# PHOENIX GLOBAL EMPLOYEE ATTRITION

*RESEARCH PROJECT BY – ALI AHMAD*

# CONTENTS

Objective

Univariate Analysis

Bivariate Analysis

Multivariate Analysis

Statistical Testing

Model Building

Observations

Different Classifier Models

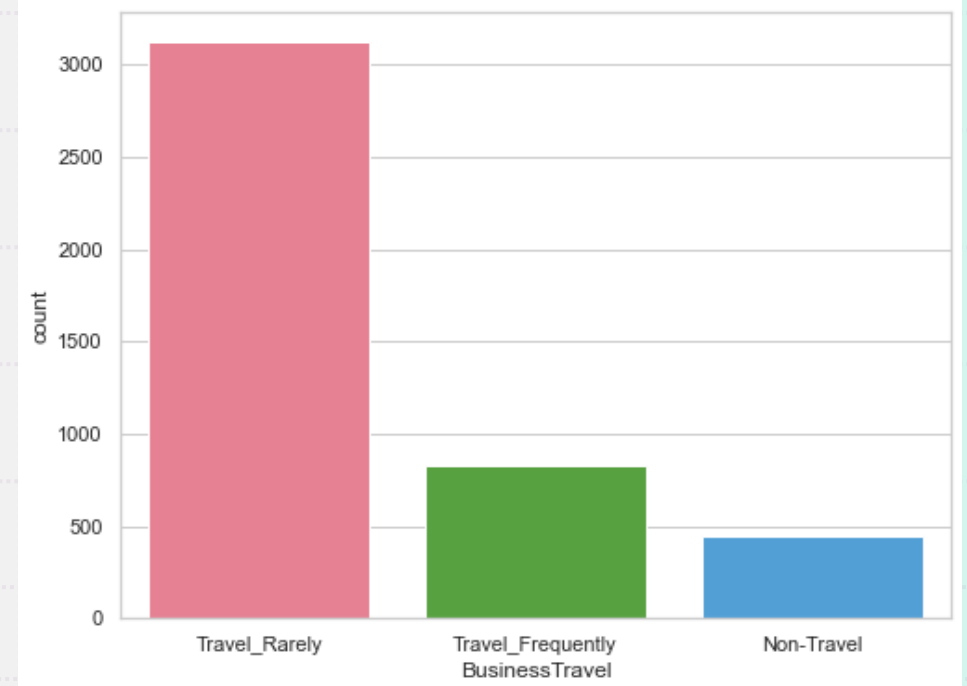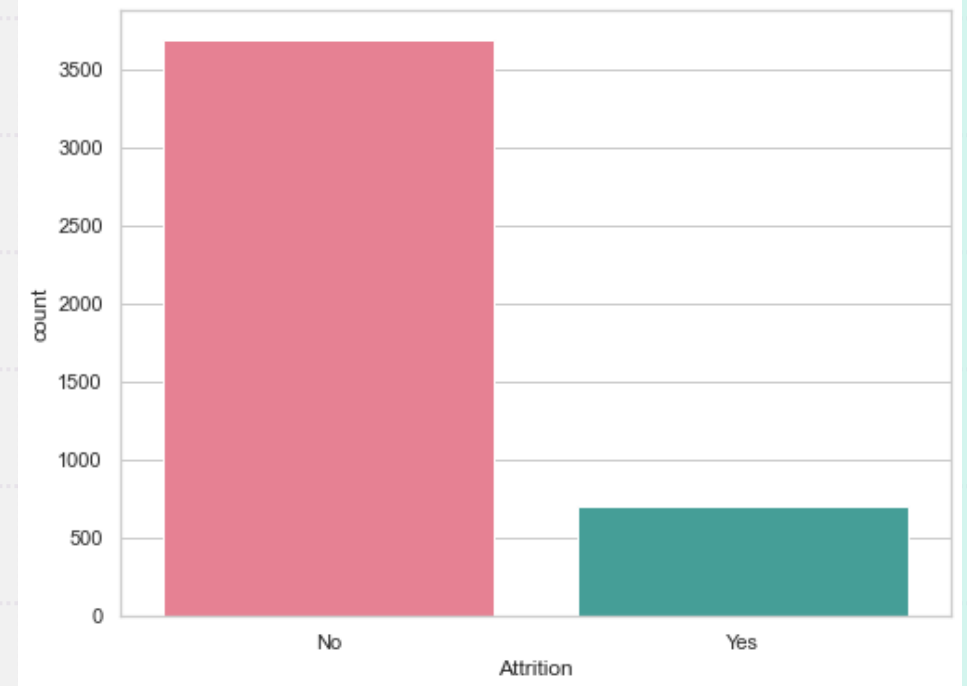Insights & Business Recommendations

# OBJECTIVE

➢ **Business Problem** : Phoenix Global company has been facing the issue of constant employee **attrition.** Close to 15 % of its employees leave the company and need to be replaced with the talent pool available in the job market. The management believes that this level of attrition (employees leaving, either on their own or because they got fired) is bad for the company.

