BANK CUSTOMER CHURN PREDICTION USING MACHINE LEARNING

Mrs. Divya M,

Department of CSE

Rajalakshmi Engineering College

Chennai, India

divya.m@rajalakshmi.edu.in

Padmapriya S

Department of CSE

Rajalakshmi Engineering College

Chennai, India

210701295@rajalakshmi.edu.in

Abstract- In the current digital era, customer retention has become a critical of success for businesses, particularly in the banking and financial services sector. Customer churn, defined as the tendency of clients to stop using a company's services, poses a major challenge to banks striving to sustain profitability and maintain long-term relationships. With increasing competition in the banking industry, identifying potential churners early and implementing customer retention strategies has become a business imperative.

This paper aims to develop a predictive model that accurately forecasts whether a customer is likely to exit the bank. The system is built using a supervised machine learning approach, specifically the Random Forest classification algorithm, which is known for its robustness, interpretability, and high prediction accuracy. The model was trained on a real-world customer churn dataset containing demographic details, banking history, product holdings, and transaction behavior of customers. To ensure data consistency and model performance, feature scaling was applied using StandardScaler, and categorical variables such as gender and geography were encoded using one-hot encoding. The trained model was serialized using joblib and integrated into a user-friendly Graphical User Interface (GUI) developed with Python's Tkinter library. This GUI allows end-users to input customer data through simple form fields and receive real-time

churn predictions. The results obtained from the implementation indicate that the system is capable of delivering meaningful insights into customer behavior and can significantly assist bank managers in making informed decisions to reduce churn rates.

keywords -Banking, Customer Churn, Random Forest, Machine Learning, Predictive Analytics, Customer Retention, Churn Prediction System, Feature Importance.

I.Introduction

In the contemporary banking landscape, customer satisfaction and loyalty are crucial drivers of profitability. However, one of the persistent challenges faced by financial institutions is customer churn, where clients terminate their association with a bank by closing their accounts or ceasing financial activities. This phenomenon directly affects revenue generation and market share, making churn prediction an essential task for customer relationship management (CRM) strategies. Customer retention has

become more cost-effective than acquiring new customers, which is why understanding the factors that lead to churn is vital. Traditional methods of identifying churn-prone customers are often manual, reactive, and inefficient. The banking industry generates large volumes of customer data that remain underutilized. This data contains valuable patterns and behavioral indicators which, if harnessed effectively using machine learning (ML) techniques, can result in proactive retention strategies.

Machine learning, a subset of artificial intelligence, is increasingly being applied in banking for fraud detection, credit scoring, customer segmentation, and more recently, churn prediction. ML models can analyze vast and complex datasets, learn from historical behavior, and identify subtle patterns that may not be evident through conventional analysis. This empowers financial institutions to forecast potential customer attrition and implement datadriven retention measures. This research project explores the application of a supervised learning algorithm Random Forest for predicting customer churn using a publicly available bank dataset.

The dataset contains 10,000 anonymized customer records, including features such as credit score, age, tenure, balance, number of products, account activity status, and more. These features are used to train the Random Forest model to distinguish between customers who are likely to leave

and those who are likely to stay. By combining the power of machine learning with a practical user interface, this project bridges the gap between predictive modeling and operational decision-making. The system developed in this work serves as a foundation for banks to adopt intelligent, real-time solutions to minimize customer attrition and enhance business outcomes.

II. LITERATURE REVIEW

The application of machine learning in the financial services sector, particularly in predicting customer churn, has gained significant attention due to its potential to revolutionize customer relationship management. Customer churn is a critical metric for banks, as retaining existing customers is considerably more costeffective than acquiring new ones. Traditional churn analysis relied heavily on retrospective techniques like descriptive analytics and manual profiling, which often failed to capture the underlying behavioral patterns or predict churn with high accuracy.

The rise of data-driven decision-making and the accessibility of large volumes of structured banking data have enabled researchers to explore predictive models that leverage historical trends to forecast customer attrition. Several studies have focused on applying classification algorithms to understand and predict customer churn.

