(stringr)
library(lmtest)
library(sampler)
library(gmodels)
library(readr)
library(ROSE)
library(tidyr)
library(gplots)
library(lsr)
library(stargazer)
library(car)
```

Importing Data:

```
# importing the bank marketing data csv file
bank = read.csv("bank-additional-full.csv", sep = ";")
```

After that it would be better to treat the dataset and preparing some variables for further analysing as explained in following part:

1. Creating two new variables related to the project question:
The age_label variable which is dividing the age variable in five categories (1. age > 60 for senior-citizen, 2. age >45 for mid-old 3. age>30 for Mid-age 4. age>15 for Young 5. age<15 for Children)

```
bank= bank%>% mutate(age_label= if_else(age >
60,"senior-citizen",if_else(age>45,"mid-old",if_else(age>30,"Mid-age",if_else(age>15,"Young
","Children"))))
```

2. Finding the unknown values then doing the best solution for each variable based on their correlation by the target variable by using the cross table function("CrossTable function | R Documentation", 2021).

```
#searching for unknowns throughout the dataset  
colSums(bank == "unknown")
```

```
   age      job    marital  education  default  housing    loan  
    0     330        80      1731     8597      990     990  
contact    month  day_of_week  duration  campaign    pdays  previous  
    0         0         0         0         0         0         0  
poutcome  emp.var.rate  cons.price.idx  cons.conf.idx  euribor3m  nr.employed    y  
    0         0         0         0         0         0         0
```

Figure 1

It was found the *unknown* values for the default (8579), education (1731), housing (990), loan (990) and job (330) and marital (80) variables as it shows in Figure 1.

It was decided to remove the *unknown values* for job, marital and housing then change the *unknown* value for education variable to *university.degree* value based on the correlation with target variable(y) as it is obvioused in Figure 2 and because of the number unknown value for default is a lot, the column was removed.

At the end it was checked again for unknown value ,the result was shown in Figure 3.

```
CrossTable(bank$job,bank$y_new)  
CrossTable(bank$marital,bank$y_new)  
CrossTable(bank$education,bank$y_new)  
CrossTable(bank$housing,bank$y_new)  
  
bank = bank %>% filter( job!= "unknown")  
bank = bank %>% filter( marital!= "unknown")  
bank$education[bank$education=="unknown"]="university.degree"  
bank = bank %>% filter( housing!= "unknown")  
#removing column default & duration  
bank = bank[-c(5)]  
  
#Checking the unknown values again  
colSums(bank == "unknown")
```

bank\$education	bank\$y	no	yes	Row Total
basic.4y		3748	428	4176
		0.486	0.029	
		0.698	0.182	0.181
		0.583	0.092	
		0.691	0.018	
basic.6y		2184	188	2292
		2.423	19.088	
		0.918	0.002	0.056
		0.658	0.041	
		0.651	0.005	
basic.9y		5572	473	6045
		6.465	63.527	
		0.922	0.078	0.147
		0.152	0.002	
		0.135	0.011	
high.school		8484	1031	9515
		0.198	1.561	
		0.092	0.008	0.231
		0.232	0.222	
		0.286	0.025	
illiterate		14	4	18
		0.244	1.018	
		0.778	0.222	0.000
		0.000	0.001	
		0.000	0.000	
professional.course		4646	595	5243
		0.004	0.032	
		0.687	0.113	0.127
		0.227	0.228	
		0.113	0.014	
university.degree		10498	1678	12168
		0.292	65.317	
		0.063	0.137	0.295
		0.287	0.368	
		0.235	0.041	
unknown		1488	251	1731
		2.041	16.079	
		0.055	0.145	0.042
		0.048	0.054	
		0.036	0.006	
Column Total		36548	4648	41188
		0.687	0.113	

Figure 2

Figure 3

- After analysing it was found that the campaign column which is about the number of contacted time, it was found the customer after more than 10 times calling didn't subscribed the term deposit so it would be better to filter the dataset with less than 10 times.

```
#Filter the dataset by campaign of less than 10
bank= bank%>%filter(campaign<10)
```

The variables of interest used in this research are shown below in Figure 4:

Concept	Variable Name	Statistical Type	Possible Values
last contact duration, in seconds	duration	Numerical	Range 0 to 4918
client age	age_label	Categorical	1.age>60(senior-citizen) 2.age>45(mid-old) 3.age>30(Mid-age) 4. age>15(Young) 5.age<15(Children)
client subscribed a term deposit?	y	Categorical	yes and no
contact communication type	contact	Categorical	cellular, telephone
client education	education	Categorical	basic.4y to university.degree

Figure 4

In this step the sample should be chosen, the sample was chosen by Sampling Design & Analysis advised by (Ziegel & Lohr, 2000).

Choosing the sample came in the following part with the name of sbank2 which is the final data that we want to use for further analysis. And the final data balancing came in Figure 4.

```
size <- rsampcalc(nrow(sbank), e=5, ci=95, p=0.5, over=0.1)
sbank2<-ssamp(sbank, size, y_new, over=0.1)
sbank2 %>% count(y_new)
```

Figure 4

Statistical Measurement

Each of the variables of interest would be inspected. The numeric variables of interest, that represent the last contact duration with the customer, in seconds. it would be inspected for normality by checking standardised scores for skewness and kurtosis and considering the percentage of standardised scores for the variables fell outside of expected boundaries and creating histograms and QQ plots. Decisionsing about the skewness and the kurtosis came from the advice of (George & Mallory, 2011) that categorizing the distribution as normal when the relevant standardised scores of the skewness and the kurtosis fall in the range ± 2 with the advice of (Field, Miles & Field, 2012) which categorizing the distribution as normal when 95% of the the variable standardised scores fall within the ± 3.29 bounds for a dataset larger than 80 cases. The categorical variables summary statistics would be identified for analysis.

Duration variable:

Inspecting the Duration variable and its Normality by code:

```
# Descriptive statistics
# getting summary statistics for duration variable
duration<- sbank2$duration
duration_d<-describe(duration,omit = TRUE, IQR = TRUE)
duration_s<-list(pastecs::stat.desc(duration, basic = FALSE))

