



King Saud University
College of Computer and Information Sciences
Information Technology department

**IT 326: Data Mining
Course Project**

Case Study: Applicants for a Gold Digger position

Project Report: Data Summarization and Preprocessing

Group#:	6	
Section#:	64063	
Group Members	Name	ID
	Seba alahmadi	439201336
	Noura alsultan	441201306
	Nouf alfulaij	441201159
	Reem Almutairi	437201303

كيف حلينا مشكلتنا بالكلاسификаشن والكلسترنق ؟

مشكلتنا هـ

1 Problem

مشكلتنا كانت ان وظيفة ال قولد داير وظيفه حساسه وهو اي احد يتوظف فيها ف الكلاسيفيكيشن توقع لنا الموظفين الجدار اذا بيتوظفون او لا بيسد اون ترايننـج سـيـت ، والكلـستـرنـق خـلـانـا نـلـاقـي السـيـمـيـلـارـيـتيـ بيـنـ كلـ كـلـسـتر او قـرـوبـ عـشـانـ نـعـرـفـ وـنـمـيـزـ بيـنـ القـرـوـبـاتـ واـيـشـ كلـ قـرـوبـ يـتـميـزـ عنـ الثـانـيـ

The gold-digger job is an important and sensitive job that needs people with high skills and specialists in this field. In this report, data mining techniques are used to predict whether a candidate will be hired or not based on set of information of people who applied for gold digger position. It is important because it will show summary of candidates that will help and serve people who are interested in hiring people with gold digger experience.

2 Data Mining Task

The data mining tasks are clustering and classification.

ايش راح يـفـيـدـناـ الـكـلـسـتـرنـقـ
فـيـ الدـاـتـاـ سـيـتـ حـقـتـناـ ؟

In clustering we will analyze the Applicants for a Gold Digger position dataset, find the similarity to identify distinct groups and what are the factors that effects and distinguishes one applicant from another.

In classification, it is a binary classification problem where the class attribute is the **embauche(hiring)** which has two values: 1 which indicates that the applicant has been hired and 0 which indicates that the applicant has not been hired.

ايش يـفـيـدـناـ الـكـلـاسـيـفـيـكـيشـنـ فـيـ الدـاـتـاـ سـيـتـ حـقـتـناـ ؟

The goal of the data mining task is to classify an applicant, based on years of relevant experience, salary expectation, highest qualification diploma, specialties, note(grade (out of 100) for gold digging exam) into either a candidate will be hired or not.

↓ Page 26
↓ Phase 2 بدـاـيـةـ

3 Data

We choose this data from Kaggle website <https://www.kaggle.com/bryanb/applicants-for-a-gold-digger-position>

The dataset as shown in *figure1* includes Applicants for a Gold Digger position of 20000 applicants(object) between 1/1/2010 and 30/12/2014 it consists of the following 12 attributes as shown in *figure2* with 396 missing values but we minimize the data to 999 rows as shown in *figure3* with 27 missing values as shown in *figure5*. The attributes names, description, data type, and possible values of features are given in *Table 1*.

Data set rows before

```
> #print NO of objects  
> nrow(data)  
[1] 20000
```

Figure1 number of rows

Data set rows after

```
> #print No of objects  
> nrow(data)  
[1] 999  
> # print NUM of features
```

Figure3 number of rows minimized

```
> # print NUM of features  
> ncol(data)  
[1] 12  
>
```

Figure2 number of columns

```
> #print names of the features  
> names(data)  
[1] "X"           "date"        "cheveux"      "age"          "exp"  
[6] "salaire"     "sexe"        "diplome"      "specialite"   "note"  
[11] "dispo"       "embauche"  
>
```

Figure4 names of features

Attributes name	description	Data type	possible values
x	Index of applicant	Numeric	Range between 1 - 999
date	date of the application	nominal	Range between 1/1/2010 – 30/12/2014
cheveux	hair color	nominal	(chatain,brun,blond,roux)
age	age of the candidate	numeric	Range between [-3.....74]
exp	Years of relevant experience	numeric	Range between [-2.....23]
salaire	Salary expectation	numeric	Range between [14.1k...54k]
sexe	Female or male	binary	F :Female OR M: Male
diplome	Highest qualification diploma	nominal	(bac, licence, master, doctorat)
specialite	minor of the diploma	nominal	(geologie, forage, detective, archeologie)
note	Grade out of 100 for gold digging exam	numeric	Range between [0...100]
dispo	directly available	binary	directly available: oui not directly available: non
embauche	Has the candidate been hired	binary	0: not hired 1: hired

Table 1 Description

Missing values:

```
> #missing values  
> sum(is.na(data))  
[1] 27  
> |
```

Figure5 missing values

The five-number summary of the dataset:

Example: It can be noticed that the dataset contains candidates with ages in the range between -3 and 74 years old. The average age of the candidates in the dataset is 35 years old. As shown in figure6.

```
> # print five number summary  
> #summary(data)  
> summary(data$X)  
   Min. 1st Qu. Median Mean 3rd Qu. Max.  
 0.0 248.5 504.0 500.9 750.5 998.0  
> |  
  
> summary(data$age)  
  Min. 1st Qu. Median Mean 3rd Qu. Max.  
10.00 29.00 35.00 34.85 41.00 59.00  
> |  
  
> summary(data$exp)  
  Min. 1st Qu. Median Mean 3rd Qu. Max.  
4.000 8.000 10.000 9.593 11.000 15.000  
> |  
  
> summary(data$salaire)  
  Min. 1st Qu. Median Mean 3rd Qu. Max.  
21735 31462 35132 35075 38528 46944  
> |  
  
> summary(data$note)  
  Min. 1st Qu. Median Mean 3rd Qu. Max.  
29.97 63.65 74.56 75.10 87.08 121.14  
> |  
  
> summary(data$embauche)  
  Min. 1st Qu. Median Mean 3rd Qu. Max.  
0.0000 0.0000 0.0000 0.1294 0.0000 1.0000  
> |
```

Figure6 five number summary

Outlier Detection and box plots:

Outliers in the dataset affect the effectiveness of the prediction models and therefore it's important to detect and remove outliers from the dataset. *Figure 7,8,9 and 10.*

age boxplot before removing outliers

```
> boxplot(data$age)
> boxplot(data$age)$out
[1] 61 65  8 62  8  8 5 63  5 62 62 62  8
> outliers1 <- boxplot(data$age, plot=FALSE)$out
> print(outliers1)
[1] 61 65  8 62  8  8 5 63  5 62 62 62  8
> data[which(data$age %in% outliers1),]
```

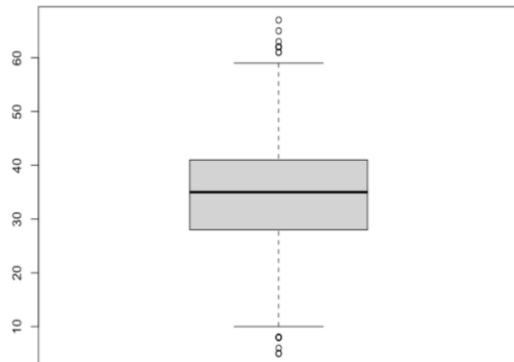


Figure7

The figure7 show us that there are outliers created in 60 and above, and between almost 5-10

note boxplot before removing outliers

```
> boxplot(data$note)
> boxplot(data$note)$out
[1] 127.58 125.47 12.63
> outliers3 <- boxplot(data$note, plot=FALSE)$out
> print(outliers3)
[1] 127.58 125.47 12.63
> data[which(data$note %in% outliers3),]
```

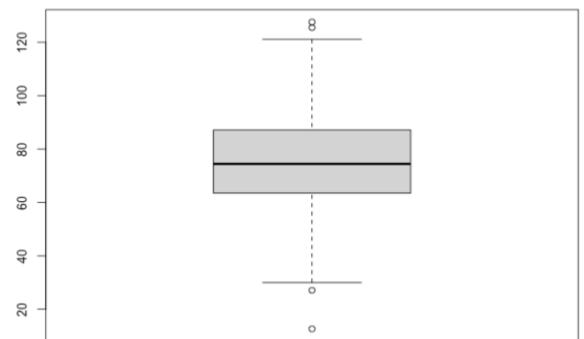
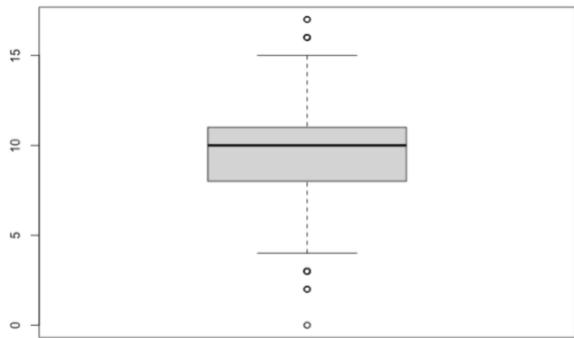


Figure8

The figure8 show us that there are outliers created in above 120 and in almost 15 and 30

Exp boxplot before removing outliers

```
> #data preprccising
> #removing outliers
> library(outliers)
> #data$exp[which(data$exp>420)] <- c(data$exp[which(data$exp>420)]*2)
> boxplot(data$exp)
> boxplot(data$exp)$out
[1] 3 16 3 3 3 3 2 3 17 17 3 3 0 16 3 3 3 3 0 16 16 2 2 3 3 3 2 2 3 3 3 16 17 16 16 16 16 16 3
[41] 16 16 17 16
> outliers <- boxplot(data$exp, plot=FALSE)$out
> print(outliers)
[1] 3 16 3 3 3 3 2 3 17 17 3 3 0 16 3 3 3 3 0 16 16 2 2 3 3 3 2 2 3 3 3 16 17 16 16 16 16 16 3
[41] 16 16 17 16
> data[which(data$exp %in% outliers).]
```

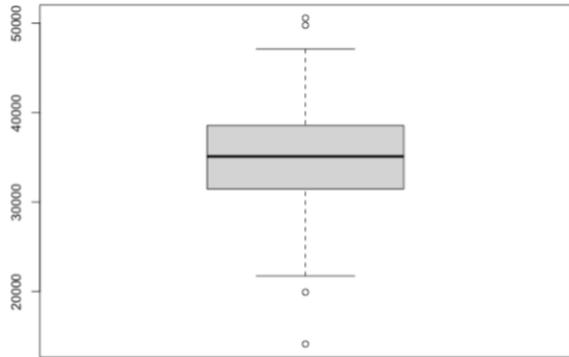


The figure9 show us that there are outliers created in 16 and above and between 0-4

Figure9

salarie boxplot before removing outliers

```
> boxplot(data$salaire)
> boxplot(data$salaire)$out
[1] 19925 50575 14128 49761
> outliers2 <- boxplot(data$salaire, plot=FALSE)$out
> print(outliers2)
[1] 19925 50575 14128 49761
> data[which(data$salaire %in% outliers2),]
```



The figure10 show us that there are outliers created in 5000 and above and between almost 1500 and 2000

Figure10

Description: After box plotting which provided us as box plots it shows that all attributes have outliers that need to be removed later on.

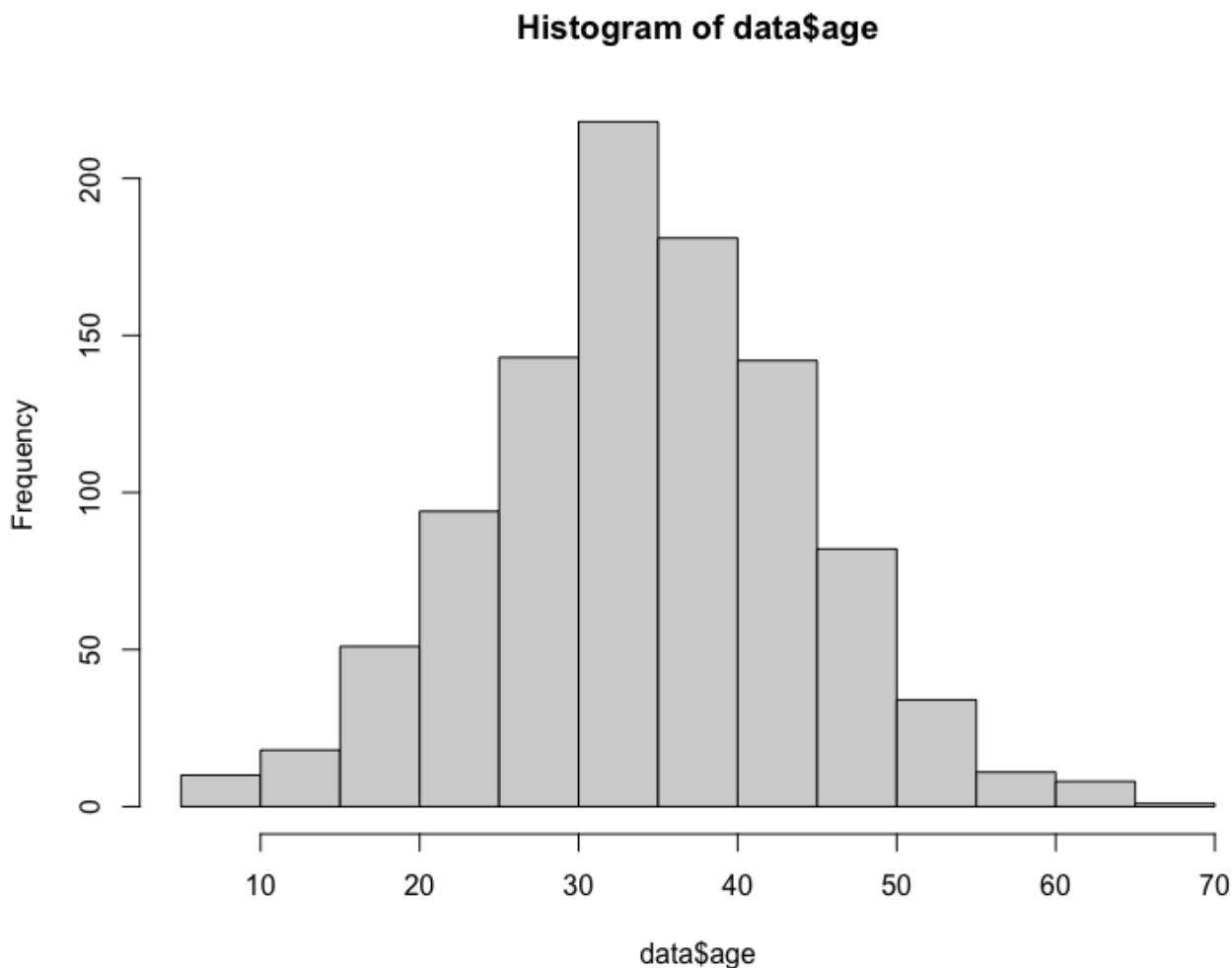


Figure 11

Description: The age distribution *Figure 11* over candidates in the dataset can be illustrated easily using a histogram as it is observed that the age of most of the candidates in the dataset ranges from 30 to 45 years old.