➢ **Goal** : To understand what factors they should focus on, in order to curb attrition in the company. In other words, they want to know what changes they should make to their workplace, in order to get most of their employees to stay.

➢ We are required to model the probability of attrition using a **Logistic Regression**. The results obtained will be used by the management to understand what changes they should make to their workplace. We will also model the probability of attrition using different classifiers like Random Forest, SVM, etc for better results.
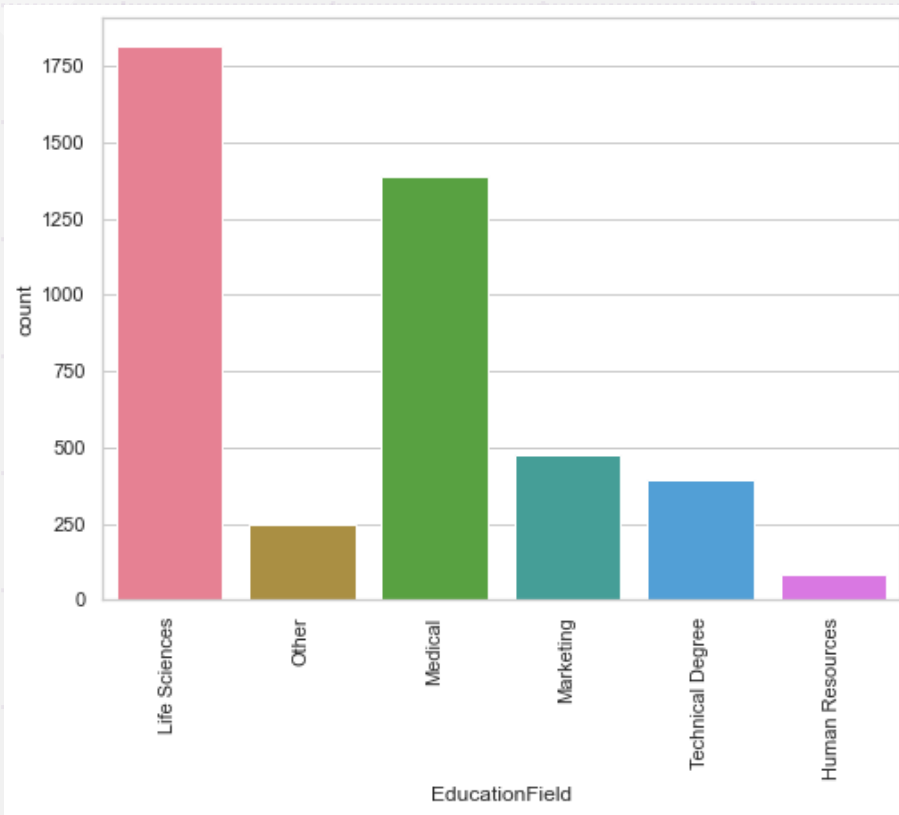
# VARIABLES IN DATASET

```
'EmployeeID', 'hrs', 'JobInvolvement', 'PerformanceRating',
'EnvironmentSatisfaction', 'JobSatisfaction',
'WorkLifeBalance', 'Age', 'Attrition', 'BusinessTravel',
'Department', 'DistanceFromHome', 'Education',
'EducationField', 'EmployeeCount', 'Gender', 'JobLevel',
'JobRole', 'MaritalStatus', 'MonthlyIncome',
'NumCompaniesWorked', 'Over18', 'PercentSalaryHike',
'StandardHours', 'StockOptionLevel',
'TotalWorkingYears', 'TrainingTimesLastYear',
'YearsAtCompany', 'YearsSinceLastPromotion',
'YearsWithCurrManager'
```
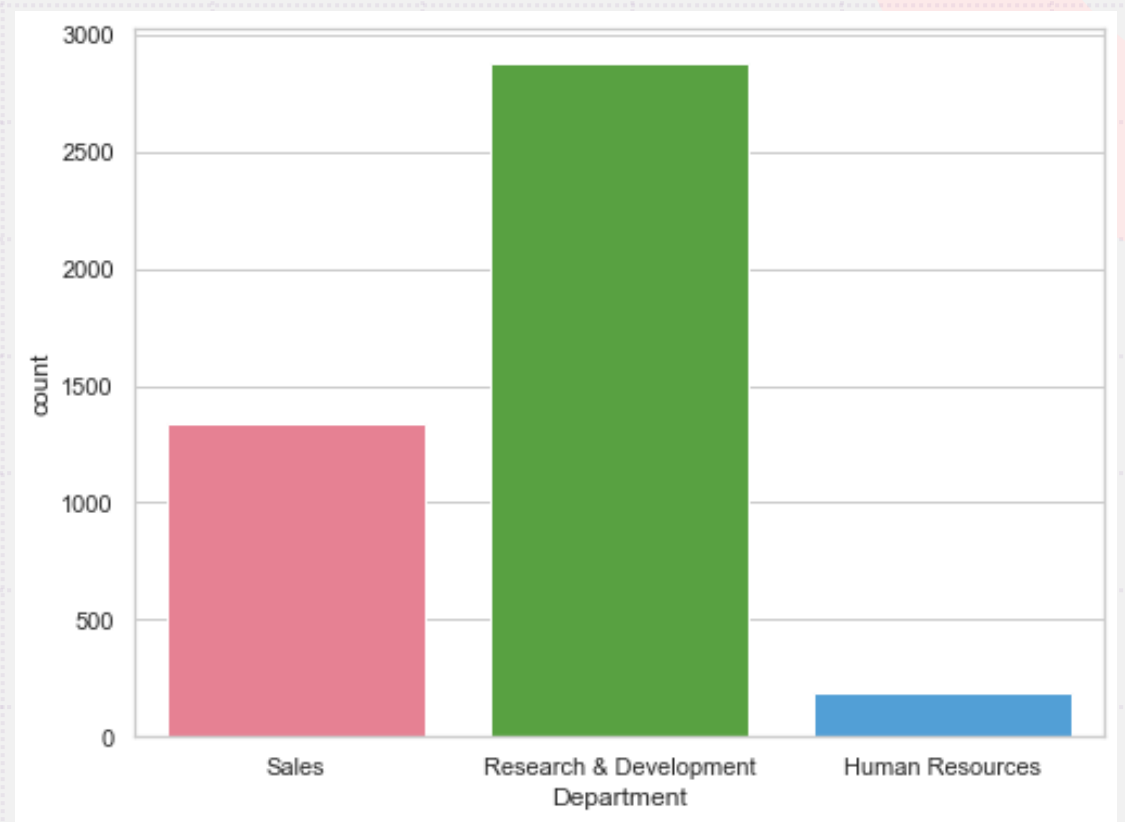
# UNIVARIATE ANALYSIS