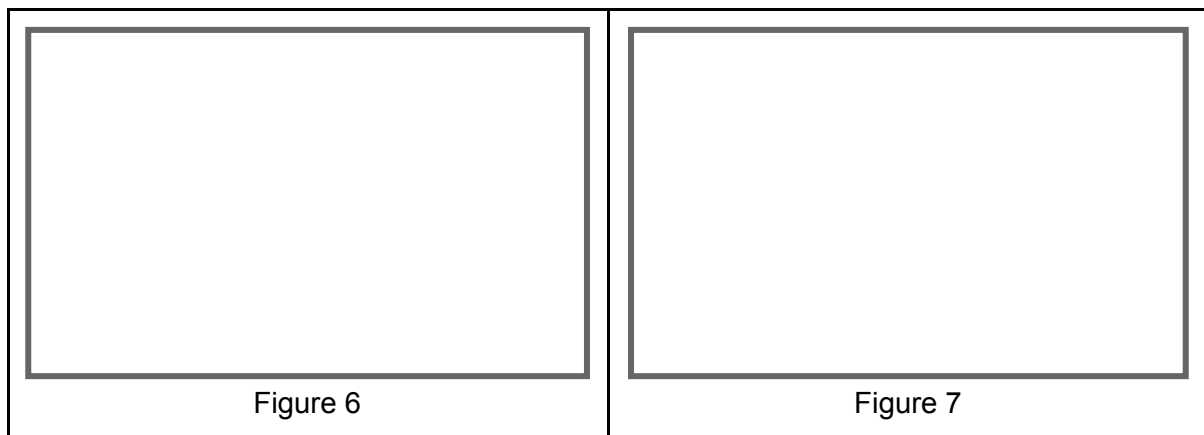
skew      <- semTools::skew(duration)
kurt      <- semTools::kurtosis(duration)
stdskew   <- skew[1] / skew[2]
stdkurt    <- kurt[1] / kurt[2]
zscore    <- abs(scale(duration))
gt196     <- FSA::perc(as.numeric(zscore), 1.96, "gt")
gt329     <- FSA::perc(as.numeric(zscore), 3.29, "gt")

duration_s$skew <- skew
duration_s$kurt <- kurt
duration_s$std.skew <- stdskew
duration_s$std.kurt <- stdkurt
duration_s$gt.196 <- gt196
duration_s$gt.329 <- gt329
duration_s

# Distribution of Age variable with visualization
ggplot(sbank2,aes(x=duration))+labs(x='duration', y='Density')+
  geom_histogram(binwidth = 30,colour='black',aes(y=..density..,fill=..count..))+
```

```
scale_fill_gradient("Count",low="#132B43", high="#56B1F7")+  
stat_function(fun = dnorm,color="red",args = list(mean=mean(duration,na.rm =  
TRUE),sd=sd(duration,na.rm = TRUE)))+  
ggtitle('Figure 1: Distribution of Duration')+  
theme(plot.title = element_text(size=10))  
  
# Create QQ Plot  
qqnorm(duration, main = "Figure 2: QQplot of Duration")  
qqline(duration,col=2) # show line on the plot
```

Figure 5



Report of normality analysis for duration variable:

The Duration is represented by a numeric variable in the dataset. Inspection of the standardised scores for skewness and kurtosis reveal that the kurtosis score (*kurtosis* = 44.5, *SE* = .227) and the skewness score (*skewness* = 24.95, *SE* = .113) is out of that range which is the range of range of -2 and 2 . This implies that kurtosis and skewness is not normal. For further inspection using plots such as histogram and normality plot (figure 8 and figure 8), we found that the distribution is positively skewed and not normalized. there are various values deviating from the normality line. On inspection of the count of outliers, there was found 2.1% standardised scores were outside the acceptable range of [-3.29, +3.29] that shows none of the values is outside the 95% Confidence Interval. In total, based on all the tests, it can be said that the data for duration variable will not be treated as a normal within this analysis (*Median*=176, *IQR*=185, *M* = 247.99, *SD* = 248.02, *N* = 462).

Transforming the duration variable:

After inspecting the duration variable, it was found that it was treated as non normal. In this analysis, the Log transformation was performed to this non normal data to convert it into the normal data before doing the linear parametric tests. This transformation approach was chosen because the Parametric test transformed normal data is considered more powerful compared to non parametric test on untransformed non normal data. So, the transformed duration variable would be used for result analysis and manipulation.

Transforming the duration variable and Checking Normality by code:

```
#Using the Log function for transforming duration
tduration<- log(sbank2$duration+10)

#adding the duration transformation to dataset table as new column by the name tduration
sbank2 <- sbank2 %>% mutate(tduration)

tduration_d<-describe(tduration,omit = TRUE, IQR = TRUE)
tduration_s<-list(pastecs::stat.desc(tduration, basic = FALSE))

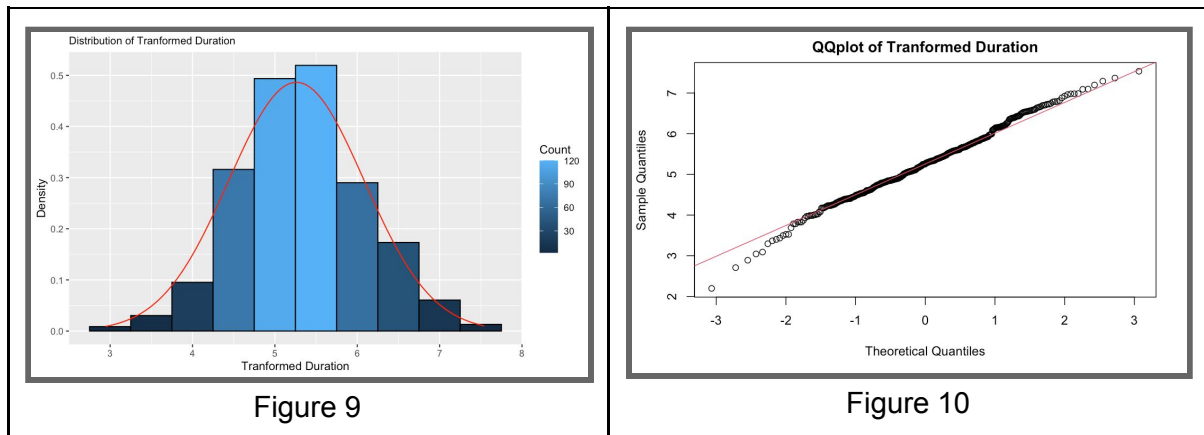
skew      <- semTools::skew(tduration)
kurt      <- semTools::kurtosis(tduration)
stdskew   <- skew[1] / skew[2]
stdkurt   <- kurt[1] / kurt[2]
zscore    <- abs(scale(tduration))
gt196     <- FSA::perc(as.numeric(zscore), 1.96, "gt")
gt329     <- FSA::perc(as.numeric(zscore), 3.29, "gt")

tduration_s$skew <- skew
tduration_s$kurt <- kurt
tduration_s$std.skew <- stdskew
tduration_s$std.kurt <- stdkurt
tduration_s$gt.196 <- gt196
tduration_s$gt.329 <- gt329
tduration_s

# Distribution of transformed duration(tduration) variable with visualization After
Transformation
ggplot(sbank2,aes(x=tduration))+labs(x='Tranformed Duration', y='Density')+
  geom_histogram(binwidth = 0.5,colour='black',aes(y=..density..,fill=..count..))+
  scale_fill_gradient("Count",low="#132B43", high="#56B1F7")+
  stat_function(fun = dnorm,color="red",args = list(mean=mean(tduration,na.rm =
TRUE),sd=sd(tduration,na.rm = TRUE)))+
  ggtitle('Distribution of Tranformed Duration')+
  theme(plot.title = element_text(size=10))