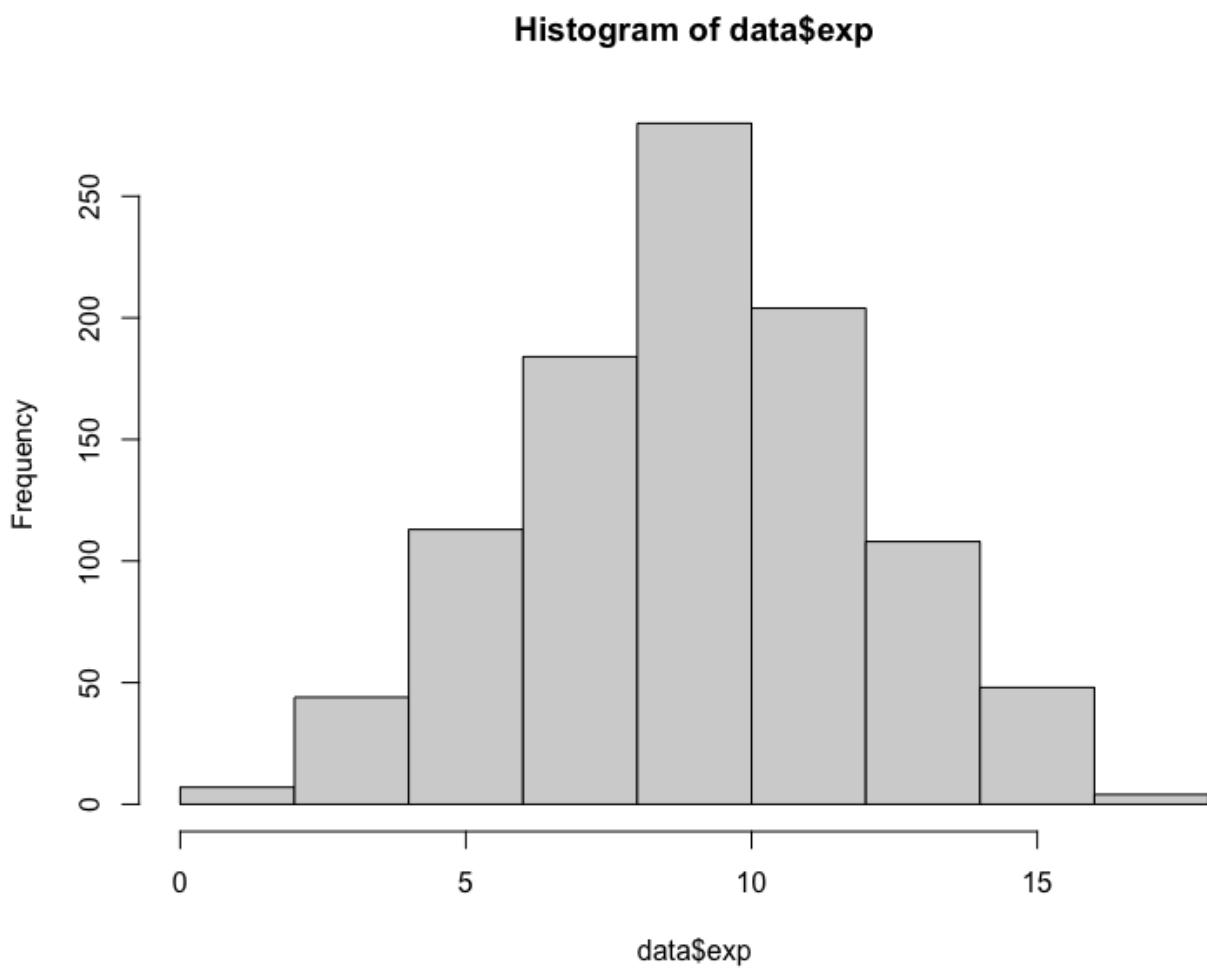


Figure12

Description: The exp distribution *Figure 12* over candidates in the dataset shows that most of the candidates have years of experience between (5-15) years

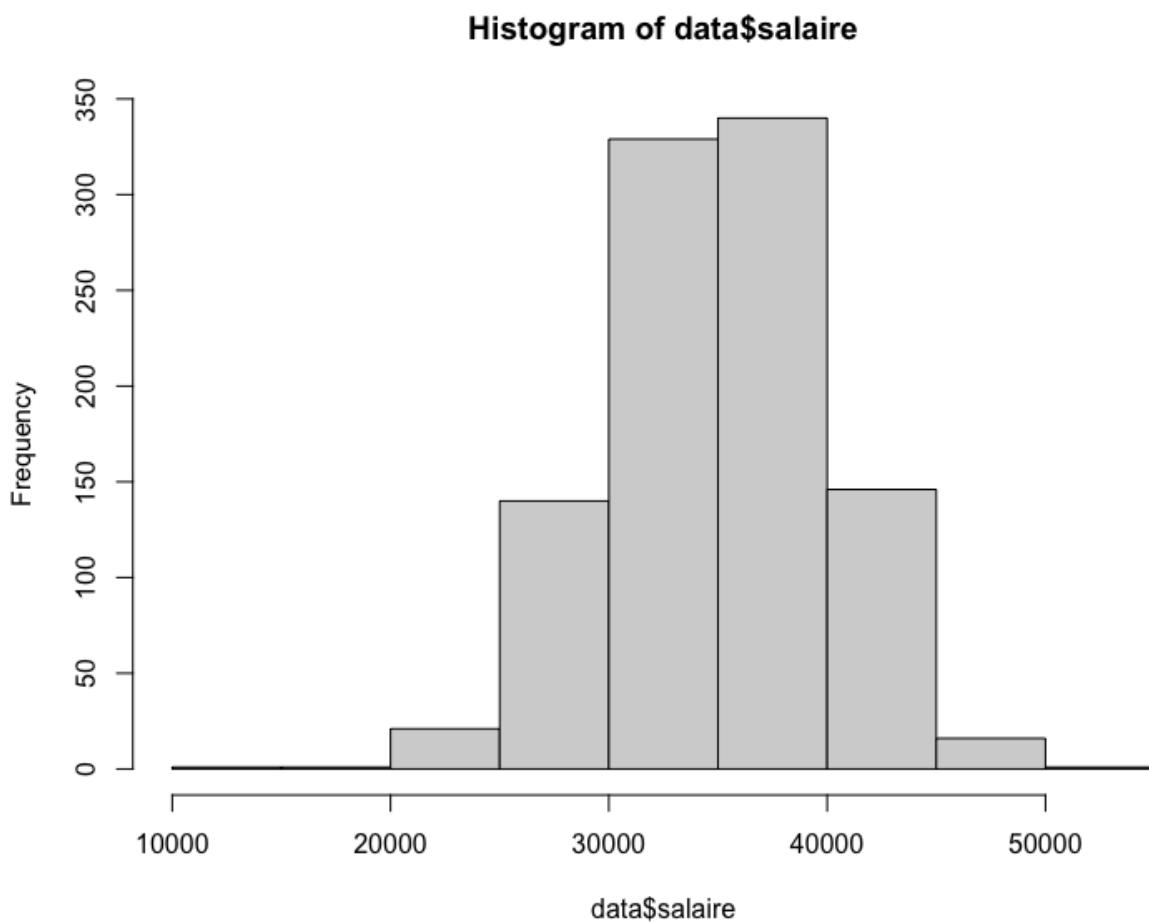


Figure 13

Description: The salaire distribution *Figure 13* over candidates in the dataset shows that most of the candidates receive salaries between (30000 and 40000)

Histogram of data\$note

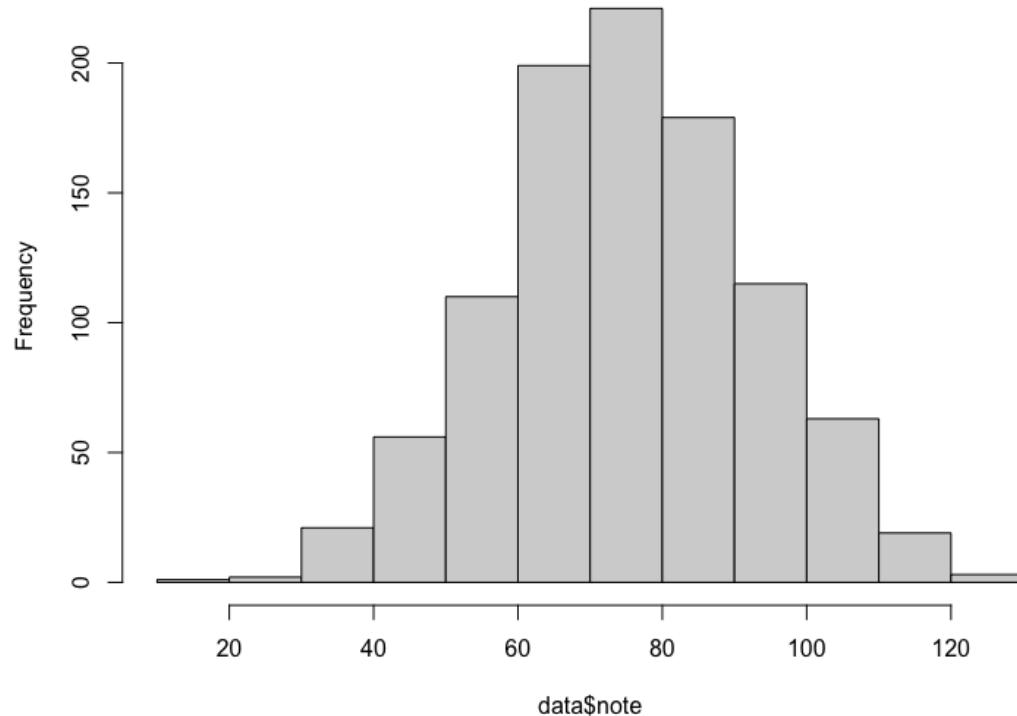


Figure 14

Description: The note distribution *Figure 14* over candidates in the dataset shows that most of the candidates have grade between(60-80)

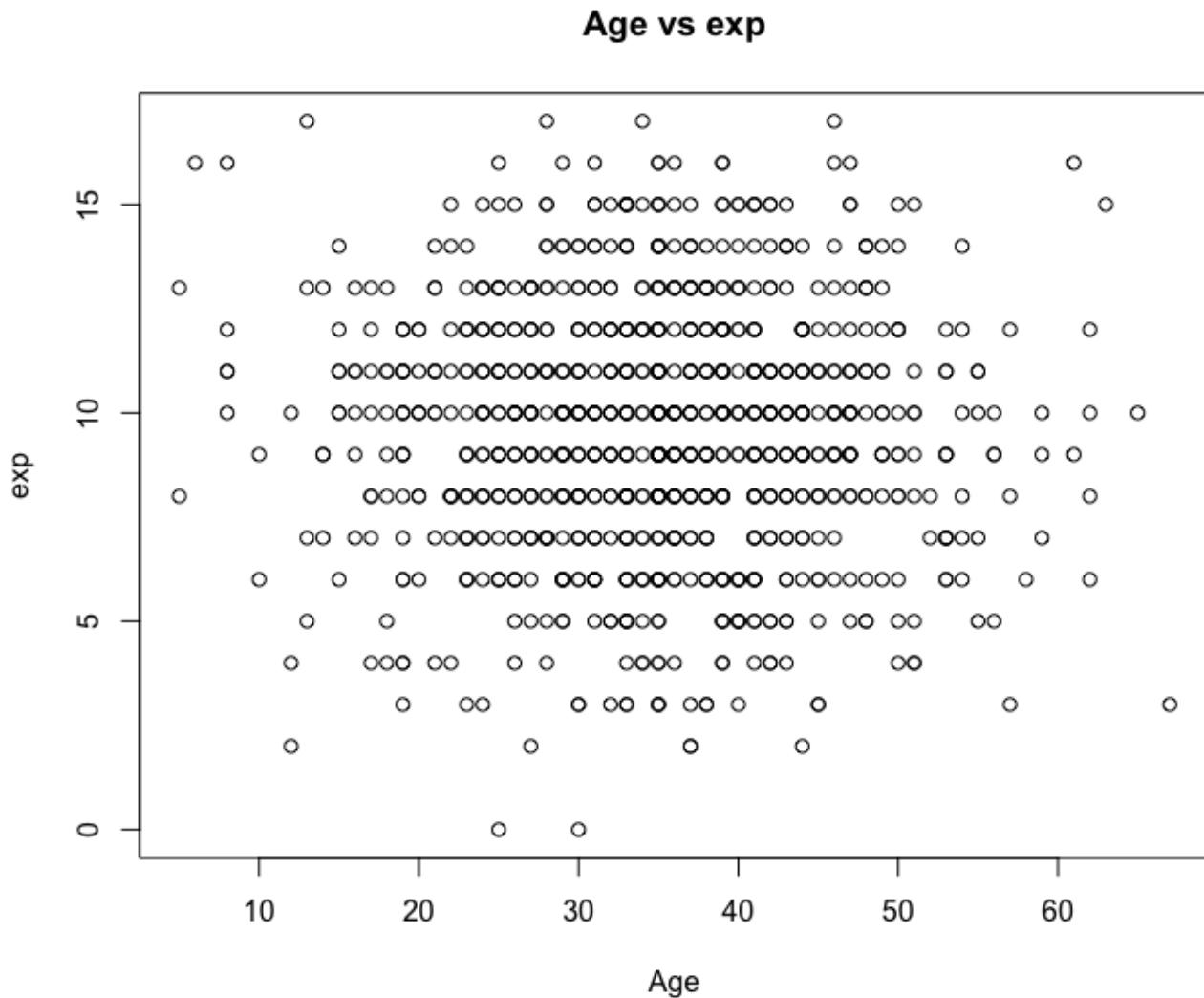


Figure 15

Description: This scatter plot *Figure 15* shows that there is no relationship between age and exp so there is no correlation between age and exp . Because the value range is spread through the scatter plot.

