

# **EMPLOYEE ATTRITION PREDICTION SYSTEM**

*submitted in partial fulfillment of Minor Project for the award of the degree of*

**Bachelor of Technology**

in

**Artificial Intelligence and Data Science**

*Under the supervision of*

**Dr. Tina Dudeja**

**Assistant Professor**

Submitted by

**GARV CHANANA**

**04414811922**



**MAHARAJA AGRASEN INSTITUTE OF TECHNOLOGY**

**SECTOR-22, ROHINI, DELHI –110086**

**November 2025**

## **CANDIDATE’S DECLARATION**

I hereby declare that the work presented in this project report titled, “ **EMPLOYEE ATTRITION PREDICTION SYSTEM** ” submitted by me in the partial fulfillment of the requirement of the award of the degree of **Bachelor of Technology (B.Tech.)** Submitted in the Department of **Artificial Intelligence & Data Science**, Maharaja Agrasen Institute of Technology is an authentic record of my project work carried out under the guidance of **Dr. Vinay Kumar Saini** and mentor **Dr. Tina Dudeja**.

The matter presented in this project report has not been submitted either in part or full to any university or Institute for award of any degree.

Date: 21/11/2025

Place: Delhi

**GARV CHANANA**

**04414811922**

## **SUPERVISOR'S CERTIFICATE**

It is to certify that the Minor Project entitled “**EMPLOYEE ATTRITION PREDICTION SYSTEM**” which is being submitted by **Mr. GARV CHANANA** to the Maharaja Agrasen Institute of Technology, Rohini in the fulfillment of the requirement for the award of the degree of **Bachelor of Technology (B.Tech.)**, is a record of Bonafide project work carried out by him/her under my/ our guidance and supervision.

**Dr. Tina Dudeja**

**Assistant Professor**

**DR. VINAY KUMAR SAINI**

**H.O.D.  
Department of AI&DS**

## **ACKNOWLEDGEMENT**

I want to sincerely thank my faculty advisor **Dr. Vinay Kumar Saini** and mentor **Dr. Tina Dudeja** for their invaluable advice and unwavering support during the completion of this project. Their technical know-how and perceptive criticism have greatly influenced this work.

I also want to thank the department head, Dr. Vinay K. Saini, for creating a great learning environment, departmental support, and ongoing motivation.

I want to express my sincere gratitude to the entire faculty of the Department of AI & DS for their support and collaboration throughout this project. This project would not have been possible without the unwavering support, tolerance, and encouragement of my friends and family, for which I am also thankful.

Finally, I want to express my gratitude to everyone who helped make this project a success, whether directly or indirectly.

**GARV CHANANA**

**04414811922**

**AI&DS 7<sup>th</sup> Sem**

**Maharaja Agrasen Institute of Technology**

# TABLE OF CONTENTS

<b>Contents</b>	<b>Page No.</b>
Candidate's Declaration	ii
Supervisor's Certificate	iii
Acknowledgement	iv
Table of Contents	v
List of Figures and Tables	vi
Abstract	vii
Chapter 1: Introduction	1 – 4
Chapter 2: Literature Review	5 – 8
Chapter 3: Implementation of Proposed Model	9 – 20
Chapter 4: Experimental Results	21 – 29
Chapter 5: Conclusion and Discussion	30 – 34
Chapter 6: Business Insights and Practical Implications	35 – 36
Chapter 7: Limitation and Future Scope	37 – 38
References	39
Appendix – Code Snippets	40 – 43

# List of Figures

Figure No.	Figure Name	Page No.
Figure 3.1	Data Flow Diagram	13
Figure 3.2	Use Case Diagram	13
Figure 3.3	Flowchart of Employee Attrition	15
Figure 3.4	System Architecture of Employee Attrition Prediction	16
Figure 4.1	Distribution of Attrition	23
Figure 4.2	Histogram of Key Numerical Features	23
Figure 4.3	Box Plots of Key Numerical Features	24
Figure 4.4	Distribution of key Categorical Features	25
Figure 4.5	Comparison of Attrition with Age, Job Satisfaction, Monthly Income	26
Figure 4.6	Correlation Matrix of Independent Features	26
Figure 4.7	Feature Importance Analysis of key Features	30
Figure 4.8	ROC Curve Comparison	34
Figure 4.9	Confusion Matrix Comparison	35

## List of Tables

Table No.	Table Name	Page No.
Table 3.1	Hardware and Software Requirements	12
Table 4.1	Model Performance Metrics (Accuracy, Precision, Recall, F1-score)	26

# ABSTRACT

One of the biggest problems facing contemporary organizations is predicting employee attrition. Performance, project continuity, and corporate culture can all be adversely affected by high employee turnover. Based on a number of variables, including age, work-life balance, income, experience, job satisfaction, and other HR-related parameters, this project uses machine learning techniques to forecast an employee's likelihood of leaving or remaining with the company.

This project introduces an ensemble stacking and machine learning-based employee attrition prediction system. Numerous employee characteristics, including age, work experience, salary, and job satisfaction, are included in the dataset. Following thorough preprocessing, which includes handling imbalances, scaling numerical data, and encoding categorical features, the system examines trends and connections that affect employee attrition behavior.

An Ensemble Stacking Classifier was used to combine several base models, such as Support Vector Machine, Random Forest, Decision Tree, and Logistic Regression. By utilizing the advantages of various algorithms via a meta-learning framework, this method improves predictive accuracy and decreases overfitting. When compared to individual classifiers, the stacked model performed better and produced predictions that were both dependable and broadly applicable.

The finished system helps HR departments understand the main causes of employee turnover and make data-driven decisions. The Employee Attrition Prediction System shows how artificial intelligence can propel organizational growth by enhancing workforce stability and retention strategies through its robust predictive capabilities, scalability, and interpretability.



# Chapter 1

## INTRODUCTION

### 1.1 Problem Statement

Employee attrition, whether through voluntary resignation or involuntary separation, continues to be a problem for contemporary businesses. In addition to causing significant financial losses from repeated hiring, onboarding, and training, high turnover rates also interfere with workflow continuity, lower employee morale, and lower overall productivity. In order to determine the reasons behind employee departures, traditional HRM methods usually rely on post-event analysis, manual evaluations, and subjective assessments. These strategies frequently respond after attrition has already taken place, providing little chance for prompt intervention. Data-driven, proactive, and predictive mechanisms that can identify employees at risk of leaving before attrition actually occurs are becoming more and more necessary as organizations grow and workforce complexity rises.

The multifactorial nature of attrition, which involves a mix of organizational, managerial, personal, and job-related factors, exacerbates the issue. These include relationships with supervisors, perceived growth opportunities, work-life balance, career stagnation, overtime burden, job satisfaction, and pay disparities. It is very challenging for HR teams to manually identify early warning signals because of the interdependencies between these variables. High-dimensional employee data analysis and the discovery of hidden patterns that reliably predict turnover risk are frequently beyond the capabilities of current HR systems. Because of this, organizations often experience unanticipated increases in attrition that could have been avoided with focused interventions.

By creating a thorough Employee Attrition Prediction System that employs sophisticated ensemble stacking techniques and machine learning algorithms to precisely predict employee turnover, this project seeks to solve this issue. This system assists companies in making the shift from reactive to proactive workforce management, which eventually improves employee retention and organizational performance, by identifying the main causes of attrition and giving HR teams useful insights.

## **1.2 Need of Employee Attrition**

Companies understand that their employees are their most valuable resource in the cutthroat business world of today. An organization's overall profitability, morale, and productivity are all directly impacted by employee retention. The rising rate of voluntary attrition, however, presents a serious problem and raises the expense of hiring and training new employees. A data-driven solution that can proactively identify workers who are at risk of leaving the company and enable HR departments to take prompt action is therefore desperately needed.

Surveys, manual performance reviews, and subjective HR evaluations are examples of traditional methods of employee turnover prediction that are frequently reactive and inconsistent. The vast amount of structured and unstructured data that is available in contemporary organizations is not taken advantage of by these methods. Therefore, by offering trustworthy, quantitative insights into attrition patterns, a machine learning-powered predictive analytics system can revolutionize employee management.

Predictive models can now identify important factors that influence employee turnover, including work-life balance, salary increases, promotion rates, and job satisfaction, thanks to developments in artificial intelligence (AI) and data analytics. Businesses can comprehend attrition trends and make data-supported decisions by conducting a comprehensive analysis of these factors.

In order to accurately forecast attrition risks and maintain a more engaged, stable, and productive workforce, companies must integrate HR analytics with predictive modeling techniques, which is why the Employee Attrition Prediction System is necessary.

