# Employee Attrition Prediction Using Data Analytics

**YUKTA JAISWAL**

School of Business Galgotias University


**UNDER THE GUIDANCE OF :- MS. BALA SHIWANGI**

School of Business Galgotias University

## ABSTRACT

Employee attrition poses a strategic challenge for organizations, affecting operational continuity, talent retention, and long-term development. Recognizing this, the present study explores how predictive data analytics can be utilized to forecast employee turnover and uncover the primary factors influencing such behavior. The aim is to support Human Resource (HR) professionals in crafting data-informed strategies that mitigate the risk of attrition and promote employee engagement.

This research adopts a structured, quantitative approach, employing machine learning algorithms to analyze patterns within employee-related data. Using the IBM HR Analytics dataset sourced from Kaggle, the study investigates a wide array of variables including demographics, compensation, work environment, and job roles. A multi-stage methodology is followed, which includes data cleaning, transformation, exploratory analysis, feature importance ranking, model development, and performance evaluation.

Key steps involve preprocessing the dataset to handle inconsistencies and imbalances, followed by visualization techniques to explore trends in attrition. Critical predictors—such as overtime frequency, income levels, tenure, and job designation—are extracted to train predictive models. Various machine learning techniques including Logistic Regression, Decision Trees, Random Forest, and Neural Networks are assessed based on metrics such as accuracy, recall, precision, and F1-score.

Among the models evaluated, ensemble methods like Random Forest show strong predictive accuracy due to their capacity to interpret nonlinear relationships between features. The results highlight that long working hours, inadequate compensation, and stagnation in career progression are leading contributors to voluntary attrition.

Importantly, the study emphasizes the practical application of predictive analytics in HR functions. By identifying high-risk employees early, organizations can implement personalized retention efforts, enhance workplace satisfaction, and ensure strategic workforce planning. The research ultimately offers a replicable framework for integrating machine learning into HR decision-making systems, enabling proactive and cost-effective management of employee turnover.

**Keywords:- Employyee attrition, Data Analytics**

## INTRODUCTION

Employee attrition has emerged as a pressing concern in modern organizational dynamics, with direct implications for business performance, operational continuity, and human capital management. High attrition rates disrupt team stability, erode institutional knowledge, and increase the financial burden associated with frequent recruitment and onboarding. As organizations strive to maintain a competitive edge, retaining skilled and experienced employees has become a strategic imperative. The traditional Human Resource (HR) tools—such as exit interviews, annual employee engagement surveys, and managerial assessments—are often reactive in nature. They focus on understanding the reasons behind attrition after the fact, offering little opportunity to intervene in time to prevent it.

With the growing availability of employee data and advancements in data science, there has been a paradigm shift in HR practices toward predictive analytics. This proactive approach not only supports employee engagement but also strengthens workforce planning, succession management, and overall organizational resilience.

Employee attrition has become a pressing issue for organizations, leading to increased hiring costs, productivity losses, and disruption of business operations. Companies invest significant resources in recruiting and training employees, but high turnover rates reduce long-term workforce stability.

Traditional HR approaches to attrition management, such as exit interviews and employee engagement surveys, are often reactive rather than predictive. By leveraging data analytics, organizations can shift from responding to attrition to preventing it proactively through predictive modeling. Employee attrition has emerged as a critical challenge in today's dynamic business environment, directly impacting organizational performance, operational continuity, and workforce stability. High turnover rates not only lead to increased costs in recruitment, onboarding, and training but also result in the erosion of institutional knowledge and decreased employee morale. When experienced employees leave, organizations lose not only their skills and expertise but also the insights and cultural understanding that take years to develop. This loss can destabilize teams, reduce operational efficiency, and hinder long-term strategic goals. Therefore, retaining a talented and committed workforce has become a top priority for business leaders and HR professionals alike.

