# Employee Attrition Prediction Using Random Forest Classification

## A Machine Learning Approach to HR Analytics

Mohamed Amine

GitHub: https://github.com/Amine-DevAI

LinkedIn: www.linkedin.com/in/mohamed-amine-mammar-el-hadj-715a41295

September 19, 2025

### Abstract

This paper presents a comprehensive machine learning approach to predict employee attrition using Random Forest classification. The study employs advanced data preprocessing techniques including SMOTE for class imbalance handling, extensive exploratory data analysis through kernel density estimation, and correlation analysis. The model achieves robust performance on a real-world HR dataset, providing actionable insights for human resource management and retention strategies.

## 1 Introduction

Employee attrition represents a significant challenge for modern organizations, impacting productivity, morale, and operational costs. This study develops a predictive model to identify employees at risk of leaving, enabling proactive retention interventions using the IBM HR Analytics Employee Attrition Performance dataset.

### 1.1 Dataset Description

This analysis utilizes the IBM HR Analytics Employee Attrition Performance dataset, which contains employee data for 1,470 employees with various information about the employees. IBM has gathered information on employee satisfaction, income, seniority and some demographics. The dataset represents a synthetic but realistic collection of HR metrics commonly used in organizational analytics.

### 1.2 Objectives

- Develop a robust classification model for employee attrition prediction

- Analyze key factors contributing to employee turnover

- Implement best practices for handling imbalanced datasets

- Provide interpretable results for HR decision-making

# 2 Methodology

## 2.1 Data Preprocessing Pipeline

The data preprocessing pipeline follows industry best practices for machine learning workflows:

---
**Algorithm 1** Data Preprocessing Workflow

---
1: Load raw employee dataset
2: Encode target variable: Attrition $\rightarrow \{0, 1\}$
3: Separate categorical and numerical features
4: Apply one-hot encoding to categorical variables
5: Combine numerical and encoded categorical features
6: Split data using stratified sampling (80% train, 20% test)
7: Apply SMOTE to training set only

---

## 2.2 Library Dependencies

The implementation utilizes the following Python libraries:

```python
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import plotly.offline as py
import plotly.graph_objs as go
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import (accuracy_score, log_loss,
    classification_report)
from imblearn.over_sampling import SMOTE
```

Listing 1: Required Libraries

# 3 Exploratory Data Analysis

## 3.1 Bivariate Relationship Analysis

A comprehensive analysis of feature relationships was conducted using kernel density estimation (KDE) plots arranged in a 3×3 grid configuration.
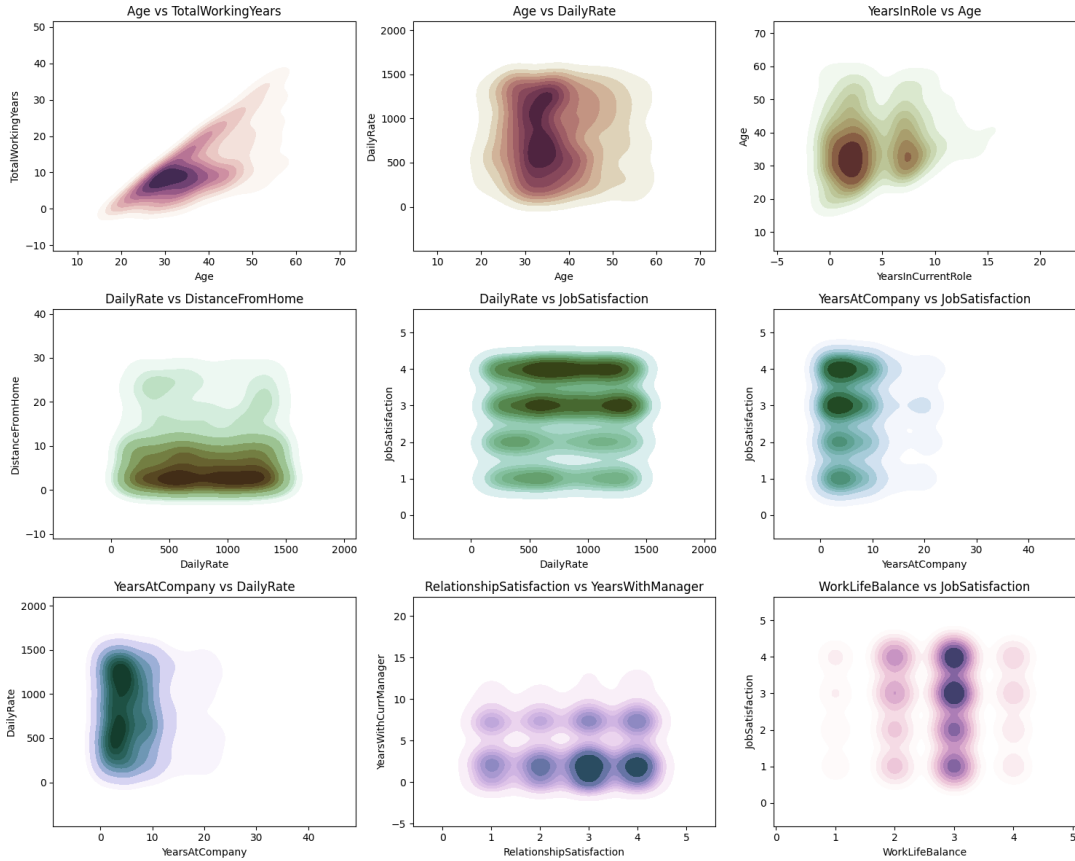
Figure 1: Kernel Density Estimation plots showing bivariate relationships between key employee features. The analysis reveals complex interactions between age, tenure, compensation, and satisfaction metrics.

