



**UNIVERSITAT POLITÈCNICA
DE CATALUNYA
BARCELONATECH**

Final Project

Ettore Falde, Samoussa Fofana, Federico Basaglia

11/29/2021

Contents

1	Introduction	4
2	Setup the software	4
3	Importing Data	5
3.1	Variables analysis	5
4	Cleaning Data	6
4.1	Names	6
4.2	Dimensions	6
4.3	Head and Tail	6
4.4	Removing	7
4.5	Checking n's	7
4.6	Dropping na	7
5	Analysis	7
5.1	Introduction Analysis	8
5.1.1	Plot job role frequency	8
5.2	Gender Analysis	9
5.2.1	Monthly Income	9
5.2.2	Department	13
5.2.3	Job Role	16
5.2.4	Seniority	18
5.2.5	Education field	23
5.2.6	Regression	28
5.3	Attrition Analysis	31
5.3.1	Gender vs Attrition	35
5.3.2	Marital status vs Attrition	36
5.3.3	Regression	42
5.4	Satisfaction analysis	45
5.4.1	1-Job_satisfaction	45
5.4.2	Differences between the two cities in term of job satisfaction	46
6	Adding dataset	60
6.1	UK Comparation	61
6.1.1	Carrer satisfaction	62
6.2	Barcelona Comparation	63
6.2.1	Introduction Analysis	63

7	Conclusions	66
8	Referemces	68

1 Introduction

In this report we consider that the new CEO of a specific IT company has contacted us because she wants us to **analyze the current Human Resources status** of the company. She has just sent a data set with all available employee information. As we can see, the **company has two locations**: the first one in **London**, and the second one in **Barcelona**.

The new CEO is concerned about several issues. She truly believes in gender equality in organizations as it implies a signal to society. On the other hand, she is concerned that the offices in Barcelona do not follow a similar structure to the one in London. **In her opinion, the structure of the Barcelona offices should tend towards the London structure.** In her meeting with us, she also told us that she would like to know the attitudes (e.g., satisfaction) of the employees across the different departments and if anything could be done to improve them. Finally, she commented that she is very concerned about the company's succession strategy and in particular some positions in certain departments.

Let's consider that the new CEO of a specific IT company has contacted us because she wants us to analyze the current Human Resources status of the company. She has just sent a data set with all available employee information. This information is in the attached data set. As we can see, the company has two locations: the first one in London, and the second one in Barcelona.

The new CEO is concerned about several issues. She truly believes in gender equality in organizations as it implies a signal to society. On the other hand, she is concerned that the offices in Barcelona do not follow a similar structure to the one in London. In her opinion, the structure of the Barcelona offices should tend towards the London structure. In her meeting with us, she also told us that she would like to know the attitudes (e.g., satisfaction) of the employees across the different departments and if anything could be done to improve them. Finally, she commented that she is very concerned about the company's succession strategy and in particular some positions in certain departments.

Based on this information, we need to carry out an exploratory data analysis and prepare a technical report (with Rmarkdown) and a technical presentation (5-10 minutes).

Note: It is highly recommended to seek external sources of information (either in dataset or report formats) for the analysis and the reporting.

Based on this information, we will to carry out an exploratory data analysis and prepare a technical report (with Rmarkdown) and a technical presentation (5-10 minutes).

2 Setup the software

The software used for the development of the study and the writing of the report is R[1]. The first step is to define the work directory and to load the libraries needed:

```
library(tidyverse)
library(ggplot2)
library(GGally)
library(gridExtra)
library(yardstick)
library(broom)
library(janitor)
library(caTools)
library(ROCR)
library(corrplot)
library(tidytext)
library(glue)
library(scales)
```

```
library(plotly)
library(patchwork)
library(skimr)
library(RColorBrewer)
```

3 Importing Data

The first step is to load the dataset in the system, and check the names of the variables.

```
source("functions_script.R")
```

```
mydb <- read.csv2("dataset.csv")
# web_db <- read.csv("WA_Fn-UseC_-HR-Employee-Attrition.csv")
# names(web_db)
```

3.1 Variables analysis

We got the dataset from the website of Atenea, it is composed by 1506 observations of 36 variables. The variables selected for this dataset are:

1. **Age:** Variable that represent the age of the employee
2. **Attrition:** variable that represent the departure of employees from the organization for any reason
3. **BusinessTravel:** Represent how often an employee travel for work purpose
4. **DailyRate:** The amount of money the employees are paid per day
5. **Department:** Department of the company at which the employee belong
6. **DistanceFromHome:** Employee home distance from the workplace
7. **Education:** Educational level of the employee (1=Below College, 2=College, 3=Bachelor, 4=Master, 5= Doctor)
8. **EducationField:** Education field of employee (Human Resources, Life Sciencies, Marketing, Medical, Technical Degree, Other)
9. **EmployeeCount:** Coolumn all equal to 1 to count the total number of employee in the data set
10. **EmployeeNumber:** unique number to identify the employee
11. **EnvironmentSatisfaction:** level of environment satisfaction (1=Low, 2=Medium, 3=High, 4=Very High)
12. **Gender:** Gender of the employee (Male, Female)
13. **HourlyRate:** The amount of money the employees are paid per hour
14. **JobInvolvement:** Level of involvement of the employee (1=Low, 2=Medium, 3=High, 4=Very High)
15. **JobLevel:** Is a category of authority in the company (1=low, 5=High)
16. **JobRole:** Represent the role cover by the employee (Sales Executive, Research Scientist, Laboratory Technician, Manufacturing Director, Healthcare Representative, Manager, Sales Representative, Research Director, Human Resources)
17. **JobSatisfaction:** Level of satisfaction of the employee (1=Low, 2=Medium, 3=High, 4=Very High)
18. **MaritalStatus:** Marital status of the employee (Divorced, Married, Single)
19. **MonthlyIncome:** Monthly income of the employee
20. **MonthlyRate:** Monthly rate of employee
21. **NumCompaniesWorked:** Number of companies for ehich the employee worked
22. **Over18:** If the age of the employee is higher than 18 (Y = yes, N = no)
23. **OverTime:** If the employee perform over time (Yes, No)
24. **PercentSalaryHike:** Represent the percentage inrease of a salary
25. **PerformanceRating:** Performance rating of the employee (1=Low, 2=Good, 3=Excellent, 4=Outstanding)

26. **RelationshipSatisfaction**: Relationship satisfaction of the employee (1=Low, 2=Medium, 3=High, 4=Very High)
27. **StandardHours**: Standard working hour per **week?** (80 for everyone)
28. **StockOptionLevel**: Stock option level
29. **TotalWorkingYears**: Total years of working
30. **TrainingTimesLastYear**: Training hours of the last year
31. **WorkLifeBalance**: the amount of time you spend doing your job compared with the amount of time you spend with your family and doing things you enjoy (1=Bad, 2=Good, 3=Better, 4=Best)
32. **YearsAtCompany**: Total years of working at the company
33. **YearsInCurrentRole**: Total years spent in the current position
34. **YearsSinceLastPromotion**: How many year ago the employee had the last promotion
35. **YearsWithCurrManager**: How many years the employee is with the actual manager
36. **City**: where the employee works (London, Barcelona)

4 Cleaning Data

4.1 Names

In this sub-point we are going to change the names of the variables in order to have all the names of the variables with the same layout.

```
mydb <- mydb %>% clean_names(., "snake")
```

4.2 Dimensions

First of all, we are going to check the actual dimension of our dataset. Hence, from the following code we can understand that there are 36 variables in total and

```
mydb %>% dim()
```

```
## [1] 1506 36
```

```
mydb %>% nrow()
```

```
## [1] 1506
```

```
mydb %>% ncol()
```

```
## [1] 36
```

4.3 Head and Tail

Here, we are going to check the first 10 elements at the beginning and at the end of the dataset. Consecutively, we are going to check the top and the bottom values of the main relevant variables to catch some errors.

```
mydb %>% head(10)
mydb %>% tail(10)
mydb <- rename(mydb, age = age)
mydb %>% arrange(desc(age)) %>% top_n(10, age)
mydb %>% arrange(age) %>% top_n(-10, age)
```

4.4 Removing

In this part of the data cleaning we are going to remove all the blank rows, the duplicates and strange values that may affect our analysis.

```
# Remove blank rows and columnsn
mydb <- mydb %>% remove_empty(c("rows", "cols"))

# Removing entries with too high and too low age
mydb <- mydb %>% filter(age <= 80 & age >= 16)
mydb <- mydb %>% filter(job_involvement <= 4)
mydb <- mydb %>% filter(num_companies_worked >= 0)
```

Therefore, as we can see, this line of code did not affected our dataset. So, this mean that there are no rows or columns that are empty.

Now, we are going to pass to the study of duplicates, by the *employee_number* variable that we suggest it is the key.

```
# Duplicates removal
mydb %>% get_dupes(employee_number)
mydb <- mydb %>% distinct(employee_number, .keep_all= TRUE)
```

4.5 Checking n's

Hence, now it is time to check the n's.

4.6 Dropping na

To conclude, the cleaning of the dataset, we are going to remove every line with at least one empty gap.

```
mydb <- mydb %>% drop_na()
```

From now on we can easily proceed with our analysis.

5 Analysis

Our analysis consists in different parts.

1. First of all, we want to describe the gender equality inside the company and understand if there are discrepancies and how we can solve those problems.
2. Secondly, we want to understand the attrition factor and have a clear comprehension of what is the structure of the offices between Barcelona and London. This would be very helpful in order to decrease the percentage of people who want to leave and how the company can improve the level of comfort of its employees.
3. Thirdly, we think is very important to describe the attitude of our employees and what is their satisfaction level. So, combining this solution with the previous question we can improve our decision and suggest to the company where it can adopt new HR methodologies.

4. Fourthly, we want to understand how the company can find its successor and what could be the best solution.
5. To conclude, there is also the good practice to import in our project new datasets that could improve the decisions that we can take.

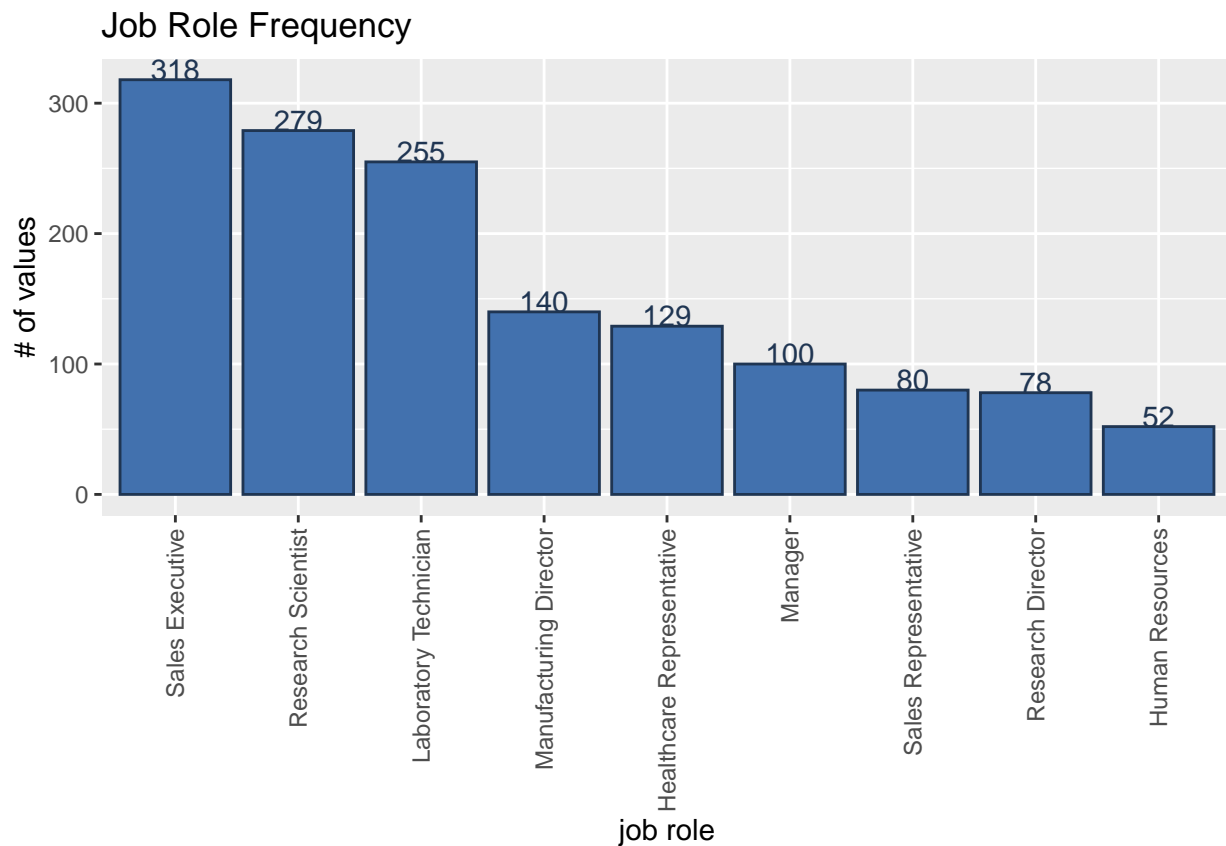
5.1 Introduction Analysis

5.1.1 Plot job role frequency

This would help to analyse the number of types of jobs and compare with types

```
k <- mydb %>% tabyl(job_role)
```

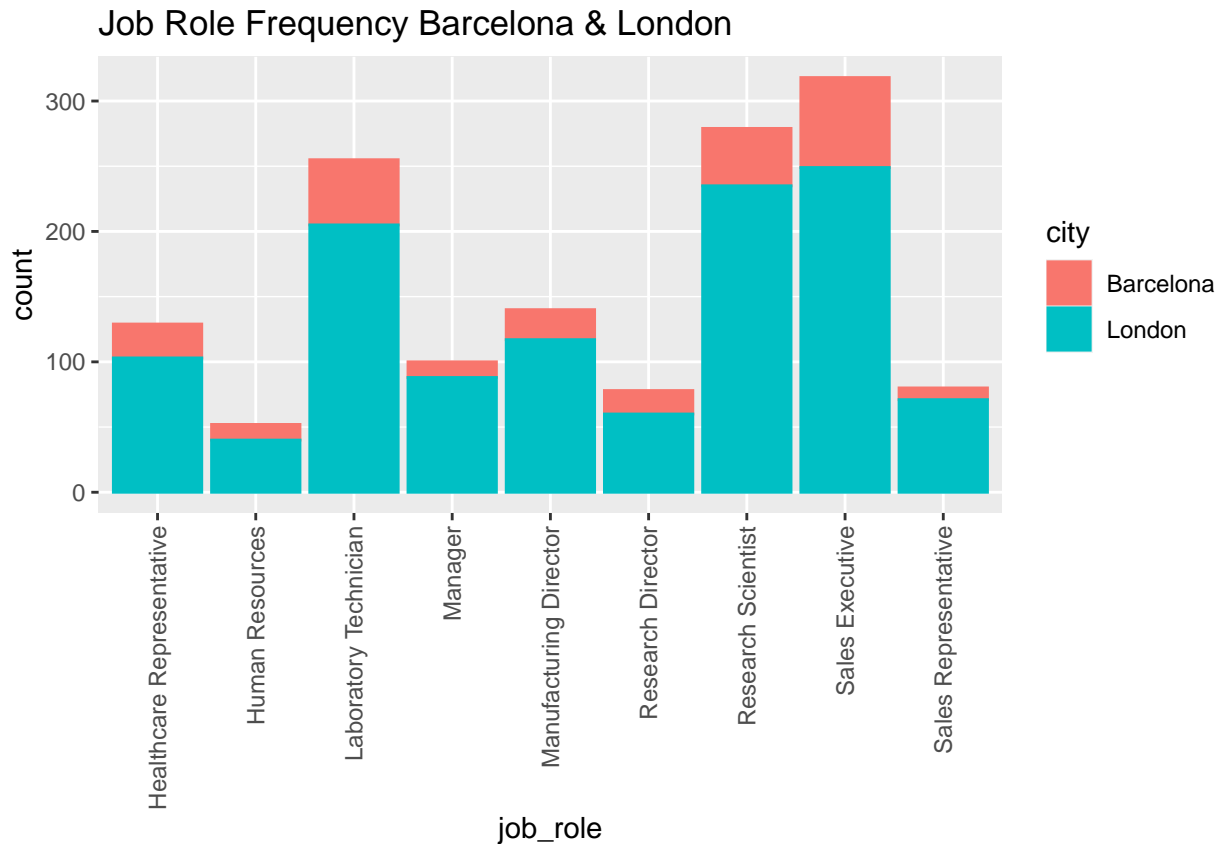
```
ggplot(k, aes(x = reorder(job_role, -n), y = n)) +
  geom_bar(stat = "identity", fill = "#4271AE", colour = "#1F3552") +
  theme(axis.text.x=element_text(angle=90,hjust=1,vjust=0.5)) +
  geom_text(aes(label = n), vjust = 0, colour = "#1F3552")+
  labs(y= "# of values", x = "job role") +
  ggtitle("Job Role Frequency")
```



```
rm("k")
```



```
ggplot(mydb, aes(x = job_role, color = city, fill = city)) +
  geom_bar(stat = "count") +
  theme(axis.text.x=element_text(angle=90,hjust=1,vjust=0.5)) +
  ggtitle("Job Role Frequency Barcelona & London")
```



5.2 Gender Analysis

First of all, we want to understand the salary that sex has. This would give us a better overview on how the salary is distributed inside the company. Moreover, we will also take into account the possible differences that we have in the Barcelona and London headwaters.

5.2.1 Monthly Income

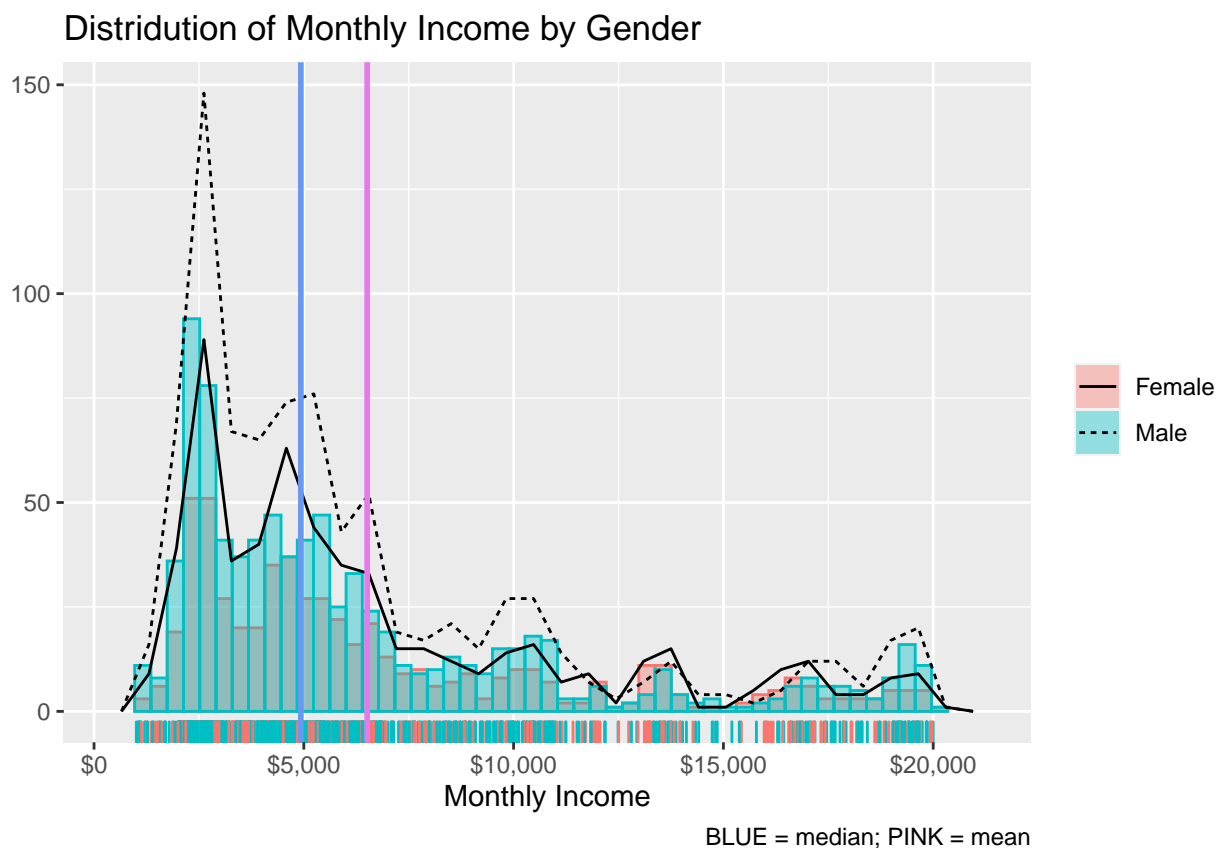
It is pretty clear that the monthly income could be one of the most important variable that will determine the differences inside a company taking into account the gender analysis.

Introduction

```
summary(mydb$monthly_income)
```

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	1009	2909	4930	6510	8386	19999

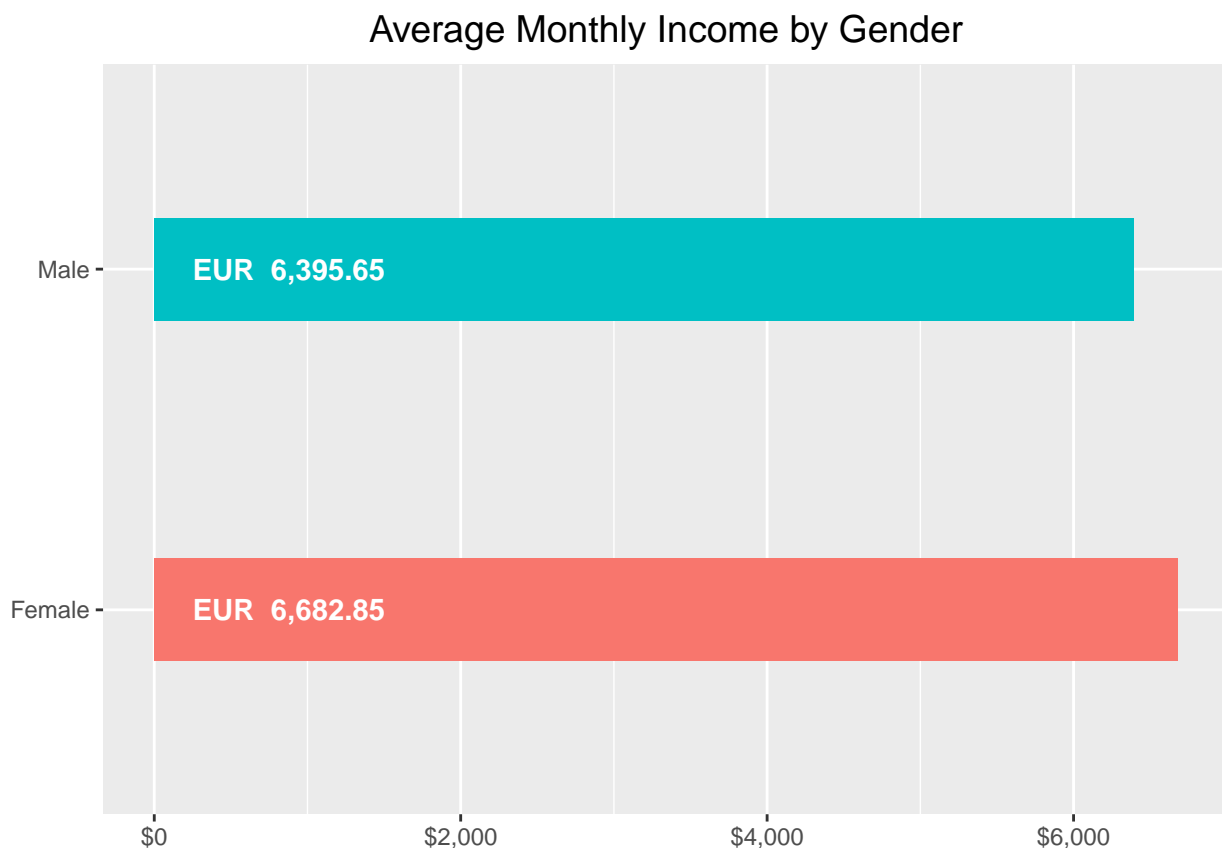
```
mydb %>% select(monthly_income, gender) %>%
  ggplot(aes(monthly_income)) +
  geom_histogram(
    aes(monthly_income, color = gender, fill = gender),
    alpha = 0.4,
    position = "identity",
    bins = 50,
  ) +
  geom_freqpoly(aes(linetype = gender), bins = 30) +
  geom_rug(aes(color = gender)) +
  scale_x_continuous(labels = label_dollar()) +
  geom_vline(aes(xintercept = mean(monthly_income)), color = "#E67AEC", size = 1) +
  geom_vline(aes(xintercept = median(monthly_income)), color = "#6899F1", size = 1) +
  guides(color = "none") +
  labs(title = "Distridution of Monthly Income by Gender",
    caption = "BLUE = median; PINK = mean",
    x = "Monthly Income",
    y = NULL,
    fill = NULL,
    linetype = NULL)
```



Hence, from this graph we can see that most of the employee have a salary lower than the average. In addition, we can also see that generally we have the dotted male line almost always above the female line, but this can be considered acceptable, due to the fact that in our dataset we have more men and women. To conclude, we can also see that the major discrepancies are given from the lower range of salary, while the higher is the salary, the lower are the differences between the female and male monthly income.

