# SC1015 Mini Project

DSF3
HUNG KUO-CHEN
JODI TAY SEOW XUAN
YANG XIAOYUE

# Table of Contents

**01**
## Motivation
Problem definition and dataset

**02**
## EDA
Insights and tools used

**03**
## ML
Why and how the techniques used help achieve the objective(s)

**04**
## Outcome
Interesting discoveries and recommendations

# Motivation

- Employee attrition: the turnover rate in various job roles

- High attrition rate brings about problems

  - Hard-to-replace employees leave → lower productivity and profits

  - High costs incurred in training and hiring new employees

- Aim: uncover reasons for an employee's resignation and recommend improvements made within IBM to retain its employees

# IBM HR Attrition Dataset

- Used data from 1470 employees in IBM

  - 16.1% left IBM, 83.9% stayed in IBM

- Includes 34 independent variables based on an employee's profile

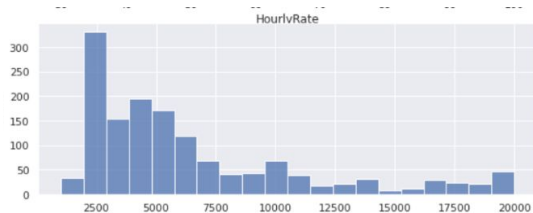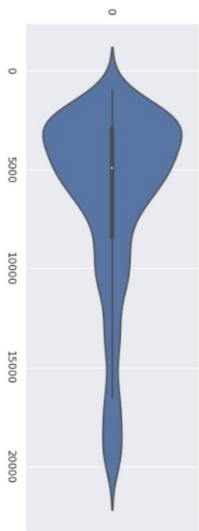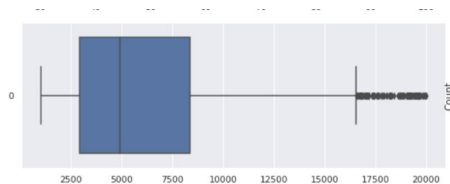  - Each contributed to whether an employee decided to leave or stay

# Cleaning the Data

- Dropped insignificant columns

  - EmployeeCount, Over18, StandardHours, EmployeeNumber

- Checked for missing values and duplicates

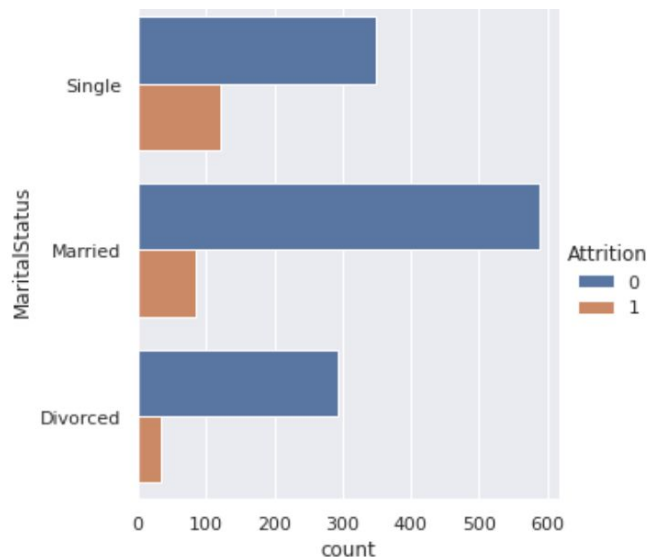  - None found

# Breakdown of Variables

## Numeric

- 14 regular numeric variables

- 9 factor numeric variables

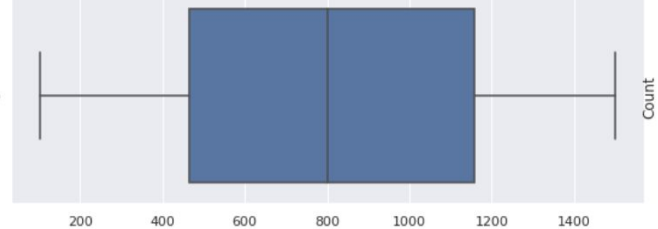- A variety of data visualization methods: box plot, histogram plot, and violin plot