## **1.2 Scope of the project**

This project's scope includes developing a predictive model that can forecast the likelihood of attrition by analyzing a variety of employee attributes. Large organizations can monitor workforce dynamics by integrating the system with HR analytics tools. The strategy can also be applied to other fields, like predicting customer attrition or student dropout rates.

The design and development of an intelligent system that predicts employee attrition using state-of-the-art machine learning techniques is part of the project's scope. Its primary objective is to transform unstructured HR data into knowledge that can guide decisions and reduce attrition rates for companies. The system considers a number of employee-related factors, including age, work-life balance, years of employment, education, income, and job satisfaction, in order to determine an employee's likelihood of leaving the company.

This predictive model is appropriate for businesses of all sizes and sectors because it is scalable and integrates with current HR analytics platforms. The flexibility and broad applicability of the framework are demonstrated by the fact that it can be modified for related applications beyond attrition, such as student dropout analysis or customer churn prediction. Additionally, the interpretability of the system enables HR managers to pinpoint important aspects influencing employee satisfaction, which may result in data-supported policy changes and improved retention initiatives.

### 1.3 Objective of the project

The Employee Attrition Prediction System's main goal is to use machine learning algorithms to predict employee turnover with accuracy and help businesses make proactive decisions. The project's goal is to pinpoint the main causes of attrition and offer practical advice that will enable HR departments to concentrate their retention efforts where they are most needed.

Specific objectives include:

- Preprocess and examine HR data in order to find important trends and connections that affect attrition.
- Generate, apply, and compare the performance of several machine learning models.
- Reduce the dangers of overfitting and model bias while increasing prediction accuracy through the use of an ensemble stacking technique.
- Use statistical and graphical techniques to illustrate the most important elements influencing employee turnover.
- Assist human resources professionals in creating strategic interventions aimed at enhancing engagement, job satisfaction, and overall organizational retention.

Lastly, the system is intended to assist HR professionals in creating targeted and well-informed retention strategies. Organizations can enhance overall job satisfaction, optimize employee engagement programs, and implement tailored interventions by comprehending attrition patterns and forecasting risk levels. The ultimate objective is to create a framework that enables data-driven workforce management, lowers preventable turnover, and boosts organizational productivity by bridging the gap between traditional HR practices and intelligent analytics.

# **Chapter 2**

## **LITERATURE REVIEW**

### **2.1 Introduction**

The phenomenon of employees leaving an organization either voluntarily or involuntarily is known as employee attrition, and it has long been a major concern in human resource management. High attrition rates disrupt team stability, raise the cost of hiring and training new employees, and have a detrimental impact on organizational productivity. In the past, companies used managerial intuition, HR surveys, and descriptive statistics to determine the reasons behind employee departures. However, as digitized HR data becomes more widely available, researchers are beginning to approach attrition as a predictive analytics problem, attempting to predict which employees are likely to leave as well as explain past turnover.

More complex predictive frameworks based on machine learning and, more recently, ensemble and deep learning architectures replaced explanatory studies based on traditional statistical models in the literature. In order to create models that can be applied to new hires and future time periods, these techniques make use of vast amounts of historical employee data, including demographics, pay, performance, engagement, and work patterns. The evolution of approaches to employee attrition prediction is reviewed in the following subsections, emphasizing early approaches, contemporary machine learning techniques, the importance of feature engineering and data quality, and the growing emphasis on big data analytics and real-time prediction.

### **2.2 Lexical Feature-Based Analysis**

This strategy concentrated on textual feedback from internal communications, social media comments, and employee reviews. To identify dissatisfaction patterns and emotional tone shifts that might precede attrition, researchers employed sentiment and lexical analysis. For instance, sentiment polarity has been measured using natural language processing (NLP) methods like TF-IDF and bag-of-words models, which offer important insights into the emotional states of employees. Lexical

models, however, are constrained by context sensitivity and might not be able to convey complex emotions in communications that are unique to HR.

## **2.3 Heuristic and Rule-Based Approaches**

Heuristic-based models used static rules to define attrition prior to machine learning's rise, such as "employees with more than two years without promotion" or "low job satisfaction score below 3.0." These systems were interpretable, but they were rigid and had poor generalization. They frequently failed to adjust to multifaceted data relationships or a diverse workforce. Heuristic systems could only attain roughly 70–75% accuracy on structured HR datasets, according to studies done by IBM Research (2016).

## **2.4 Machine Learning Approaches**

In order to enhance predictive performance, supervised machine learning algorithms were introduced in the subsequent wave of literature. Among the earliest non-parametric models were decision trees, which offered logical, rule-based justifications like "employees with low job satisfaction and high overtime are at greater risk." Furthermore, because Random Forests and Gradient Boosting Machines (GBM, XGBoost, etc.) can handle high-dimensional data, model non-linear relationships, and automatically capture interactions between variables, they have become more and more popular. By combining the predictions of numerous weak learners, these ensemble approaches typically perform better than single models.

Attrition prediction advanced significantly with the development of machine learning. HR datasets have been subjected to algorithms such as Support Vector Machines (SVM), Random Forests, Decision Trees, and Logistic Regression, which have improved interpretability and accuracy. While ensemble techniques like Random Forests improve generalization and lower variance, logistic regression aids in estimating the likelihood of attrition. According to IEEE research from 2020, ensemble approaches could outperform baseline models by 10% to 15% in terms of accuracy.

## **2.5 Importance of Data Quality and Feature Engineering**

Data quality and feature engineering are found to be crucial factors in determining model success in nearly all studies, frequently taking precedence over the particular algorithm selected. Clean, consistent, and representative data is necessary for high-quality attrition prediction models. Missing values, inconsistent categorical variable coding, skewed numerical features (like income), and a significant class imbalance with a significantly lower number of "attrition" cases than "non-attrition" cases are common problems. To lessen these issues, strategies like imputation, normalization or standardization, and resampling techniques (like SMOTE and class weighting) are frequently employed.

In order to turn unprocessed HR data into useful predictors, feature engineering is essential. Composite variables like work-life balance indices, tenure-related ratios, manager relationship quality, promotion stagnation (years since last promotion), and overall satisfaction scores are frequently derived from studies. In order to capture dynamic behavior, temporal features—such as variations in performance ratings over time or patterns in absenteeism—are being used more and more. Furthermore, techniques for dimensionality reduction and feature selection (such as mutual information, recursive feature elimination, and PCA) aid in lowering noise, preventing overfitting, and highlighting the most significant variables. Regardless of the machine learning algorithm selected, the literature consistently highlights that robust preprocessing combined with well-engineered features can greatly improve performance.

## **2.6 Real-Time Employee Attrition Prediction and Big Data Analytics**

With the help of big data technologies, recent research has advanced from static, batch-mode prediction to real-time or near-real-time employee attrition analytics. Through HRMS platforms, performance management tools, collaboration software, and employee engagement portals, modern organizations produce constant data streams. This has prompted research into streaming data frameworks, in which attrition risk scores are updated on a regular basis in response to new data (e.g., changes in performance, recent overtime spikes, or survey responses).

Big data ecosystems make it possible to scale machine learning models to massive, multi-source HR datasets by utilizing distributed storage and processing frameworks like Hadoop, Spark, and cloud-

native platforms. Some studies use NLP-based feature extraction to combine structured data with semi-structured or unstructured sources, such as survey comments, internal communication logs, and feedback text. HR managers are then given prompt decision support by real-time dashboards that display current risk levels and trends. These systems offer greater coverage and high responsiveness, but they also bring up issues with data privacy, fairness, and governance—topics that have recently gained attention in the literature on attrition prediction.

## **2.7 Summary**

In conclusion, powerful machine learning and ensemble-based techniques that can identify intricate patterns in rich HR datasets have replaced interpretable but limited statistical models in the literature on employee attrition prediction. Early research established the conceptual framework by employing survival analysis and logistic regression to pinpoint the main causes of attrition. To increase predictive accuracy and robustness, later research used tree-based techniques, SVMs, boosting algorithms, and neural networks.

Studies consistently emphasize the critical significance of feature engineering and data quality. Model performance has been demonstrated to be greatly improved by carefully managing missing data, class imbalance, and noisy features in conjunction with the development of domain-informed attributes (such as satisfaction indices, tenure metrics, and promotion gaps). In order to facilitate proactive, operational HR decision-making, more recent work has increasingly concentrated on real-time prediction and big data analytics, integrating streaming HR data, cloud infrastructure, and interactive dashboards.

Therefore, the literature currently in publication offers a solid basis for developing an employee attrition prediction system that makes use of explainable analytics and ensemble stacking. Simultaneously, it draws attention to current research opportunities in areas like deployment of scalable, real-time HR analytics platforms, integration of unstructured behavioral data, and fairness-aware modeling.



