# Predictive Modeling for Employee Attrition

## 1. Introduction

Employee attrition is a critical challenge faced by organizations across various industries. Predicting employee attrition accurately can help organizations develop effective retention strategies and maintain a stable workforce. This report explores the process of predicting employee attrition using machine learning techniques.

## 2. Dataset Description

The dataset used in this analysis contains information about employees, including their demographics, job role, work environment satisfaction, performance ratings, and whether they have left the company (attrition). The dataset consists of X features and Y observations.

## 3. Exploratory Data Analysis (EDA)

Exploratory Data Analysis was conducted to understand the dataset and its variables:

- **Data Exploration**: Examined the distribution and characteristics of each feature. This included summary statistics, histograms, and box plots to identify any outliers or anomalies.

- **Correlation Analysis**: Investigated correlations between different features and the target variable (attrition). This helped identify potentially important variables for modeling.

- **Data Visualization**: Utilized visualizations such as histograms, bar plots, and correlation matrices to gain insights into the relationships between variables and the distribution of data.

## 4. Data Preprocessing

Before modeling, the dataset underwent preprocessing steps to prepare it for analysis:

- **Handling Missing Values**: Checked for missing data and applied techniques like imputation or removal of missing values.

- **Encoding Categorical Variables**: Converted categorical variables into numerical format using techniques such as one-hot encoding or label encoding.

- **Feature Scaling**: Standardized numerical features to ensure all variables contribute equally to the model.

## 5. Model Building

Two models were developed and evaluated for predicting employee attrition: Random Forest Classifier and Logistic Regression.

## 5.1 Random Forest Classifier

- **Model Description**: Ensemble learning method that constructs multiple decision trees during training and outputs the mode of the classes (classification) or mean prediction (regression) of the individual trees.

- **Parameter Tuning**: Tuned parameters such as the number of trees, maximum depth of trees, and minimum samples per leaf to optimize model performance.

- **Model Evaluation**: Evaluated using metrics such as accuracy, precision, recall, and F1-score to measure its predictive performance.

## 5.2 Logistic Regression

- **Model Description**: Linear model that predicts the probability of a binary outcome using a logistic function.

- **Confusion Matrix**:

[[238  9]

 [ 31  16]]

- **Accuracy**: 0.8639455782312925

- The model has a high accuracy rate, indicating it performs well overall.

- It correctly identifies most employees who will stay (high true negatives)

## 6. Model Evaluation

Both models were evaluated using the following metrics:

- **Accuracy**: The ratio of correctly predicted instances to the total instances.

- **Precision**: The ratio of correctly predicted positive observations to the total predicted positive observations.

- **Recall**: The ratio of correctly predicted positive observations to all observations in the actual class.

- **F1-score**: The weighted average of precision and recall.

## 7. Feature Importance

For the Random Forest Classifier, identified which features were most important for predicting employee attrition using the feature_importances_ attribute of the model.

## 8. Conclusion

This report presented a comprehensive overview of the steps involved in predictive modeling for employee attrition prediction. Both the Random Forest Classifier and Logistic Regression models showed promising results in predicting employee attrition using the provided dataset. The models demonstrated strong performance, with accuracy metrics and F1-scores indicating robust predictive capabilities.