# University of Exeter Business School

**BEM2031 Introduction to Business Analytics**

**Assignment**

**Academic Year 2023 / 2024**

Student ID: 710032265

# Contents

# 1 Understanding the Business

## 1.1 The Goals of the Project

The overarching goal of this project entails conducting an extensive evaluation of employee statistics, particularly focusing on instances of employees leaving the company. This involves the meticulous examination of data to discern patterns and draw conclusive insights. Then, to provide the employers a detailed breakdown of steps to help reduce the number of employees leaving. This includes what steps to take, how to implement them and which departments to focus on first .

## 1.2 The Analysis

The data analysis is thorough yet incomplete and inaccurate. The report identifies trends that mostly use satisfaction level, which is not the most useful variable and is the only purely subjective result that is skewed for bad employees. Further, and more importantly, the report has no breakdown of department, it simply searches for trends in the data across the whole company. Then still offers no actionable solutions to the tailored to the wider company or specific departments.

The costs of this analysis of employee retention primarily include the resources and time invested in collecting, organising, and analysing the data. Additionally, there would inevitably be costs associated with implementing any recommendations derived from the analysis, had there been any recommendations. These could include training programs or policy changes. On the other hand, the benefits of this analysis could potentially be significant. By identifying patterns and trends in employee turnover, the company can take proactive measures to address underlying issues, ultimately improving employee retention. This can lead to higher employee morale, increased productivity, and ultimately, cost savings associated with recruiting and training. However, the benefits are to be limited as the analysis fails to provide any actionable insights tailored to the company or the specific departments.

From this specific report it is unlikely that anyone would benefit or be potentially harmed as the report writer makes no actionable recommendations. However, an employer may take their own insights from the report. For example as the Random Forest implies that 'Satisfaction Level' and 'Time spent in the Company' are the strongest variables. The employer may then be more inclined to let go employees they deem to be unhappy, or that spend what they perceive to be too much time at work.

Furthermore, it is suggested that 'Salary' is not an important factor to employee retention, this statistic would only drive employers to give fewer raises, and offer lower starting salaries, increasing the harm to employees. Moreover, given deeper analysis, this data trend may not be accurate for all departments and could actually lead to an increase in employee turnover.

Overall, this report is unlikely to produce any benefits to anyone as it has not shown any actionable solutions to be implemented.

## 2 Understanding the Data

### 2.1 The Data

The data is somewhat appropriate, however it lacks some important additions and useful terminologies. Firstly, there is no information on when any of the data was collected. This is important as had this data been available, trends over time would have been able to be analysed. This would help contribute towards making actionable changes to potential managers and staff. Further, the term "work accident" may not fully encompass all work-related incidents, such as cases involving abuse or assault. This suggests potential gaps in the coverage of workplace safety concerns. Finally, while the inclusion of salary and department data is valuable, converting these string variables into numerical values could enhance the analytical capabilities robustness of the analysis.

There is a multitude of information that is unknown to the reader of the report. From the source of the data to the year it was collected there are some serious gaps in information.
Understanding whether the data was collected from a manager, employer or HR significantly impacts the data and biases that lie within it. As, if it were collected from HR, this would invalidate the data for the HR department as it would have been collected by their superiors. which makes for poor data reliability. Furthermore, regarding the time frame during which the data was collected. Understanding when these results were compiled and whether there are discernible trends over time could offer valuable context for interpreting the findings.

Some data that should have been included could include more specific performances metrics, such as date of departure from the company or even location. This additional data could show potential geographical trends amongst different branches of the company, or similar trends over time. For example, showing 'an increase in sales employees leaving in the London branch over the last 2 years'. This sort of data is actionable and can lead to clear recommendations to help reduce turnover.

### 2.2 The Visualisations in the Report

The visualisations throughout the report are dubious at best. For example Figure(1) shows the large distribution of data amongst the different departments. Yet the report goes on to view the average salary of all the data, which now knowing the department distribution, holds little relevance. What more, after this point the report never returns to to concept of viewing data based on department, despite knowing the lack of even representation.
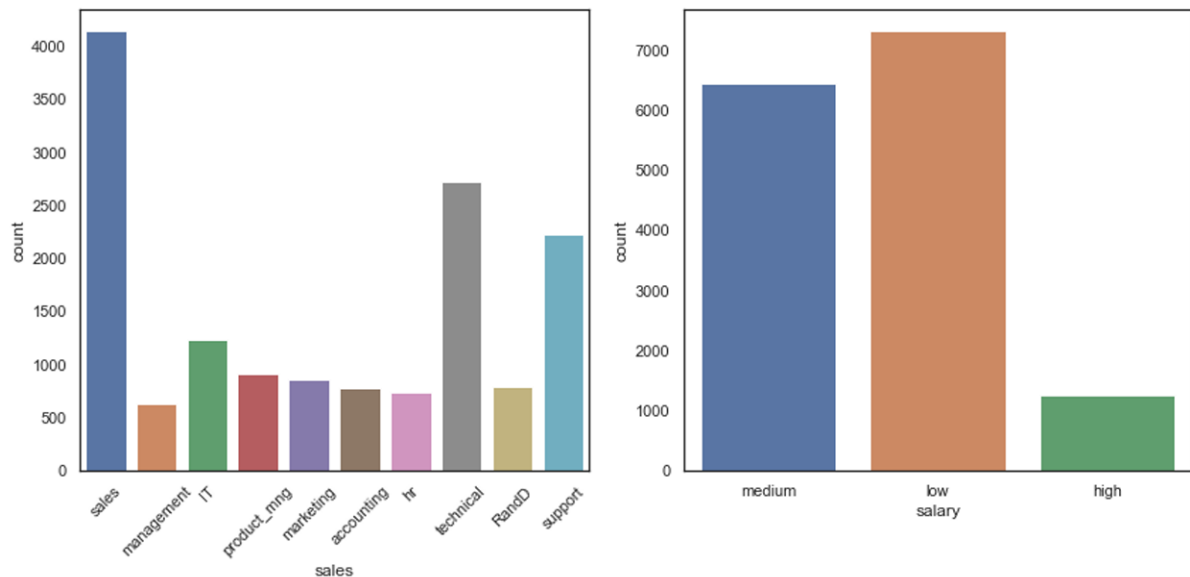
Figure 1: The First Visualisation from the Report

This represents how poor the visualisations are at conveying a meaningful narrative, fundamentally the report falls at the first hurdle in this sense.