Now, we will take a closer look to the monthly income between male and female. So, we can observe that the male have a worse results. While, the female has in medium an higher salary than man.

```
mydb %>%
  select(gender, monthly_income) %>%
  group_by(gender) %>%
  summarise(avg_income = round(mean(monthly_income), 2), .groups = "drop") %>%
  ggplot(aes(x = gender, y = avg_income)) +
  geom_col(aes(fill = gender), width = 0.3, show.legend = FALSE) +
  geom_text(
    aes(
      x = gender,
      y = 0.01,
      label = dollar(avg_income, prefix = "EUR ")
    ),
    hjust = -0.2,
    size = 4,
    colour = "white",
    fontface = "bold"
  ) +
  coord_flip() +
  scale_y_continuous(labels = label_dollar()) +
  theme(plot.title = element_text(size = 14, hjust = 0.5)) +
  labs(title = "Average Monthly Income by Gender",
       x = NULL,
       y = NULL)
```



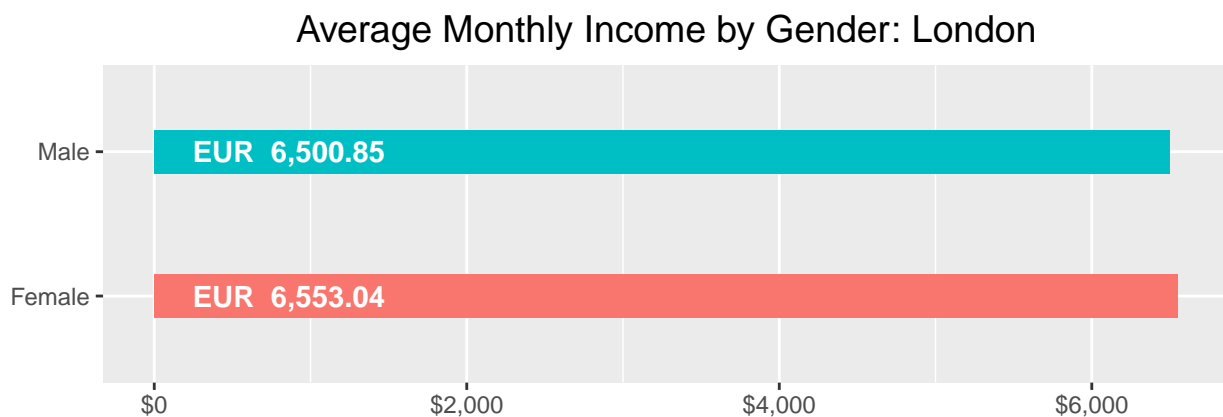
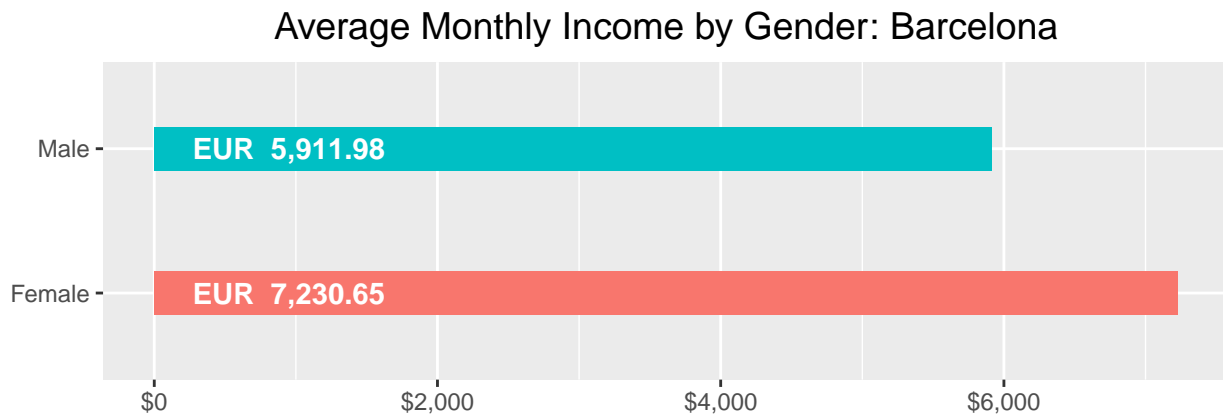
Clearly, this difference is so minimal that we need to go further in detail to clearly understand if we can adopt any strategy to balance the gender inside a company or not.

So, due to we also need to understand where is located this difference, we will filter by country.

```
# Barcelona
g1 <- mydb %>% filter(city == "Barcelona") %>%
  select(gender, monthly_income) %>%
  group_by(gender) %>%
  summarise(avg_income = round(mean(monthly_income), 2), .groups = "drop") %>%
  ggplot(aes(x = gender, y = avg_income)) +
  geom_col(aes(fill = gender), width = 0.3, show.legend = FALSE) +
  geom_text(
    aes(
      x = gender,
      y = 0.01,
      label = dollar(avg_income, prefix = "EUR ")
    ),
    hjust = -0.2,
    size = 4,
    colour = "white",
    fontface = "bold"
  ) +
  coord_flip() +
  scale_y_continuous(labels = label_dollar()) +
  theme(plot.title = element_text(size = 14, hjust = 0.5)) +
  labs(title = "Average Monthly Income by Gender: Barcelona",
       x = NULL,
       y = NULL)

# London
g2 <- mydb %>% filter(city == "London") %>%
  select(gender, monthly_income) %>%
  group_by(gender) %>%
  summarise(avg_income = round(mean(monthly_income), 2), .groups = "drop") %>%
  ggplot(aes(x = gender, y = avg_income)) +
  geom_col(aes(fill = gender), width = 0.3, show.legend = FALSE) +
  geom_text(
    aes(
      x = gender,
      y = 0.01,
      label = dollar(avg_income, prefix = "EUR ")
    ),
    hjust = -0.2,
    size = 4,
    colour = "white",
    fontface = "bold"
  ) +
  coord_flip() +
  scale_y_continuous(labels = label_dollar()) +
  theme(plot.title = element_text(size = 14, hjust = 0.5)) +
  labs(title = "Average Monthly Income by Gender: London",
       x = NULL,
       y = NULL)
```

```
grid.arrange(g1, g2, nrow = 2)
```



```
rm(g1, g2)
```

In this differentiation is definitely more clear how in Barcelona the monthly salary seems to privilege the female gender despite the male one. While, in London even if women has slightly higher average monthly salary, this is so small that can be omitted. In conclusion, we can focus to Barcelona and try to identify here the causes of this differences.

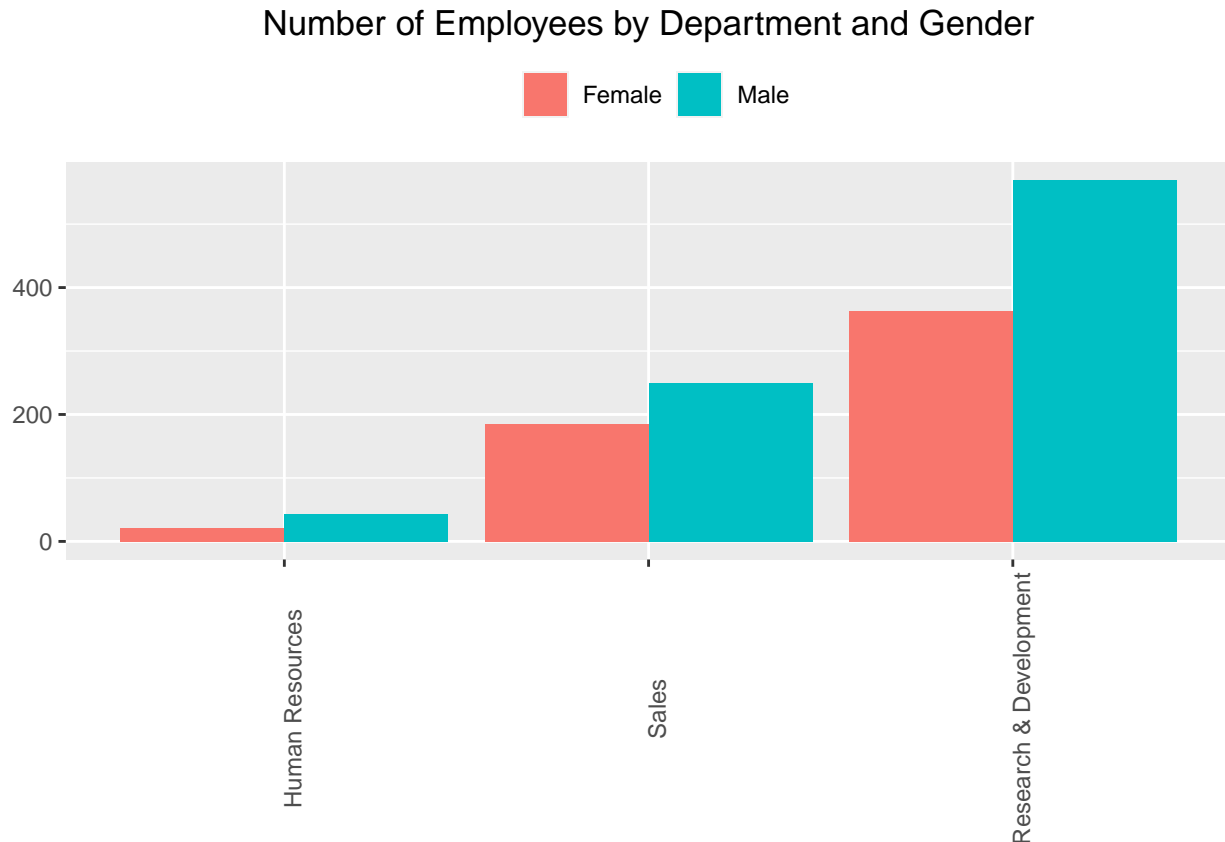
5.2.2 Department

```
mydb %>%
  group_by(department, gender) %>%
  summarise(amount = n(), .groups = "drop") %>%
  ggplot(aes(
    x = fct_reorder(department, amount),
    y = amount,
    fill = gender
  )) +
  geom_col(position = "dodge") +
  theme(axis.text.x = element_text(angle = 90),
        plot.title = element_text(hjust = 0.5),
```

```

    legend.position = "top") +
  labs(title = "Number of Employees by Department and Gender",
x = NULL,
y = NULL,
fill = NULL)

```



Therefore, in this graph we can see how inside the company we have the majority of employees in the field of research and development. Moreover, also here we can see that in general the male sex has an higher frequency in each department compared to female one.

Also here we found pretty important to stat that the differences between Barcelona and London to decide ultimately where to take decisions.

```

# Barcelona
g1 <- mydb %>% filter(city == "Barcelona") %>%
  group_by(department, gender) %>%
  summarise(amount = n(), .groups = "drop") %>%
  ggplot(aes(
    x = fct_reorder(department, amount),
    y = amount,
    fill = gender
  )) +
  geom_col(position = "dodge") +
  theme(axis.text.x = element_text(angle = 90),
    plot.title = element_text(hjust = 0.5),
    legend.position = "top") +
  labs(title = "Barcelona",

```

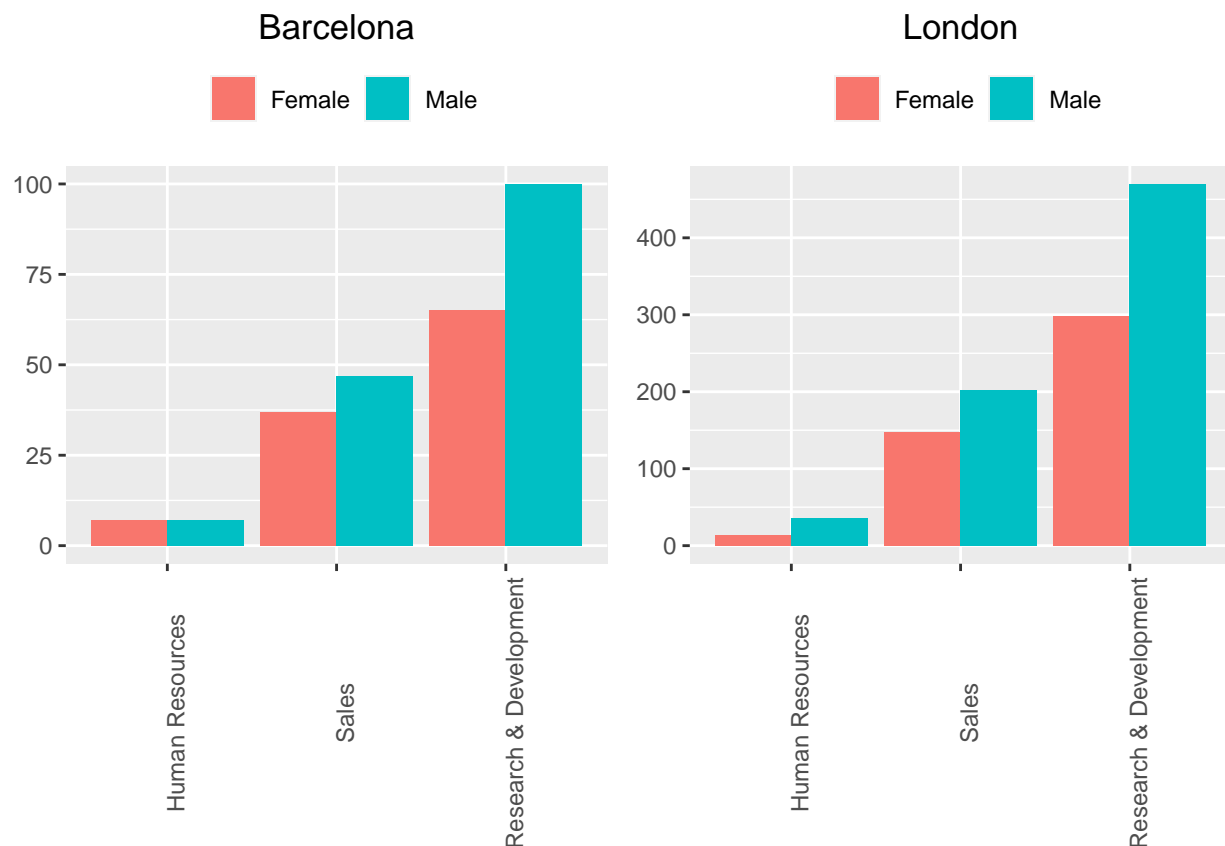
```

x = NULL,
y = NULL,
fill = NULL)

# London
g2 <- mydb %>% filter(city == "London") %>%
  group_by(department, gender) %>%
  summarise(amount = n(), .groups = "drop") %>%
  ggplot(aes(
    x = fct_reorder(department, amount),
    y = amount,
    fill = gender
  )) +
  geom_col(position = "dodge") +
  theme(axis.text.x = element_text(angle = 90),
        plot.title = element_text(hjust = 0.5),
        legend.position = "top") +
  labs(title = "London",
x = NULL,
y = NULL,
fill = NULL)

grid.arrange(g1, g2, nrow = 1)

```



```
rm(g1, g2)
```

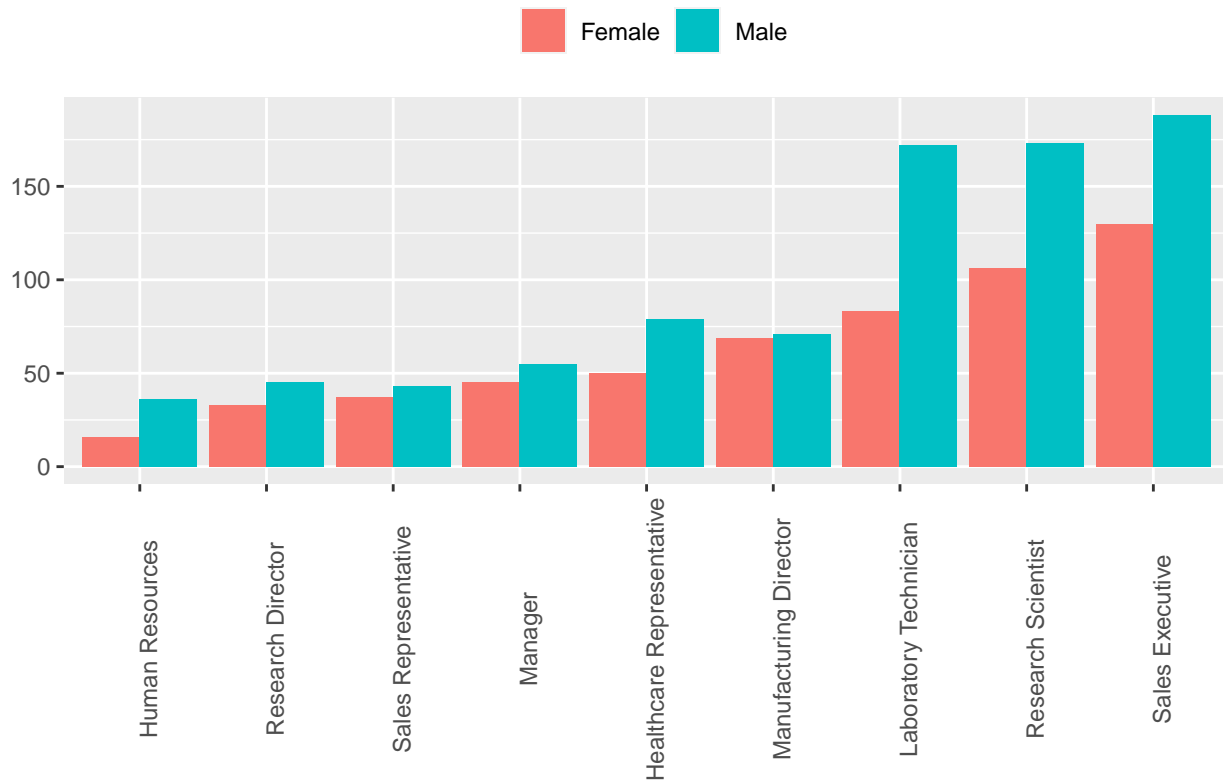
Now, we can start seeing that the major differences about the number of employees and their distribution is mainly in London, while in Barcelona for instance, In the HR department we have almost the same number of employee per gender. This balance inside the company is a good sign of gender equality, even if there are some department, such as sales and R&D where the male sex is predominant In both London and Barcelona.

5.2.3 Job Role

The same analysis as above is represented below taking into account the *job_role* variable.

```
mydb %>%
  group_by(job_role, gender) %>%
  summarise(amount = n(), .groups = "drop") %>%
  ggplot(aes(
    x = fct_reorder(job_role, amount),
    y = amount,
    fill = gender
  )) +
  geom_col(position = "dodge") +
  theme(axis.text.x = element_text(angle = 90),
        plot.title = element_text(hjust = 0.5),
        legend.position = "top") +
  labs(title = "Number of Employees by Job Role and Gender",
x = NULL,
y = NULL,
fill = NULL)
```


Number of Employees by Job Role and Gender

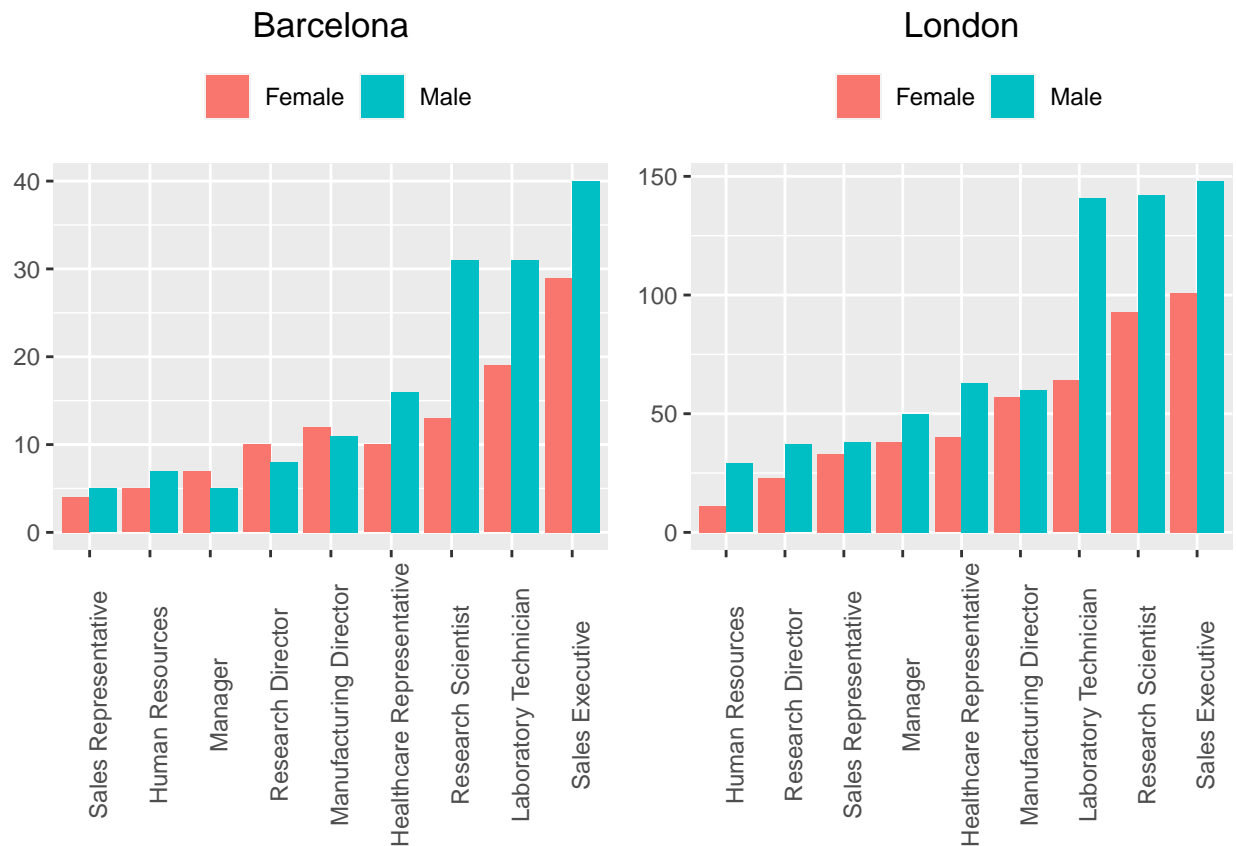


```
# Barcelona
g1 <- mydb %>% filter(city == "Barcelona") %>%
  group_by(job_role, gender) %>%
  summarise(amount = n(), .groups = "drop") %>%
  ggplot(aes(
    x = fct_reorder(job_role, amount),
    y = amount,
    fill = gender
  )) +
  geom_col(position = "dodge") +
  theme(axis.text.x = element_text(angle = 90),
        plot.title = element_text(hjust = 0.5),
        legend.position = "top") +
  labs(title = "Barcelona",
       x = NULL,
       y = NULL,
       fill = NULL)

# London
g2 <- mydb %>% filter(city == "London") %>%
  group_by(job_role, gender) %>%
  summarise(amount = n(), .groups = "drop") %>%
  ggplot(aes(
    x = fct_reorder(job_role, amount),
    y = amount,
    fill = gender
  ))
```

```
)) +
  geom_col(position = "dodge") +
  theme(axis.text.x = element_text(angle = 90),
        plot.title = element_text(hjust = 0.5),
        legend.position = "top") +
  labs(title = "London",
       x = NULL,
       y = NULL,
       fill = NULL)

grid.arrange(g1, g2, nrow = 1)
```



```
rm(g1, g2)
```

This two graphs seems to be relevant in some job roles. In fact, as we can see in Barcelona, women play an essential role in the two Director job roles.

