

EMPLOYEE CHURN PREDICTION

A Project Report

submitted in partial fulfillment of the requirements

of

fundamentals with cloud computing and gen AI

by

K.SURYAPRAKASH

Prakashsurya3003@gmail.com

1A0E77C91FE6AEDB0914C7FF8E0899F5 (aut2291240012)

912421114306

Under the Guidance of

P.RAJA

Master trainer, Edunet Foundation

ACKNOWLEDGEMENT

We would like to take this opportunity to express our deep sense of gratitude to all individuals who helped us directly or indirectly during this thesis work.

Firstly, we would like to thank my supervisor, P.Raja (Master Trainer of Edunet foundation) and Dr.M. Anto Alosius (Assistant professor) for being a great mentor and the best adviser I could ever have. His advice, encouragement and the critics are a source of innovative ideas, inspiration and causes behind the successful completion of this project. The confidence shown in me by him was the biggest source of inspiration for me. It has been a privilege working with him for the last one year. He always helped me during my project and many other aspects related to the program. His talks and lessons not only help in project work and other activities of the program but also make me a good and responsible professional.

ABSTRACT of the Project

Employee churn refers to the voluntary or involuntary departure of employees from an organization, posing significant challenges for businesses. Understanding and predicting employee churn is crucial for maintaining a stable workforce, enhancing productivity, and reducing recruitment costs. Unlike customer churn, where businesses cannot select their customers, employee retention is a choice made by organizations. This project aims to analyze employee churn through a structured approach that includes “*exploratory analysis, data visualization, cluster analysis, model building, and performance evaluation.*” The process begins with a thorough examination of the dataset to identify patterns and trends related to employee turnover. Visualization techniques will be employed to highlight key factors influencing churn and to present insights effectively.

Next, cluster analysis will be conducted to group employees with similar characteristics, providing deeper understanding of the employee segments at risk of leaving. Following this, a predictive model will be developed using machine learning algorithms to forecast churn based on identified features. The model's performance will be evaluated using metrics such as “*accuracy, precision, recall, and F1-score, ensuring its reliability for practical application.*” By leveraging data-driven insights, organizations can implement targeted strategies to enhance employee retention, ultimately leading to a more engaged workforce and reduced operational disruptions.

TABLE OF CONTENTS

Abstract	i
List of Figures	v
List of Tables	v
Chapter 1. Introduction	1
1.1 Importance of Employee Churn Prediction.....	1
1.2 Key Steps in Employee Churn Prediction	1
1.3 Problem Statement	4
1.4 Motivation	4
1.5 Objectives	5
1.6 Scope of the Project	5
Chapter 2. Literature Survey	7
2.1 Review relevant literature or previous work in this domain.	7
2.2 Mention the existing models, techniques, or methodologies related to the problem.	8
2.3 Highlight the gaps or limitations in existing solutions and my project will address them.	9
Chapter 3. Proposed Methodology.....	11
3.1 System Design.	11
3.2 Module Used.	12
3.3 Data Flow Diagram.	13
3.4 Advantages.	16
3.5 Requirement Specification.	17
Chapter 4. Implementation and Results	19
4.1 Results of Employee churn prediction.....	19

Chapter 5. Discussion and Conclusion	23
5.1 Key Finding	23
5.2 Git Hub Link of the Project.	23
5.3 Video Recording of Project Demonstration.	23
5.4 Limitations.	23
5.5 Future Work.	24
5.6 Conclusion	24
References	26

LIST OF FIGURES

FIGURE NO	TITLE	PAGE NO.
1	Data Flow Diagram	15
2	Display the first few lines.	19
3	Summary Statistics	19
4	Check Missing Values	20
5	Churn Distribution	20
6	Correlation Matrix	21
7	Elbow curve	21
8	Confusion Matrix	22
9	ROC Curve	22

LIST OF TABLES

TABLE NO	TITLE	PAGE NO.
1	A comparison table of model performance	3
2	Differences Between Employee and Customer Churn	3
3	Scope of the Project	6

CHAPTER 1

INTRODUCTION

Employee churn, or employee turnover, refers to the departure of employees from an organization. This can be due to a variety of factors, such as dissatisfaction with their role, better opportunities elsewhere, or personal reasons. Unlike customer churn, where businesses don't have much control over who becomes their customer, employee churn presents a deeper challenge since the company actively selects its employees. Losing an employee is not just a matter of replacing a headcount—each employee represents a significant investment in training, knowledge, and organizational culture.

The cost of replacing an employee can be high, both in terms of time and money. It involves advertising job roles, conducting interviews, onboarding new hires, and giving them time to adjust and become productive. For this reason, businesses focus on minimizing employee churn as much as possible. Identifying the factors leading to employee churn and creating predictive models to anticipate and prevent churn is essential for long-term success.

1.1 Importance of Employee Churn Prediction

Employee churn has far-reaching effects on any organization. It impacts team morale, increases costs, and disrupts productivity. Predicting churn can enable organizations to intervene and address underlying issues, ultimately retaining talent. By using data-driven approaches, organizations can preemptively identify employees at risk of leaving and take targeted actions to improve retention.

1.2 Key Steps in Employee Churn Prediction

1. Exploratory Analysis

Exploratory data analysis (EDA) helps us uncover patterns, anomalies, and trends in the employee dataset. We look at various factors like employee demographics, job satisfaction, work experience, and salary to understand what might lead to an employee leaving. This phase often involves:

- Checking data distributions.
- Identifying missing or inconsistent data.

- Analyzing relationships between variables, such as age, experience, and likelihood of leaving.

2. Data Visualization

Data visualization allows us to see trends more clearly and understand the impact of various factors on employee churn. Charts, graphs, and heatmaps help in visualizing correlations and patterns that might not be immediately obvious. For example:

- **Bar charts** can show which departments have the highest churn.
- **Heatmaps** can show correlations between factors like salary, years at the company, and job satisfaction.

3. Cluster Analysis

Cluster analysis groups employees based on similar characteristics, such as salary, job satisfaction, and years of service. This method helps in identifying which groups of employees are more likely to leave. For instance:

- Employees in the lower salary range with low job satisfaction might form a high-risk cluster.

Using cluster analysis, we can create targeted retention strategies for different employee segments.

4. Building a Prediction Model

Once we've gathered insights through EDA and clustering, we can build a machine learning model to predict employee churn. Common algorithms used for this purpose include:

- **Logistic Regression:** A simple yet effective model for binary classification tasks.
- **Decision Trees and Random Forests:** These help in understanding how various features impact churn decisions.
- **Support Vector Machines (SVM):** Good for distinguishing between employees who are likely to leave and those who aren't.