## 2.3   Improvements to the Visualisations

To improve the effectiveness of the visualisations the 'Salary' data should have been further refined. Searching for inter-departmental trends would lead to accurate data analysis that can lead to actionable, effective suggestions. For example 'Sales' is the largest department, and 'Management' is the smallest, so one could assume there would be a discrepancy between the two. Hence why below I have done precisely this and will continue in the 'Improved Analysis' Section.
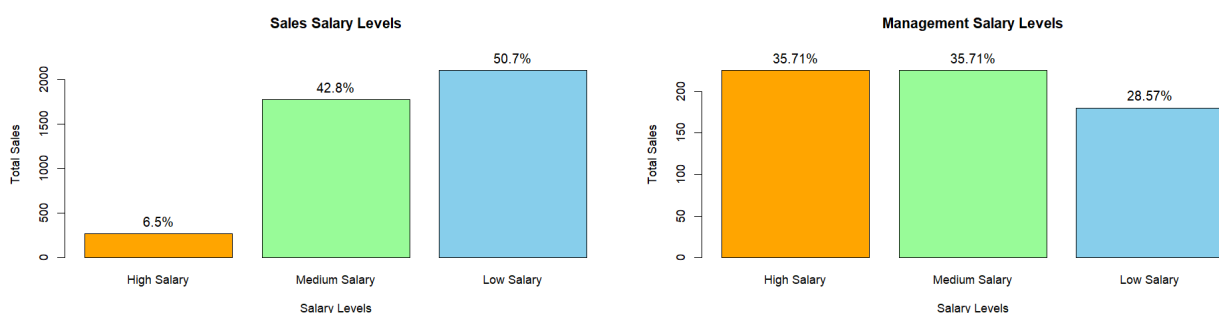


Figure 2: Sales and Management Salary Breakdown

Figure(2) shows precisely why the original report should have broken down the data by department. It is evident that the 'Salary' distributions are very different when comparing 'Sales' and 'Management'.

# 3 Data Preparation

## 3.1 Data Corruption

The report adequately addressed the presence of missing values within the dataset, demonstrating a proactive approach to data analysis. However, the report failed to examine other possible forms of data corruption. For example no attempt was made to check for duplicate or incorrectly formatted data entries. Failing to verify this makes the report vulnerable to producing incorrect visualisations and undermining the accuracy of the findings.

Further, the report could have delved into detecting and managing outliers within the dataset. Looking at the Interquartile Range (IQR) and the z score would have helped the report identify the potential anomalous data. Additionally, given the lack of numeration in the 'Salary' and 'Sales' columns, the data the should have been reshaped to have new 'Dummy Columns' containing Boolean values (1s and 0s).

## 3.2 Trust in the Data

As to weather we can trust the data, there are many different questions to be asked:

- What branch/location was the data was collected from?

- What is the source of the information?

- What year is the data from?

- Who collected the data?

This information is crucial to having reliable trustworthy data. For example the data would be much more credible had it been collected by a third party as apposed to an internal figure. Or even if there were dates attached to each entry so as to be able to truly verify if there were any duplicates.

# 4 Modelling

## 4.1 Analytical Choices

The choice to use a Random Forest and a Decision Tree were sensible options as they both offer different types of information in varying different formats. However, the main issue arises as to how they were implemented. The report used both predictive models to evaluate the company as a whole, this did find some general patterns, but never evaluates if these are true inter-departmentally. The report should have created new data frames for specific departments to evaluate the differing trends, so as to allow them to make actionable suggestions based on informed data. Failing to do this forces the report to make large generalisations not appropriate for a professional company.

## 4.2 Reliability of the Models

Looking at the Model evaluation we see that for the Random Forest the report identified a Variable Importance plot was and found that it's mean accuracy was calculated as 98.99% accurate. However this was only evaluated on the 'training data' and no 'test data' was used to verify the fit of the model, hence we cannot know whether the model was good at generalising unseen data.
As for the Decision Tree, it was plotted to view the outcomes and similarly to the Random Forest had it's mean accuracy calculated as 95%. This, like the previous model suggest a good fit however, given no further analysis using unseen data was performed, so we cannot know truly how it's accuracy.

Overfitting occurs when a model learns noise or random fluctuations in the training data instead of the underlying pattern. There were no explicit signs of overfitting mentioned in the analysis. However, the analysis required to capture such outcomes was not performed. For example, the report could have utilised ROC curves to identify the AUC of the Random Forest. Comparing the AUC for the 'training data' and the 'test data' would inform us whether the model was over or underfitting the data. On top of this, inspecting the confusion matrix of the Random Forest to evaluate is it was prone to 'False Positives' or 'False Negatives' could hugely add to the report.

## 4.3 What do the Models tell us

In the report a Variable Importance plot was introduced to show how much each variable affects the outcome of employee retention. This showed that 'Satisfaction Level' by far was the most important variable to someone leaving the company, with next being 'Time Spent' in the company and the 'Number of Projects'. However, it is important to remember that this is generalised for the whole company and may not be correct for specific departments.
Although not explicitly expressed in the report, from inspection on the decision tree we can deduce that a higher satisfaction level correlates to a higher likelihood of staying and a higher time spent in the company correlates to a higher likelihood of staying. Nevertheless, this is again a generalised estimate. It would be much more reliable and clear if a Logistic Regression Model was introduced to evaluate *how* these variable affect the outcome, as shown in an upcoming section.