5.2.4 Seniority

Taking in consideration both company location In this graphs in order to show that up in a better way, we decided to remove the confidence interval and the points.

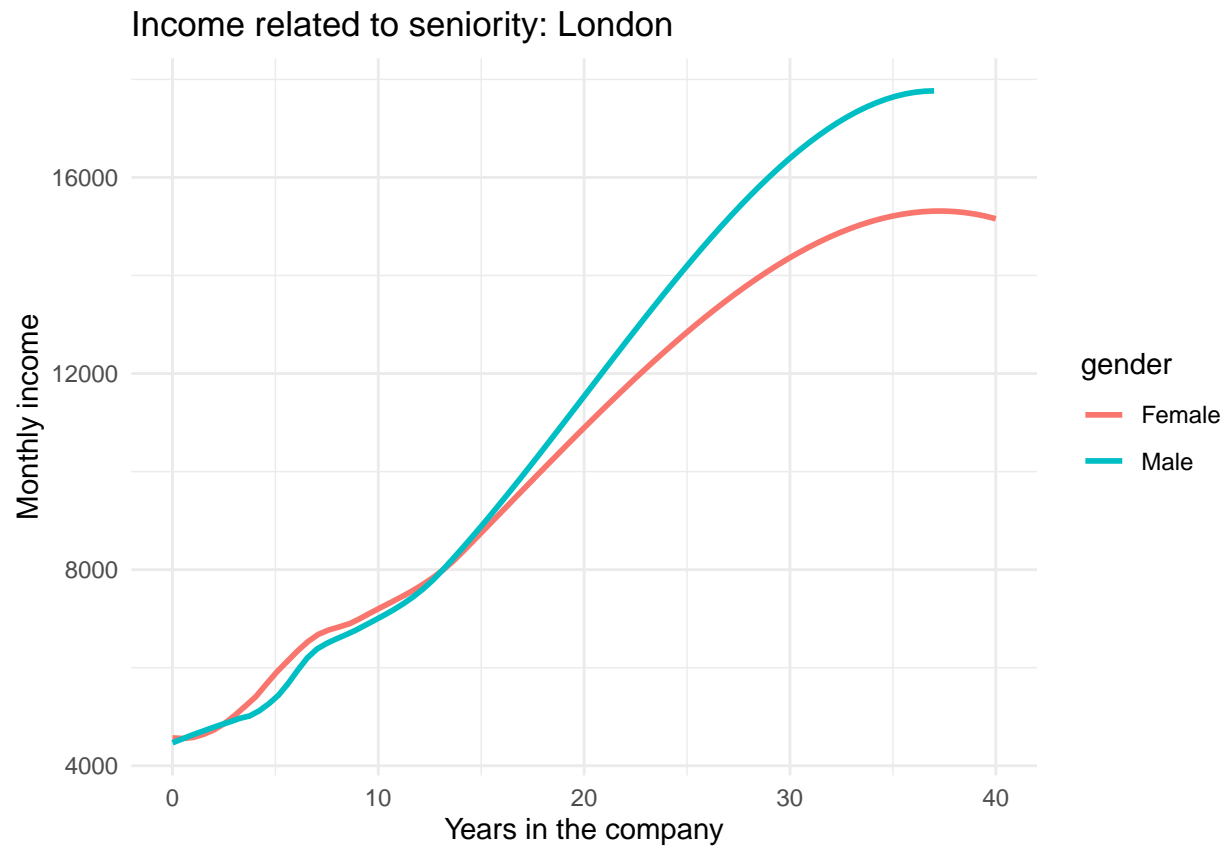
```
ggplot(mydb, aes(x=years_at_company, y=monthly_income, colour=gender))+
  geom_smooth(method = 'loess', formula = y ~ x, se=F) +
  labs(title = "Income related to seniority to all company location",
       x = "Years in the company",
       y = "Monthly income") +
  theme_minimal()
```



So, as we can see here, generally speaking the female gender has an higher salary compare to the male gender. This is true till approximately the 20th year of seniority inside the company where the male gender will overcome the female monthly income in a dramatic way.

Taking in consideration London

```
ggplot(mydb%>%filter(city=="London"),
       aes(x=years_at_company, y=monthly_income, colour=gender)) +
  geom_smooth(method = 'loess', formula = y ~ x, se=F) + labs(
    title = "Income related to seniority: London",
    x = "Years in the company",
    y = "Monthly income"
  ) +
  theme_minimal()
```



In the specific case of London we can stat that generally speaking the monthly income is well balanced between male and female work positions, but around the 15th year of seniority the male gender will have a boost in therm of monthly income.

Taking in consideration Barcelona

```
ggplot(mydb%>%filter(city=="Barcelona"),
  aes(x=years_at_company, y=monthly_income, colour=gender)) +
  geom_smooth(method = 'loess', formula = y ~ x, se=F) + labs(
    title = "Income related to seniority: Barcelona",
    x = "Years in the company",
    y = "Monthly income"
  ) +
  theme_minimal()
```

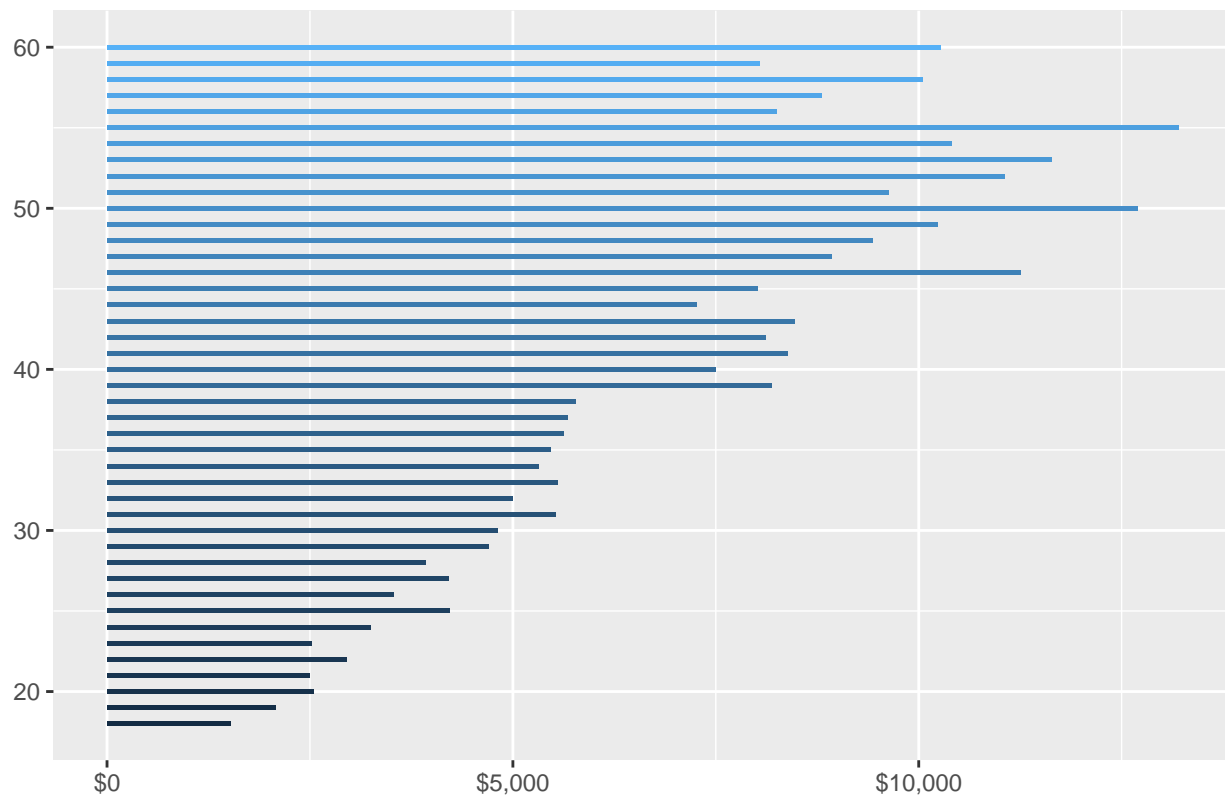


Now, in Barcelona we can see that the results we obtained are pretty different than before. In fact, the female gender has generally a higher income compared to the male gender.

Other interesting graphs we can draw are the one related to the age.

```
mydb %>%
  select(age, monthly_income) %>%
  group_by(age) %>%
  summarise(avg_income = round(mean(monthly_income), 2), .groups = "drop") %>%
  ggplot(aes(x = age, y = avg_income)) +
  geom_col(aes(fill = age), width = 0.3, show.legend = FALSE) +
  coord_flip() +
  scale_y_continuous(labels = label_dollar()) +
  theme(plot.title = element_text(size = 14, hjust = 0.5)) +
  labs(title = "Average Monthly Income by Gender",
       x = NULL,
       y = NULL)
```

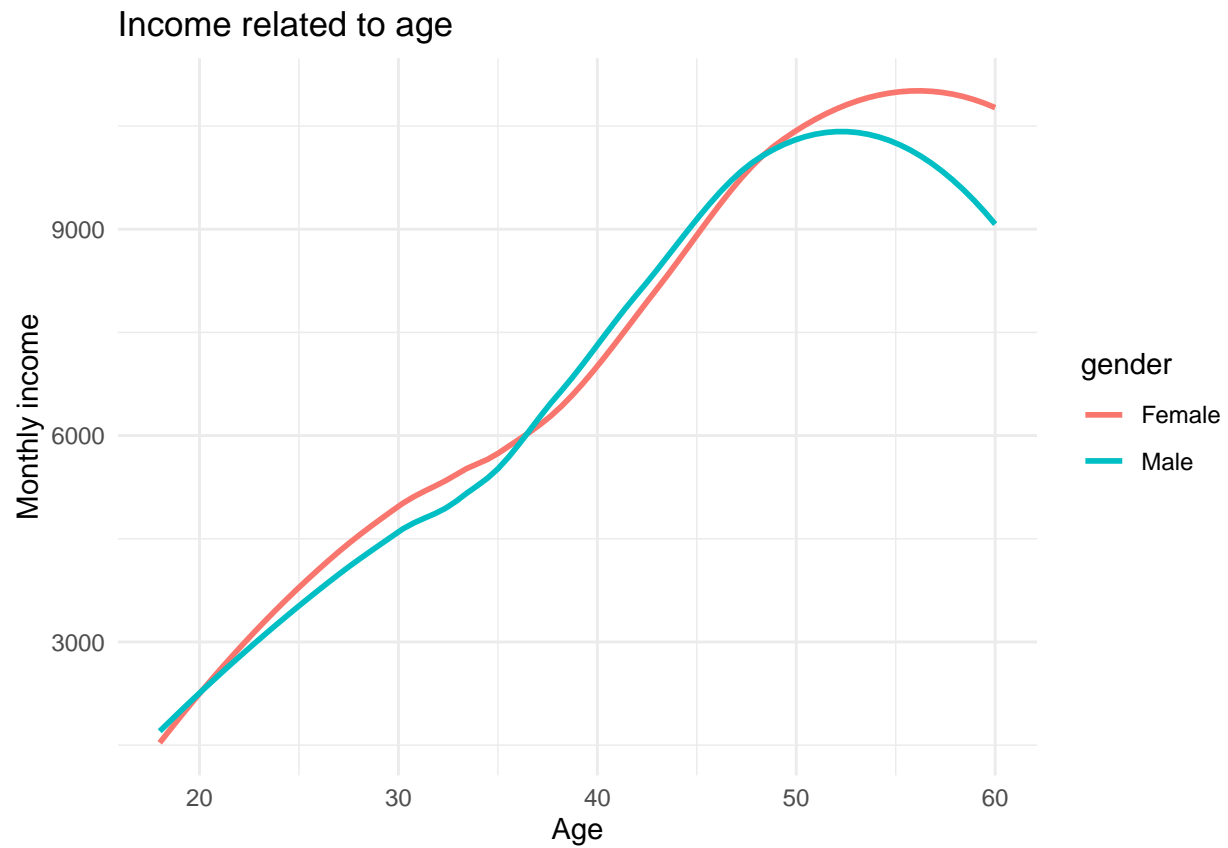
Average Monthly Income by Gender



```

ggplot(mydb,
  aes(x=age, y=monthly_income, colour=gender)) +
  geom_smooth(method = 'loess', formula = y ~ x, se=F) + labs(
    title = "Income related to age",
    x = "Age",
    y = "Monthly income"
  ) +
  theme_minimal()

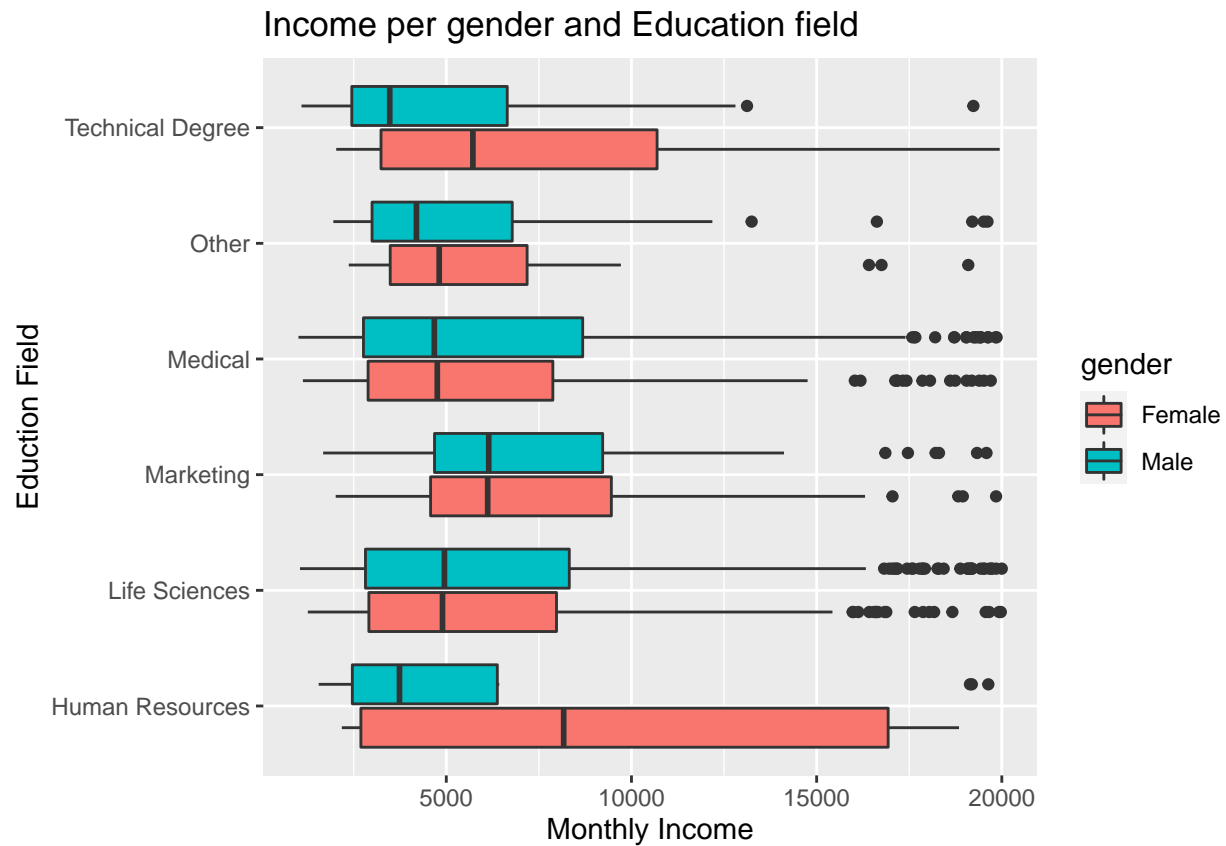
```



5.2.5 Education field

Performance rating education field.

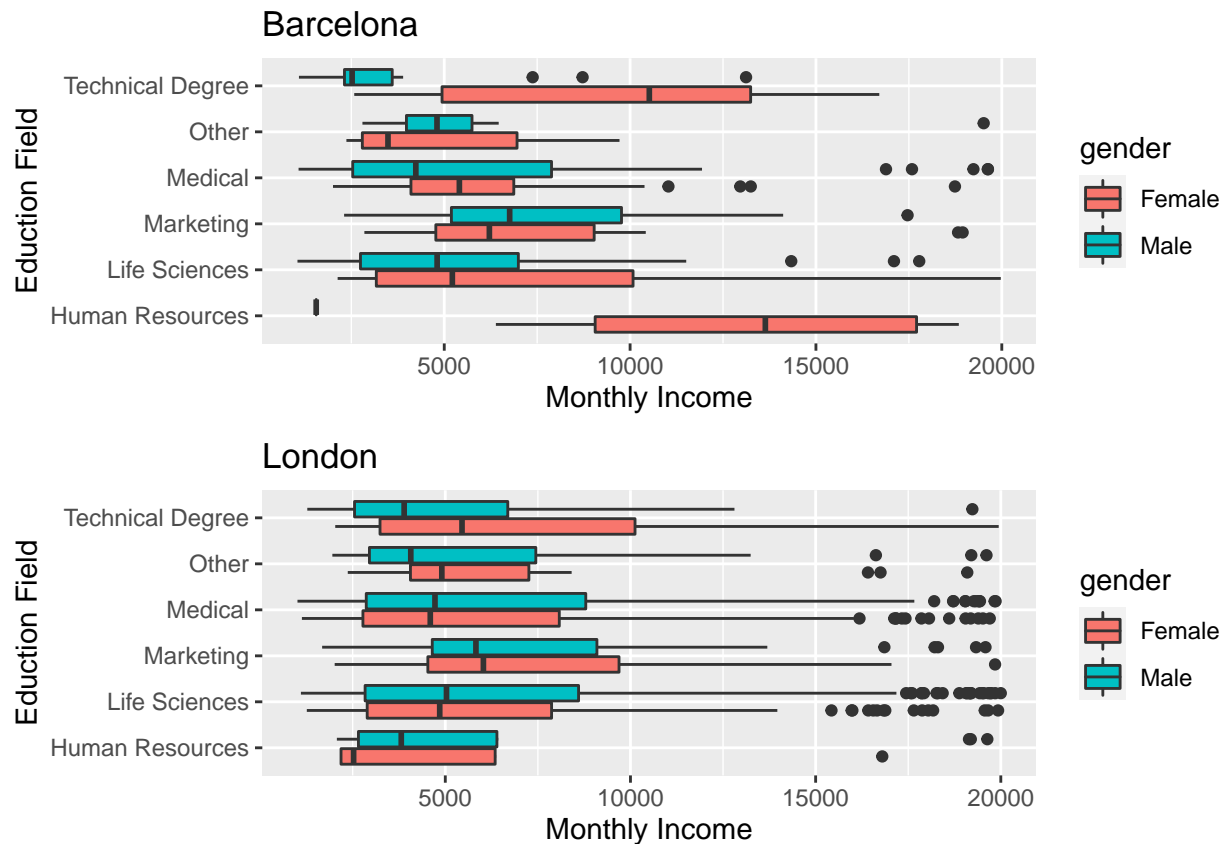
```
ggplot(mydb, aes(y= education_field, x = monthly_income, fill = gender))+  
  labs(y = "Education Field", x = "Monthly Income", title = "Income per gender and Education field") +  
  geom_boxplot()
```



```
# Barcelona
g1 <- ggplot(mydb%>%filter(city == "Barcelona"), aes(y= education_field, x = monthly_income, fill = gender)) +
  labs(y = "Education Field", x = "Monthly Income", title = "Barcelona") +
  geom_boxplot()

# London
g2 <- ggplot(mydb%>%filter(city == "London"), aes(y= education_field, x = monthly_income, fill = gender)) +
  labs(y = "Education Field", x = "Monthly Income", title = "London") +
  geom_boxplot()

grid.arrange(g1, g2, nrow = 2)
```

```
rm(g1, g2)
```

So, in the first graph we can see that there are some differences. Hence we decided to locate them and understood that the major problems are in Barcelona. In fact, in the HR sector, women have an higher salary than men, also for technical degree we can see how the female gender seems to have a clear advantage on that.

Clearly, per education field we can see also that in some cases also the male gender seems to have a slightly higher salary, but the main focus are those from which the discrepancy are very high, such as HR in Barcelona and Technical Degree, always in Barcelona.

Moreover, we want also to understand if for the same job role the gender has a crucial role in the determination of the monthly income of each employee.

```
ggplot(mydb, aes(x = monthly_income, y = job_role, fill = gender)) + geom_boxplot() +
  theme(axis.text.x = element_text(angle = 90),
        plot.title = element_text(hjust = 0.5),
        legend.position = "top") +
  labs(title = "Number of Employees by Job Role and Gender",
       x = "monthly income",
       y = "job role",
       fill = NULL)
```

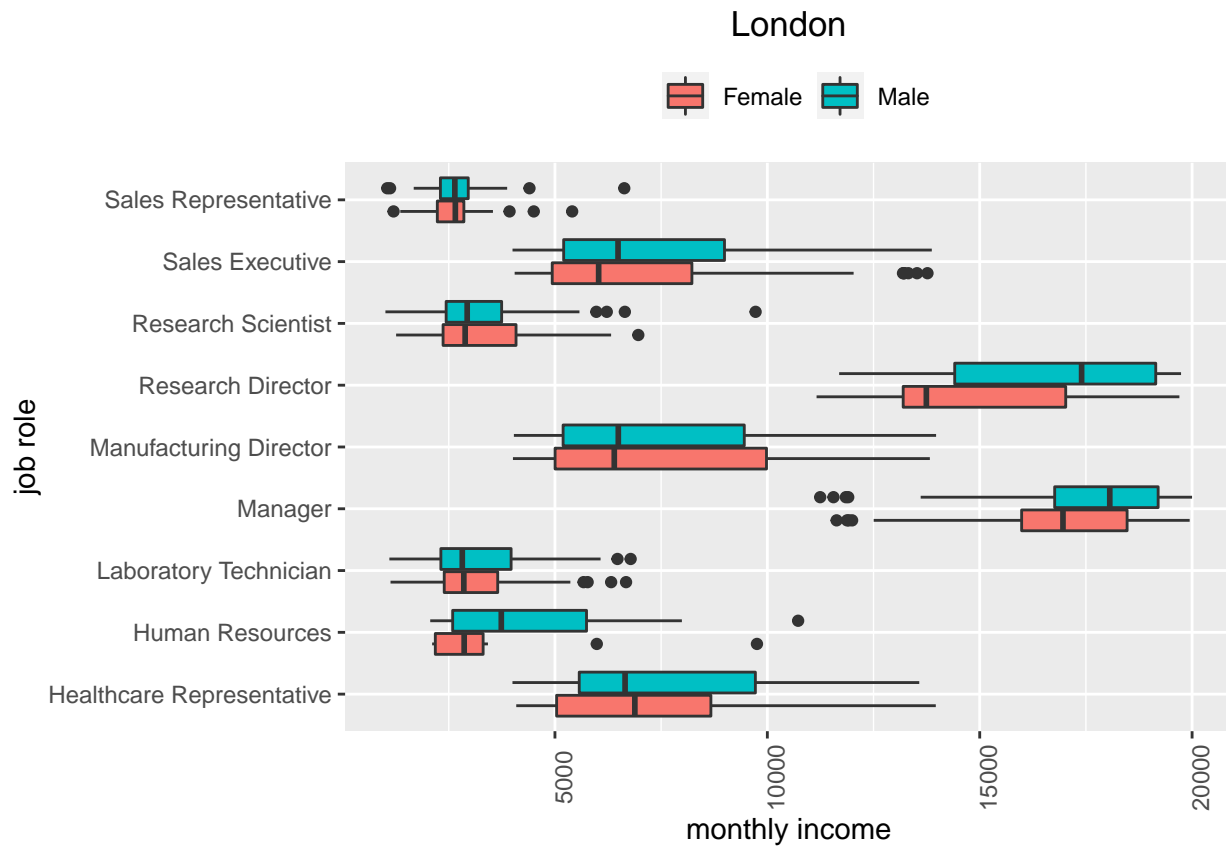
Number of Employees by Job Role and Gender



```
# Barcelona
ggplot(mydb%>%filter(city == "Barcelona"), aes(x = monthly_income, y = job_role, fill = gender)) + geom_boxplot()
  theme(axis.text.x = element_text(angle = 90),
        plot.title = element_text(hjust = 0.5),
        legend.position = "top") +
  labs(title = "Barcelona",
       x = "monthly income",
       y = "job role",
       fill = NULL)
```



```
# London
ggplot(mydb%>%filter(city == "London"), aes(x = monthly_income, y = job_role, fill = gender)) + geom_boxplot()
  theme(axis.text.x = element_text(angle = 90),
        plot.title = element_text(hjust = 0.5),
        legend.position = "top") +
  labs(title = "London",
       x = "monthly income",
       y = "job role",
       fill = NULL)
```



Hence, we can understand that also in this case the majority of the problem related to gender are located in Barcelona.