- Over here we see that maximum employees stayed at the company (No) and less employees have left the company

- After checking the plot, it implies that this Phoenix Global data is imbalanced. We will use SMOTE technique to balance the data.

- We can see that employees who left the company are mostly travelling rarely followed by frequent travel employees
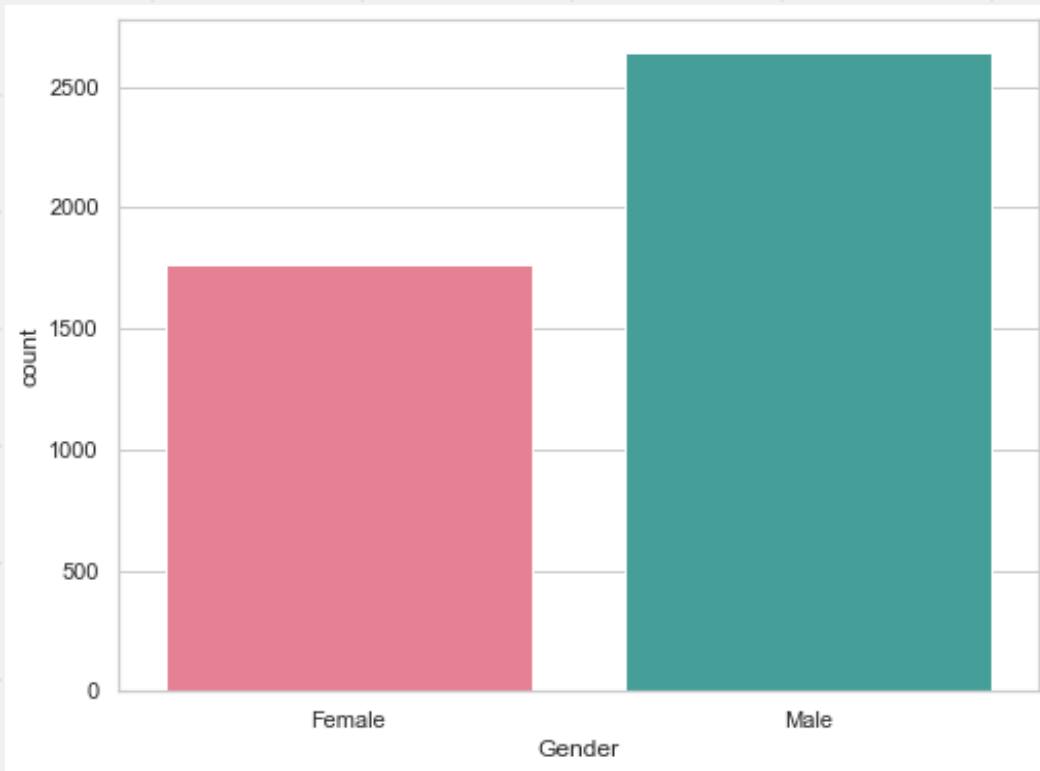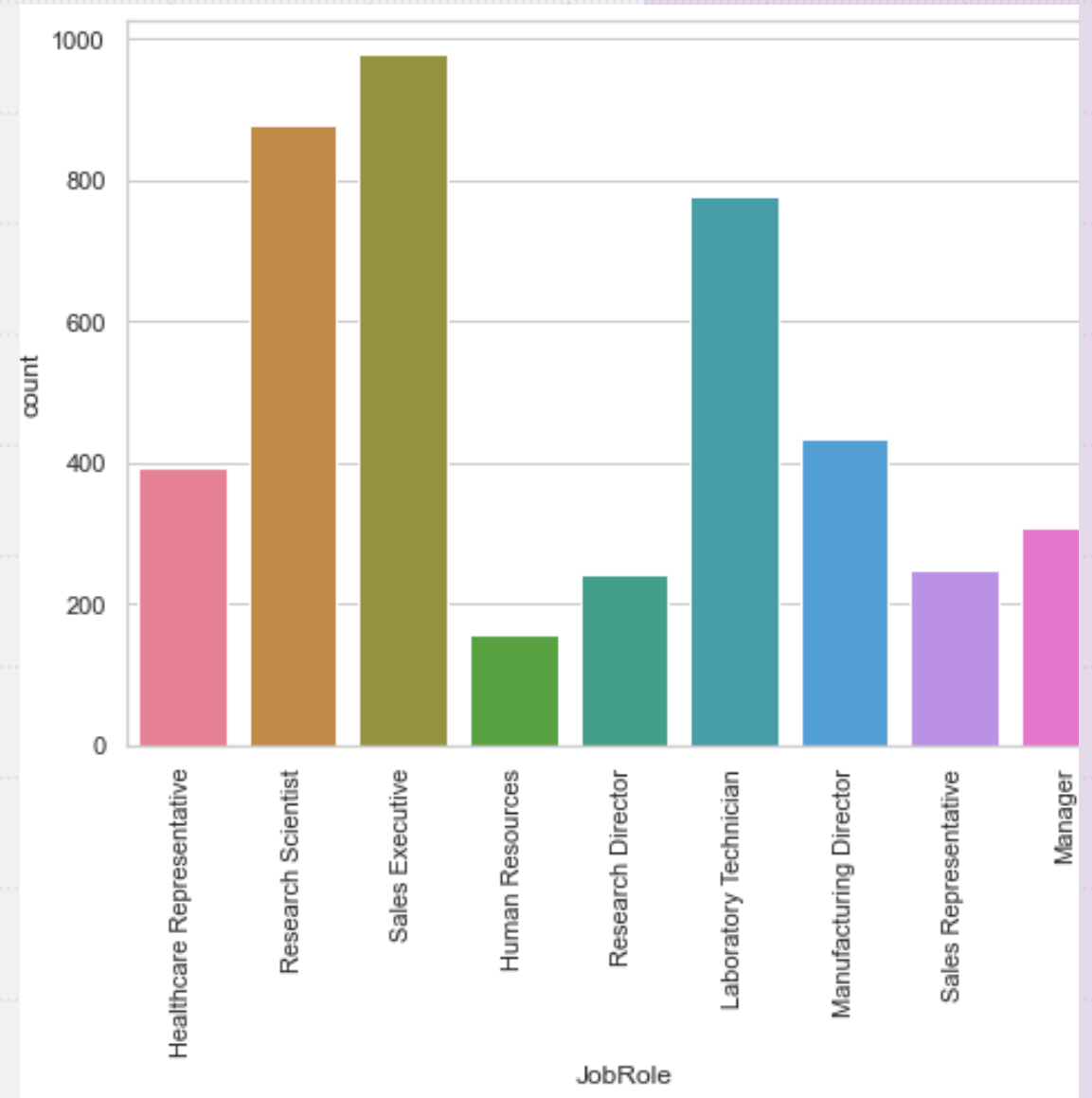
- The maximum attritions is seen where the employee's have the education field in Life Sciences followed by Medical.