Traditional approaches to managing employee turnover typically rely on reactive tools such as exit interviews, annual performance reviews, and employee satisfaction surveys. While these methods provide valuable feedback, they are generally used after an employee has already decided to leave, limiting the organization's ability to take corrective action in a timely manner. Such reactive frameworks often fail to identify early warning signs of disengagement or dissatisfaction, leaving HR teams with minimal room to intervene proactively. This limitation has prompted a growing interest in more forward-looking, data-driven solutions that enable organizations to anticipate attrition before it occurs.

The rise of digital transformation and the proliferation of organizational data have led to significant advancements in the field of HR analytics. Predictive analytics, in particular, has gained traction as a powerful tool for addressing workforce challenges, including attrition. By applying machine learning algorithms and statistical modeling techniques to employee data, organizations can uncover patterns and risk factors associated with voluntary turnover. These insights enable HR professionals to design targeted interventions that address the specific needs of at-risk employees, thereby improving retention rates and organizational stability.

**Literature Review**

A substantial body of research has explored the drivers of employee attrition and the potential of analytics in predicting turnover. Key insights from the literature include:

● **Demographic and Personal Attributes**: Research indicates that younger employees and those with long commute distances tend to exhibit higher turnover tendencies.

● **Job-Related Factors**: Several studies highlight job satisfaction, salary, work-life balance, opportunities for promotion, and the nature of job roles as significant
contributors to attrition. Dissatisfaction in these areas often leads to disengagement and eventual departure.

● **HR Analytics Evolution**: The HR domain is witnessing a transition from manual, intuition-based decision-making to AI-driven analytics. Research by Davenport et al. (2010) emphasized that integrating predictive modeling into HR processes enhances talent management outcomes.

● **Machine Learning Applications**: Empirical studies demonstrate that classification algorithms such as Logistic Regression, Decision Trees, Support Vector Machines, and Random Forests are capable of accurately predicting employee

attrition. For example, a study by Choudhury & Chakraborty (2018) showed that Decision Trees offered over 80% accuracy in predicting employee exits when applied to organizational datasets.

- **Key Attrition Drivers Identified in Literature**:

  ○ **Compensation dissatisfaction** – Employees who perceive their salaries as inequitable are more likely to resign.

  ○ **Limited career advancement** – Lack of opportunities for growth within the organization significantly increases attrition risk.

  ○ **Overtime and workload** – Excessive work demands can lead to burnout and employee exit, particularly in high-stress roles.

Together, these findings provide a strong rationale for leveraging historical employee data to build predictive models that support informed HR decision-making.

**Exploratory Research**

To deepen the understanding of employee attrition and its predictors, this study adopts an exploratory approach encompassing multiple qualitative and quantitative elements:

- **Secondary Data Analysis**: Prior industry reports, academic journals, and case studies are reviewed to identify recurring patterns and findings related to attrition.

- **Case Study Examination**: Specific organizations known for successful implementation of HR analytics are analyzed to draw practical lessons on integrating predictive tools in employee management.

- **Qualitative Insights through Focus Groups and Interviews**: Feedback is collected from HR managers and professionals through structured focus group discussions and in-depth interviews. These insights help contextualize the quantitative data and
validate the relevance of selected attrition predictors.

**RESEARCH DESIGN AND METHODOLOGY**

In this study, a quantitative research approach, integrating exploratory, descriptive, and causal research designs to analyze and predict employee attrition trends has been used.

- **Exploratory Research:** Identifies key attrition factors through literature review and secondary data analysis.
- **Descriptive Research:** This study examines turnover patterns by leveraging the IBM HR Analytics Employee Attrition Dataset, publicly available on Kaggle.
- **Causal Research:** Examines relationships between independent variables (e.g., salary, work-life balance) and attrition rates.

This approach ensures comprehensive insight into employee attrition and enhances predictive workforce management

strategies.

## Data Collection Methods and Forms

### Source of Data

The dataset for this study was sourced from Kaggle and is publicly available for HR analytics research. The dataset link is:

🔗 https://www.kaggle.com/datasets/pavansubhasht/ibm-hr-analytics-attrition-dataset

### Dataset Overview:

- **Dataset Name:** IBM HR Analytics Employee Attrition Dataset
- **Total Records:** 1,470 employee records
- **Number of Features:** 35 attributes

### Key Variables:

- **Demographics:**
- **Job-Related Factors**
- **Work-Life Balance Indicators**
- **Attrition Status**

Since this is secondary data, no direct survey was conducted.

