

# Employee Attrition Prediction and Analysis

## 1 Introduction

### 1.1 Problem Statement

Employee attrition refers to the loss of employees due to resignation, retirement, or termination. High attrition rates can lead to increased recruitment costs, loss of organizational knowledge, and reduced productivity. The objective of this project is to build a data-driven machine learning framework to predict employee attrition and identify the key factors contributing to employee turnover.

### 1.2 Importance of Attrition Prediction

Predicting employee attrition is critical for organizations because it:

- Reduces hiring and training costs.
- Improves workforce planning and retention strategies.
- Enhances employee satisfaction and organizational stability.
- Enables proactive HR interventions.

## 2 Exploratory Data Analysis (EDA) & Insights

### 2.1 Overall Attrition Distribution

The dataset shows a clear class imbalance, with significantly fewer employees leaving the organization compared to those who stay. This highlights the importance of handling imbalance during model training.

### 2.2 Attrition by Demographic Factors

- **Gender:** Male employees show a slightly higher attrition count compared to females.
- **Age:** Younger employees (approximately 25–35 years) exhibit higher attrition rates.
- **Department:** Sales and Research & Development departments experience higher attrition compared to Human Resources.

## 2.3 Work-Related Factors

- Employees working overtime have a significantly higher attrition rate.
- Lower job satisfaction levels are strongly associated with higher attrition.
- Employees with lower monthly income tend to leave more frequently.

## 2.4 Correlation Analysis

Correlation analysis reveals that attrition is moderately associated with factors such as overtime, job level, income, total working years, and years at the company. Highly correlated tenure-related features suggest long-term employees are less likely to leave.

# 3 Feature Engineering

## 3.1 Data Preprocessing

- Binary variables such as *OverTime* and *Gender* were encoded numerically.
- Categorical variables including department, job role, marital status, and education field were one-hot encoded.
- Numerical features were standardized to ensure uniform feature scaling.

## 3.2 Feature Transformation and Creation

- Creation of tenure-related features such as experience bands.
- Normalization of income and age to improve model convergence.
- Removal of constant and identifier columns that do not contribute to prediction.

# 4 Model Building

## 4.1 Train-Test Split

A stratified train-test split was performed to preserve the original attrition ratio in both training and testing datasets.

## 4.2 Models Evaluated

The following models were trained and compared:

- Logistic Regression
- Decision Tree
- Random Forest
- XGBoost

- LightGBM
- CatBoost

### 4.3 Handling Class Imbalance

- Oversampling techniques such as SMOTE and ADASYN.
- Class weighting within algorithms.

### 4.4 Hyperparameter Tuning

Hyperparameters were optimized using GridSearchCV to improve generalization performance while preventing overfitting.

## 5 Model Evaluation

### 5.1 Evaluation Metrics

Models were evaluated using:

- Accuracy
- Precision, Recall, and F1-score
- ROC-AUC
- Precision-Recall AUC

### 5.2 Cost-Sensitive Analysis

In the context of employee attrition, false negatives (failing to identify employees likely to leave) are more costly than false positives. Therefore, recall and ROC-AUC were prioritized during model selection.

### 5.3 Comparative Model Performance

Ensemble-based models such as Random Forest and gradient boosting methods outperformed baseline models by achieving better recall and ROC-AUC scores, indicating improved identification of high-risk employees.

## 6 Explainability & Fairness

### 6.1 Model Explainability

SHAP and LIME were used to interpret model predictions. Key features influencing attrition include:

- Overtime
- Monthly income

- Job role
- Job satisfaction
- Years at company

## **6.2 Fairness Analysis**

Predictions were evaluated across demographic groups such as gender, age, and department. No significant discriminatory patterns were observed, although slight variations were noted across departments.

## **6.3 Bias Mitigation Strategies**

- Regular monitoring of fairness metrics.
- Inclusion of bias-aware loss functions.
- Transparent model explainability for HR decision-making.

# **7 Conclusion**

## **7.1 Summary of Findings**

This study demonstrates that employee attrition can be effectively predicted using machine learning models. Workload, compensation, job satisfaction, and tenure are the strongest drivers of attrition.