## 4.4 Improved Analysis

*N.B* - (The visualisations in the following section were produced in R-Studio as improvements and additions to the original report, the code can be found in the appendix.)

When producing these visualisations, I took the original dataset, then randomly separated it into 'test' (30%) and 'training' (70%) dataset. These will be use throughout to evaluate the predictive models.

### 4.4.1 Salary Breakdown

After identifying the 'Salary' distribution in Figure(2), it is suitable to then break this down to two columns 'Left' and 'Stayed' to further evaluated any underlying trends.
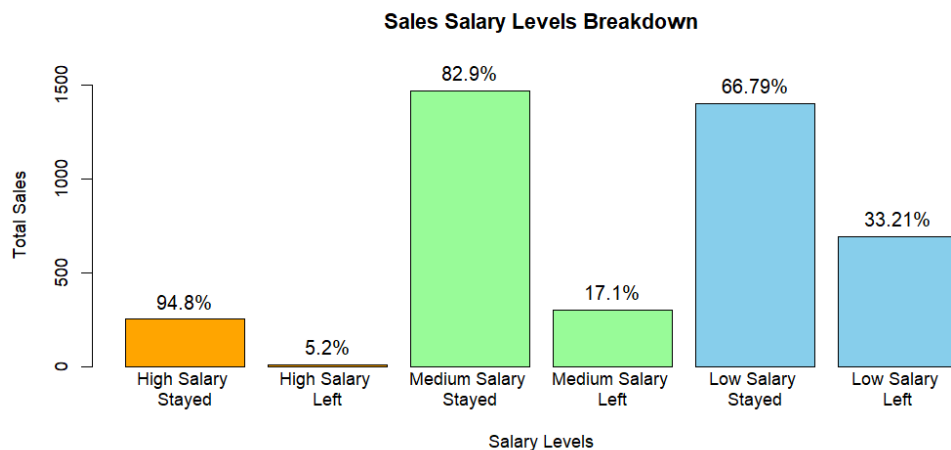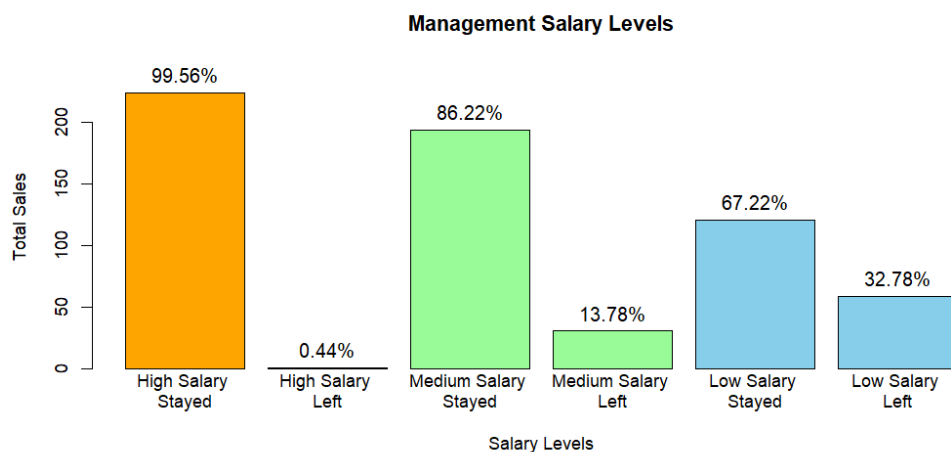


Figure 3: Sales Salary Breakdown



Figure 4: Management Salary Breakdown

Figure(4) shows us a more detailed breakdown of whether employees decided to stay or leave, based on their pay. We see that 'High Salary' Managers were more likely to stay than 'High Salary' Sales workers, and that 'Low Salary' workers in both departments are roughly equally likely to stay. 'High Salary' Managers are likely more probable to stay due to the nature of the position as it often require years of promotion, hence we could assume that the 0.44% of them who are leaving is likely due to retirement.
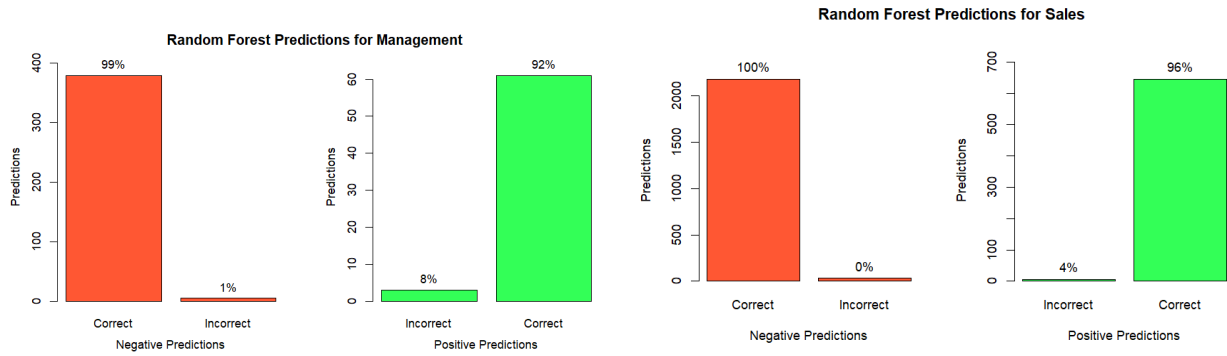
### 4.4.2 Random Forest



Figure 5: Random Forests for both Departments

Looking at the Random Forest breakdown in Figure(5)we see that for both 'Sales' and 'Management' it favours 'False Positives'. Although, this becomes twice as likely for the 'Management' data compared to the 'Sales' data. This allow us to understand potential inconsistencies with the predictions and then value providing more data from 'Sales' to the Random Forest to improve its accuracy.