5. Evaluating Model Performance

The success of any predictive model depends on how accurately it can forecast churn. We use evaluation metrics like:

- **Accuracy:** How many predictions were correct?
- **Precision:** Of those predicted to leave, how many actually left?
- **Recall:** Of all who left, how many were correctly predicted?
- **F1-Score:** A balance between precision and recall.

A comparison table of model performance might look like this:

Model	Accuracy	Precision	Recall	F1-Score
Logistic Regression	0.85	0.81	0.78	0.79
Random Forest	0.89	0.86	0.84	0.85
SVM	0.88	0.84	0.83	0.83

Table 1: A comparison table of model performance

Key Differences Between Employee and Customer Churn

Aspect	Employee Churn	Customer Churn
Selection	The company chooses who to hire.	The company cannot select its customers.
Impact on Operations	Employees drive the day-to-day business.	Customers impact revenue and brand image.
Cost of Replacement	High cost and time for hiring and training.	High cost of acquiring new customers.
Long-term Impact	Affects morale, culture, and internal growth.	Affects sales and market share.

Table 2: Differences Between Employee and Customer Churn

1.3 Problem Statement:

You are tasked to perform Employee Churn prediction in Python. Employee churn can be defined as a leak or departure of an intellectual asset from a company or organization. or in simple words, you can say, when employees leave the organization is known as churn. The following points help you to understand, employee and customer churn in a better way:

- The business chooses the employee to hire someone while in marketing you don't get to choose your customers.
- Employees will be the face of your company, and collectively, the employees produce everything your company does.
- Losing a customer affects revenues and brand image. acquiring new customers is difficult and costly compared to retaining existing customers. Employee churn is also painful for companies in organizations. It requires time and effort to find and train a replacement. You are going to cover the following steps

- Exploratory Analysis
- Data Visualization
- Cluster Analysis
- Building Prediction Model
- Evaluating Model Performance this statement

1.4 Motivation:

Employee churn, or the departure of valuable employees, significantly disrupts operations, reduces productivity, and incurs high costs for recruitment and training. Predicting and preventing churn is critical for organizations aiming to maintain a stable and experienced workforce.

This project leverages data analytics and machine learning to identify patterns and factors leading to employee churn, offering insights that businesses can use to retain employees and minimize turnover.

By understanding the causes of churn, companies can take proactive steps to mitigate it, safeguarding both their operational continuity and employee investment. The potential applications of this project are broad, particularly in human resource management. Organizations can use the churn prediction model to identify employees at risk of leaving and implement targeted retention strategies, such as enhancing career development opportunities or improving work conditions. This predictive capability enables HR teams to prioritize interventions, reduce replacement costs, and retain talent. Moreover, by addressing the factors contributing to churn, companies can enhance overall employee satisfaction, fostering a more engaged workforce and contributing to long-term organizational success.

1.5 Objective:

The primary objective of this project is to predict employee churn using Python, helping organizations identify employees at risk of leaving and take proactive steps to retain them. The project aims to explore and analyze employee-related data to uncover patterns and factors that contribute to churn, such as job satisfaction, compensation, work environment, and performance metrics. By conducting exploratory analysis, data visualization, and cluster analysis, the goal is to gain insights into the key drivers of employee turnover. Additionally, the project seeks to develop a machine learning prediction model that accurately forecasts which employees are most likely to leave the organization.

The model will be evaluated based on various performance metrics to ensure its effectiveness in predicting churn. Ultimately, the project aims to equip organizations with a data-driven tool to reduce employee turnover, improve retention strategies, and enhance overall workforce stability and productivity.

1.6 Scope of the Project:

This project focuses on predicting employee churn, a critical issue that impacts business productivity, morale, and incurs hiring and training costs.

The goal is to identify patterns that lead to churn and develop a machine learning model to predict which employees are most likely to leave. The process begins with exploratory data analysis (EDA) to review factors like demographics, job satisfaction, and performance, followed by data cleaning and handling inconsistencies.

Data visualization, using techniques like histograms and heat maps, will reveal relationships between variables and churn. Cluster analysis, through methods like K-means, will group employees based on similar characteristics to identify high-risk groups. The core of the project is building a predictive model using algorithms such as logistic regression, decision trees, and Random Forest. The model will be fine-tuned for accuracy and interpretability, making it useful for HR teams.

Finally, the model's performance will be evaluated using metrics like accuracy, precision, recall, and ROC-AUC. Cross-validation will ensure robustness, and a confusion matrix will offer insights into its predictive capabilities. This project aims to provide a comprehensive solution for predicting and reducing employee churn, helping organizations retain talent through data-driven insights.

Step	Objectives	Tasks
1. Exploratory Analysis	Understand the dataset, identify patterns, and detect anomalies.	<ul style="list-style-type: none"> - Load and clean the dataset. - Analyze basic statistics (mean, median, etc.). - Identify missing values.
2. Data Visualization	Create visual representations to highlight trends and insights.	<ul style="list-style-type: none"> - Use histograms, bar charts, and pie charts to visualize distributions.
3. Cluster Analysis	Identify segments of employees based on characteristics.	<ul style="list-style-type: none"> - Normalize the data. - Apply clustering algorithms (e.g., K-means, hierarchical clustering).
4. Building Prediction Model	Develop a predictive model to forecast employee churn.	<ul style="list-style-type: none"> - Split the dataset into training and testing sets. - Train the model on the training set.
5. Evaluating Model Performance	Assess the accuracy and reliability of the prediction model.	<ul style="list-style-type: none"> - Use metrics such as accuracy, precision, recall, F1-score, and AUC-ROC curve. - Fine-tune the model parameters for better performance.

Table3: Scope of the project.

CHAPTER 2

Literature Survey

Employee churn prediction is a crucial aspect of human resource analytics, aimed at identifying the reasons employees leave organizations and developing models to predict churn. While businesses have control over whom they hire, employee retention poses challenges, similar to customer retention in marketing. Losing employees is not only costly in terms of recruitment, training, and onboarding but also disrupts productivity and affects team morale. Research in this area highlights key factors influencing churn, such as job satisfaction, compensation, career growth, and company culture. Early prediction models, like logistic regression and decision trees, focused on these factors to estimate churn risk, while more recent studies employ ensemble techniques like Random Forest, XGBoost, and neural networks to handle complex datasets and improve accuracy. Exploratory data analysis (EDA) and data visualization are key to understanding trends, while clustering techniques, such as K-means, help segment employees based on shared characteristics, revealing which groups are more likely to churn. Evaluating prediction models using metrics like accuracy, precision, recall, and ROC-AUC ensures that these tools provide actionable insights to help organizations retain their valuable talent.

2.1 Review relevant literature or previous work in this domain.

Employee churn prediction is a vital area of HR analytics, aimed at identifying factors that cause employees to leave and predicting those likely to churn. This is crucial for companies as losing employees impacts operational efficiency and incurs significant costs in recruitment and training. Factors like job satisfaction, compensation, work-life balance, and career growth opportunities are key predictors of churn.