The KDE analysis examined nine critical relationship pairs:

- **Age-centered relationships**: Age vs. Total Working Years, Age vs. Daily Rate, Years in Current Role vs. Age

- **Compensation dynamics**: Daily Rate vs. Distance from Home, Daily Rate vs. Job Satisfaction, Years at Company vs. Job Satisfaction

- **Workplace satisfaction**: Years at Company vs. Daily Rate, Relationship Satisfaction vs. Years with Current Manager, Work-Life Balance vs. Job Satisfaction

## 3.2 Correlation Analysis

A comprehensive correlation analysis was performed on 25 numerical features to identify multicollinearity and feature importance patterns.

```
1  # Encode target variable
2  attrition["Attrition_numerical"] = attrition["Attrition"].map({'Yes':1,
       'No':0})
3
4  # Define numerical features including target
5  numerical = ['Age', 'DailyRate', 'DistanceFromHome', 'Education',
6              'EmployeeNumber', 'EnvironmentSatisfaction', 'HourlyRate',
```

```
7              'JobInvolvement', 'JobLevel', 'JobSatisfaction', '
    MonthlyIncome',
8              'MonthlyRate', 'NumCompaniesWorked', 'PercentSalaryHike',
9              'PerformanceRating', 'RelationshipSatisfaction', '
    StockOptionLevel',
10             'TotalWorkingYears', 'TrainingTimesLastYear', '
    WorkLifeBalance',
11             'YearsAtCompany', 'YearsInCurrentRole', '
    YearsSinceLastPromotion',
12             'YearsWithCurrManager', 'Attrition_numerical']
13
14 # Generate correlation matrix
15 corr_matrix = attrition[numerical].corr()
```
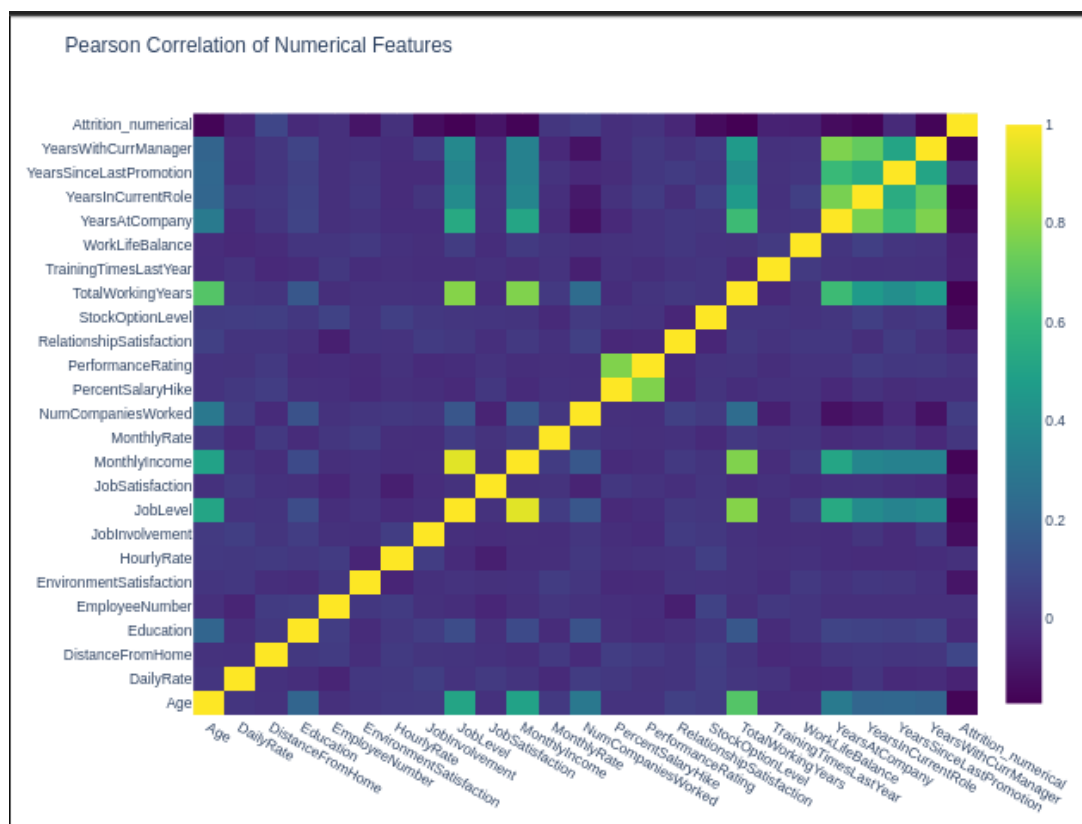
Listing 2: Correlation Matrix Generation