### 4.4.3 ROC of the Random Forest



Figure 6: ROC Plot for training data

Looking at the ROC curves, we can identify that the Random Forest is an extremely good fit for the training data for both 'Sales' and 'Management', having achieved an AUC = 1 for both.



Figure 7: ROC Plot for test data

Then, comparing Figure(6) to the ROC of the 'test data', we see the AUC has slightly dropped to 0.9991 and 0.9892. However, these are so close to 1 we can confidently say that the Random Forest shows no signs of over-fitting or under-fitting and would be suitable to use on any new unseen data.

### 4.4.4 Variable Importance of Random Forest



Figure 8: Variable Importance Plots

Looking at this visualisation we can see that the order of variable importance doesn't change inter-departmentally. However, the Importance of them does. For exam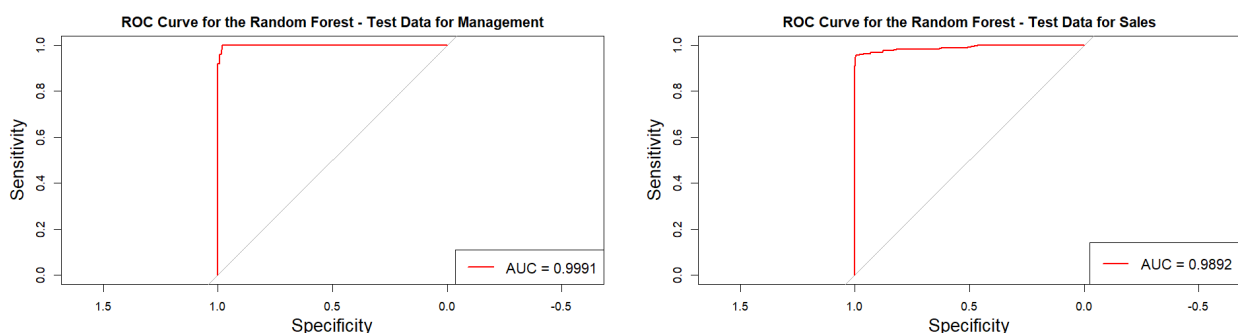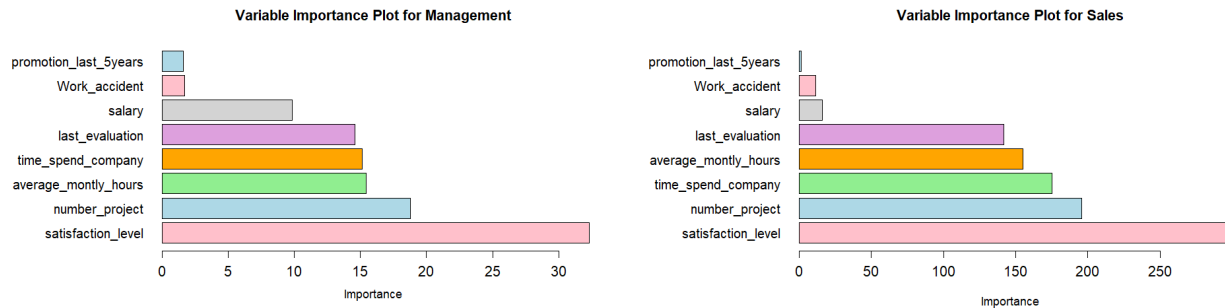ple looking at 'Salary' we identify that it is of much more importance for employees from 'Management' than from 'Sales'. This is likely due to the nature of the pay scales of both departments. But, ultimately this still doesn't tell us *how* they impact employee retention.

### 4.4.5 Variable Effects from the Logistic Regression Model



Figure 9: Variable Importance Effect Plot

To understand, *how* the different variables impact if an employee leaves, we have to perform a Logistic Regression on the data. This now show significance on a Positive and Negative scale, which allows us to understand the implications of different variables. As above the red bars represent variables which, if increased would increase the chance of an employee leaving. Along with green bars representing variables which if increased would increase the chance of an employee staying. The biggest discrepancy being 'last_evaluation' which shows as red for 'Sales' and green for 'Management'.

### 4.4.6 ROC of the Logistic Regression Model



Figure 10: ROC Plot for the Logistic Regression

To check that the Logistic Regression in Figure(9) is accurate, it is good practice to verify this using the ROC curve. This produced AUC of 0.8901 and 0.8235 for the training data and test data respectively. Both these scores suggest the model is a good fit for the data and would be suitable on new, unseen data. We can confidently say that it shows no signs of over or underfitting.

# 5 Evaluation

## 5.1 The Process, Data and Visualisations

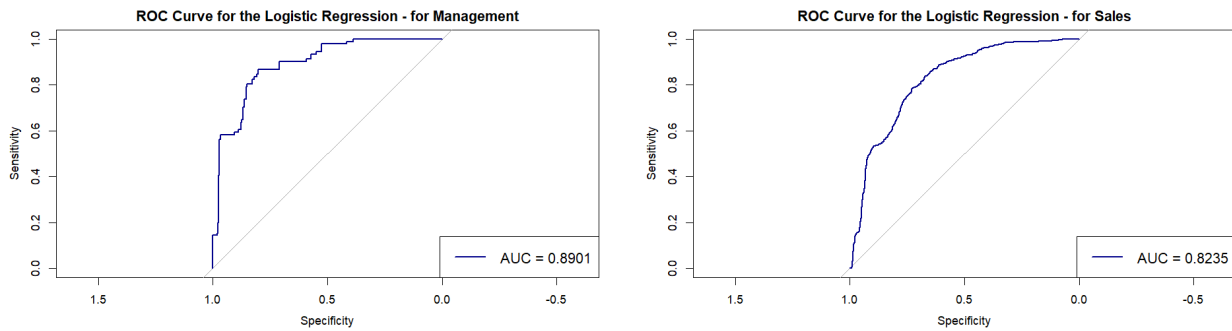The analytical process was fairly poor and did not endeavour to discover underlying trends within the company, as it did not even simply separate the data amongst departments. Further, the data was an inadequate choice, as contains no time related data, meaning no trends over time could be assessed. Not to mention that the data contained many duplicates, likely from the 'Management' column as this is not a department, but simply a section from every other department. Finally, the data contained simple spelling errors, for example 'average_montly_hours' should have been 'Average_Monthly_Hours'.