5.2.6 Regression

To conclude, we wanted also to see if the gender is a relevant variable that affect *monthly_income* or other variables relevant for our study.

```
mod1 <- lm(monthly_income ~ ., data = mydb)
summary(mod1)
```

```
##
## Call:
## lm(formula = monthly_income ~ ., data = mydb)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3908.5   -708.4    -2.6    671.1   4386.4
##
## Coefficients: (2 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3.470e+02  7.669e+02   0.453  0.650964
## age          -5.641e+00  4.709e+00  -1.198  0.231111
## attritionYes  3.134e+01  9.423e+01   0.333  0.739467
## business_travelTravel_Frequently 1.522e+02  1.179e+02   1.291  0.196892
## business_travelTravel_Rarely    1.809e+02  1.007e+02   1.797  0.072541 .
```

```
## daily_rate 8.238e-02 7.538e-02 1.093 0.274661
## departmentResearch & Development 3.289e+02 3.880e+02 0.848 0.396800
## departmentSales 1.243e+02 4.096e+02 0.303 0.761617
## distance_from_home -4.345e+00 3.731e+00 -1.165 0.244375
## education -1.657e+01 3.035e+01 -0.546 0.585303
## education_fieldLife Sciences -1.192e+02 2.915e+02 -0.409 0.682715
## education_fieldMarketing -2.094e+01 3.114e+02 -0.067 0.946408
## education_fieldMedical -1.418e+02 2.928e+02 -0.484 0.628357
## education_fieldOther -2.207e+02 3.142e+02 -0.702 0.482566
## education_fieldTechnical Degree -4.758e+01 3.056e+02 -0.156 0.876305
## employee_count NA NA NA NA
## employee_number 7.576e-02 5.037e-02 1.504 0.132803
## environment_satisfaction -3.513e+00 2.804e+01 -0.125 0.900322
## genderMale 1.033e+02 6.209e+01 1.664 0.096398 .
## hourly_rate 1.065e+00 1.492e+00 0.714 0.475601
## job_involvement -9.031e+01 4.298e+01 -2.101 0.035832 *
## job_level 2.758e+03 6.892e+01 40.018 < 2e-16 ***
## job_roleHuman Resources -4.542e+01 4.133e+02 -0.110 0.912508
## job_roleLaboratory Technician -5.865e+02 1.413e+02 -4.150 3.53e-05 ***
## job_roleManager 4.247e+03 2.105e+02 20.177 < 2e-16 ***
## job_roleManufacturing Director -4.501e+01 1.394e+02 -0.323 0.746829
## job_roleResearch Director 4.072e+03 1.846e+02 22.064 < 2e-16 ***
## job_roleResearch Scientist -5.202e+02 1.399e+02 -3.718 0.000209 ***
## job_roleSales Executive 1.204e+02 2.746e+02 0.439 0.661048
## job_roleSales Representative -5.310e+02 3.059e+02 -1.736 0.082827 .
## job_satisfaction 8.567e-01 2.776e+01 0.031 0.975385
## marital_statusMarried 1.722e+01 8.113e+01 0.212 0.831897
## marital_statusSingle -3.919e+01 1.119e+02 -0.350 0.726291
## monthly_rate -4.583e-03 4.246e-03 -1.079 0.280650
## num_companies_worked 1.128e+01 1.361e+01 0.829 0.407411
## over18Y -2.533e+01 4.363e+02 -0.058 0.953715
## over_timeYes 7.032e+01 7.004e+01 1.004 0.315615
## percent_salary_hike 1.746e+01 1.306e+01 1.337 0.181600
## performance_rating -1.627e+02 1.316e+02 -1.237 0.216350
## relationship_satisfaction 2.005e+01 2.817e+01 0.712 0.476750
## standard_hours NA NA NA NA
## stock_option_level -4.319e+01 4.853e+01 -0.890 0.373648
## total_working_years 4.752e+01 8.490e+00 5.598 2.62e-08 ***
## training_times_last_year -1.715e+01 2.369e+01 -0.724 0.469296
## work_life_balance -1.562e+01 4.295e+01 -0.364 0.716218
## years_at_company 3.395e+00 1.057e+01 0.321 0.748023
## years_in_current_role 5.206e+00 1.356e+01 0.384 0.701107
## years_since_last_promotion 2.366e+01 1.219e+01 1.940 0.052544 .
## years_with_curr_manager -3.138e+01 1.416e+01 -2.216 0.026842 *
## cityLondon -1.143e+01 7.823e+01 -0.146 0.883837
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1129 on 1383 degrees of freedom
## Multiple R-squared: 0.9445, Adjusted R-squared: 0.9426
## F-statistic: 500.8 on 47 and 1383 DF, p-value: < 2.2e-16
```

```
mod2 <- lm(monthly_income ~ gender, data = mydb)
summary(mod2)
```

```
##
## Call:
## lm(formula = monthly_income ~ gender, data = mydb)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -5554   -3586   -1555    1924   13603
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   6682.8      197.5   33.837  <2e-16 ***
## genderMale    -287.2      254.5   -1.129    0.259
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4711 on 1429 degrees of freedom
## Multiple R-squared:  0.0008906, Adjusted R-squared:  0.0001914
## F-statistic: 1.274 on 1 and 1429 DF, p-value: 0.2593
```

```
names(mydb)
```

```
## [1] "age" "attrition"
## [3] "business_travel" "daily_rate"
## [5] "department" "distance_from_home"
## [7] "education" "education_field"
## [9] "employee_count" "employee_number"
## [11] "environment_satisfaction" "gender"
## [13] "hourly_rate" "job_involvement"
## [15] "job_level" "job_role"
## [17] "job_satisfaction" "marital_status"
## [19] "monthly_income" "monthly_rate"
## [21] "num_companies_worked" "over18"
## [23] "over_time" "percent_salary_hike"
## [25] "performance_rating" "relationship_satisfaction"
## [27] "standard_hours" "stock_option_level"
## [29] "total_working_years" "training_times_last_year"
## [31] "work_life_balance" "years_at_company"
## [33] "years_in_current_role" "years_since_last_promotion"
## [35] "years_with_curr_manager" "city"
```

```
mod3 <- lm(monthly_income ~ age + gender + business_travel + education_field + education + distance_from_home, data = mydb)
summary(mod3)
```

```
##
## Call:
## lm(formula = monthly_income ~ age + gender + business_travel +
##      education_field + education + distance_from_home + job_level +
##      job_role + marital_status + total_working_years + performance_rating +
##      years_at_company, data = mydb)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3891.7   -705.5       2.1    660.6   4258.4
```

```
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)      193.964    450.615   0.430 0.666940
## age              -4.646     4.638  -1.002 0.316737
## genderMale       104.622     61.651   1.697 0.089917 .
## business_travelTravel_Frequently 140.283    116.457   1.205 0.228563
## business_travelTravel_Rarely     166.238    100.023   1.662 0.096735 .
## education_fieldLife Sciences    -44.626    268.716  -0.166 0.868124
## education_fieldMarketing     30.145    287.329   0.105 0.916458
## education_fieldMedical    -57.443    269.850  -0.213 0.831460
## education_fieldOther    -132.605    291.814  -0.454 0.649599
## education_fieldTechnical Degree   31.189    283.180   0.110 0.912314
## education        -17.829     30.166  -0.591 0.554597
## distance_from_home    -4.104     3.703  -1.108 0.268013
## job_level          2772.607     68.512  40.469 < 2e-16 ***
## job_roleHuman Resources   -328.629    221.961  -1.481 0.138945
## job_roleLaboratory Technician -608.474    140.212  -4.340 1.53e-05 ***
## job_roleManager        4114.167    182.436  22.551 < 2e-16 ***
## job_roleManufacturing Director  -91.199    138.511  -0.658 0.510374
## job_roleResearch Director  4006.173    183.527  21.829 < 2e-16 ***
## job_roleResearch Scientist  -528.001    139.542  -3.784 0.000161 ***
## job_roleSales Executive   -105.359    126.509  -0.833 0.405089
## job_roleSales Representative -701.603    182.652  -3.841 0.000128 ***
## marital_statusMarried      48.096     77.150   0.623 0.533118
## marital_statusSingle      23.255     83.002   0.280 0.779380
## total_working_years      49.009     8.232   5.953 3.32e-09 ***
## performance_rating     -27.222     82.507  -0.330 0.741496
## years_at_company         -4.357     6.479  -0.672 0.501379
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1130 on 1405 degrees of freedom
## Multiple R-squared:  0.9435, Adjusted R-squared:  0.9425
## F-statistic: 937.9 on 25 and 1405 DF,  p-value: < 2.2e-16
```

```
rm(mod1, mod2, mod3)
```

Hence, after this tree linear regression models we finally understood that actually the gender variable isn't significative to determine the monthly income. Anyway, this doesn't means that there are no discrepancy between gender, but simply that we need to develop an internal company analysis understand how in Barcelona for some kind of roles women and men have higher salaries compared to the opposite sex.

5.3 Attrition Analysis

First of all, we want to understand how the percentage of attrition is distributed inside the company. Then, we will differentiate per country.

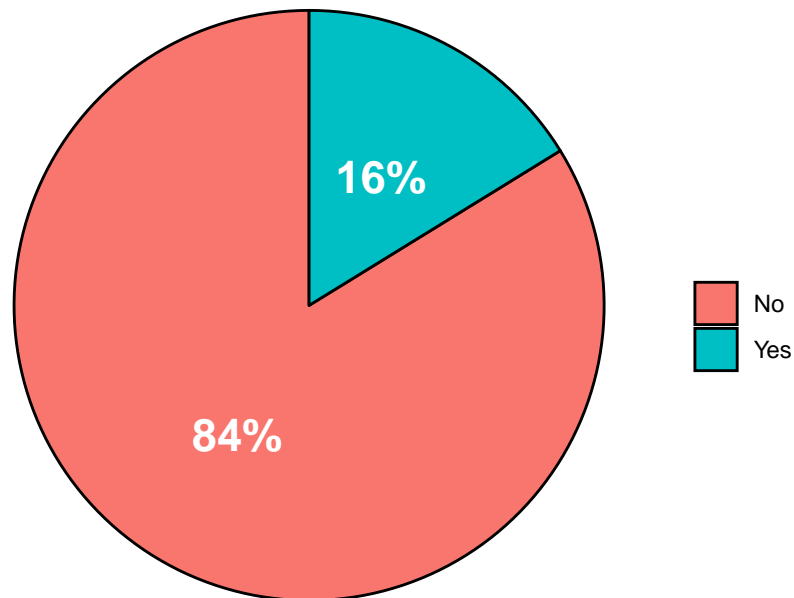
```
# Create a temporary dataset to get the percentage of employees that would leave the company
temp <- mydb %>%
  group_by(attrition) %>%
  summarize(counts = n()) %>%
  mutate(percent = percentage_func(counts)) %>%
```

```
arrange(desc(percent))

# Create a pie chart to see the percentage of people that will leave the company
pie_chart_func(dataset = temp,
               counts_var = temp$counts,
               var_interest = temp$attrition,
               title = "Are the Attrition var Balanced?",
               subtitle = "Pie Plot,percentortion of YES to NO in Attrition Var",
               caption = "UPC")
```

Are the Attrition var Balanced?

Pie Plot,percentortion of YES to NO in Attrition Var



UPC

```
# Removing the un-used variables
rm("temp")
```

```
# Barcelona
# Create a temporary dataset to get the percentage of employees that would leave the company
t1 <- mydb %>% filter(city == "Barcelona") %>%
  group_by(attrition) %>%
  summarize(counts = n()) %>%
  mutate(percent = percentage_func(counts)) %>%
  arrange(desc(percent))

g1 <- pie_chart_func(dataset = t1,
                    counts_var = t1$counts,
                    var_interest = t1$attrition,
                    title = "Barcelona",
```



```

        subtitle = "Pie Plot,percentortion of YES to NO in Attrition Var",
        caption = "")

# London
t2 <- mydb %>% filter(city == "London") %>%
  group_by(attrition) %>%
  summarize(counts = n()) %>%
  mutate(percent = percentage_func(counts)) %>%
  arrange(desc(percent))

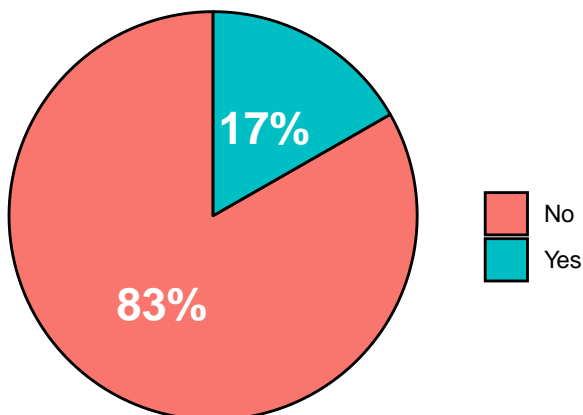
# Create a pie chart to see the percentage of people that will leave the company
g2 <- pie_chart_func(dataset = t2,
  counts_var = t2$counts,
  var_interest = t2$attrition,
  title = "London",
  subtitle = "Pie Plot,percentortion of YES to NO in Attrition Var",
  caption = "")

grid.arrange(g1, g2, nrow = 1)

```

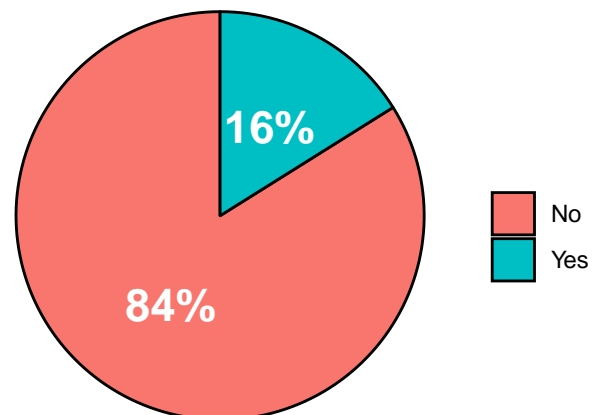
Barcelona

Pie Plot,percentortion of YES to NO in Attrition Var



London

Pie Plot,percentortion of YES to NO in Attrition Var



```

# Removing the un-used variables
rm(g1, g2, t1, t2)

```

Fortunately, as we can see the percentage of attrition between Barcelona and London is almost the same, so

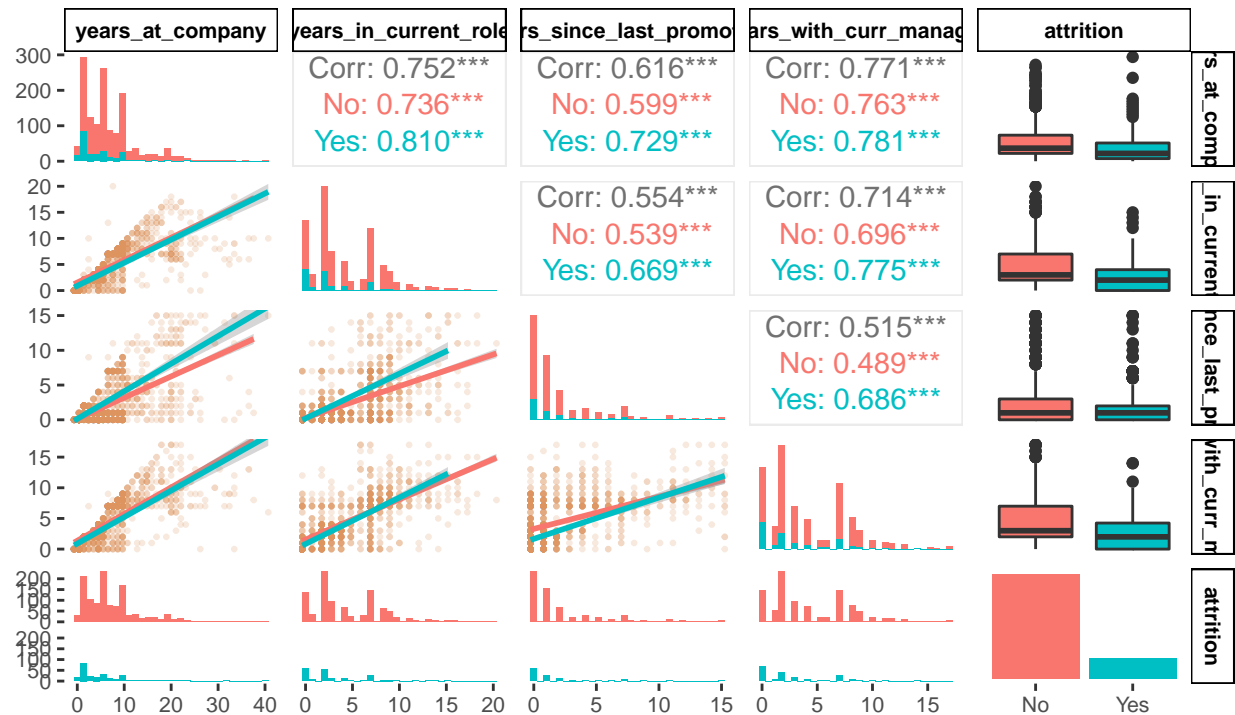
we can proceed with a generalized analysis.

```
mydb %>% select(starts_with("years"), attrition) %>%

ggpairs(
  aes(color = attrition),
  lower = list(continuous = wrap(
    "smooth",
    alpha = 0.2,
    size = 0.5,
    color = "#DE945E"
  )),
  diag = list(continuous = "barDiag"),
  upper = list(continuous = wrap("cor", size = 4))
) +
  theme(
    axis.text = element_text(size = 8),
    panel.background = element_rect(fill = "white"),
    strip.background = element_rect(fill = "white"),
    strip.background.x = element_rect(colour = "black"),
    strip.background.y = element_rect(colour = "black"),
    strip.text = element_text(color = "black", face = "bold", size = 8)
  ) +
  labs(
    title = "Pair plot by attrition Var",
    subtitle = "Pair Plot, scatter plot, Histogram and Correlation coefficient",
    caption = "",
    x = NULL,
    y = NULL
  )
```

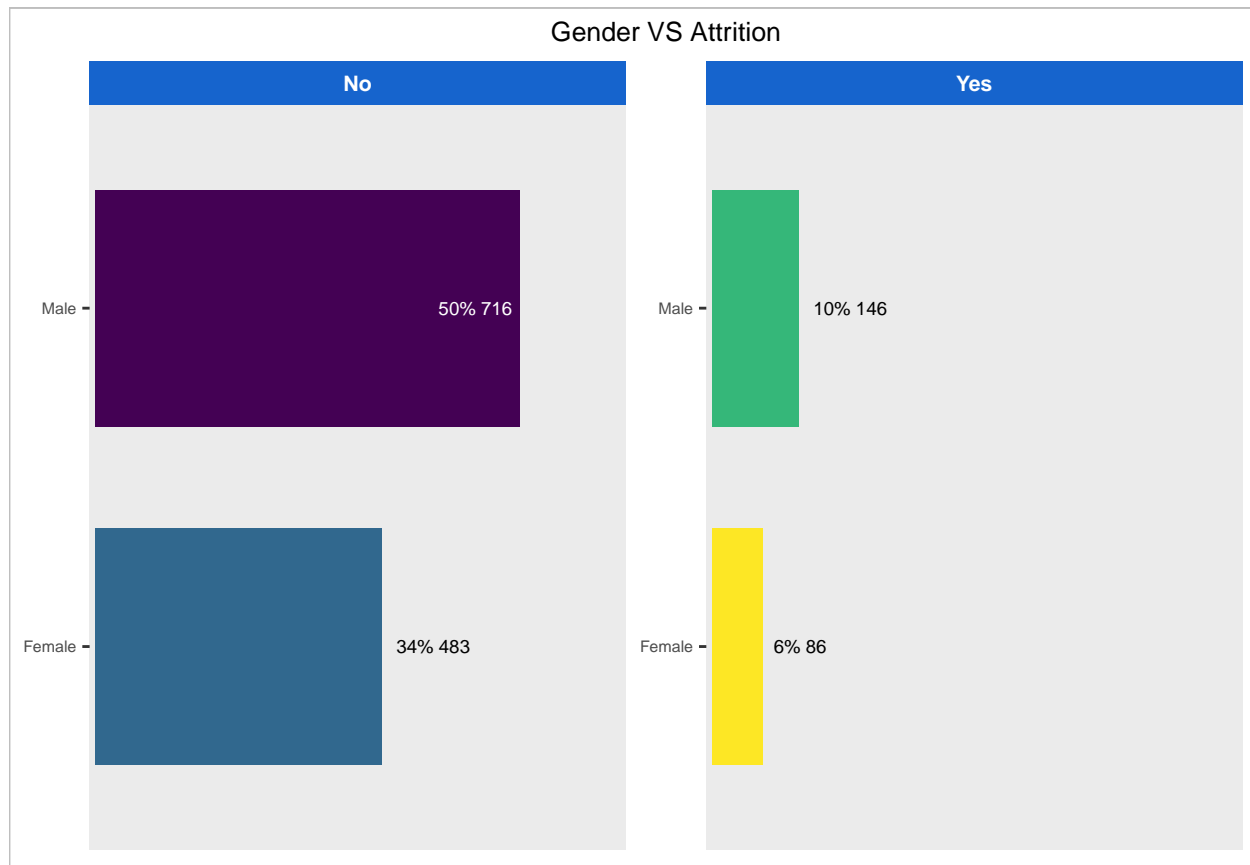
Pair plot by attrition Var

Pair Plot, scatter plot, Histogram and Correlation coefficient



5.3.1 Gender vs Attrition

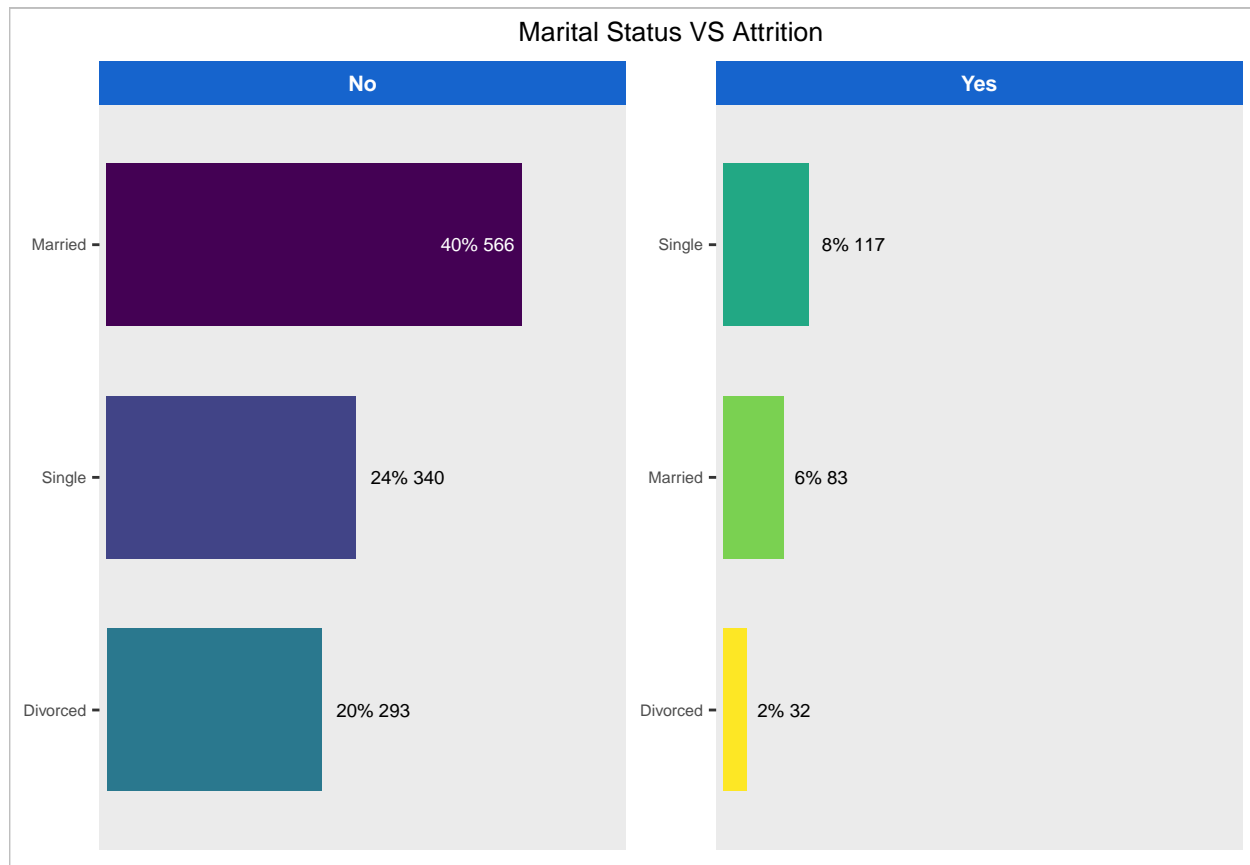
```
bar_plot_proportions(gender, attrition)
```



So, as we can see male tend to leave the company more the women does, but we also have to take into account the fact that there are more men in the company. Hence, the relation seems to be pretty balanced.

