

# CSV Data Analysis in R - project00- dataset HR-attribution

```
setwd("/Users/sinapc/Desktop/Aidapt/project00/2-R project")  
options(repos = c(CRAN = "https://cloud.r-project.org/"))  
  
install.packages("tidyverse")
```

The downloaded binary packages are in  
/var/folders/hc/kctppyk503b3hrnh27w5qmw0000gn/T//RtmpNwAYWa/downloaded\_packages

```
install.packages("summarytools")
```

The downloaded binary packages are in  
/var/folders/hc/kctppyk503b3hrnh27w5qmw0000gn/T//RtmpNwAYWa/downloaded\_packages

```
install.packages("knitr")
```

The downloaded binary packages are in  
/var/folders/hc/kctppyk503b3hrnh27w5qmw0000gn/T//RtmpNwAYWa/downloaded\_packages

```
install.packages("corrplot")
```

The downloaded binary packages are in  
/var/folders/hc/kctppyk503b3hrnh27w5qmw0000gn/T//RtmpNwAYWa/downloaded\_packages

```
library(tidyverse)
```

```
-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
v dplyr      1.1.4      v readr      2.1.5
v forcats    1.0.0      v stringr    1.5.1
v ggplot2    3.5.2      v tibble     3.2.1
v lubridate  1.9.4      v tidyr      1.3.1
v purrr      1.0.4
-- Conflicts ----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()     masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become
```

```
library(summarytools)
```

Attaching package: 'summarytools'

The following object is masked from 'package:tibble':

view

```
library(knitr)
```

```
data <- read_csv("original.csv")
```

Rows: 1470 Columns: 35

```
-- Column specification -----
```

Delimiter: ","

chr (9): Attrition, BusinessTravel, Department, EducationField, Gender, Job...

dbl (26): Age, DailyRate, DistanceFromHome, Education, EmployeeCount, Employ...

i Use `spec()` to retrieve the full column specification for this data.

i Specify the column types or set `show\_col\_types = FALSE` to quiet this message.

```
#dim(data)
```

```
#head(data, 10)
```

```
glimpse(data)
```

Rows: 1,470

Columns: 35

```
$ Age <dbl> 41, 49, 37, 33, 27, 32, 59, 30, 38, 36, 35, 2~
$ Attrition <chr> "Yes", "No", "Yes", "No", "No", "No", "No", "~
$ BusinessTravel <chr> "Travel_Rarely", "Travel_Frequently", "Travel~
$ DailyRate <dbl> 1102, 279, 1373, 1392, 591, 1005, 1324, 1358,~
$ Department <chr> "Sales", "Research & Development", "Research ~
$ DistanceFromHome <dbl> 1, 8, 2, 3, 2, 2, 3, 24, 23, 27, 16, 15, 26, ~
$ Education <dbl> 2, 1, 2, 4, 1, 2, 3, 1, 3, 3, 3, 2, 1, 2, 3, ~
$ EducationField <chr> "Life Sciences", "Life Sciences", "Other", "L~
$ EmployeeCount <dbl> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
$ EmployeeNumber <dbl> 1, 2, 4, 5, 7, 8, 10, 11, 12, 13, 14, 15, 16,~
$ EnvironmentSatisfaction <dbl> 2, 3, 4, 4, 1, 4, 3, 4, 4, 3, 1, 4, 1, 2, 3, ~
$ Gender <chr> "Female", "Male", "Male", "Female", "Male", "~
$ HourlyRate <dbl> 94, 61, 92, 56, 40, 79, 81, 67, 44, 94, 84, 4~
$ JobInvolvement <dbl> 3, 2, 2, 3, 3, 3, 4, 3, 2, 3, 4, 2, 3, 3, 2, ~
$ JobLevel <dbl> 2, 2, 1, 1, 1, 1, 1, 1, 3, 2, 1, 2, 1, 1, 1, ~
$ JobRole <chr> "Sales Executive", "Research Scientist", "Lab~
$ JobSatisfaction <dbl> 4, 2, 3, 3, 2, 4, 1, 3, 3, 3, 2, 3, 3, 4, 3, ~
$ MaritalStatus <chr> "Single", "Married", "Single", "Married", "Ma~
$ MonthlyIncome <dbl> 5993, 5130, 2090, 2909, 3468, 3068, 2670, 269~
$ MonthlyRate <dbl> 19479, 24907, 2396, 23159, 16632, 11864, 9964~
$ NumCompaniesWorked <dbl> 8, 1, 6, 1, 9, 0, 4, 1, 0, 6, 0, 0, 1, 0, 5, ~
$ Over18 <chr> "Y", "Y", "Y", "Y", "Y", "Y", "Y", "Y", "Y", "~
$ OverTime <chr> "Yes", "No", "Yes", "Yes", "No", "No", "Yes",~
$ PercentSalaryHike <dbl> 11, 23, 15, 11, 12, 13, 20, 22, 21, 13, 13, 1~
$ PerformanceRating <dbl> 3, 4, 3, 3, 3, 3, 4, 4, 4, 3, 3, 3, 3, 3, 3, ~
$ RelationshipSatisfaction <dbl> 1, 4, 2, 3, 4, 3, 1, 2, 2, 2, 3, 4, 4, 3, 2, ~
$ StandardHours <dbl> 80, 80, 80, 80, 80, 80, 80, 80, 80, 80, 80, 8~
$ StockOptionLevel <dbl> 0, 1, 0, 0, 1, 0, 3, 1, 0, 2, 1, 0, 1, 1, 0, ~
$ TotalWorkingYears <dbl> 8, 10, 7, 8, 6, 8, 12, 1, 10, 17, 6, 10, 5, 3~
$ TrainingTimesLastYear <dbl> 0, 3, 3, 3, 3, 2, 3, 2, 2, 3, 5, 3, 1, 2, 4, ~
$ WorkLifeBalance <dbl> 1, 3, 3, 3, 3, 2, 2, 3, 3, 2, 3, 3, 2, 3, 3, ~
$ YearsAtCompany <dbl> 6, 10, 0, 8, 2, 7, 1, 1, 9, 7, 5, 9, 5, 2, 4,~
$ YearsInCurrentRole <dbl> 4, 7, 0, 7, 2, 7, 0, 0, 7, 7, 4, 5, 2, 2, 2, ~
$ YearsSinceLastPromotion <dbl> 0, 1, 0, 3, 2, 3, 0, 0, 1, 7, 0, 0, 4, 1, 0, ~
$ YearsWithCurrManager <dbl> 5, 7, 0, 0, 2, 6, 0, 0, 8, 7, 3, 8, 3, 2, 3, ~
```

```
#str(data)
```

```
summary(data) #descriptive statistics
```

Age

Attrition

BusinessTravel

DailyRate

Min. :18.00	Length:1470	Length:1470	Min. : 102.0
1st Qu.:30.00	Class :character	Class :character	1st Qu.: 465.0
Median :36.00	Mode :character	Mode :character	Median : 802.0
Mean :36.92			Mean : 802.5
3rd Qu.:43.00			3rd Qu.:1157.0
Max. :60.00			Max. :1499.0
Department	DistanceFromHome	Education	EducationField
Length:1470	Min. : 1.000	Min. :1.000	Length:1470
Class :character	1st Qu.: 2.000	1st Qu.:2.000	Class :character
Mode :character	Median : 7.000	Median :3.000	Mode :character
	Mean : 9.193	Mean :2.913	
	3rd Qu.:14.000	3rd Qu.:4.000	
	Max. :29.000	Max. :5.000	
EmployeeCount	EmployeeNumber	EnvironmentSatisfaction	Gender
Min. :1	Min. : 1.0	Min. :1.000	Length:1470
1st Qu.:1	1st Qu.: 491.2	1st Qu.:2.000	Class :character
Median :1	Median :1020.5	Median :3.000	Mode :character
Mean :1	Mean :1024.9	Mean :2.722	
3rd Qu.:1	3rd Qu.:1555.8	3rd Qu.:4.000	
Max. :1	Max. :2068.0	Max. :4.000	
HourlyRate	JobInvolvement	JobLevel	JobRole
Min. : 30.00	Min. :1.00	Min. :1.000	Length:1470
1st Qu.: 48.00	1st Qu.:2.00	1st Qu.:1.000	Class :character
Median : 66.00	Median :3.00	Median :2.000	Mode :character
Mean : 65.89	Mean :2.73	Mean :2.064	
3rd Qu.: 83.75	3rd Qu.:3.00	3rd Qu.:3.000	
Max. :100.00	Max. :4.00	Max. :5.000	
JobSatisfaction	MaritalStatus	MonthlyIncome	MonthlyRate
Min. :1.000	Length:1470	Min. : 1009	Min. : 2094
1st Qu.:2.000	Class :character	1st Qu.: 2911	1st Qu.: 8047
Median :3.000	Mode :character	Median : 4919	Median :14236
Mean :2.729		Mean : 6503	Mean :14313
3rd Qu.:4.000		3rd Qu.: 8379	3rd Qu.:20462
Max. :4.000		Max. :19999	Max. :26999
NumCompaniesWorked	Over18	OverTime	PercentSalaryHike
Min. :0.000	Length:1470	Length:1470	Min. :11.00
1st Qu.:1.000	Class :character	Class :character	1st Qu.:12.00
Median :2.000	Mode :character	Mode :character	Median :14.00
Mean :2.693			Mean :15.21
3rd Qu.:4.000			3rd Qu.:18.00
Max. :9.000			Max. :25.00
PerformanceRating	RelationshipSatisfaction	StandardHours	StockOptionLevel
Min. :3.000	Min. :1.000	Min. :80	Min. :0.0000