Next the analytical choices and visualisations. They did not convey any form of narrative towards identifying or solving a problem. The visualisations did identify patterns in the wider data, however, this offered no help in terms of suggestions for the company to take action.
So to sum up, no this project should not move forward to deployment as it contains no clear actionable suggestions, and if an employer were to act on some of the visualisations, it is likely that many people would be harmed and the company would not succeed as a result.

## 5.2 Changes

To continue this analysis into the future I have broken down the variable importance from the Random Forest into different departments. This allows us to identify more specific trends and find out what variables the different departments value more. For example 'Management' value their salary quite highly, whereas for 'Sales' it is barely a significant factor. This is likely due to 'Sales' employees working off of commission, hence making their base salary less important.

Further to identify *how* these factor affect the probability of an employee leaving the company, I have performed a Logistic Regression Model as shown in Figure(9). Doing this for separate departments is vital, as we can understand which areas to target specifically, and then make informed decisions.

## 5.3 Additions

Some additions to the report could include, some analysis on the Random Forest to verify it's accuracy. This could include ROC plots as I used in Figure(6) or perhaps PRC. The addition of these would make the results more reliable, adding to the robustness of the report allowing for trusted recommendations.

# 6 Deployment

## 6.1 Actions

**From the Report:**

I would recommend that the company aims to improve the general satisfaction of all the employees. However, as the report has no actionable suggestions on how to implement this, I could make no recommendations to elaborate on this.

**From my Own Analysis:**

I would recommend that the company aims to increase the rate at which it promotes its 'Sales' employees so as to improve their retention. This is because they value promotions significantly. This could be due to promotions being rarer in 'Sales' or some other underlying factors. To understand this further, deeper analysis would nee to be carried out. It is feasible to increase the promotion rate of 'Sales' employees as there would likely be other departments that do not uphold promotions as highly. Hence the resources could be allocated differently.

Further the company could aim to slightly reduce the total time employees spend at the company for both 'Management' and 'Sales' employees as this increases their likelihood of staying as shown in Figure(9) . The task of reducing 'Sales' employee's time spent at the company is fairly straight forward. Although, it is understandable that employers may be sceptical about this reduction and it's effect on profits/turnover. This model does not explicitly tell us if the reduction in work hours would directly translate to increased profits, but it does inform us that it would lead to a higher employee retention.

A final suggestion could be for the company to aim to reduce the rate at which it gives 'Sales' employees raises, as it has a lesser effect on their retention. Although, we have seen that promotions have a strong impact, from Figure(9) we see that without the title, a pay raise for 'Sales' employees does little for their retention. This money could be distributed to other departments to improve their retention depending on how they value raises. However, this needs to be implemented carefully so as not to harm the 'Sales' employees. Obviously, reducing the raises for 'Sales' employees is not a difficult task for employers. But, the difficulty comes in moderating this, so as not to disproportionately target them and create further pay imbalances.

Overall, I would suggest that the company invests in collecting more detailed data. More particularly some data matching to the date at which the current data was collected. Then also I would suggest having more detailed data, such as when an employee last asked for a raise and if it was successful, as well as whether they received compensation for their work place accident. This last one is particularly important as more deeper analysis is needed to look into 'work_accident' as in Figure(9) some employees are more likely to stay given they had a work accident. I believe this is a good area for the company to investigate further. Further, it would be useful to have the employees previous evaluations, this way we could identify if 'Hard Working' employees are being adequately compensated. I believe this data analysis would be the most impactful on driving company profits.

## 6.2 Deploying an Ongoing Process

To effectively deploy this strategy, I would suggest that the company implements the previous recommendations, including driving my analysis further by examining each department individually and assessing the variables they prioritise the most through additional Variable Importance plots. Furthermore, enhancing the analysis by incorporating more data to track ongoing trends over time would be hugely beneficial. In addition to these steps, fostering a culture of data-driven decision-making across all levels of the organisation is essential. This involves

providing training and resources to employees to enable them to understand and utilise the data effectively in their day-to-day operations. Regular communication and feedback loops should be established to ensure that insights from data analysis are effectively integrated into decision-making processes. Ultimately, the company should consistently evaluate their actions and their impact to ensure continuous improvement in company turnover.

## 6.3 Different types of Deployment

A **Descriptive Deployment** requires identifying important metrics and trends in the data. For example recognising the importance of 'satisfaction_level' in the data and understanding how it affects employee retention.

A **Predictive Deployment** requires a consistent approach to the model, ensuring it is regularly updated with new information to improve the accuracy for new predictions. In the case if the current data, this could look like providing the Random Forest with updated information on existing employees, such as updating their salary status or even new information such as their 'Last_Request_for_Raise' or 'Date_of_Last_Evaluation'.

A **Prescriptive Deployment** requires suggesting feasible recommendations, to help improve the valued outcomes. Such as the ones previously suggested, like reducing 'Sales' employees work hours to improve employee retention.

