

Final Project Abstract

Matthew Baker, Erinda Budo, Don Padmaperuma, Subhalaxmi Rout

Abstarct

HR Analytics finds out the people-related trends in the data and helps the HR Department take the appropriate steps to keep the organization running smoothly and profitably. Attrition is a corporate setup is one of the complex challenges that the people managers and the HRs personnel have to deal with it.

In this research assignment, we investigated data on employee attrition of a company. This is a fictional data set created by IBM data scientists.

We have collected this dataset from Kaggle, using the below link: <https://www.kaggle.com/pavansubhasht/ibm-hr-analytics-attrition-dataset>

Problem Statement The first question comes into our mind, what is attrition?

Attrition is a process in which the workforce dwindles at a company, following a period in which a number of people retire or resign, and are not replaced.

Second question, what are the reason for attrition?

This can happen for many reasons:

- Employees looking for better opportunities
- A negative working environment
- Bad management
- Excessive working hours

We will do the analysis based on Gender, Education, Income, Working Environment, and lastly, build a predictive model to determine whether an employee is going to quit or not.

Approaches follow The first part will consider the data provided and attempt to identify trends and patterns. The data is then split into training and testing sets and using machine learning techniques identify the individuals that are more likely to leave the organization.

Response or Target feature is the Attrition which is going to be our feature of interest for the prediction - based on the independent features.

Dataset Overview The dataset consists of 1470 observations (rows), 35 features (variables). There is no missing data! this will make it easier to work with the dataset. We only have two datatypes in this dataset: factors and integers.