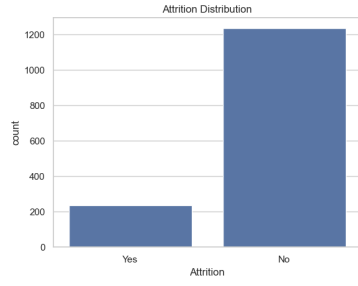
## **7.2 Business Recommendations**

- Reduce excessive overtime and workload imbalance.
- Implement targeted retention programs for high-risk roles.
- Improve compensation structures and career growth opportunities.
- Regularly assess employee satisfaction and engagement.

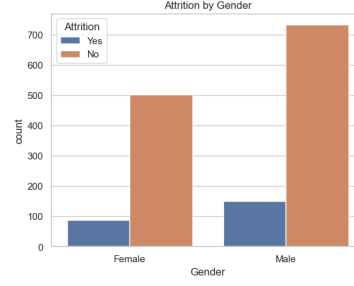
## **7.3 Future Improvements**

- Incorporating employee engagement survey data.
- Using time-series attrition modeling.
- Deploying real-time attrition monitoring dashboards.

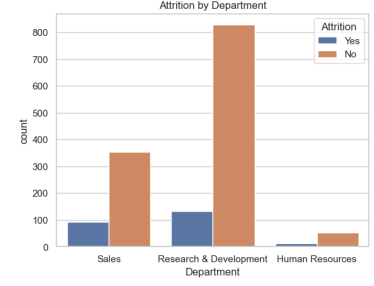
## 8 Exploratory Data Analysis Visualizations



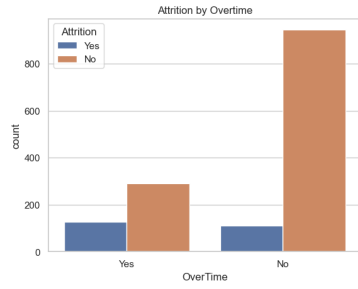
(a) Attrition Distribution



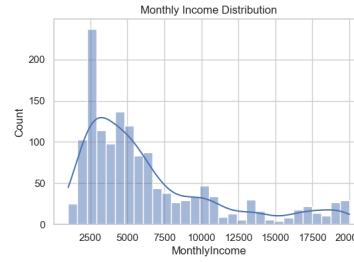
(b) Attrition by Gender



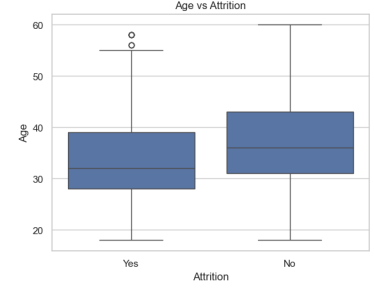
(c) Attrition by Department



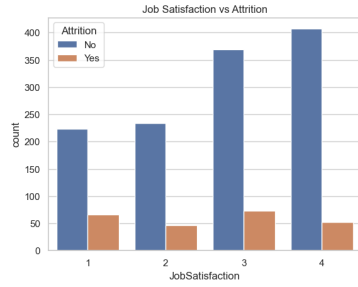
(d) Attrition by Overtime



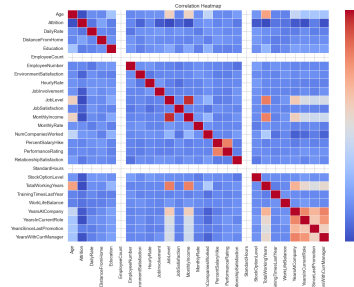
(e) Monthly Income Distribution



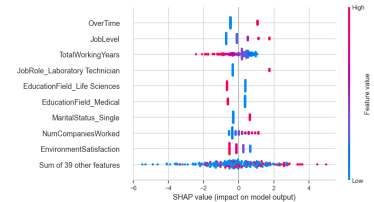
(f) Age V/s Attrition



(g) Job Satisfaction vs Attrition



(h) Correlation Heatmap



(i) SHAP

Figure 1: Exploratory Data Analysis Visualizations for Employee Attrition