

LB/TH/41/2025
TH6007

**ENHANCING EMPLOYEE RETENTION:
AN IMPROVED HYBRID MODEL FOR PREDICTING
EMPLOYEE ATTRITION**

Madhubhashini J M M

239333B

MSc in Computer Science

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
FACULTY OF ENGINEERING

University of Moratuwa
Sri Lanka

June 2025

**ENHANCING EMPLOYEE RETENTION:
AN IMPROVED HYBRID MODEL FOR PREDICTING
EMPLOYEE ATTRITION**

Madhubhashini J M M

239333B

Thesis/Dissertation submitted in partial fulfillment of the requirements for the degree
MSc in Computer Science

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
FACULTY OF ENGINEERING

University of Moratuwa
Sri Lanka

June 2025

DECLARATION

I confirm that this thesis/dissertation is entirely my own work. It does not include, without proper acknowledgment, any content that has previously been submitted for any degree or diploma at any other university or higher education institution. To the best of my knowledge, it contains no material previously published or written by another individual, except where references are clearly provided. I reserve the right to reuse this work in whole or in part in future publications such as books or journal articles.

Date:27-June-2025

The above candidate has carried out research for the PhD/MPhil/Masters thesis/dissertation under my supervision. I confirm that the declaration made above by the student is true and correct.

Name of Supervisor: Prof. Chathura R. De Silva

Signature of the Supervisor:

Date:27-June-2025

DEDICATION

I would like to dedicate the research thesis,
“Enhancing Employee Retention: Hybrid Model for Predicting Employee Attrition”
To,
The Dean of Faculty of Engineering, Prof. Manatunge J.M.A
The Head of Department of Computer Science and Engineering, Dr. Uthayasanker
Thayasivam
Supervisor of the Project, Prof. Chathura R. De Silva,
all the academic and non-academic staff in the Faculty of Engineering who
inspired
and motivated us in various ways.

ACKNOWLEDGEMENT

Here I owe my deep gratitude to the below mentioned people who have supported
and
guided me in different
ways to meet success of the project.
My supervisor Prof. Chathura R. De Silva, for the constant encouragement and
guidance
and enormous support,
that he has always given to me towards fulfilling the aim of the project.
All the academic staff at the Faculty of Engineering of University of Moratuwa,
who
have support, encouragement and guidance throughout my postgraduate years.
All the non-academic staff of Faculty of Engineering of University of Moratuwa,
members of my family, my batch mates and friends who have given support to me in
numerous ways.

ABSTRACT

Employee attrition poses a significant challenge to organizational stability, directly impacting productivity, operational continuity, and long-term strategic goals. While traditional methods of attrition management rely on retrospective analyses, machine learning (ML) offers a proactive approach to predict and mitigate workforce turnover. This study addresses the critical need for robust predictive frameworks by evaluating eight ML classifiers; Support Vector Machines with Radial Basis Function kernel, Fisher's Linear Discriminant, Logistic Regression, Multi-Layer Perceptron, Random Forests, Naive Bayes, Adaboost and XGBoost with IBM HR Analytics Employee Attrition & Performance dataset to forecast employee attrition. A novel hybrid stacking ensemble model is proposed to enhance prediction accuracy by integrating the strengths of individual classifiers. The research investigates the impact of dataset balancing and feature optimization on model performance, emphasizing the role of data preprocessing in improving predictive reliability. Models underwent hyperparameter tuning and stratified 5-fold cross-validation, with the hybrid ensemble (Logistic Regression meta-learner) achieving superior metrics: 95.81% accuracy, 96.69% Precision, 95.77% F1-score, and 99.23% ROC-AUC. Individual models exhibited strong performance, notably XGBoost and LR (93.38% accuracy), while dataset balancing improved minority-class recall by 18–22% and feature selection reduced redundancy by 30%. Key findings reveal that dataset balancing and feature engineering significantly enhance model precision, particularly for minority class identification. The stacking model's superior performance underscores the value of meta-learning in synthesizing heterogeneous classifier outputs. Practical implications include scalable frameworks for HR analytics, enabling early risk identification and personalized retention strategies. The stacking model's adaptability allows integration with HR platforms, transforming reactive practices into proactive decision-making. Methodologically, the study advances predictive human capital management by demonstrating the synergy of hybrid ML approach and interpretable outputs. By bridging technical innovation with organizational psychology, this work offers a replicable blueprint for sustainable workforce management, emphasizing real-time analytics and employee-centric interventions.