To implement an ongoing prescriptive dashboard system I would need to develop systems to automatically input the newly incoming data into the Random Forest after being cleaned. Alongside this, new algorithms to verify the accuracy and generate feasible recommendations would be put in place. These would help to create a clear and simple system for employers and partners to interact with, and then make quality informed decisions.

## 6.4 Consequences

From this analysis, if my suggestions were implemented correctly I would expect employee turnover to decrease dramatically in all departments, as well as employee satisfaction to increase or at least remain constant. In turn this would drive up company profits, improving the quality of life for the employees and the profit margins for the employers.

# 7   Code Appendix

```r
1  rm(list = ls())
2  library(tidyverse)
3  library(ggplot2)
4  library(outliers)
5  library(fastDummies)
6  library(factoextra)
7  library(FactoMineR)
8  library(survival)
9  library(survminer)
10 library(randomForest)
11 library(fmsb)
12 library(rpart)
13 library(rpart.plot)
14
15
16 # Import the data
17 hr <- read.csv('kaggle_hr_analytics.csv')
18
19 #summary(hr)
20 glimpse(hr)
21
22 # dummy columns
23 hr <- dummy_cols(hr, c('salary', 'sales'))
24
25 # rename columns
26 hr <- hr %>%
27   rename(department =  sales) %>%
28   mutate(department = fct_infreq(department))
29
30 # add an ID column
31 hr$ID <- as.character(1:nrow(hr))
32 glimpse(hr)
33
34
35
36
37
38 ##sales salary breakdown
39
40 # select the data
41 hr_sales_low =  sum(hr$sales_sales[hr$salary_low == 1])
42 hr_sales_med =  sum(hr$sales_sales[hr$salary_medium == 1])
43 hr_sales_high =  sum(hr$sales_sales[hr$salary_high == 1])
44
45
46 # Create a vector containing the sums
47 sales_sums <- c(hr_sales_high, hr_sales_med, hr_sales_low)
48 total_sales <- sum(sales_sums)
49 # Calculate percentages
50 percentages <- round((sales_sums / total_sales) * 100, 2)
51
52
53
54
```

```r
55  # Assigning colors
56  bar_colors <- c("orange", "palegreen","skyblue")
57
58  # Create bar plot with specified colors
59  barplot(sales_sums,
60          beside = TRUE,
61          names.arg = c( "High Salary", "Medium Salary","Low Salary"),
62          main = "Sales Salary Levels",
63          xlab = "Salary Levels",
64          ylab = "Total Sales",
65          col = bar_colors,
66          ylim = c(0, max(sales_sums) * 1.1))
67
68  # Add percentages on the bars
69  text(x = barplot(sales_sums, beside = TRUE, plot = FALSE),
70       y = sales_sums + 0.5,
71       labels = paste0(percentages, "%"),
72       pos = 3,
73       cex = 1.1,
74       col = "black")
75
76
77
78
79
80
81
82  ## managment salary breakdown
83
84  # select new data
85  hr_managment_low =  sum(hr$sales_management[hr$salary_low == 1])
86  hr_managment_med =  sum(hr$sales_management[hr$salary_medium == 1])
87  hr_managment_high =  sum(hr$sales_management[hr$salary_high == 1])
88
89
90  # Create a vector containing the sums
91  hr_sums <- c(hr_managment_high,  hr_managment_med, hr_managment_low)
92
93  # Calculate total hr
94  total_hr <- sum(hr_sums)
95
96  # Calculate percentages
97  percentages <- round((hr_sums / total_hr) * 100, 2)
98
99  # Assigning colors
100 bar_colors <- c("#FF5733", "#33FF57", "#3366FF")
101 bar_colors <- c("orange", "palegreen","skyblue")
102
103 # Create bar plot with specified colors
104 barplot(hr_sums,
105         beside = TRUE,
106         names.arg = c( "High Salary", "Medium Salary","Low Salary"),
107         main = "Management Salary Levels",
108         xlab = "Salary Levels",
109         ylab = "Total Sales",
110         col = bar_colors,
```

```r
111          ylim = c(0, max(hr_sums) * 1.1))
112
113
114 # Add percentages on the bars
115 text(x = barplot(hr_sums, beside = TRUE, plot = FALSE),
116      y = hr_sums + 0.5,
117      labels = paste0(percentages, "%"),
118      pos = 3,
119      cex = 1.1,
120      col = "black")
121
122
123
124
125
126
127 ## sales salary left breakdown
128 hr_sales_low_stay =  sum(hr$sales_sales[hr$salary_low == 1 & hr$left == 0])
129 hr_sales_low_left =  sum(hr$sales_sales[hr$salary_low == 1 & hr$left == 1])
130
131 hr_sales_med_stay =  sum(hr$sales_sales[hr$salary_medium == 1 & hr$left == 0])
132 hr_sales_med_left =  sum(hr$sales_sales[hr$salary_medium == 1 & hr$left == 1])
133