# Create QQ Plot
qqnorm(tduration, main = "QQplot of Tranformed Duration")
qqline(tduration,col=2) # show line on the plot
```


Figure 8



Report of normality analysis for transformed duration variable:

Transformed duration variable is represented by a numeric variable which was calculated by doing a Log function on the duration variable in the dataset. Inspecting the standardized scores for skewness (*skewness* = $-.62$, *SE* = $.113$) and kurtosis (*kurtosis* = 1.48 , *SE* = $.227$) shows that both of the skewness and the kurtosis value fall within the standardized score range of -2 and 2 which implies that both skewness and kurtosis are normal. For further inspection by using plots such as histogram and normality plot (Figure 9 and Figure 10), it was found that the distribution is normal. On inspection of the count of outliers, we found the 0.2% standardised scores were outside the acceptable range of $[-3.29, +3.29]$. which shows that none of the values is outside of the 95% Confidence Interval. In total, it was found base all the test that was done the data for transformed duration variable has a normal distribution by this analysis (*M* = 5.23 , *SD* = $.77$, *N* = 462).

Has the client subscribed to a term deposit?

Has the client subscribed a term deposit(y) is a nominal variable in the bank marketing dataset. The sample dataset contains data from 53 clients who subscribed 'Yes' and 409

clients who did not subscribe 'No'. The variable is representative of a sample which the clients will subscribe to the term deposit.

Inspecting the variable by code:

```
y<-table(sbank2$y)
y

# report basic summary statistics by a grouping variable
describeBy(tduration, sbank2$y)

#remove NA from if there is exist
dy<-data.frame(sbank2$y, tduration)
dy<-na.omit(dy)
names(dy)<-c("y", "duration")

# Create box plot for the variable
medhelp_graph<-ggplot(dy, aes(y, duration))
medhelp_graph+stat_summary(fun.y = mean, geom = "bar", fill="blue", colour="black", na.rm =
TRUE)+stat_summary(fun.data = mean_cl_normal, geom = "pointrange", na.rm =
TRUE)+labs(x="Client subscription of term deposit", y="Total Anxiety", title="Mean call
duration by Client subscription of term deposit")
```

Figure 14



Figure 15

Age of client:

Age of client is a categorical variable in the bank marketing dataset. The sample dataset has data from 9 clients in the senior-citizen group(age>60), 265 clients in the Mid-age group(age

>45), 117 clients in the mid-old group(age>30), 71 clients in the young group(age>15). The variable is representative of a sample which is a client from a Portuguese banking institution.

Inspecting the variable by code:

```
age_label<-table(sbank2$age_label)
age_label

# report basic summary statistics by a grouping variable
describeBy(tduration,sbank2$age_label)

#remove NA from if there is exist
dage<-data.frame(sbank2$age_label,tduration)
dage<-na.omit(dage)
names(dage)<-c("age_label","duration")

# Create box plot for the variable
medhelp_graph<-ggplot(dage,aes(age_label,duration))
medhelp_graph+stat_summary(fun.y = mean,geom = "bar",fill="#56B1F7",colour="black",na.rm =
TRUE)+stat_summary(fun.data = mean_cl_normal,geom = "pointrange",na.rm =
TRUE)+labs(x="Clients Age",y="Last call duration",title="Mean Last call duration by Clients
age")
```

Descriptive statistics by group													
group: Mid-age													
	vars	n	mean	sd	median	trimmed	mad	min	max	range	skew	kurtosis	se
X1	1	247	5.22	0.8	5.18	5.19	0.75	2.83	7.46	4.62	0.32	0.3	0.05

group: mid-old													
	vars	n	mean	sd	median	trimmed	mad	min	max	range	skew	kurtosis	se
X1	1	112	5.34	0.73	5.34	5.34	0.59	3.53	7.25	3.72	-0.03	0.21	0.07

group: senior-citizen													
	vars	n	mean	sd	median	trimmed	mad	min	max	range	skew	kurtosis	se
X1	1	11	5.58	0.74	5.26	5.58	0.64	4.43	6.74	2.31	0.16	-1.52	0.22

group: Young													
	vars	n	mean	sd	median	trimmed	mad	min	max	range	skew	kurtosis	se
X1	1	92	5.12	0.77	5.17	5.12	0.78	2.94	6.79	3.84	-0.14	-0.01	0.08

Figure 16

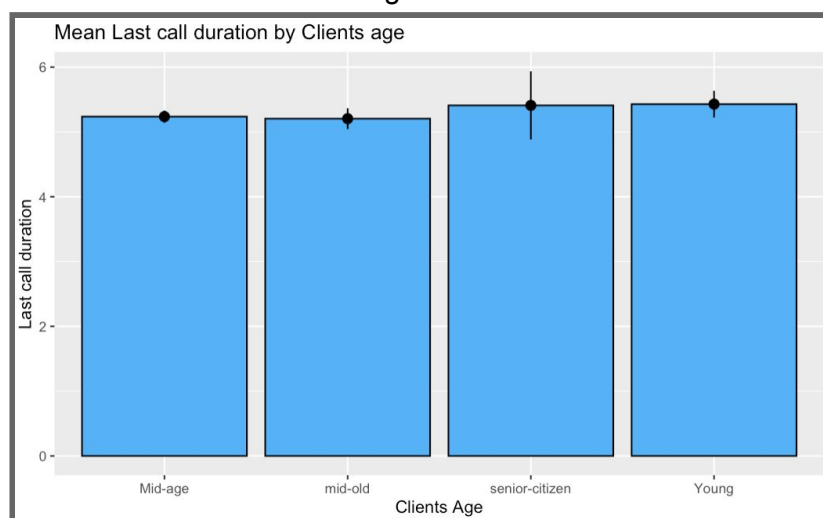


Figure 17

Clients Education:

Client Education is a categorical variable in the bank marketing dataset. The sample dataset has data from 42 clients in basic.4y group education, 30 clients in the basic.6y education group, 81 clients in the basic.9y education group, 110 clients in the high.school education group, 58 clients in the professional.course education group, 141 clients in the university.degree education group. The variable is representative of a sample which is a client from a Portuguese banking institution.

Inspecting the variable by code:

```
education<-table(sbank2$education)
education

# report basic summary statistics by a grouping variable
describeBy(tduration,sbank2$education)

#remove NA from if there is exist
dedu<-data.frame(sbank2$education,tduration)
dedu<-na.omit(dedu)
names(dedu)<-c("education","duration")

# Create box plot for the variable
medhelp_graph<-ggplot(dedu,aes(education,duration))
medhelp_graph+stat_summary(fun.y = mean,geom = "bar",fill="#56B1F7",colour="black",na.rm =
TRUE)+stat_summary(fun.data = mean_cl_normal,geom = "pointrange",na.rm =
TRUE)+labs(x="Clients Education",y="Last call duration",title="Mean Last call duration by
Clients Education")
```