Disengaged or dissatisfied employees are more prone to leave, making proactive retention strategies essential for businesses. Traditional models like logistic regression and decision trees were initially used for predicting churn based on employee demographics, job role, salary, and satisfaction scores. More recent

approaches include advanced ensemble models like Random Forest and XG Boost, which handle complex datasets more effectively.

Additionally, clustering techniques like K-means help group employees with similar traits, allowing businesses to target specific segments for retention efforts. Deep learning models, such as neural networks, have also been explored for predicting churn, particularly in dynamic environments with time-series data.

A typical approach involves exploratory data analysis (EDA) to uncover patterns and relationships between variables, followed by data visualization techniques like histograms and heat maps to identify trends. Cluster analysis groups employees based on common traits to pinpoint which segments are at higher risk of churn. Machine learning models such as logistic regression, Random Forest, and XG Boost are then used to predict churn, with performance evaluated using metrics like accuracy, precision, recall, and ROC-AUC. Cross-validation ensures model robustness and avoids over fitting.

In summary, employee churn prediction combines data exploration, clustering, machine learning, and model evaluation to provide businesses with actionable insights. By leveraging modern machine learning techniques, organizations can more accurately predict churn and implement effective strategies to improve employee retention.

2.2 Mention the existing models, techniques, or methodologies related to the problem.

Several existing models, techniques, and methodologies have been applied to the problem of employee churn prediction, drawing from both human resource analytics and machine learning. Logistic regression has been widely used as a traditional method to estimate the likelihood of employee churn based on various factors such as job satisfaction, salary, and performance.

Decision trees and Random Forest models are also common, offering the advantage of interpretability by identifying key features that influence churn. Ensemble methods like Gradient Boosting Machines (GBM) and XG Boost have proven highly effective due to their ability to handle large datasets and capture complex patterns between variables.

Clustering techniques, such as K-means clustering and hierarchical clustering, are often used to segment employees into groups based on similarities in characteristics, helping to better understand the specific groups that are more prone to leaving. More advanced methodologies, such as deep learning with neural networks, have also been applied to employee churn prediction, particularly in organizations dealing with large-scale or time-series data. These methods allow organizations to develop more accurate and tailored models for predicting churn while leveraging past data to intervene and improve employee retention strategies.

2.3 Highlight the gaps or limitations in existing solutions and my project will address them.

Existing solutions for employee churn prediction, while effective, have certain limitations and gaps that this project aims to address. Traditional models like logistic regression and decision trees, though simple and interpretable, often fail to capture the complexity of factors that influence churn, especially. When dealing with non-linear relationships between features like job satisfaction, salary, and career progression. Additionally, many existing models lack scalability and struggle with large datasets, limiting their effectiveness in organizations with diverse employee profiles. Ensemble methods like Random Forest and XG Boost improve accuracy but can be computationally intensive, making them difficult to implement in real-time prediction scenarios.

Another gap in current approaches is that they often focus on a narrow set of features, ignoring important dimensions such as employee engagement, sentiment from performance reviews, or external factors like market conditions that could impact churn. Furthermore, many churn prediction models don't incorporate clustering to identify distinct employee groups with varying churn risks, resulting in one-size-fits-all retention strategies that are less effective.

This project will address these gaps by combining advanced machine learning techniques with clustering to offer a more nuanced approach to employee churn prediction. By exploring multiple features in-depth, using advanced algorithms like Gradient Boosting, and incorporating cluster analysis, the model will provide more personalized insights into employee segments at risk. This will allow organizations to develop targeted retention strategies, addressing the limitations of current models and making the solution more applicable to large-scale, dynamic environments. Additionally, the model's performance will be optimized for accuracy and efficiency, ensuring it can handle real-time predictions and scale to different organizational needs.

CHAPTER 3

Proposed Methodology

The proposed methodology for employee churn prediction in Python involves a systematic approach to identify the factors contributing to churn and develop a predictive model to forecast which employees are likely to leave. The process begins with **Exploratory Data Analysis (EDA)**, where the dataset is examined for patterns, outliers, and relationships between factors like job satisfaction, salary, and performance that influence churn. Next, **Data Visualization** techniques such as bar charts, histograms, and heatmaps are used to visually uncover trends and correlations. Following this, **Cluster Analysis** (e.g., K-means clustering) will group employees based on similar traits, helping to pinpoint segments with higher churn risk. The core of the project involves building a **Prediction Model** using machine learning algorithms like logistic regression, Random Forest, or Gradient Boosting, trained on historical data. Finally, the model's effectiveness will be evaluated using performance metrics like accuracy, precision, recall, and ROC-AUC to ensure it reliably predicts churn and provides actionable insights for HR teams.

3.1 System Design

3.1.1 Registration: The registration component serves to collect and store essential data about employees that will be instrumental for churn analysis. It includes a user-friendly interface, such as a web or mobile application, where HR personnel can input relevant employee information, including name, age, department, job title, salary, joining date, performance ratings, and job satisfaction levels. This data is securely stored in a relational or No SQL database, ensuring proper validation rules are applied to maintain data integrity and accuracy.

Access control mechanisms, such as role-based access control (RBAC), ensure that only authorized personnel can register or modify employee data, while logging features track all registration activities for auditing purposes.

3.1.2.1 Recognition: The Recognition component focuses on analyzing the registered employee data to identify patterns that may indicate potential churn. This involves a data preprocessing module that cleans and prepares the data for analysis, followed by exploratory data analysis (EDA) to uncover insights into employee characteristics and churn trends. Clustering algorithms, like K-means or hierarchical clustering, are employed to segment employees into groups based on shared attributes, highlighting which segments may be more prone to churn. Subsequently, machine learning models—such as logistic regression, decision trees, and Random Forest—are built to predict churn likelihood based on employee data, with rigorous cross-validation to evaluate performance. Finally, the system includes dashboards and reporting tools to present churn predictions and actionable recommendations to HR teams, enabling them to proactively address churn and enhance retention strategies. This comprehensive system design integrates effective registration processes with advanced recognition mechanisms, allowing organizations to better understand and mitigate employee churn, ultimately fostering a more stable and productive workforce.

3.2 Modules Used

3.2.2 Face Detection: To perform employee churn prediction in Python, we will begin with **exploratory analysis** to gain a deeper understanding of the dataset. This includes identifying and handling missing values, detecting outliers, and performing summary statistics on key features such as salary, tenure, and job satisfaction. We'll explore correlations between these features and the target variable (churn) to uncover potential drivers of employee turnover. The goal is to develop an initial understanding of the data, which can be further refined through the application of data cleaning techniques and feature engineering.