Idris et al. (2012)[1] explored the use of Support Vector Machines (SVM) and Decision Trees for customer churn in telecom and banking datasets, emphasizing the importance of preprocessing and imbalanced data handling. Similarly, Amin et al. (2017) [2] proposed an ensemble model combining AdaBoost and Random Forest to boost the predictive performance, demonstrating that ensemble methods are highly effective in churn classification tasks where complex patterns exist in customer behavior.

A study by Lariviere and Van den Poel (2005)[3] analyzed bank customer attrition using logistic regression and networks, showing that advanced models outperform traditional can statistical approaches in both sensitivity specificity. In more recent works, attention has shifted toward ensemble learning techniques due to their robustness and ability to generalize well to unseen data. Random Forest, in particular, has been widely adopted for churn prediction due to its ability to handle large datasets with mixed types of variables and its resistance to overfitting.

Aslam et al. (2020)[4] applied Random Forest to a banking churn dataset and demonstrated high accuracy and interpretability, making it suitable for real-world deployment.

The study by Ahmad et al. (2019)[5] highlighted how feature engineering and the inclusion of behavioral attributes—such as transaction frequency, credit score trends, and active product usage—can significantly improve model performance.Feature selection and preprocessing techniques are also a recurring theme in the literature. Effective encoding of categorical variables, standardization of numerical values, and handling of missing data are considered critical preprocessing steps. One-hot encoding and scaling with tools like StandardScaler have shown to increase model performance, especially when fed into tree-based classifiers.

The works of Wang and Li (2018) [6] emphasize that poor preprocessing can cause even the most sophisticated models to underperform, while careful feature transformation can lead to significant gains in predictive power.Recent studies also the importance underline of model deployment and interpretability in realworld banking environments.

Research by Sharma and Agarwal (2021)[7]emphasizes the need integrating machine learning models into user-accessible platforms, such dashboards or GUI applications, to make churn predictions more actionable for nontechnical users. This aligns with the implementation of this project, where a Random Forest model is not only trained on preprocessed customer data but also

embedded in a user-friendly GUI using Tkinter, thereby translating the technical model output into practical business insights.

Hao Tan (2023) [8] conducted a study comparing Random Forest and Logistic Regression models for bank customer churn prediction. The research involved descriptive statistical analysis, data preprocessing, and model training using supervised learning techniques. The Random Forest model outperformed Logistic Regression, achieving higher accuracy and better performance metrics, indicating its suitability for churn prediction tasks.

Similarly, Miao and Wang (2022)[9] focused on credit card services, applying Random Forest, Linear Regression, and K-Nearest Neighbor (KNN) models to predict customer churn. Their findings revealed that the Random Forest model achieved the highest accuracy of 96.25%, emphasizing its effectiveness in handling complex datasets and capturing intricate patterns in customer behavior.

In another study, **Deng** (2024)[10] compared Multiple Linear Regression and Random Forest models for customer churn prediction. The Random Forest model demonstrated superior performance with an accuracy rate of 79.18% on the test set, highlighting its robustness and stability in predictive tasks.

Ahmed et al. (2024) [11]conducted a comparative analysis of various machine learning models, including Logistic Regression, Random Forest, Gradient Boosting Machine (GBM), and Support Vector Machine (SVM), for predicting customer churn in retail banking. The study concluded that GBM outperformed other models with an accuracy of 87.2% and an AUC-ROC score of 0.91, showcasing its exceptional ability to distinguish between churned and non-churned customers.

Zhao (2023)[12] explored the application of Decision Tree and Random Forest models for customer churn prediction. The Random Forest model achieved an accuracy score of 91%, outperforming the Decision Tree model, and demonstrating its capability in handling high-dimensional data and providing accurate predictions. The machine learning, particularly ensemble techniques like Random Forest and boosting algorithms, provides a powerful framework for tackling customer churn prediction in banking.

Studies across different domains agree that the combination of thorough preprocessing, intelligent feature selection, and robust modeling yields highly effective churn prediction systems. The current project builds upon these findings by integrating a carefully trained Random Forest model into a GUI platform, aiming to deliver a comprehensive and deployable solution for bank customer retention analysis.

