# Unveiling the Secrets Behind Employee Attrition at IBM: A Data-Driven Approach



## The Motivation Behind the Project

In the competitive landscape of modern businesses, retaining talent is critical. IBM, a global tech giant, has faced a surge in employee attrition recently. As a People Data Analyst intern, I embarked on a mission to uncover the factors contributing to this trend, using R and advanced statistical techniques.

## Why This Analysis Matters

Understanding the root causes of employee turnover can save organizations substantial costs associated with hiring and training new employees. This article delves into the intricacies of attrition at IBM, revealing insights that could help not only IBM but also other companies facing similar challenges. You'll learn about the data analysis process, key findings, and actionable takeaways to mitigate attrition.

# Exploring the Dataset

The dataset used for this analysis is a well-known HR dataset created by IBM data scientists. It contains 1470 rows representing employees and 35 columns with various attributes, including demographics, job-related information, and the critical "Attrition" column. The dataset is available for public use and can be found [here](#).
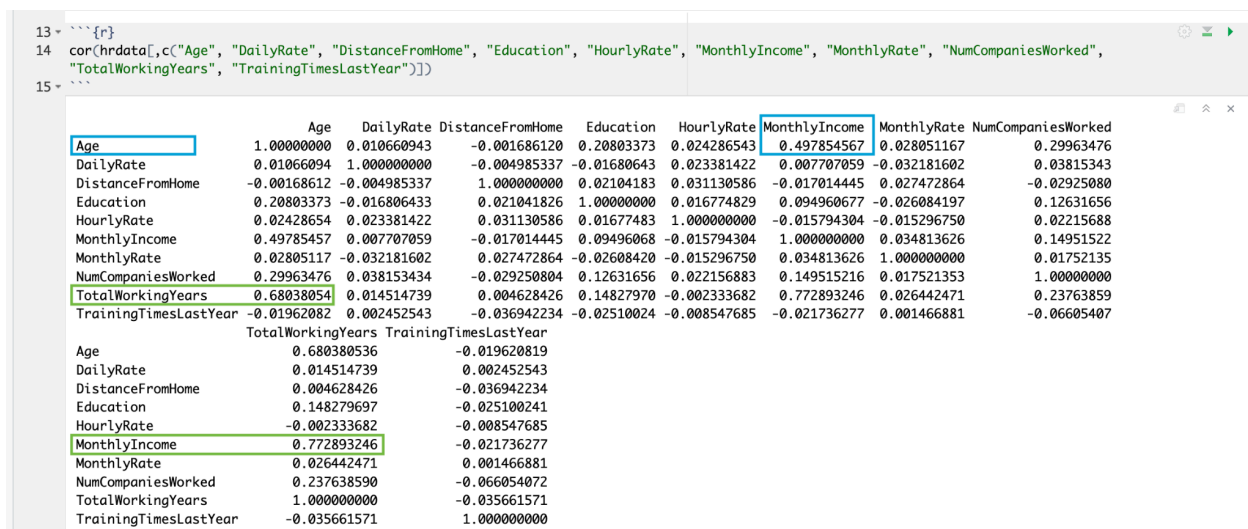
# In-depth Data Analysis

## Preparing and Cleaning the Data

First, I imported the dataset into R and conducted an initial exploration to understand the structure and contents. The key variable, "Attrition," was analyzed to differentiate between employees who stayed and those who left.

## Correlation Matrix Insights

A correlation matrix was generated for important demographic and job-related variables such as Age, Monthly Income, and Total Working Years.
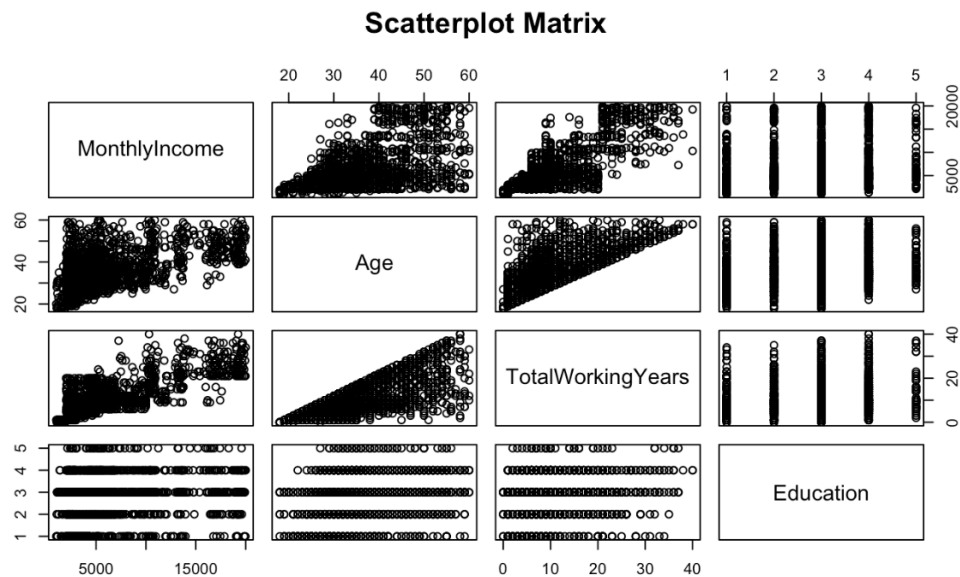