```
Descriptive statistics by group
group: basic.4y
  vars  n mean  sd median trimmed mad min max range skew kurtosis se
X1     1 44  5.4 0.78   5.29   5.36 0.6 3.93 7.25  3.32 0.37  -0.15 0.12
-----
group: basic.6y
  vars  n mean  sd median trimmed mad min max range skew kurtosis se
X1     1 24  5.09 0.78   5.2   5.12 0.94 3.37 6.25  2.88 -0.27  -1.04 0.16
-----
group: basic.9y
  vars  n mean  sd median trimmed mad min max range skew kurtosis se
X1     1 57  5.04 0.63   5.12   5.09 0.68 2.83 6.19  3.36 -0.86   1.14 0.08
-----
group: high.school
  vars  n mean  sd median trimmed mad min max range skew kurtosis se
X1     1 91  5.35 0.88   5.33   5.34 0.87 3.18 7.46  4.28 0.12  -0.33 0.09
-----
group: illiterate
  vars  n mean sd median trimmed mad min max range skew kurtosis se
X1     1 1 5.23 NA   5.23   5.23  0 5.23 5.23   0  NA    NA NA
-----
group: professional.course
  vars  n mean  sd median trimmed mad min max range skew kurtosis se
X1     1 51  5.33 0.81   5.18   5.3 0.82 3.85 6.94  3.09 0.33  -0.78 0.11
-----
group: university.degree
  vars  n mean  sd median trimmed mad min max range skew kurtosis se
X1     1 194  5.2 0.75   5.2   5.19 0.71 2.94 7.43  4.48 0.1   0.46 0.05
```

Figure 18



Figure 19

Section 3 - Result

An alpha level of .05 was adopted for the Pearson Correlation Test and Cohen's rules on effect size (coefficient which is r) were adopted. For correlation, according to Cohen, the effect size is low if the r value varies around .1, medium if r varies around .3, and high if r varies more than .5 (Field, Miles & Field, 2012). A p -value lower than the level of significance of .05 indicates that the null hypothesis is clear evidence to reject it. The effect size is low if the value of eta-squared varies around .01, medium if eta-squared varies around .06, and high if eta-squared varies more than .14, as per Cohen, for difference (Field, Miles & Field, 2012).

Hypothesis 1:

- H0: There will be no significant predictor for last call duration that leads clients to subscribed term deposits in different age groups?
- H1: There will be a significant predictor for last call duration that leads clients to subscribed term deposits in different age groups?

Statistical Evidence

Check the difference in the last call duration for clients by different age groups.

```
# Descriptive statistics by clients age group
durage<-na.omit(data.frame(sbank2$tduration, sbank2$age_label))
names(durage)<-c('tduration', 'age_label')

# Descriptive statistics by client age groups
describeBy(durage$tduration, durage$age_label)
mean_durage<-round(tapply(durage$tduration, durage$age_label, FUN =mean), digits = 2)
mean_durage
```

Descriptive statistics by group

group: Mid-age													
	vars	n	mean	sd	median	trimmed	mad	min	max	range	skew	kurtosis	se
X1	1	247	5.22	0.8	5.18	5.19	0.75	2.83	7.46	4.62	0.32	0.3	0.05

group: mid-old													
	vars	n	mean	sd	median	trimmed	mad	min	max	range	skew	kurtosis	se
X1	1	112	5.34	0.73	5.34	5.34	0.59	3.53	7.25	3.72	-0.03	0.21	0.07

group: senior-citizen													
	vars	n	mean	sd	median	trimmed	mad	min	max	range	skew	kurtosis	se
X1	1	11	5.58	0.74	5.26	5.58	0.64	4.43	6.74	2.31	0.16	-1.52	0.22

group: Young													
	vars	n	mean	sd	median	trimmed	mad	min	max	range	skew	kurtosis	se
X1	1	92	5.12	0.77	5.17	5.12	0.78	2.94	6.79	3.84	-0.14	-0.01	0.08

Figure 20

Mid-age	mid-old	senior-citizen	Young
5.22	5.34	5.58	5.12

Figure 21

The plot (Figure 22) shows how mean values of last call duration changes with different client age groups and the number of people belonging to each group as well.

```
plotmeans(durage$tduration~durage$age_label,digits = 2,
          ccol = 'red',mean.labels = TRUE,xlab = 'Clients age group',ylab = 'Last call
duration',
          main='Plot of Last call duration Mean by Clients age group')
```

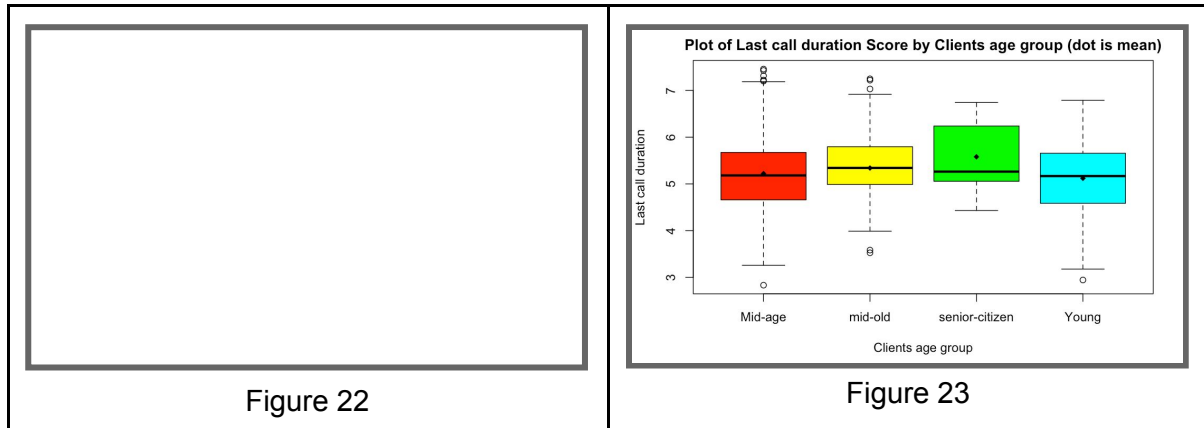
From the graph(Figure 22), we can understand that the mean value of last call duration differs for different groups.

Last call duration with 'young' group having the lowest mean and Group 'senior-citizen' having the highest mean.

the boxplot analysis for further hypothesis testing was performed.

```
boxplot(durage$tduration~durage$age_label,
        main='Plot of Last call duration Score by Clients age group (dot is mean)',
        xlab = 'Clients age group',ylab = 'Last call duration',col=rainbow(6))
points(mean_durage,col='black',pch=18)
```

As it was obvious in boxplots(Figure 23), it was inferred that each client in different age groups has a different amount of variation in last call duration and there is a lot of overlap among values for different groups. But this information is not enough to provide evidence to simply affirm or reject null hypothesis as it does not give information whether the differences are statistically significant. To determine statistical significance, we need to assess the confidence intervals for the differences of means. We further investigate the difference in mean values, considering there is a lot of overlap of last call duration for different clients groups, just because of variation within groups or variation among the groups. The ANOVA Test would be done.



Doing the Bartlett test of homogeneity of variances for last call duration and clients in different age groups.

```
bartlett.test(durage$tduration,durage$age_label)
```

As it was aboviused in Figure 23 the p-value = .71 > .05 so the null hypothesis of the test is accepted and should be said that the variance of different groups can be assumed to be equal.

Figure 24

Doing the assumption of homoscedasticity for Anova test for last call duration and clients in different age groups.

