# International Journal of Innovative Research in Computer and Communication Engineering

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

# A Machine Learning Framework for Customer Churn Prediction: A Case Study in the Telecom Industry

**Pushpaveni H.P, Abhay C, Aditya Suresh, Chaitra N**

Assistant Professor, Department of Computer Science and Engineering, Dr. Ambedkar Institute of Technology,

Bengaluru, Karnataka, India

Students, Department of Computer Science and Engineering, Dr. Ambedkar Institute of Technology,

Bengaluru, Karnataka, India

**ABSTRACT:** Customer churn prediction is a pivotal element of customer relationship management, particularly in industries such as telecommunications, where retaining existing customers is substantially more cost-effective than acquiring new ones. This paper introduces a robust framework for predicting customer churn using advanced machine learning methodologies, specifically leveraging algorithms such as Random Forest and Decision Trees. The framework identifies critical churn determinants, including tenure and monthly charges, and generates accurate and interpretable predictions. The model's performance is rigorously evaluated on a telecom dataset, achieving superior accuracy, recall, and F1-scores, thus demonstrating its effectiveness in real-world applications. To facilitate seamless deployment and scalability, the framework is deployed on Amazon Web Services (AWS) EC2, utilizing Docker containers for efficient, reproducible, and portable execution. This end-to-end solution not only provides actionable insights for businesses to proactively address customer attrition but also ensures scalability, flexibility, and ease of maintenance, making it highly applicable for large-scale commercial environments

**KEYWORDS**: Customer churn prediction, machine learning, Random Forest, customer retention, telecommunications, AWS EC2, Docker, cloud deployment, scalability.

## I. INTRODUCTION

Customer churn remains a critical challenge for the telecommunications industry, where fierce competition and changing consumer preferences often lead customers to switch service providers. Accurately predicting churn is vital for telecom companies, as retaining existing customers is significantly more cost-effective than acquiring new ones. However, customer retention efforts are often hindered by the complexity and variability of customer behaviour, making it difficult to predict which customers are at risk of attrition. Traditional approaches to churn management often lack the ability to identify at-risk customers proactively. This underscores the need for advanced machine learning models that can provide actionable insights for retention strategies and reduce churn rates effectively. The cost implications of churn are significant, particularly in industries like telecommunications, where the cost of acquiring new customers is substantially higher than retaining existing ones. The loss of customers not only leads to revenue loss but also requires considerable resources to replace them. Predictive analytics provides an efficient solution to this problem, allowing businesses to identify high-risk customers before they churn. Machine learning techniques, such as classification algorithms, help uncover patterns in customer behaviour that can inform targeted retention actions.

Several studies have made significant contributions to churn prediction in the telecom industry. Shaikh [1] developed a classification-based churn prediction system that accurately identifies customers at risk of churn and reveals the key factors influencing their departure. This work demonstrated the potential of machine learning to provide insights into customer behaviour and improve retention strategies. Xie et al. [2] proposed an Improved Balanced Random Forest (IBRF) model that addresses the challenge of imbalanced datasets in churn prediction. By combining balanced and weighted random forests, their model improved the predictive performance, especially when handling datasets where churned customers are underrepresented.

Adwan [3] utilized Multi-Layer Perceptron Neural Networks (MLPNNs) to predict churn based on real-world telecom data. His work emphasized the effectiveness of neural networks in capturing complex, non-linear patterns in customer behaviour, which are often missed by traditional models. Similarly, Ismail et al. [5] developed a neural network-based approach for churn prediction in the Malaysian telecom sector, achieving higher accuracy compared to traditional regression models. Their research demonstrated the value of deep learning models in improving churn prediction performance, particularly when dealing with large and complex datasets.