1st Qu.:3.000	1st Qu.:2.000	1st Qu.:80	1st Qu.:0.0000
Median :3.000	Median :3.000	Median :80	Median :1.0000
Mean :3.154	Mean :2.712	Mean :80	Mean :0.7939
3rd Qu.:3.000	3rd Qu.:4.000	3rd Qu.:80	3rd Qu.:1.0000
Max. :4.000	Max. :4.000	Max. :80	Max. :3.0000
TotalWorkingYears	TrainingTimesLastYear	WorkLifeBalance	YearsAtCompany
Min. : 0.00	Min. :0.000	Min. :1.000	Min. : 0.000
1st Qu.: 6.00	1st Qu.:2.000	1st Qu.:2.000	1st Qu.: 3.000
Median :10.00	Median :3.000	Median :3.000	Median : 5.000
Mean :11.28	Mean :2.799	Mean :2.761	Mean : 7.008
3rd Qu.:15.00	3rd Qu.:3.000	3rd Qu.:3.000	3rd Qu.: 9.000
Max. :40.00	Max. :6.000	Max. :4.000	Max. :40.000
YearsInCurrentRole	YearsSinceLastPromotion	YearsWithCurrManager	
Min. : 0.000	Min. : 0.000	Min. : 0.000	
1st Qu.: 2.000	1st Qu.: 0.000	1st Qu.: 2.000	
Median : 3.000	Median : 1.000	Median : 3.000	
Mean : 4.229	Mean : 2.188	Mean : 4.123	
3rd Qu.: 7.000	3rd Qu.: 3.000	3rd Qu.: 7.000	
Max. :18.000	Max. :15.000	Max. :17.000	

```
#names(data)
```

```
# categorical value check with Knitr:Kable frames
```

```
data %>% distinct(Attrition) %>% kable()
```

---

Attrition

Yes

No

---

```
data %>% distinct(BusinessTravel) %>% kable()
```

---

BusinessTravel

Travel\_Rarely

Travel\_Frequently

Non-Travel

---

```
data %>% distinct(Department) %>% kable()
```

Department
Sales
Research & Development
Human Resources

```
data %>% distinct(Education) %>% kable()
```

Education
2
1
4
3
5

```
data %>% distinct(EducationField) %>% kable()
```

EducationField
Life Sciences
Other
Medical
Marketing
Technical Degree
Human Resources

```
data %>% distinct(EmployeeCount) %>% kable()
```

EmployeeCount
1

```
data %>% distinct(Gender) %>% kable()
```

Gender
Female
Male

```
data %>% distinct(JobRole) %>% kable()
```

JobRole
Sales Executive
Research Scientist
Laboratory Technician
Manufacturing Director
Healthcare Representative
Manager
Sales Representative
Research Director
Human Resources

```
data %>% distinct(MaritalStatus) %>% kable()
```

MaritalStatus
Single
Married
Divorced

```
data %>% distinct(Over18) %>% kable()
```

Over18
Y

```
data %>% distinct(OverTime) %>% kable()
```

OverTime
Yes
No

```
data %>% distinct(PerformanceRating) %>% kable()
```

PerformanceRating
3
4

```
data %>% distinct(StandardHours) %>% kable()
```

StandardHours
80

```
#missing value check
colSums(is.na(data))
```

Age	Attrition	BusinessTravel
0	0	0
DailyRate	Department	DistanceFromHome
0	0	0
Education	EducationField	EmployeeCount
0	0	0
EmployeeNumber	EnvironmentSatisfaction	Gender
0	0	0
HourlyRate	JobInvolvement	JobLevel
0	0	0
JobRole	JobSatisfaction	MaritalStatus
0	0	0
MonthlyIncome	MonthlyRate	NumCompaniesWorked
0	0	0
Over18	Overtime	PercentSalaryHike
0	0	0
PerformanceRating	RelationshipSatisfaction	StandardHours
0	0	0
StockOptionLevel	TotalWorkingYears	TrainingTimesLastYear
0	0	0
WorkLifeBalance	YearsAtCompany	YearsInCurrentRole
0	0	0
YearsSinceLastPromotion	YearsWithCurrManager	
0	0	



```
#comprehensive overveiw of data from package
dfSummary(data)
```

Data Frame Summary  
data  
Dimensions: 1470 x 35  
Duplicates: 0

No	Variable	Stats / Values	Freqs (% of Valid)	Graph
1	Age [numeric]	Mean (sd) : 36.9 (9.1) min < med < max: 18 < 36 < 60 IQR (CV) : 13 (0.2)	43 distinct values	: : . .
2	Attrition [character]	1. No 2. Yes	1233 (83.9%) 237 (16.1%)	IIIIIIII III
3	BusinessTravel [character]	1. Non-Travel 2. Travel_Frequently 3. Travel_Rarely	150 (10.2%) 277 (18.8%) 1043 (71.0%)	II III IIIIIIII
4	DailyRate [numeric]	Mean (sd) : 802.5 (403.5) min < med < max: 102 < 802 < 1499 IQR (CV) : 692 (0.5)	886 distinct values	. : : : :
5	Department [character]	1. Human Resources 2. Research & Development 3. Sales	63 ( 4.3%) 961 (65.4%) 446 (30.3%)	IIIIIIII IIIIIIII IIIIIIII
6	DistanceFromHome [numeric]	Mean (sd) : 9.2 (8.1) min < med < max: 1 < 7 < 29 IQR (CV) : 12 (0.9)	29 distinct values	: : : : :
7	Education [numeric]	Mean (sd) : 2.9 (1) min < med < max:	1 : 170 (11.6%) 2 : 282 (19.2%)	II III

		1 < 3 < 5 IQR (CV) : 2 (0.4)	3 : 572 (38.9%) 4 : 398 (27.1%) 5 : 48 ( 3.3%)	IIIIIIII IIIIII
8	EducationField [character]	1. Human Resources 2. Life Sciences 3. Marketing 4. Medical 5. Other 6. Technical Degree	27 ( 1.8%) 606 (41.2%) 159 (10.8%) 464 (31.6%) 82 ( 5.6%) 132 ( 9.0%)	IIIIIIII II IIIIIIII I I
9	EmployeeCount [numeric]	1 distinct value	1 : 1470 (100.0%)	IIIIIIII
10	EmployeeNumber [numeric]	Mean (sd) : 1024.9 (602) min < med < max: 1 < 1020.5 < 2068 IQR (CV) : 1064.5 (0.6)	1470 distinct values	: . : : : : : : : : : : : : :
11	EnvironmentSatisfaction [numeric]	Mean (sd) : 2.7 (1.1) min < med < max: 1 < 3 < 4 IQR (CV) : 2 (0.4)	1 : 284 (19.3%) 2 : 287 (19.5%) 3 : 453 (30.8%) 4 : 446 (30.3%)	III III IIIIII IIIIII
12	Gender [character]	1. Female 2. Male	588 (40.0%) 882 (60.0%)	IIIIIIII IIIIIIII
13	HourlyRate [numeric]	Mean (sd) : 65.9 (20.3) min < med < max: 30 < 66 < 100 IQR (CV) : 35.8 (0.3)	71 distinct values	. . . : : : : : : : : : : : :
14	JobInvolvement [numeric]	Mean (sd) : 2.7 (0.7) min < med < max: 1 < 3 < 4 IQR (CV) : 1 (0.3)	1 : 83 ( 5.6%) 2 : 375 (25.5%) 3 : 868 (59.0%) 4 : 144 ( 9.8%)	I IIIIII IIIIIIII I
15	JobLevel [numeric]	Mean (sd) : 2.1 (1.1) min < med < max: 1 < 2 < 5 IQR (CV) : 2 (0.5)	1 : 543 (36.9%) 2 : 534 (36.3%) 3 : 218 (14.8%) 4 : 106 ( 7.2%)	IIIIIIII IIIIIIII II I