134 hr_sales_high_stay =  sum(hr$sales_sales[hr$salary_high == 1 & hr$left == 0])
135 hr_sales_high_left =  sum(hr$sales_sales[hr$salary_high == 1 & hr$left == 1])
136
137
138 # Create a vector containing the sums
139 sales_sums <- c(  hr_sales_high_stay,hr_sales_high_left,
140                   hr_sales_med_stay, hr_sales_med_left,
141                   hr_sales_low_stay,hr_sales_low_left)
142
143
144 total_sales <- sum(sales_sums)
145
146
147
148 sum1 = sum(sales_sums[1] , sales_sums[2])
149 sum2 = sum(sales_sums[3] , sales_sums[4])
150 sum3 = sum(sales_sums[5] , sales_sums[6])
151
152 percentages = c(round((sales_sums[1] / sum1) * 100, 2), round((sales_sums[2] / sum1) *
        100, 2),
153                 round((sales_sums[3] / sum2) * 100, 2), round((sales_sums[4] / sum2) *
        100, 2),
154                 round((sales_sums[5] / sum3) * 100, 2), round((sales_sums[6] / sum3) *
        100, 2))
155
156
157
158
159 # Assigning colors
160 bar_colors <- c("orange", "orange", "palegreen","palegreen", "skyblue", "skyblue")
161
162 # Create bar plot with specified colors
163 barplot(sales_sums,
```

```r
164         beside = TRUE,
165         names.arg = c("High Salary \nStayed", "High Salary \nLeft",
166                       "Medium Salary \nStayed", "Medium Salary \nLeft",
167                       "Low Salary \nStayed", "Low Salary \nLeft"),
168
169         main = "Sales Salary Levels Breakdown",
170         xlab = "Salary Levels",
171         ylab = "Total Sales",
172         col = bar_colors,
173         ylim = c(0, max(sales_sums) * 1.1))
174
175
176 # Add percentages on the bars
177 text(x = barplot(sales_sums, beside = TRUE, plot = FALSE),
178      y = sales_sums,
179      labels = paste0(percentages, "%"),
180      pos = 3,
181      cex = 1.1,
182      col = "black")
183
184
185
186 ## Management salary left breakdown
187
188 hr_sales_low_stay =  sum(hr$sales_management[hr$salary_low == 1 & hr$left == 0])
189 hr_sales_low_left =  sum(hr$sales_management[hr$salary_low == 1 & hr$left == 1])
190
191 hr_sales_med_stay =  sum(hr$sales_management[hr$salary_medium == 1 & hr$left == 0])
192 hr_sales_med_left =  sum(hr$sales_management[hr$salary_medium == 1 & hr$left == 1])
193
194 hr_sales_high_stay =  sum(hr$sales_management[hr$salary_high == 1 & hr$left == 0])
195 hr_sales_high_left =  sum(hr$sales_management[hr$salary_high == 1 & hr$left == 1])
196
197
198
199 # Create a vector containing the sums
200 sales_sums <- c(  hr_sales_high_stay,hr_sales_high_left,
201                   hr_sales_med_stay, hr_sales_med_left,
202                   hr_sales_low_stay,hr_sales_low_left)
203
204
205 total_sales <- sum(sales_sums)
206 # Calculate percentages
207 percentages <- round((sales_sums / total_sales) * 100, 2)
208
209
210 sum1 = sum(sales_sums[1]  , sales_sums[2])
211 sum2 = sum(sales_sums[3]  , sales_sums[4])
212 sum3 = sum(sales_sums[5]  , sales_sums[6])
213
214 percentages = c(round((sales_sums[1] / sum1) * 100, 2), round((sales_sums[2] / sum1) *
    100, 2),
215               round((sales_sums[3] / sum2) * 100, 2), round((sales_sums[4] / sum2) *
    100, 2),
216               round((sales_sums[5] / sum3) * 100, 2), round((sales_sums[6] / sum3) *
    100, 2))
```

```r
# Assigning colors
bar_colors <- c("orange", "orange", "palegreen","palegreen", "skyblue", "skyblue")

# Create bar plot with specified colors
barplot(sales_sums,
        beside = TRUE,
        names.arg = c("High Salary \nStayed", "High Salary \nLeft",
                      "Medium Salary \nStayed", "Medium Salary \nLeft",
                      "Low Salary \nStayed", "Low Salary \nLeft"),

        main = "Management Salary Levels",
        xlab = "Salary Levels",
        ylab = "Total Sales",
        col = bar_colors,
        ylim = c(0, max(sales_sums) * 1.1))


# Add percentages on the bars
text(x = barplot(sales_sums, beside = TRUE, plot = FALSE),
     y = sales_sums,
     labels = paste0(percentages, "%"),
     pos = 3,
     cex = 1.1,
     col = "black")








## Random Forest
library(dplyr)
set.seed(42)

# CHANGE HERE!!!
variable = "Management"
hr_sales = hr[hr$sales_management == 1, ]




hr_sales$left <- factor(hr_sales$left)



#use 70% of dataset as training set and 30% as test set
sample <- sample(c(TRUE, FALSE), nrow(hr_sales), replace=TRUE, prob=c(0.7,0.3))
hr_train  <- hr_sales[sample, ]
hr_test   <- hr_sales[!sample, ]