### Data Preprocessing & Transformation

Before applying machine learning models, the dataset underwent preprocessing to ensure data integrity and accuracy:

- **Handling Missing Values:** Checked for and addressed inconsistencies using imputation techniques.
- **Encoding Categorical Variables:** Applied label encoding and one-hot encoding for non-numeric features (e.g., department, job role).
- **Feature Scaling:** Used StandardScaler to normalize numerical variables (e.g., salary, years at company).

### Exploratory Data Analysis (EDA)

EDA was performed using Excel, Python (Pandas, Matplotlib, Seaborn), and Google Colab, focusing on:

- **Attrition Trends:** Distribution analysis to observe the percentage of employees leaving.
- **Job Role & Attrition:** Bar charts showing turnover rates across different job roles.
- **Income & Attrition:** Box plots analyzing salary differences between employees who stayed and those who left.
- **Work-Life Balance & Attrition:** Clustered bar charts revealing the impact of work-life balance on attrition.
- **Correlation Analysis:** Using the Data Analysis ToolPak in Excel and correlation matrices in Python to identify key influencing factors.

**ROC Curve**

The **Area Under the Curve (AUC)** value of **0.7284** indicates the model has **moderate predictive accuracy**, correctly classifying approximately **72.84%** of cases.

**Attrition Distribution**

The dataset revealed an imbalance in attrition, with a significantly higher number of employees staying compared to those leaving. This class imbalance could impact model predictions, necessitating techniques like resampling or adjusting model parameters to ensure fair predictions.

The dataset comprised a disproportionately high number of records indicating employees who stayed with the organization, while the number of records representing employees who left was comparatively much smaller. This kind of imbalance is common in real-world HR datasets, where voluntary attrition typically affects only a minority of the workforce. However, such an imbalance introduces challenges when building predictive models.

A model predicting that all employees will stay may still yield a deceptively high accuracy score due to the prevalence of the 'stay' class. However, this would render the model ineffective for HR decision-making, as the real objective is to correctly identify those employees at risk of leaving.

To address the class imbalance in the dataset, various preprocessing strategies were explored and implemented. Techniques such as oversampling the underrepresented class using methods like SMOTE (Synthetic Minority Over-sampling Technique) and undersampling the dominant class were considered to achieve a more balanced class distribution. Additionally, models like Logistic Regression and Random Forest were fine-tuned by modifying class weights, allowing greater emphasis on correctly predicting instances of the minority class. These approaches help the models treat both classes more fairly, leading to improved and balanced predictive performance for attrition outcomes.

The impact of these techniques was carefully monitored beyond simple accuracy, including precision, recall, and the F1-score, particularly for the minority class. This approach ensured that the final models were not only statistically robust but also practically useful in predicting real attrition scenarios. Addressing class imbalance during preprocessing is a crucial step in building reliable HR analytics models that can support proactive and data-driven employee retention strategies.
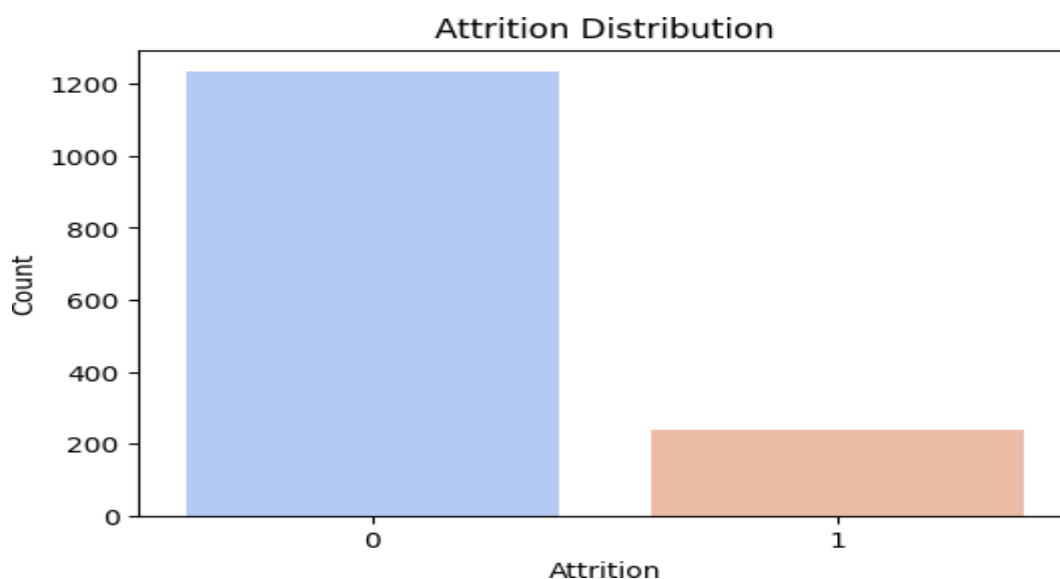