```
aov_durage<-aov(durage$tduration~durage$age_label)
summary(aov_durage)
```

Figure 25

Since the value of F-statistic=2.05 > 1 (significant) and p-value=0.105 > 0.05 in Figure(24), this shows that the variation among the groups and the variation within groups is high, so the mean values for different groups are not significantly different. Therefore, we could not reject the null hypothesis of the test which means the values for different groups are equal. In total, it was concluded that for the confidence interval the null hypothesis can not be rejected that there is no significant difference in last call duration for clients in different age groups.

Doing the calculation of eta square, the result came in Figure 25.

```
etaSquared(aov_durage)
```

Figure 25

Report of Difference Analysis:

A Bartlett's test was done and the equality of variance for Last call duration for all clients in different age group was indicated *K-squared* = 1.35, *P* = 0.7.

A one-way between-groups analysis of variance was conducted to last call duration for clients in different age groups. clients were divided into four groups according to their age (Group 1 : senior-citizen, Group 2 : Mid-age, Group 3 : mid-old, Group 4 : young).

There was no statistically significant difference level in last call duration mean for different clients age (*F(3, 458)* = 2.05, *p* = .105).

The effect size, calculated using eta squared was .01.

The test results indicate there is evidence to support accepting the null hypothesis that there is no difference in last call duration for clients in different age groups.

Checking the difference in last duration call leads clients to subscribed a term deposit.

```
# Descriptive statistics
contact<-as.factor(sbank2$y)

#Conduct Levene's test for homogeneity of variance in library car
ltest<-car::leveneTest(tduration ~ y, data=sbank2)
#Pr(F) is the probability
ltest

#Conduct the t-test from package stats
#You can use the var.equal = TRUE option to specify equal variances and a pooled variance
estimate
stats::t.test(tduration~y,var.equal=TRUE,data=sbank2)
#Effect Size

effsize::cohen.d(tduration,y, alpha = 0.05, na.rm=TRUE)
# effet size=-19.23 large difference
```

Report of Difference Analysis

A Levene's test was conducted and indicated equality of variance for Last call duration for clients who subscribed the term deposit (*F-value* = 3.02, *P* = .08). A t-test analysis of variance was conducted to explore last call duration for clients who subscribed term deposit. Participants were divided into groups according to which clients will subscribe to the term deposit(Group 1 : Yes, Group 2 : No). There was a statistically significant difference in last call duration mean scores for clients who subscribed to a term deposit. The p-value for two

sample tests is very small which means it is significant and can reject the null hypothesis. The effect size, calculated using Cohen's d was -19.23 which implies there is a strong standardised mean difference for both groups. The test results indicate there is evidence to reject null alternative hypothesis which is no difference in Last call duration for clients who subscribed term deposit. The result came in Figure 27.

```
Levene's Test for Homogeneity of Variance (center = median)
      Df F value Pr(>F)
group  1   3.027 0.08256 .
      460
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

      Two Sample t-test

data:  tduration by y
t = -7.9077, df = 460, p-value = 1.963e-14
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -1.0548170 -0.6349074
sample estimates:
mean in group no mean in group yes
      5.140269      5.985131

Cohen's d

d estimate: -19.23453 (large)
95 percent confidence interval:
      lower      upper
-21.09966 -17.36939
```

Figure 27

Model 1

Building the linear regression models

Baseline Model last call duration predicted that clients in different age groups will subscribe to term deposit.

```
#dummy code
df<-data.frame(sbank2)
which(names(df)=='age_label')
which(names(df)=='y')
which(names(df)=='education')
df<-df[,c(1,21,5)]
df$tduration<-tduration
df<-na.omit(df)
df
df$age_label=recode(df$age_label,"senior-citizen"="1","Mid-age"="2","mid-old"="3","Young"="4")
df$education=recode(df$education,'basic.4y'='1','basic.6y'='2','basic.9y'='3','high.school'='4','professional.course'='5','university.degree'='6')
df$y=recode(df$y,'yes'='1','no'='2')
df
```

Creating model1

```
model1=lm(df$tduration~df$y+df$age_label)
anova(model1)
summary(model1)
```

Figure 28

Figure 29

```
stargazer(model1, type="text") #Tidy output of all the required stats  
# Figure 29  
plot(model1)
```

Figure 30

Probability and Statistical Inference
Continuous Assessment Part II

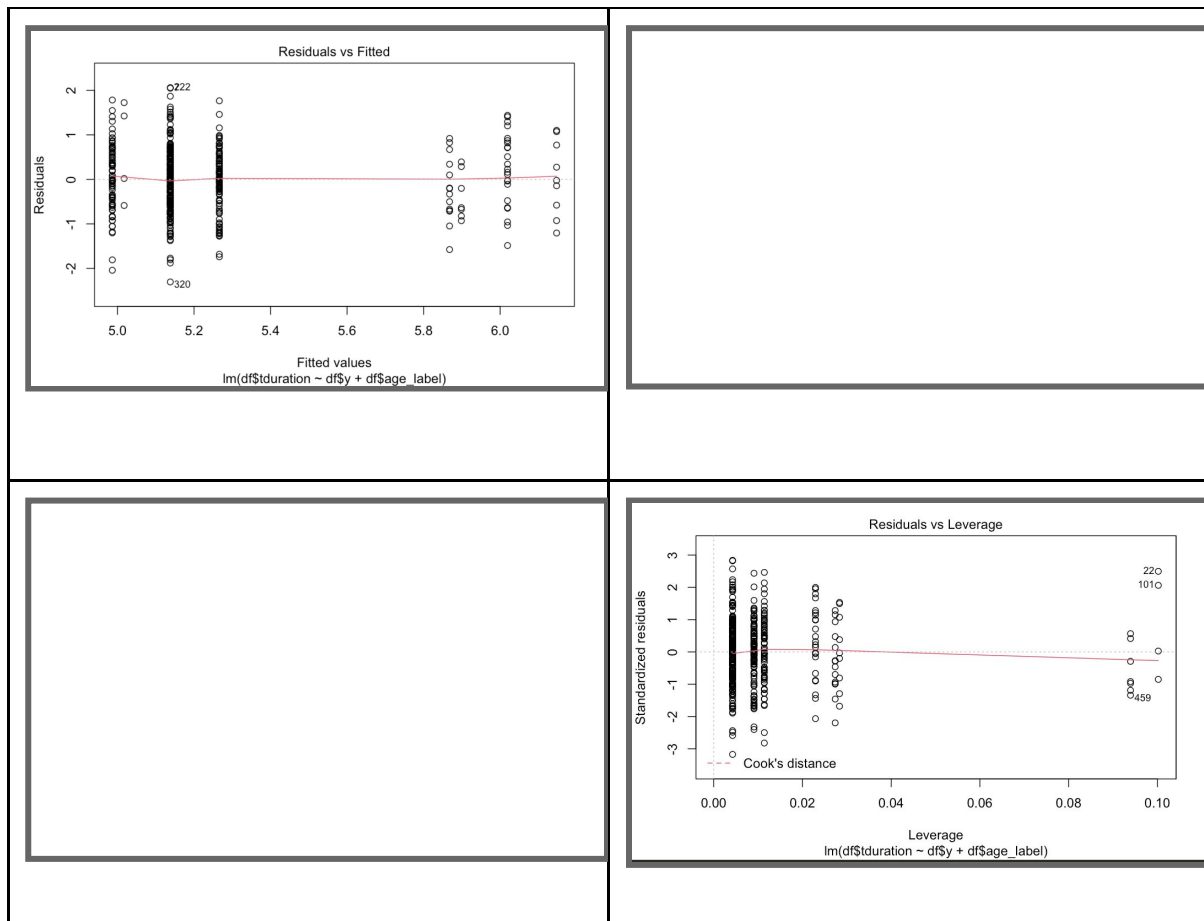


Figure31