			5 : 69 ( 4.7%)	
16	JobRole [character]	1. Healthcare Representative 2. Human Resources 3. Laboratory Technician 4. Manager 5. Manufacturing Director 6. Research Director 7. Research Scientist 8. Sales Executive 9. Sales Representative	131 ( 8.9%) 52 ( 3.5%) 259 (17.6%) 102 ( 6.9%) 145 ( 9.9%) 80 ( 5.4%) 292 (19.9%) 326 (22.2%) 83 ( 5.6%)	I  III I I I III IIII I
17	JobSatisfaction [numeric]	Mean (sd) : 2.7 (1.1) min < med < max: 1 < 3 < 4 IQR (CV) : 2 (0.4)	1 : 289 (19.7%) 2 : 280 (19.0%) 3 : 442 (30.1%) 4 : 459 (31.2%)	III III IIIIII IIIIII
18	MaritalStatus [character]	1. Divorced 2. Married 3. Single	327 (22.2%) 673 (45.8%) 470 (32.0%)	IIII IIIIII IIIIII
19	MonthlyIncome [numeric]	Mean (sd) : 6502.9 (4708) min < med < max: 1009 < 4919 < 19999 IQR (CV) : 5468 (0.7)	1349 distinct values	: : : : : : : : : : : : : :
20	MonthlyRate [numeric]	Mean (sd) : 14313.1 (7117.8) min < med < max: 2094 < 14235.5 < 26999 IQR (CV) : 12414.5 (0.5)	1427 distinct values	. : : : : : : : : : : : : : :
21	NumCompaniesWorked [numeric]	Mean (sd) : 2.7 (2.5) min < med < max: 0 < 2 < 9 IQR (CV) : 3 (0.9)	0 : 197 (13.4%) 1 : 521 (35.4%) 2 : 146 ( 9.9%) 3 : 159 (10.8%) 4 : 139 ( 9.5%) 5 : 63 ( 4.3%) 6 : 70 ( 4.8%) 7 : 74 ( 5.0%) 8 : 49 ( 3.3%) 9 : 52 ( 3.5%)	II IIIIII I II I   I

22	Over18 [character]	1. Y	1470 (100.0%)	IIIIII
23	OverTime [character]	1. No 2. Yes	1054 (71.7%) 416 (28.3%)	IIIIII IIII
24	PercentSalaryHike [numeric]	Mean (sd) : 15.2 (3.7) min < med < max: 11 < 14 < 25 IQR (CV) : 6 (0.2)	15 distinct values	: : : . : : : : : : :
25	PerformanceRating [numeric]	Min : 3 Mean : 3.2 Max : 4	3 : 1244 (84.6%) 4 : 226 (15.4%)	IIIIII III
26	RelationshipSatisfaction [numeric]	Mean (sd) : 2.7 (1.1) min < med < max: 1 < 3 < 4 IQR (CV) : 2 (0.4)	1 : 276 (18.8%) 2 : 303 (20.6%) 3 : 459 (31.2%) 4 : 432 (29.4%)	III IIII IIIIII IIII
27	StandardHours [numeric]	1 distinct value	80 : 1470 (100.0%)	IIIIII
28	StockOptionLevel [numeric]	Mean (sd) : 0.8 (0.9) min < med < max: 0 < 1 < 3 IQR (CV) : 1 (1.1)	0 : 631 (42.9%) 1 : 596 (40.5%) 2 : 158 (10.7%) 3 : 85 ( 5.8%)	IIIIII IIIIII II I
29	TotalWorkingYears [numeric]	Mean (sd) : 11.3 (7.8) min < med < max: 0 < 10 < 40 IQR (CV) : 9 (0.7)	40 distinct values	: : : : : : . : : :
30	TrainingTimesLastYear [numeric]	Mean (sd) : 2.8 (1.3) min < med < max: 0 < 3 < 6 IQR (CV) : 1 (0.5)	0 : 54 ( 3.7%) 1 : 71 ( 4.8%) 2 : 547 (37.2%) 3 : 491 (33.4%) 4 : 123 ( 8.4%) 5 : 119 ( 8.1%) 6 : 65 ( 4.4%)	IIIIII IIIIII IIIIII IIIIII I I

31	WorkLifeBalance [numeric]	Mean (sd) : 2.8 (0.7) min < med < max: 1 < 3 < 4 IQR (CV) : 1 (0.3)	1 : 80 ( 5.4%) 2 : 344 (23.4%) 3 : 893 (60.7%) 4 : 153 (10.4%)	I IIII IIIIII II
32	YearsAtCompany [numeric]	Mean (sd) : 7 (6.1) min < med < max: 0 < 5 < 40 IQR (CV) : 6 (0.9)	37 distinct values	: : : : : : : : .
33	YearsInCurrentRole [numeric]	Mean (sd) : 4.2 (3.6) min < med < max: 0 < 3 < 18 IQR (CV) : 5 (0.9)	19 distinct values	: : : : : : : .
34	YearsSinceLastPromotion [numeric]	Mean (sd) : 2.2 (3.2) min < med < max: 0 < 1 < 15 IQR (CV) : 3 (1.5)	16 distinct values	: : : : : : .
35	YearsWithCurrManager [numeric]	Mean (sd) : 4.1 (3.6) min < med < max: 0 < 3 < 17 IQR (CV) : 5 (0.9)	18 distinct values	: : : : : : : .

```

#mean(data$MonthlyIncome, na.rm = TRUE)

# Summary of all numeric columns
summary_table <- data %>%
  summarise(across(where(is.numeric),list(mean = mean,median = median, sd = sd), na.rm = TRUE))

```

```

Warning: There was 1 warning in `summarise()`.
i In argument: `across(...)` .
Caused by warning:
! The `...` argument of `across()` is deprecated as of dplyr 1.1.0.
Supply arguments directly to `.fns` through an anonymous function instead.

```

```

# Previously
across(a:b, mean, na.rm = TRUE)

# Now
across(a:b, \(x) mean(x, na.rm = TRUE))

# write a summary table in format_csv()
summary_table %>%
  summarise(across(where(is.numeric), list(mean = mean, median = median, sd = sd), na.rm = TRUE))

# A tibble: 1 x 234
  Age_mean_mean Age_mean_median Age_mean_sd Age_median_mean Age_median_median
      <dbl>         <dbl>         <dbl>         <dbl>         <dbl>
1      36.9         36.9           NA           36           36
# i 229 more variables: Age_median_sd <dbl>, Age_sd_mean <dbl>,
#   Age_sd_median <dbl>, Age_sd_sd <dbl>, DailyRate_mean_mean <dbl>,
#   DailyRate_mean_median <dbl>, DailyRate_mean_sd <dbl>,
#   DailyRate_median_mean <dbl>, DailyRate_median_median <dbl>,
#   DailyRate_median_sd <dbl>, DailyRate_sd_mean <dbl>,
#   DailyRate_sd_median <dbl>, DailyRate_sd_sd <dbl>,
#   DistanceFromHome_mean_mean <dbl>, DistanceFromHome_mean_median <dbl>, ...

write_csv(summary_table, "summary_table.csv")

# Attrition rate by Job Role
data %>%
  group_by(JobRole, Attrition) %>%
  summarise(Count = n()) %>%
  mutate(Percent = round(Count / sum(Count) * 100, 1))

```

`summarise()` has grouped output by 'JobRole'. You can override using the `.groups` argument.

```

# A tibble: 18 x 4
# Groups:   JobRole [9]
  JobRole      Attrition Count Percent
  <chr>         <chr>    <int>   <dbl>
1 Healthcare Representative No      122    93.1
2 Healthcare Representative Yes        9     6.9

```

3	Human Resources	No	40	76.9
4	Human Resources	Yes	12	23.1
5	Laboratory Technician	No	197	76.1
6	Laboratory Technician	Yes	62	23.9
7	Manager	No	97	95.1
8	Manager	Yes	5	4.9
9	Manufacturing Director	No	135	93.1
10	Manufacturing Director	Yes	10	6.9
11	Research Director	No	78	97.5
12	Research Director	Yes	2	2.5
13	Research Scientist	No	245	83.9
14	Research Scientist	Yes	47	16.1
15	Sales Executive	No	269	82.5
16	Sales Executive	Yes	57	17.5
17	Sales Representative	No	50	60.2
18	Sales Representative	Yes	33	39.8

```
# Select numeric variables
num_data <- data %>% select(where(is.numeric))

# Correlation matrix
cor_matrix <- cor(num_data)
```

Warning in cor(num\_data): the standard deviation is zero

```
round(cor_matrix, 2)
```