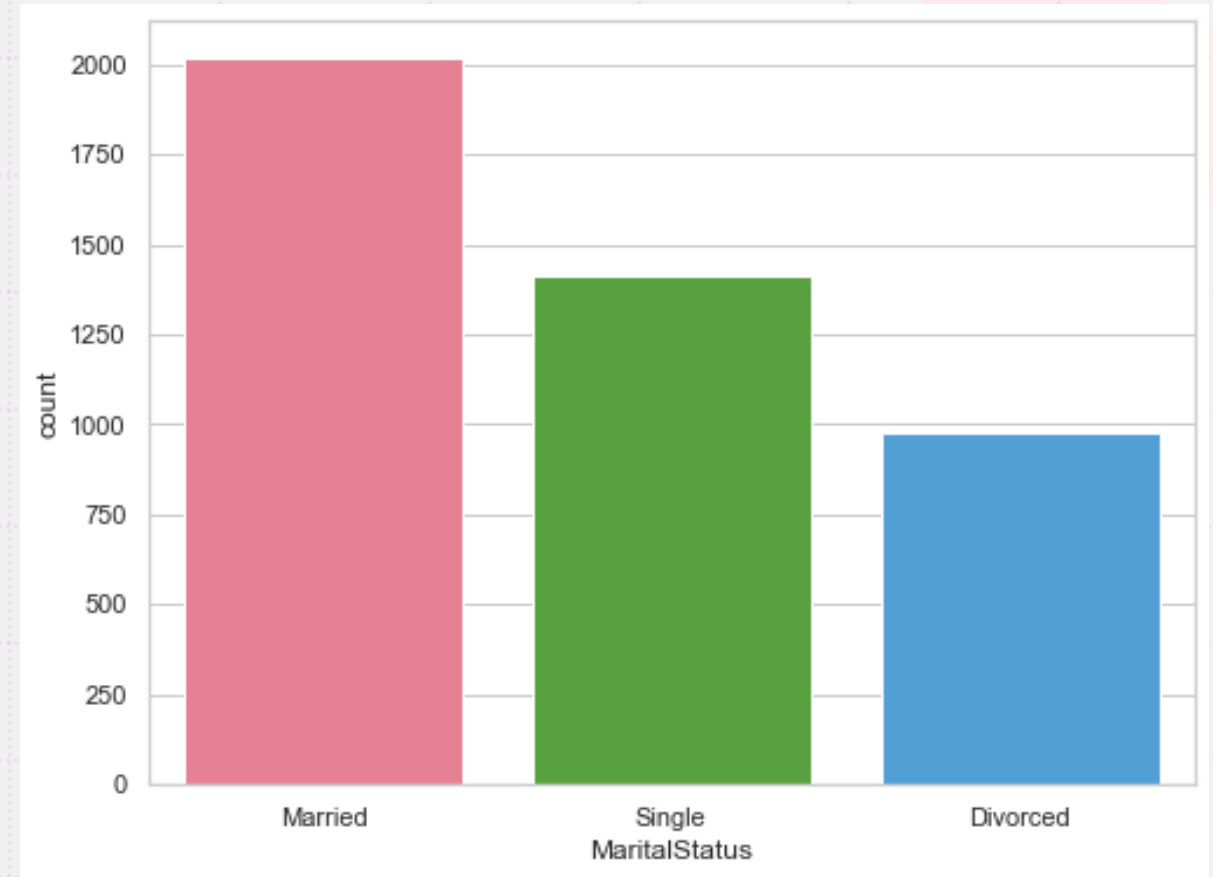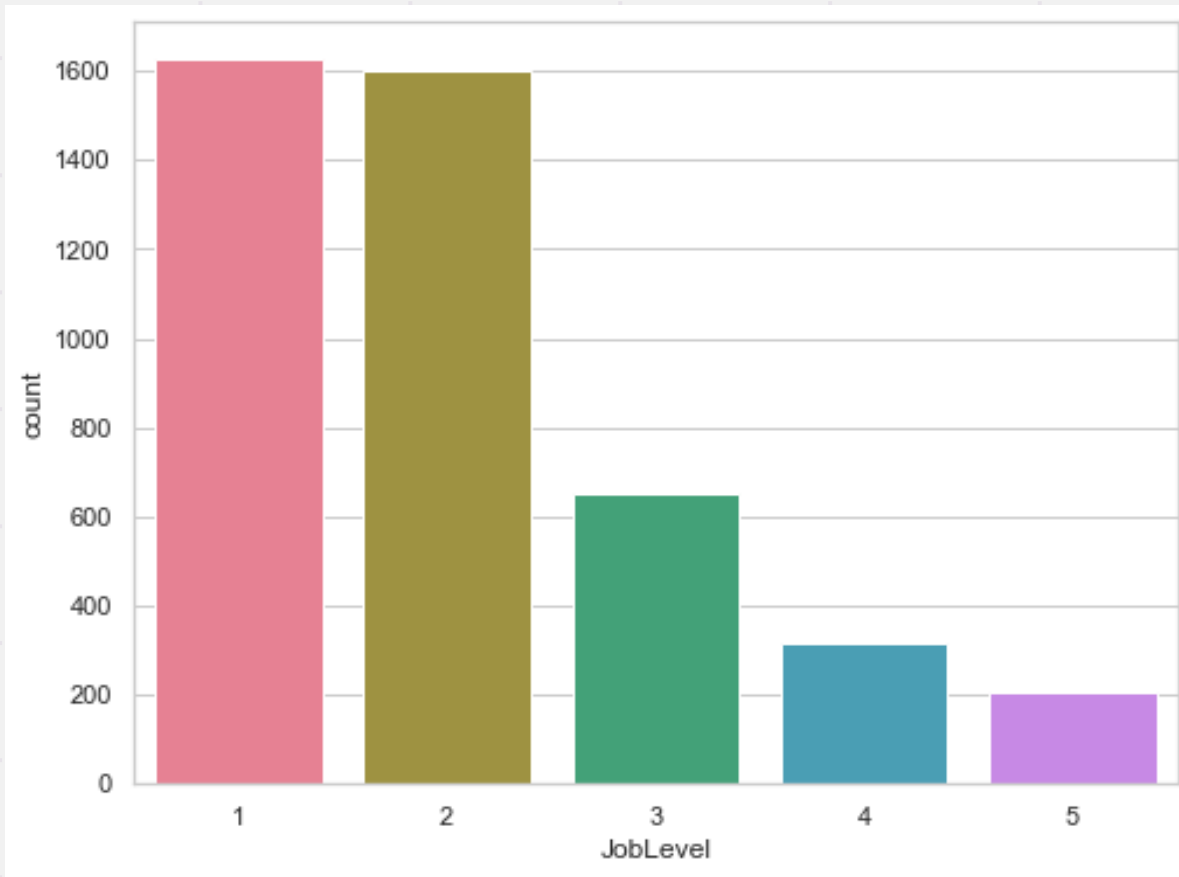
- Over here the employees that worked in the Research & Development Department have most attritions and very less attritions are seen from the Human Resources department