```
#Check assumptions
# List of residuals
resid(model1)
#A density plot of the residuals
plot(density(resid(model1)))
```



Figure 32

```
# leverage plots
leveragePlots(model1)
```



Figure 33

```
#Cook's distance
cooks.distance(model1)
#Plot Cook's distance
plot(cooks.distance(model1), ylab="Cook's statistic")
# none of the values is greater than 1 so no influential values
```

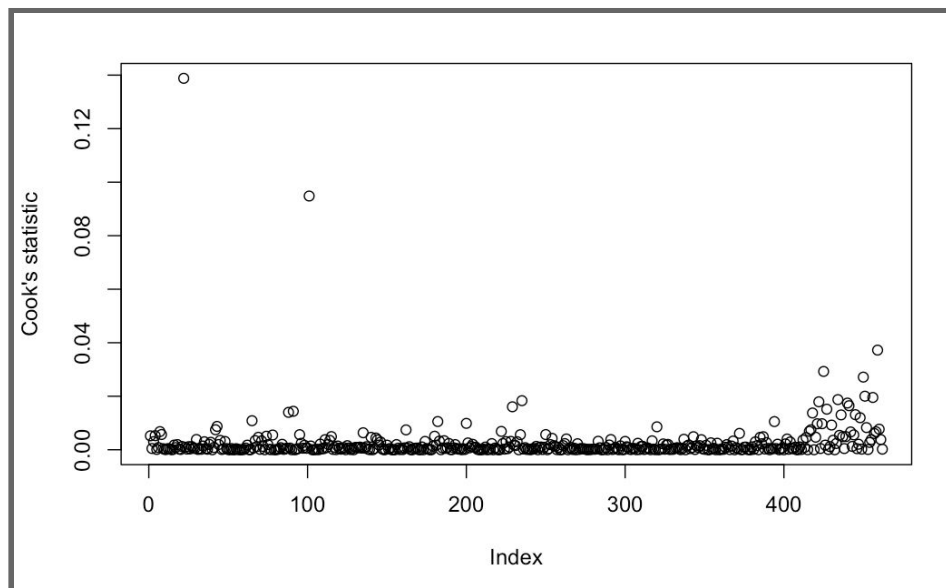


Figure 34

```
# Collinearity
vifmodel<-vif(model1)
vifmodel
# value < 2.5 not problem
```

```
1/(vifmodel)
# values > .4 not problem
```

Figure 35

Report of Linear Modelling Analysis

A multiple linear regression analysis was conducted to determine whether last call duration and clients age lead the customer to subscribed the term deposit. A significant regression equation was found ($F(4,457) = 17.7$, $p = 1.612e-13$), with an Multiple R-squared=.1341.

Examination of the histogram, normal P-P plot of standardised residuals and the scatterplot of the dependent variable, last call duration, and standardised residuals showed that no outliers existed and the residuals followed normal distribution. Also, examination of the standardised residuals showed that none of the values was outside the standard range (95% within limits of -3.29 to +3.29) as the minimum and maximum values are -2.3 and 2.06 respectively further affirming that there were no outliers. Also, none of the Cook's distances were found to be more than 1, hence there are no influential values.

Examination for multicollinearity showed that the tolerance and variance influence factor measures were within acceptable levels (tolerance >0.4, VIF <2.5) as outlined in Tarling (2008). The scatter plot of standardised residuals showed that the data met the assumptions of homogeneity of variance and linearity. The data also meets the assumption of non-zero variances of the predictors.

Because all the assumptions for the model 1 have been proven true and 13.41% of the variance in last call duration is explained by the considered predictors. On checking the significance levels for each of the main terms (in this case the coefficients associated with the constant, y_2), we found that there is evidence that each of these terms are adding something to the model (they are statistically significant as $p < .05$). Hence, these statistical values provide enough evidence to reject that null hypothesis which is no significant prediction of last call duration that lead clients to subscribed term deposits in different age groups.

Hypothesis 2:

- H0: There will be no significant predictor for last call duration that leads clients to subscribed term deposits by different age and different education?
- H1: There will be a significant predictor for last call duration that leads clients to subscribed term deposits by different age and different education?

Statistical Evidence

Check difference in the last call duration for clients with different education

```
# Descriptive statistics
duredu<-na.omit(data.frame(tduration,sbank2$education))
names(duredu)<-c('tduration','education')
# Descriptive statistics by education
describeBy(duredu$tduration,duredu$education)
# check mean for each education group
mean_duredu<-round(tapply(duredu$tduration,duredu$education,FUN =mean),digits = 2)
mean_duredu
```

Figure 36

basic.4y	basic.6y	basic.9y	high.school	professional.course	university.degree
5.40	5.09	5.04	5.35	5.33	5.20

Figure 37

The plot (Figure 38) shows how mean values of last call duration changes with different client age groups and the number of people belonging to each group as well.

```
plotmeans(duredu$tduration~duredu$education,digits = 2,
          ccol = 'red',mean.labels = TRUE,xlab = 'Clients education',ylab = 'Last call
duration',main='Plot of Last call duration Mean by Clients educationh')
```

From the graph (Figure 38), it was founded that the mean value of last call duration differs for different education groups.

Last call duration with 'basic.9y' group having the lowest mean and group 'basic.4y' having the highest mean.

the boxplot analysis for further hypothesis testing was performed.

```
boxplot(duredu$tduration~duredu$education,
        main='Plot of last call duration by Clients educationh (dot is mean)',xlab = '
Clients educationh',ylab = 'Last call duration',col=rainbow(11))
```