5.3.2 Marital status vs Attrition

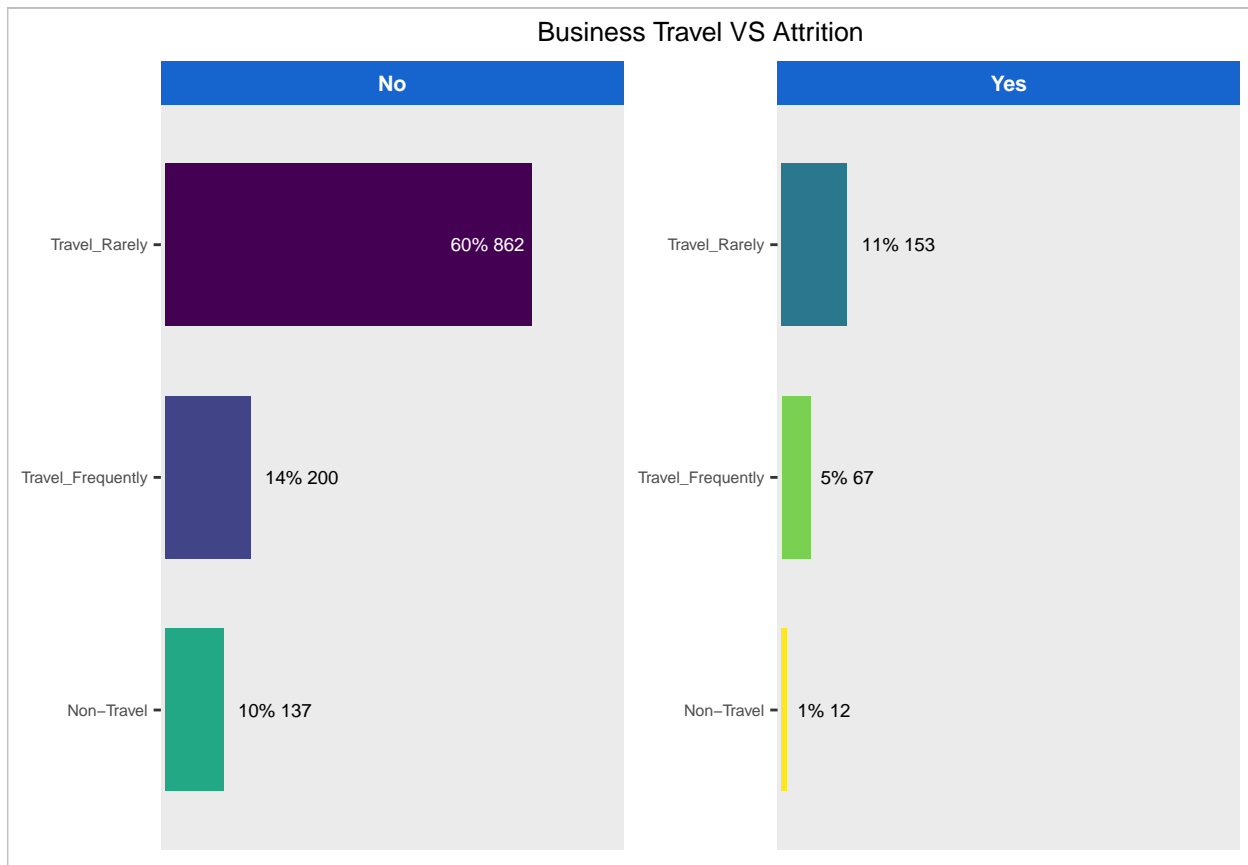
```
bar_plot_proportions(marital_status, attrition)
```



Here, we can see how single people tend to leave the company more, compared to married and divorced.

The *business_travel* variable could be a very important variable that will tell us if the travels affect the attrition variable.

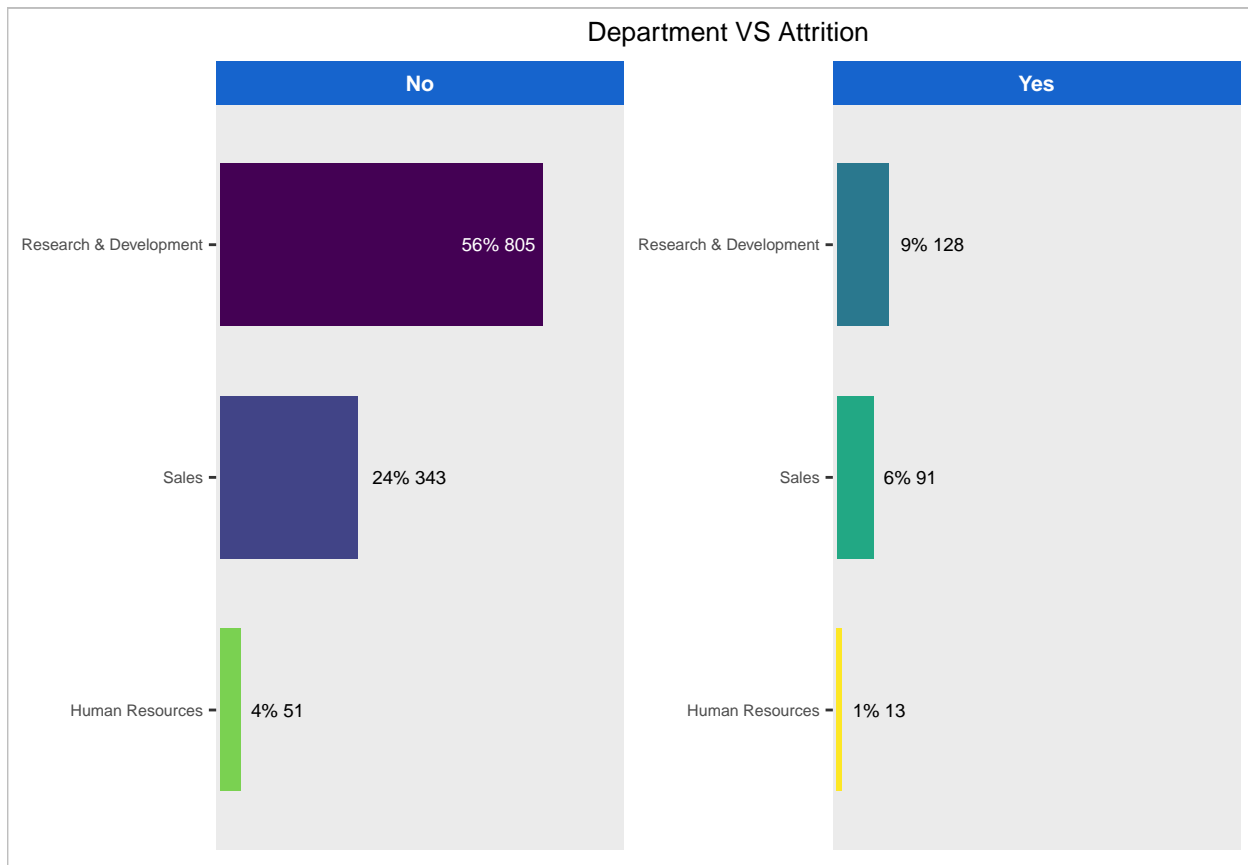
```
bar_plot_proportions(business_travel, attrition)
```



At the end, we have that the main obs that affect the attrition is the `travel_rarely`, which could be that the non routine and the change of plans for the single employee could cause an higher attrition. Hence, we suggest to clarify with higher advance if that employee have to travel or not, and finally take the new results and compute a further analysis.

The `department` variable could be a very important variable that will tell us if the department affect the attrition variable.

```
bar_plot_proportions(department, attrition)
```



Here instead, we can see how the R&D department seems the one with higher attrition. Anyway, we have also have to take into consideration that the majority of the people inside the company work in this department.

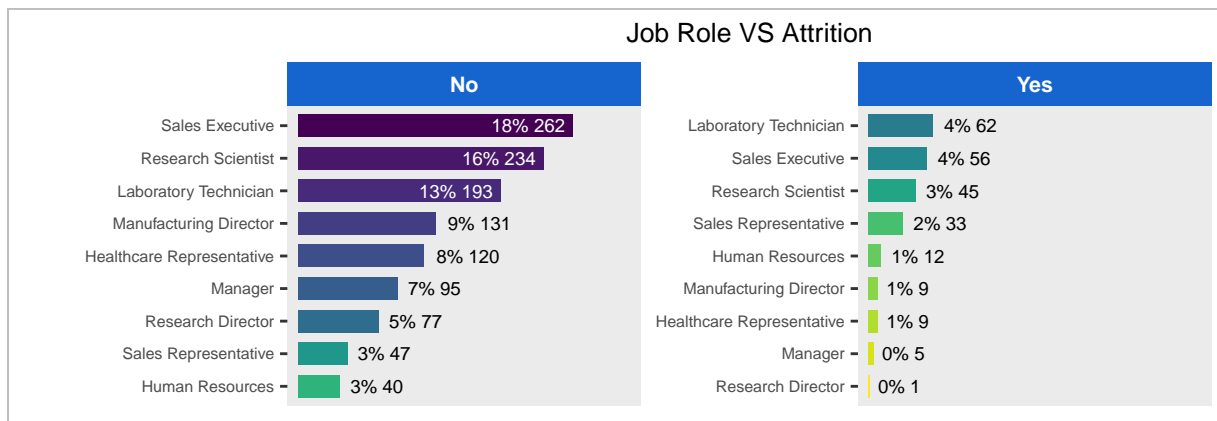
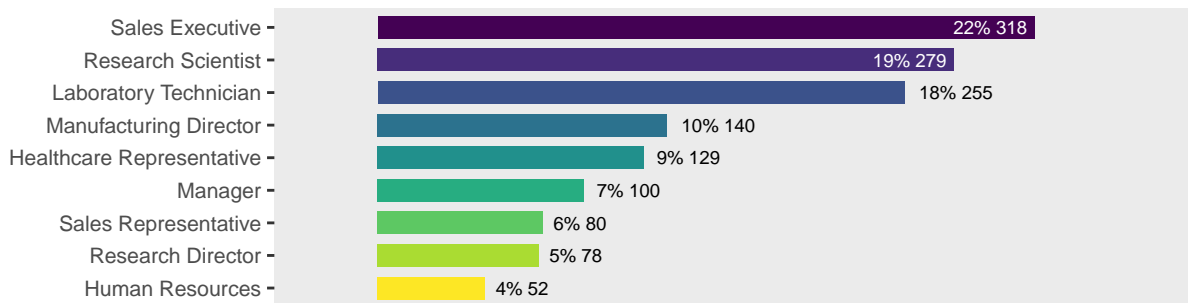
Now, we want to study also the *job_role* versus the *attrition* variable. This would tell us what are the most common job positions that will affect this variable and where we can work in order to reduce the percentage of people who leave the company in the specific job role.

```

plt_job_role <- bar_plot_proportions(job_role)
plt_job_role_att <- bar_plot_proportions(job_role, attrition)
(plt_job_role /
  plt_job_role_att) +
  plot_annotation(
    title = "Proportions of Job role VS Attrition",
    caption = ""
  ) &
  theme(plot.caption = element_text(color = "#969696", size = 7))

```

Proportions of Job role VS Attrition



In this case, Laboratory Technician, Sales Executive, Research Scientist and Sales Representative, tend to leave the organisation more than others.

Another interesting graph that we can draw is the proportion of education version attrition.

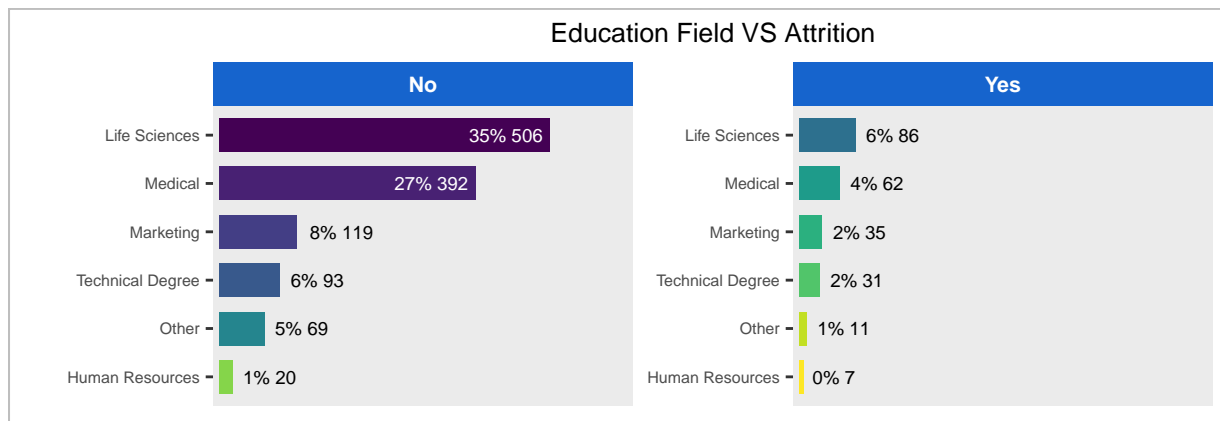
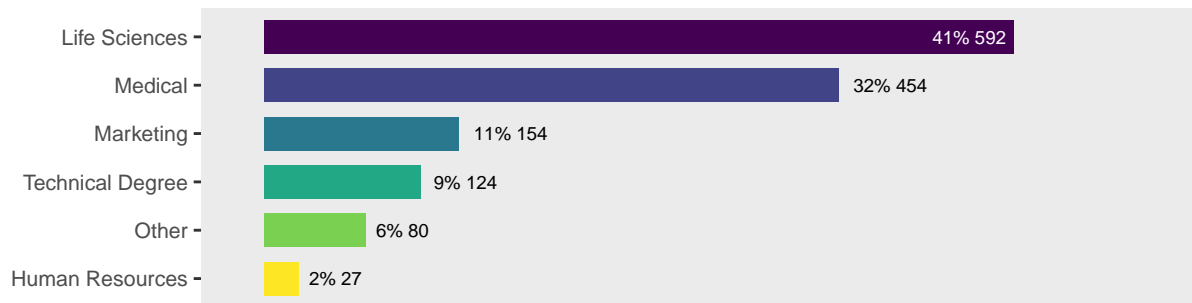
```

plt_education_field <- bar_plot_proportions(education_field)
plt_education_field_att <- bar_plot_proportions(education_field, attrition)

(plt_education_field /
  plt_education_field_att) +
  plot_annotation(
    title = "Proportions of Education field VS Attrition",
    caption = ""
  ) &
  theme(plot.caption = element_text(color = "#969696", size = 7))

```


Proportions of Education field VS Attrition



Hence, life sciences, medical will leave the company more frequently than the others.

Another variable that can affect the attrition is the dummy variable if the employees work over time or not.

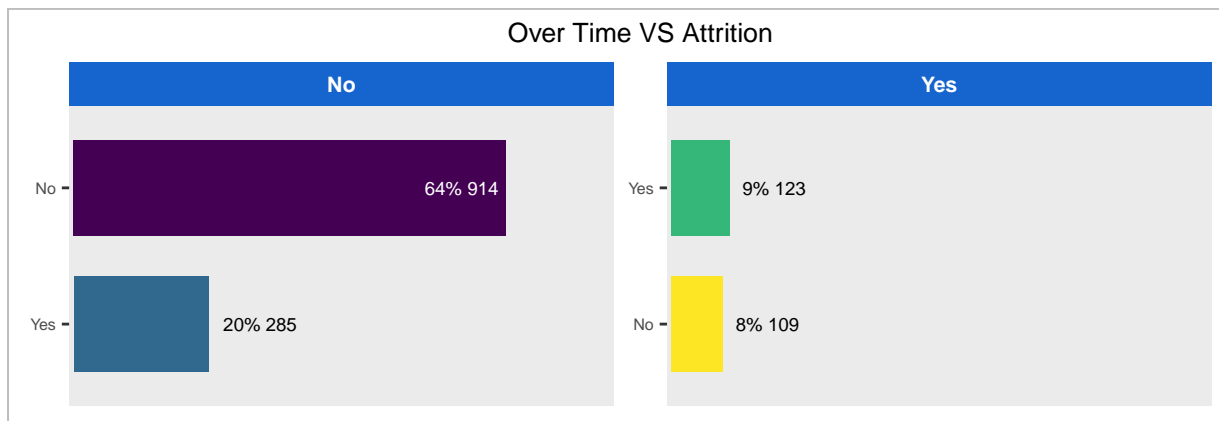
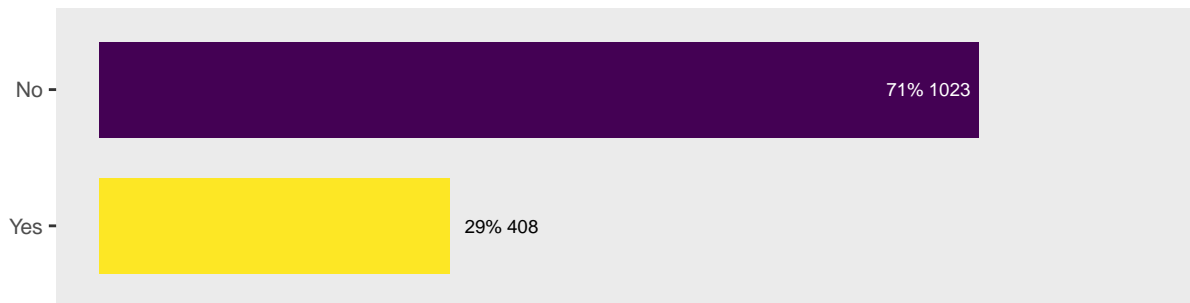
```

plt_education_field <- bar_plot_proportions(over_time)
plt_education_field_att <- bar_plot_proportions(over_time, attrition)

(plt_education_field /
  plt_education_field_att) +
  plot_annotation(
    title = "Proportions of over time work VS Attrition",
    caption = ""
  ) &
  theme(plot.caption = element_text(color = "#969696", size = 7))

```

Proportions of over time work VS Attrition



Therefore, we can conclude that if employees work more overtime they would tend to leave the company more than the people who do not. This variable is extremely significant because as we can see the majority of the workers don't work overtime.

5.3.3 Regression

In order to deal with *attrition* we have to convert the yes or no variable into a binary variable. Once we did that we can work with the logistic regression model. Obviously, we are dealing now with a dummy variable. Hence, we can try to understand what are the variables that will affect our attrition.

```
mydb$attrition <- ifelse(mydb$attrition=="Yes",1,0)
```

```
mod1 <- glm(attrition ~ ., data = mydb)
summary(mod1)
```

```
##
## Call:
## glm(formula = attrition ~ ., data = mydb)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -0.56380  -0.20926  -0.08298   0.08646   1.13122
##
## Coefficients: (2 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    4.916e-01  2.185e-01   2.250 0.024588 *
```

```
## age -3.542e-03 1.341e-03 -2.641 0.008358 **
## business_travelTravel_Frequently 1.517e-01 3.341e-02 4.542 6.06e-06 ***
## business_travelTravel_Rarely 6.205e-02 2.872e-02 2.161 0.030875 *
## daily_rate -2.586e-05 2.151e-05 -1.202 0.229386
## departmentResearch & Development -8.383e-03 1.107e-01 -0.076 0.939666
## departmentSales -7.603e-03 1.169e-01 -0.065 0.948150
## distance_from_home 3.733e-03 1.060e-03 3.521 0.000445 ***
## education 4.978e-04 8.662e-03 0.057 0.954184
## education_fieldLife Sciences -8.436e-02 8.315e-02 -1.015 0.310491
## education_fieldMarketing -3.531e-02 8.887e-02 -0.397 0.691157
## education_fieldMedical -9.266e-02 8.353e-02 -1.109 0.267465
## education_fieldOther -1.043e-01 8.962e-02 -1.164 0.244614
## education_fieldTechnical Degree 8.256e-03 8.722e-02 0.095 0.924596
## employee_count NA NA NA NA
## employee_number -1.202e-05 1.438e-05 -0.836 0.403552
## environment_satisfaction -4.232e-02 7.920e-03 -5.343 1.07e-07 ***
## genderMale 3.227e-02 1.772e-02 1.821 0.068761 .
## hourly_rate -1.807e-04 4.258e-04 -0.424 0.671339
## job_involvement -5.786e-02 1.219e-02 -4.748 2.27e-06 ***
## job_level -5.677e-03 2.889e-02 -0.197 0.844246
## job_roleHuman Resources 9.030e-02 1.179e-01 0.766 0.443924
## job_roleLaboratory Technician 1.374e-01 4.041e-02 3.400 0.000693 ***
## job_roleManager 3.070e-02 6.833e-02 0.449 0.653228
## job_roleManufacturing Director 5.641e-04 3.978e-02 0.014 0.988689
## job_roleResearch Director -1.384e-02 6.124e-02 -0.226 0.821273
## job_roleResearch Scientist 3.917e-02 4.011e-02 0.977 0.328942
## job_roleSales Executive 7.343e-02 7.833e-02 0.937 0.348721
## job_roleSales Representative 2.369e-01 8.716e-02 2.718 0.006640 **
## job_satisfaction -3.807e-02 7.855e-03 -4.847 1.39e-06 ***
## marital_statusMarried 1.582e-02 2.315e-02 0.684 0.494339
## marital_statusSingle 1.034e-01 3.182e-02 3.249 0.001186 **
## monthly_income 2.552e-06 7.673e-06 0.333 0.739467
## monthly_rate 4.119e-07 1.212e-06 0.340 0.734037
## num_companies_worked 1.743e-02 3.856e-03 4.521 6.69e-06 ***
## over18Y 2.014e-01 1.244e-01 1.619 0.105642
## over_timeYes 2.057e-01 1.921e-02 10.705 < 2e-16 ***
## percent_salary_hike -2.603e-03 3.729e-03 -0.698 0.485281
## performance_rating 1.798e-02 3.756e-02 0.479 0.632330
## relationship_satisfaction -2.225e-02 8.018e-03 -2.775 0.005596 **
## standard_hours NA NA NA NA
## stock_option_level -2.024e-02 1.384e-02 -1.462 0.143914
## total_working_years -4.342e-03 2.447e-03 -1.774 0.076240 .
## training_times_last_year -1.410e-02 6.750e-03 -2.089 0.036913 *
## work_life_balance -3.240e-02 1.223e-02 -2.650 0.008153 **
## years_at_company 2.993e-03 3.014e-03 0.993 0.320805
## years_in_current_role -7.562e-03 3.864e-03 -1.957 0.050563 .
## years_since_last_promotion 1.128e-02 3.471e-03 3.251 0.001179 **
## years_with_curr_manager -7.716e-03 4.043e-03 -1.908 0.056548 .
## cityLondon 1.143e-02 2.232e-02 0.512 0.608768
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 0.1037191)
##
```

```
## Null deviance: 194.39 on 1430 degrees of freedom
## Residual deviance: 143.44 on 1383 degrees of freedom
## AIC: 867.43
##
## Number of Fisher Scoring iterations: 2
```

```
mod2 <- glm(attrition ~ age + business_travel + distance_from_home + environment_satisfaction + job_inv
summary(mod2)
```

```
##
## Call:
## glm(formula = attrition ~ age + business_travel + distance_from_home +
## environment_satisfaction + job_involvement + job_role + job_satisfaction +
## num_companies_worked + over_time + relationship_satisfaction +
## training_times_last_year + work_life_balance + years_since_last_promotion,
## data = mydb)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -0.60053  -0.21285  -0.08711   0.05264   1.15940
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      0.636633   0.085509   7.445 1.68e-13 ***
## age             -0.005727   0.001110  -5.160 2.83e-07 ***
## business_travelTravel_Frequently  0.153883   0.033865   4.544 5.99e-06 ***
## business_travelTravel_Rarely      0.065628   0.029036   2.260 0.02396 *
## distance_from_home  0.003394   0.001079   3.147 0.00168 **
## environment_satisfaction -0.040395   0.008028  -5.032 5.48e-07 ***
## job_involvement    -0.064178   0.012312  -5.213 2.14e-07 ***
## job_roleHuman Resources  0.143701   0.054727   2.626 0.00874 **
## job_roleLaboratory Technician  0.157089   0.036353   4.321 1.66e-05 ***
## job_roleManager      0.028496   0.044866   0.635 0.52544
## job_roleManufacturing Director -0.010194   0.040356  -0.253 0.80062
## job_roleResearch Director -0.048350   0.047634  -1.015 0.31027
## job_roleResearch Scientist  0.066199   0.035785   1.850 0.06454 .
## job_roleSales Executive  0.095625   0.034582   2.765 0.00576 **
## job_roleSales Representative  0.307871   0.048330   6.370 2.55e-10 ***
## job_satisfaction     -0.037746   0.007946  -4.750 2.24e-06 ***
## num_companies_worked  0.017607   0.003705   4.752 2.22e-06 ***
## over_timeYes         0.203682   0.019525  10.432 < 2e-16 ***
## relationship_satisfaction -0.016991   0.008121  -2.092 0.03661 *
## training_times_last_year -0.014758   0.006849  -2.155 0.03136 *
## work_life_balance    -0.032849   0.012403  -2.649 0.00817 **
## years_since_last_promotion  0.003856   0.002876   1.341 0.18029
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 0.1085529)
##
## Null deviance: 194.39 on 1430 degrees of freedom
## Residual deviance: 152.95 on 1409 degrees of freedom
## AIC: 907.27
##
```

```
## Number of Fisher Scoring iterations: 2
```

```
rm(mod1, mod2)
```

To sum up, we can see that in the first regression we can see the different types of variables and their correlations, while in the second model we took only those variable with a low p-value in the *mod1*.