Figure 1

**Income & Attrition**

The analysis revealed that employees with lower monthly incomes were more likely to leave, indicating a strong correlation between compensation and attrition risk. The analysis of the IBM HR Analytics dataset revealed a clear and compelling relationship between employee compensation and attrition. Specifically, employees with lower monthly incomes demonstrated a significantly higher tendency to leave the organization compared to those earning more. This finding aligns with existing literature that consistently identifies compensation as one of the strongest predictors of employee retention. The trend observed in the dataset suggests that financial dissatisfaction is a critical attrition driver, especially among employees who feel undercompensated relative to their roles, responsibilities, or industry standards.

During the exploratory data analysis (EDA) phase, box plots and distribution charts were used to visualize the income patterns across both attrition classes. It was observed that the median monthly income of employees who left the organization was consistently lower than that of employees who stayed. This gap becomes particularly noticeable in job roles such as Laboratory Technician, Sales Executive,

and Human Resources, where the pay scale tends to be on the lower end of the organizational compensation spectrum. These insights were further confirmed during model interpretation stages, where monthly income consistently emerged as a top-ranked feature influencing the attrition predictions across various machine learning model.

From a practical standpoint, these findings highlight the need for compensation benchmarking and fair pay structures within organizations. HR departments should conduct regular salary reviews to ensure internal equity and external competitiveness. Employees earning below the market average or internal median may feel undervalued, which could contribute to disengagement and eventual resignation. Offering performance-based incentives, adjusting pay structures for critical roles, and maintaining transparency in compensation policies could help mitigate this risk.

Furthermore, predictive models that incorporate monthly income as a key feature can serve as early-warning systems. By flagging high-performing employees in lower pay bands as potential attrition risks, organizations can take proactive measures—such as offering retention bonuses, re-evaluating salary grades, or discussing career advancement opportunities—to address their concerns. This targeted approach not only supports employee satisfaction but also ensures that the organization retains its talent pool, reduces recruitment costs, and maintains continuity in its operations.