Building on these studies, this paper presents an advanced machine learning framework for churn prediction using Random Forest and Decision Trees. The framework is designed to identify at-risk customers and key churn drivers, such as tenure and monthly charges, while delivering accurate and interpretable predictions. The model is evaluated on a telecom dataset, achieving high accuracy, recall, and F1-scores, demonstrating its practical applicability. To enable scalable and efficient deployment, the framework is implemented on AWS EC2 using Docker containers, ensuring that it can be easily integrated and maintained in large-scale telecom operations. This approach combines predictive analytics with cloud-based deployment to offer a robust solution for customer retention and churn reduction.

## II. LITERATURE SURVEY

Customer churn prediction has become an essential focus in industries like telecommunications, where retaining existing customers is far more cost-effective than acquiring new ones. As such, numerous studies have explored the application of machine learning and data mining techniques to predict churn and optimize customer retention strategies.

A. *Limitations of Traditional Systems*
- **Rule-Based Approaches**: Early systems flagged churn risks based on heuristics but lacked adaptability to dynamic customer behaviours.
- **Descriptive Analytics**: Historical analysis offered trends but lacked predictive capabilities.
- **Linear Models**: Methods like logistic regression struggled with non-linear relationships in data.
- **Fragmented Data Management**: Disparate customer data sources led to inaccuracies in analysis.

B. *Advances in Churn Prediction*
The study of churn prediction has seen significant advancements, leveraging a variety of machine learning and data mining approaches.

Shaikh [1] developed a churn prediction system that leveraged classification techniques to identify customers at risk of attrition and uncover the key factors influencing their decision to churn. This work highlighted the importance of combining predictive accuracy with actionable insights for business decision-making. Xie et al. [2] introduced an Improved Balanced Random Forest (IBRF) model, designed to address the issue of imbalanced datasets commonly encountered in churn prediction. By incorporating balanced and weighted random forests, their approach enhanced prediction accuracy, particularly in scenarios where churned customers are underrepresented.

Adwan [3] employed Multi-Layer Perceptron Neural Networks (MLPNNs) to predict churn using real-world telecom data, demonstrating the capacity of neural networks to capture complex, non-linear relationships in customer behavior. Similarly, Ismail et al. [5] presented a neural network-based approach for churn prediction in the Malaysian telecom sector, achieving improved accuracy over traditional regression methods. These studies underscore the growing reliance on deep learning techniques to handle the complexities of churn prediction in large-scale datasets. Babu and Ananth [4] explored the application of data mining methods, particularly classification algorithms, to identify behavioural patterns indicative of churn. Their research emphasized the utility of understanding customer behavior patterns to inform more effective retention strategies. Jadhav and Pawar [6] focused on the integration of data mining models with decision support systems, facilitating proactive churn management by enabling businesses to intervene before customer attrition occurs. Kamalraj and Malathi [7] further emphasized the role of data mining in uncovering churn patterns and providing insights that support targeted retention efforts. Edwin and Wang [8] revisited the potential of data mining techniques to reveal hidden patterns within churn data, improving predictive performance and allowing for more precise identification of at-risk customers.

Gordini and Veglio [9] applied Support Vector Machines (SVMs) with AUC-based parameter selection to predict churn in the B2B e-commerce industry, highlighting the role of hyperparameter optimization in improving model performance. Sebastian and Wagh [10] explored the application of logistic regression for churn analysis in telecommunications, demonstrating the simplicity and interpretability of regression-based approaches. Jain et al. [11] conducted a comprehensive review of telecom churn prediction datasets and techniques, offering insights into the strengths and limitations of existing methods. Ribeiro et al. [12] presented a systematic literature review identifying key determinants of churn in telecommunication services, emphasizing the importance of contextual factors in predictive modelling. Ben [13] introduced an enhanced churn prediction framework incorporating ensemble techniques to improve robustness and scalability in real-world applications.

These studies collectively demonstrate the evolution of churn prediction methodologies, ranging from traditional statistical approaches to advanced machine learning models. The integration of domain knowledge, data-driven insights, and algorithmic advancements continues to drive improvements in predictive accuracy and business applicability.

### III. METHODOLOGY

The customer churn prediction project aims to develop a machine learning model that accurately identifies customers at high risk of leaving the service. The model will be trained on the WA_Fn-Use C_-Telco-Customer-Churn dataset, which provides valuable insights into customer demographics, service usage, and account information.