In conclusion the *mod2* shows all the variables that have a significative effect in the the definition of the attrition variable.

5.4 Satisfaction analysis

In this section we want to analyze the level of satisfaction of the employees and of the workplace more in general to see if there are some differences between the factors and characteristics that we are going to consider. For this reason our studyings are focused on the variable and their trends:

1-job_satisfaction

2-environment_satisfaction

5.4.1 1-Job_satisfaction

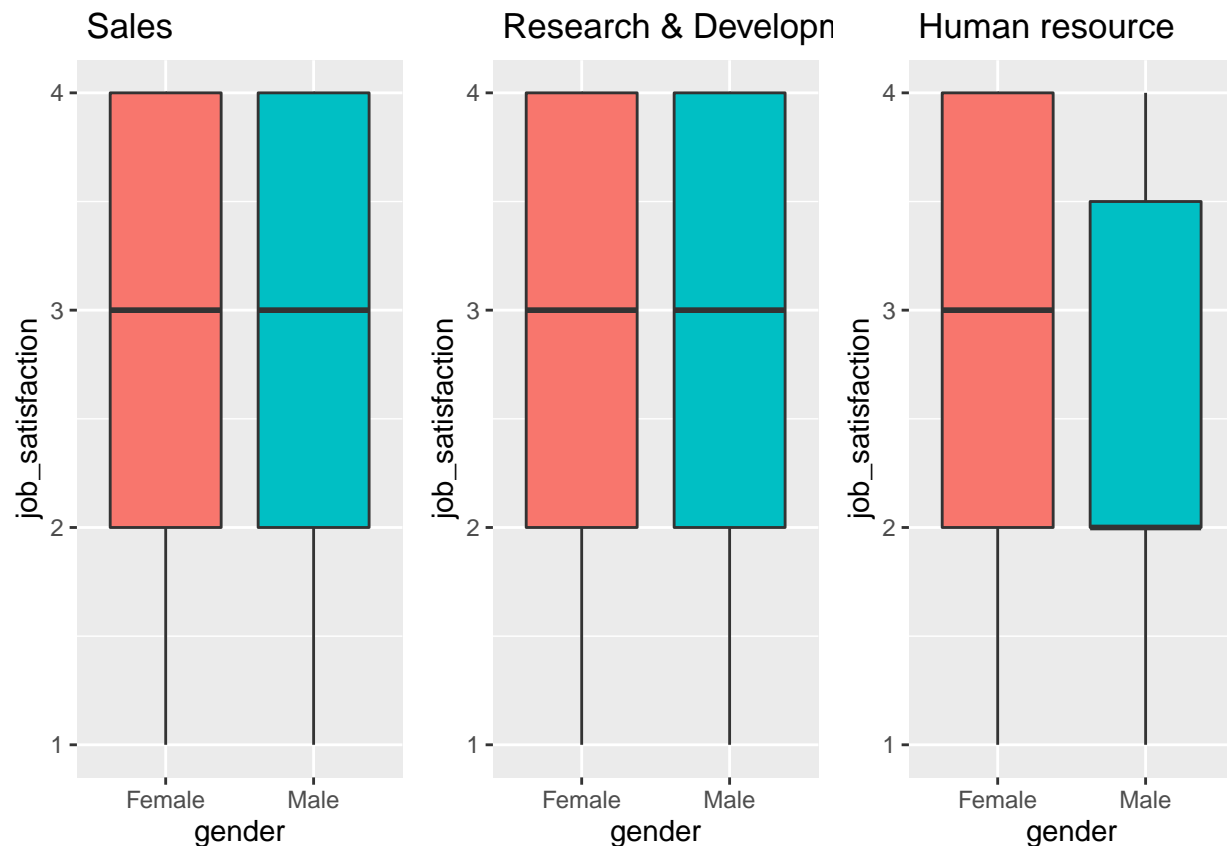
5.4.1.1 Gender and Department Considering the satisfaction of our employees to make a first summary of the data we represent its value in each department differentiated per gender. For this section we use some boxplot to explain our results as to have an idea on the main distribution of the values between different factor we will consider later.

```
##Sales
g1 <- ggplot(mydb%>%filter(department == "Sales"),
             aes(x= gender, y=job_satisfaction, fill= gender))+
  geom_boxplot(show.legend = FALSE)+
  labs(title = " Sales")

##Research & Development
g2 <- ggplot(mydb%>%filter(department == "Research & Development"),
             aes(x= gender, y=job_satisfaction, fill= gender))+
  geom_boxplot(show.legend = FALSE)+
  labs(title = " Research & Development ")

##Human Resource
g3 <- ggplot(mydb%>%filter(department == "Human Resources"),
             aes(x= gender, y=job_satisfaction, fill= gender))+
  geom_boxplot(show.legend = FALSE)+
  labs(title = " Human resource ")

grid.arrange(g1, g2, g3, ncol = 3)
```



In this first general overview we can observe a constant trend per each department, also due to the fact that we have a limited choice of value for the variable “job_satisfaction” and so it’s normal to see similar pattern. It’s relevant the distribution for male in the Human resource business unit, where we see that the median and the 1st quantile are overlapped and then presents the lowest value of satisfaction.

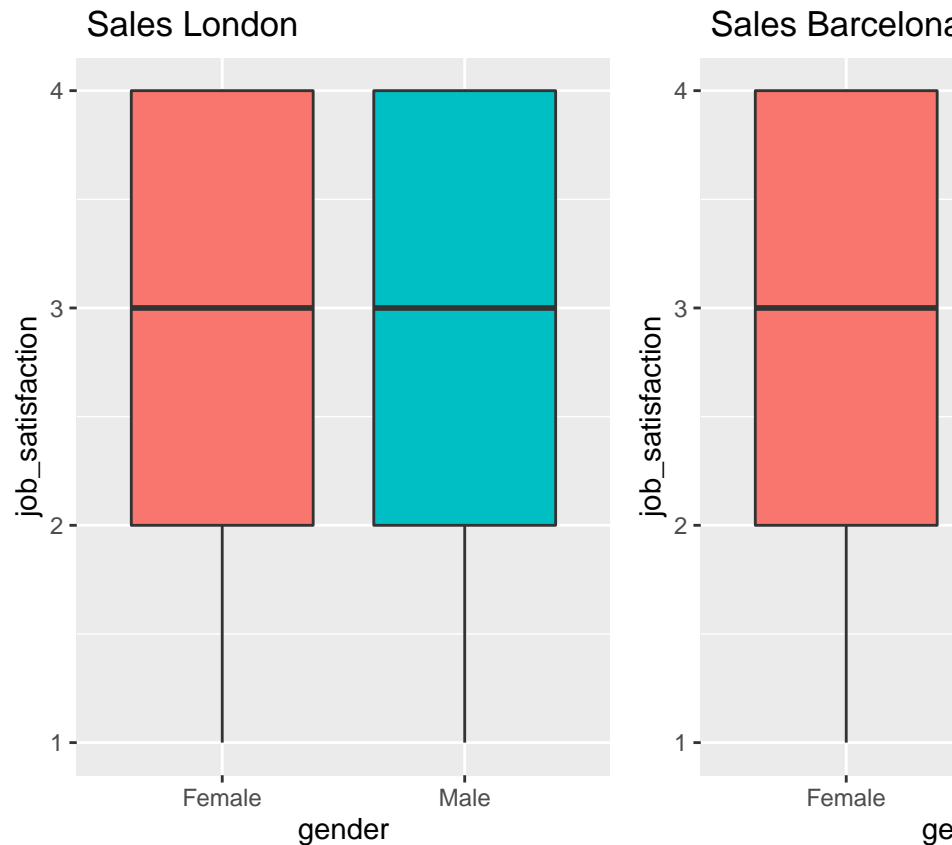
5.4.2 Differences between the two cities in term of job satisfaction

A possible way to go further in our analysis is to study the possible differences in the feelings of the workers between London and Barcelona to understand if there are some parameters or other point to go into detail. In this way we compare for each department the level of satisfaction per gender.

```
##Sales London
g4 <- ggplot(mydb%>%filter(department == "Sales" & city == "London"),
  aes(x= gender, y= job_satisfaction, fill= gender))+
  geom_boxplot(show.legend = FALSE)+
  labs(title = " Sales London ")

##Sales Barcelona
g5 <- ggplot(mydb%>%
  filter(department == "Sales" & city == "Barcelona"),
  aes(x= gender, y=job_satisfaction, fill= gender))+
  geom_boxplot(show.legend = FALSE)+
```

```
labs(title = " Sales Barcelona ")
grid.arrange(g4,g5, ncol = 2)
```



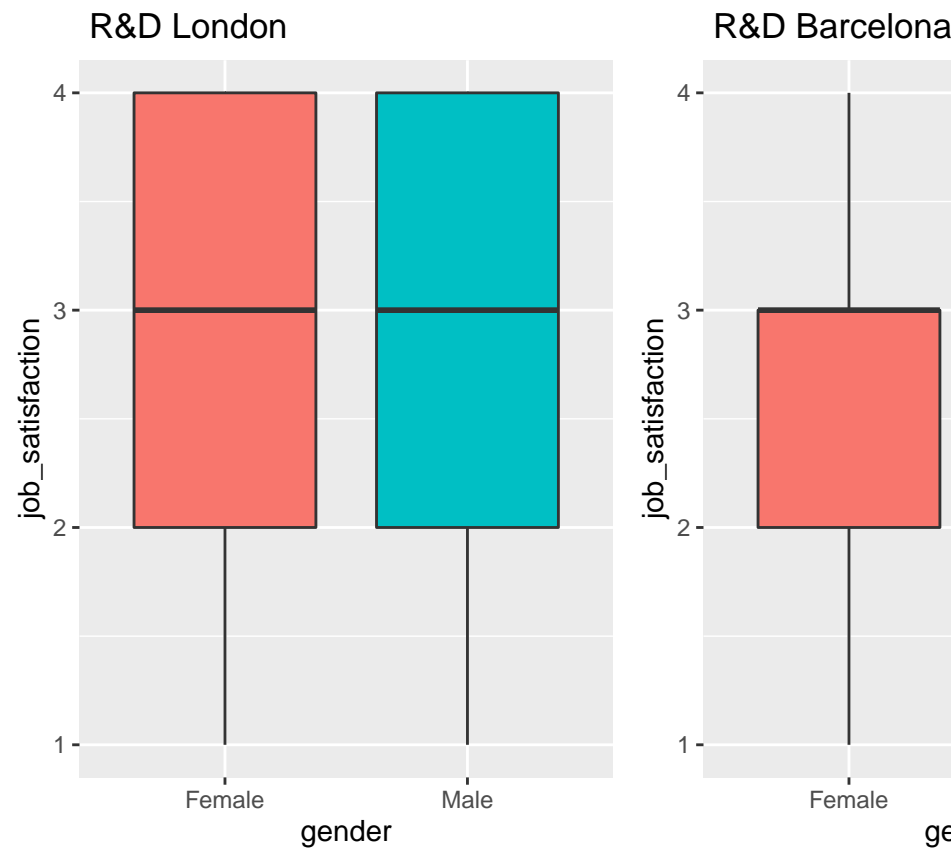
5.4.2.1 Sales department , Gender

Starting from Sales department we do not observe any peculiarity as the 2 graph are equal ; in general we can assume a medium level of satisfaction of the workers in this field, which stands at the value “3”.

```
##Research & Development London
g6 <- ggplot(mydb%>%
  filter(department == "Research & Development" & city == "London"),
  aes(x = gender, y = job_satisfaction, fill = gender)) +
  geom_boxplot(show.legend = FALSE) +
  labs(title = " R&D London")

##Research & Development Barcelona
g7 <- ggplot(mydb%>%
  filter(department == "Research & Development" & city == "Barcelona"),
  aes(x = gender, y = job_satisfaction, fill = gender)) +
  geom_boxplot(show.legend = FALSE) +
  labs(title = " R&D Barcelona")

grid.arrange(g6,g7, ncol = 2)
```



5.4.2.2 R&D department , Gender

For R&D we observe another time a regular trend close to 3 as average value of satisfaction. The only relevant note is the equality between the median and the 3rd quantile at the value 3, but except for that we can consider this as a medium satisfaction level department.

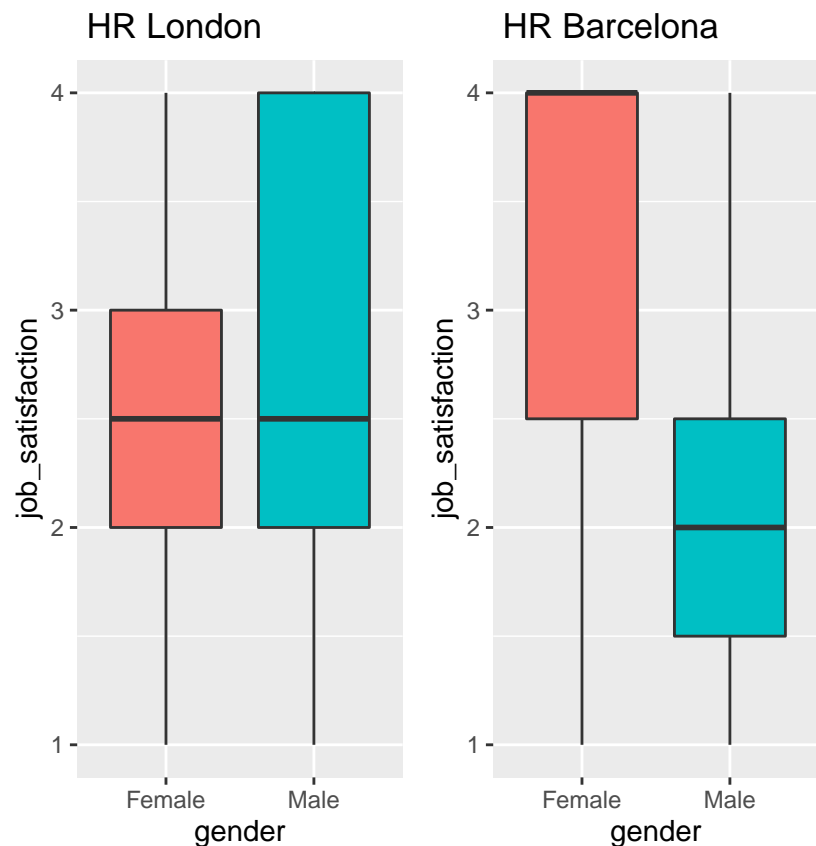
```
##Human Resource London
```

```
g8 <- ggplot(mydb%>%
  filter(department == "Human Resources" & city == "London"),
  aes(x = gender, y = job_satisfaction, fill = gender)) +
  geom_boxplot(show.legend = FALSE) +
  labs(title = "HR London")
```

```
##Human Resource Barcelona
```

```
g9 <- ggplot(mydb%>%
  filter(department == "Human Resources" & city == "Barcelona"),
  aes(x = gender, y = job_satisfaction, fill = gender)) +
  geom_boxplot(show.legend = FALSE) +
  labs(title = "HR Barcelona")
```

```
grid.arrange(g8, g9, ncol = 3)
```

5.4.2.3 HR department , Gender

Going deep in H&R department we see that the previous graph present inside it some interest sides ; in fact it is clearly evident from this division some difference between the 2 cities. For London size we observe that the median is located between 2 and 3 for both male and female gender , with maximum values for male that reach 4 , that is also the 3rd quantile. On the other side in Barcelona we see a big discrepancy inside gender : in fact considering male they present lowest values, a median which is equal to 2 (medium low value) and also with 1 as 1st quantile. Opposite is the trend for female in which we can observe a median equal to 4 which indicates a good level of satisfaction. That is can due to the fact that , doing further analysis on the number of employees, there is a clear predominance of female in this department in Barcelona and so this a general course could be not the best for the integration of male in this sector.

However to have a main view of the graph presented as now we resume them in this picture :

```
grid.arrange(g1, g2, g3, g4, g5, g6, g7, g8, g9, ncol = 3, top = "Comparison between job satisfaction a
```

Comparison between job satisfaction and gender



```
# clean variables
```

```
rm(g1, g2, g3, g4, g5, g6, g7, g8, g9)
```

```
mean_mi <- round(mean(mydb$monthly_income),2)
```

```
median_mi <- round(median(mydb$monthly_income),2)
```

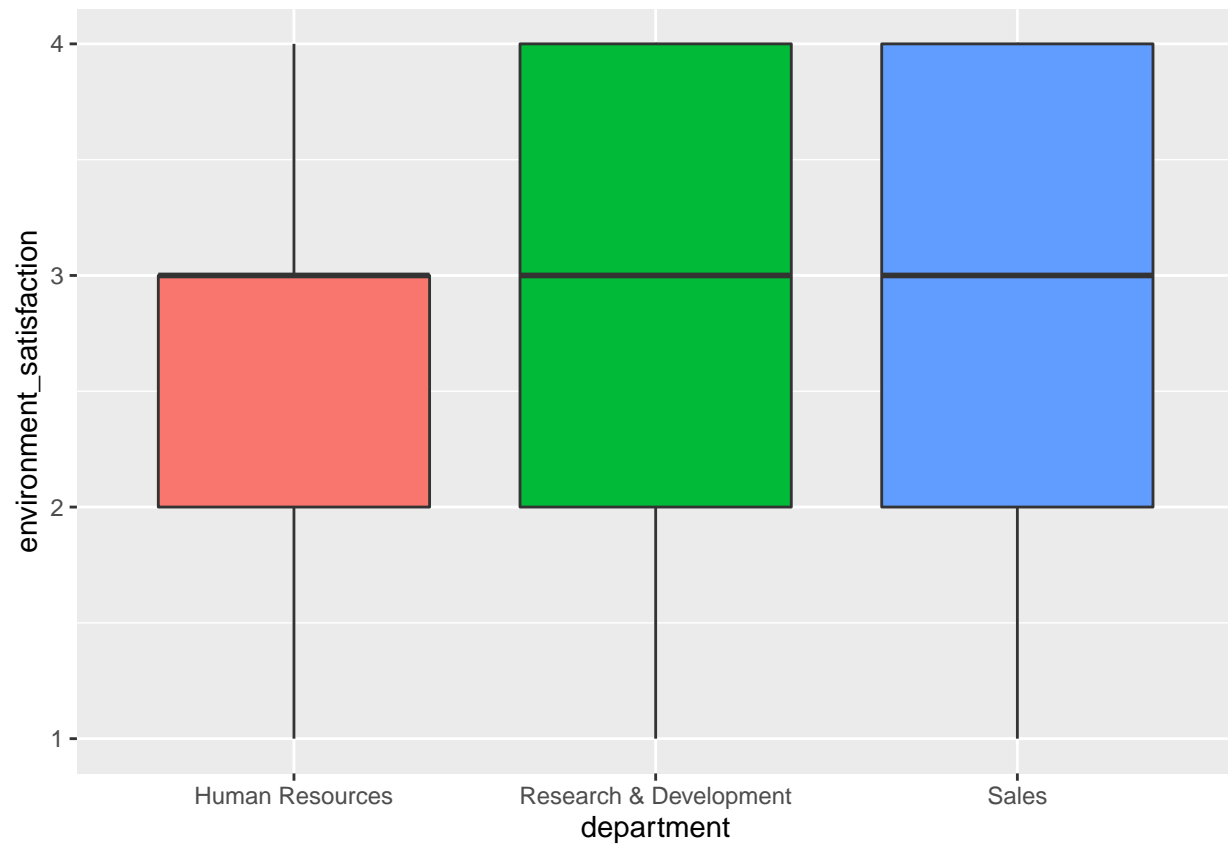
```
ggplot(mydb,
  aes(x = gender,
      y = monthly_income,
      fill = gender)) +
  geom_boxplot(show.legend = FALSE) +
  coord_flip() +
  scale_y_continuous(labels = label_dollar()) +
  labs(title = "The gender gap in monthly income",
       caption = "Data Source: Kaggle IBM HR Employee Attrition",
       x = "Gender",
       y = "Monthly Income")
```



Environemnt_satisfaction

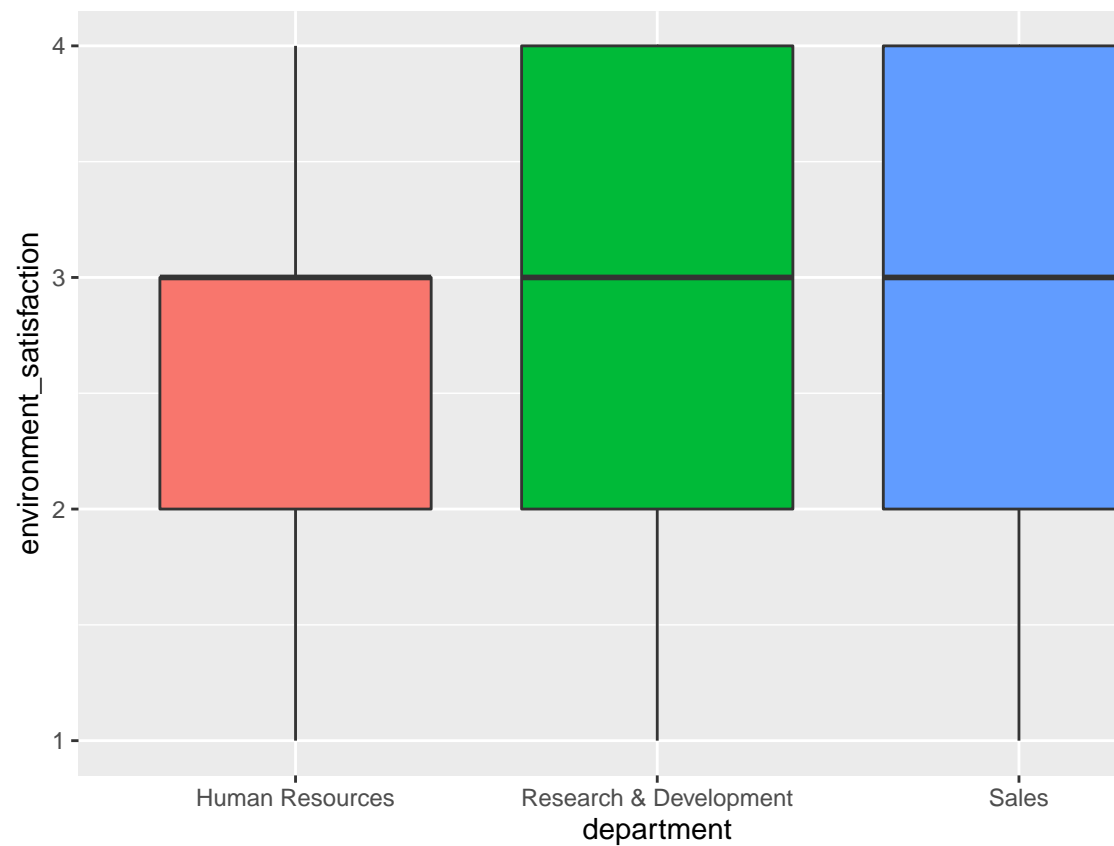
5.4.2.4 Environment Satisfaction per department Now we focus on the satisfaction on the work-place and we try to go deepen if there are some interesting aspects into the single department.

```
ggplot(mydb, aes(x= department, y= environment_satisfaction, fill = department))+
  geom_boxplot(show.legend = FALSE)
```



Generally there is a medium level of satisfaction for R&D and Sales department while we see a medium negative attitude in Human Resources. However this pattern is not so outstanding, reason for what we decide to go deeper in considerations per gender and per city.