Hairbarplot

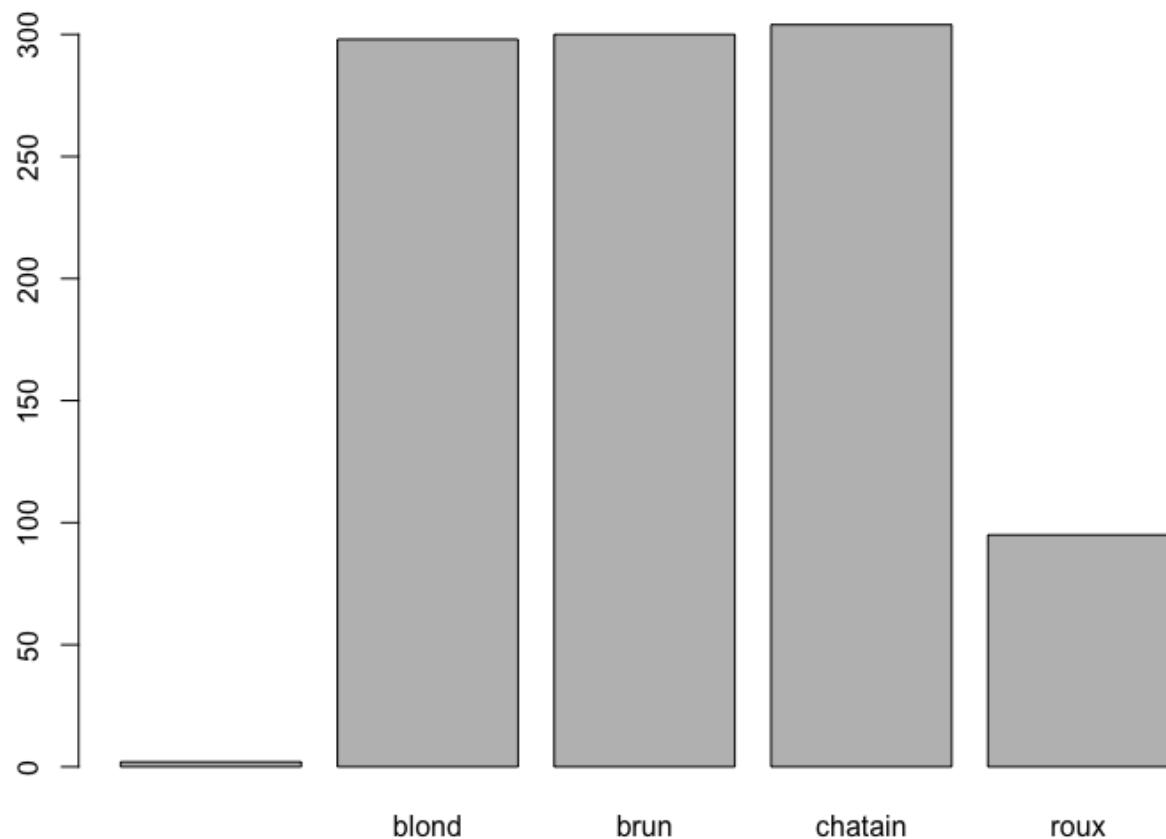


Figure16

Description: From the bar plot *Figure16* it shows that there are almost 300 blond and 300 brun and more than 300 chatain and 100 roux

Date barplot

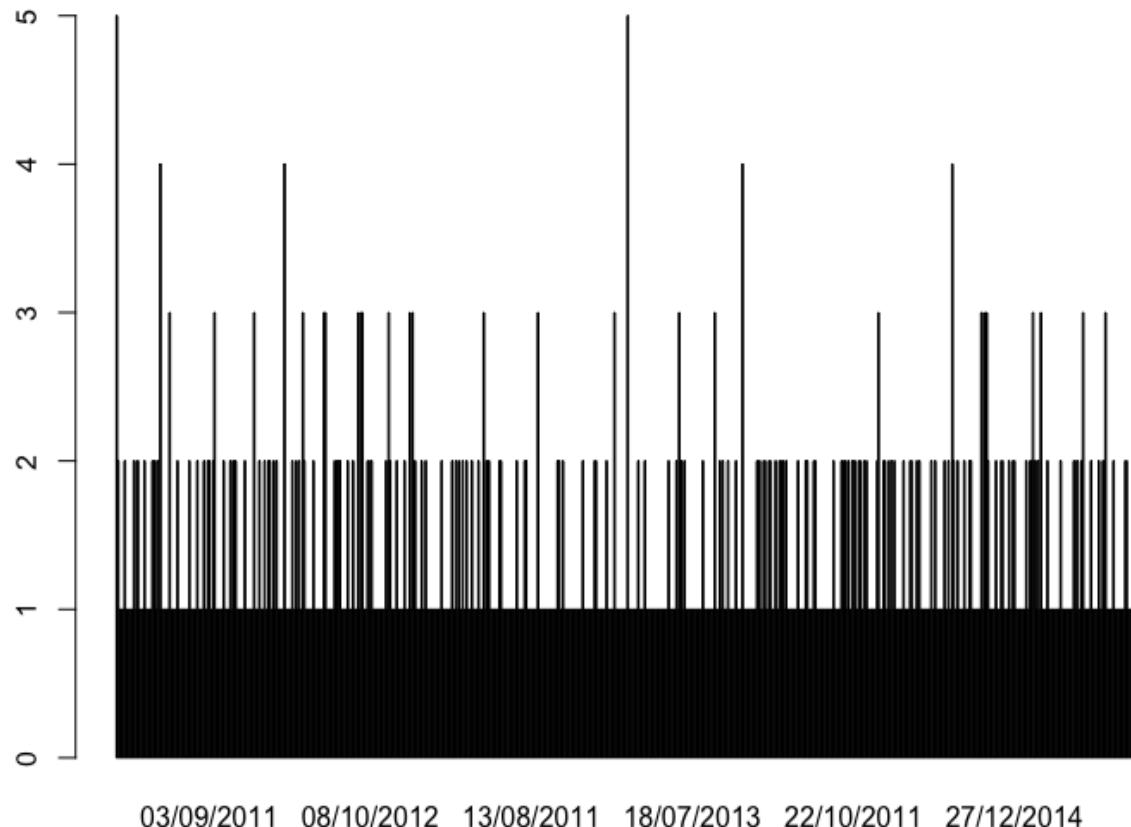


Figure17

Description: From the bar plot *Figure17* it shows that there are random distribution for dates from 2011-2014

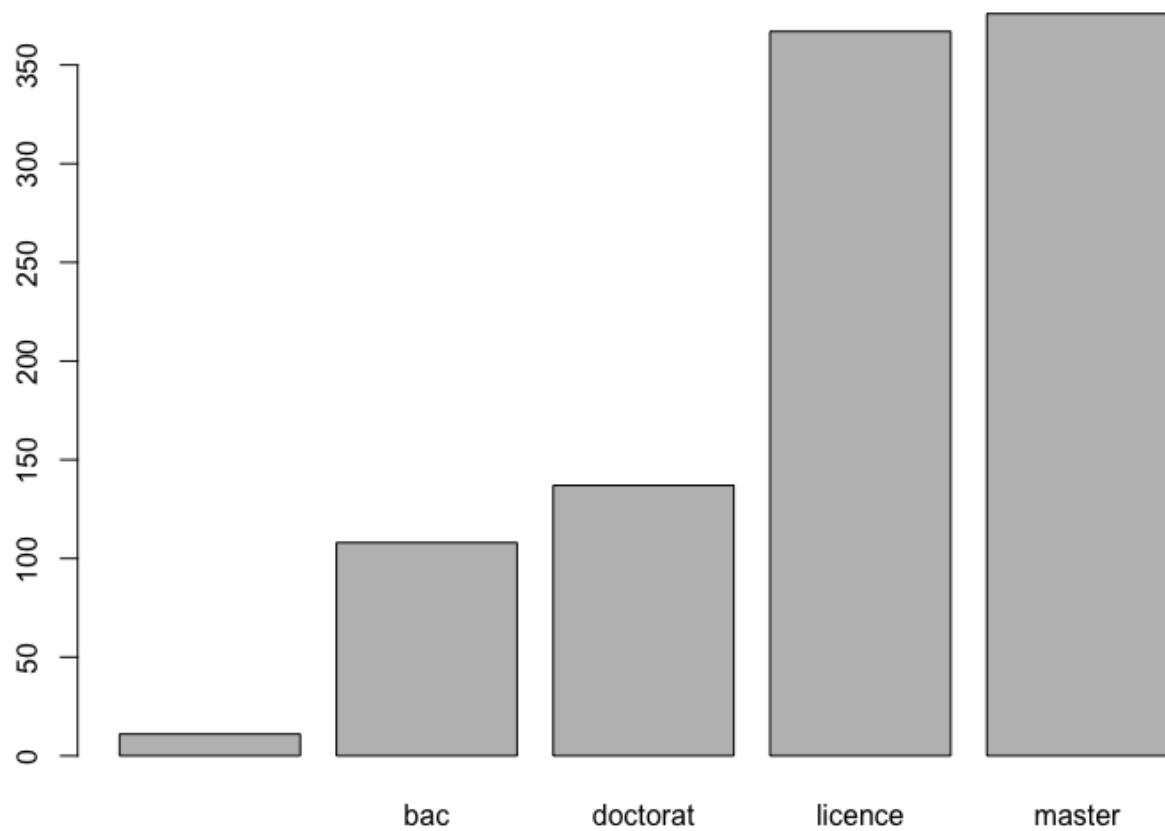


Figure18

Description: From the bar plot *Figure 18* it shows that there are more than 100 bac and doctorat and 50 licence and more than 350 master.

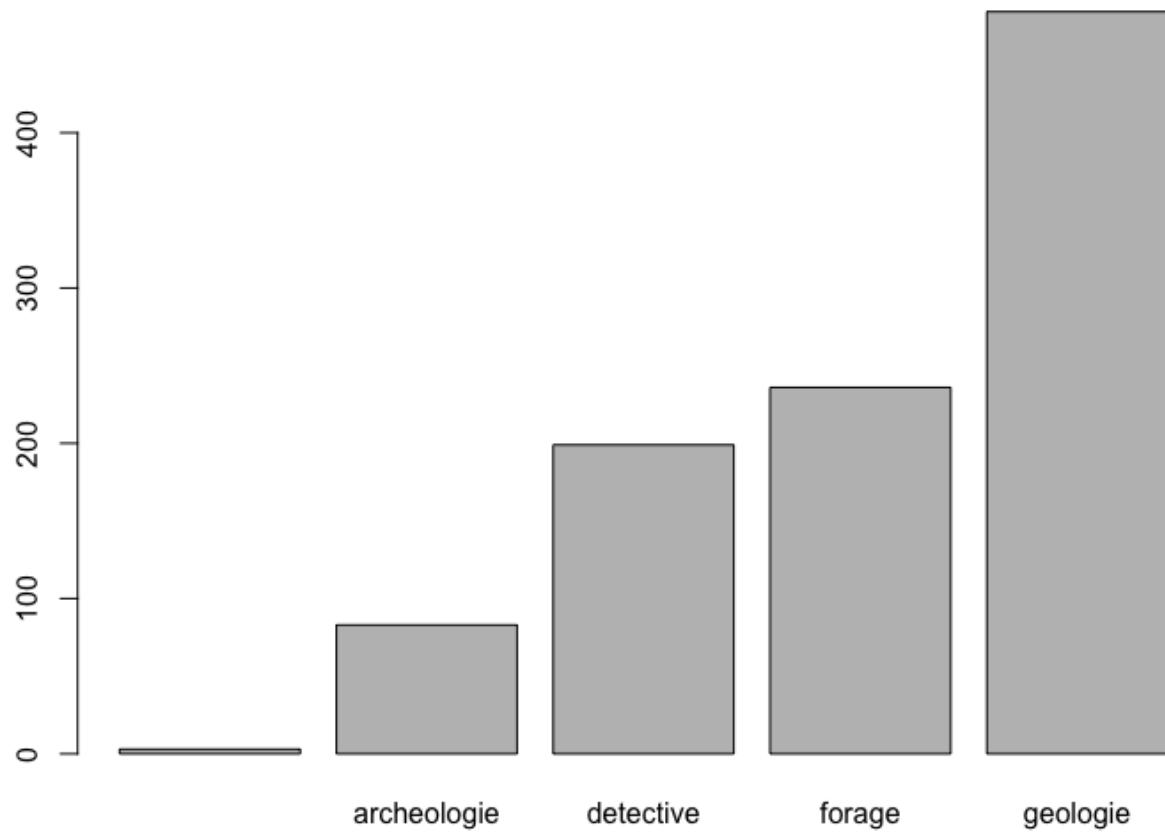


Figure19

Description: From the bar plot *Figure19* it shows that there are less than 100 archeologie and 200 detective and more than 200 forage and more than 400 geologie.

4 Data preprocessing

Our dataset before preprocessing *Figure20* : is contain 999 object and 12 attribute.

The screenshot shows the RStudio Source Editor with a data frame titled 'data'. The table has 25 rows and 12 columns. The columns are labeled: X, date, cheveux, age, exp, salaire, sexe, diplome, specialite, note, dispo, and embauche. The data includes various personal and professional information such as date of birth, hair color, age, experience, salary, gender, education level, specialization, grade, availability, and employment status. Some cells contain NA or specific codes like 'geologie' or 'forage'.

	X	date	cheveux	age	exp	salaire	sexe	diplome	specialite	note	dispo	embauche
1	0	2012-06-02	roux	25	9	26803	F	licence	geologie	97.08	non	0
2	1	2011-04-21	blond	35	13	38166	M	licence	forage	63.86	non	0
3	2	2012-09-07	blond	29	13	35207	M	licence	geologie	78.50	non	0
4	3	2011-07-01	brun	NA	12	32442	M	licence	geologie	45.09	non	0
5	4	2012-08-07	roux	35	6	28533	F	licence	detective	81.91	non	0
6	5	2014-02-12	chatain	37	8	38558	M	master	geologie	63.46	non	1
7	6	2013-11-11	brun	33	12	39476	M	master	geologie	50.20	oui	0
8	7	2012-03-10	roux	31	10	42392	M	licence	forage	62.20	oui	0
9	8	2014-10-17	chatain	43	10	28625	M	doctorat	geologie	65.17	non	1
10	9	2011-06-04	chatain	28	11	32454	M	master	forage	66.93	non	1
11	10	2014-08-06	brun	50	12	38516	F	licence	geologie	58.93	non	0
12	11	2010-04-22	blond	43	8	38719	M	master	archeologie	80.07	oui	0
13	12	2011-10-23	brun	44	9	28688	M	doctorat	forage	115.76	oui	0
14	13	2014-06-25	brun	39	4	30822	F	licence	geologie	66.85	non	0
15	14	2011-10-15	roux	23	14	33882	M	licence	geologie	61.46	non	0
16	15	2010-07-12	brun	38	3	40889	M	master	geologie	62.85	oui	0
17	16	2012-02-05	brun	31	10	28507	M	doctorat	forage	104.26	oui	0
18	17	2010-07-03	blond	37	13	41512	M	bac	geologie	55.91	non	0
19	18	2010-07-21	blond	35	14	30359	M	licence	geologie	94.58	non	0
20	19	2014-07-22	chatain	44	8	28956	M	master	geologie	74.82	non	0
21	20	2010-02-12	brun	31	12	34483	M	master	forage	66.27	oui	0
22	21	2014-05-22	chatain	30	10	39069	F	licence	forage	65.19	non	0
23	22	2010-01-04	brun	39	13	34632	M	master	geologie	81.35	non	0
24	23	2012-03-20	brun	31	8	41582	M	licence	geologie	58.94	non	0
25	24	2012-01-16	blond	35	8	30782	F	licence	geologie	90.31	non	0

Figure20

Data cleaning:

1-Filling in Missing values:

One of the most important preprocessing steps in data mining is finding instances that contain missing values and treating these missing values.