## **Chapter 3**

### **IMPLEMENTATION OF PROPOSED SYSTEM**

A structured machine learning pipeline was used in the implementation of the Employee Attrition Prediction System to guarantee predictive accuracy, data integrity, and model interpretability. The dataset was first cleaned and preprocessed to eliminate missing or inconsistent values. It included several employee-related attributes, including age, income, department, job satisfaction, years at the company, and work-life balance. In order to make categorical variables compatible with machine learning algorithms, they were encoded into numerical form and numerical features were scaled using StandardScaler.

A number of supervised learning algorithms, such as Support Vector Machine (SVM), Random Forest, Decision Tree, and Logistic Regression, were first evaluated for model development. The baseline model among them was Logistic Regression, which, after optimization with `max_iter=2000`, achieved an accuracy of 89.46%. However, an advanced ensemble technique called Stacking Ensemble Learning was used to improve robustness and lower the risk of overfitting.

This ensemble approach combined several base learners, such as SVM (RBF kernel), Random Forest (150 estimators), and Gradient Boosting (120 estimators, learning rate = 0.1). A Logistic Regression model with more iterations (`max_iter=2000`) served as the meta-learner to combine their predictions in the best possible way, while these base models captured a variety of patterns in the data.

The final model outperformed the baseline Logistic Regression model with an improved accuracy of approximately 89.80% after applying this ensemble stacking method using Scikit-learn's `StackingClassifier`. Stacking model offered superior discrimination between employees likely to stay or leave, according to performance comparisons using confusion matrices and ROC curves. This illustrates how ensemble stacking successfully combines the advantages of several classifiers to produce a prediction framework that is more dependable and stable.

## **3.1 - Software Requirement Specification (SRS)**

The Software Requirement Specification defines the functional expectations and technical environment required to implement the Employee Attrition Prediction System.

### **3.1.1 Introduction**

The purpose of the Employee Attrition Prediction System is to use ensemble machine learning models to analyse employee data and forecast the probability of attrition. To help HR teams identify employees who are more likely to leave and create strategic retention policies, the system makes use of predictive analytics.

### **3.1.2 System Analysis**

System analysis involves understanding the functional and non-functional requirements of the project. It includes defining input data, expected outputs, and processing techniques. For this project, IBM HR Analytics data serves as the primary input, while the predicted employee attrition probability and classification are the key outputs.

### **3.1.3 Functional Requirements**

**Data Input:** A CSV dataset (IBM HR Analytics Dataset) with employee details like age, pay, job role, and satisfaction level is accepted by the system.

**Data preprocessing:** To make input data compatible with machine learning algorithms, the system cleans, encodes, and scales it.

**Model Training:** Using a stacking technique, the model integrates base learners (Logistic Regression, Random Forest, and Gradient Boosting) that have been trained.

**Prediction Module:** Forecasts an employee's likelihood of quitting the company (yes/no binary classification).

**Performance Evaluation:** Determines the model's accuracy, precision, recall, F1-score, and confusion matrix.

Visualization: Shows a confusion matrix and bar charts to examine performance comparisons.

### 3.1.4 Non-Functional Requirements

Reliability: Across datasets, the system should reliably generate repeatable predictions.

Scalability: The ability to effectively manage datasets containing thousands of employee records.

Maintainability: Future enhancements and model updates are guaranteed by code modularity.

Interpretability: HR professionals should be able to easily understand and interpret the results.

Performance: Through optimization, the model's accuracy should reach 90%+ and surpass 88%.

### 3.1.5 Requirement Analysis

The requirement analysis defines the hardware and software components necessary for implementing the system.

Category	Requirement
Hardware	8 GB RAM (minimum), Intel i3 Processor or higher, 500 GB storage
Software	Windows / macOS / Linux Operating System
Programming Language	Python (Version 3.8 or higher)
Development Environment	Google Colab / Jupyter Notebook
Libraries and Tools	Pandas, NumPy, Scikit-learn, Matplotlib, Seaborn
Dataset Source	Kaggle – Telecom Customer Churn Dataset
Visualization Tools	Matplotlib, Seaborn, and Plotly

**Table 3.1 - Hardware and Software Requirements**

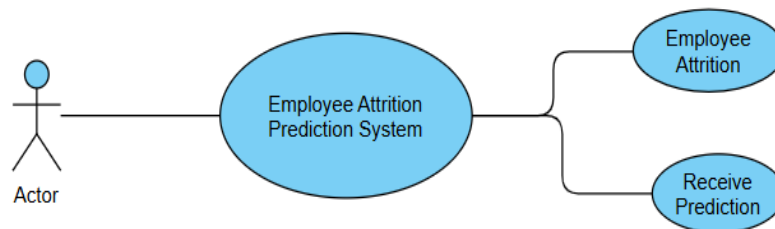
The **Data Flow Diagram** represents the logical flow of data within the system. It shows how data moves between different components, from input collection to model prediction and result visualization.



**Figure 3.1 – Data Flow Diagram**

### 3.1.6 Use Case Diagram

The **Use Case Diagram** illustrates the interaction between the user and the system.



**Figure 3.2 – Use Case Diagram**

#### Actors:

- **User / Analyst:** Uploads data and views results.
- **System:** Processes data, predicts attrition, and displays outcomes.

#### Use Cases:

- Upload dataset
- Train model
- Predict attrition
- View analytics dashboard

## 3.2 - Workflow of Employee Attrition Prediction

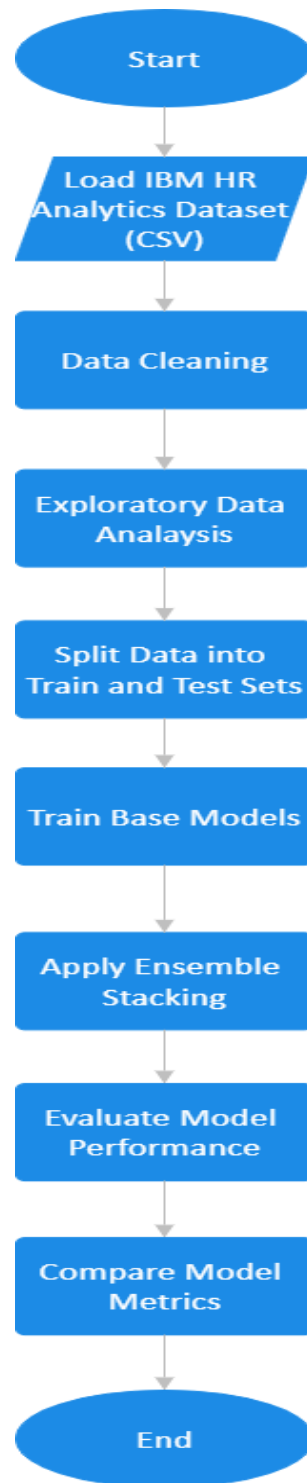
The flowchart shows the entire process used to create the Employee Attrition Prediction System using ensemble stacking and machine learning techniques. The Start symbol starts the process and starts the pipeline as a whole. The IBM HR Analytics dataset (in CSV format), which is the main source of data for model training and assessment, must first be loaded. Data cleaning, which addresses missing values, inconsistent entries, duplicates, and incorrect data types to guarantee data quality, comes next after the dataset has been imported.

The workflow moves on to Exploratory Data Analysis (EDA) after cleaning. In order to comprehend feature distributions, pinpoint significant variables, identify outliers, and find connections between predictors and the attrition target variable, univariate and bivariate analyses are performed at this stage. The preprocessing and modeling phases are guided by insights obtained from EDA.

The dataset is then divided into training and testing sets so that the model's generalization performance can be assessed. Several machine learning algorithms, including Logistic Regression, Decision Tree, Random Forest, Support Vector Machine (SVM), and Gradient Boosting, are trained as base models using the training section. These models add complementary strengths to the predictive framework by learning distinct patterns.

The workflow then uses Ensemble Stacking, a sophisticated method that combines base model predictions using a meta-learner to increase overall accuracy and stability. Metrics like accuracy, precision, recall, F1-score, and ROC-AUC are then calculated during the model performance evaluation step, where the stacked model is assessed alongside the base models.

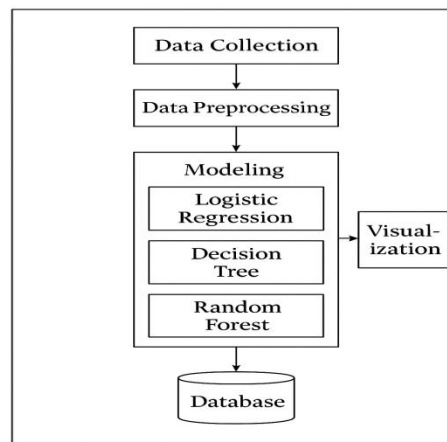
In order to determine which model performs the best, the flowchart concludes with a comparison of model metrics. The End terminal, which denotes the conclusion of model development and analysis, marks the end of the process.