```

```r
273
274
275  set.seed(42)
276
277  # Fit the random forest model
278  rf <- randomForest(left ~ number_project +
279                         time_spend_company +
280                         last_evaluation +
281                         average_montly_hours +
282                         salary +
283                         Work_accident +
284                         promotion_last_5years +
285                         satisfaction_level
286                     ,
287
288                     data = hr_train,
289                     type = "classification")
290
291
292
293  # print confusion matrix
294  print(rf$confusion)
295
296
297  ## look at the rf
298
299
300
301  # Extract confusion matrix data with percentages
302  conf_data <- rf$confusion[, 1:3]
303
304  # Separate negative and positive predictions
305  conf_neg <- conf_data[, 1]
306  conf_pos <- conf_data[, 2]
307
308  # Set up the plotting area
309  par(mfrow = c(1, 2))
310
311  # Create bar plot for negative predictions
312  neg_bar <- barplot(conf_neg,
313                     xlab = "Negative Predictions",
314                     ylab = "Predictions",
315                     col = "#FF5733",
316                     names.arg = c("Correct", "Incorrect"),
317                     ylim = c(0, max(conf_neg) * 1.1))
318
319
320  # Add percentage labels
321  text(x = neg_bar[2],
322       y = conf_neg[2],
323       label = paste0(round(conf_data[1, 3] * 100), "%"),
324       pos = 3,
325       col = "black")
326
327  text(x = neg_bar[1],
328       y = conf_neg[1],
```

```r
329        label = paste0(round((1 - conf_data[1, 3]) * 100), "%"),
330        pos = 3,
331        col = "black")
332
333
334
335 # Create bar plot for positive predictions
336 pos_bar <- barplot(conf_pos,
337                    xlab = "Positive Predictions",
338                    ylab = "Predictions",
339                    col = "#33FF57",
340                    names.arg = c("Incorrect", "Correct"),
341                    ylim = c(0, max(conf_pos) * 1.1))
342
343
344 # Add percentage labels
345 text(x = pos_bar[2],
346      y = conf_pos[2],
347      label = paste0(round((1 - conf_data[2, 3]) * 100), "%"),
348      pos = 3,
349      col = "black")
350
351 text(x = pos_bar[1],
352      y = conf_pos[1],
353      label = paste0(round(conf_data[2, 3] * 100), "%"),
354      pos = 3,
355      col = "black")
356
357
358 # Add a title
359 title(main = paste0("Random Forest Predictions for ", variable),
360       outer = TRUE,
361       line = -1.5)
362
363
364
365
366
367
368
369 ## Look at the accuracy
370
371
372 library(pROC)
373
374
375 # plot size
376 par(mfrow = c(1, 1))
377
378 # Predictions for training data
379 pred_train <- predict(rf,
380                       newdata = hr_train,
381                       type = "prob")[, 2]
382
383 # ROC curve for training data
384 mailing_train_roc <- roc(hr_train$left, pred_train)
```

```r
385
386 # Plot ROC curve for training data
387 plot(mailing_train_roc,
388      main = paste0("ROC Curve for the Random Forest - Training Data for ",  variable),
389      col = "darkblue",
390      lwd = 2,
391      cex.lab = 1.5)
392
393 legend("bottomright",
394        legend = paste("AUC =", round(auc(mailing_train_roc), 4)),
395        col = "darkblue",
396        lwd = 2,
397        cex = 1.2)
398
399
400
401
402
403 # plot size
404 par(mfrow = c(1, 1))
405
406 # Predictions for test data
407 pred_test <- predict(rf, newdata = hr_test, type = "prob")[, 2]
408
409 # ROC curve for test data
410 mailing_test_roc <- roc(hr_test$left, pred_test)
411
412
413 # Plot ROC curve for test data
414 plot(mailing_test_roc,
415      main = paste0("ROC Curve for the Random Forest - Test Data for ",  variable),
416      col = "red",
417      lwd = 2,
418      cex.lab = 1.5)
419
420 # display a legend
421 legend("bottomright",
422        legend = paste("AUC =", round(auc(mailing_test_roc), 4)),
423        col = "red",
424        lwd = 2,
425        cex = 1.2)
426
427
428
429
430
431
432
433 ## Variable Importance
434
435 # plot dimensions
436 par(mar = c(5, 12, 4, 2), mfrow = c(1,1))
437
438 # Get variable importance data from the random forest model
439 var_imp <- data.frame(Variables = row.names(rf$importance), Importance = rf$importance
     [,1])
```

```r
440
441  # put the data in decreasing order
442  var_imp <- var_imp[order(var_imp$Importance, decreasing = TRUE), ]
443
444
445  # define the colours be used
446  col <- c("pink",
447           "lightblue",
448           "lightgreen",
449           "orange",
450           "plum",
451           "lightgrey")
452
453
454
455  # Create a bar plot with sorted values
456  barplot(var_imp$Importance,
457          names.arg = var_imp$Variables,
458          xlab = "Importance",
459          main = paste0("Variable Importance Plot for ", variable),
460          col = col,
461          cex.names = 1.1,
462          cex.axis = 1.2,
463          las = 1,
464          horiz = TRUE)  # <- this line makes the bars horizontal
465
466
467
468
469
470
471
472
473  ## Logistic Regression
474
475  # Fit logistic regression model
476  logit_model <- glm(left ~ number_project +
477                         time_spend_company +
478                         last_evaluation +
479                         average_montly_hours +
480                         salary +
481                         Work_accident +
482                         promotion_last_5years +
483                         satisfaction_level,
484
485                     data = hr_sales,
486                     family = binomial)
487
488  # Summary of the model
489  log_sum = summary(logit_model)
490
491  # Interpret coefficients and their significance
492  log_sum$coefficients
493
494
495
```

```r
496  ## plot it
497
498  # Extract coefficients from the summary
499  coefficients <- log_sum$coefficients[2:10, 0:1]
500
501
502
503  # Plot
504  par(mar = c(5, 10, 4, 2), mfrow = c(1,1))  # Adjust margin to accommodate longer variable
           names
505
506  barplot(coefficients,
507          horiz = TRUE,
508          col = ifelse(coefficients < 0, "#33FF57", "#FF5733"),
509          las = 1,
510          xlab = "Significance",
511          main = paste0("Variable Effect Plot for ", variable),
512  )
513
514  legend("bottomleft",
515         legend = c("More Likely to Stay if increased", "More Likely to Leave if increased"
           ),
516         fill = c("green", "red"))
517
518
519
520
521  ## check the accuracy:
522
523
524  library(pROC)
525
526  # Compute predicted probabilities
527  predicted_probs <- predict(logit_model, type = "response")
528
529  # Compute ROC curve
530  roc_curve <- roc(hr_sales$left, predicted_probs)
531
532  # Plot ROC curve
533  plot(roc_curve,
534       col = "darkblue",
535       main = paste0("ROC Curve for the Logistic Regression - for ",  variable))
536
537
538  legend("bottomright",
539         legend = paste("AUC =", round(auc(roc_curve), 4)),
540         col = "darkblue",
541         lwd = 2,
542         cex = 1.2)
```