## Categorical

- 6 categorical variables

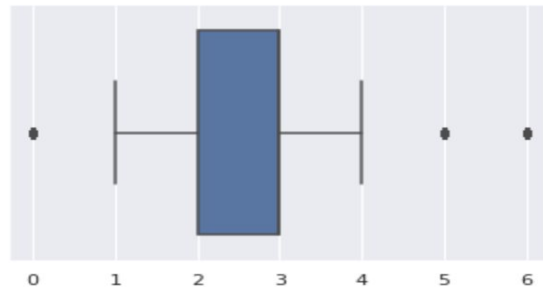- Categorical bar plot of each variable against attrition using GroupBy
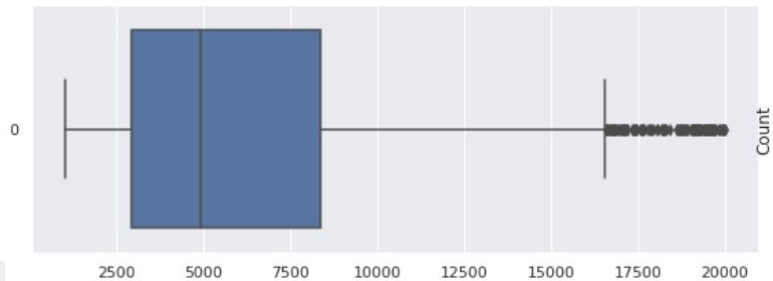
# Outliers (Box Plot)

Different variables have different degrees of outliers (close to none, moderate, and large)