Correlation Matrix

The results revealed significant correlations between:

- Age and Total Working Years
- Age and Monthly Income
- Age and Number of Companies Worked

## Visualizing Relationships with Scatter Plots

To visualize these relationships, I created scatter plots using the pairs function in R.

```r
16 ▾ ```{r}
17    pairs(~MonthlyIncome+Age+TotalWorkingYears+Education,data = hrdata, main = "Scatterplot Matrix")
18 ▾ ```
```
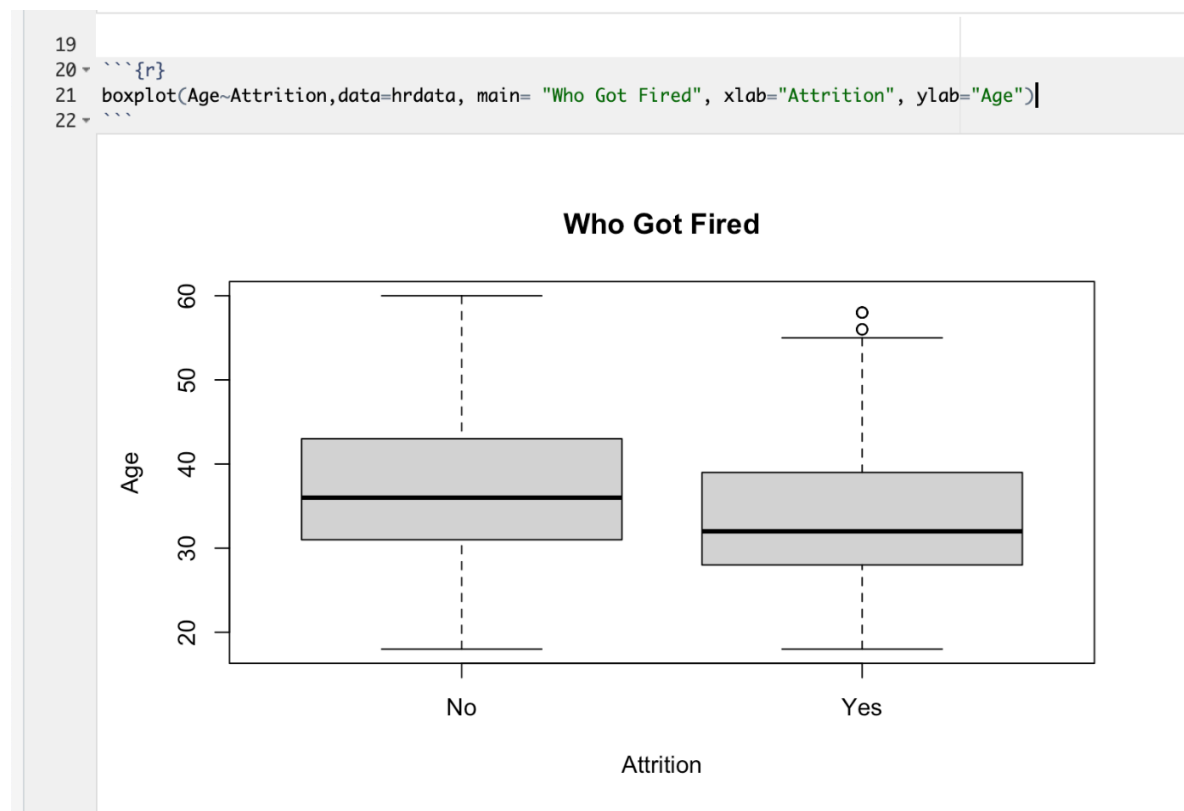


Demographic scatterplot

Here are some noteworthy observations from the scatterplot matrix:

- There is a positive correlation between age and total working years, indicating that experience tends to increase with age.
- Monthly income tends to rise with age, although there are a few individuals in their 30s who earn more than those in the subsequent age groups.