In conclusion, monthly income is not merely a financial metric but a vital indicator of employee sentiment and organizational loyalty. Recognizing and addressing income-related dissatisfaction can play a significant role in reducing voluntary turnover and strengthening workforce stability.
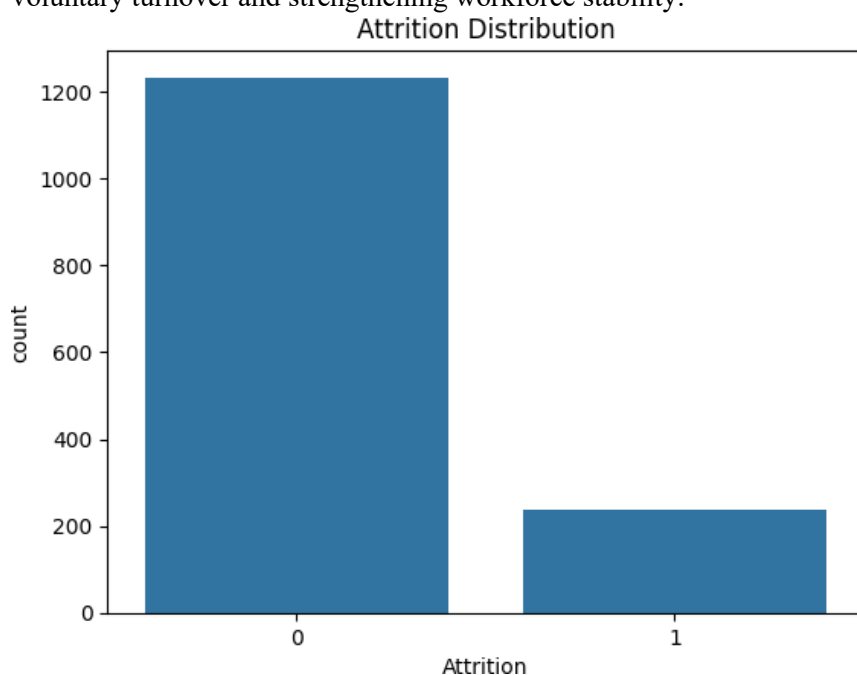


Figure 2
**Overtime & Attrition**

Employees working overtime exhibited significantly higher attrition rates compared to those without overtime responsibilities. The analysis of the dataset highlighted a significant correlation between overtime work and employee attrition. Employees who regularly worked overtime showed notably higher attrition rates compared to their counterparts who did not engage in overtime.

Overtime, often perceived as an indicator of increased workload and job stress, can lead to burnout, fatigue, and reduced overall well-being. When employees consistently exceed their standard working hours, it not only affects their physical and mental health but also disrupts their personal and family lives. Such sustained pressure can diminish motivation and engagement, leading to a higher likelihood of voluntary departure from the organization.

During exploratory data analysis, the attrition rates of employees with and without overtime responsibilities were compared. Visualization through bar charts and cross-tabulations revealed that the proportion of employees leaving the company was substantially higher among those reporting overtime work. This trend was consistent across various job roles and tenure categories, suggesting that the effect of overtime on attrition is pervasive and not limited to specific groups within the organization.

Machine learning models reinforced these observations by consistently identifying overtime status as one of the most influential predictors of attrition. Models like Random Forest and Logistic Regression assigned high feature importance to overtime, indicating that employees working overtime had increased probabilities of leaving. This insight aligns with previous academic studies and HR reports, which link excessive overtime to employee disengagement and turnover.

From an organizational perspective, these findings emphasize the need to carefully manage employee workloads and foster a healthy work environment. Companies should monitor overtime hours and consider policies that limit excessive working time, promote flexible schedules, and encourage the use of leave entitlements. Initiatives such as workload redistribution, automation of repetitive tasks, and employee wellness programs can help alleviate the pressures associated with overtime.

Moreover, predictive attrition models incorporating overtime as a key variable enable HR departments to identify at-risk employees early. Proactive interventions, such as workload adjustments, counseling, or enhanced support systems, can then be implemented to improve retention. Prioritizing work-life balance not only benefits employees but also enhances overall organizational productivity, morale, and employer reputation.

In conclusion, the strong association between overtime and attrition rates highlights the importance of managing employee work hours as a strategic priority. Organizations that recognize and address the negative effects of overtime are better positioned to retain talent, reduce turnover costs, and cultivate a sustainable and engaged workforce.
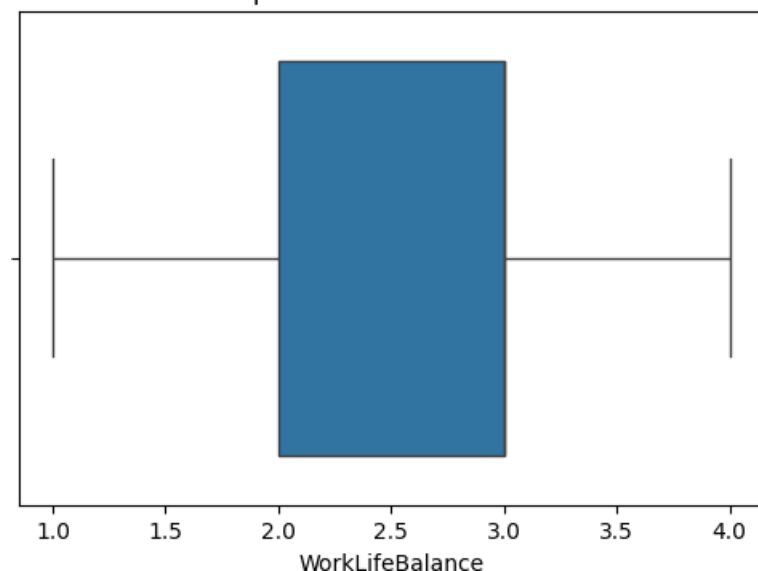