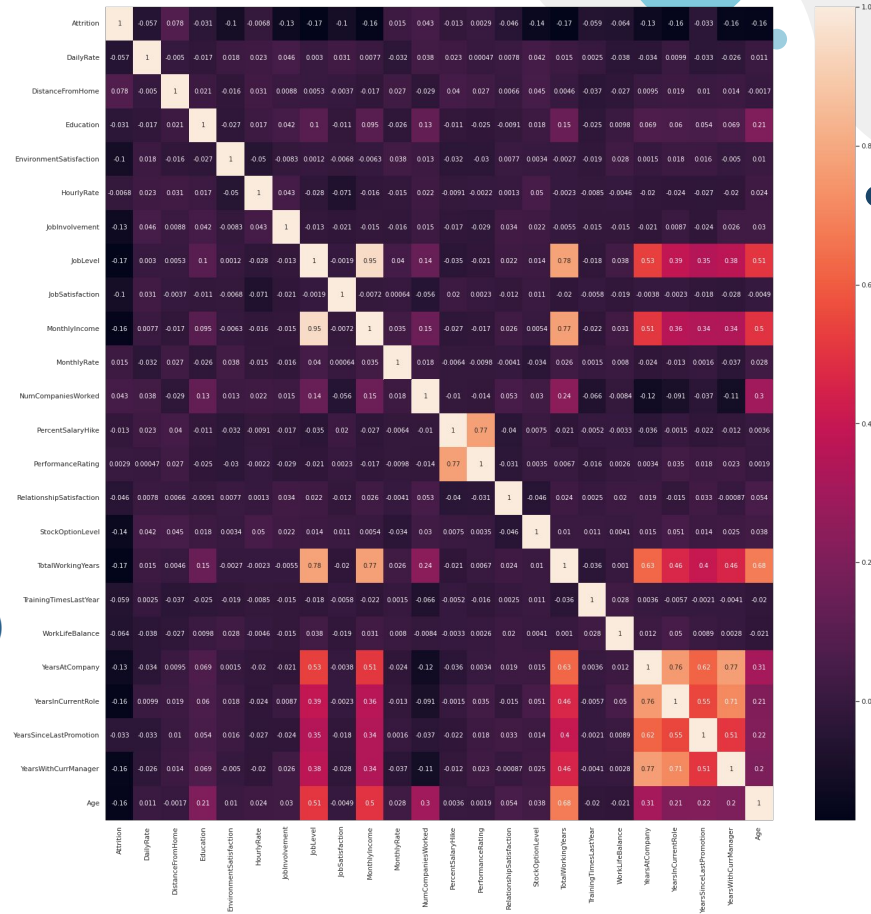

Box plot with close to no outliers



Box plot with a moderate number of outliers



Box plot with a large number of outliers
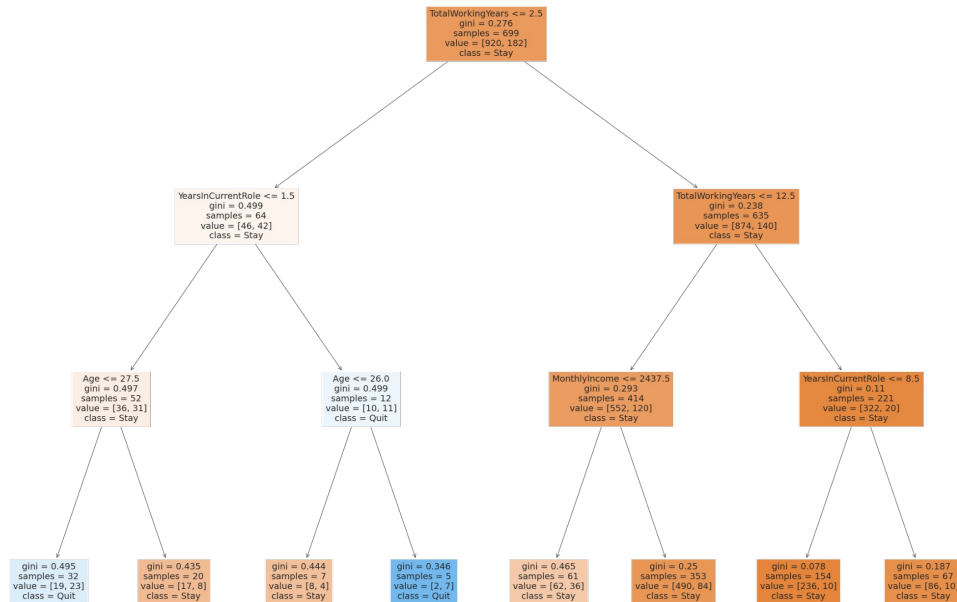
# Correlation Matrix

- Reclassify attrition as a numeric variable
- Plot the correlation matrix for attrition against all other numeric variables
- Interesting Findings:
  - MonthlyIncome and JobLevel (0.95)
  - TotalWorkingYears and JobLevel (0.78)
  - TotalWorkingYears and MonthlyIncome (0.77)

# Random Forest

- Extract variables with relatively high correlation
- **Accuracy of random forest: 84.5%**
- Tune hyperparameters with random search
  - Random combinations of the hyperparameters are used to find the optimal solution for the built model

# Logistic Regression

- Determine level of influence of categorical variables on attrition
- Convert each variable into numeric indicator variables with get_dummies
- Accuracy: 0.84 (train), 0.86 (test)
- Ineffective due to **extremely low precision**
  - The model only classifies 38% of employees that quit correctly

| | OverTime | Gender | MaritalStatus | Department | EducationField |
|---|---|---|---|---|---|
| 0 | Yes | Female | Single | Sales | Life Sciences |
| 1 | No | Male | Married | Research & Development | Life Sciences |
| 2 | Yes | Male | Single | Research & Development | Other |
| 3 | Yes | Female | Married | Research & Development | Life Sciences |
| 4 | No | Male | Married | Research & Development | Medical |
| ... | ... | ... | ... | ... | ... |
| 1465 | No | Male | Married | Research & Development | Medical |
| 1466 | No | Male | Married | Research & Development | Medical |
| 1467 | Yes | Male | Married | Research & Development | Life Sciences |
| 1468 | No | Male | Married | Sales | Medical |
| 1469 | No | Male | Married | Research & Development | Medical |

1470 rows × 5 columns

# Neural Network

- Accuracy:

# Summary of Findings

## 1. Random Forest

Numeric variables that have relatively high correlation with attrition can be used to predict attrition

## 2. Logistic Regression

Not recommended using categorical variables to predict attrition
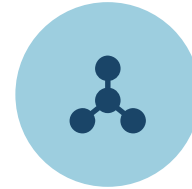
## 3. Neural Network

# Data Driven Insights

## Reasons for Leaving

- Low salary
- Low chance for career progression
- Lack of opportunities

## Improvements for IBM

- Provide more salary incentives
- Open up spots for changes in senior management

# Thank You