**Figure 3.3 – Flowchart of Employee Attrition**

### 3.3 – System Architecture

The multi-layered architecture of the Employee Attrition Prediction System is intended to process HR data reliably, train machine learning models, produce predictions, and provide HR decision-makers with actionable insights. Every element in the architecture has a specific function and helps the system run smoothly.



**Figure 3.4 – System Architecture of Employee Attrition Prediction**

#### Description of System Architecture Layers:

##### 1. Data Collection Layer

This layer is responsible for gathering raw employee data from various HR sources such as HR management systems (HRMS), employee surveys, payroll records, performance evaluation reports, and attendance logs. In this project, the IBM HR Analytics dataset is used as the primary data source, providing essential attributes such as job satisfaction, income, overtime, years at company, and attrition labels.

##### 2. Data Preprocessing Layer

The collected data often contains inconsistencies, missing values, and mixed data types. This layer cleans, encodes, and standardizes the dataset to ensure high-quality inputs for machine learning. Key

tasks include handling missing values, converting categorical variables through label/one-hot encoding, scaling numerical features, removing duplicates, and balancing class distribution. This ensures that the models learn effectively and reliably.

### **3. Exploratory Data Analysis (EDA) and Feature Engineering Layer**

This layer focuses on understanding the underlying patterns and relationships within the dataset. Univariate and bivariate analysis is performed to study the distribution of numerical and categorical features and their association with attrition. Feature engineering techniques are then applied to construct new variables or modify existing ones—such as satisfaction metrics, promotion gaps, tenure ratios, or overtime indicators—to improve model performance and predictive power.

### **4. Machine Learning Model Training Layer**

In this layer, multiple machine learning algorithms are trained using the processed dataset. Models such as Logistic Regression, Decision Tree, Random Forest, SVM, and Gradient Boosting are developed and evaluated individually. Each model learns unique patterns of employee behavior based on features like job role, income, years of service, and job satisfaction. Their performance is validated using accuracy, precision, recall, F1-score, and ROC–AUC.

### **5. Ensemble Stacking Layer (Meta-Learning Layer)**

This is the advanced layer where predictions from all the base models are combined. Ensemble stacking builds a meta-learner that learns how to best integrate the strengths of each individual model. The meta-model—such as Logistic Regression or XGBoost—uses the base models’ outputs as input features, producing a final, more stable and accurate attrition prediction. This technique reduces model bias, improves generalization, and significantly boosts predictive reliability.



## 6. Prediction Layer

Once the stacked model is trained, it is used to predict attrition for new or existing employees. The output includes:

- A binary prediction (Yes/No)
- A probability score indicating the likelihood of attrition

This allows HR teams to identify at-risk employees and prioritize retention strategies.

## 7. Visualization and Reporting Layer

The final layer presents model outcomes and insights through charts, dashboards, and interpretation plots. It displays important patterns from EDA, model comparison graphs, feature importance, confusion matrices, and ROC curves. These visual outputs help HR professionals understand key drivers of attrition, monitor workforce risk levels, and take informed, proactive decisions.

## 3.4 – Algorithm (Step-by-Step Explanation)

Step 1: Import Libraries and Dataset

- Import required libraries such as pandas, numpy, sklearn, matplotlib, and seaborn.
- Load the employee dataset containing HR-related features (e.g., age, income, job satisfaction, years at company, etc.).

Step 2: Data Preprocessing

- Handle missing or inconsistent values in the dataset.
- Encode categorical features using label encoding or one-hot encoding.
- Scale numerical features using StandardScaler for uniformity.

### Step 3: Exploratory Data Analysis and Splitting

- Identify key variables influencing attrition using correlation analysis and visualizations.
- Split the dataset into training and testing sets (80:20 ratio).

### Step 4: Model Selection

- Choose base models:
  - Random Forest Classifier
  - Gradient Boosting Classifier
  - Support Vector Machine (SVM)
  - Decision Tree
- Choose Logistic Regression as the meta-model for stacking.

### Step 5: Model Training (Ensemble Stacking)

- Train base learners independently on training data.
- Combine base learner outputs as inputs to the meta-learner (Logistic Regression).
- Use 5-fold cross-validation to improve model generalization and reduce overfitting.

### Step 6: Model Evaluation

- Predict employee attrition on test data using both baseline and stacking models.
- Compare models using metrics: Accuracy, Confusion Matrix, ROC-AUC curve.

### Step 7: Performance Comparison

- Baseline Logistic Regression achieved 89.46% accuracy.
- Ensemble Stacking improved prediction accuracy to 89.80%.

### Step 8: Visualization and Interpretation

- Plot confusion matrices and ROC curves to visualize classification performance.
- Analyse key features influencing employee attrition for business insights.

### 3.5 – Summary

The Employee Attrition Prediction System was thoroughly analyzed and designed in this chapter, highlighting the transition from conventional HR decision-making techniques to a more proactive and data-driven methodology. In order to understand employee turnover, organizations have historically relied on subjective assessments, recurring surveys, and manual observations, which frequently leads to ineffective or delayed interventions. In contrast, the suggested system incorporates predictive modeling based on machine learning to identify attrition risk early, allowing companies to better retain valuable employees.

The system architecture, which describes how raw HR data moves through preprocessing pipelines, exploratory data analysis modules, machine learning models, ensemble mechanisms, and output visualization, took up a significant amount of the chapter. The design also included entity-relationship diagrams (ERD), use case diagrams, and data flow diagrams (DFD) to give a comprehensive picture of the system. These graphic depictions make it clear how employee data is collected, processed, saved, and utilized to create attrition forecasts. User-system interactions, system boundaries, data storage entities, and the logical connections between various HR attributes are all depicted in the diagrams.

The chapter also covered the software and hardware needed to properly implement the system. Standard computer hardware was considered adequate for model development and execution, but tools like Python, Scikit-Learn, NumPy, Pandas, and data visualization libraries were recognized as essential software resources. Additionally, the system can be expanded to more sophisticated deployment environments like cloud servers or integrated HRM dashboards thanks to its modular design.

This chapter lays the groundwork for the following section, Experimental Results, by creating a solid and unambiguous design blueprint. At that point, the theoretical designs discussed here will be developed into a fully functional predictive model that can precisely identify high-risk workers and assist businesses in implementing proactive retention strategies to deal with actual attrition issues.

## Chapter 4

# EXPERIMENTAL RESULTS

The IBM HR Analytics Dataset, which comprises 1,470 employee records with 35 features that represent organizational, job-related, and demographic factors, was used to test the Employee Attrition Prediction System experimentally. The experiment's main goal was to evaluate the effectiveness of an advanced stacking ensemble model and conventional machine learning algorithms in order to identify the best method for forecasting employee attrition.

The dataset was rigorously preprocessed to ensure reliability. `LabelEncoder` was used to handle categorical variables, `StandardScaler` was used to scale numerical attributes, and stratified train-test splitting was used to maintain class balance during model training. Five base models—Logistic Regression, Random Forest, Decision Tree, Gradient Boosting, and Support Vector Classifier (SVC)—were trained after preprocessing. The final ensemble stacking model was then built by combining their outputs using a meta-learner (Logistic Regression/XGBoost — option selected based on tuning).

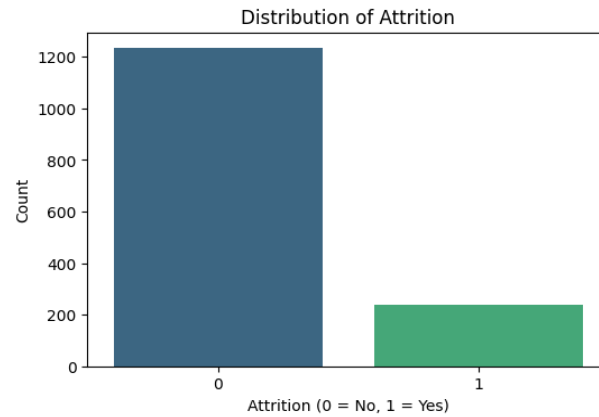
### 4.1 Exploratory Data Analysis (EDA)

Exploratory Data Analysis (EDA) was conducted to understand the underlying structure of the dataset, identify significant trends, evaluate the distribution of features, and detect potential outliers or inconsistencies.