- We can see that most employees' roles who left is Sales Executive followed by Research Scientist and Laboratory Technician
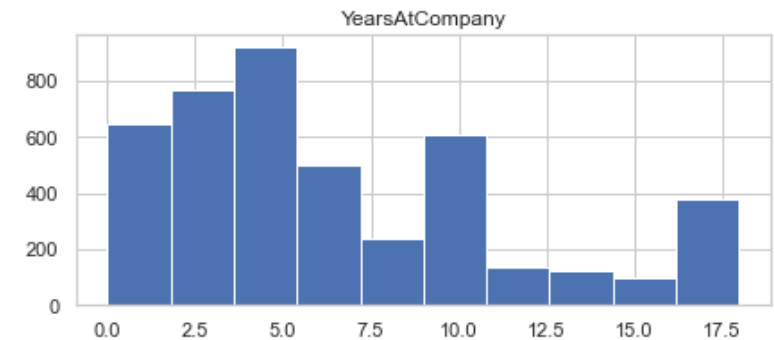- Most employees who left the company is male

- The most attritions of the employees can be seen in Job Level 1 followed by job level 2.

- The Married employees are most likely to be seen leaving the company.

# UNIVARIATE ANALYSIS



- Age Feature Distribution is almost Normal Distribution.

- Most of the Columns are in Skew Distribution form except Age.

# BIVARIATE ANALYSIS

# BIVARIATE ANALYSIS :
## using Boxplots
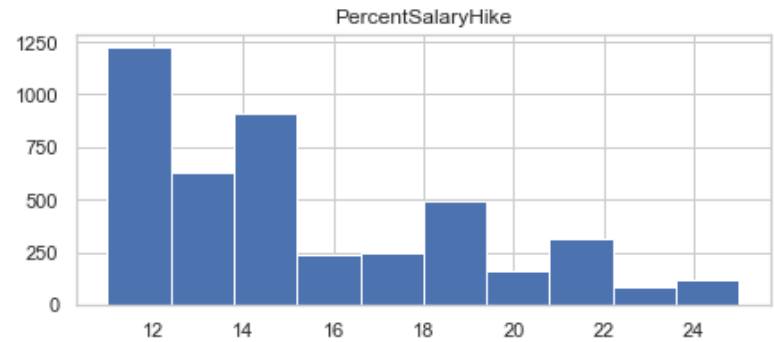
Boxplot shows the distribution of quantitative data that facilitates comparisons between variables or across levels of a categorical variable.

**Multivariate Analysis:** Using Pair plots for numerical variables

# Correlation: Heatmap



Correlation matrix among variables

# ANOVA TEST (Analysis of Variance)

It gives the relationship between continuous and categorical columns whether they are statistically significant or not.

H0- Null hypothesis means the variables are not correlated:

1. If Small P-Value < 5% (0.05) means the variables are correlated. Null hypothesis H0 is rejected.

2. If Large P-Value > 5% (0.05) means the variables are not correlated and the Null hypothesis H0 is accepted.

Result for Anova test : Variables which are not correlated with Attrition

• Hrs
• JobInvolvement
• PerformanceRating
• DistanceFromHome
• JobLevel
• StockOptionLevel