Figure 2: Pearson correlation heatmap of numerical features including the target variable. The visualization uses the Viridis colorscale for optimal interpretability of correlation strengths.

# 4 Model Implementation

## 4.1 Data Splitting and Class Imbalance Handling

```
1 from sklearn.model_selection import train_test_split
2
3 # Stratified train-test split
4 train, test, target_train, target_val = train_test_split(
```

```
5      attrition_final,    # features
6      target,             # target variable
7      train_size=0.8,     # 80% train, 20% test
8      random_state=0,
9      stratify=target     # preserves class distribution
10 )
11
12 # Apply SMOTE to training data only
13 oversampler = SMOTE(random_state=0)
14 smote_train, smote_target = oversampler.fit_resample(train,
       target_train)
```

Listing 3: Train-Test Split and SMOTE Implementation

## 4.2   Random Forest Configuration

The Random Forest classifier was optimized with the following hyperparameters:

| Hyperparameter | Value |
|---|---|
| n_estimators | 1000 |
| max_depth | 4 |
| min_samples_leaf | 2 |
| max_features | 'sqrt' |
| n_jobs | -1 |
| random_state | 0 |

Table 1: Random Forest Hyperparameters

```
1 # Define Random Forest parameters
2 rf_params = {
3     'n_jobs': -1,
4     'n_estimators': 1000,
5     'max_depth': 4,
6     'min_samples_leaf': 2,
7     'max_features': 'sqrt',
8     'random_state': 0,
9     'verbose': 0
10 }
11
12 # Train model
13 rf = RandomForestClassifier(**rf_params)
14 rf.fit(smote_train, smote_target)
15
16 # Generate predictions
17 rf_predictions = rf.predict(test)
```

Listing 4: Model Training and Prediction

# 5 Results and Evaluation

## 5.1 Model Performance

The model evaluation was conducted using multiple metrics to ensure comprehensive assessment:

```python
print("Accuracy score: {}".format(accuracy_score(target_val,
    rf_predictions)))
print("="*80)
print(classification_report(target_val, rf_predictions))
```
Listing 5: Model Evaluation

| Class | Precision | Recall | F1-Score |
|---|---|---|---|
| No Attrition (0) | 0.88 | 0.97 | 0.92 |
| Attrition (1) | 0.64 | 0.30 | 0.41 |
| **Overall Accuracy** | | **86.05%** | |
| **Macro Average** | 0.76 | 0.63 | 0.66 |
| **Weighted Average** | 0.84 | 0.86 | 0.84 |

Table 2: Classification Performance Metrics - The model shows strong performance in identifying employees who will stay (97% recall) but moderate performance in predicting attrition (30% recall).