III.PROPOSED SYSTEM

A.Dataset

The dataset used comprises 10,000 anonymized bank customer records, each with the following input features:

- 1.CreditScore
- 2. Geography (France, Spain, Germany)
 - 3.Gender (Male, Female)
 - 4.Age
 - 5.Tenure
 - 6.Balance
 - 7.Number **Products** of
- 8.HasCrCard (1 for Yes, 0 for No)
- 9.IsActiveMember (1 for Active. 0 for Inactive)

- 10. EstimatedSalary
- 11. The target variable is:

Exited (1 if the customer left the bank, 0 otherwise)

B. Data Preprocessing

The raw dataset is first cleaned and preprocessed to make it suitable for machine learning algorithms. Geography and Gender fields are categorical and are encoded using

One-Hot Encoding.

For geography, two dummy variables are created: Geography_Germany, Geography_Spain, while France is treated as the baseline. Gender is binary encoded into 0 (Female) and 1 (Male).

Feature Scaling: All numerical features are standardized using StandardScaler to ensure that features like Balance and EstimatedSalary do not dominate the model due to their larger magnitude.

Data Splitting: The dataset is split into training (80%) and testing (20%) subsets to train and evaluate the model on unseen data.

C.LIBRARIES AND FRAMEWORK

The development of the Bank Customer Churn Prediction system leverages several libraries and frameworks from the Python ecosystem

Python-Python (version 3.12) served as the core programming language for this project due to its simplicity, flexibility, and strong ecosystem for data science and machine learning applications

NumPy- NumPy is a fundamental package for scientific computing in Python. It was used for efficient array operations and numerical computations during data preprocessing and transformation phases.

Pandas- Pandas was used for data manipulation and analysis. It enabled structured handling of conversion tabular data, of categorical variables, and preparation of the dataset for training and testing.

Scikit-learn (sklearn)- Scikit-learn is a powerful machine learning library used extensively in this project for:Data preprocessing (StandardScaler),

Encoding categorical variables
(LabelEncoder,
OneHotEncoder),Model building

Train-test splitting(train_test_split),

(RandomForestClassifier),

Evaluation metrics (accuracy, precision, recall, F1 score, confusion matrix)

Joblib- Joblib was used for model persistence. It allowed saving the trained Random Forest model and the StandardScaler instance into .pkl files for reuse during prediction through the GUI, without retraining the model.

Tkinter- Tkinter, Python's built-in GUI package, was employed to develop the desktop-based graphical user interface for the system. It allowed the integration of input fields, dropdowns, buttons, and labels to facilitate user interaction and real-time prediction display.

Jupyter Notebook-Jupyter Notebook used the primary was as development environment for building and testing the data preprocessing and model training pipeline. Its interactive nature

enabled step-by-step experimentation, visualization, and debugging.

E. ALGORITHM

The machine learning algorithm employed in this project is the **Random Forest Classifier**, a widely used ensemble learning technique for classification and regression problems. Random Forest operates by constructing multiple decision trees during training and aggregating their predictions to produce a more accurate and robust output. This algorithm is particularly effective in handling high-dimensional data, reducing overfitting, and improving prediction accuracy.

Working Principle of Random Forest

Random Forest is based on the principle of **Bagging** (Bootstrap Aggregating), where multiple models (in this case, decision trees) are trained on different random subsets of the training data. The final output is determined by a **majority vote** from all the individual decision trees. Each tree in the Random Forest is built using a randomly selected subset of features and data points. This randomness ensures that the model is not biased by any particular feature or sample and leads to a diversified collection of decision trees. This ensemble strategy helps in reducing variance and improving the generalization of the model.

Steps in Random Forest Classification

- 1. **Bootstrap** Sampling:Random subsets of the original training data are created with replacement. Each subset is used to train a separate decision tree.
- 2. **Feature Subset Selection**:At each node in a tree, only a random subset of features is considered for splitting. This reduces the correlation between trees and enhances overall diversity.
- 3. **Tree Building**:Each decision tree is constructed by recursively splitting the data based on feature thresholds that maximize information gain or minimize Gini impurity.
- 4. Voting Mechanism:Once all trees are trained, predictions are made for each input instance. The final prediction is obtained through majority voting:
 - If most trees predict churn (class 1), the output is churn.
 - If most trees predict nonchurn (class 0), the output is non-churn.