	Age	DailyRate	DistanceFromHome	Education
Age	1.00	0.01	0.00	0.21
DailyRate	0.01	1.00	0.00	-0.02
DistanceFromHome	0.00	0.00	1.00	0.02
Education	0.21	-0.02	0.02	1.00
EmployeeCount	NA	NA	NA	NA
EmployeeNumber	-0.01	-0.05	0.03	0.04
EnvironmentSatisfaction	0.01	0.02	-0.02	-0.03
HourlyRate	0.02	0.02	0.03	0.02
JobInvolvement	0.03	0.05	0.01	0.04
JobLevel	0.51	0.00	0.01	0.10
JobSatisfaction	0.00	0.03	0.00	-0.01
MonthlyIncome	0.50	0.01	-0.02	0.09
MonthlyRate	0.03	-0.03	0.03	-0.03

NumCompaniesWorked	0.30	0.04	-0.03	0.13
PercentSalaryHike	0.00	0.02	0.04	-0.01
PerformanceRating	0.00	0.00	0.03	-0.02
RelationshipSatisfaction	0.05	0.01	0.01	-0.01
StandardHours	NA	NA	NA	NA
StockOptionLevel	0.04	0.04	0.04	0.02
TotalWorkingYears	0.68	0.01	0.00	0.15
TrainingTimesLastYear	-0.02	0.00	-0.04	-0.03
WorkLifeBalance	-0.02	-0.04	-0.03	0.01
YearsAtCompany	0.31	-0.03	0.01	0.07
YearsInCurrentRole	0.21	0.01	0.02	0.06
YearsSinceLastPromotion	0.22	-0.03	0.01	0.05
YearsWithCurrManager	0.20	-0.03	0.01	0.07
	EmployeeCount	EmployeeNumber	EnvironmentSatisfaction	
Age	NA	-0.01	0.01	
DailyRate	NA	-0.05	0.02	
DistanceFromHome	NA	0.03	-0.02	
Education	NA	0.04	-0.03	
EmployeeCount	1	NA	NA	
EmployeeNumber	NA	1.00	0.02	
EnvironmentSatisfaction	NA	0.02	1.00	
HourlyRate	NA	0.04	-0.05	
JobInvolvement	NA	-0.01	-0.01	
JobLevel	NA	-0.02	0.00	
JobSatisfaction	NA	-0.05	-0.01	
MonthlyIncome	NA	-0.01	-0.01	
MonthlyRate	NA	0.01	0.04	
NumCompaniesWorked	NA	0.00	0.01	
PercentSalaryHike	NA	-0.01	-0.03	
PerformanceRating	NA	-0.02	-0.03	
RelationshipSatisfaction	NA	-0.07	0.01	
StandardHours	NA	NA	NA	
StockOptionLevel	NA	0.06	0.00	
TotalWorkingYears	NA	-0.01	0.00	
TrainingTimesLastYear	NA	0.02	-0.02	
WorkLifeBalance	NA	0.01	0.03	
YearsAtCompany	NA	-0.01	0.00	
YearsInCurrentRole	NA	-0.01	0.02	
YearsSinceLastPromotion	NA	-0.01	0.02	
YearsWithCurrManager	NA	-0.01	0.00	
	HourlyRate	JobInvolvement	JobLevel	JobSatisfaction
Age	0.02	0.03	0.51	0.00
DailyRate	0.02	0.05	0.00	0.03



DistanceFromHome	0.03	0.01	0.01	0.00
Education	0.02	0.04	0.10	-0.01
EmployeeCount	NA	NA	NA	NA
EmployeeNumber	0.04	-0.01	-0.02	-0.05
EnvironmentSatisfaction	-0.05	-0.01	0.00	-0.01
HourlyRate	1.00	0.04	-0.03	-0.07
JobInvolvement	0.04	1.00	-0.01	-0.02
JobLevel	-0.03	-0.01	1.00	0.00
JobSatisfaction	-0.07	-0.02	0.00	1.00
MonthlyIncome	-0.02	-0.02	0.95	-0.01
MonthlyRate	-0.02	-0.02	0.04	0.00
NumCompaniesWorked	0.02	0.02	0.14	-0.06
PercentSalaryHike	-0.01	-0.02	-0.03	0.02
PerformanceRating	0.00	-0.03	-0.02	0.00
RelationshipSatisfaction	0.00	0.03	0.02	-0.01
StandardHours	NA	NA	NA	NA
StockOptionLevel	0.05	0.02	0.01	0.01
TotalWorkingYears	0.00	-0.01	0.78	-0.02
TrainingTimesLastYear	-0.01	-0.02	-0.02	-0.01
WorkLifeBalance	0.00	-0.01	0.04	-0.02
YearsAtCompany	-0.02	-0.02	0.53	0.00
YearsInCurrentRole	-0.02	0.01	0.39	0.00
YearsSinceLastPromotion	-0.03	-0.02	0.35	-0.02
YearsWithCurrManager	-0.02	0.03	0.38	-0.03

	MonthlyIncome	MonthlyRate	NumCompaniesWorked
Age	0.50	0.03	0.30
DailyRate	0.01	-0.03	0.04
DistanceFromHome	-0.02	0.03	-0.03
Education	0.09	-0.03	0.13
EmployeeCount	NA	NA	NA
EmployeeNumber	-0.01	0.01	0.00
EnvironmentSatisfaction	-0.01	0.04	0.01
HourlyRate	-0.02	-0.02	0.02
JobInvolvement	-0.02	-0.02	0.02
JobLevel	0.95	0.04	0.14
JobSatisfaction	-0.01	0.00	-0.06
MonthlyIncome	1.00	0.03	0.15
MonthlyRate	0.03	1.00	0.02
NumCompaniesWorked	0.15	0.02	1.00
PercentSalaryHike	-0.03	-0.01	-0.01
PerformanceRating	-0.02	-0.01	-0.01
RelationshipSatisfaction	0.03	0.00	0.05
StandardHours	NA	NA	NA

StockOptionLevel	0.01	-0.03	0.03
TotalWorkingYears	0.77	0.03	0.24
TrainingTimesLastYear	-0.02	0.00	-0.07
WorkLifeBalance	0.03	0.01	-0.01
YearsAtCompany	0.51	-0.02	-0.12
YearsInCurrentRole	0.36	-0.01	-0.09
YearsSinceLastPromotion	0.34	0.00	-0.04
YearsWithCurrManager	0.34	-0.04	-0.11

	PercentSalaryHike	PerformanceRating
Age	0.00	0.00
DailyRate	0.02	0.00
DistanceFromHome	0.04	0.03
Education	-0.01	-0.02
EmployeeCount	NA	NA
EmployeeNumber	-0.01	-0.02
EnvironmentSatisfaction	-0.03	-0.03
HourlyRate	-0.01	0.00
JobInvolvement	-0.02	-0.03
JobLevel	-0.03	-0.02
JobSatisfaction	0.02	0.00
MonthlyIncome	-0.03	-0.02
MonthlyRate	-0.01	-0.01
NumCompaniesWorked	-0.01	-0.01
PercentSalaryHike	1.00	0.77
PerformanceRating	0.77	1.00
RelationshipSatisfaction	-0.04	-0.03
StandardHours	NA	NA
StockOptionLevel	0.01	0.00
TotalWorkingYears	-0.02	0.01
TrainingTimesLastYear	-0.01	-0.02
WorkLifeBalance	0.00	0.00
YearsAtCompany	-0.04	0.00
YearsInCurrentRole	0.00	0.03
YearsSinceLastPromotion	-0.02	0.02
YearsWithCurrManager	-0.01	0.02

	RelationshipSatisfaction	StandardHours
Age	0.05	NA
DailyRate	0.01	NA
DistanceFromHome	0.01	NA
Education	-0.01	NA
EmployeeCount	NA	NA
EmployeeNumber	-0.07	NA
EnvironmentSatisfaction	0.01	NA