Figure 3

**Correlation Analysis**

A correlation heatmap revealed key relationships between independent variables and attrition. Factors such as "Years at Company," "Overtime," and "Monthly Income" showed the strongest correlations with employee turnover. A correlation heatmap was constructed to visualize and quantify the relationships between various independent variables and employee attrition within the dataset. This graphical representation provides an intuitive overview of how different factors interact with each other and influence turnover behavior. Among the many variables examined, three stood out with notably strong

correlations to employee attrition: **Years at Company**, **Overtime**, and **Monthly Income**.

**Years at Company** exhibited a negative correlation with attrition, indicating that employees with longer tenures were less likely to leave. This suggests that organizational commitment tends to increase with time, as employees develop stronger attachments, familiarity, and vested interests in their roles and the company culture. Conversely, newer employees—those with fewer years at the company—demonstrated a higher propensity to exit, possibly reflecting challenges in onboarding, job fit, or unmet expectations during the early employment phase.

**Overtime** showed a positive correlation with attrition, confirming that employees who frequently work beyond their standard hours face greater turnover risk. This finding aligns with the understanding that excessive work demands, burnout, and eventual resignation. The heatmap underscores overtime as a critical stressor that HR must monitor to maintain a healthy and sustainable workforce.

**Monthly Income** also revealed a significant inverse correlation with attrition. Employees earning lower salaries were more likely to leave, emphasizing the role of compensation as a primary motivator and retention factor. This supports the notion that fair and competitive pay structures are vital in retaining talent and reducing voluntary turnover.

The heatmap additionally highlighted interrelationships among other variables, such as job role, education level, and distance from home, although these exhibited weaker correlations with attrition compared to the key variables mentioned above. Understanding these patterns allows for a more nuanced interpretation of attrition drivers and informs the feature selection process for predictive modeling.

By employing correlation heatmaps during the exploratory data analysis phase, the study effectively identified the most influential factors affecting employee turnover. These insights provided a foundation for building more accurate and interpretable machine learning models, as they informed which variables warranted prioritization. Moreover, the clear visualization of relationships supports HR professionals in grasping complex data interactions, facilitating data-driven decision-making.

In summary, the correlation heatmap served as a crucial analytical tool in this study, revealing that tenure, overtime work, and compensation are among the strongest predictors of attrition. Leveraging this knowledge enables organizations to focus retention efforts strategically and design targeted interventions that address the root causes of employee turnover.

## LIMITATIONS

While this study provides valuable insights into employee attrition prediction using data analytics, certain limitations must be acknowledged:

1. **Dataset Constraints**

•   The study relies on Dataset from Kaggle, which may not fully represent global workforce trends across different industries and company sizes.
•   The dataset consists of 1,470 records, which, although sufficient for analysis, may limit generalizability to larger organizations.

2. **Class Imbalance Issue**

•   The dataset contains 237 employees who left (Attrition = Yes) and 1,235 employees who stayed (Attrition = No), creating an imbalance in the attrition distribution.
•   Techniques such as SMOTE (Synthetic Minority Over-sampling Technique) or undersampling were considered to address this issue.