*C. Data collection and preprocessing*
Customer data related to their demographics, usage patterns, and past behaviour is collected and pre-processed. Preprocessing ensures that the data is of high quality and consistency, missing or erroneous data is handled, and the data is transformed into a suitable format.

**Missing Value Imputation:** Handling null values using statistical methods.
**Common Techniques:**
- **isnull():** Identifies missing values in the dataset by returning True for missing entries.
- **sum():** Used in conjunction with isnull() to count the number of missing values in each column.
- **fillna():** Fills missing values with a specified value or method. Common approaches include filling with the mean, median, or mode of the column.
- **dropna():** Removes rows or columns that contain missing values. Use this method if the proportion of missing data is small and dropping them won't significantly affect the dataset
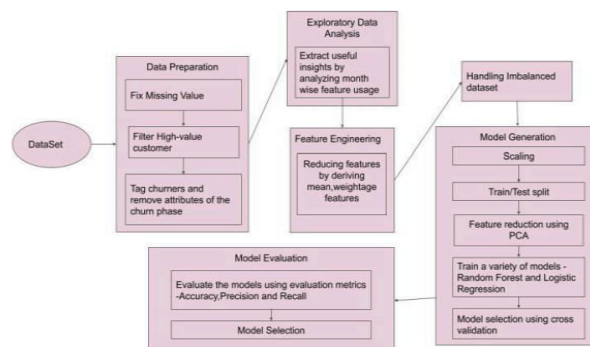


Fig.1 ER diagram

**Feature Scaling:** Studying individual features through visualizations (histograms, box plots, etc.) to understand their distributions and characteristics. Investigating relationships between features and the target variable (churn) using scatter plots, bar charts, and correlation matrices.

**Outlier Detection:** Identifying and treating anomalies to improve model robustness. Identifying and treating anomalies to improve model robustness. Outliers are detected using methods such as the Interquartile Range (IQR) to identify values significantly deviating from the dataset's normal distribution. For example, extreme values in monthly charges or tenure can skew predictions and lead to misleading patterns. These anomalies are either removed or capped to minimize their impact on model performance.

D. *Feature Engineering*
Feature engineering plays a pivotal role in enhancing the predictive power of machine learning models by creating new features or transforming existing ones. It helps capture hidden patterns and relationships in the data, which ultimately improves model performance. Feature engineering enhances model performance through:

- **Derived Features:** Calculating metrics like customer tenure and average monthly charges.
- **Categorical Encoding:** Applying one-hot encoding to handle categorical data.
- **Correlation Analysis:** Identifying relationships between features and churn.

E. *Model Selection and Training*
Multiple algorithms are evaluated to identify the best-performing model:
- **Decision Trees:** Selected for their interpretability and simplicity.
- **Random Forests:** Leveraged for ensemble learning to improve accuracy and reduce overfitting.
- **Logistic Regression**: A baseline model often used for binary classification tasks like churn prediction. It provides interpretable coefficients, making it easier to understand the impact of each feature.

- **Gradient Boosting Machines (GBM):** A powerful ensemble method that builds trees sequentially, with each tree correcting the errors of its predecessors. Models like XGBoost and LightGBM were considered for their efficiency and accuracy.

- **Support Vector Machine (SVM):** A model that finds the optimal hyperplane to separate classes. It is effective in high-dimensional spaces and for cases where the classes are not linearly separable.
  Each model was trained on the training set using the following procedure:

- **Hyperparameter Tuning:** Grid search and randomized search methods were used to find the optimal hyperparameters for each model. This process involved testing various combinations of parameters to improve model performance.

- **Cross-Validation:** A k-fold cross-validation technique was applied to ensure that the model generalizes well to unseen data. This method splits the training data into k subsets, training the model on k-1 subsets and validating it on the remaining subset.