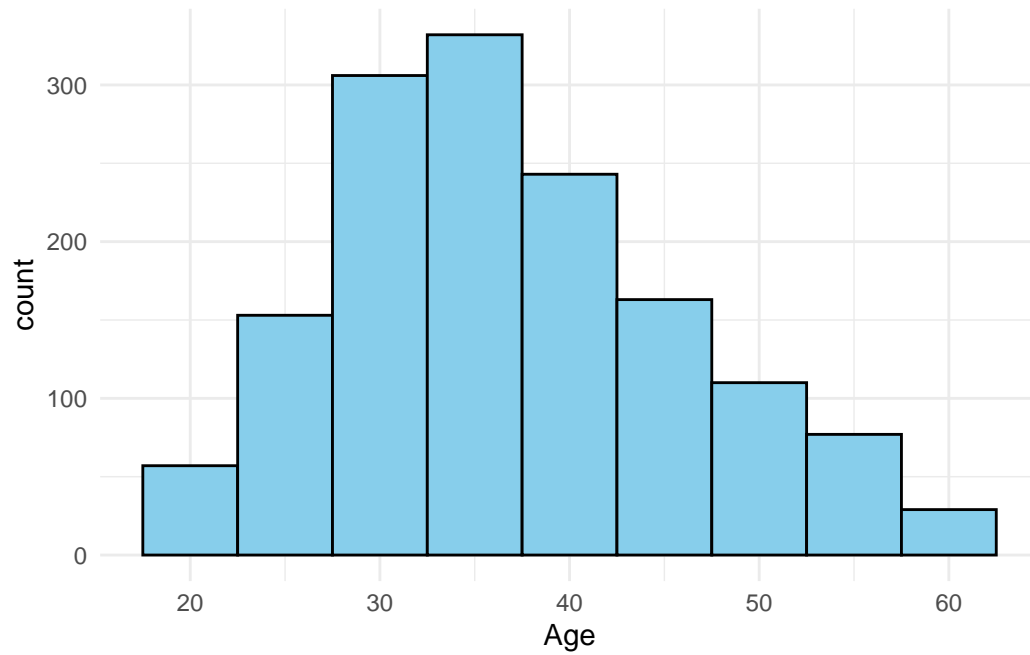
HourlyRate	0.00	NA
JobInvolvement	0.03	NA
JobLevel	0.02	NA
JobSatisfaction	-0.01	NA
MonthlyIncome	0.03	NA
MonthlyRate	0.00	NA
NumCompaniesWorked	0.05	NA
PercentSalaryHike	-0.04	NA
PerformanceRating	-0.03	NA
RelationshipSatisfaction	1.00	NA
StandardHours	NA	1
StockOptionLevel	-0.05	NA
TotalWorkingYears	0.02	NA
TrainingTimesLastYear	0.00	NA
WorkLifeBalance	0.02	NA
YearsAtCompany	0.02	NA
YearsInCurrentRole	-0.02	NA
YearsSinceLastPromotion	0.03	NA
YearsWithCurrManager	0.00	NA

	StockOptionLevel	TotalWorkingYears
Age	0.04	0.68
DailyRate	0.04	0.01
DistanceFromHome	0.04	0.00
Education	0.02	0.15
EmployeeCount	NA	NA
EmployeeNumber	0.06	-0.01
EnvironmentSatisfaction	0.00	0.00
HourlyRate	0.05	0.00
JobInvolvement	0.02	-0.01
JobLevel	0.01	0.78
JobSatisfaction	0.01	-0.02
MonthlyIncome	0.01	0.77
MonthlyRate	-0.03	0.03
NumCompaniesWorked	0.03	0.24
PercentSalaryHike	0.01	-0.02
PerformanceRating	0.00	0.01
RelationshipSatisfaction	-0.05	0.02
StandardHours	NA	NA
StockOptionLevel	1.00	0.01
TotalWorkingYears	0.01	1.00
TrainingTimesLastYear	0.01	-0.04
WorkLifeBalance	0.00	0.00
YearsAtCompany	0.02	0.63

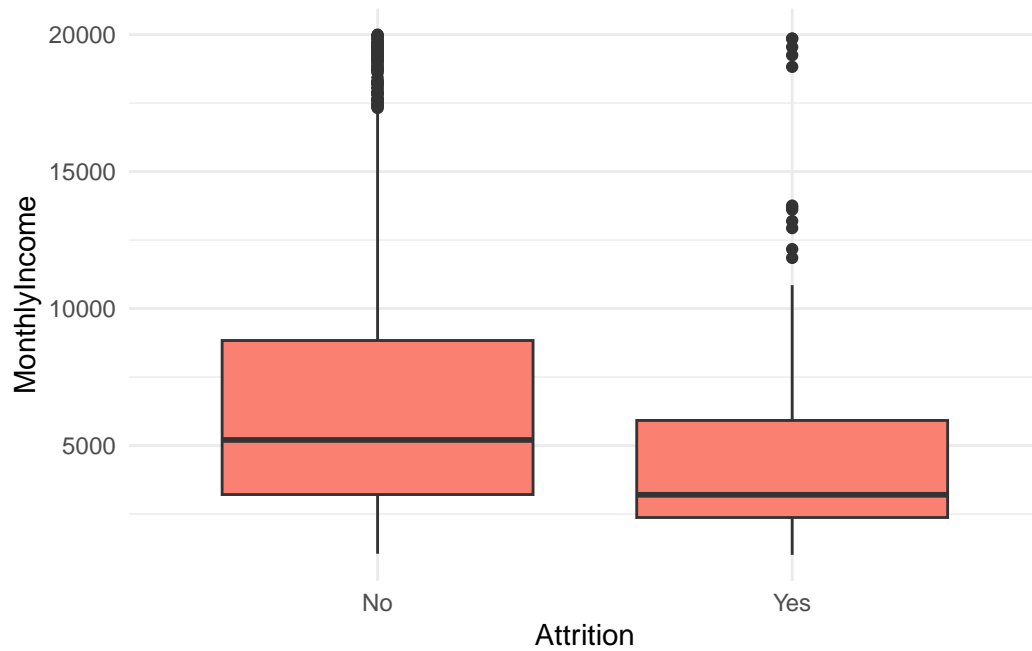
YearsInCurrentRole	0.05	0.46	
YearsSinceLastPromotion	0.01	0.40	
YearsWithCurrManager	0.02	0.46	
	TrainingTimesLastYear	WorkLifeBalance	YearsAtCompany
Age	-0.02	-0.02	0.31
DailyRate	0.00	-0.04	-0.03
DistanceFromHome	-0.04	-0.03	0.01
Education	-0.03	0.01	0.07
EmployeeCount	NA	NA	NA
EmployeeNumber	0.02	0.01	-0.01
EnvironmentSatisfaction	-0.02	0.03	0.00
HourlyRate	-0.01	0.00	-0.02
JobInvolvement	-0.02	-0.01	-0.02
JobLevel	-0.02	0.04	0.53
JobSatisfaction	-0.01	-0.02	0.00
MonthlyIncome	-0.02	0.03	0.51
MonthlyRate	0.00	0.01	-0.02
NumCompaniesWorked	-0.07	-0.01	-0.12
PercentSalaryHike	-0.01	0.00	-0.04
PerformanceRating	-0.02	0.00	0.00
RelationshipSatisfaction	0.00	0.02	0.02
StandardHours	NA	NA	NA
StockOptionLevel	0.01	0.00	0.02
TotalWorkingYears	-0.04	0.00	0.63
TrainingTimesLastYear	1.00	0.03	0.00
WorkLifeBalance	0.03	1.00	0.01
YearsAtCompany	0.00	0.01	1.00
YearsInCurrentRole	-0.01	0.05	0.76
YearsSinceLastPromotion	0.00	0.01	0.62
YearsWithCurrManager	0.00	0.00	0.77
	YearsInCurrentRole	YearsSinceLastPromotion	
Age	0.21	0.22	
DailyRate	0.01	-0.03	
DistanceFromHome	0.02	0.01	
Education	0.06	0.05	
EmployeeCount	NA	NA	
EmployeeNumber	-0.01	-0.01	
EnvironmentSatisfaction	0.02	0.02	
HourlyRate	-0.02	-0.03	
JobInvolvement	0.01	-0.02	
JobLevel	0.39	0.35	
JobSatisfaction	0.00	-0.02	
MonthlyIncome	0.36	0.34	

MonthlyRate	-0.01	0.00
NumCompaniesWorked	-0.09	-0.04
PercentSalaryHike	0.00	-0.02
PerformanceRating	0.03	0.02
RelationshipSatisfaction	-0.02	0.03
StandardHours	NA	NA
StockOptionLevel	0.05	0.01
TotalWorkingYears	0.46	0.40
TrainingTimesLastYear	-0.01	0.00
WorkLifeBalance	0.05	0.01
YearsAtCompany	0.76	0.62
YearsInCurrentRole	1.00	0.55
YearsSinceLastPromotion	0.55	1.00
YearsWithCurrManager	0.71	0.51
	YearsWithCurrManager	
Age	0.20	
DailyRate	-0.03	
DistanceFromHome	0.01	
Education	0.07	
EmployeeCount	NA	
EmployeeNumber	-0.01	
EnvironmentSatisfaction	0.00	
HourlyRate	-0.02	
JobInvolvement	0.03	
JobLevel	0.38	
JobSatisfaction	-0.03	
MonthlyIncome	0.34	
MonthlyRate	-0.04	
NumCompaniesWorked	-0.11	
PercentSalaryHike	-0.01	
PerformanceRating	0.02	
RelationshipSatisfaction	0.00	
StandardHours	NA	
StockOptionLevel	0.02	
TotalWorkingYears	0.46	
TrainingTimesLastYear	0.00	
WorkLifeBalance	0.00	
YearsAtCompany	0.77	
YearsInCurrentRole	0.71	
YearsSinceLastPromotion	0.51	
YearsWithCurrManager	1.00	