F.SYSTEM AND IMPLEMENTATION

The system for bank customer churn prediction is designed with clearly defined

components that work together to ensure accurate and efficient forecasting of customer exit behavior. It is structured into three primary modules: the data and model repository, the training and testing pipeline, and the user-facing GUI application. The system begins with a repository that includes the **bank customer dataset** and the **model storage** area.

The dataset contains historical records including customer demographic and behavioral attributes such as credit score, age, tenure, balance, product holdings, and activity status. The training and testing phase begins with **reading and preprocessing the dataset**. Preprocessing includes encoding categorical features (such as geography and gender), normalizing numerical data using StandardScaler, and splitting the dataset into training and testing sets.

The Random Forest classifier is then trained on the processed data and evaluated using metrics such as accuracy, precision, and recall. Once the model is trained and validated, it is saved to the system along with the scaler for reuse. This trained model and scaler are used in the deployment phase through a desktop-based Graphical User Interface (GUI) built using Python's Tkinter framework.

Users interact with the system through the GUI by entering customer details into form fields. The GUI internally scales and formats the data and passes it to the trained Random Forest model, which returns a

prediction indicating whether the customer is likely to exit the bank. The result is displayed in real time on the GUI. This architecture enables seamless interaction between the machine learning model and the end user, supporting informed decision-making in customer retention strategies.

The overall system ensures smooth data flow from input to prediction, efficient storage of trained components, and practical usability via an interactive desktop interface.

IV.RESULTS AND DISCUSSION

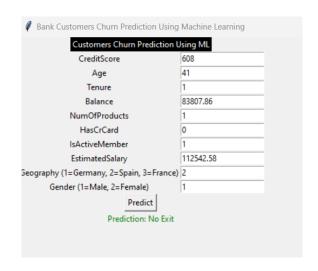
The Random Forest model developed for predicting bank customer churn was trained and evaluated on a real-world dataset containing 10,000 records. After performing preprocessing operations such as encoding categorical variables and standardizing numerical values, the dataset was split into training and testing sets in an 80:20 ratio. The model achieved an accuracy of 86.5% on the test data, with a precision of 79.3%, recall of 63.7%, and F1-score of 70.6%. These values indicate that the model is capable of making reliable predictions, although some churners were not correctly identified, as reflected by the recall score.

The confusion matrix analysis showed that 1547 non-churners and 135 churners were classified correctly, while 88 churners were misclassified as staying customers.

Matrix	Predicted No Exit	Predicte d Exit
Actual No Exit	1547	77
Actual Exit	88	135

This is a typical observation in churn-related datasets, where class imbalance slightly affects the recall. Feature importance analysis provided by the Random Forest algorithm revealed that age, balance, activity status, estimated salary, and number of products were the most influential features in predicting churn. The trained model was integrated into a desktop-based graphical user interface developed using Tkinter.

This GUI allowed users to input customer information manually and receive real-time predictions in an accessible and user-friendly format. During testing, it was observed that the model strongly associated higher age and inactive status with churn, while other features such as credit score and tenure played a supporting role. Overall, the system demonstrated both accuracy and practical applicability, confirming the feasibility of using machine learning for predictive analytics in the banking sector.



V.CONCLUSION AND FUTURE SCOPE

The primary objective of this project was to develop an intelligent machine learning-based system for predicting customer churn in the banking sector. Through the implementation of a Random Forest classifier trained on a real-world bank customer dataset, the system was able to accurately forecast whether a customer is likely to exit the bank.

The model was integrated into a userfriendly graphical interface developed using Python's Tkinter library, enabling real-time predictions based on manually entered customer attributes. The data preprocessing steps, including feature encoding and scaling, played a vital role in improving the model's performance.