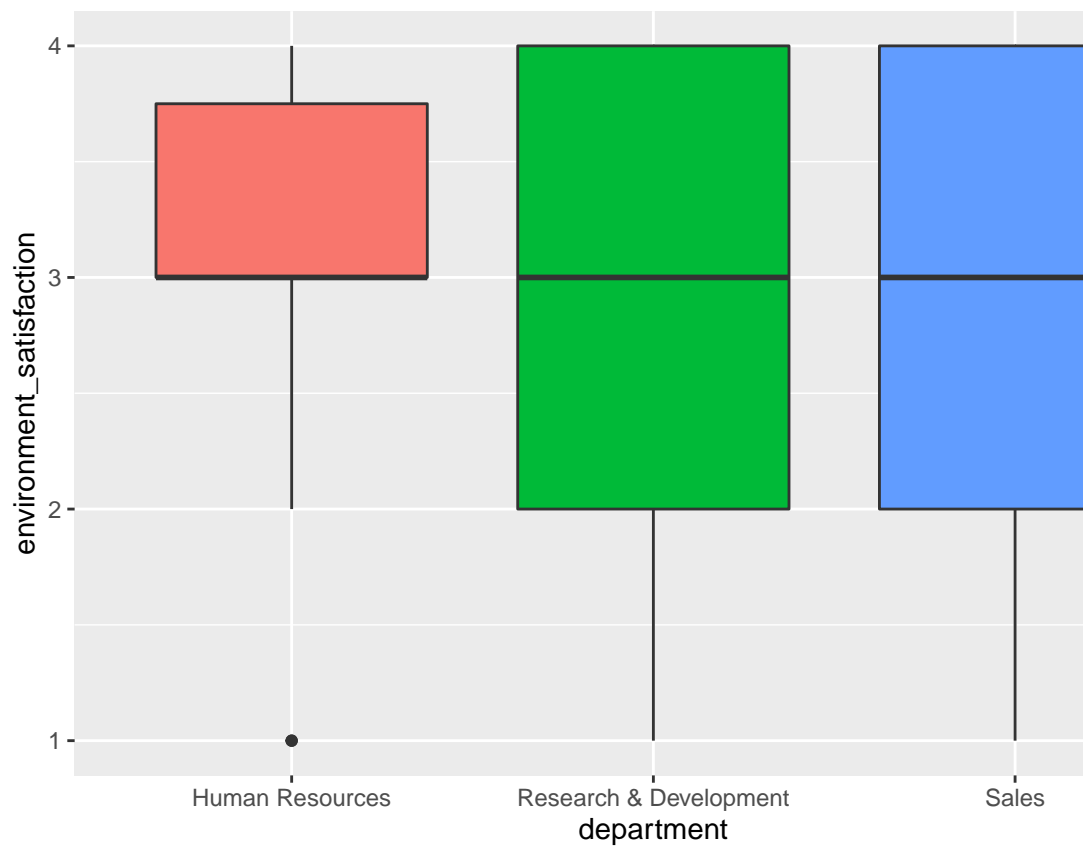
```
ggplot(mydb%>%filter(city == "London", gender == "Male"), aes(x= department, y= environment_satisfaction))  
  geom_boxplot(show.legend = FALSE)
```



5.4.2.5 London , Male

Firstly we have filtered per “Male” and “London” and as in the previous analysis we observe a medium low trend for male in Human Resources department. For the other two business unit we observe a constant trend that we have found also considering job_satisfaction.

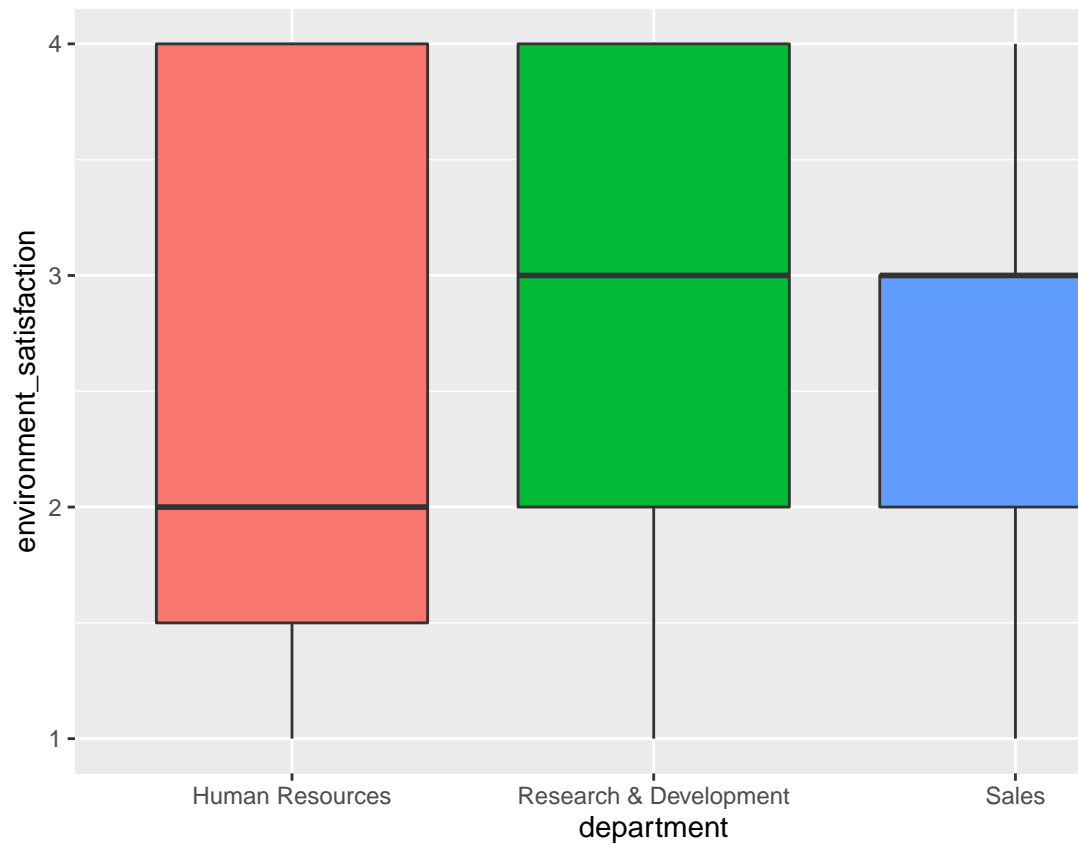
```
ggplot(mydb%>%filter(city == "London", gender== "Female"), aes(x= department, y= environment_satisfaction))  
  geom_boxplot(show.legend = FALSE)
```



5.4.2.6 London , Female

As we were supposing to see the medium low level in Human Resources for male is balanced in the female as to reach in general the medium trend (=3) of environment satisfaction which we can also find for the other 2 department.

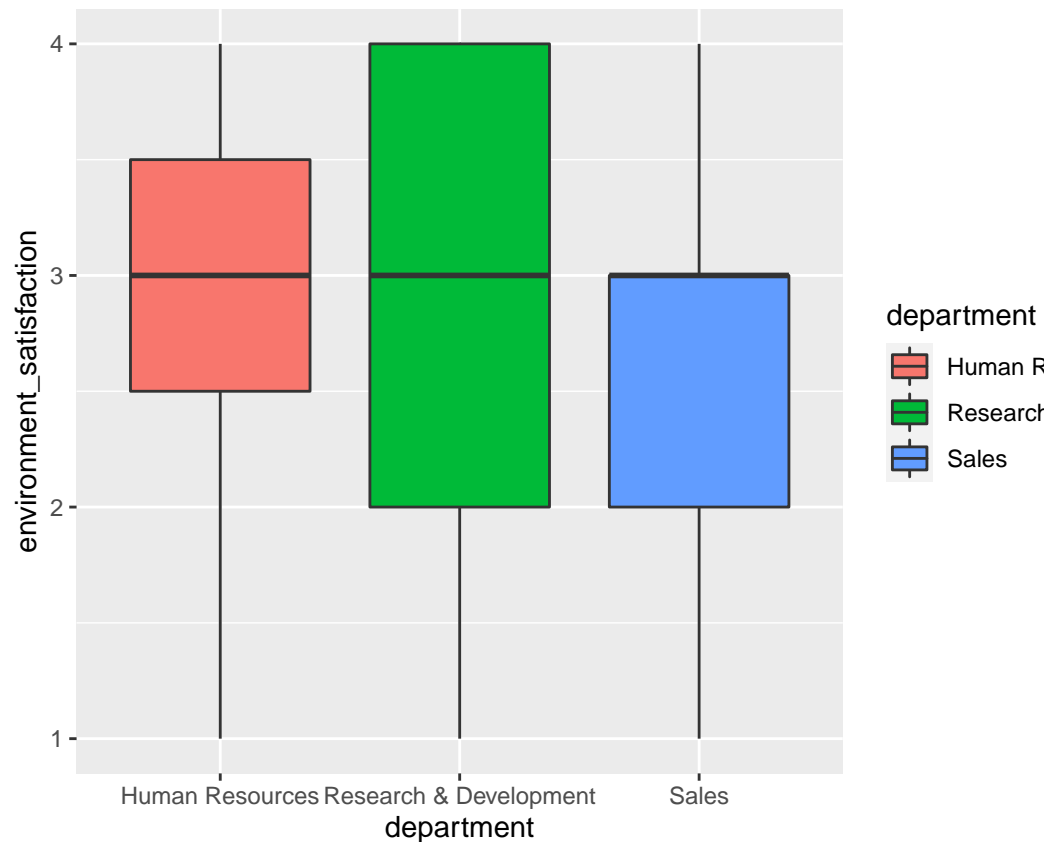
```
ggplot(mydb%>%filter(city == "Barcelona", gender== "Male"), aes(x= department, y= environment_satisfaction))  
  geom_boxplot(show.legend = FALSE)
```



5.4.2.7 Barcelona , Male

For the city of Barcelona it's confirmed the unsatisfaction(in this case related if the environment) for male in Human Resource department : in fact even if we have 1st quantile = 4 , we see the median equal to 2 , value that can not be appreciated by the company. R&D presents a positive average of values with the 1st quantile =4 and the median=3 while the Sales department has lowest value distribution that stay most in a range between 2 and 3.

```
ggplot(mydb%>%filter(city == "Barcelona", gender== "Female"), aes(x= department, y= environment_satisfaction)) +  
  geom_boxplot()+ labs(show.legend = FALSE)
```



5.4.2.8 Barcelona , Female

In this last graph , we do not see differences from the boxplot of R&D and Sales which are the same as the one of male. However, as we were expecting , Human Resources department follows a pretty positive course , thus going against the trend of the other sex.

This series of graph so confirms the big differences between male and female satisfaction level just for the city of Barcelona.

5.4.2.9 Regression We firstly try to interpretate the variable “job_satisfaction” and “environment_satisfaction” using the linear regression model :

```
mod4 <- lm(job_satisfaction ~ ., data = mydb)
summary(mod4)
```

```
##
## Call:
## lm(formula = job_satisfaction ~ ., data = mydb)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.2016 -0.8105  0.1685  1.0232  1.9179
##
## Coefficients: (2 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   4.235e+00  7.342e-01   5.768 9.88e-09 ***
## age           4.050e-03  4.562e-03   0.888  0.37481
## attrition     -4.388e-01  9.052e-02  -4.847 1.39e-06 ***
```



```
## business_travelTravel_Frequently 9.581e-02 1.142e-01 0.839 0.40168
## business_travelTravel_Rarely -2.949e-02 9.765e-02 -0.302 0.76270
## daily_rate 6.281e-05 7.303e-05 0.860 0.38993
## departmentResearch & Development -1.273e-01 3.759e-01 -0.339 0.73484
## departmentSales 2.951e-02 3.968e-01 0.074 0.94072
## distance_from_home -1.702e-04 3.616e-03 -0.047 0.96245
## education -7.746e-03 2.940e-02 -0.263 0.79227
## education_fieldLife Sciences 1.193e-01 2.824e-01 0.422 0.67282
## education_fieldMarketing -4.812e-02 3.017e-01 -0.159 0.87331
## education_fieldMedical 1.072e-02 2.837e-01 0.038 0.96987
## education_fieldOther 9.001e-02 3.044e-01 0.296 0.76749
## education_fieldTechnical Degree 9.608e-03 2.961e-01 0.032 0.97412
## employee_count NA NA NA NA
## employee_number -8.635e-05 4.878e-05 -1.770 0.07691 .
## environment_satisfaction -2.127e-02 2.716e-02 -0.783 0.43362
## genderMale 9.140e-02 6.016e-02 1.519 0.12892
## hourly_rate -3.858e-03 1.442e-03 -2.676 0.00754 **
## job_involvement -5.069e-02 4.168e-02 -1.216 0.22419
## job_level 2.230e-02 9.808e-02 0.227 0.82019
## job_roleHuman Resources -2.292e-01 4.003e-01 -0.573 0.56700
## job_roleLaboratory Technician -4.056e-02 1.378e-01 -0.294 0.76849
## job_roleManager -9.080e-02 2.319e-01 -0.391 0.69552
## job_roleManufacturing Director -1.168e-01 1.350e-01 -0.865 0.38707
## job_roleResearch Director -7.072e-02 2.079e-01 -0.340 0.73374
## job_roleResearch Scientist 2.680e-02 1.362e-01 0.197 0.84407
## job_roleSales Executive -1.261e-01 2.660e-01 -0.474 0.63561
## job_roleSales Representative -1.346e-01 2.966e-01 -0.454 0.65015
## marital_statusMarried 8.875e-02 7.856e-02 1.130 0.25875
## marital_statusSingle 2.567e-01 1.082e-01 2.372 0.01784 *
## monthly_income 8.039e-07 2.605e-05 0.031 0.97538
## monthly_rate -6.230e-07 4.115e-06 -0.151 0.87968
## num_companies_worked -1.949e-02 1.317e-02 -1.479 0.13924
## over18Y -7.884e-01 4.221e-01 -1.868 0.06199 .
## over_timeYes 1.366e-01 6.778e-02 2.015 0.04405 *
## percent_salary_hike 6.169e-03 1.266e-02 0.487 0.62609
## performance_rating -6.804e-02 1.275e-01 -0.534 0.59369
## relationship_satisfaction -1.733e-02 2.729e-02 -0.635 0.52553
## standard_hours NA NA NA NA
## stock_option_level 6.979e-02 4.699e-02 1.485 0.13768
## total_working_years -8.124e-03 8.314e-03 -0.977 0.32866
## training_times_last_year -2.176e-02 2.294e-02 -0.949 0.34300
## work_life_balance -3.828e-02 4.160e-02 -0.920 0.35767
## years_at_company 3.407e-03 1.023e-02 0.333 0.73929
## years_in_current_role 1.521e-02 1.313e-02 1.159 0.24682
## years_since_last_promotion -2.323e-03 1.183e-02 -0.196 0.84434
## years_with_curr_manager -2.717e-02 1.372e-02 -1.980 0.04788 *
## cityLondon -6.203e-04 7.578e-02 -0.008 0.99347
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.093 on 1383 degrees of freedom
## Multiple R-squared:  0.04775,    Adjusted R-squared:  0.01539
## F-statistic: 1.476 on 47 and 1383 DF,  p-value: 0.02088
```

```
mod5 <- lm(environment_satisfaction ~ ., data = mydb)
summary(mod5)
```

```
##
## Call:
## lm(formula = environment_satisfaction ~ ., data = mydb)
##
## Residuals:
```

	Min	1Q	Median	3Q	Max
	-2.2781	-0.8165	0.1687	0.9918	2.0542

```
##
## Coefficients: (2 not defined because of singularities)
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	2.585e+00	7.322e-01	3.530	0.000429 ***
age	8.935e-04	4.518e-03	0.198	0.843243
attrition	-4.779e-01	8.944e-02	-5.343	1.07e-07 ***
business_travelTravel_Frequently	1.118e-02	1.131e-01	0.099	0.921274
business_travelTravel_Rarely	5.562e-03	9.667e-02	0.058	0.954121
daily_rate	4.980e-05	7.231e-05	0.689	0.491088
departmentResearch & Development	-3.833e-01	3.720e-01	-1.030	0.302966
departmentSales	-5.335e-01	3.926e-01	-1.359	0.174390
distance_from_home	-4.296e-04	3.579e-03	-0.120	0.904482
education	-3.578e-02	2.909e-02	-1.230	0.218898
education_fieldLife Sciences	-3.662e-03	2.795e-01	-0.013	0.989551
education_fieldMarketing	1.659e-01	2.986e-01	0.556	0.578632
education_fieldMedical	-5.093e-03	2.808e-01	-0.018	0.985532
education_fieldOther	3.091e-01	3.012e-01	1.026	0.304931
education_fieldTechnical Degree	1.434e-01	2.931e-01	0.489	0.624727
employee_count	NA	NA	NA	NA
employee_number	2.365e-05	4.834e-05	0.489	0.624666
genderMale	3.187e-02	5.960e-02	0.535	0.592853
hourly_rate	-2.701e-02	1.429e-03	-1.890	0.058943 .
job_involvement	-3.629e-02	4.127e-02	-0.879	0.379441
job_level	5.397e-03	9.709e-02	0.056	0.955683
job_roleHuman Resources	-5.135e-01	3.961e-01	-1.296	0.195083
job_roleLaboratory Technician	1.802e-02	1.364e-01	0.132	0.894897
job_roleManager	4.089e-03	2.296e-01	0.018	0.985794
job_roleManufacturing Director	1.332e-01	1.336e-01	0.997	0.319091
job_roleResearch Director	-2.512e-01	2.057e-01	-1.221	0.222118
job_roleResearch Scientist	-5.472e-02	1.348e-01	-0.406	0.684917
job_roleSales Executive	3.037e-02	2.633e-01	0.115	0.908205
job_roleSales Representative	1.971e-01	2.936e-01	0.671	0.502239
job_satisfaction	-2.084e-02	2.661e-02	-0.783	0.433615
marital_statusMarried	-3.291e-02	7.780e-02	-0.423	0.672355
marital_statusSingle	4.698e-02	1.073e-01	0.438	0.661689
monthly_income	-3.230e-06	2.579e-05	-0.125	0.900322
monthly_rate	6.286e-06	4.070e-06	1.545	0.122684
num_companies_worked	2.070e-02	1.304e-02	1.587	0.112636
over18Y	9.346e-01	4.176e-01	2.238	0.025383 *
over_timeYes	2.652e-01	6.681e-02	3.970	7.57e-05 ***
percent_salary_hike	-8.637e-03	1.253e-02	-0.689	0.490737
performance_rating	-3.473e-02	1.262e-01	-0.275	0.783249
relationship_satisfaction	-8.075e-03	2.702e-02	-0.299	0.765108

```
## standard_hours          NA          NA          NA          NA
## stock_option_level      2.809e-03  4.655e-02  0.060 0.951891
## total_working_years    -6.205e-03  8.231e-03 -0.754 0.451076
## training_times_last_year -1.409e-02  2.272e-02 -0.620 0.535120
## work_life_balance       2.448e-02  4.119e-02  0.594 0.552377
## years_at_company        3.227e-03  1.013e-02  0.318 0.750155
## years_in_current_role   1.600e-03  1.300e-02  0.123 0.902086
## years_since_last_promotion 1.507e-02  1.170e-02  1.288 0.197950
## years_with_curr_manager -8.518e-03  1.360e-02 -0.626 0.531302
## cityLondon              5.384e-02  7.500e-02  0.718 0.472989
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.082 on 1383 degrees of freedom
## Multiple R-squared:  0.0517, Adjusted R-squared:  0.01947
## F-statistic: 1.604 on 47 and 1383 DF, p-value: 0.006356
```

```
mod6 <-lm(job_satisfaction ~ gender + city, data = mydb)
summary(mod6)
```

```
##
## Call:
## lm(formula = job_satisfaction ~ gender + city, data = mydb)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.774 -0.754  0.246  1.246  1.315
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.70492    0.07637  35.418  <2e-16 ***
## genderMale   0.06887    0.05953   1.157   0.248
## cityLondon  -0.01978    0.07522  -0.263   0.793
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.102 on 1428 degrees of freedom
## Multiple R-squared:  0.0009779, Adjusted R-squared: -0.0004213
## F-statistic: 0.6989 on 2 and 1428 DF, p-value: 0.4973
```

```
mod7 <- lm(environment_satisfaction ~ gender + city, data = mydb)
summary(mod7)
```

```
##
## Call:
## lm(formula = environment_satisfaction ~ gender + city, data = mydb)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.7334 -0.7334  0.2666  1.2666  1.3278
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
##
```

```
## (Intercept) 2.672169    0.075783   35.261    <2e-16 ***
## genderMale  0.007919    0.059073    0.134    0.893
## cityLondon  0.053339    0.074645    0.715    0.475
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.094 on 1428 degrees of freedom
## Multiple R-squared:  0.0003723, Adjusted R-squared:  -0.001028
## F-statistic: 0.2659 on 2 and 1428 DF,  p-value: 0.7665
```

We can observe that both “gender” and “city” are not relevant to explain the two variables studied “job_satisfaction” and “environment_satisfaction”. Anyway is it clear from the previous part that there is a more pronounced dissatisfaction among Spanish male workers and so in this sense would be useful to understand the motivations of this general mood and also the difference between them and London workers in their workplace. A more deepen analysis on the behavior of the workers in this company, the approach between males and females and , more than all , their growth prospects within the company : in our opinion this could be the key in order to solve this problem.

6 Adding dataset

```
# UK
london_data<-read.csv("survey_results_public2.csv")
glimpse(london_data)

# Barcelona
barcelona_data <- read.csv("barcelona_dataset.csv")

# Clean Barcelona
barcelona_data <- rename(barcelona_data, age = db_bar.EDAT1899_1A7)
barcelona_data <- subset(barcelona_data, select = -X)
tabyl(barcelona_data$sex)
barcelona_data$sex <- ifelse(barcelona_data$sex ==
                             "DONA",
                             "Female", barcelona_data$sex)
barcelona_data$sex <- ifelse(barcelona_data$sex ==
                             "HOME",
                             "Male", barcelona_data$sex)

tabyl(barcelona_data$monthly_income)
barcelona_data$monthly_income <- ifelse(barcelona_data$monthly_income ==
                                         "DE 1.001 A 1.500 EUROS",
                                         (1001+1500)/2, barcelona_data$monthly_income)
barcelona_data$monthly_income <- ifelse(barcelona_data$monthly_income ==
                                         "DE 2.001 A 2.500 EUROS",
                                         (2001+2500)/2, barcelona_data$monthly_income)
barcelona_data$monthly_income <- ifelse(barcelona_data$monthly_income ==
                                         "DE 2.501 A 3.000 EUROS",
                                         (2501+3000)/2, barcelona_data$monthly_income)
barcelona_data$monthly_income <- ifelse(barcelona_data$monthly_income ==
                                         "DE 3.001 A 5.000 EUROS",
                                         (3001+5000)/2, barcelona_data$monthly_income)
```

```
barcelona_data$monthly_income <- ifelse(barcelona_data$monthly_income ==
  "DE 5.001 A 7.000 EUROS",
  (5001+7000)/2, barcelona_data$monthly_income)
barcelona_data$monthly_income <- ifelse(barcelona_data$monthly_income ==
  "DE 500 A 1.000",
  (500+1000)/2, barcelona_data$monthly_income)
barcelona_data$monthly_income <- ifelse(barcelona_data$monthly_income ==
  "DE 7.001 A 9.000 EUROS",
  (7001+9000)/2, barcelona_data$monthly_income)
barcelona_data$monthly_income <- ifelse(barcelona_data$monthly_income ==
  "MENYS DE 500 EUROS",
  (500)/2, barcelona_data$monthly_income)
barcelona_data$monthly_income <- ifelse(barcelona_data$monthly_income ==
  "DE 1.501 A 2.000 EUROS",
  (1.501+2000)/2, barcelona_data$monthly_income)
barcelona_data$monthly_income <- ifelse(barcelona_data$monthly_income ==
  "MÉS DE 9.000 EUROS",
  (10000), barcelona_data$monthly_income)
barcelona_data$monthly_income <- ifelse(barcelona_data$monthly_income ==
  "DE 500 A 1.000 EUROS",
  (500+1000)/2, barcelona_data$monthly_income)