P-value > 0.05 for these variables, hence we can reject these variable for further analysis.

# Chi-Square Test

Chi-square is a quantitative measure used to determine whether relationship exists between two categorical variables.

H0- Null hypothesis means the variables are not correlated:

1. If Small P-Value < 5% (0.05) means the variables are correlated. Null hypothesis H0 is rejected
2. If Large P-Value > 5% (0.05) means the variables are not correlated and the Null hypothesis H0 is accepted

The variables (p-value < 0.05) that are found significant from Chi-square test:

```
['BusinessTravel_Non-Travel', 'BusinessTravel_Travel_Frequently',
'BusinessTravel_Travel_Rarely', 'Department_Human Resources',
 'EducationField_Human Resources', 'EducationField_Technical Degree',
'JobRole_Manufacturing Director', 'JobRole_Research Director',
 'MaritalStatus_Divorced', 'MaritalStatus_Married', 'MaritalStatus_Single']
```

The variables that were found with large p-value i.e. p-value > 0.05. So, we can reject these variables.

```
['JobRole_Sales Representative', 'JobRole_Sales Executive', 'JobRole_Research
Scientist', 'JobRole_Manager',
 'JobRole_Laboratory Technician', 'JobRole_Human Resources', 'JobRole_Healthcare
Representative', 'Gender_Male', 'Gender_Female',
'EducationField_Other', 'EducationField_Medical', 'EducationField_Marketing',
 'EducationField_Life Sciences', 'Department_Sales', 'Department_Research &
Development']
```
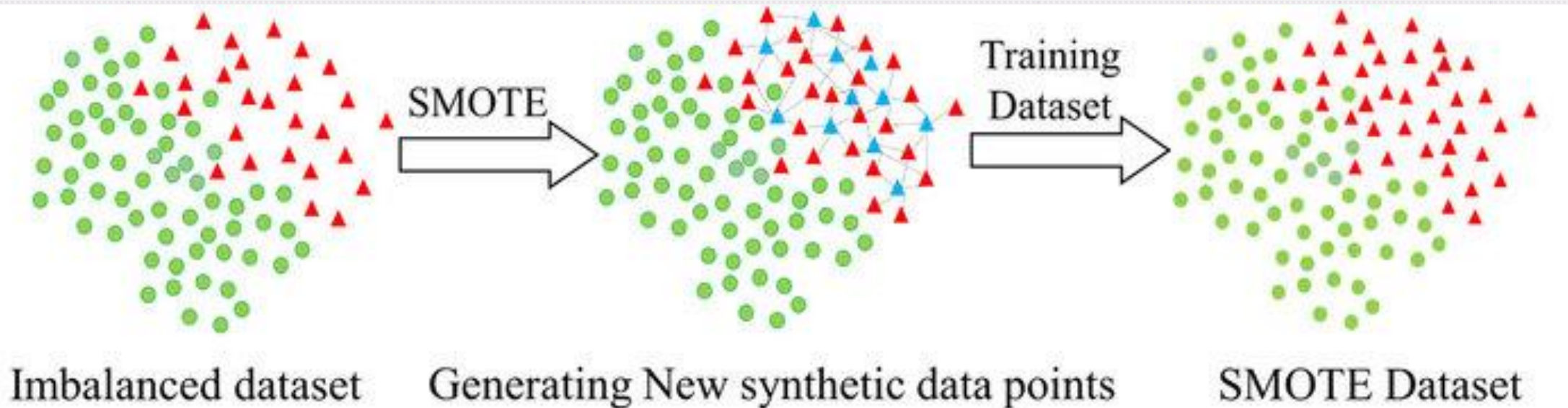
# LOGISTIC REGRESSION MODEL BUILDING