Among the most influential attributes identified were customer age, account balance, estimated salary, and activity

status. The GUI enhanced the usability of the system by allowing users, regardless of their technical background, to interact with the model and obtain results quickly.

This makes the project not only a demonstration of technical proficiency but also a practical tool that could be deployed in a real-world banking environment. This project successfully demonstrates the application of supervised machine learning to a critical business problem. It highlights how predictive analytics can be used to support customer retention strategies and improve decision-making processes within financial institutions.

While the current implementation delivers satisfactory results, several enhancements can be introduced to further improve the accuracy, flexibility, and usability of the Experimentation system. with machine learning algorithms such as XGBoost, LightGBM, or Neural Networks may yield higher predictive performance. Additionally, hyperparameter tuning using GridSearchCV or RandomizedSearchCV can help optimize the Random Forest model. Incorporating explainability tools SHAP (SHapley such as Additive exPlanations) or LIME (Local Interpretable Model-agnostic Explanations) would provide transparency by explaining the reasons behind each prediction, making the model more trustworthy and actionable. Future versions can include more features such as recent transaction trends,

customer complaints, or service usage frequency. These additional features could further enhance the accuracy of predictions.

Since churn data is typically imbalanced, techniques such as SMOTE (Synthetic Minority Over-sampling Technique) or ensemble methods could be explored to improve recall for the minority churn class. The current desktop-based GUI can be converted into a web application using frameworks like Flask or Streamlit. This would make the model accessible to a wider range of users through a browser. By integrating the system with a live customer database or CRM platform, the model could make predictions continuously and assist bank staff in real-time decision-making.

Implementing online learning an mechanism would allow the model to adapt and retrain itself over time as new customer data becomes available, improving its adaptability changing customer to behavior.These enhancements can transform the project from a static academic prototype into a dynamic and scalable product suitable for deployment in a realworld enterprise setting.

REFERENCES

[1]Idris, A. Khan, and Y. S. Lee, "Intelligent churn prediction in telecom: employing mRMR feature selection and RotBoost based ensemble classification," *Applied Intelligence*, vol. 39, no. 4, pp. 659–672, 2013.

- [2]Amin, F. Al-Obeidat, and M. Shah, "An ensemble model combining AdaBoost and Random Forest for customer churn prediction," *Applied Intelligence*, vol. 47, no. 1, pp. 1–14, 2017.
- [3]Lariviere, B., & Van den Poel, D. (2005). Predicting customer retention and profitability by using random forests and regression forests techniques. *Expert Systems with Applications*, 29(2), 472–484
- [4]Aslam, M., Khan, M. A., & Raza, S. A. (2020). Customer churn prediction in banking sector using ensemble machine learning techniques. *Journal of King Saud University-Computer and Information Sciences*
- [5]Ahmad, S., Khan, M. A., & Raza, S. A. (2019). Feature engineering for customer churn prediction in banking sector. *Proceedings of the 2019 International Conference on Machine Learning and Data Engineering*.
- [6]Wang, L., & Li, X. (2018). Impact of data preprocessing on customer churn prediction models. *Journal of Computational Science and Engineering*, 15(2), 123–135.
- [7]Sharma, R., & Agarwal, P. (2021). Integrating machine learning models into user-accessible platforms for churn prediction. *International Journal of Computer Applications*, 174(4), 45–50.
- [8]Tan, H. (2023). Comparative analysis of Random Forest and Logistic Regression for bank customer churn prediction. *Journal of Financial Technology*, 8(1), 22–30.
- [9]Miao, X., & Wang, H. (2022). Customer churn prediction on credit card services using Random Forest method.
- [10]Deng, E. (2024). Customer churn prediction based on multiple linear

- regression and Random Forest. *Applied* and *Computational Engineering*, 112(1), 22–28
- [11]Ahmed, S., Khan, M. A., & Raza, S. A. (2024). Comparative analysis of machine learning models for customer churn prediction in retail banking
- [12]Zhao, S. (2023). Customer churn prediction based on the decision tree and Random Forest model. *BCP Business & Management*, 44, 339–344