In this project, methods for checking the existence of missing values are first used as follows: *Figure21*

```
> #missing values  
> sum(is.na(data))  
[1] 27  
>
```

Figure 21

Where the function “IS.NA ()” was used to check if there are missing values or not in the dataset and the “SUM” was used to sum all missing values in the dataset.

total number of missing values in the dataset is 27 which indicates that there are a lot of missing values that we need to fill in the dataset.

```
> #replace missing data with average value
> data$date=ifelse(is.na(data$date),ave(data$date,FUN = function(x)mean(x,na.rm=TRUE)),data$date)
> data$cheveux=ifelse(is.na(data$cheveux),ave(data$cheveux,FUN = function(x)mean(x,na.rm=TRUE)),data$cheveux)
> data$age=ifelse(is.na(data$age),ave(data$age,FUN = function(x)mean(x,na.rm=TRUE)),data$age)
> data$exp=ifelse(is.na(data$exp),ave(data$exp,FUN = function(x)mean(x,na.rm=TRUE)),data$exp)
> data$salaire=ifelse(is.na(data$salaire),ave(data$salaire,FUN = function(x)mean(x,na.rm=TRUE)),data$salaire)
> data$sexe=ifelse(is.na(data$sexe),ave(data$sexe,FUN = function(x)mean(x,na.rm=TRUE)),data$sexe)
> data$diplome=ifelse(is.na(data$diplome),ave(data$diplome,FUN = function(x)mean(x,na.rm=TRUE)),data$diplome)
> data$specialite=ifelse(is.na(data$specialite),ave(data$specialite,FUN = function(x)mean(x,na.rm=TRUE)),data$specialite)
> data$note=ifelse(is.na(data$note),ave(data$note,FUN = function(x)mean(x,na.rm=TRUE)),data$note)
> sum(is.na(data))
[1] 0
```

Figure22

We replaced the missing values with the average value and the sum of the missing values at the end of the replacement process became zero. *Figure22*

2-Outliers Removal

Outliers should be removed from the dataset since they can affect on the accuracy of the data mining models. In this preprocessing step, outliers were detected and removed from the dataset.

Removing outlier from salaire:

```
> boxplot(data$salaire)
> boxplot(data$salaire)$out
[1] 19925 50575 14128 49761
> outliers2 <- boxplot(data$salaire, plot=FALSE)$out
> print(outliers2)
[1] 19925 50575 14128 49761
> data[which(data$salaire %in% outliers2),]
```

Figure23

salaire boxplot after removing outlier

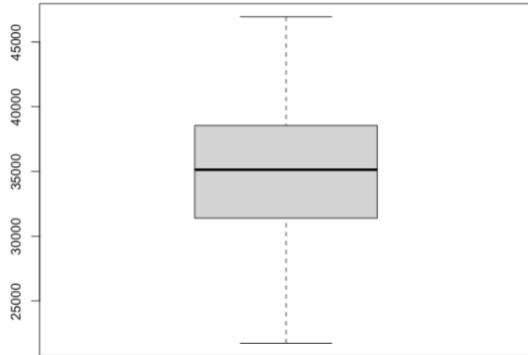


Figure25

note boxplot after removing outlier

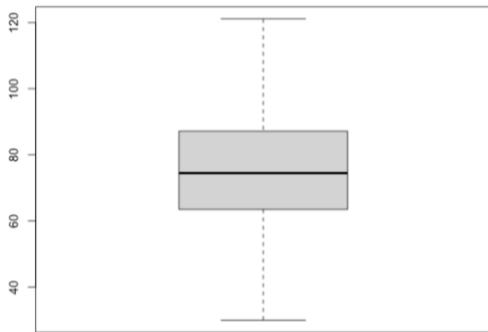


Figure27

Removing outlier from age:

```
> boxplot(data$age)
> boxplot(data$age)$out
[1] 61 65 8 62 8 8 5 63 5 62 62 62 8
> outliers1 <- boxplot(data$age, plot=FALSE)$out
> print(outliers1)
[1] 61 65 8 62 8 8 5 63 5 62 62 62 8
> data[which(data$age %in% outliers1),]
```

Figure24

age boxplot after removing outlier

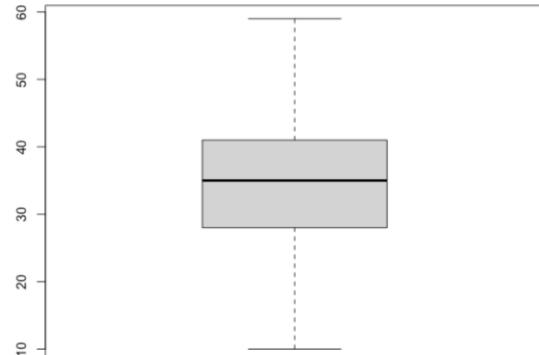


Figure26

Removing outlier from note

```
> boxplot(data$note)
> boxplot(data$note)$out
[1] 127.58 125.47 12.63
> outliers3 <- boxplot(data$note, plot=FALSE)$out
> print(outliers3)
[1] 127.58 125.47 12.63
> data[which(data$note %in% outliers3),]
```

Figure28

Removing outlier from exp:

```
> #data preprccising
> #removing outliers
> library(outliers)
> #data$exp[which(data$exp>420)] <- c(data$exp[which(data$exp>420)]*2)
> boxplot(data$exp)
> boxplot(data$exp)$out
[1] 3 16 3 3 3 3 2 3 17 17 3 3 0 16 3 3 3 3 0 16 16 2 2 3 3 2 2 3 3 3 3 16 17 16 16 16 16 16 3
[41] 16 16 17 16
> outliers <- boxplot(data$exp, plot=FALSE)$out
> print(outliers)
[1] 3 16 3 3 3 3 2 3 17 17 3 3 0 16 3 3 3 3 0 16 16 2 2 3 3 2 2 3 3 3 3 16 17 16 16 16 16 16 3
[41] 16 16 17 16
> data[which(data$exp %in% outliers),]
```

Figure29

Exp boxplot after removing

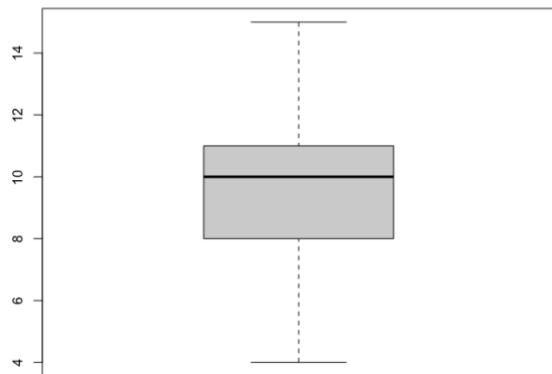


Figure30

Data transformation:

3-Normalization

Normalization of features can speed up training of some data mining models since we have attributes in different scale, we need to normalize to bring the dataset to a common scale (between 0 and 1) while keeping the distribution of the variables the same

the dataset before normalization *Figure31*:

X	date	cheveux	age	exp	salaire	sexe	diplome	specialite	note	dispo	embauche
0	2012-06-02	roux	25	9	26803	F	licence	geologie	97.08	non	0
1	2011-04-21	blond	35	13	38166	M	licence	forage	63.86	non	0
2	2012-09-07	blond	29	13	35207	M	licence	geologie	78.50	non	0
3	2011-07-01	brun	NA	12	32442	M	licence	geologie	45.09	non	0
4	2012-08-07	roux	35	6	28533	F	licence	detective	81.91	non	0
5	2014-02-12	chatain	37	8	38558	M	master	geologie	63.46	non	1
6	2013-11-11	brun	33	12	39476	M	master	geologie	50.20	oui	0
7	2012-03-10	roux	31	10	42392	M	licence	forage	62.20	oui	0
8	2014-10-17	chatain	43	10	28625	M	doctorat	geologie	65.17	non	1
9	2011-06-04	chatain	28	11	32454	M	master	forage	66.93	non	1
10	2014-08-06	brun	50	12	38516	F	licence	geologie	58.93	non	0
11	2010-04-22	blond	43	8	38719	M	master	archeologie	80.07	oui	0
12	2011-10-23	brun	44	9	28688	M	doctorat	forage	115.76	oui	0
13	2014-06-25	brun	39	4	30822	F	licence	geologie	66.85	non	0
14	2011-10-15	roux	23	14	33882	M	licence	geologie	61.46	non	0
15	2010-07-12	brun	38	3	40889	M	master	geologie	62.85	oui	0
16	2012-02-05	brun	31	10	28507	M	doctorat	forage	104.26	oui	0
17	2010-07-03	blond	37	13	41512	M	bac	geologie	55.91	non	0
18	2010-07-21	blond	35	14	30359	M	licence	geologie	94.58	non	0
19	2014-07-22	chatain	44	8	28956	M	master	geologie	74.82	non	0
20	2010-02-12	brun	31	12	34483	M	master	forage	66.27	oui	0
21	2014-05-22	chatain	30	10	39069	F	licence	forage	65.19	non	0
22	2010-01-04	brun	39	13	34632	M	master	geologie	81.35	non	0
23	2012-03-20	brun	31	8	41582	M	licence	geologie	58.94	non	0
24	2012-01-16	blond	35	8	30782	F	licence	geologie	90.31	non	0

Figure31

Normalization code : *Figure32*

```
#####
# Normalization of features

str(data)
normalize<- function(x){(x-min(x))/(max(x)-min(x))}

dataN = data
dataN$age <- normalize(data$age)
dataN$salaire <- normalize(data$salaire)
dataN$exp <- normalize(data$exp)
dataN$note <- normalize(data$note)

str(dataN)
```

Figure32

the dataset after normalization: *Figure33*

The screenshot shows the RStudio Source Editor window with the tab 'dataN' selected. The data is presented in a table format with 25 rows and 14 columns. The columns are labeled: X, date, cheveux, age, exp, salaire, sexe, diplome, specialite, note, dispo, and embauche. The data includes various categorical and numerical values, such as dates ranging from 2010-01-04 to 2014-07-22, and numerical values for age and experience.

X	date	cheveux	age	exp	salaire	sexe	diplome	specialite	note	dispo	embauche	
1	0	2012-06-02	roux	0.3733333	0.4545455	0.3180757	F	licence	geologie	0.6570537	non	0
2	1	2011-04-21	blond	0.5066667	0.6363636	0.6032272	M	licence	forage	0.4101382	non	0
3	2	2012-09-07	blond	0.4266667	0.6363636	0.5289719	M	licence	geologie	0.5189535	non	0
4	3	2011-07-01	brun	0.5067269	0.5909091	0.4595849	M	licence	geologie	0.2706258	non	0
5	4	2012-08-07	roux	0.5066667	0.3181818	0.3614896	F	licence	detective	0.5442991	non	0
6	5	2014-02-12	chatain	0.5333333	0.4090909	0.6130643	M	master	geologie	0.4071652	non	1
7	6	2013-11-11	brun	0.4800000	0.5909091	0.6361013	M	master	geologie	0.3086071	oui	0
8	7	2012-03-10	roux	0.4533333	0.5000000	0.7092775	M	licence	forage	0.3977999	oui	0
9	8	2014-10-17	chatain	0.6133333	0.5000000	0.3637983	M	doctorat	geologie	0.4198751	non	1
10	9	2011-06-04	chatain	0.4133333	0.5454545	0.4598861	M	master	forage	0.4329567	non	1
11	10	2014-08-06	brun	0.7066667	0.5909091	0.6120103	F	licence	geologie	0.3734949	non	0
12	11	2010-04-22	blond	0.6133333	0.4090909	0.6171046	M	master	archeologie	0.5306229	oui	0
13	12	2011-10-23	brun	0.6266667	0.4545455	0.3653793	M	doctorat	forage	0.7958971	oui	0
14	13	2014-06-25	brun	0.5600000	0.2272727	0.4189315	F	licence	geologie	0.4323621	non	0
15	14	2011-10-15	roux	0.3466667	0.6818182	0.4957213	M	licence	geologie	0.3922997	non	0
16	15	2010-07-12	brun	0.5466667	0.1818182	0.6715601	M	master	geologie	0.4026312	oui	0
17	16	2012-02-05	brun	0.4533333	0.5000000	0.3608372	M	doctorat	forage	0.7104207	oui	0
18	17	2010-07-03	blond	0.5333333	0.6363636	0.6871942	M	bac	geologie	0.3510480	non	0
19	18	2010-07-21	blond	0.5066667	0.6818182	0.4073126	M	licence	geologie	0.6384718	non	0
20	19	2014-07-22	chatain	0.6266667	0.4090909	0.3721047	M	master	geologie	0.4916010	non	0
21	20	2010-02-12	brun	0.4533333	0.5909091	0.5108033	M	master	forage	0.4280511	oui	0
22	21	2014-05-22	chatain	0.4400000	0.5000000	0.6258877	F	licence	forage	0.4200238	non	0
23	22	2010-01-04	brun	0.5600000	0.6363636	0.5145424	M	master	geologie	0.5401368	non	0
24	23	2012-03-20	brun	0.4533333	0.4090909	0.6889508	M	licence	geologie	0.3735692	non	0
25	24	2012-01-16	blond	0.5066667	0.4090909	0.4179277	F	licence	geologie	0.6067341	non	0

Figure33

4-Encoding

We discovered that we have 4 nominal attributes(cheveux,diplome,dispo,specialite) that needs to be encoded to numeric values so we can do calculations on it.

