

# Early Prediction of Employee Attrition Using Data Mining Techniques

Md. Shamsul rahat Chy  
Computer Science and Engineering  
Brac University  
Dhaka, Bangladesh  
shamsul.rahata.chy@g.bracu.ac.bd

**Abstract**—This research paper focuses on predicting employee attrition using data mining techniques, which is crucial for reducing turnover costs and improving organizational stability. Utilizing a dataset of 1,470 employees, which includes 35 attributes such as age, job role, and monthly income, various predictive models like Logistic Regression, Support Vector Machine, and Neural Networks are employed. The methodology involves preprocessing the data to handle numeric and categorical variables and implementing a split of 70% training and 30% testing data. Initial model evaluations are based on accuracy, ROC AUC, and F1 scores. The goal is to identify the most effective model for predicting attrition, enabling proactive retention strategies. This approach not only aids in understanding the factors influencing attrition but also serves as a guide for implementing targeted interventions to retain valuable human resources.

**Index Terms**—Employee Attrition, Data Mining, Logistic Regression, SVM, Neural Networks, Predictive Analytics

## I. INTRODUCTION

Employee attrition—the phenomenon of workers leaving an organization—is a natural but complex issue that affects all businesses, regardless of size or industry. It encompasses voluntary departures, involuntary terminations, and retirements, each of which can destabilize operational continuity and incur substantial costs related to recruitment and training of new personnel. Understanding and predicting employee attrition is therefore critical for maintaining workforce stability and optimizing resource allocation.

Recent advancements in data mining and predictive analytics offer powerful tools for tackling this challenge. By analyzing historical data on employee behavior and characteristics, organizations can identify potential patterns and triggers for attrition. This research paper applies various data mining techniques, including Logistic Regression, Support Vector Machines (SVM), and Neural Networks, to a dataset of 1,470 employees to predict the likelihood of employee departures. Through this analysis, the study aims to highlight significant predictors of attrition and evaluate the effectiveness of different models in forecasting such events.

By proactively predicting attrition, organizations can implement targeted retention strategies, ultimately reducing turnover and enhancing overall organizational performance. This paper explores these possibilities, aiming to contribute valuable insights into the strategic management of human resources.

## II. DATA ANALYSIS

The dataset used in this study encompasses detailed records of 1,470 employees, covering a diverse range of 35 attributes, which include demographic details like age and gender, job-specific information such as role, daily rate, and monthly income, and psychological metrics like job satisfaction. This comprehensive dataset allows for a multi-faceted analysis of factors that could potentially influence employee attrition.

### A. Descriptive Statistics

The age of employees in the dataset is fairly normally distributed, with a mean around the mid-30s, reflecting a workforce that ranges from young entrants to seasoned professionals. This spread suggests different generational attitudes and values, which can influence attrition differently. Other numeric data such as daily rates and monthly income are also analyzed to understand their correlation with employee turnover.

### B. Attrition Distribution

The attrition variable, which is the focus of this study, shows that a significant minority of the workforce has left the company. A binary classification ('Yes' for attrition, 'No' for non-attrition) helps in clearly distinguishing between employees who have left and those who remain, serving as a foundation for further predictive modeling.

### C. Job Satisfaction Analysis

Job satisfaction levels, categorized from 1 to 4, reveal that a majority of employees report high satisfaction (levels 3 and 4), suggesting a generally positive work environment. However, the presence of employees with lower satisfaction levels (1 and 2) indicates areas where there might be room for improvement to reduce potential attrition.

The initial data analysis provides a robust understanding of the dataset's structure and the key variables that may influence attrition. It establishes a clear path for the preprocessing steps needed to prepare the data for predictive modeling. By identifying key trends and patterns, this analysis sets the stage for developing effective models to predict and, consequently, mitigate employee attrition.

### III. DATA PREPROCESSING

The data preprocessing stage is crucial for preparing the raw dataset for effective modeling and analysis. This study's dataset, consisting of 1,470 employees across 35 attributes, required several preprocessing steps to ensure data quality and relevance for predicting employee attrition.

#### A. Handling Missing Values

Initially, the dataset was thoroughly examined for missing values across all columns. Fortunately, this particular dataset did not contain any missing data, eliminating the need for imputation techniques which can sometimes introduce bias or affect the dataset's integrity.

#### B. Encoding Categorical Variables

Many attributes in the dataset, such as job role and department, are categorical. These were converted into numerical format using one-hot encoding. This method transforms each categorical value into a new binary column, ensuring that the machine learning algorithms can effectively interpret and process the data.

#### C. Feature Selection and Data Normalization

Feature selection was performed to identify the most relevant attributes that impact attrition. This involved analyzing correlation coefficients between each feature and the target variable (attrition), as well as employing automated feature selection techniques like Recursive Feature Elimination (RFE) to refine the feature set. This step is vital to enhance model performance by reducing complexity and avoiding overfitting.

Given the variety in the scale of the data—ranging from salary figures to binary variables—normalization was essential. Using the StandardScaler, numeric features were scaled to have zero mean and unit variance. This normalization ensures that no single feature dominates the model due to its scale, allowing for a fair evaluation of feature importance.

#### D. Splitting the Dataset

The final step involved splitting the data into training and testing sets. A stratified split was used to maintain the same proportion of attrition occurrences in both sets. Typically, the data was divided into 70% training and 30% testing segments. This division allows for the comprehensive training of the models while also setting aside a portion of the data to evaluate model performance unbiasedly.

### IV. METHODOLOGY

The methodology adopted for predicting employee attrition involves a structured approach that begins with the careful preparation of data to ensure its suitability for the subsequent analytical processes. Initially, the dataset undergoes a thorough review to confirm all data types are appropriate for analysis, which includes encoding categorical variables to numerical formats. This is critical for seamless integration into machine learning algorithms.

Following data preparation, the dataset is bifurcated into independent features (X) and the dependent target variable (y),

where the target is 'attrition'—indicating whether an employee has left the organization. The data is then split into a training set and a testing set using a stratified sampling method to maintain a uniform distribution of the target variable across both sets, typically allocating 70% of the data to training and 30% to testing. This step is essential for training the models on a comprehensive set of data while reserving a portion for unbiased evaluation.

Scaling of features is performed next using the StandardScaler to normalize the data, ensuring that no single feature disproportionately influences the model's outcomes due to varying scales. This normalization is especially crucial for distance-based models like the Support Vector Machine.

Once the data is scaled, various predictive models such as Logistic Regression, Random Forest, Support Vector Machine, and XGBoost are initialized. Each model is then trained on the scaled training data to learn from the patterns that suggest possible attrition. Initial evaluations of these models are conducted using accuracy and ROC AUC metrics on the testing set to gauge preliminary effectiveness.

The promising models from this initial evaluation undergo a phase of hyperparameter tuning, which optimizes their performance by adjusting parameters such as learning rates or tree depths. Techniques like grid search or random search are employed to explore the best combinations of these parameters.

After tuning, the models are trained once more with these optimized parameters to refine their predictive accuracy. This final training phase includes a detailed evaluation using a broader set of metrics, including the F1 score, to assess each model comprehensively. The outcomes are then visually represented through various methods like confusion matrices and ROC curves, facilitating a straightforward comparison of each model's effectiveness.

The methodology culminates in a detailed conclusion that draws from the comparative analysis of the models, discussing the implications of the findings for practical applications in human resource management and suggesting avenues for future research. This systematic approach ensures that the study not only identifies the most effective model for predicting employee attrition but also provides actionable insights that can be applied in real-world scenarios.

### REFERENCES