```
points(mean_durededu,col='black',pch=18)
```

As it was obvious in boxplots(Figure 39), it was inferred that each client in different age groups has a different amount of variation in last call duration and there is a lot of overlap among values for different groups. But this information is not enough to provide evidence to simply affirm or reject null hypothesis as it does not give information whether the differences are statistically significant. To determine statistical significance, we need to assess the confidence intervals for the differences of means. We further investigate the difference in mean values, considering there is a lot of overlap of last call duration for different clients groups, just because of variation within groups or variation among the groups. The ANOVA Test would be done.

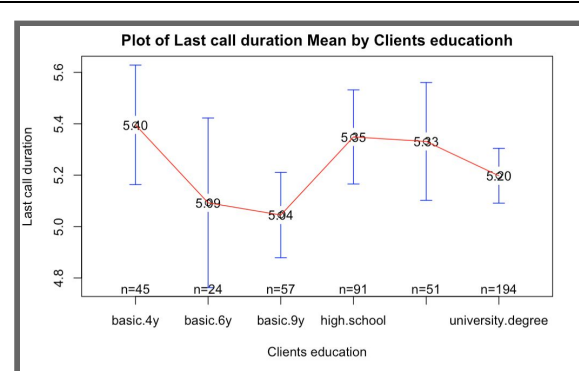


Figure 38

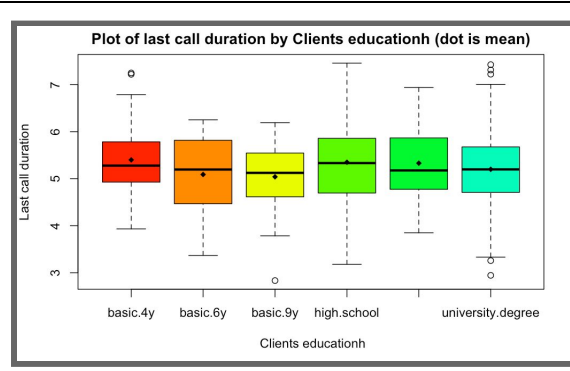


Figure 39

Doing the Bartlett test of homogeneity of variances for last call duration and clients with different education.

```
bartlett.test(durededu$tduration,durededu_new$education)
```

As it was aboviused in (Figure 40) the p-value = .14 > .05 so the null hypothesis of the test is accepted and should be said that the variance of different groups can be assumed to be equal.

Figure 40

Doing the assumption of homoscedasticity for Anova test for last call duration and clients in different age groups.

```
aov_durage<-aov(durage$tduration~durage$age_label)
summary(aov_durage)
```

Figure 41

Since the value of F-statistic=1.87 > 1 (significant) and p-value=0.097 > 0.05 in (Figure41), this shows that the variation among the groups and the variation within groups is high, so the mean values for different groups are not significantly different. Therefore, we could not reject the null hypothesis of the test which means the values for different groups are equal. In total, it was concluded that for the confidence interval the null hypothesis can not be rejected that there is no significant difference in last call duration for clients with different education.

Doing the calculation of eta square, the result came in Figure 41.

```
etaSquared(aov_durage)
```

```
eta.sq eta.sq.part  
durededu$education 0.02010697 0.02010697
```

Figure 42

Report of Difference Analysis:

A Bartlett's test was done and the equality of variance for Last call duration for all clients in different age group was indicated *K-squared* = 8.12, *P* = .14.

A one-way between-groups analysis of variance was conducted to last call duration for clients with different education. clients were divided into six groups according to their education (Group 1 : basic.4y, Group 2 : basic.6y, Group 3 : basic.9y, Group 4 : high school, Group 5 : professional.course, Group 6 :university.degree).

There was no statistically significant difference level in last call duration mean for different clients education (*F(5, 456)* = 1.87, *p* = .09).

The effect size, calculated using eta squared was .02.

The test results indicate there is evidence to support accepting the null hypothesis that there is no difference in last call duration for clients with different education.

Model 2

Building the linear regression models

Model last call duration predicted that clients by different age groups and different education will subscribe to term deposit.

Creating model 2

```
model2<-lm(df1$tduration~df1$y+df1$age_label+df1$education)  
anova(model2)  
summary(model2)
```


Figure 43

Figure 44

```
stargazer(model2, type="text") #Tidy output of all the required stats  
stargazer(model1, model2, type="text") #Quick model comparison  
plot(model2)
```

Probability and Statistical Inference
Continuous Assessment Part II

```

=====
                        Dependent variable:
                        -----
                        tduration
=====
y2                      -0.853***
                        (0.112)

age_label2              0.109
                        (0.234)

age_label3              0.231
                        (0.239)

age_label4             -0.045
                        (0.240)

education2             -0.136
                        (0.186)

education3             -0.170
                        (0.148)

education4              0.038
                        (0.136)

education5              0.010
                        (0.151)

education6             -0.103
                        (0.123)

Constant               5.950***
                        (0.239)

=====
Observations            462
R2                      0.143
Adjusted R2             0.126
Residual Std. Error    0.729 (df = 452)
F Statistic            8.368*** (df = 9; 452)
=====
Note:                  *p<0.1; **p<0.05; ***p<0.01

```

Figure 45

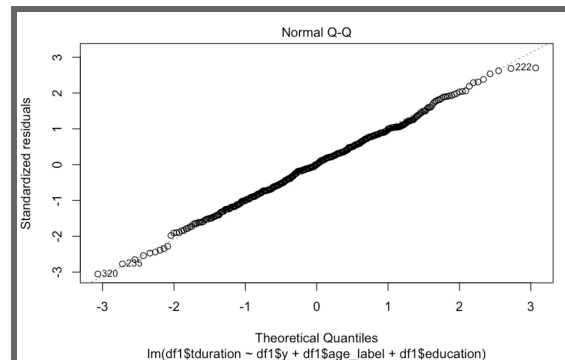
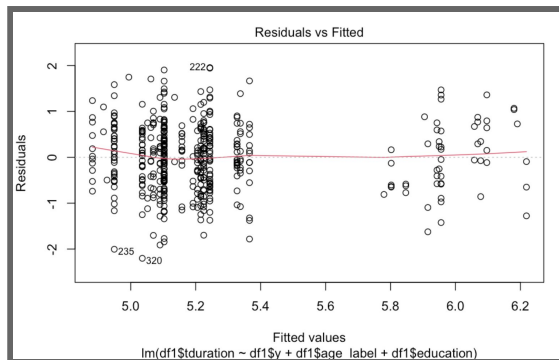


Figure 46

```
#Check assumptions  
# List of residuals  
resid(model1)  
#A density plot of the residuals  
plot(density(resid(model1)))
```

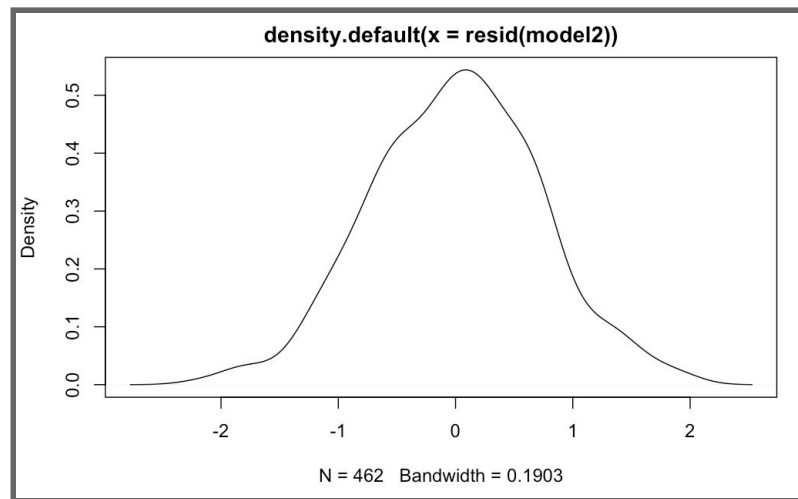


Figure 47