Keywords: Machine Learning (ML) Support Vector Machines (SVM). Radial Basis Function (RBF) kernel, Fisher's Linear Discriminant (FLD), Logistic Regression (LR), Multi-Layer Perceptron (MLP), Random Forests (RF), Naive Bayes (NB)

TABLE OF CONTENTS

DECLARATION	i
DEDICATION	ii
ACKNOWLEDGEMENT	iii
ABSTRACT	iv
LIST OF FIGURES	vii
LIST OF TABLES	viii
CHAPTER 1	1
INTRODUCTION	1
CHAPTER 2	2
RESEARCH PROBLEM	2
CHAPTER 3	3
RESEARCH AIM AND OBJECTIVES	3
CHAPTER 4	4
LITERATURE REVIEW	4
CHAPTER 5	8
METHODOLOGY	8
5.1 The Preprocessing Phase	8
5.1.1 Dataset.....	9
5.1.2 Remove Missing and Duplicate Values	10
5.1.3 Remove Unnecessary Features	10
5.1.4 Apply Appropriate Feature Encoding	11
5.1.5 Remove Correlated Feature.....	10
5.1.6 Feature Scaling.....	11
5.1.7 Identify Significant Features with PCA	11
5.2 Model	12
5.2.1 Support Vector Machine (RBF kernel):.....	12
5.2.2 Fisher Linear Discriminant Function:	12
5.2.3 Logistic Regression:.....	12
5.2.4 Multi-Layer Perceptron (MLP) classifier:	13
5.2.5 Random Forest:	13
5.2.6 Naive Bayes:	13
5.2.7 AdaBoost Classifier:	14
5.2.8 XGBoost Classifier:	14
5.3 Performance Evaluation	15

5.3.1	Accuracy:	15
5.3.2	Precision:.....	15
5.3.3	Recall:	15
5.3.4	F1-score:.....	15
5.3.5	The Area Under the Curve (AUC).....	16
5.4	Cross Validation	16
CHAPTER 6		17
RESULTS AND DISCUSSION		17
6.1	Performance Evaluation of Preprocessing Phase	17
6.2	Performance Evaluation of Different Models	19
6.3	Performance Evaluation of Stacked Ensemble Model	24
CHAPTER 7		27
CONCLUSION		27
REFERENCES		28

LIST OF FIGURES

Figure	Description	Page
Figure 5.1	Preprocessing Steps	8
Figure 5.2	Attrition Yes, No cases distribution	10
Figure 6.1	SVM – Imbalanced data set	22
Figure 6.2	SVM – Balanced data set	22
Figure 6.3	FLD – Imbalanced data set	22
Figure 6.4	FLD – Balanced data set	22
Figure 6.5	LR – Imbalanced data set	22
Figure 6.6	LR – Balanced data set	22
Figure 6.7	MLP – Imbalanced data set	22
Figure 6.8	MLP – Balanced data set	22
Figure 6.9	RF – Imbalanced data set	23
Figure 6.10	RF – Balanced data set	23
Figure 6.11	NB – Imbalanced data set	23
Figure 6.12	NB – Balanced data set	23
Figure 6.13	AdaBoost – Imbalanced data set	23
Figure 6.14	AdaBoost – Balanced data set	23
Figure 6.15	XGBoost – Imbalanced data set	23
Figure 6.16	XGBoost – Balanced data set	23
Figure 6.17	The proposed model's architecture	26

LIST OF TABLES

Table	Description	Page
Table 5.1	Data Description of the data set	9
Table 6.1	Accuracy of different models	19
Table 6.2	Precisions of different models	19
Table 6.3	Recall of different models	20
Table 6.4	F1 Score of different models	20
Table 6.5	Different combination of base models for building Stacked Models and their performance	25
Table 6.6	Results of different stacked model approaches	25
Table 6.7	Performance of the current work and the related work	26