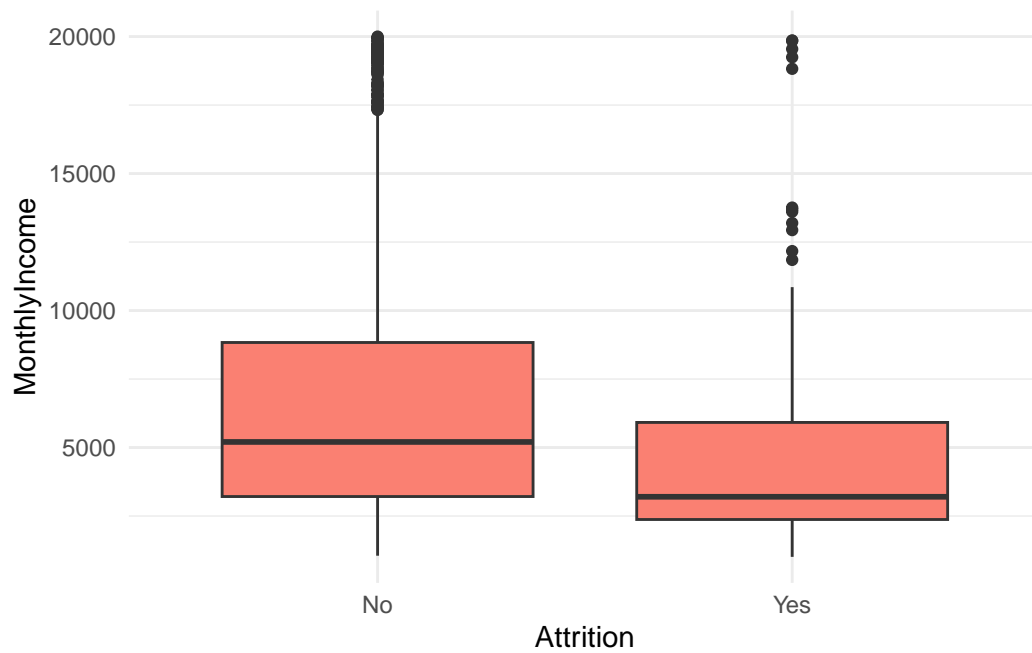
```
ggplot(data, aes(x = Age)) +  
  geom_histogram(binwidth = 5, fill = "skyblue", color = "black") +  
  theme_minimal()
```



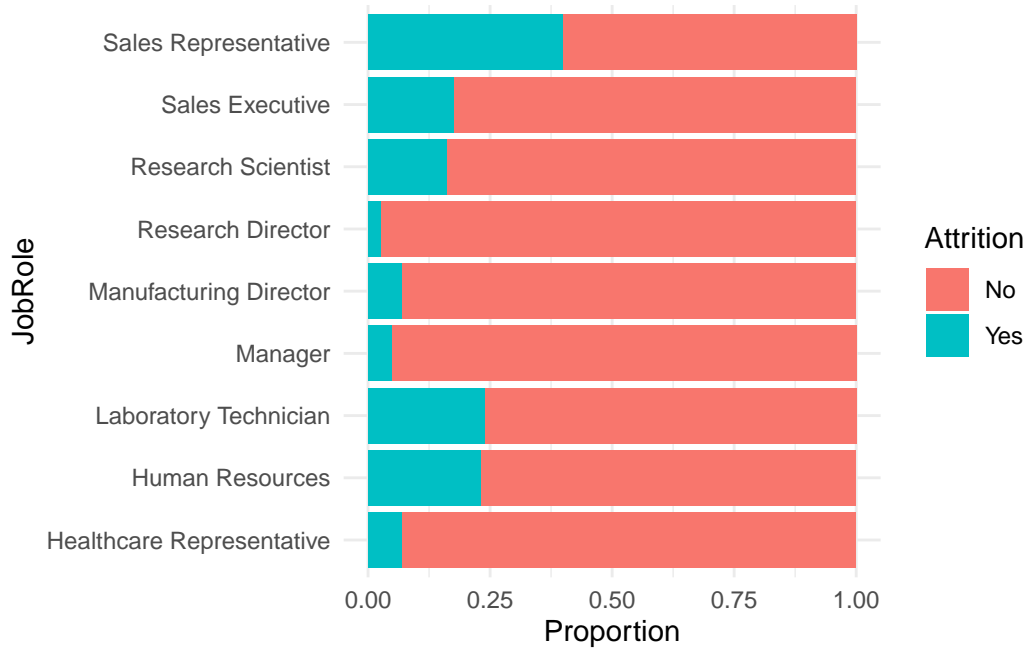
```
ggplot(data, aes(x = Attrition, y = MonthlyIncome)) +  
  geom_boxplot(fill = "salmon") +  
  theme_minimal()
```



```
ggplot(data, aes(x = Attrition, y = MonthlyIncome)) +  
  geom_boxplot(fill = "salmon") +  
  theme_minimal()
```



```
ggplot(data, aes(x = JobRole, fill = Attrition)) +
  geom_bar(position = "fill") +
  coord_flip() +
  labs(y = "Proportion") +
  theme_minimal()
```

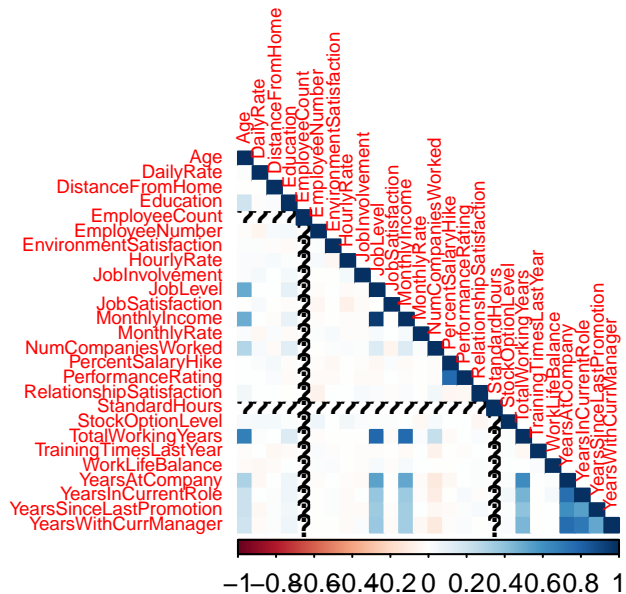


```
library(corrplot)
```

corrplot 0.95 loaded

```
corrplot(cor_matrix, method = "color", type = "lower", tl.cex = 0.6)
```





```
# Convert Attrition to binary
data$Attrition <- ifelse(data$Attrition == "Yes", 1, 0)

# Basic model
model <- glm(Attrition ~ Age + JobRole + MonthlyIncome + OverTime + YearsAtCompany,
data = data, family = "binomial")

summary(model)
```

Call:

```
glm(formula = Attrition ~ Age + JobRole + MonthlyIncome + OverTime +
    YearsAtCompany, family = "binomial", data = data)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )	
(Intercept)	-1.994e+00	5.549e-01	-3.594	0.000326	***
Age	-3.260e-02	1.021e-02	-3.194	0.001403	**
JobRoleHuman Resources	1.465e+00	5.146e-01	2.847	0.004414	**
JobRoleLaboratory Technician	1.519e+00	4.255e-01	3.570	0.000357	***
JobRoleManager	-3.464e-01	7.151e-01	-0.484	0.628117	
JobRoleManufacturing Director	-2.418e-02	4.871e-01	-0.050	0.960414	

```

JobRoleResearch Director      -1.264e+00  8.871e-01  -1.425  0.154156
JobRoleResearch Scientist     8.113e-01  4.302e-01   1.886  0.059293 .
JobRoleSales Executive        1.006e+00  3.864e-01   2.605  0.009198 **
JobRoleSales Representative    2.103e+00  4.744e-01   4.433  9.27e-06 ***
MonthlyIncome                 4.592e-05  4.667e-05   0.984  0.325156
OverTimeYes                   1.485e+00  1.576e-01   9.423  < 2e-16 ***
YearsAtCompany                -3.891e-02  1.869e-02  -2.081  0.037395 *

```

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 1298.6 on 1469 degrees of freedom  
Residual deviance: 1102.9 on 1457 degrees of freedom  
AIC: 1128.9

Number of Fisher Scoring iterations: 6