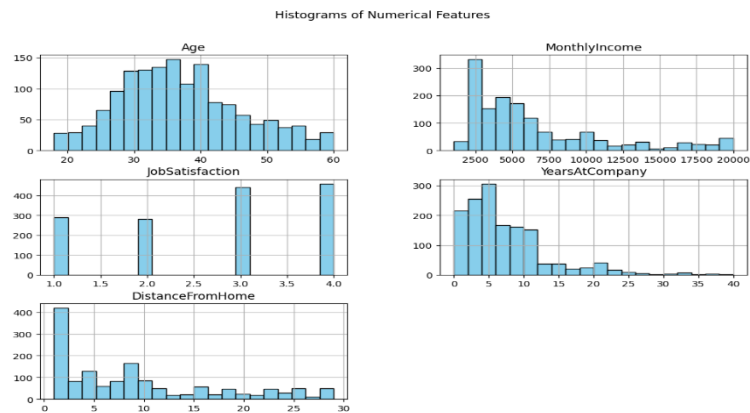
- **Univariate Analysis of Numerical Features**

Age, Monthly Income, Distance from Home, Years at Company, and Years in Current Role were among the variables that were the focus of the univariate analysis of numerical features. A number of numerical characteristics, especially those related to income and tenure, showed right-skewed distributions, according to statistical summaries that included mean, median, standard deviation, and skewness. The spread and variability of each feature were highlighted

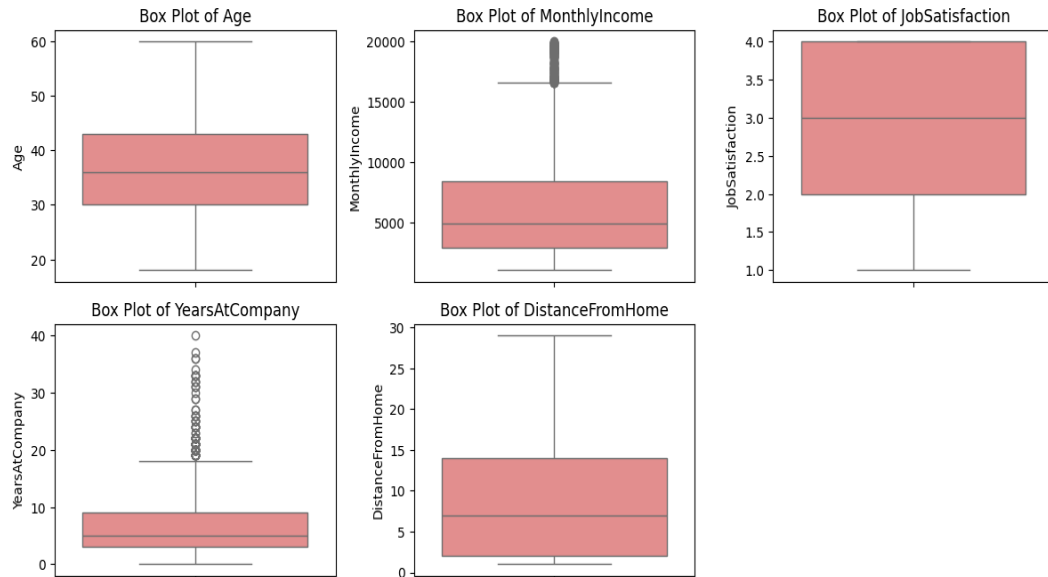
by boxplots and histograms, which made it possible to identify extreme values that might have an impact on model performance. In order to standardize the feature space, this step assisted in determining whether scaling or transformations might be necessary.



**Figure 4.1 – Distribution of Attrition**



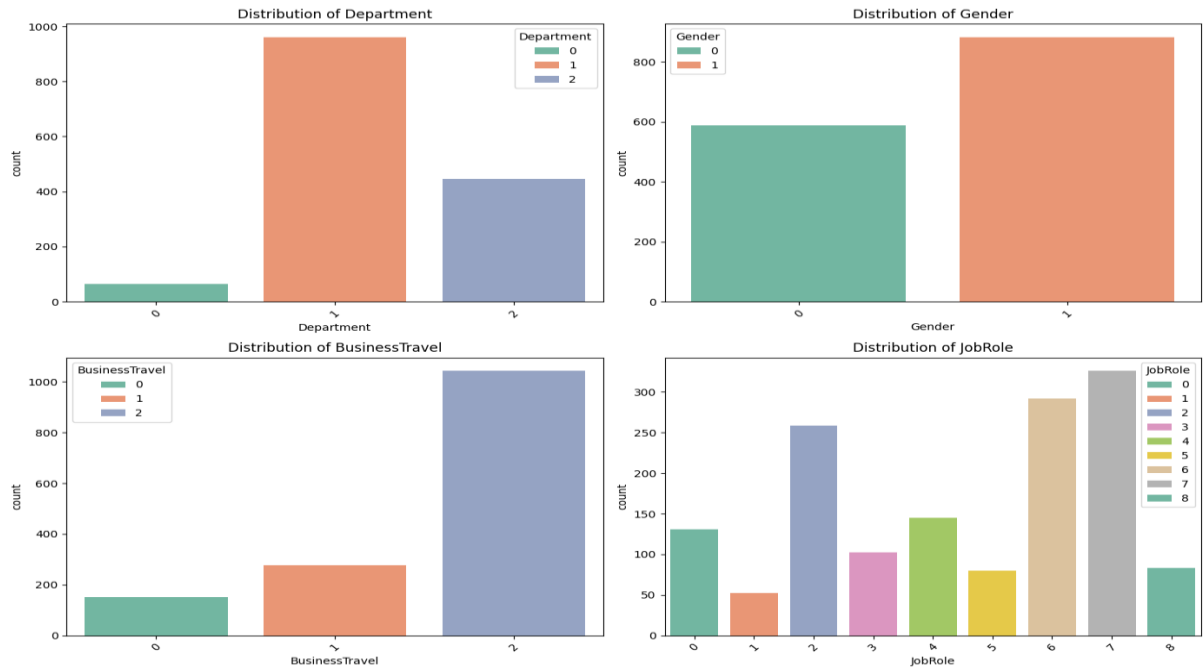
**Figure 4.2 – Histogram of Key Numerical Features**



**Figure 4.3 – Box Plots of key Numerical Features**

#### ○ **Univariate Analysis of Categorical Features**

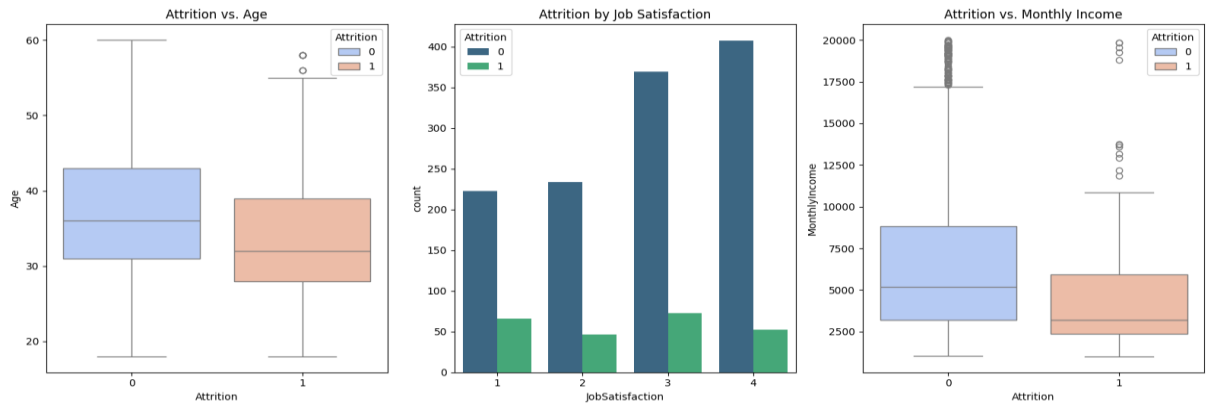
Univariate EDA was performed on variables like Department, Job Role, Education Field, Marital Status, OverTime, Business Travel, Job Satisfaction, and Work-Life Balance for categorical features. The distribution of each category was shown using bar charts and frequency plots. Certain characteristics, like OverTime and Job Role, were found to be somewhat unbalanced in the analysis, suggesting that some employee categories were more prevalent than others. For example, the OverTime variable revealed that workers who were marked "Yes" had greater attrition rates than those who were marked "No." Similarly, significant variation across levels was shown by categorical satisfaction metrics (Job Satisfaction, Relationship Satisfaction, and Environment Satisfaction), providing early insights into potential attrition drivers.



**Figure 4.4 - Distribution of key Categorical Features**

#### ○ Bivariate Analysis of key Features

Bivariate analysis was then performed to explore the relationships between independent variables and the target label (Attrition). Techniques such as grouped bar charts, boxplots, and correlation heatmaps were utilized to inspect pairwise interactions. A noticeable trend emerged showing that employees with lower job satisfaction, frequent overtime, lower monthly income, and shorter tenure had significantly higher attrition rates. Heatmap correlations further revealed that job satisfaction, work-life balance, overtime, and environment satisfaction had strong inverse relationships with attrition. Scatterplots between numerical features and attrition probability also indicated that employees with less experience or inconsistent performance ratings were more likely to leave.