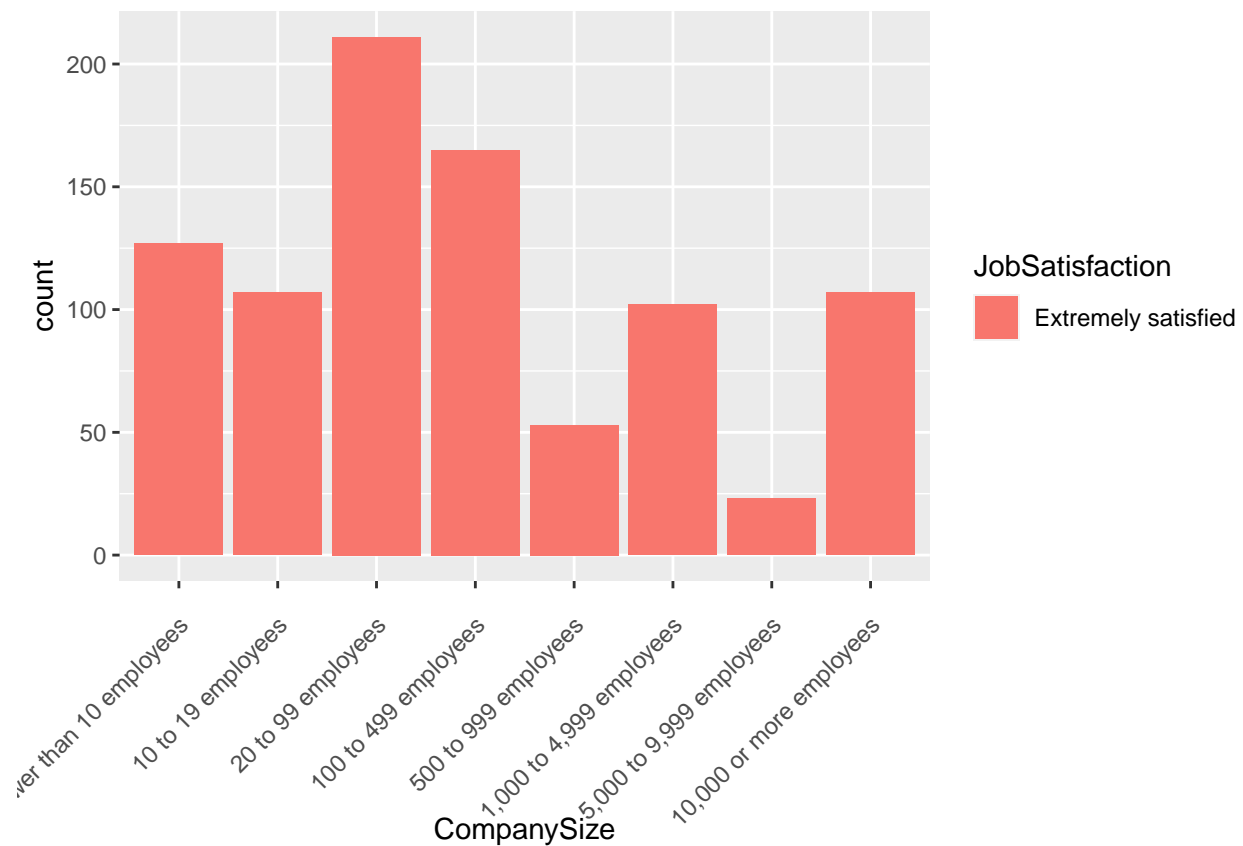
barcelona_data <- barcelona_data %>% filter(monthly_income != "NO CONTESTA")
barcelona_data <- barcelona_data %>% filter(monthly_income != "NO HO SAP")

barcelona_data$monthly_income <- as.numeric(barcelona_data$monthly_income)
```

6.1 UK Comparation

```
ggplot(london_data%>%filter(JobSatisfaction == "Extremely satisfied"), aes(x= CompanySize, fill = JobSa
  theme( axis.text.x = element_text(angle=45, hjust=1, vjust=0.9))+ scale_x_discrete(limits = c("Fewer

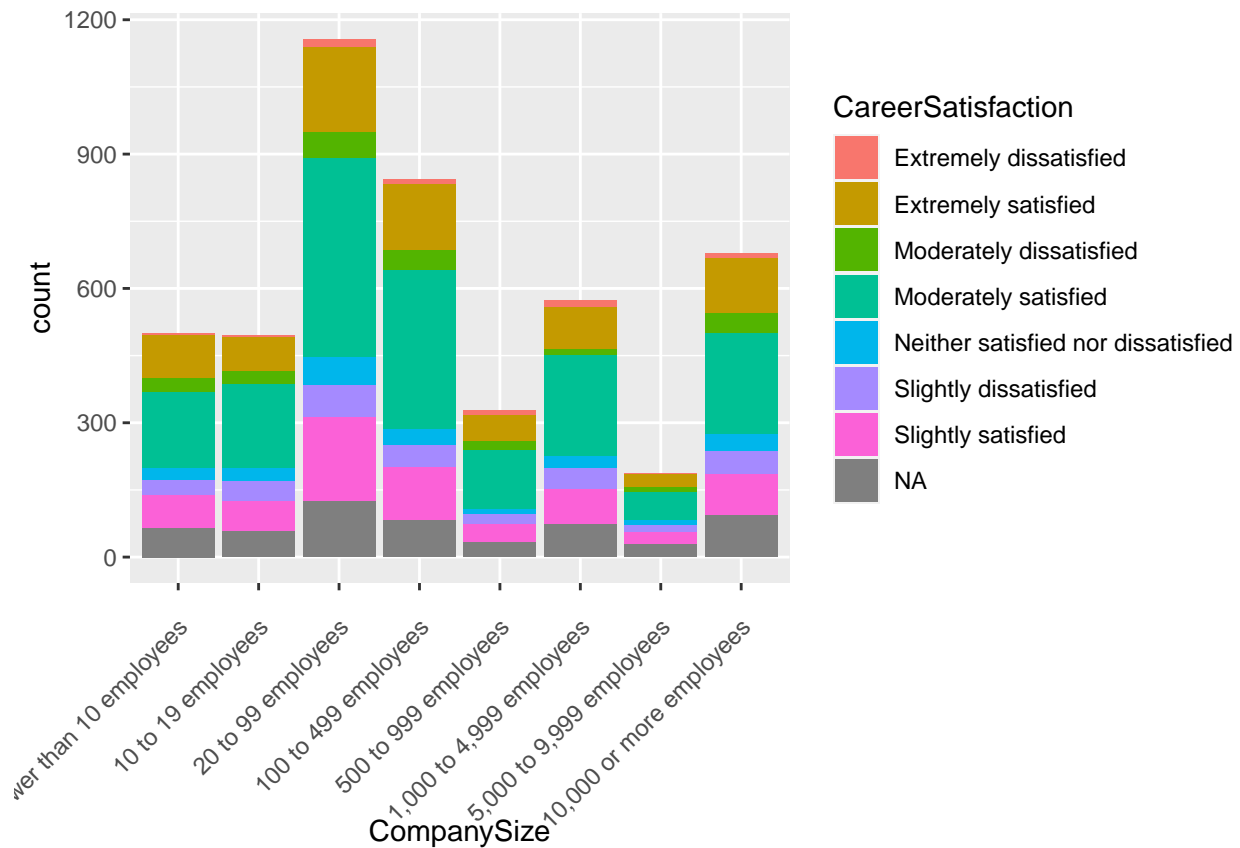
## Warning: Removed 150 rows containing non-finite values (stat_count).
```



As we can see from this graph in the Uk if people work in a company which is small the satisfaction of the employee tends to increase with respect to a big vompany

6.1.1 Carrer satisfaction

```
ggplot(london_data, aes(x= CompanySize , fill = CareerSatisfaction))+ geom_bar()+
  theme( axis.text.x = element_text(angle=45, hjust=1, vjust=0.9))+ scale_x_discrete(limits = c("Fewer
```



6.2 Barcelona Comparison

6.2.1 Introduction Analysis

```
summary(barcelona_data)
```

```
##      year      sex      age      monthly_income
##  Min.   :2021   Length:2266   Length:2266   Min.    : 250
##  1st Qu.:2021   Class :character   Class :character   1st Qu. : 1001
##  Median :2021   Mode  :character   Mode  :character   Median  : 2250
##  Mean   :2021                                     Mean   : 2686
##  3rd Qu.:2021                                     3rd Qu.: 4000
##  Max.   :2021                                     Max.   :10000
```

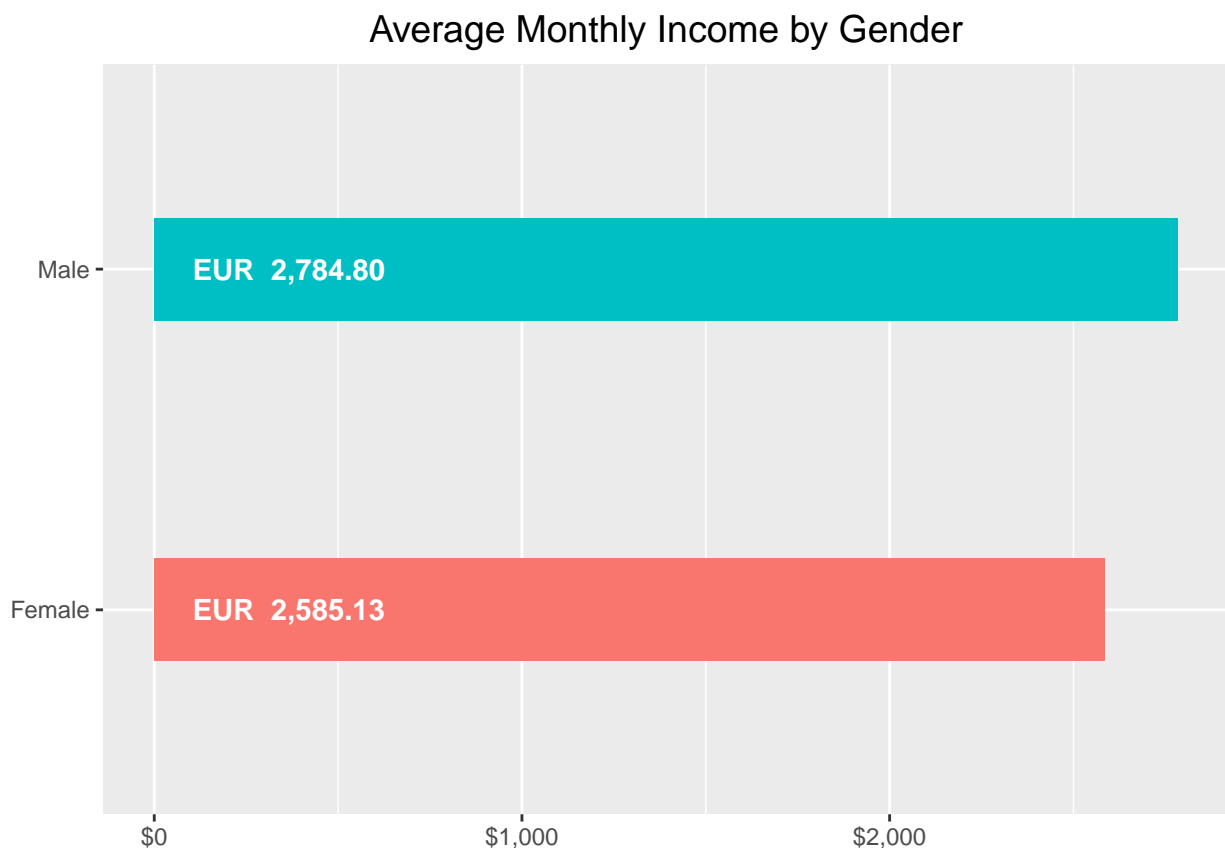
```
tabyl(barcelona_data$sex)
```

```
## barcelona_data$sex    n    percent
##           Female 1122 0.4951456
##           Male  1144 0.5048544
```

From this data we can see how the general mean of the monthly_income is €2686. Moreover, we can see how well this dataset is distributed between Male and Female.

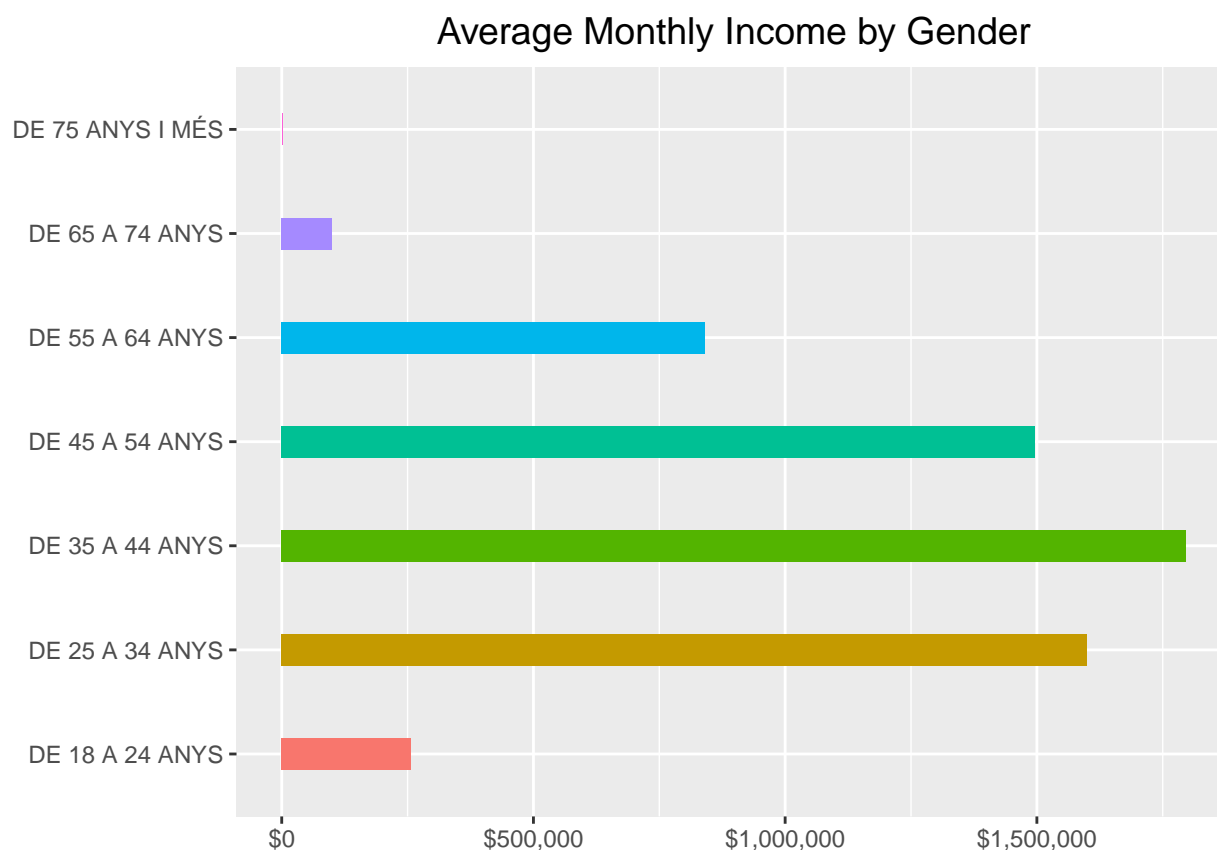
With the cleaned new dataset, we can make some analysis in order to have a general overview and can do a comparison between the starting dataset solution and the new one.

```
barcelona_data %>%
  select(sex, monthly_income) %>%
  group_by(sex) %>%
  summarise(avg_income = round(mean(monthly_income), 2), .groups = "drop") %>%
  ggplot(aes(x = sex, y = avg_income)) +
  geom_col(aes(fill = sex), width = 0.3, show.legend = FALSE) +
  geom_text(
    aes(
      x = sex,
      y = 0.01,
      label = dollar(avg_income, prefix = "EUR ")
    ),
    hjust = -0.2,
    size = 4,
    colour = "white",
    fontface = "bold"
  ) +
  coord_flip() +
  scale_y_continuous(labels = label_dollar()) +
  theme(plot.title = element_text(size = 14, hjust = 0.5)) +
  labs(title = "Average Monthly Income by Gender",
       x = NULL,
       y = NULL)
```



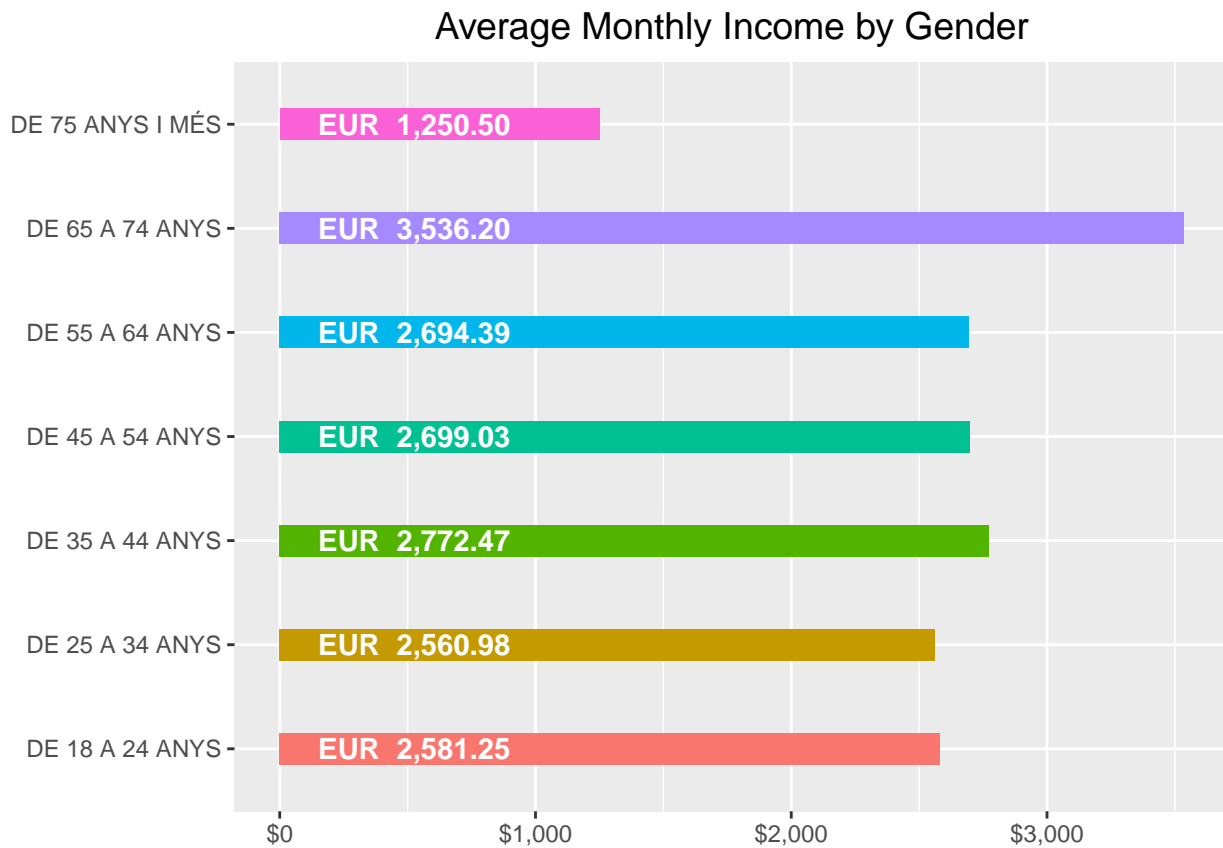
In this case, we can see how in Barcelona Male has generally an higher salary compared to women.

```
barcelona_data %>% ggplot(aes(x = age, y = monthly_income)) +
  geom_col(aes(fill = age), width = 0.3, show.legend = FALSE) +
  coord_flip() +
  scale_y_continuous(labels = label_dollar()) +
  theme(plot.title = element_text(size = 14, hjust = 0.5)) +
  labs(title = "Average Monthly Income by Gender",
       x = NULL,
       y = NULL)
```



```
barcelona_data %>%
  select(age, monthly_income) %>%
  group_by(age) %>%
  summarise(avg_income = round(mean(monthly_income), 2), .groups = "drop") %>%
  ggplot(aes(x = age, y = avg_income)) +
  geom_col(aes(fill = age), width = 0.3, show.legend = FALSE) +
  geom_text(
    aes(
      x = age,
      y = 0.01,
      label = dollar(avg_income, prefix = "EUR ")
    ),
    hjust = -0.2,
    size = 4,
```

```
colour = "white",
fontface = "bold"
) +
coord_flip() +
scale_y_continuous(labels = label_dollar()) +
theme(plot.title = element_text(size = 14, hjust = 0.5)) +
labs(title = "Average Monthly Income by Gender",
x = NULL,
y = NULL)
```



To conclude, we can understand here how the people between 65 and 74 years old have higher monthly_income.

7 Conclusions

In the analysis we reach some conclusion let's take a look on what we reach.

We started by analyzing the frequency of the different role that we have in the dataset, we obtained that looking at the whole data set without differentiating by city the most frequent role are:

- Sales Executive
- Research Scientist
- Laboratory technician

The next step is relate to see if this path is followed also in the two city. The second graph show us that the distribution of the role is the same in both city.

Proceeding forward, we shift on the gender analysis to understand if both company reach the goal of the gender equality.

We first calculate some statistics indicator of the distribution of the monthly income, obtaining an average monthly income of 6510 euro.

We plotted a graph that show the general situation of the whole company, and there we can see that on average the monthly income of the female employee tend to be higher. In fact, the average monthly income for male is 6395.65 while for the female is 6682.85 euro.

The next step will be to analyze the same point but after a distinction between the two different city.

In both the city the monthly income of the female is higher than the male, but the gap in Barcelona is more significant.

Let's shift to the result that take in consideration also the department to which the employee belong.

In general the number of male per department is always higher with respect to the female, and the majority of personnel work in the Research & Development.

As before we will take a look also in the two different city. The gap between the gender per department is very high in London, but in Barcelona the HR department respect the gender equality, so our CEO need to work firstly on the department of Sales and Research & Development.

Let's shift now on the job role, in all the role we have more male than female. From the graph where we don't differentiate between the city, all the role has more male than girl except the Manufacturing director role where the number of male and female are almost the same. We can state the same conclusion also if we analyze only London. In the case of Barcelona we have a slight difference, the female are higher in three role.

- Manager
- Research Director
- Manufacturing director

Now, we want to write down the result obtained once we put in relation the monthly income versus the years in the company differentiating by gender and then by city also.

The whole picture show us that in the initial stage the female tend to get more money with respect to the male, then in the middle stage from 8 years till 20 the income is almost the same, but in the final stage so after 20 years of seniority the gap between female and male increase in the favor of the male.

Looking to the two city, we can state the same things only the duration of the middle stage change, because now the male start earn more money 15 years of seniority. In Barcelona the picture is totally different, female employee tend to earn more almost always.

Let's consider the monthly income in relation to the different department and differentiate also by gender.

In the graph where we consider at the same time both the city we can see that the differences in income are present in the education field of technical degree and Human resource cause the mean of female income with respect to the male is more bigger. Considering the two city the problem are still present but in Barcelona is more bigger.

We carried out the same type of analysis but now focusing on the job role instead of education field. We obtained that in general the income as the same path for both male and female, but in two role which are Research Director and Manager the male tend to earn more than the female. In Barcelona the discrepancy is present for both cause in some role male earn in average more than female and also the contrary happen. While in London only the research director show a big gap in favor of the male.

8 References

1. *M., L. (2004). Moneyball: The art of winning an unfair game. New york: John Wiley and sons.*
2. Wickham et al., (2019). Welcome to the tidyverse. Journal of Open Source Software, 4(43), 1686, <https://doi.org/10.21105/joss.01686>
3. H. Wickham. ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New York, 2016.
4. David Robinson, Alex Hayes and Simon Couch (2021). broom: Convert Statistical Objects into Tidy Tibbles. R package version 0.7.9. <https://CRAN.R-project.org/package=broom>