```

# Calculate attrition percentages by BusinessTravel
attrition_by_businesstravel <- data %>%
  group_by(BusinessTravel, Attrition) %>%
  summarise(Count = n(), .groups = 'drop') %>%
  group_by(BusinessTravel) %>%
  mutate(Percentage = round(Count / sum(Count) * 100, 1))

# Display table
attrition_by_businesstravel %>% kable()

```

BusinessTravel	Attrition	Count	Percentage
Non-Travel	0	138	92.0
Non-Travel	1	12	8.0
Travel_Frequently	0	208	75.1
Travel_Frequently	1	69	24.9
Travel_Rarely	0	887	85.0
Travel_Rarely	1	156	15.0

```

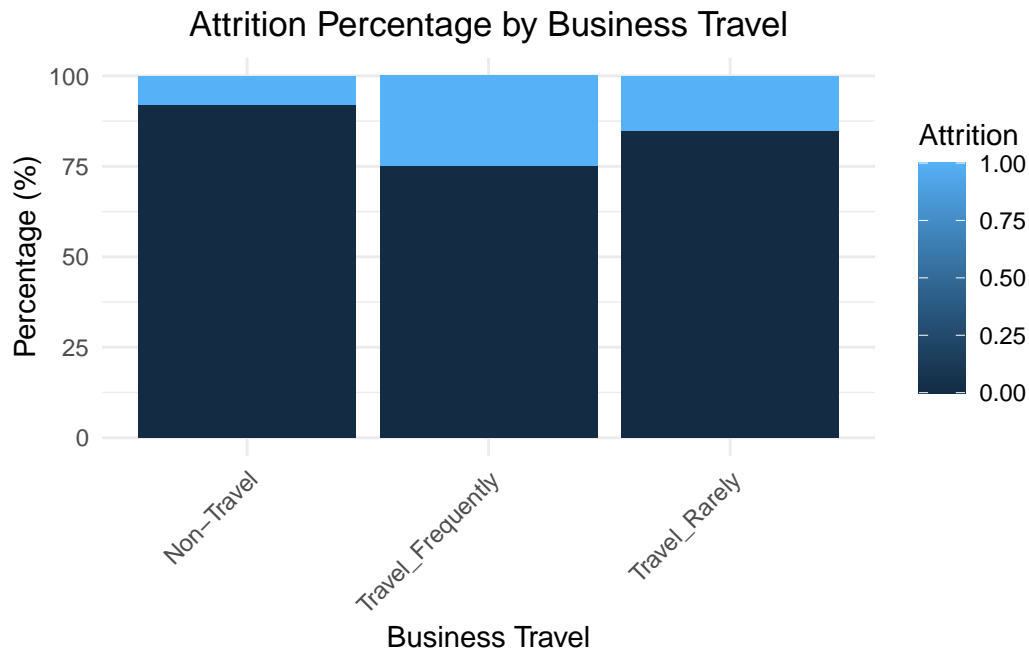
# Visualize attrition by BusinessTravel
ggplot(attrition_by_businesstravel, aes(x = BusinessTravel, y = Percentage, fill = Attrition)) +
  geom_bar(stat = "identity", position = "stack") +
  labs(title = "Attrition Percentage by Business Travel",
       x = "Business Travel",

```

```

y = "Percentage (%)" +
theme_minimal() +
theme(axis.text.x = element_text(angle = 45, hjust = 1), # Rotate labels for readability
      plot.title = element_text(hjust = 0.5))

```



```

# Calculate attrition percentages by Department
attrition_by_department <- data %>%
  group_by(Department, Attrition) %>%
  summarise(Count = n(), .groups = 'drop') %>%
  group_by(Department) %>%
  mutate(Percentage = round(Count / sum(Count) * 100, 1))

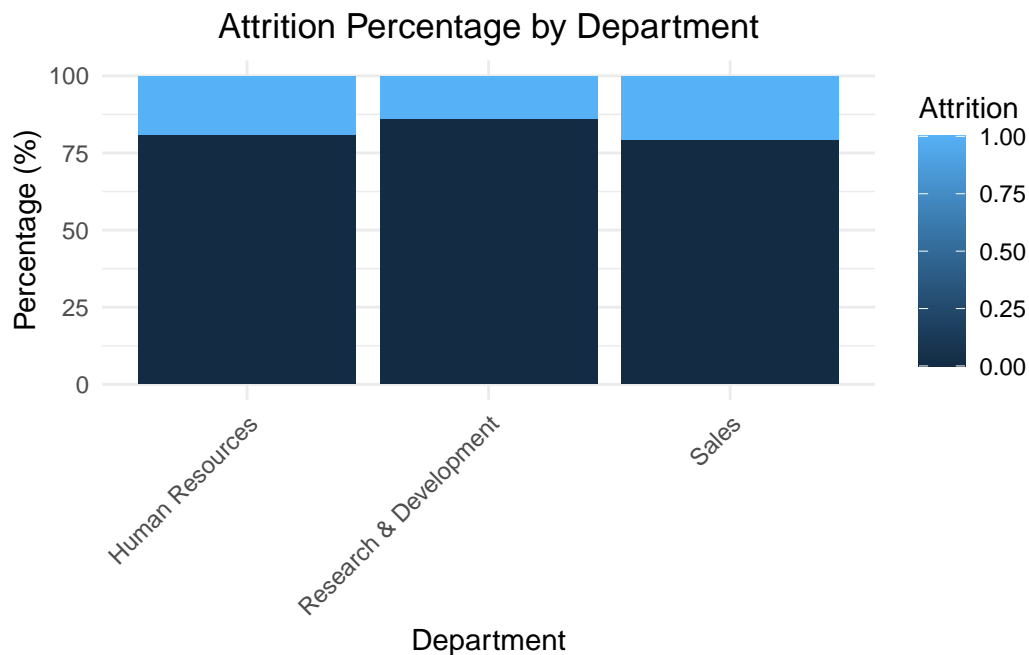
# Display table
attrition_by_department %>% kable()

```

Department	Attrition	Count	Percentage
Human Resources	0	51	81.0
Human Resources	1	12	19.0
Research & Development	0	828	86.2
Research & Development	1	133	13.8
Sales	0	354	79.4

Department	Attrition	Count	Percentage
Sales	1	92	20.6

```
# Visualize attrition by Department
ggplot(attrition_by_department, aes(x = Department, y = Percentage, fill = Attrition)) +
  geom_bar(stat = "identity", position = "stack") +
  labs(title = "Attrition Percentage by Department",
       x = "Department",
       y = "Percentage (%)") +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 45, hjust = 1),
        plot.title = element_text(hjust = 0.5))
```

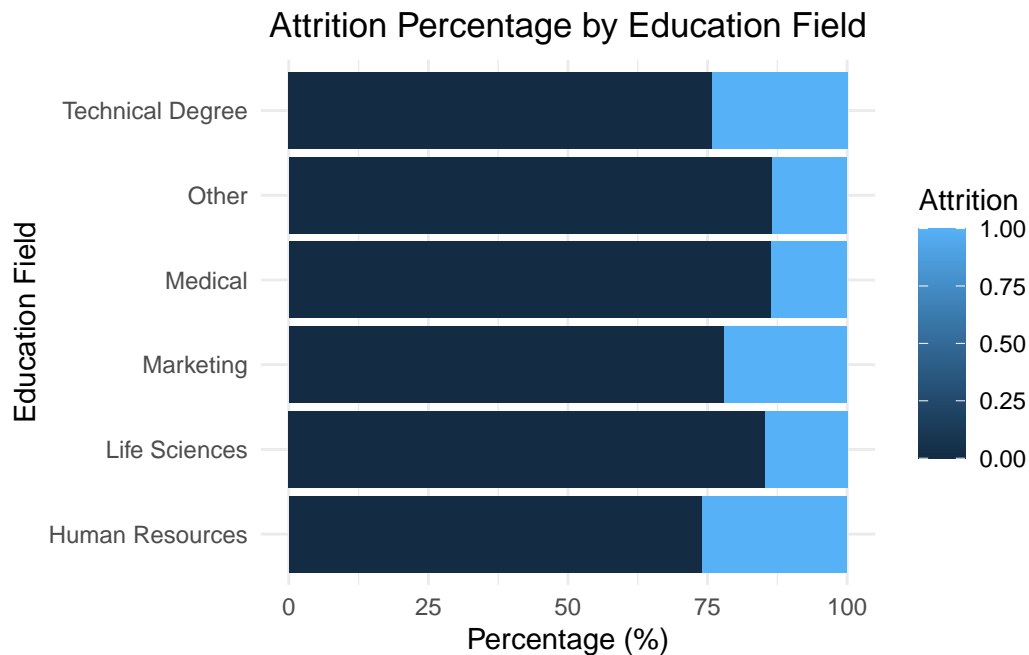