For model building, Scikit-Learn model is used to predict the attrition rate of the employees.

SMOTE technique with oversampling is used to balance the dataset.



Imbalanced dataset    Generating New synthetic data points    SMOTE Dataset

# OBSERVATIONS

- ❖ The accuracy of logistic regression classifier is coming around 70%.

- ❖ Root Mean Square Error (RMSE) of logistic regression is 0.54.

- ❖ The ROC AUC score for logistic regression model is approximately 62%



```
CLASSIFICATION REPORT
==========================================================================

                precision      recall    f1-score     support

          0         0.92        0.71        0.80        1109
          1         0.31        0.70        0.43         214

   accuracy                                 0.70        1323
  macro avg         0.62        0.70        0.62        1323
weighted avg        0.83        0.70        0.74        1323


==========================================================================

Accuracy of logistic regression classifier on test set: 0.70
==========================================================================

RMSE: 0.543636
==========================================================================
```

# CLASSIFIER MODELS TO ACHIEVE BETTER ACCURACY

❖ *Support Vector Machine Classifier (SVM)* : SVC is a linear model and creates a line or a hyperplane which separates the data into classes.

❖ *Decision Tree Classifier* : DTs are a non-parametric supervised learning method used for classification and regression.

❖ *Gradient Boosting* : Gradient boosting is another CART ensemble algorithm which relies on the assumption that the best possible next model.

❖ *Random Forest Classifier* : The random forest is an ensemble learning algorithm which builds upon the Classification And Regression Tree (CART) paradigm. It uses bagging and feature randomness

Below are the results from different classification models to achieve a better accuracy and rmse :-

| MODEL | ACCURACY | RMSE |
|---|---|---|
| Logistic Regression Classifier | 70% | 0.54% |
| Support Vector Machine Classifer | 82% | 0.41% |
| Decision Tree Classifier | 97% | 0.17% |
| Gradient Boosting | 78% | 0.46% |
| Random Forest Classifier | 99% | 0.08% |

# Insights and Business Recommendations

**INSIGHTS :**

❖ From feature importance, we find that Age appears to be the most important feature followed by the employee's Total Working Years, Monthly Income and Years At Company.

❖ Digging further we found also that the employees who left the company worked in the Research & Development department and Sales department and more specifically those who had Sales executive, Research scientist and Laboratory technicians as their Job role.

❖ We observed that most of the employees who left the company has normal Work Life Balance and employees with education level 3 and 4 had most attritions and the least attritions were seen with employees having education level 5.

❖ Most attritions is seen where the employee is travelling rarely.

❖ If we go by Random forest or Decision tree approach for model building, we will yield a nice accuracy and less root mean square error. Here we are getting a 99% accuracy and less than 10 % Rmse from Random forest model approach.

**RECOMMENDATIONS :**

❖ So my recommendation is to check the job satisfaction and environment satisfaction to lower the risk of attrition.

❖ Employees who didn't get promoted since last 5 years have higher risk of attrition.

❖ More years an employee spends with the manager lesser the risk of attrition.

❖ If the employee travels frequently then it can lead to a higher risk of attrition.

# THANK YOU