the dataset before encoding: *Figure34*

X	date	cheveux	age	exp	salaire	sexe	diplome	specialite	note	dispo	embauche
0	2012-06-02	roux	25	9	26803	F	licence	geologie	97.08	non	0
1	2011-04-21	blond	35	13	38166	M	licence	forage	63.86	non	0
2	2012-09-07	blond	29	13	35207	M	licence	geologie	78.50	non	0
3	2011-07-01	brun	NA	12	32442	M	licence	geologie	45.09	non	0
4	2012-08-07	roux	35	6	28533	F	licence	detective	81.91	non	0
5	2014-02-12	chatain	37	8	38558	M	master	geologie	63.46	non	1
6	2013-11-11	brun	33	12	39476	M	master	geologie	50.20	oui	0
7	2012-03-10	roux	31	10	42392	M	licence	forage	62.20	oui	0
8	2014-10-17	chatain	43	10	28625	M	doctorat	geologie	65.17	non	1
9	2011-06-04	chatain	28	11	32454	M	master	forage	66.93	non	1
10	2014-08-06	brun	50	12	38516	F	licence	geologie	58.93	non	0
11	2010-04-22	blond	43	8	38719	M	master	archeologie	80.07	oui	0
12	2011-10-23	brun	44	9	28688	M	doctorat	forage	115.76	oui	0
13	2014-06-25	brun	39	4	30822	F	licence	geologie	66.85	non	0
14	2011-10-15	roux	23	14	33882	M	licence	geologie	61.46	non	0
15	2010-07-12	brun	38	3	40889	M	master	geologie	62.85	oui	0
16	2012-02-05	brun	31	10	28507	M	doctorat	forage	104.26	oui	0
17	2010-07-03	blond	37	13	41512	M	bac	geologie	55.91	non	0
18	2010-07-21	blond	35	14	30359	M	licence	geologie	94.58	non	0
19	2014-07-22	chatain	44	8	28956	M	master	geologie	74.82	non	0
20	2010-02-12	brun	31	12	34483	M	master	forage	66.27	oui	0
21	2014-05-22	chatain	30	10	39069	F	licence	forage	65.19	non	0
22	2010-01-04	brun	39	13	34632	M	master	geologie	81.35	non	0
23	2012-03-20	brun	31	8	41582	M	licence	geologie	58.94	non	0
24	2012-01-16	blond	35	8	30782	F	licence	geologie	90.31	non	0

Figure34

The encoding code: *Figure35*

```
#####
#ENCODING

data$dispo = factor(data$dispo, levels = c("non", "oui"), labels = c(0, 1))

data$cheveux = factor(data$cheveux, levels = c("chatain", "brun", "blond", "roux"), labels = c(1,2,3,4))

data$diplome = factor(data$diplome, levels = c("doctorat", "master", "bac", "licence"), labels = c(1,2,3,4))
data$specialite = factor(data$specialite, levels = c("geologie", "forage", "detective", "archeologie"), labels = c(1,2,3,4))
str(data)
```

Figure35

the dataset after encoding: *Figure36*

	X	date	cheveux	age	exp	salaire	sexe	diplome	specialite	note	dispo	embauche
1	0	02/06/2012	4	25.00000	9.000000	26803.00	F	4	1	97.08000	0	0
2	1	21/04/2011	3	35.00000	13.000000	38166.00	M	4	2	63.86000	0	0
3	2	07/09/2012	3	29.00000	13.000000	35207.00	M	4	1	78.50000	0	0
4	3	01/07/2011	2	34.84839	12.000000	32442.00	M	4	1	45.09000	0	0
5	4	07/08/2012	4	35.00000	6.000000	28533.00	F	4	3	81.91000	0	0
6	5	12/02/2014	1	37.00000	8.000000	38558.00	M	2	1	63.46000	0	1
7	6	11/11/2013	2	33.00000	12.000000	39476.00	M	2	1	50.20000	1	0
8	7	10/03/2012	4	31.00000	10.000000	42392.00	M	4	2	62.20000	1	0
9	8	17/10/2014	1	43.00000	10.000000	28625.00	M	1	1	65.17000	0	1
10	9	04/06/2011	1	28.00000	11.000000	32454.00	M	2	2	66.93000	0	1
11	10	06/08/2014	2	50.00000	12.000000	38516.00	F	4	1	58.93000	0	0
12	11	22/04/2010	3	43.00000	8.000000	38719.00	M	2	4	80.07000	1	0
13	12	23/10/2011	2	44.00000	9.000000	28688.00	M	1	2	115.76000	1	0
14	13	25/06/2014	2	39.00000	4.000000	30822.00	F	4	1	66.85000	0	0
15	14	15/10/2011	4	23.00000	14.000000	33882.00	M	4	1	61.46000	0	0
17	16	05/02/2012	2	31.00000	10.000000	28507.00	M	1	2	104.26000	1	0
18	17	03/07/2010	3	37.00000	13.000000	41512.00	M	3	1	55.91000	0	0
19	18	21/07/2010	3	35.00000	14.000000	30359.00	M	4	1	94.58000	0	0
20	19	22/07/2014	1	44.00000	8.000000	28956.00	M	2	1	74.82000	0	0
21	20	12/02/2010	2	31.00000	12.000000	34483.00	M	2	2	66.27000	1	0
22	21	22/05/2014	1	30.00000	10.000000	39069.00	F	4	2	65.19000	0	0
23	22	04/01/2010	2	39.00000	13.000000	34632.00	M	2	1	81.35000	0	0
24	23	20/03/2012	2	31.00000	8.000000	41582.00	M	4	1	58.94000	0	0
25	24	16/01/2012	3	35.00000	8.000000	30782.00	F	4	1	90.31000	0	0
26	25	13/05/2014	1	31.00000	12.000000	33238.00	F	2	3	70.39000	0	0
27	26	12/04/2013	2	33.00000	8.000000	35075.23	F	4	3	56.99000	0	0
28	27	18/06/2013	2	42.00000	8.000000	30641.00	F	4	1	78.89000	0	0
29	28	02/01/2010	1	27.00000	13.000000	31505.00	F	2	4	88.33000	0	0
30	29	25/07/2010	1	19.00000	9.000000	31357.00	M	1	3	88.69000	1	0
31	30	21/03/2010	2	19.00000	4.000000	39031.00	M	2	1	57.48000	1	0
32	31	29/05/2014	1	40.00000	13.000000	37495.00	M	1	1	80.86000	1	0
33	32	05/02/2012	1	43.00000	7.000000	37579.00	M	1	1	91.15000	1	0
34	33	03/09/2012	3	28.00000	13.000000	34485.00	F	4	2	76.94000	0	0
35	34	21/06/2010	3	38.00000	6.000000	41909.00	M	4	1	87.76000	1	1
36	35	04/11/2011	1	43.00000	8.000000	46245.00	M	4	1	49.39000	0	0
37	36	01/08/2010	2	49.00000	8.000000	38020.00	M	2	2	60.63000	1	0
38	37	10/01/2013	2	26.00000	11.000000	31759.00	F	2	2	102.66000	0	1
39	38	14/06/2013	3	59.00000	7.000000	40397.00	M	4	2	60.42000	1	0
40	39	08/11/2014	1	37.00000	6.000000	27002.00	M	2	2	55.83000	0	1

Figure36

5 Data Mining Technique

1

Classification:

النكتيك الذي استخدمناها عشان نحل المسألة

We used the Decision tree classifier technique to classify candidates who applied for gold digger position into hired (`embauche = 1`) or not hired (`embauche = 0`) using the target value (`embauche`) since we're doing classification the dataset class label and number of classes are known. to do classification we need to split the data into training set and test set. Using `training data` set to determine the classification model. Then, Using the `test data` with the classification model to estimate accuracy rate of the classification model. We used Decision tree because it generates understandable rules. It performs classification without requiring much computation.

How we used technique?

Package:
الباكيجز والمكتوب
الي المستخدمون في
الكتور

party "for recursive partitioning"

caret "loads packages as needed and assume they are installed"

rpart "builds the classification tree"

rpart.plot "print / plot the tree"

Methods:

`nrow` " method to return the number of rows of data"

`predict` "makes prediction for new data"

`confusionMatrix` "prediction table "

②

Clustering:

التكنيك اللي استخدمنا هي عشان تحل المسئله

How we used technique?

→ using 3 different sizes (2,4,6)

We used K-means clustering technique algorithm that first partition data into groups where data in the same group are similar, K-mean work by randomly select objects as cluster centers and assigns other objects to the nearest cluster center and then improves the clustering by iteratively updating the cluster centers and reassigning the objects to new centers. We used k-mean because it relatively simple to implement , it suitable to large data sets it also guarantees convergence, can warm-start the positions of centroids. K-mean generalizes to clusters of different shapes and sizes.

يعني
وتحمن لنا أن المطلوب في كل كلسستره تغير
متقاربة

We used these packages:

centroids ④

cluster

⑤ ملائمه للبيانات الكبيرة

factoextra(visualize the clusters)



NbClust(determine the optimal number)



GGally, plotly (plotting the cluster)

also this libraries: NbClust , cluster , factoextra, GGally, plotly

في البيانات حققتنا كل يوم من نوموك والكلسترنق عشان يتسمى لازم كل الأعمة نوموك خولنا كل اعدهنا نوموك

* and because we have non numeric attributes, we convert our attributes to be numeric and we removed the class label form our data set.

```
> datacl$embauche<-NULL
> datacl$diplome<-as.numeric(as.character(datacl$diplome))
> datacl$sexe<-as.numeric(as.character(datacl$sexe))
> datacl$cheveux<-as.numeric(as.character(datacl$cheveux))
> datacl$specialite<-as.numeric(as.character(datacl$specialite))
> datacl$dispo<-as.numeric(as.character(datacl$dispo))
> datacl$specialite<-as.numeric(as.character(datacl$specialite))
> sapply(datacl,class)
  cheveux    age      exp    salaire     sexe   diplome specialite      note
"numeric" "numeric" "numeric" "numeric" "numeric" "numeric" "numeric"
  dispo
"numeric"
```

لأن الكلسترنق احنا ننسوي اذا ما عندنا كلاس ليبل
فهنا عشان ننسوي كلسترنق نفترض انه ما عندنا
كلاس ليبل فنحطه العاود له null اللي هو اسم
الكلاس ليبل عندنا

Methods:

- Set.seed(): Reproduce results.
- Scale(): Scaling and centering of matrix-like objects.

- Kmeans(): Classify observations into k groups, based on their similarity. Each group is represented by the mean value of points in the group, known as the cluster centroid.
- Fviz_cluster(): Visualize clustering results.
→ Always must be between $(-1, 1)$, A good average Silhouette is when average silhouette is = 1
- silhouette(): find the average silhouette for each cluster and all cluster
- Fviz_nbclust(): Determining and visualizing the optimal number of clusters.
- ggparcoord(): interpretation clusters relation by choosing specific attributes.

6 Evaluation and Comparison

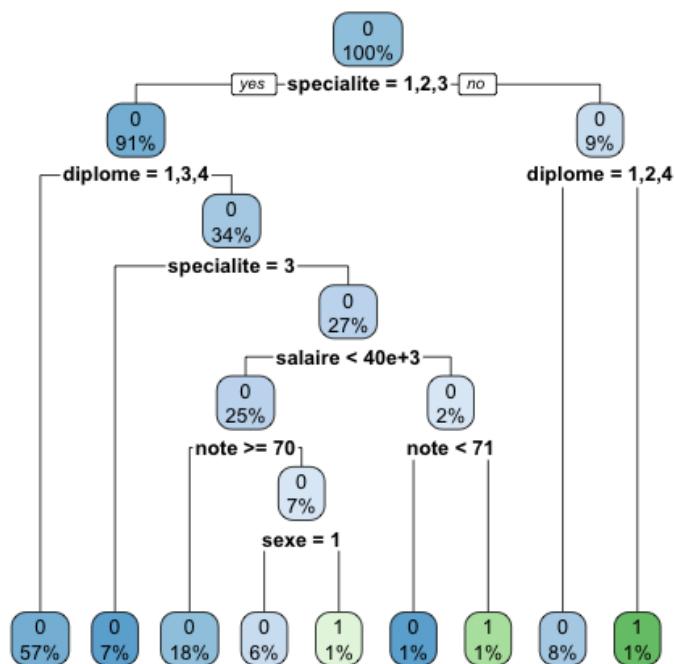


Classification:

Decision tree (1): first Siz e

In Decision tree (1) we divided the data set into training set with percentage (90%) and test set with percentage (10%).

Decision tree (1)



اذنات وش
فتحوا من الراوية
و خواه الراوية



Using `rpart.plot` we plot the Decision tree (1) it shows us that the root is `specialite` and that indicated that it has the highest information gain of all attributes.

is:

The best/ most valuable attribute in the tree and is selected to split the data

Confusion Metrix for (1)

```
R 4.1.1 : ~/Desktop/326 project/ > confusionMatrix(xtab)
Confusion Matrix and Statistics

predicted1  0  1
      0 86 18
      1  3  1
      Accuracy : 0.8056
      95% CI : (0.7183, 0.8754)
No Information Rate : 0.8241
P-Value [Acc > NIR] : 0.74153

      Kappa : 0.0274

McNemar's Test P-Value : 0.00225

      Sensitivity : 0.96629
      Specificity : 0.05263
Pos Pred Value : 0.82692
Neg Pred Value : 0.25000
      Prevalence : 0.82407
Detection Rate : 0.79630
Detection Prevalence : 0.96296
Balanced Accuracy : 0.50946

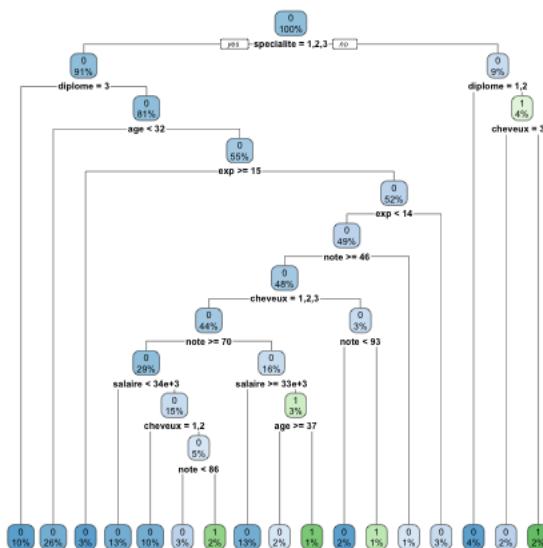
      'Positive' Class : 0

> precision(xtab)
[1] 0.8269231
>
```

نفسي أول وحدة Decision tree (2):

In Decision tree (2) we divided the data set into training set with percentage (80%) and test set with percentage (20%).