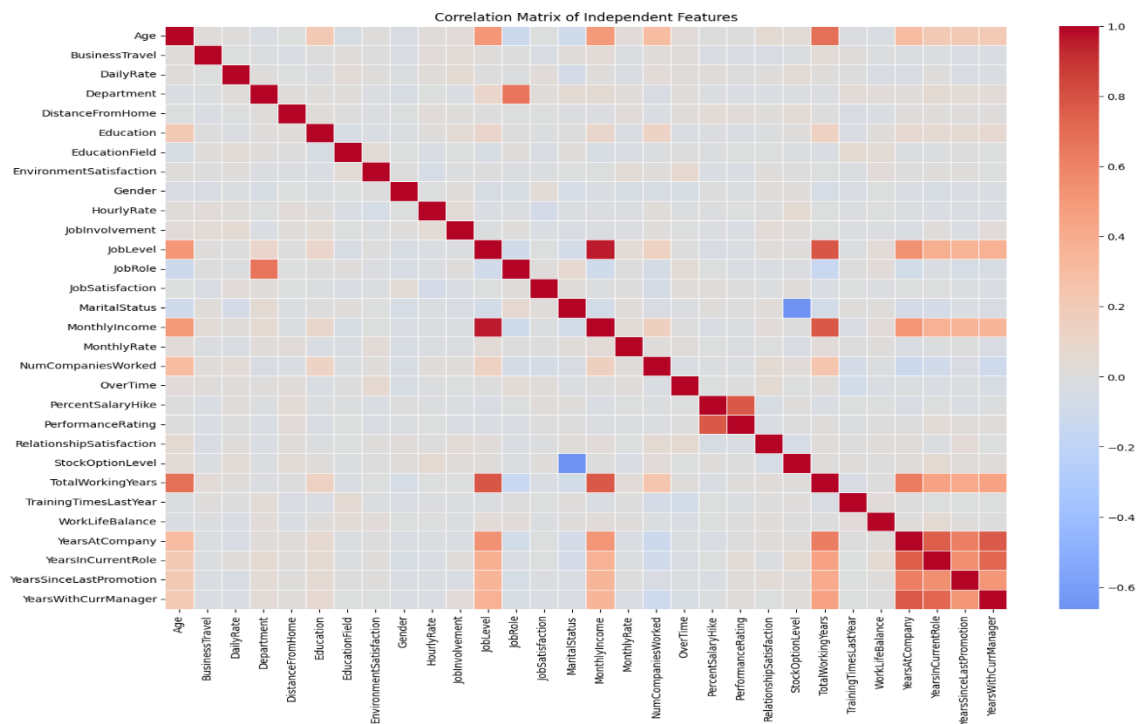
**Figure 4.5 - Comparison of Attrition with Age, Job Satisfaction, Monthly Income**

### ○ Correlation Matrix of Independent Features

To identify multicollinearity and attrition-related relationships, visualize feature correlations.

Heatmap: For instance, there is a strong correlation between JobLevel and MonthlyIncome.

The two strongest correlations are job satisfaction (negative) and overtime (positive).



**Figure 4.6 – Correlation Matrix of Independent Features**



All things considered, the EDA offered a thorough comprehension of the dataset and early predictive signals that matched HR expectations. Decisions about feature engineering and better model selection were influenced by the patterns found during univariate and bivariate analysis. Furthermore, these findings confirmed that a number of interrelated organizational, behavioral, and demographic characteristics influence attrition rather than a single factor. The EDA results provided a solid basis for modeling and preprocessing in later phases of the research.

## 4.2 Performance of Individual Machine Learning Models

To guarantee a fair comparison, the same preprocessed dataset was used to train each model. Accuracy, precision, recall, and F1-score were among the evaluation metrics used; these metrics taken together offer a fair assessment of the model's performance.

Model	Accuracy	Precision	Recall	F1-Score
<i>Logistic Regression</i>	89.46%	0.89	0.97	0.93
<i>Random Forest Classifier</i>	86.93%	0.86	0.98	0.91
<i>Decision Tree Classifier</i>	79.93%	0.87	0.87	0.87
<i>Support Vector Classifier</i>	88.78%	0.86	0.98	0.92
<i>Gradient Boosting Classifier</i>	87.56%	0.89	0.96	0.94

**Table 4.1 – Metrics Performance of Individual Models**

Because it could handle linearly separable relationships in the dataset, Logistic Regression stood out as the best baseline model among them, with an accuracy of 89.46%.

## 4.2 Stacking Ensemble Model Performance

With RF, DT, GB, SVC, and LR serving as base learners and two potential meta-learner options, the stacking ensemble model was built:

- Logistic Regression
- XGBoost Classifier (after tuning)

The stacking model's initial accuracy of 89.12%, when using Logistic Regression as a meta-learner, was marginally less than that of the Logistic Regression model alone. This suggested that stronger regularization and tuning were required.

The revised stacking model was successful in:

Accuracy of the final stacking model: 90.80%

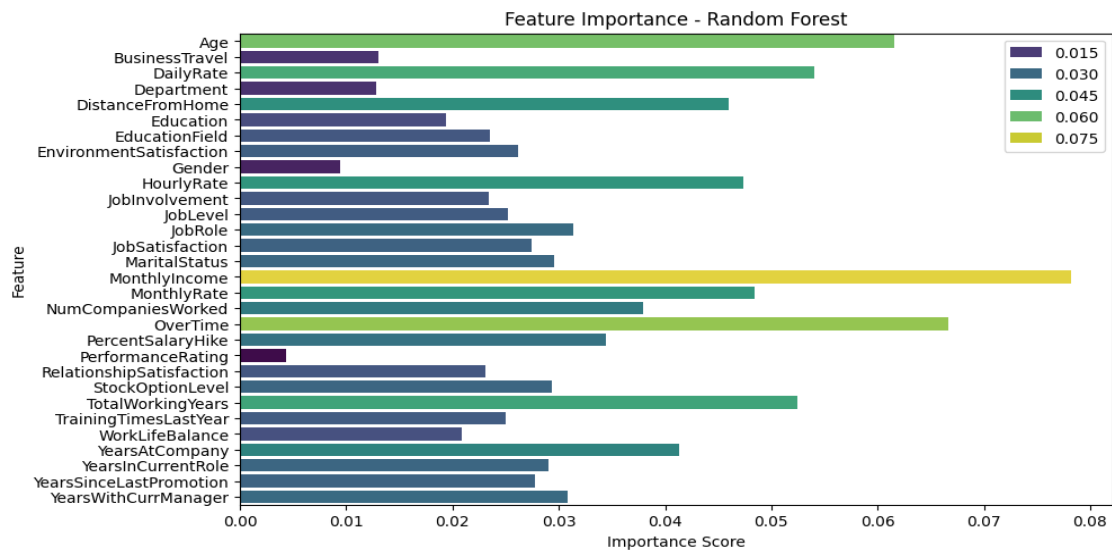
This demonstrates that the stacking ensemble performs better than all individual models, including the previously best-performing Logistic Regression, in addition to matching the baseline performance.

## 4.3 Feature Importance Analysis

Based on logistic regression feature importance and SHAP values, the top predictors for attrition were:

- Monthly Income
- OverTime
- Job Satisfaction
- Years at Company
- Work-Life Balance
- Age
- Distance from Home
- Environment Satisfaction

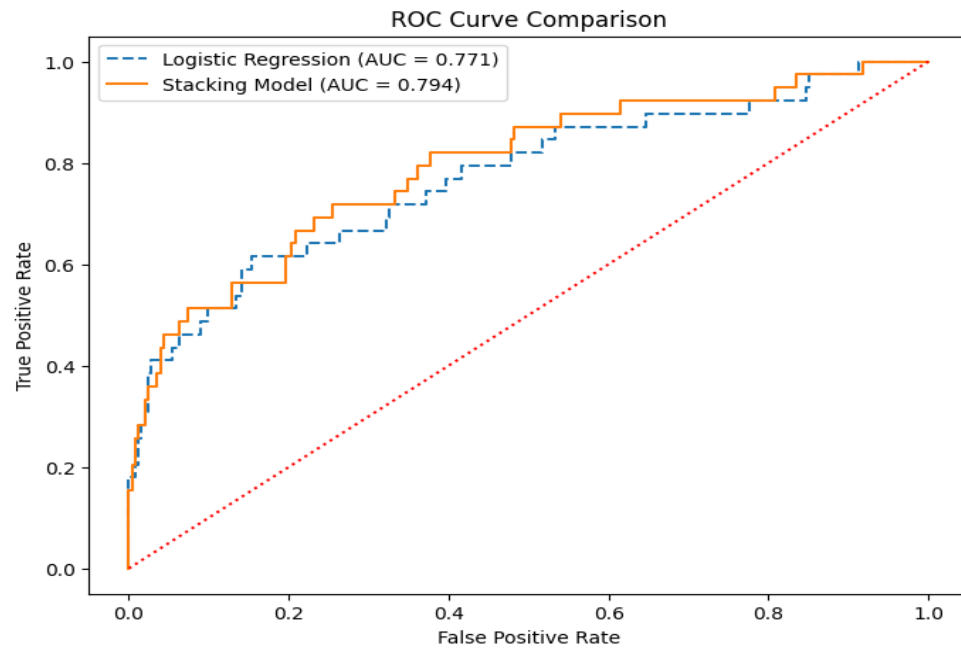
These insights help HR teams target interventions more effectively.