Let's have an overview of the dataset.

```
hr_data <-  
  read.csv("https://raw.githubusercontent.com/mharrisonbaker/DATA621_GroupWork2/main/Final%20Project/WA_Fn-UseC_-HR-Employee-Attrition.csv")  
dim(hr_data)
```

```
## [1] 1470    35
```

```
str(hr_data)
```

```
## 'data.frame':    1470 obs. of  35 variables:  
## $ Age           : int  41 49 37 33 27 32 59 30 38 36 ...  
## $ Attrition      : chr  "Yes" "No" "Yes" "No" ...  
## $ BusinessTravel : chr  "Travel_Rarely" "Travel_Frequently" "Travel_Rarely" "Travel_Frequently" ...  
## $ DailyRate      : int  1102 279 1373 1392 591 1005 1324 1358 216 1299 ...  
## $ Department     : chr  "Sales" "Research & Development" "Research & Development" "Research & Development" ...  
## $ DistanceFromHome : int  1 8 2 3 2 2 3 24 23 27 ...  
## $ Education       : int  2 1 2 4 1 2 3 1 3 3 ...  
## $ EducationField  : chr  "Life Sciences" "Life Sciences" "Other" "Life Sciences" ...  
## $ EmployeeCount   : int  1 1 1 1 1 1 1 1 1 1 ...  
## $ EmployeeNumber  : int  1 2 4 5 7 8 10 11 12 13 ...  
## $ EnvironmentSatisfaction : int  2 3 4 4 1 4 3 4 4 3 ...  
## $ Gender          : chr  "Female" "Male" "Male" "Female" ...  
## $ HourlyRate       : int  94 61 92 56 40 79 81 67 44 94 ...  
## $ JobInvolvement   : int  3 2 2 3 3 3 4 3 2 3 ...  
## $ JobLevel        : int  2 2 1 1 1 1 1 1 3 2 ...  
## $ JobRole         : chr  "Sales Executive" "Research Scientist" "Laboratory Technician" "Research Scientist" ...  
## $ JobSatisfaction  : int  4 2 3 3 2 4 1 3 3 3 ...  
## $ MaritalStatus    : chr  "Single" "Married" "Single" "Married" ...  
## $ MonthlyIncome    : int  5993 5130 2090 2909 3468 3068 2670 2693 9526 5237 ...  
## $ MonthllyRate     : int  19479 24907 2396 23159 16632 11864 9964 13335 8787 16577 ...  
## $ NumCompaniesWorked : int  8 1 6 1 9 0 4 1 0 6 ...  
## $ Over18          : chr  "Y" "Y" "Y" "Y" ...  
## $ OverTime         : chr  "Yes" "No" "Yes" "Yes" ...  
## $ PercentSalaryHike : int  11 23 15 11 12 13 20 22 21 13 ...  
## $ PerformanceRating : int  3 4 3 3 3 3 4 4 4 3 ...  
## $ RelationshipSatisfaction: int  1 4 2 3 4 3 1 2 2 2 ...  
## $ StandardHours     : int  80 80 80 80 80 80 80 80 80 80 ...  
## $ StockOptionLevel  : int  0 1 0 0 1 0 3 1 0 2 ...  
## $ TotalWorkingYears : int  8 10 7 8 6 8 12 1 10 17 ...  
## $ TrainingTimesLastYear : int  0 3 3 3 3 2 3 2 2 3 ...  
## $ WorkLifeBalance    : int  1 3 3 3 3 2 2 3 3 2 ...  
## $ YearsAtCompany     : int  6 10 0 8 2 7 1 1 9 7 ...  
## $ YearsInCurrentRole : int  4 7 0 7 2 7 0 0 7 7 ...  
## $ YearsSinceLastPromotion : int  0 1 0 3 2 3 0 0 1 7 ...  
## $ YearsWithCurrManager : int  5 7 0 0 2 6 0 0 8 7 ...
```

```
# missing values  
hr_data[!complete.cases(hr_data),]
```

```
## [1] Age           Attrition      BusinessTravel  
## [4] DailyRate      Department    DistanceFromHome  
## [7] Education       EducationField EmployeeCount  
## [10] EmployeeNumber EnvironmentSatisfaction Gender  
## [13] HourlyRate      JobInvolvement JobLevel  
## [16] JobRole         JobSatisfaction MaritalStatus  
## [19] MonthlyIncome    MonthlyRate    NumCompaniesWorked  
## [22] Over18          OverTime       PercentSalaryHike  
## [25] PerformanceRating RelationshipSatisfaction StandardHours  
## [28] StockOptionLevel TotalWorkingYears TrainingTimesLastYear  
## [31] WorkLifeBalance  YearsAtCompany YearsInCurrentRole  
## [34] YearsSinceLastPromotion YearsWithCurrManager  
## <0 rows> (or 0-length row.names)
```

```
colnames(hr_data)
```

```
## [1] "Age" "Attrition"
## [3] "BusinessTravel" "DailyRate"
## [5] "Department" "DistanceFromHome"
## [7] "Education" "EducationField"
## [9] "EmployeeCount" "EmployeeNumber"
## [11] "EnvironmentSatisfaction" "Gender"
## [13] "HourlyRate" "JobInvolvement"
## [15] "JobLevel" "JobRole"
## [17] "JobSatisfaction" "MaritalStatus"
## [19] "MonthlyIncome" "MonthlyRate"
## [21] "NumCompaniesWorked" "Over18"
## [23] "OverTime" "PercentSalaryHike"
## [25] "PerformanceRating" "RelationshipSatisfaction"
## [27] "StandardHours" "StockOptionLevel"
## [29] "TotalWorkingYears" "TrainingTimesLastYear"
## [31] "WorkLifeBalance" "YearsAtCompany"
## [33] "YearsInCurrentRole" "YearsSinceLastPromotion"
## [35] "YearsWithCurrManager"
```

```
summary(hr_data)
```

```
##      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
```

Attrition Distribution Below plot shows the distribution of attrition, the employee leave the company is 236 out of 1470. This is not a balanced dataset.

```
counts <- table(hr_data$Attrition)
print(counts)
```

```
##
## No  Yes
## 1233 237
```

```
counts_per <- round(counts / 1470 * 100,1)
print(counts_per)
```

```
##
## No  Yes
## 83.9 16.1
```

```
bp <- barplot(counts_per, main="Attrition Distribution",
  xlab="Attrition True or False",
  ylab = "Count",
  col="darkblue", legend = rownames(counts))
text(bp, 0, counts_per,cex=1,pos=3, col = 'white')
```