Following this, we will proceed with **data visualization** to visually assess relationships in the data.

Using bar plots, histograms, box plots, scatter plots, and heat maps, we will analyze the distribution of variables and the correlation between employee attributes like department, promotions, or salary brackets, and churn rates.

Visualizing these patterns helps identify key factors that influence employee decisions to leave. Next, we will conduct cluster analysis (e.g., k-means clustering) to group employees into clusters based on similar characteristics such as job role, performance, and satisfaction levels. This helps identify unique patterns within employee groups and how these clusters correlate with churn, providing more context to the churn behavior.

Moving on, we will focus on **building a prediction model**, starting with data preprocessing, which involves converting categorical variables into numerical formats, scaling the data, and possibly using techniques like SMOTE to handle any class imbalances. We will experiment with various machine learning models such as Logistic Regression, Decision Trees, Random Forests, Gradient Boosting, and XG Boost to predict employee churn. After training the models, we'll evaluate their performance using metrics such as accuracy, precision, recall, F1-score, and the ROC-AUC curve. Cross-validation will be applied to ensure that the model is robust and generalizes well to new data. Additionally, feature importance analysis will help us understand which factors contribute most to churn, offering actionable insights for business decisions.

3.3 Data Flow Diagram

Employee churn refers to the departure of employees from an organization, which can be costly and disruptive, as it involves not only losing talent but also the time and resources required to hire and train replacements. Just like customer churn, where retaining existing customers is often more cost-effective than acquiring new ones, retaining employees is critical for businesses.

In contrast to customer churn, however, businesses can choose which employees to hire, and these employees are integral to the company's operations and reputation.

To predict employee churn, the process will involve several key steps: **Exploratory Analysis** to understand the dataset and identify key factors influencing churn, **Data Visualization** to illustrate patterns and relationships between features like salary, job satisfaction, and churn, **Cluster Analysis** to group employees with similar characteristics.

Building a Prediction Model using machine learning algorithms like Logistic Regression or Random Forests, and finally, **Evaluating Model Performance** using metrics such as accuracy, precision, recall, and ROC-AUC to ensure the model effectively predicts churn. A **Data Flow Diagram (DFD)** would illustrate the flow of information throughout these steps, from data collection and preprocessing to model training, testing, and evaluation.

Block Diagram Description

1. Start

- Open the google colab.
- Initiates the process.

2. Load Data

- Load the dataset (e.g., employee_churn.csv).

3. Exploratory Analysis

- Check for missing values.
- Basic statistics overview.
- Distribution of the target variable (Churn).

4. Data Visualization

- Visualize numerical features against churn.
- Generate a correlation heat map.

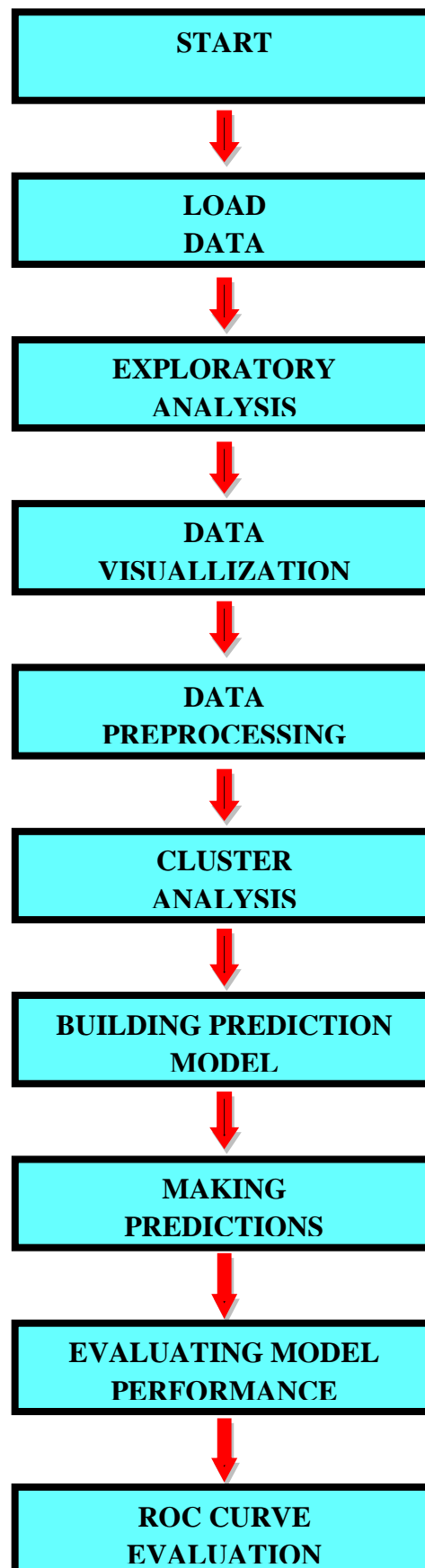


Fig 1: Data Flow Diagram

5. Data Preprocessing

- One-hot encoding for categorical variables.
- Define features (X) and target (y).
- Split the dataset into training and testing sets.
- Scale the features using Standard Scaler.

6. Cluster Analysis

- Determine optimal clusters using the Elbow method.
- Fit KMeans clustering and assign clusters.

7. Building Prediction Model

- Create and train a Random Forest Classifier.

8. Making Predictions

- Use the trained model to predict employee churn on the test dataset.

9. Evaluating Model Performance

- Generate a classification report.
- Plot confusion matrix.
- Calculate and print accuracy.

10. ROC Curve Evaluation

- Calculate predicted probabilities.
- Compute ROC curve points.
- Plot the ROC curve and calculate AUC.

3.4 Advantages

- **Cost-Saving:** Reduces recruitment and training costs associated with replacing departing employees.

- **Stability:** Maintains a stable workforce, preserving institutional knowledge and experience.
- **Satisfaction:** Improves overall employee satisfaction by addressing issues that contribute to churn.
- **Data-Driven:** Utilizes data analytics to make informed decisions and tailor strategies based on empirical evidence.
- **Optimization:** Allocates resources efficiently by focusing on high-risk employee groups identified through analysis.
- **Performance:** Enhances organizational performance through better teamwork and continuity among retained employees.