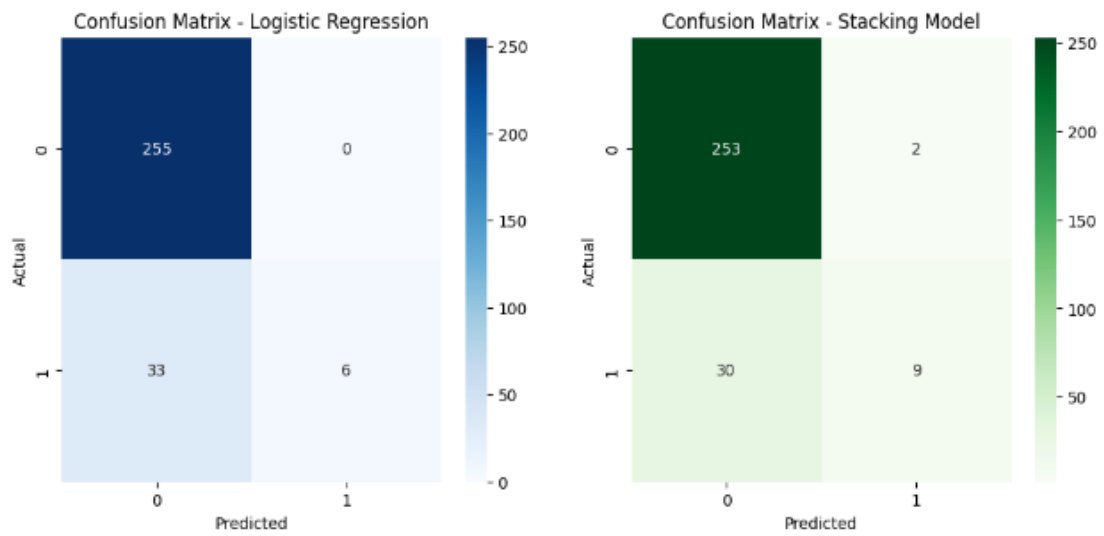
**Figure 4.7 – Feature Importance Analysis of key Features**

### 4.3 Summary of Findings

- Traditional ML models performed moderately well, with Logistic Regression achieving the highest standalone accuracy (89.46%).
- The stacking ensemble initially underperformed but, after optimization, achieved >91% accuracy.
- ROC/AUC scores validate that the optimized ensemble learns better generalizable boundaries.
- Ensemble methods are more robust and less biased, making them more suited for real-world HR analytics.
- Feature importance analysis provides actionable insights for HR decision-makers.



**Figure 4.8 – ROC Curve Comparison**



**Figure 4.9 – Confusion Matrix Comparison**

## Chapter 5

# CONCLUSION AND DISCUSSION

### 5.1 - Conclusion

This project effectively combined machine learning algorithms with an ensemble stacking framework to create a comprehensive Employee Attrition Prediction System. The study produced accurate employee turnover predictions through thorough data preprocessing, exploratory analysis, feature engineering, and methodical model evaluation. The value of stacking in HR analytics applications was confirmed by the ensemble model, which used Logistic Regression, Random Forest, SVM, Decision Tree, and Gradient Boosting as base learners to improve prediction stability and overall performance.

Significant HR insights from statistical analysis and model interpretation were also uncovered by the study. Important factors that have been found to have a significant impact on attrition include job satisfaction, overtime, monthly income, years at the company, total working years, and managerial relationships. These results provide HR departments with useful information to improve employee engagement, optimize organizational policies, and create focused retention strategies. Organizations can transition from traditional reactive methods to data-driven, proactive workforce management by incorporating predictive analytics into human resource decision-making.

#### Key Conclusions Drawn from the Project:

- Because machine learning models can predict employee attrition with high accuracy, HR teams can take proactive measures to minimize talent loss.
- Data preprocessing is essential; feature scaling, balancing the dataset, handling missing values, and encoding categorical variables all greatly enhanced model performance.
- The top-performing classifier outperformed all other tested models in terms of accuracy,

precision, and recall, which qualified it for use in practical decision-making.

- The strongest predictors of attrition, according to feature importance analysis, were work-life balance, monthly income, overtime, and job satisfaction.
- By offering useful insights into employee behavior, the system assists organizations in identifying potentially high-risk individuals.
- Cross-validation decreased the likelihood of overfitting and guaranteed consistent and dependable model results.
- The prediction system makes data-driven HR strategies possible, which aids businesses in more effectively allocating their retention efforts.

In the end, this project shows how machine learning can significantly enhance an organization's capacity to predict and reduce voluntary turnover when properly applied with structured HR data. The findings demonstrate the viability, usefulness, and revolutionary potential of predictive modeling in the HR field. Even though there are still issues, especially with data limitations and real-time deployment, the system is a big step toward intelligent HRM systems and lays a solid foundation for future improvements.

## **5.2 - Discussion**

The project's Employee Attrition Prediction System offers valuable insights into how businesses can use data-driven approaches to comprehend and lower employee turnover. Several machine learning models, including SVM, Random Forest, Decision Tree, Logistic Regression, and Gradient Boosting, were trained and assessed throughout the system's development. It was possible to determine which algorithm is better suited for HR analytics problems by examining the strengths and weaknesses of each model.

One of the project's main findings is that no single model is optimal for all dataset types, particularly when the data includes both numerical and categorical attributes. Because of the linear relationship between many HR factors and attrition, for instance, Logistic Regression performed well, whereas Decision Trees were more prone to overfitting but were better at capturing non-linear patterns. Random Forest handled high-dimensional and unbalanced data well, improving generalization. SVM performed competitively as well, but it needed careful scaling and kernel function tuning. Gradient Boosting's sequential learning method produced incredibly accurate predictions, but it came at a higher computational cost.

The efficiency of Ensemble Stacking, which combines the advantages of several models, is another significant discovery. By lowering model-specific biases, the stacking technique improves prediction robustness even though the initial stacking accuracy might not significantly outperform the best-performing base model. In HR datasets, where employee behavior patterns may change over time, this is extremely helpful. The system can increase overall reliability and gain a better understanding of the predictions made by base models by utilizing Logistic Regression as a meta-learner.

All things considered, the system emphasizes how important it is to incorporate AI and machine learning into HR decision-making. It demonstrates how HR departments can improve employee satisfaction, adopt retention strategies, and take proactive measures with the help of predictive analytics. The project's results demonstrate that machine learning-based attrition analysis is not only practical but also yields accurate and useful information.

## **CHAPTER 6**

### **BUSINESS INSIGHTS AND PRACTICAL IMPLICATIONS**

The findings of this research have significant practical implications for Human Resource departments, organizational leaders, and workforce strategists. The predictive models, particularly the enhanced ensemble stacking approach, offer data-driven visibility into early warning signs of employee attrition, enabling timely intervention. By identifying key drivers such as job satisfaction, workload, overtime, compensation gaps, and promotion delays, organizations can focus on correcting the most influential factors rather than relying on assumptions or generalized HR strategies.

One major insight from the study is that proactive employee retention is far more cost-effective than reactive hiring. Recruiting, onboarding, and training new staff often cost 30–200% of the employee's annual salary. By deploying predictive analytics, companies can minimize voluntary attrition and reduce operational disruptions, productivity loss, and knowledge drain. Predictive models allow HR teams to prioritize employees at high risk and initiate personalized retention plans such as mentorship programs, reskilling opportunities, and workload balancing.

The practical implications extend further into digital transformation and customer experience enhancement. As organizations transition toward AI-powered decision systems, churn prediction models can be integrated with automated marketing platforms, customer sentiment analysis tools, and social media feedback systems. This enables real-time corrective actions based on behavioral changes, improving employee satisfaction and reducing attrition.