Decision tree (2)



Using rpart.plot we plot the Decision tree (2) it shows us that the root is also specialite and that indicated that it has the highest information gain of all attributes.

Confusion Metrix for (2)

```
R 4.1.1 · ~/Desktop/326 project/ 
> confusionMatrix(xtab )
Confusion Matrix and Statistics

predicted\0\1
\0 152 20
\1 5 3

Accuracy : 0.8611
95% CI : (0.8018, 0.9081)
No Information Rate : 0.8722
P-Value [Acc > NIR] : 0.71807

Kappa : 0.1366

McNemar's Test P-Value : 0.00511

Sensitivity : 0.9682
Specificity : 0.1304
Pos Pred Value : 0.8837
Neg Pred Value : 0.3750
Prevalence : 0.8722
Detection Rate : 0.8444
Detection Prevalence : 0.9556
Balanced Accuracy : 0.5493

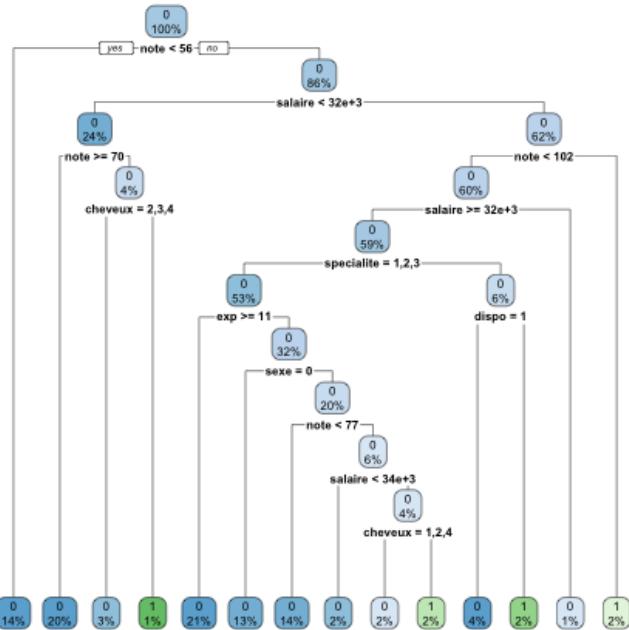
'Positive' Class : 0

> precision(xtab)
[1] 0.8837209
>
```

Decision tree (3): نفس أول ووحده

In Decision tree (3) we divided the data set into training set with percentage (70%) and test set with percentage (30%).

Decision tree (3)



Using rpart.plot we plot the Decision tree (3) it shows us that the root is note and that indicated that it has the highest information gain of all attributes.

Confusion Metrix for (3)

```
R 4.1.1 · ~/Desktop/326 project/
> confusionMatrix(xtab)
Confusion Matrix and Statistics

predicted1   0   1
      0 227 28
      1 16 14

Accuracy : 0.8456
95% CI : (0.7983, 0.8855)
No Information Rate : 0.8526
P-Value [Acc > NIR] : 0.66780

Kappa : 0.3033

McNemar's Test P-Value : 0.09725

Sensitivity : 0.9342
Specificity : 0.3333
Pos Pred Value : 0.8902
Neg Pred Value : 0.4667
Prevalence : 0.8526
Detection Rate : 0.7965
Detection Prevalence : 0.8947
Balanced Accuracy : 0.6337

'Positive' Class : 0

> precision(xtab)
[1] 0.8901961
>
```

Evaluation matrix table

Decision trees	Decision tree (1): 90% - 10%	Decision tree (2) 80% - 20% <i>نحوه اول</i> <i>نحوه اول</i>	Decision tree (3) 70% - 30%
Accuracy	80%	86% <i>highest accuracy</i>	84%
precision	82%	88%	89%
sensitivity	96%	96%	93%
specificity	5%	13%	33%

Accuracy: number of classifications a model correctly predicts

Precision: Number of true positive divided by the total number of positive prediction

Sensitivity: predict true positives which is people that hired
Specificity: predict true negative which is people that's not hired

Clustering:

جربنا 3 سایزات مختلفة

Size (1):

In size (1) we chose number of clusters to be k=2.

Size (2):

In size (2) we chose number of clusters to be k=4.

Size (3):

In size (3) we chose number of clusters to be k=6.

	Size 1	Clustering table	Size 2	Size 3
Clusters	K=2	K=4		K=6
Silhouette width for each cluster	<pre>> avrsil<-silhouette(kmeans, results\$cluster) > fviz_silhouette(avrsil)</pre> <p>Figure2</p>	<pre>> fviz_silhouette(avrsil) cluster size ave.sil.width 1 1 176 0.53 2 2 156 0.54 3 3 284 0.54 4 4 298 0.51 > </pre> <p>Figure4</p>	<pre>> fviz_silhouette(avrsil) cluster size ave.sil.width 1 1 136 -0.48 2 2 67 0.54 3 3 61 0.51 4 4 246 0.53 5 5 216 0.52 6 6 188 0.53 > </pre> <p>Figure6</p>	
Silhouette width for all clusters	<p>average silhouette</p> <p>Cluster 1 width=0.56</p> <p>Cluster 2 width=0.55</p> <p>$\frac{0.56+0.55}{2}$</p>	<p>average silhouette</p> <p>width=0.53</p>	<p>average silhouette</p> <p>width=0.52</p>	
Visualization	Figure1	Figure3	Figure5	

For size 1 (K=2)

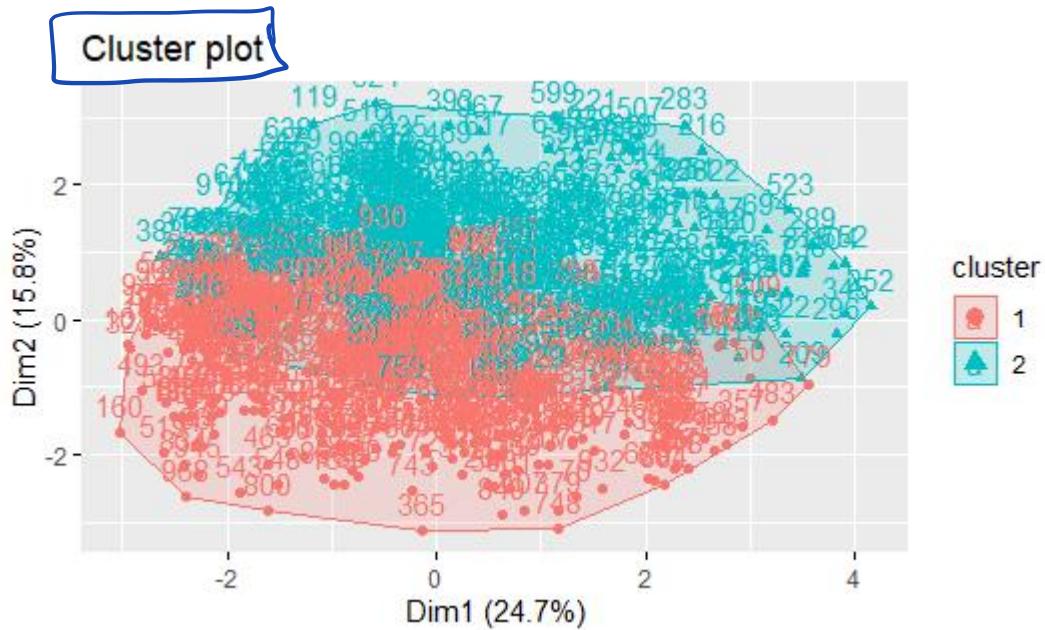


Figure1

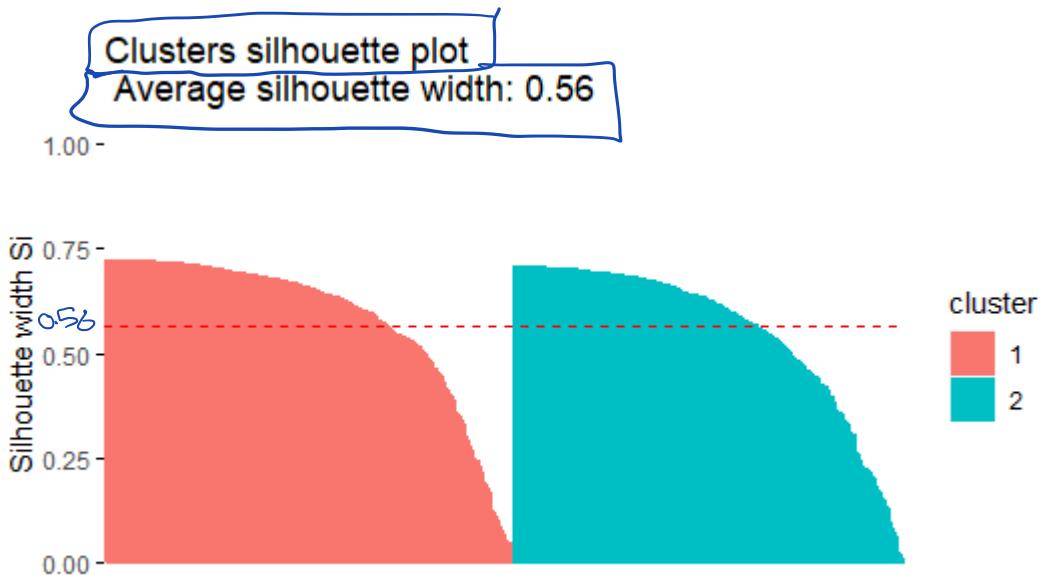


Figure2 K=2

For Sizec2) ($K=4$)

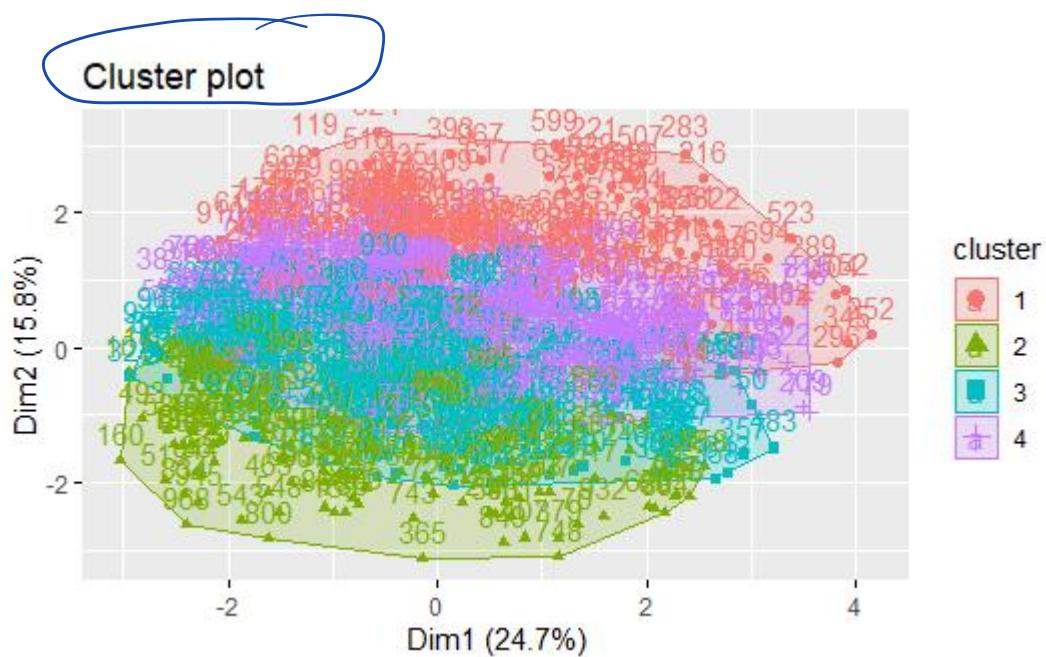


Figure3

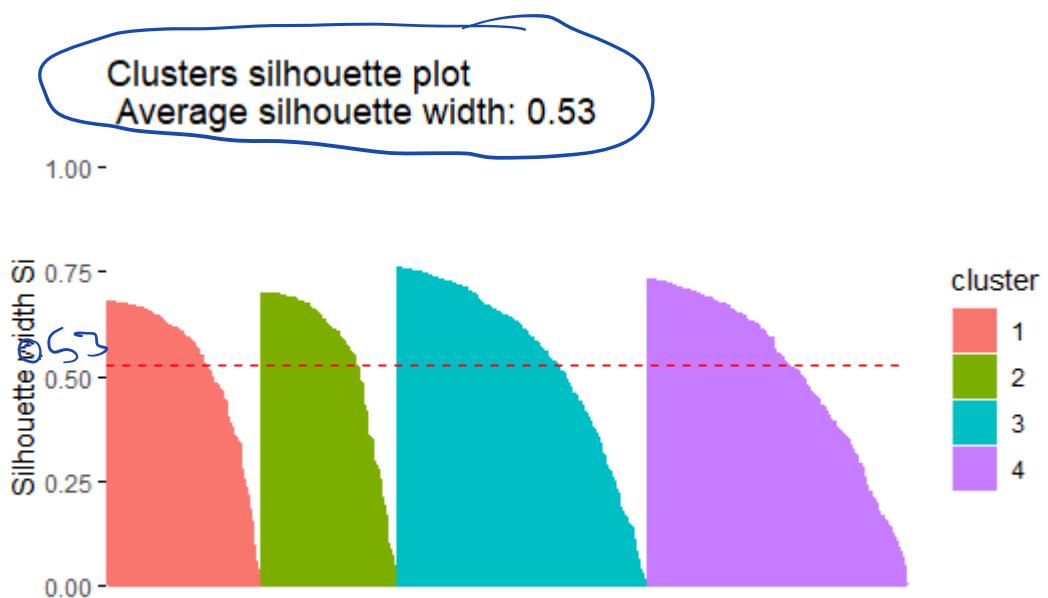


Figure4 $K=4$

For Size 3 ($K=6$)