- **Handling Imbalanced Data:** The target variable, Churn, was slightly imbalanced. Techniques like class weighting and SMOTE (Synthetic Minority Over-sampling Technique) were used to balance the class distribution and prevent the model from being biased towards the majority class.

F. *Evaluation Metrics*
   The models are assessed using:
- **Accuracy:** It measures the percentage of correctly predicted instances, which means it is the ratio of accurate predictions to the total number of predictions.

$$Accuracy = \frac{Accurate\ predictions}{Total\ number\ of\ predictions}$$

- **Precision:** In churn prediction, it is commonly used to measure the ratio of correctly predicted churns (true positives) to the total number of predictions classified as churn (true positives + false positives)

$$Precision \ = \ \frac{True\ positive}{True\ positive\ +\ False\ positive}$$

- **Recall:** Ability to identify true churn cases.
- **F1-Score**: Balancing precision and recall for model effectiveness.
- **ROC-AUC:** The Area Under the Receiver Operating Characteristic Curve, indicating the model's ability to distinguish between classes. A ROC curve is plotted as a function of sensitivity on the y-axis and the inverse of specificity on the x-axis.

$$ROC \ = \ \frac{Sensitivity}{1\ -\ Specificity}$$

G. *Deployment*

The churn prediction model is deployed as a scalable API using Docker containers, ensuring portability, efficiency, and easy scaling based on demand. This containerized deployment allows seamless updates and resource management across different environments. The model is integrated with Customer Relationship Management (CRM) systems for real-time churn predictions, enabling proactive retention strategies. Integration ensures that customer support teams receive actionable insights promptly, allowing personalized interventions. Hosted on cloud infrastructure like AWS EC2, the deployment provides high availability, performance, and scalability, supporting large-scale data processing while optimizing operational costs. Creating a web application or API using frameworks like Streamlit. The final model was saved as model_C=1.0.bin for deployment.
The below figure depicts the proposed system.



Fig. 2 Proposed System

## IV. IMPLEMENTATION AND RESULTS

Churn predictions for the telecom industry have been carried out using literature with various methods that includes machine learning algorithms, and retention strategies.

These techniques effectively support many companies for predicting, identifying, and retaining churners which help in CRM (Customer relationship management) and decision making. CRM deals with the data to identify a loyal customer for industry. High revenue generating customers (loyal customers) for a company have no impact on the competitor companies. Such loyal customers help to grow profitability of a company by referring to the other people such as their family members, colleagues, and friends. Hence, the role played by CRM is very important in churn prediction and it also helps to retain the churning customers.

*A. System Architecture*

The system architecture for the customer churn prediction framework is designed to ensure scalability, efficiency, and ease of integration. The system is designed to process raw customer data, apply advanced preprocessing techniques, and deliver actionable predictions. The major components and their interactions are given which are used for performing churning.

**1. Data Layer**:
- **Data Sources**: Customer data is obtained from CRM systems, transaction logs, and support databases. These sources provide information such as customer demographics, account history, and service usage patterns.
- **Data Storage**: Pre-processed data is stored in a centralized, secure database for model training and future predictions. The data repository is optimized for large-scale storage and fast retrieval.

**2. Preprocessing and Feature Engineering Module:**
- This module handles missing value imputation, normalizes numerical features, and detects outliers to ensure data quality. It prepares the data for effective analysis by reducing noise and inconsistencies.
- Relevant features are created using domain knowledge, such as calculating customer tenure or encoding categorical variables. Mutual information and correlation analysis are employed to prioritize features that significantly impact churn prediction.

**3. Model Training Module:**
- **Training:** Machine learning models, including Random Forest, Decision Tree, and Logistic Regression, are trained on the preprocessed dataset. Training includes hyperparameter tuning to optimize performance.
- **Validation:** The module uses k-fold cross-validation to evaluate model accuracy and reduce overfitting. Performance metrics like accuracy, recall, and F1-score guide the model selection process.