```
# leverage plots  
leveragePlots(model1)
```



Figure 47

```
#Cook's distance  
cooks.distance(model1)  
#Plot Cook's distance  
plot(cooks.distance(model1), ylab="Cook's statistic")  
# none of the values is greater than 1 so no influential values
```

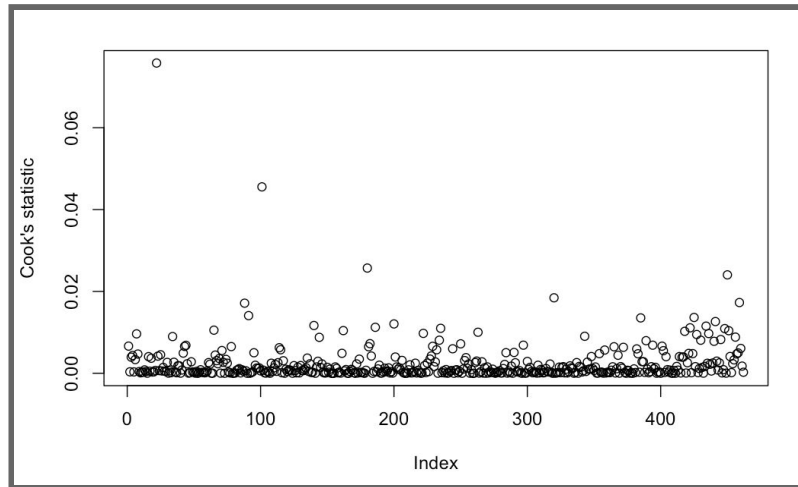


Figure 48

```
# Collinearity
vifmodel<-vif(model2)
vifmodel
# value < 2.5 not problem
1/(vifmodel)
# values > .4 not problem
```

	GVIF	Df	GVIF^(1/(2*Df))
df1\$y	1.104606	1	1.051002
df1\$age_label	1.146879	3	1.023104
df1\$education	1.091177	5	1.008764
	GVIF	Df	GVIF^(1/(2*Df))
df1\$y	0.9053000	1.0000000	0.9514726
df1\$age_label	0.8719319	0.3333333	0.9774182
df1\$education	0.9164419	0.2000000	0.9913123

Figure 49

Report of Linear Modelling Analysis

A multiple linear regression analysis was conducted to determine whether last call duration and clients age lead the customer to subscribed the term deposit. A significant regression equation was found ($F(9,452) = 8.368$, $p = 1.4e-11$), with an Multiple R-squared=.1428.

Examination of the histogram, normal P-P plot of standardised residuals and the scatterplot of the dependent variable, last call duration, and standardised residuals showed that no outliers existed and the residuals followed normal distribution. Also, examination of the standardised residuals showed that none of the values was outside the standard range (95% within limits of -3.29 to +3.29) as the minimum and maximum values are -2.2 and 1.95 respectively further affirming that there were no outliers. Also, none of the Cook's distances were found to be more than 1, hence there are no influential values.

Examination for multicollinearity showed that the tolerance and variance influence factor measures were within acceptable levels (tolerance >0.4, VIF <2.5) as outlined in Tarling (2008). The scatter plot of standardised residuals showed that the data met the assumptions of homogeneity of variance and linearity. The data also meets the assumption of non-zero variances of the predictors.

Because all the assumptions for the model 1 have been proven true and 14.28% of the variance in last call duration is explained by the considered predictors. On checking the significance levels for each of the main terms (in this case the coefficients associated with the constant, y_2), we found that there is evidence that just y_2 and constant are adding something to the model (this statistically significant as $p < .05$). Hence, these statistical values provide enough evidence to reject that null hypothesis which is no significant prediction of last call duration that lead clients to subscribed term deposits by different age and different education.

Summary Comparison

Compare Model 1 and Model 2

```
anova(model1,model2)
```

Figure 50

Model Comparison Results

It was obvious that the p-value obtained for second model ($*p = .47$) is not statistically significant (more than .05), so addition of a new variable significantly not improved the fit over model 1. therefore, we should not reject model 1.

Section 4 – Discussion/Conclusion

In this study, we wanted to determine which factors could be considered as the best predictors for determining the last call duration that can lead customers to subscribed term deposits . This was conducted by using Multiple Linear Regression where we firstly tried to establish evidence that the various predictors chosen can be used for modelling. We investigated whether a client's age, their education, and customer subscribed term deposit. We found that there were statistically significant differences in last call duration with customers with different client age, education. Furthermore, different models were built to determine which variables are the best for prediction.

The results of the baseline first model are analysed that determine whether last duration call can lead customers with different age to subscribed to the term deposit can be used as predictor for output variable. Since the p-value for the model obtained is statistically significant ($1.612e-13 < .05$), it was shown the model is good to fit as it performs better than the average score method for prediction. As per the analysis of R-squared value which is found to be .1341, we can say that 13.41% of the variance in the Last Call duration is explained by the considered predictors. On checking the significance levels for each of the

main terms (in this case the coefficients associated with the constant, y_2), we found that there is evidence that each of these terms are adding something to the model (they are statistically significant as $p < .05$). There were found no outliers, residuals, leverage points and influential values for the model. This model explained the % of variance in the Last call duration.

The results of the second model are analysed which determine whether last call duration can lead clients with different age and education to subscribed term deposit can be used as predictor for output variable. Since the p-value for the model obtained is statistically significant ($1.4e-11 < .05$), it gives us evidence to suggest that the model is good to fit as it performs better than the average score method for prediction. As per the analysis of R-squared value which is found to be .1428, we can say that 14.28% of the variance in Last call duration is explained by the considered predictors. On checking the significance levels for each of the main terms (in this case the coefficients associated with the constant, y_2), we found that there is evidence that each of these terms are adding something to the model (they are statistically significant as $p < .05$). There were no outliers, residuals, leverage points and influential values found for the model. Overall, the second model can not be considered better than the baseline model as it is statistically significant based on comparison results and addresses a few higher amounts of variance in the output variable, Last call duration which implies it is not better at making predictions.

From this analysis, there is evidence to conclude that Model 2 is the better compared to the model because of 14.28%% of the variance in Last call duration but this is not very much different between model 1 to model 2, so if we want to decide we should say the model 1 is better because this differences is not a lot.

Reference

1. George, D. (2011).
2. Ziegel, E., & Lohr, S. (2000). *Sampling: Design and Analysis*. *Technometrics*, 42(2), 223. doi: 10.2307/1271491
3. Field, A., Miles, J., & Field, Z. (2012).
4. Moro, S., Cortez, P., & Rita, P. (2014). A data-driven approach to predict the success of bank telemarketing. *Decision Support Systems*, 62, 22-31. doi: 10.1016/j.dss.2014.03.001
5. UCI Machine Learning Repository: Bank Marketing Data Set. (2021). Retrieved 6 January 2021, from <https://archive.ics.uci.edu/ml/datasets/Bank+Marketing>
6. *CrossTable function | R Documentation*. (2021). Retrieved 4 January 2021, from <https://www.rdocumentation.org/packages/gmodels/versions/2.18.1/topics/CrossTable>