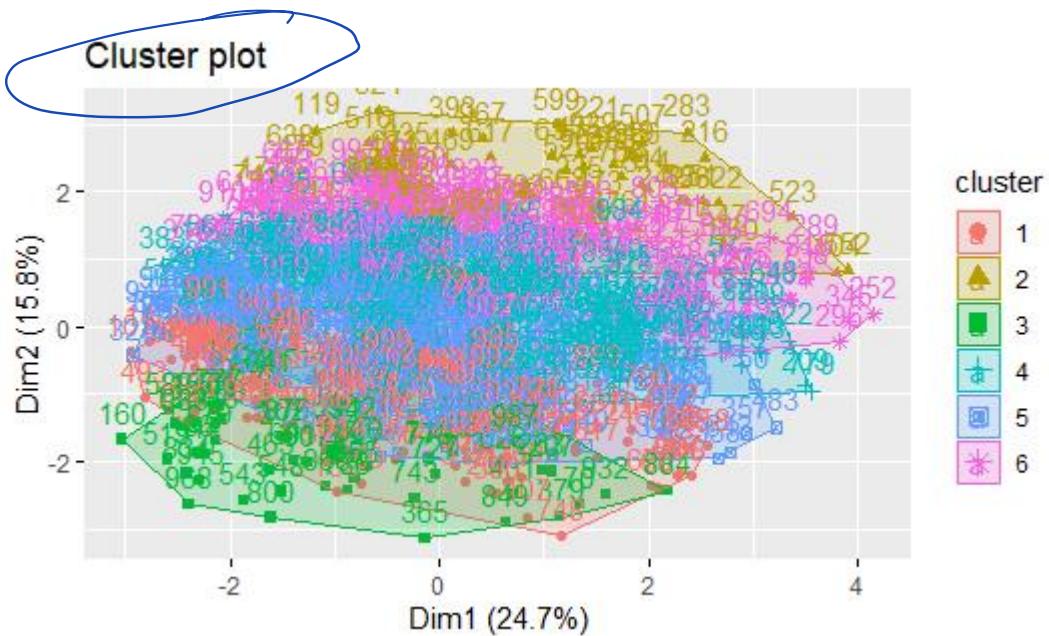


Figure5

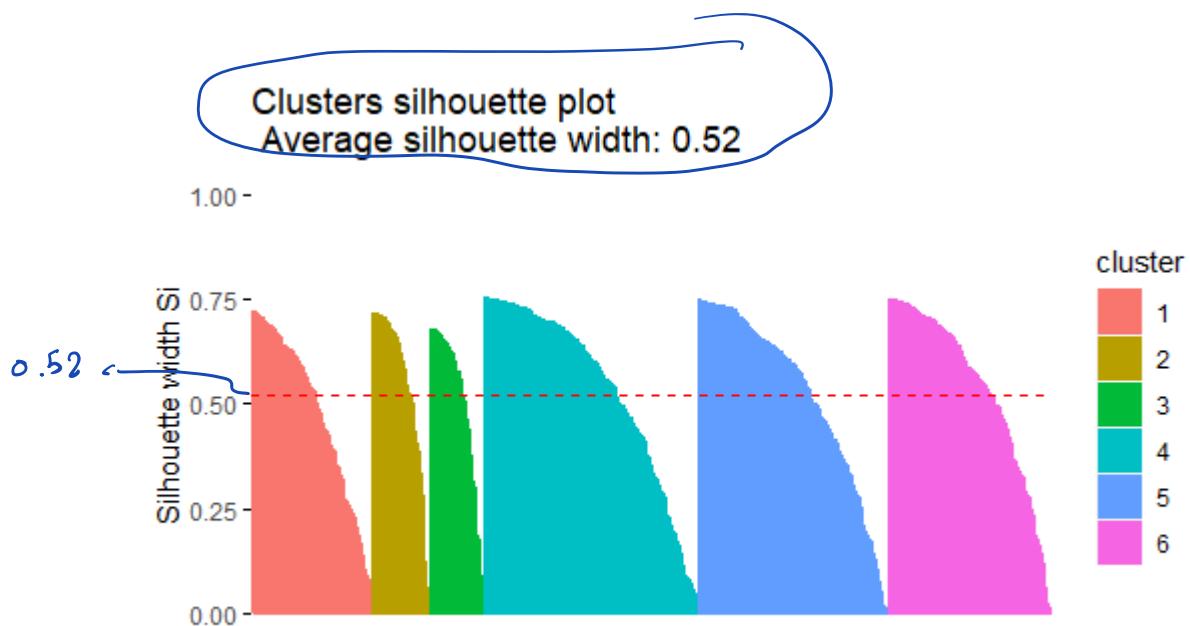


Figure 6 $K=6$

هنا زي الكلاسفيشن، لازم نسوي **evaluate** عشان نختار مين افضل سايز عندنا واعطانا قيمة اقرب لل ١ ، استخدمنا طريقتين وكلهم اشاروا ان $k=2$ افضل يعني اتنا نقصم الكستر الى مجموعتين يعطينا نتائج افضل ويسمى $k=2$ في هذه الحالة optimal number

We used two ways to find the optimal number of clusters to make sure, the first way we used: fviz_nbclust using Silhouette method, and the second way is NbClust Validation, and in both ways they indicates k=2 to be the optimal number of cluster.

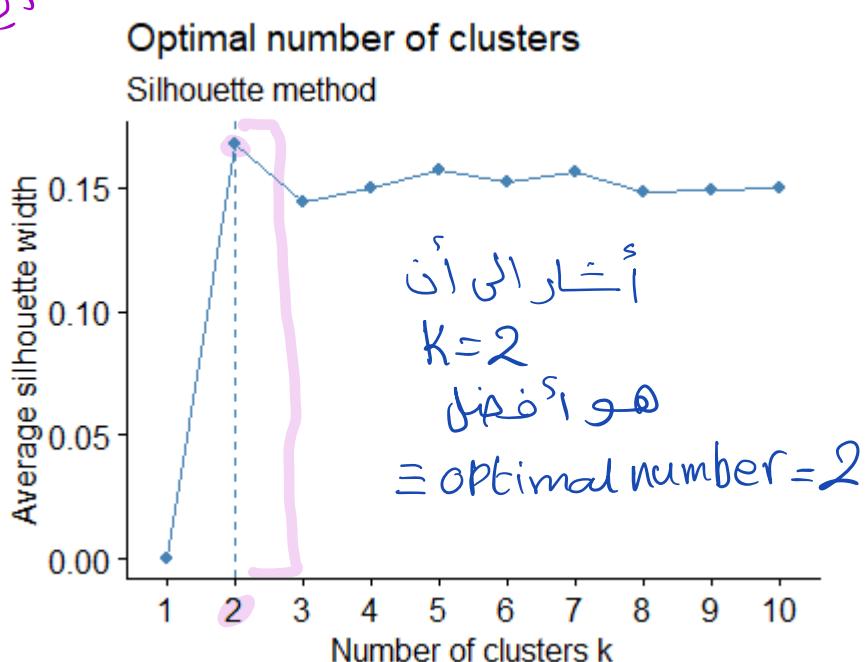
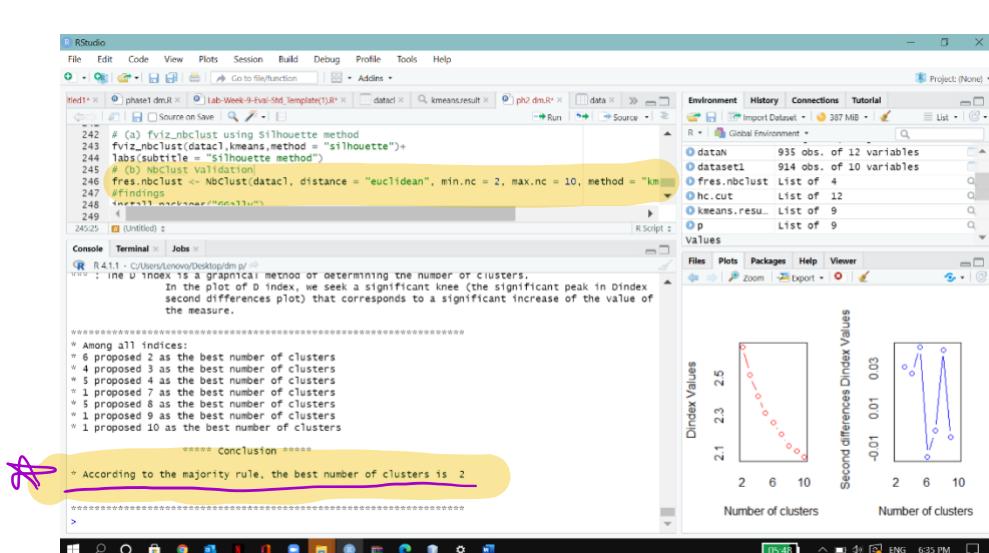


Figure 7 The best number of $k = 2$



هنا كل واحدة فيها بشرح الفايندنقر كل نقطة فيها مهمة
لازم نفهم كل شيء

7 Findings

Classification:

First, we did Mining the data by using classification technique and used tree model, use rpart, predict, confusionMatrix, We applied tree algorithm on 3 different sizes: -training set 90% and testing set 10% - training set 80% and testing set 20% and 70% training 30% testing.

after calculating the accuracy for the 3 different sizes we conclude that the best result of accuracy was “86%” which belongs to (80% training -20% testing). That means most of the tuples was classified correctly based on our class label “embauche”. From the decision tree it was split based on the “specialite” since it has the heights information gain, and it means that “specialite” has major effect on the class label “embauche” which indicate the candidates will be hired or not.

We used rpart method to plot our decision tree because it gave us more reasonable information than ctree, when we tried using ctree it only gave us one node so rpart was more understandable and suitable for our dataset.

Clustering:

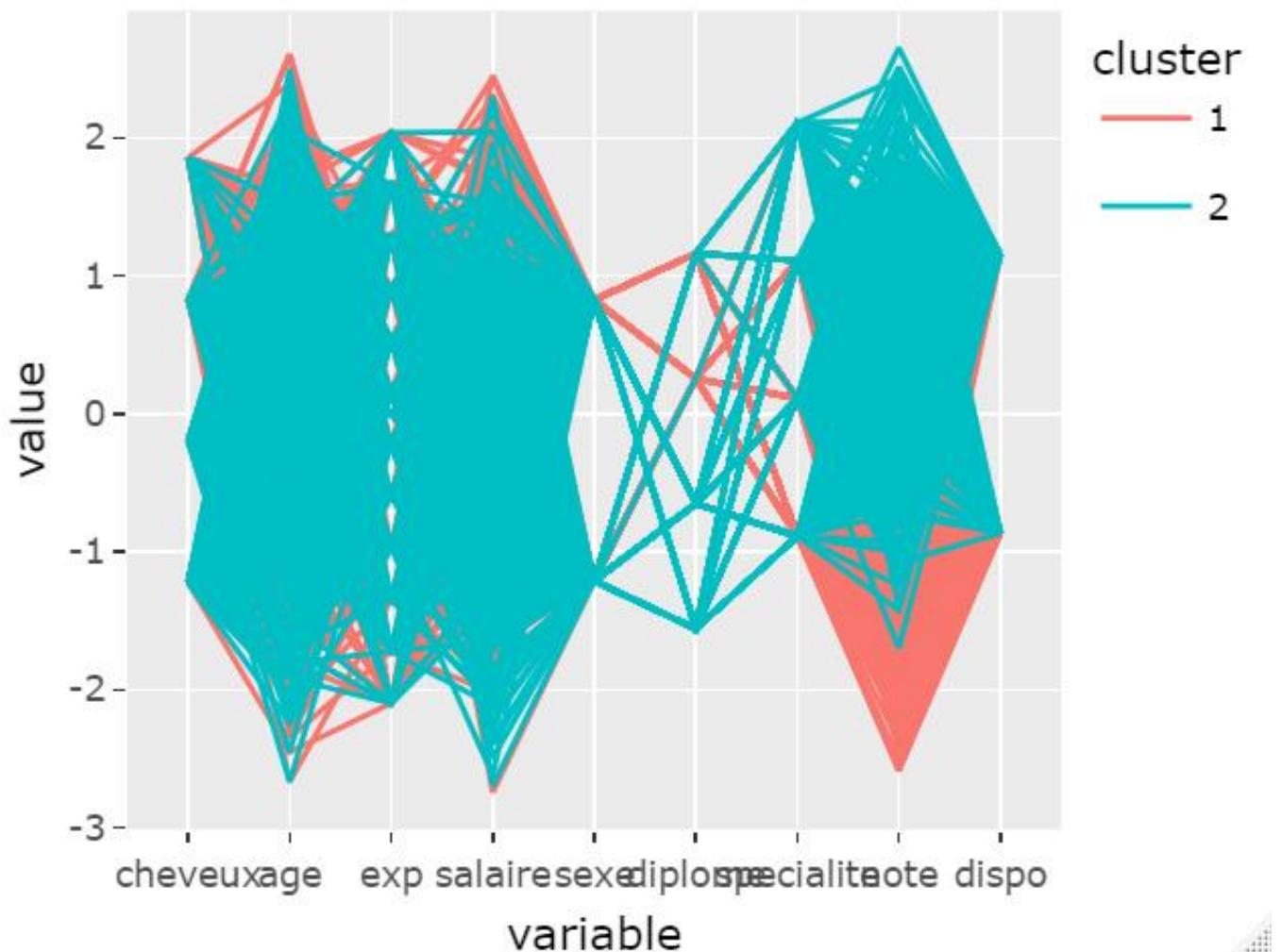
Secondly, we did Mining the data by using clustering technique by k means We want accurate results so we changed the number of clusters by choosing random k (2 ,4,6) to cluster the data based on their similarity. object in a cluster is similar to one another each object is assigned to the cluster to which it is most similar, based on the distance between the object and the cluster center for ex: in cluster one candidates with same or close years of experience and same age rang and same specialty...etc, are clustered together.

so when there is new candidate we can assign the candidate to one of our clusters. Also, we used fviz_cluster() to Visualize the clusters then we used the silhouette() method that gives the average of each cluster in the whole cluster .And we excepted the optimal number is 2 because it has the highest average among the number of k.

When the K=2 the average of cluster 0.56 which is the best average between our clusters, and in each cluster, silhouette are c1= 0.58 for first cluster and c2=0.55 for second cluster. When the K=4 the average of cluster 0.53, and in each cluster silhouette are c1= 0.53 for first cluster and c2= 0.54 for second cluster

$c_3=0.54$ for the third cluster $c_4=0.51$ for the fourth cluster. When the $K=6$ the average of cluster 0.52, and in each cluster silhouette are $c_1= 0.48$ for first cluster and $c_2=0.54$ for second cluster, $c_3=0.51$ for the third cluster $c_4=0.53$ for the fourth cluster , $c_5=0.52$ for the fifth cluster $c_6=0.53$ for the sixth .

we found that the optimal number of clusters (best number of clusters) in our data set is $k=2$ and the average width is 0.56



cluster 1: they have the highest in age, and have the lowest at note.

cluster 2: they have the highest in note, and have the lowest at salaires.