- The correlation between education and monthly income is weakly positive. This suggests that higher education doesn't always lead to higher income, although some highly educated employees do earn more.
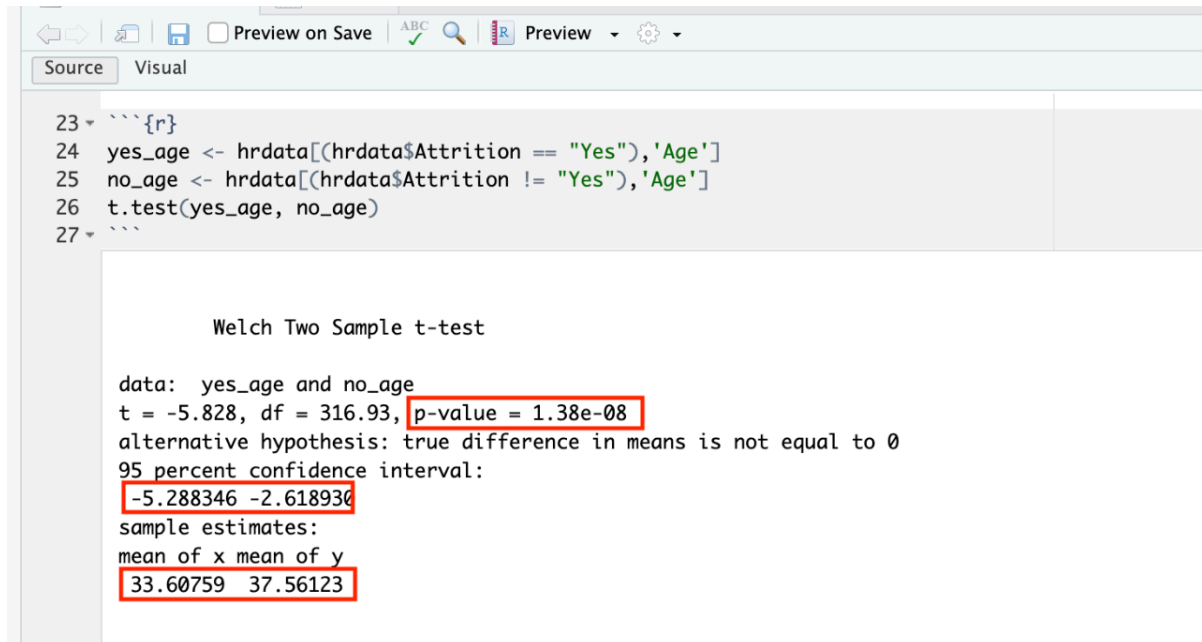
## Hypothesis Analysis: Age and Attrition

To address claims of ageism, I performed hypothesis testing using Welch Two Sample t-tests. The boxplot of age distribution for employees who stayed versus those who left suggested that younger employees were leaving at a higher rate.

```r
boxplot(Age~Attrition,data=hrdata, main= "Who Got Fired", xlab="Attrition", ylab="Age")
```



Employee Retention vs Age - Whisker Plots

"Yes" indicates employees who were laid off, while "No" represents current employees. The median age is shown by the line in the middle of the box. Both groups appear quite similar, but the average age of employees who stayed is 38 years, compared to 30 years for those who were laid off.

To confirm our hypothesis, we will perform a two-sample t-test to compare the ages of both groups and determine the p-value.

```r
23  ```{r}
24  yes_age <- hrdata[(hrdata$Attrition == "Yes"),'Age']
25  no_age <- hrdata[(hrdata$Attrition != "Yes"),'Age']
26  t.test(yes_age, no_age)
27  ```
```

```
            Welch Two Sample t-test

data:  yes_age and no_age
t = -5.828, df = 316.93, p-value = 1.38e-08
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -5.288346 -2.618930
sample estimates:
mean of x mean of y
 33.60759  37.56123
```

The t-test results confirmed this observation with a p-value less than 0.05, indicating a significant difference between the two age groups. This suggests that there is no evidential bias, and the employees who left the company are indeed younger than those who stayed.

## Key Takeaways

This project highlighted key factors influencing employee attrition at IBM, such as age and working years. Younger employees were more likely to leave, and there was a strong correlation between age and monthly income. By understanding these patterns, IBM can develop targeted strategies to improve employee retention.

If you found these insights valuable, consider how similar analysis can be applied in your organization. Feel free to reach out with questions or for collaboration on similar projects!