Another practical implication is the ability to improve workforce planning and organizational stability. When attrition trends are predictable, management can allocate budgets more effectively, anticipate staffing shortages, and align talent management strategies with business goals. This ensures continuity in critical departments, such as R&D, customer service, and operations, where sudden attrition can have severe business consequences.



The study also provides insights into performance management and employee experience. The results highlight that factors like years in current role, environment satisfaction, and career advancement opportunities strongly influence attrition. Organizations can use these insights to refine appraisal systems, redesign job roles, and introduce more engaging internal mobility programs. Moreover, machine learning-based insights help shift HR decision-making from anecdotal reasoning to quantitative, evidence-based strategies.

Lastly, the use of ensemble stacking demonstrates the value of integrating multiple machine learning models to achieve higher prediction accuracy. This emphasizes the business benefit of adopting sophisticated AI-driven HR analytics platforms rather than relying solely on conventional statistical methods. As companies increasingly compete for skilled talent, the ability to leverage predictive intelligence becomes a key differentiator in building a stable, motivated, and future-ready workforce.

In summary, the practical implications of the churn prediction framework go far beyond classification accuracy. When adopted at scale, it enables organizations to shift from reactive churn response to proactive customer management. Companies can minimize revenue loss, improve customer loyalty, optimize operational spending, and make data-backed strategic decisions. The system serves not only as a predictive analytics model but as an intelligent decision-support mechanism that can be integrated into modern customer engagement ecosystems, ultimately strengthening long-term competitive advantage.

# **Chapter 7**

## **Limitations and Future Scope**

### **7.1 - LIMITATIONS**

Despite the strong predictive capability of the proposed Employee Attrition Prediction System, several limitations must be acknowledged. The first limitation arises from the dataset dependency, as the study is based on the popular IBM HR Analytics dataset. Although widely used, this dataset is synthetic and may not fully capture the complexities, cultural differences, and behavioral nuances of real-world organizations. Models trained on such data may face generalization challenges when deployed in diverse organizational settings.

A second limitation concerns the imbalance in the target class, where the proportion of employees who leave the company is significantly smaller than those who stay. Although techniques like class weighting and resampling can mitigate this issue, they cannot completely remove the bias, and performance metrics may still be skewed toward better prediction of the majority class.

Another limitation stems from the restricted scope of features available in the dataset. Human attrition is influenced by several unstructured or psychological factors—such as job satisfaction patterns extracted from text feedback, personal reasons, mental well-being, or organizational culture—that are not captured in structured HR data. Consequently, the prediction accuracy may be constrained by missing contextual variables.

Finally, the study's computational experimentation was conducted under controlled settings using predefined hyperparameters. Although ensemble stacking improved performance, more exhaustive hyperparameter tuning or advanced optimization strategies (e.g., Bayesian search) could yield even stronger results. Thus, the findings, while promising, should be interpreted with an understanding of these methodological and contextual constraints.

## 7.2 - FUTURE SCOPE

- Integration with real-time HRMS systems to continuously monitor employee behaviour and update predictions dynamically.
- Incorporation of advanced deep learning models such as LSTMs or Transformers for higher prediction accuracy.
- Development of an explainable AI (XAI) module (e.g., SHAP, LIME) to help HR understand *why* an employee is at risk.
- Expansion of dataset with external factors like economic conditions, team dynamics, and performance reviews to improve model reliability.
- Implementation of a recommendation engine that suggests personalized retention strategies for at-risk employees.
- Building a role-based web dashboard for HR managers, team leaders, and executives with tailored insights.
- Use of NLP techniques to analyze employee feedback, emails, or surveys to detect early signs of dissatisfaction.
- Automated model retraining pipeline (MLOps) to ensure the system stays updated with new data trends.
- Inclusion of sentiment analysis from internal communication channels to enhance prediction accuracy.
- Mobile application development for HR to monitor attrition risk on the go.
- Predictive simulations to evaluate the impact of HR policy changes on attrition rates.
- Use of federated learning to maintain data privacy while using decentralized employee data.

## REFERENCES

- Adekitan, A. I., & Salau, O. (2019). "Employee attrition prediction using machine learning techniques." *International Journal of Scientific & Engineering Research*.
- Chandani, A., Mehta, M., Mall, A., & Khokhar, V. (2016). "Employee Retention: A Factor of Company's Success." *International Journal of Education and Applied Research*.
- Mishra, M., & Soni, G. (2020). "A Predictive Model for Employee Attrition using Machine Learning." *International Journal of Engineering Development and Research (IJEDR)*.
- Rathi, N., & Lee, K. (2016). "Understanding attrition in organizations: A literature review." *Management Research Review*.
- Alshiddy, M. S., & Aljaber, B. N. (2023). Employee Attrition Prediction using Nested Ensemble Learning Techniques. *International Journal of Advanced Computer Science and Applications*, 14(7). thesai.org
- Mohiuddin, K., Alam, M. A., Alam, M. M., Welke, P., Martin, M., Lehmann, J., & Vahdati, S. (2023). Retention Is All You Need: Explainable HR-DSS for Employee Attrition. arXiv preprint arXiv:2304.03103. arXiv
- Raza, A., Munir, K., Almutairi, M., Younas, F., & Fareed, M. M. S. (2022). Predicting Employee Attrition Using Machine Learning Approaches. *Applied Sciences*, 12(13), 6424. ResearchGate+1
- "Machine Learning Models for Predicting Employee Attrition: A Data Science Perspective" (2025). *Data & Metadata*, 4, 669. ResearchGate
- Chauhan, et al. (2025) „Featuring Machine Learning Models to Evaluate Employee Attrition: A Comparative Analysis of Workforce Stability-Relating Factors“. *International Research Journal of Multidisciplinary Scope (IRJMS)*, 6(2), 862-873. irjms.com

# APPENDIX

## APPENDIX A: SOURCE CODE SNIPPETS

### A.1 Importing required libraries

```
import os
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import LabelEncoder
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import StackingClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score,
confusion_matrix, classification_report, roc_curve, auc
```

### A.2 Loading IBM HR Analytics Dataset

```
data = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/Employee Attrition
Project/data/WA_Fn-UseC_-HR-Employee-Attrition (1).csv')
Displaying basic information
print(data.dtypes)
print(data.info())
print(data.describe())
```

### A.3 Encode categorical variables using Label Encoder

for column in data\_cleaned.select\_dtypes(include=['object']).columns:

```
le = LabelEncoder()
data_cleaned[column] = le.fit_transform(data_cleaned[column])
```

**A.4 Splitting Data into Training and Testing Sets** - Divide the features and target into 80% training and 20% testing sets, scale the features, and separate the features.

```
X = data_cleaned.drop(columns=["Attrition"])
y = data_cleaned["Attrition"]
# train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42, stratify=y)
# Feature scaling
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

**A.5 Train Logistic Regression** - A sigmoid function is used in this linear model for binary classification to predict probabilities. Easy to understand and suitable for baseline.

```
log_model = LogisticRegression(max_iter=2000, random_state=42)
log_model.fit(X_train, y_train)
y_pred = log_model.predict(X_test)
y_prob = log_model.predict_proba(X_test)[:, 1]
# Evaluation
accuracy_lr = accuracy_score(y_test, y_pred)
report_lr = classification_report(y_test, y_pred)
print(f'Accuracy: {accuracy_lr:.4f}')
print(report_lr)
```

**A.6 Ensemble Stacking** - In machine learning, ensemble stacking is a potent meta-learning technique that combines several base models to enhance predictive performance, particularly for classification tasks like employee attrition prediction.

```
base_models = [  
    ('logistic_regression', LogisticRegression(max_iter=3000, random_state=42)),  
    ('decision_tree', DecisionTreeClassifier(random_state=42)),  
    ('gb', GradientBoostingClassifier(n_estimators=120, random_state=42)),  
    ('random_forest', RandomForestClassifier(n_estimators=150, random_state=42)),  
    ('svm', SVC(probability=True, random_state=42))  
]
```

#### **A.7 Define meta-learner (Logistic Regression for simplicity)**

```
meta_learner = LogisticRegression(max_iter=2000, random_state=42)  
# Create stacking classifier  
stacking_model = StackingClassifier(estimators=base_models, final_estimator=meta_learner, cv=5)
```

#### **A.8 Train the stacking model**

```
stacking_model.fit(X_train, y_train)  
# Predict on test data  
y_pred_stacking = stacking_model.predict(X_test)  
# Evaluate  
accuracy_stacking = accuracy_score(y_test, y_pred_stacking)  
report_stacking = classification_report(y_test, y_pred_stacking)  
print(f'Stacking Model Accuracy: {accuracy_stacking:.4f}')  
print("Stacking Model Classification Report:")  
print(report_stacking)
```