~~inconclusion:~~

we conclude that the classification is accurate in our dataset because the classification accuracy was 86% which is close to 100% which indicates that our classification is accurate and useful and in other hand when we used clustering, the optimal number of clusters was k=2 and silhouette average was 0.56 which is close to 1 which indicates that we almost have high quality cluster.

We reach 100% high quality When
Silhouette average is = 1

8 Code

Preprocessing:

```
1 data<-read.csv("data2.csv")
2 getwd()
3
4 #print No of objects
5 nrow(data)
6
7 # print NUM of features
8 ncol(data)
9
10 #print names of the features
11 names(data)
12
13 # print values of features of first five objects
14 head(data)
15 # print data type of features
16 str(data)
17
18
19 # print five number summary for numeric
20 #summary(data)
21 summary(data$X)
22 summary(data$age)
23 summary(data$exp)
24 summary(data$salaire)
25 summary(data$note)
26 summary(data$embauche)
27
28 #missing values
29 sum(is.na(data))
30 #replace missing data with average value
31
32 #missing values
33 sum(is.na(data))
34
35 #boxplot outlier
36
```

```
35 #boxplot outlier
36
37 boxplot(data$age)
38 boxplot(data$exp)
39 boxplot(data$note)
40 boxplot(data$salaire)
41
42 # histogram
43 hist(data$age)
44 hist(data$exp)
45 hist(data$note)
46 hist(data$salaire)
47
48 barplot(table(data$cheveux))
49 barplot(table(data$date))
50 barplot(table(data$diplome))
51 barplot(table(data$specialite))
52
53 # scatter plot
54 plot(x = data$age, y = data$exp,
55       xlab = "Age",
56       ylab = "exp",
57       main = "Age vs exp"
58 )
59
60 #data preprpccising
61 #removing outliers
62 library(outliers)
63
64 #data$exp[which(data$exp>420)] <- c(data$exp[which(data$exp>420)]*2)
65 boxplot(data$exp)
66 boxplot(data$exp)$out
67 outliers <- boxplot(data$exp, plot=FALSE)$out
68 print(outliers)
69 data[which(data$exp %in% outliers),]
70 data <- data[-which(data$exp %in% outliers),]
```

```
69 data[which(data$exp %in% outliers),]
70 data <- data[-which(data$exp %in% outliers),]
71 boxplot(data$exp)
72
73 boxplot(data$age)
74 boxplot(data$age)$out
75 outliers1 <- boxplot(data$age, plot=FALSE)$out
76 print(outliers1)
77 data[which(data$age %in% outliers1),]
78 data <- data[-which(data$age %in% outliers1),]
79 boxplot(data$age)
80
81 boxplot(data$salaire)
82 boxplot(data$salaire)$out
83 outliers2 <- boxplot(data$salaire, plot=FALSE)$out
84 print(outliers2)
85 data[which(data$salaire %in% outliers2),]
86 data <- data[-which(data$salaire %in% outliers2),]
87 boxplot(data$salaire)
88
89 boxplot(data$note)
90 boxplot(data$note)$out
91 outliers3 <- boxplot(data$note, plot=FALSE)$out
92 print(outliers3)
93 data[which(data$note %in% outliers3),]
94 data <- data[-which(data$note %in% outliers3),]
95 boxplot(data$note)
96
97
98
99 #replace missing data with average value
100 data$date=ifelse(is.na(data$date),ave(data$date,FUN = function(x)mean(x,na.rm=TRUE)),data$date)
101 data$cheveux=ifelse(is.na(data$cheveux),ave(data$cheveux,FUN = function(x)mean(x,na.rm=TRUE)),data$cheveux)
102 data$age=ifelse(is.na(data$age),ave(data$age,FUN = function(x)mean(x,na.rm=TRUE)),data$age)
103
104 data$exp=ifelse(is.na(data$exp),ave(data$exp,FUN = function(x)mean(x,na.rm=TRUE)),data$exp)
```

```
103 data$exp=ifelse(is.na(data$exp),ave(data$exp,FUN = function(x)mean(x,na.rm=TRUE)),data$exp)
104 data$salaire=ifelse(is.na(data$salaire),ave(data$salaire,FUN = function(x)mean(x,na.rm=TRUE)),data$salaire)
105 data$sexe=ifelse(is.na(data$sexe),ave(data$sexe,FUN = function(x)mean(x,na.rm=TRUE)),data$sexe)
106 data$diplome=ifelse(is.na(data$diplome),ave(data$diplome,FUN = function(x)mean(x,na.rm=TRUE)),data$diplome)
107 data$specialite=ifelse(is.na(data$specialite),ave(data$specialite,FUN = function(x)mean(x,na.rm=TRUE)),data$specialite)
108 data$note=ifelse(is.na(data$note),ave(data$note,FUN = function(x)mean(x,na.rm=TRUE)),data$note)
109 sum(is.na(data))
110
111 # Normalization of features
112
113 str(data)
114 normalize<- function(x){(x-min(x))/(max(x)-min(x))}
115
116
117 dataN = data
118 dataN$age <-normalize(data$age)
119 dataN$salaire <-normalize(data$salaire)
120 dataN$exp <- normalize(data$exp)
121 dataN$note <- normalize(data$note)
122
123 str(dataN)
124
125 #####
126 #ENCODING
127
128 data$dispo = factor(data$dispo, levels = c("non", "oui"), labels = c(0, 1))
129
130 data$cheveux = factor(data$cheveux, levels = c("chatain", "brun","blond" ,"roux"), labels = c(1,2,3,4))
131
132 data$diplome = factor(data$diplome, levels = c("doctorat", "master","bac" , "licence"), labels = c(1,2,3,4))
133 data$specialite = factor(data$specialite, levels = c("geologie", "forage","detective" , "archeologie"), labels = c(1,2,3,4))
134 data$sexe = factor(data$sexe, levels = c("F", "M"), labels = c(0, 1))
135 str(data)
136
137
138 data [! is.na(data$embauche),]
139
```

a. classification

→ *النوع كل مجموعات
كود في الكود
التي فوقة*

```
142 #####classification#####
143 library(caret)
144 library(party)
145 library(rpart)
146 library(rpart.plot)
147
148
149
150 #random number generator for objects that can be reproduced
151 set.seed(1234)
152
153 ##### decision tree 1 (90% -10% )#####
154
155 #partition data to 90% to be training set
156 #nrow is a method to return the number of rows of data
157 ind <- sample(2, nrow(data) , replace = TRUE , prob = c(0.9 , 0.1))
158 #store training set value 90% in variable trainData
159 trainData <- data[ind == 1,]
160 #store testing set value 10% in variable testData
161 testData <- data[ind == 2,]
162
163 #creating/building decision tree by saving all attributes (columns) data in variable "model1"
164 #rpart is a machine learning library that is used to build the decesion tree
165 model1<- rpart(embauche ~ cheveux + age + exp + salaire + sexe + diplome + specialite + note + dispo, data = trainData, method = "class")
166
167 # plot decision tree for model1 , extra = 100 it display the number of percentage of observation in node of tree
168 rpart.plot(model1, extra = 100)
169
170 # test decision tree , predict the values based on the input we store the predeected values in variable predicted1
171 predicted1 <- predict(model1, newdata = testData, type = "class")
172
173
174
175 #confusion metrix
176 xtab <- table(predicted1, testData$embauche)
177 confusionMatrix(xtab )
```

```

174
175 #confusion matrix
176 xtab <- table(predicted1, testData$embauche)
177 confusionMatrix(xtab )
178 precision(xtab)
179
180 #####diction tree 2 (80% -20%)#####
181 #partition data to 50% to be training set
182 #nrow is a method to return the number of rows of data
183 ind <- sample(2, nrow(data) , replace = TRUE , prob = c(0.8 , 0.2))
184 #store training set value 90% in variable trainData
185 trainData <- data[ind == 1,]
186 #store testing set value 10% in variable trainData
187 testData <- data[ind == 2,]
188
189
190 # Create the decision tree
191 model1<- rpart(embauche ~ cheveux + age + exp + salaire + sexe + diplome + specialite + note + dispo, data = trainData, method='class')
192
193 # plot the decision tree
194 rpart.plot(model1, extra = 100)
195
196 # test the decision tree
197 predicted1 <- predict(model1, newdata = testData, type = "class")
198
199 #confusion metrix
200 xtab <- table(predicted1, testData$embauche)
201 confusionMatrix(xtab )
202 precision(xtab)
203
204 ##### diction tree 3 (70% -30%)#####
205 #partition data to 70% to be training set
206 #nrow is a method to return the number of rows of data
207 ind <- sample(2, nrow(data) , replace = TRUE , prob = c(0.7 , 0.3))
208 #store training set value 90% in variable trainData
209 trainData <- data[ind == 1,]
210 #store testing set value 10% in variable trainData
211 testData <- data[ind == 2,]
212
213
214 # Create decision tree, Using training data set to determine the classification model
215 model1<- rpart(embauche ~ cheveux + age + exp + salaire + sexe + diplome + specialite + note + dispo, data = trainData, method='class')
216
217 # plot decision tree
218 rpart.plot(model1, extra = 100)
219
220 # test decision tree
221 predicted1 <- predict(model1, newdata = testData, type = "class")
222
223 #confusion metrix
224 xtab <- table(predicted1, testData$embauche)
225 confusionMatrix(xtab )
226 precision(xtab)
227
228
229
230
231
232
233
234
235

```

Evaluation

هنا نسخ
 كل كود نفس اللي فوق
 فقط يغير
 اللي تغيروا

١٢

```

200 xtab <- table(predicted1, testData$embauche)
201 confusionMatrix(xtab )
202 precision(xtab)
203
204 ##### diction tree 3 (70% -30%)#####
205 #partition data to 70% to be training set
206 #nrow is a method to return the number of rows of data
207 ind <- sample(2, nrow(data) , replace = TRUE , prob = c(0.7 , 0.3))
208 #store training set value 90% in variable trainData
209 trainData <- data[ind == 1,]
210 #store testing set value 10% in variable trainData
211 testData <- data[ind == 2,]
212
213
214 # Create decision tree, Using training data set to determine the classification model
215 model1<- rpart(embauche ~ cheveux + age + exp + salaire + sexe + diplome + specialite + note + dispo, data = trainData, method='class')
216
217 # plot decision tree
218 rpart.plot(model1, extra = 100)
219
220 # test decision tree
221 predicted1 <- predict(model1, newdata = testData, type = "class")
222
223 #confusion metrix
224 xtab <- table(predicted1, testData$embauche)
225 confusionMatrix(xtab )
226 precision(xtab)
227
228
229
230
231
232
233
234
235

```

Evaluation

b. clustering

```
165 #clustering
166 datacl<-data
167 datacl$i..<-NULL
168 datacl$date<-NULL
169 datacl$embauche<-NULL
170
171 datacl$diplome<-as.numeric(as.character(datacl$diplome))
172 datacl$sexe<-as.numeric(as.character(datacl$sex))
173 datacl$cheveux<-as.numeric(as.character(datacl$cheveux))
174 datacl$specialite<-as.numeric(as.character(datacl$specialite))
175 datacl$dispo<-as.numeric(as.character(datacl$dispo))
176 datacl$specialite<-as.numeric(as.character(datacl$specialite))
177 #check if all attributes are numeric
178 sapply(datacl,class)
179 sum(is.na(datacl))
180 datacl<-na.omit(datacl)
181 #k-means set seed to random number to make result reproducible
182 set.seed(8953)
183 datacl <- scale(datacl)
184
185 library(factoextra)
186 library(cluster)
187 library(NbClust)
188 # Split the dataset into 2clusters using k-mean method
189 kmeans.result <- kmeans(datacl,2)
190
191 # Print the result
192 kmeans.result
193
194 # Visualize clusters
195 fviz_cluster(kmeans.result, data = datacl)

196 #average for each cluster
197 avrsil<-silhouette(kmeans.result$cluster,dist(datacl))
198 fviz_silhouette(avrsil)
199
200 # split the dataset into 4 clusters using k-mean method
201 kmeans.result <- kmeans(datacl,4)
202
203 # Print the result
204 kmeans.result
205
206 # Visualize clusters
207 datacl<-na.omit(datacl))
208 fviz_cluster(kmeans.result, data = datacl)
209
210 #average for each cluster
211 avrsil<-silhouette(kmeans.result$cluster,dist(datacl))
212 fviz_silhouette(avrsil)
213
214
215 # Split the dataset into 6 clusters using k-mean method
216 kmeans.result <- kmeans(datacl,6)
217
218 # Print the result
219 kmeans.result
220
221 # Visualize clusters
222 datacl<-na.omit(datacl))
223 fviz_cluster(kmeans.result, data = datacl)
224
225 #average for each cluster
226 avrsil<-silhouette(kmeans.result$cluster,dist(datacl))
227 fviz_silhouette(avrsil)
228
229 #Get the optimal number of clusters
230
231 # (a) fviz_nbclust using silhouette method
232 fviz_nbclust(datacl,kmeans,method = "silhouette")+
233 labs(subtitle = "silhouette method")
234
235 # (b) NbClust validation
236 fres.nbclust <- NbClust(datacl, distance = "euclidean", min.nc = 2, max.nc = 10, method = "kmeans",index = "all")
237 #findings
238 #install.packages("GGally")
239 #install.packages("plotly")
240 library(GGally)
241 library(plotly)
242 dataset1<-as.data.frame(datacl)
243 dataset1$cluster<-as.factor(kmeans.result$cluster)
244 p<-ggparcoord(data=dataset1, columns=c(1,2,3,4,5,6,7,8,9),groupColumn = "cluster", scale = "std")
245 labs(x="d",y=ggplotly(p))
246
```

Size 1

Size 2

size 3

Evaluation

9 References

Kaggle website <https://www.kaggle.com/bryanb/applicants-for-a-gold-digger-position>

Use IEEE format <https://libraryguides.vu.edu.au/ieeereferencing/gettingstarted>

You can use any resource for reference generating such as <https://www.citethisforme.com/citation-generator/ieee>

<https://towardsdatascience.com/clustering-analysis-in-r-using-k-means-73eca4fb7967>

10 Tasks Distribution

ID	Name	Responsibilities
439201336	Seba alahmadi	Data mining technique Clustering Report
441201306	Noura Alsultan	Data mining technique Classification Report
441201159	Nouf alfulaij	Data mining technique Clustering Report
437201303	Reem Almutairi	Data mining technique Classification Report