3.5 Requirement Specification

3.5.1. Hardware Requirements:

1. **Processor (CPU):**(Minimum: Dual-core processor, Recommended: Quad-core or higher for faster data processing and model training.)
2. **Memory (RAM):**(Minimum: 8 GB, Recommended: 16 GB or more for handling larger datasets and running complex models)
3. **Storage:**(Minimum: 10 GB of available disk space, Recommended: 50 GB or more for storing datasets, models, and outputs)
4. **Graphics Processing Unit (GPU):**(Recommended: NVIDIA GPU with CUDA support for accelerated model training, especially for deep learning tasks. Google Colab provides access to GPUs and TPUs)
5. **Network:**(Reliable internet connection to access Google Colab and datasets stored online (e.g., Google Drive, Kaggle))

3.5.2 Software Requirements:

1. **Operating System:**(Google Colab runs in a cloud-based environment, so no specific operating system is required on the user's local machine. However, a web browser (Chrome, Firefox, etc.) is needed)

2. Programming Language:(Python 3.x: Ensure the project is compatible with Python 3, which is the default in Google Colab)

3. Libraries and Frameworks:

- Data Manipulation:(Pandas, NumPy)
- Data Visualization:(Matplotlib,Seaborn)
- Machine Learning:(Scikit-learn,TensorFlow or PyTorch)

4. Development Environment:(Google Colab, which provides an interactive Jupyter notebook environment with built-in support for many libraries.)

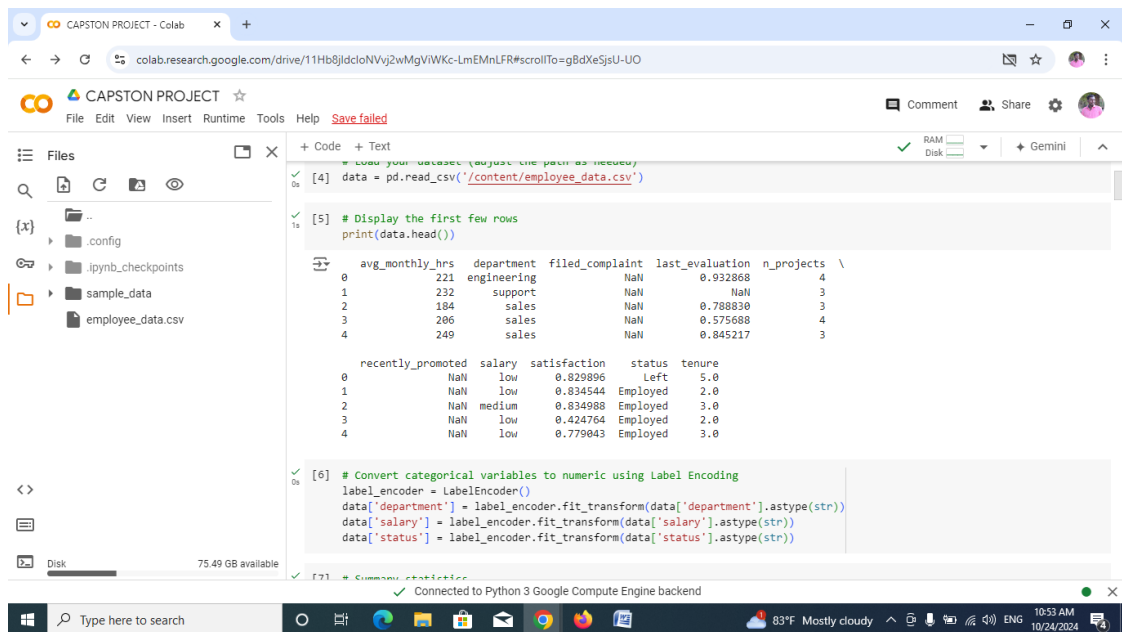
5. Version Control (Optional):(Git for tracking changes and collaborating on code (can be integrated with GitHub))

6. Data Storage Solutions:(Google Drive for storing and accessing datasets easily within Google Colab)

CHAPTER 4

Implementation and Result

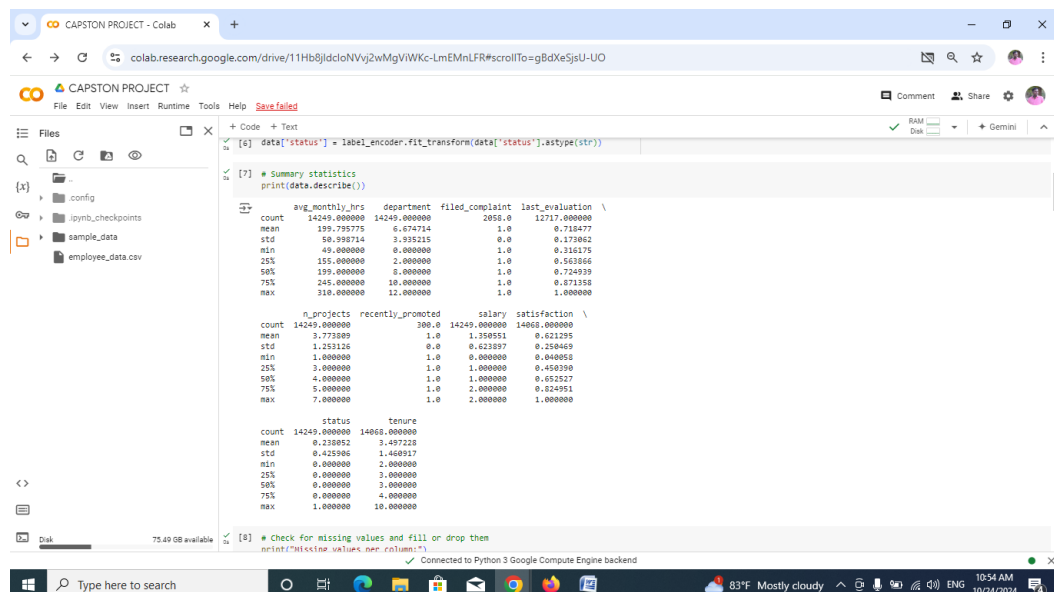
4.1 Results of Employee churn prediction



The screenshot shows a Google Colab notebook titled 'CAPSTON PROJECT'. The file explorer on the left shows a folder named 'sample_data' containing 'employee_data.csv'. The code cell [4] reads the CSV file into a DataFrame named 'data'. Cell [5] displays the first five rows of the dataset. The output shows columns: avg_monthly_hrs, department, filed_complaint, last_evaluation, n_projects, recently_promoted, salary, satisfaction, status, and tenure. The first five rows of data are as follows:

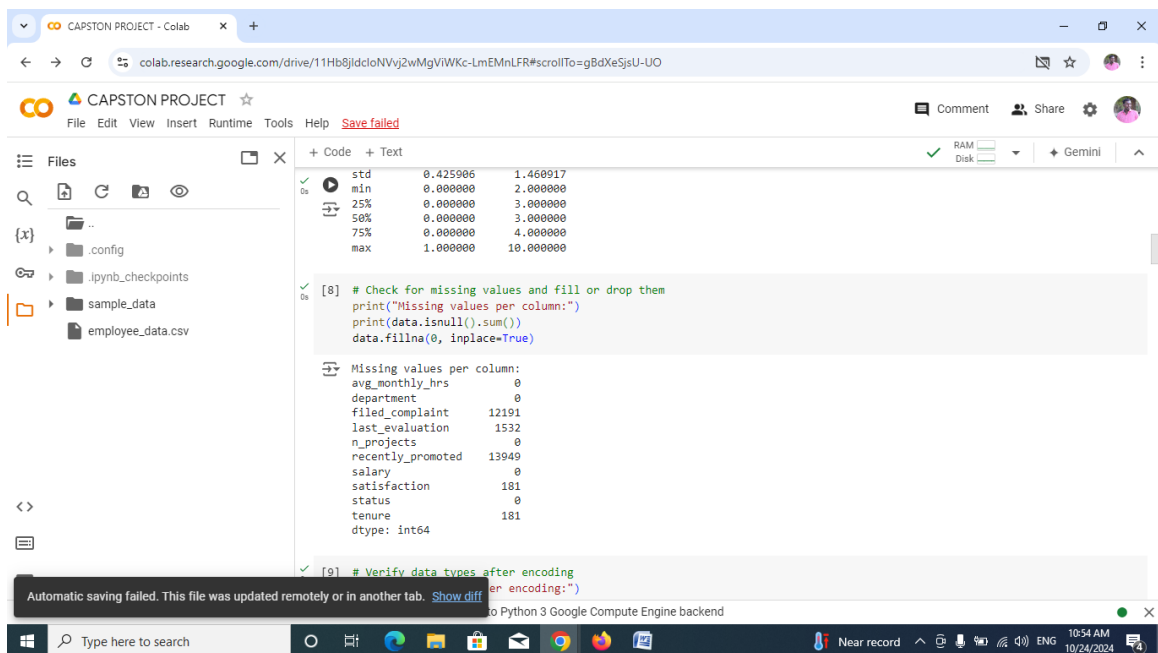
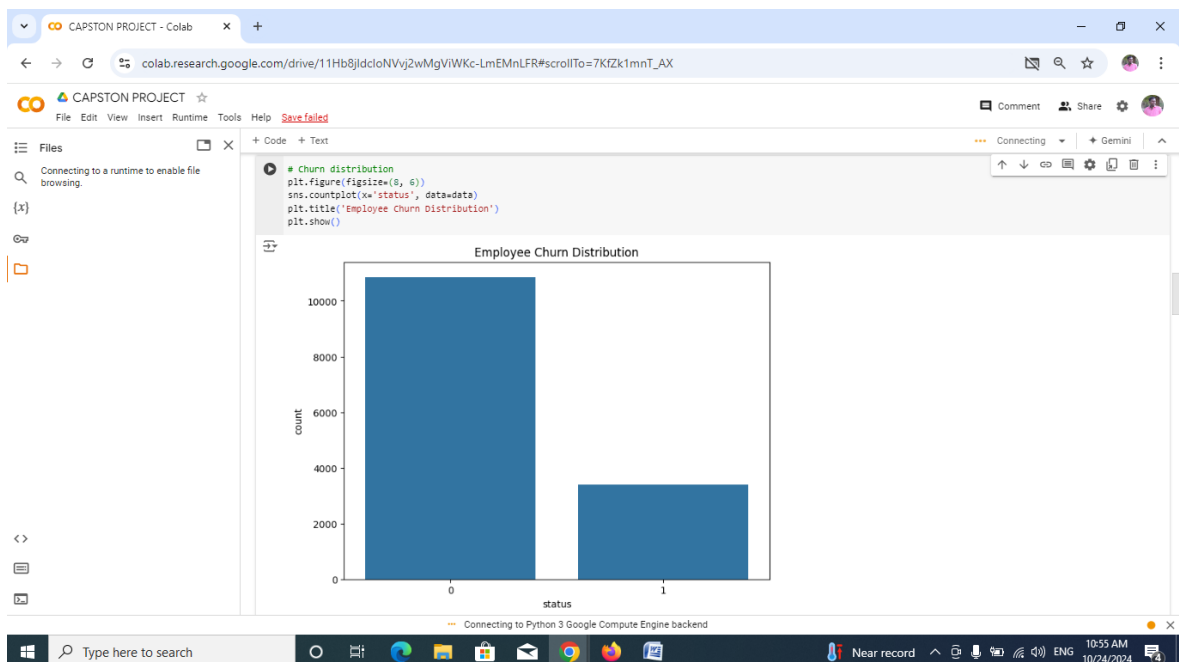
	avg_monthly_hrs	department	filed_complaint	last_evaluation	n_projects
0	221	engineering	NaN	0.932868	4
1	232	support	NaN	NaN	3
2	184	sales	NaN	0.788830	3
3	206	sales	NaN	0.575688	4
4	249	sales	NaN	0.845217	3

Fig 2. Display the first few lines.



The screenshot shows the same Google Colab notebook. Cell [6] continues the data preprocessing by converting the 'status' variable to numeric using Label Encoding. Cell [7] displays the summary statistics of the dataset using the 'describe()' method. The output shows statistics for all columns, including counts, means, standard deviations, and percentiles. The summary statistics are as follows:

	avg_monthly_hrs	department	filed_complaint	last_evaluation
count	14249.000000	14249.000000	2050.0	12717.000000
mean	199.798775	6.674714	1.0	0.718477
std	58.998714	3.935215	0.0	0.173062
min	49.000000	0.000000	1.0	0.316175
25%	155.000000	2.000000	1.0	0.503866
50%	199.000000	8.000000	1.0	0.724939
75%	245.000000	10.000000	1.0	0.871358
max	318.000000	12.000000	1.0	1.000000

Fig 3. Summary Statistics**Fig 4. Check Missing Values****Fig 5. Churn Distribution**