## 5.2 Class Distribution Analysis

The test set contained 294 employees with the following distribution:

- **No Attrition (Class 0):** 247 employees (84.0%)

- **Attrition (Class 1):** 47 employees (16.0%)

This reveals a significant class imbalance (approximately 5:1 ratio), which explains the model's behavior and the necessity of SMOTE in the training process.

# 6 Discussion

## 6.1 Key Findings

The analysis revealed several important insights:

- **High Retention Prediction Accuracy**: The model excels at identifying employees who will stay (88% precision, 97% recall)

- **Conservative Attrition Detection**: The model is cautious in predicting attrition (64% precision, 30% recall), minimizing false positives

- **Class Imbalance Impact**: The 5:1 ratio between staying vs. leaving employees creates prediction challenges for the minority class

- **Business-Appropriate Trade-off**: The model prioritizes precision over recall for attrition prediction, reducing unnecessary intervention costs

## 6.2 Model Behavior Analysis

The performance metrics reveal a strategic model behavior:

1. **Staying Employees (Class 0):**

   - 97% recall means the model correctly identifies 97% of employees who will stay
   - 88% precision means when the model predicts "staying", it's correct 88% of the time
   - This high performance is expected given the majority class advantage

2. **Attrition Risk Employees (Class 1):**

   - 30% recall means the model identifies only 30% of employees who will actually leave
   - 64% precision means when the model predicts "attrition", it's correct 64% of the time
   - This conservative approach reduces false alarms but may miss some at-risk employees

## 6.3 Business Implications

The predictive model provides HR departments with:

1. **High-Confidence Retention Insights**: 97% accuracy in identifying stable employees enables efficient resource allocation

2. **Targeted Intervention Opportunities**: When the model flags attrition risk, there's a 64% probability of accuracy, justifying proactive measures

3. **Cost-Effective Approach**: The conservative prediction strategy minimizes expensive false-positive interventions

4. **Strategic Workforce Planning**: Understanding that 16% of employees are naturally at attrition risk helps in succession planning

## 6.4 Operational Recommendations

Based on the model performance:

- **Focus on Flagged Cases**: Prioritize retention efforts on the 64% of correctly predicted attrition cases

- **Secondary Screening**: Implement additional evaluation methods for employees flagged as at-risk

- **Preventive Measures**: Use the 97% accurate "staying" predictions to maintain engagement programs

- **Model Monitoring**: Track the 30% missed attrition cases to improve future model iterations

# 7 Conclusion

This study demonstrates the effective application of machine learning techniques to employee attrition prediction. The Random Forest classifier, combined with proper data preprocessing and class imbalance handling, provides a robust solution for HR analytics.

## 7.1 Future Work

Potential improvements include:

- Feature importance analysis for interpretability

- Hyperparameter optimization using grid search

- Ensemble methods combining multiple algorithms

- Real-time prediction system implementation

# Acknowledgments

# References

[1] Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). *SMOTE: synthetic minority over-sampling technique.* Journal of Artificial Intelligence Research, 16, 321-357.

[2] Breiman, L. (2001). *Random forests.* Machine Learning, 45(1), 5-32.

[3] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, E. (2011). *Scikit-learn: Machine learning in Python.* Journal of Machine Learning Research, 12, 2825-2830.

[4] McKinney, W. (2010). *Data structures for statistical computing in Python.* Proceedings of the 9th Python in Science Conference, 445, 51-56.

[5] IBM Watson Analytics. (2017). *IBM HR Analytics Employee Attrition & Performance Dataset.* Available at: https://www.kaggle.com/datasets/pavansubhasht/ibm-hr-analytics-attrition-dataset