```
# Calculate attrition percentages by EducationField
attrition_by_educationfield <- data %>%
  group_by(EducationField, Attrition) %>%
  summarise(Count = n(), .groups = 'drop') %>%
  group_by(EducationField) %>%
  mutate(Percentage = round(Count / sum(Count) * 100, 1))

# Display table
attrition_by_educationfield %>% kable()
```

EducationField	Attrition	Count	Percentage
Human Resources	0	20	74.1
Human Resources	1	7	25.9
Life Sciences	0	517	85.3
Life Sciences	1	89	14.7
Marketing	0	124	78.0
Marketing	1	35	22.0
Medical	0	401	86.4
Medical	1	63	13.6
Other	0	71	86.6
Other	1	11	13.4
Technical Degree	0	100	75.8
Technical Degree	1	32	24.2

```
# Visualize attrition by EducationField
ggplot(attrition_by_educationfield, aes(x = EducationField, y = Percentage, fill = Attrition)) +
  geom_bar(stat = "identity", position = "stack") +
  coord_flip() + # Flip coordinates for better readability
  labs(title = "Attrition Percentage by Education Field",
        x = "Education Field",
        y = "Percentage (%)") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5))
```



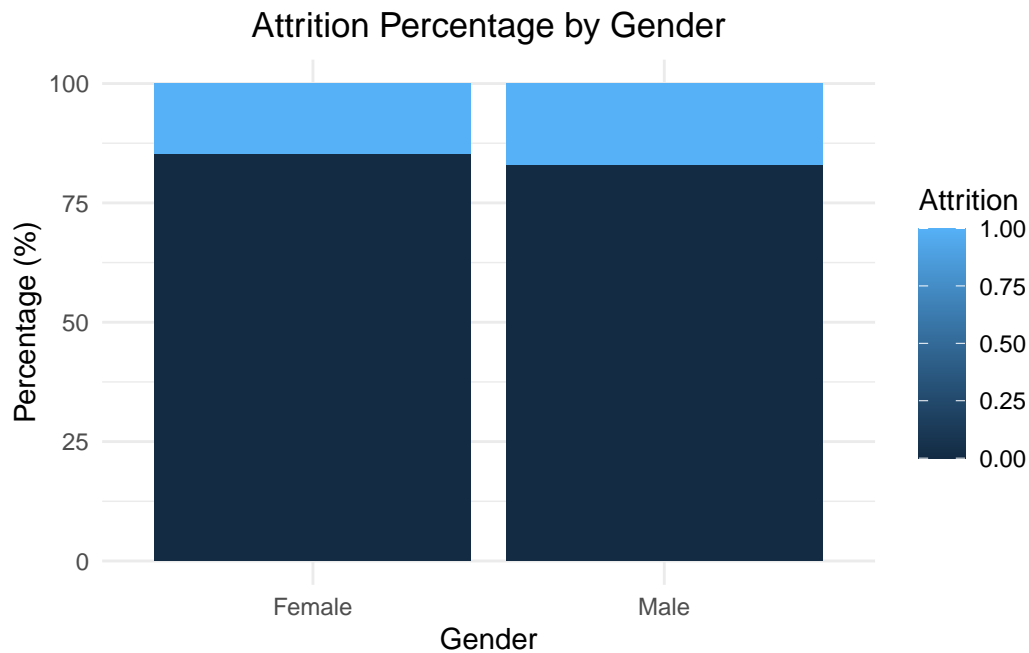
```
#Calculate attrition percentages by Gender
attrition_by_gender <- data %>%
  group_by(Gender, Attrition) %>%
  summarise(Count = n(), .groups = 'drop') %>%
  group_by(Gender) %>%
  mutate(Percentage = round(Count / sum(Count) * 100, 1))

# Display table
attrition_by_gender %>% kable()
```

Gender	Attrition	Count	Percentage
Female	0	501	85.2
Female	1	87	14.8
Male	0	732	83.0
Male	1	150	17.0

```
# Visualize attrition by Gender
ggplot(attrition_by_gender, aes(x = Gender, y = Percentage, fill = Attrition)) +
  geom_bar(stat = "identity", position = "stack") +
  labs(title = "Attrition Percentage by Gender",
       x = "Gender",
```

```
y = "Percentage (%)" +
theme_minimal() +
theme(plot.title = element_text(hjust = 0.5))
```



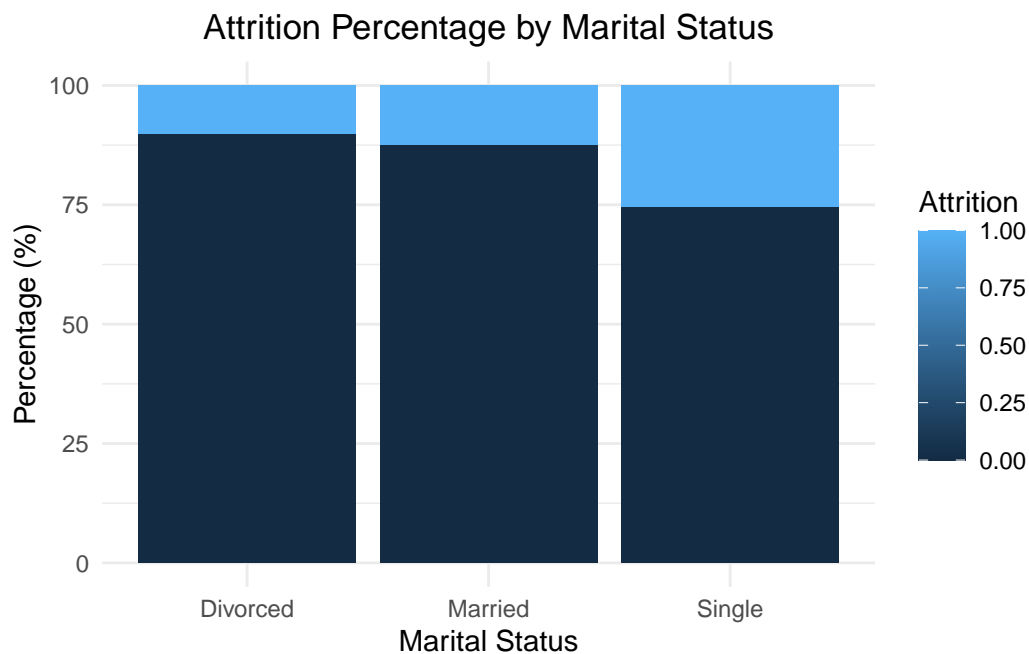
```
# Calculate attrition percentages by MaritalStatus}
attrition_by_maritalstatus <- data %>%
  group_by(MaritalStatus, Attrition) %>%
  summarise(Count = n(), .groups = 'drop') %>%
  group_by(MaritalStatus) %>%
  mutate(Percentage = round(Count / sum(Count) * 100, 1))

# Display table
attrition_by_maritalstatus %>% kable()
```

MaritalStatus	Attrition	Count	Percentage
Divorced	0	294	89.9
Divorced	1	33	10.1
Married	0	589	87.5
Married	1	84	12.5
Single	0	350	74.5
Single	1	120	25.5

MaritalStatus	Attrition	Count	Percentage
---------------	-----------	-------	------------

```
# Visualize attrition by MaritalStatus
ggplot(attrition_by_maritalstatus, aes(x = MaritalStatus, y = Percentage, fill = Attrition))
  geom_bar(stat = "identity", position = "stack") +
  labs(title = "Attrition Percentage by Marital Status",
       x = "Marital Status",
       y = "Percentage (%)") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5))
```



```
# Calculate attrition percentages by OverTime
attrition_by_overtime <- data %>%
  group_by(OverTime, Attrition) %>%
  summarise(Count = n(), .groups = 'drop') %>%
  group_by(OverTime) %>%
  mutate(Percentage = round(Count / sum(Count) * 100, 1))

# Display table
attrition_by_overtime %>% kable()
```



OverTime	Attrition	Count	Percentage
No	0	944	89.6
No	1	110	10.4
Yes	0	289	69.5
Yes	1	127	30.5

```
# Visualize attrition by OverTime
ggplot(attrition_by_overtime, aes(x = OverTime, y = Percentage, fill = Attrition)) +
  geom_bar(stat = "identity", position = "stack") +
  labs(title = "Attrition Percentage by OverTime",
       x = "OverTime",
       y = "Percentage (%)") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5))
```