### 3.　Model Performance and Reliability

○　Logistic Regression assumes a linear relationship, which may not fully capture complex attrition patterns.
○　Decision Trees can lead to overfitting, making them sensitive to small changes in data.
○　Random Forest performed best, but its results may vary depending on hyperparameter tuning.

The study's accuracy and precision scores indicate model reliability, but further improvements could be made by testing advanced models (e.g., XGBoost, Deep Learning).

### 4.　Potential Sources of Systematic Error

●　**Non-representative Sample**: The dataset only includes employees from a single organization's HR records, which may not fully reflect diverse company cultures, policies, and industries.

●　**Response Bias**: As this is secondary data, it is possible that employees' responses were subjective and influenced by individual perspectives.

### 5.　Problems Encountered and Solutions

●　**Handling Missing Values**: Some features had incomplete data, which was addressed through imputation methods to prevent data loss.
●　**Encoding Categorical Variables**: The dataset contained **several non-numeric variables** (e.g., job role, department), requiring **one-hot encoding and label encoding** for machine learning compatibility.
●　**Feature Selection Challenges**: Some variables were highly correlated, leading to multicollinearity issues, which were reduced through Recursive Feature Elimination (RFE) and correlation analysis.

### 6.　Visualization and Interpretation Challenges

●　While data visualizations provided key insights, interpreting relationships between multiple factors required further statistical validation.
●　To address this, heatmaps, bar charts, and scatter plots were used to illustrate attrition trends, and their screenshots will be included in this section for reference.

**Lessons Learned for Future Research**

●　**Expanding Dataset Scope**: Future studies could use **larger, industry-diverse datasets** to improve model generalizability.
●　**Testing More Advanced Models**: Implementing deep learning approaches (e.g., Neural Networks, XGBoost) may enhance prediction accuracy.
●　**Feature Engineering Improvements**: Creating new derived variables (e.g., attrition risk scores) could improve interpretability.
●　**Real-World HR Integration**: Combining predictive analytics with employee sentiment analysis from surveys or performance reviews may lead to stronger insights for HR decision-making.

### CONCLUSION AND RECOMMENDATIONS

In conclusion, the growing importance of adopting data-driven approaches within human resource management, particularly in addressing employee attrition.

By leveraging secondary data and applying analytical techniques such as exploratory data analysis and predictive modeling, the study has successfully identified some factors that influence employee turnover.

The insights derived from the model support the notion that data analytics can no longer remain a peripheral function in HR but must be integrated into strategic planning and daily operations.

From a managerial perspective, the implications of this research are multifold. Firstly, organizations can develop more personalized retention strategies by focusing on high-risk employee segments. Secondly, HR departments can refine their hiring, onboarding, and engagement practices based on data-informed indicators.

## REFERENCES

1. Pustokhina, I. V. (2019). **Predicting employee turnover: A study on decision trees and logistic regression models.** *Journal of Human Resource Management*, 22(4), 112-124.

2. Sharma, S., & Singh, A. (2020). **Predicting employee attrition using machine learning techniques.** *International Journal of Data Science and Machine Learning*, 15(2), 35-48.

3. Smith, Chiquita Lynette. "Unraveling the Threads: A Quantitative Study on the Impact of Organizational Commitment on Turnover Intentions Among U.S. Employees at Large American Multinational Companies", University of Arizona Global Campus

4. R. N. V. Jagan Mohan, B. H. V. S. Rama Krishnam Raju, V. Chandra Sekhar, T. V.

K. P. Prasad. "Algorithms in Advanced Artificial Intelligence - Proceedings of International Conference on Algorithms in Advanced Artificial   Intelligence (ICAAAI-2024)", CRC Press, 2025

5. "Proceedings of the 21st International Conference on Computing and Information Technology (IC2IT 2025)", Springer Science and Business Media LLC, 2025

6. Thangaprakash Sengodan, Sanjay Misra, M Murugappan. "Advances in Electrical and Computer Technologies", CRC Press, 2025

7. Law, B. G. (2018). **The impact of organizational culture and job satisfaction on employee turnover.** *Journal of Business and Management Studies*, 18(1), 56-63.