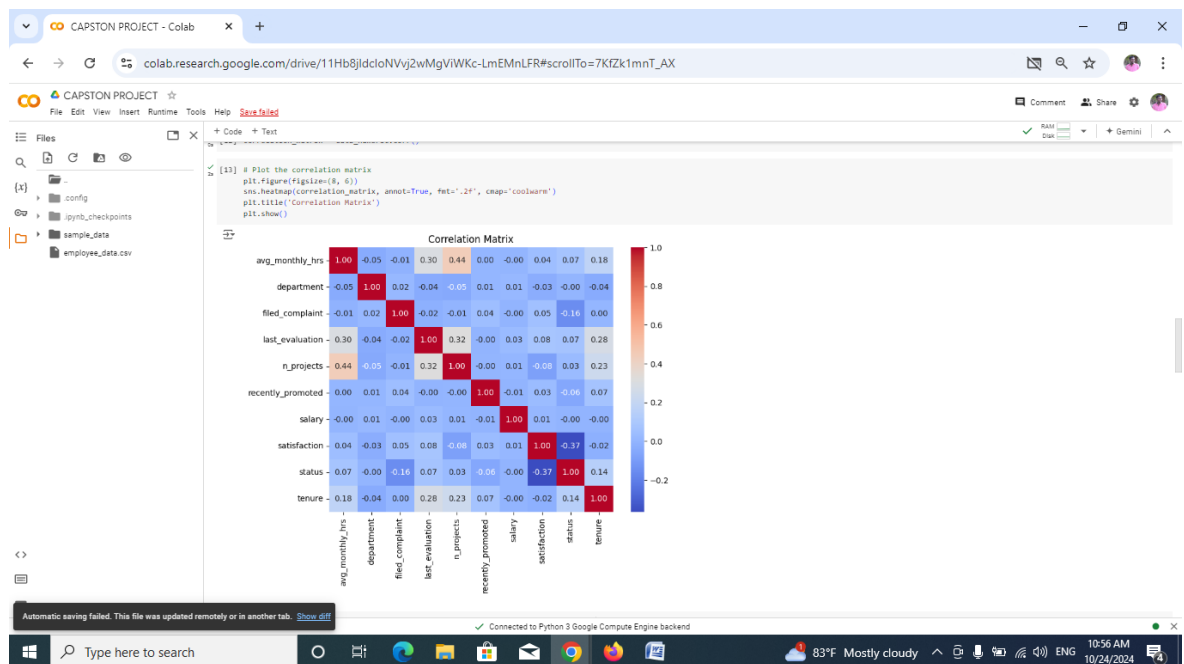


Fig 6. Correlation Matrix

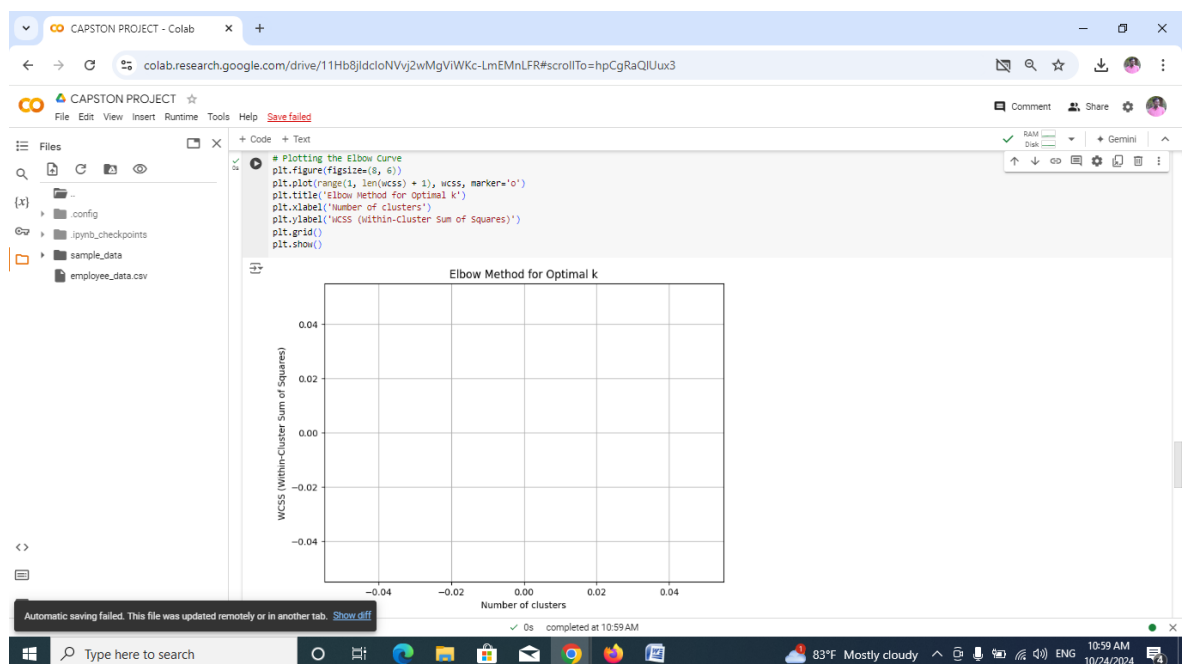


Fig 7. Elbow curve

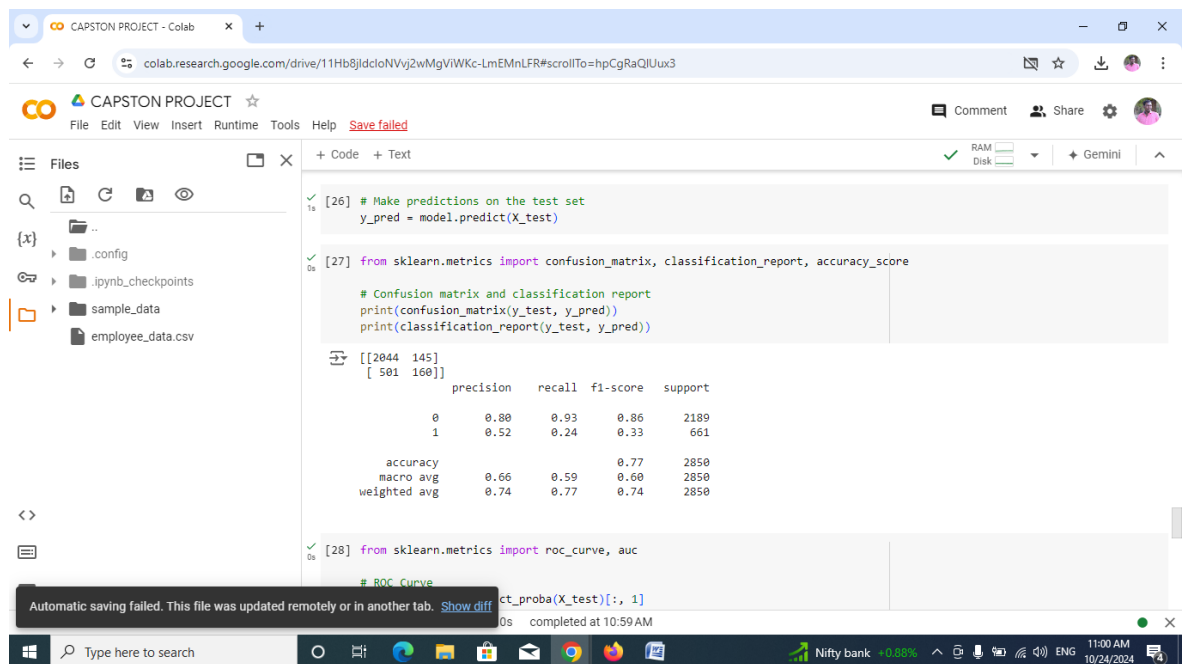


Fig 8. Confusion Matrix

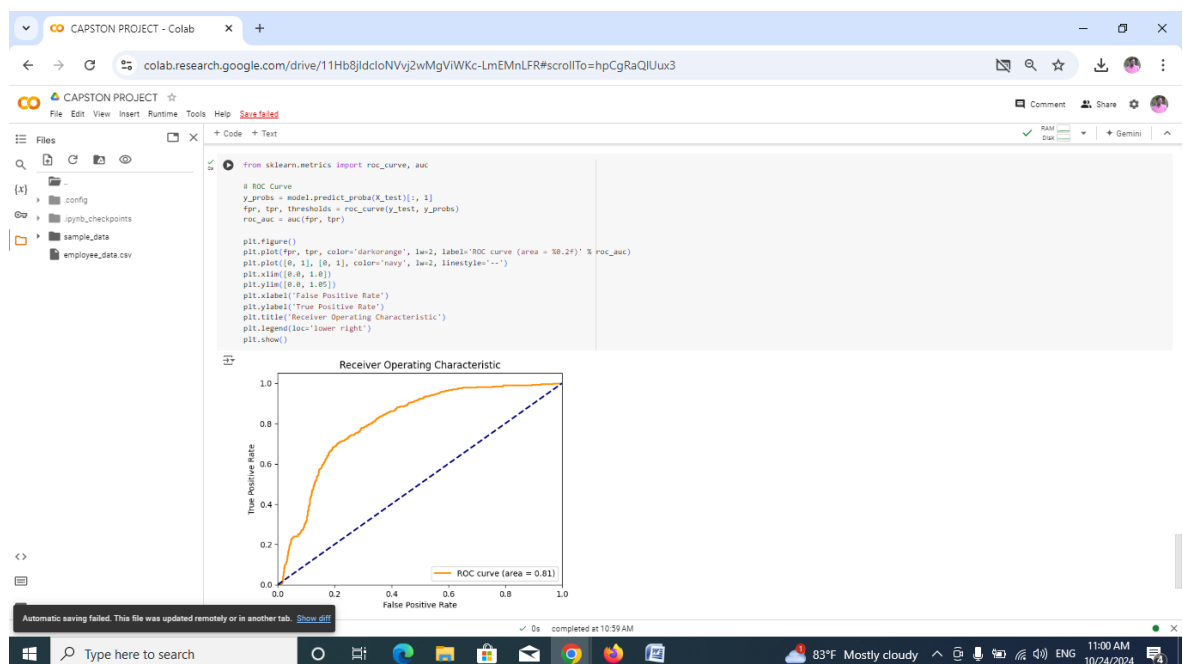


Fig 9. ROC Curve

CHAPTER 5

Discussion and Conclusion

5.1 Key Findings:

Employee churn is costly and disruptive, requiring data-driven approaches to explore, visualize, cluster, predict, and evaluate employee turnover to optimize retention strategies.

5.2 Git Hub Link of the Project:

<https://github.com/Suryaprakash3003/AI-ML-CAPSTON-PROJECT.git>

5.3 Video Recording of Project Demonstration:

<https://drive.google.com/file/d/1qu7EBkKxEmlTD9vFxrQc1CM9JLVSL4X/view?usp=drivesdk>.

5.4 Limitations:

- **Data Quality and Availability:** Incomplete, outdated, or inaccurate data can lead to misleading insights and predictions, impacting the reliability of the analysis.
- **Bias in Historical Data:** If past data reflects biases (e.g., systemic discrimination), models may perpetuate these biases, leading to unfair predictions or decision-making.
- **Dynamic Workforce Changes:** Employee behaviors and external factors (e.g., economic conditions, industry trends) can change over time, making historical data less relevant for future predictions.

- **Complex Interactions:** The reasons for employee churn can be multifaceted and influenced by numerous interconnected factors, making it challenging to model accurately.
- **Overfitting:** Predictive models may perform well on training data but fail to generalize to unseen data, especially if they are too complex or not properly validated.
- **Interpretability of Models:** Advanced machine learning models can be difficult to interpret, making it challenging to extract actionable insights for HR strategies.
- **Ethical Considerations:** Predicting churn may raise ethical concerns regarding privacy and the potential misuse of data, leading to distrust among employees.
- **Implementation Challenges:** Even with accurate predictions, implementing strategies to reduce churn requires organizational change, which can be met with resistance.

5.5 Future Work:

Future work should focus on enhancing data quality, incorporating real-time analytics, addressing model biases, and fostering cross-departmental collaboration to refine churn prediction strategies and improve employee retention.

5.6 Conclusion:

Employee churn prediction is crucial for organizations as it directly influences operational efficiency, costs, and employee morale. High turnover rates can lead to significant expenses associated with recruiting and training new staff, alongside the loss of valuable institutional knowledge that departing employees take with them. Furthermore, frequent turnover can create a destabilized work environment, impacting the morale of remaining employees and potentially exacerbating the churn problem.

Additionally, a company known for high employee turnover may struggle with its brand image, making it less appealing to potential hires and customers alike.

To effectively address employee churn, organizations can utilize various data analysis techniques. Exploratory analysis helps identify patterns and trends in employee behavior, while data visualization tools like Matplotlib and Seaborn provide intuitive insights into these findings. Implementing cluster analysis can reveal at-risk employee groups, allowing for targeted retention strategies. Predictive modeling using algorithms such as logistic regression or decision trees helps forecast churn based on numerous factors, including demographics and job satisfaction. Continuous evaluation of model performance ensures reliability, and future directions, such as integrating real-time data and personalizing retention initiatives, can further enhance these efforts, leading to a more engaged and stable workforce.

REFERENCES

- Price, J. L. (1977). *The study of turnover*. A classic work that examines organizational and individual factors affecting turnover.
- Shumaila, S. (2020). *Predicting Employee Turnover Using Logistic Regression: A Case Study*. This study explores logistic regression as a tool for churn prediction with moderate accuracy.
- Tiwari, A., & Dubey, R. (2018). *Predicting Employee Attrition Using Decision Tree Algorithm*. This paper demonstrated the use of decision trees with good accuracy.
- Singh, P., & Sharma, A. (2021). *Employee Attrition Prediction Using Random Forest Model*. A study focusing on using random forests to identify employees at risk of leaving.
- Gupta, M., & Mittal, A. (2019). *Employee Attrition Prediction Using Machine Learning Models*. This paper provides an extensive analysis of the most important predictors of churn.
- Wang, Z., & Wei, L. (2021). *Predicting Employee Turnover with Neural Networks: A Deep Learning Approach*. This paper discusses the advantages of using deep learning techniques for improving prediction accuracy.