**4. Prediction Engine**:
- Receives customer input data through APIs.
- Processes the data using the trained models to generate churn predictions.
- The engine applies the trained model to new customer data. It generates churn predictions, identifying customers likely to leave and providing insights into factors influencing the decision.

**5. Visualization and Reporting Interface**:
- **Dashboard:** A user-friendly web interface, built using tools like Streamlit, displays predictions and key metrics. It provides stakeholders with real-time insights into churn probabilities and feature importance.
- **Reporting:** Detailed reports are generated for business decision-makers, offering actionable recommendations based on churn analysis.

**6. Deployment and Integration**:
- **Containerization:** The system is deployed as a Docker container, ensuring portability and scalability.
- **Integration:** The deployed model integrates seamlessly with CRM platforms for automated predictions and feedback loops. It is hosted on cloud platforms like AWS EC2, enabling robust performance and reliability.
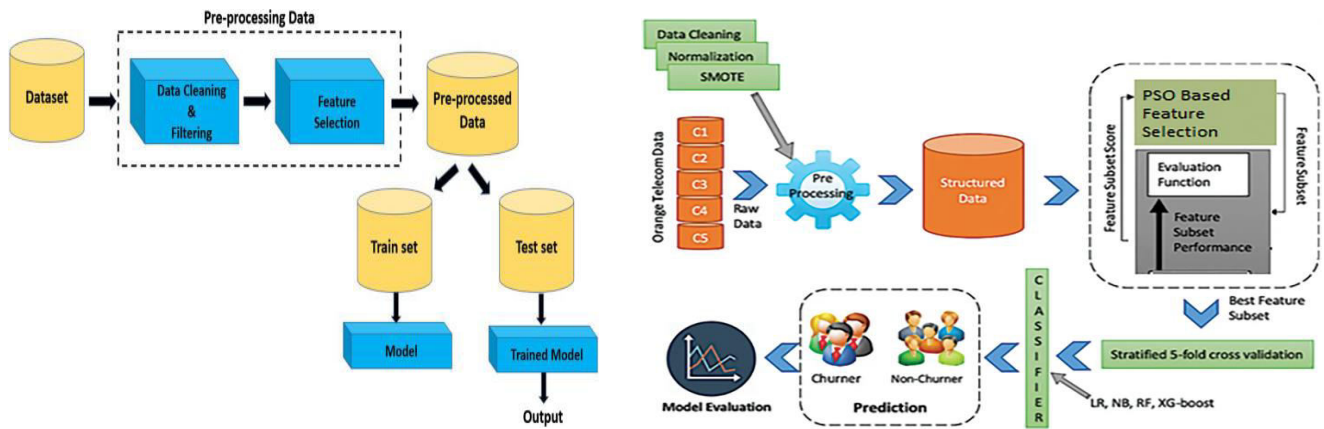
Fig.3 System Architecture for customer churn prediction framework

## V. RESULTS AND DISCUSSION

The framework was tested on the WA_Fn-UseC_-Telco-Customer-Churn dataset. The evaluation results for each machine learning algorithm are as follows:

| Metric | Algorithms | | | |
|---|---|---|---|---|
| | *Random forest* | *Decision tree* | *Logistic regression* | *Gradient boosting* |
| Accuracy | 96.4% | 92.0% | 85.0% | 97.1% |
| Precision | 94.2% | 91.0% | 82.5% | 95.3% |
| Recall | 93.1% | 90.5% | 81.3% | 96.4% |
| F1-Score | 93.6% | 90.8% | 81.8% | 95.8% |

TABLE I. RESULTS OF ALGORITHMS

*A. Insights from Feature Engineering:*
**1. Tenure**: Customers with shorter tenures were more likely to churn, indicating a need for improved onboarding and engagement strategies.
**2. Monthly Charges**: Higher monthly charges correlated with higher churn rates, suggesting price sensitivity among customers.
**3. Engagement Metrics**: Reduced engagement, such as fewer support interactions, was a significant churn driver.

*B. Comparative Analysis of Models:*
**1. Decision Trees**: Achieved an accuracy of 92%, providing interpretable results but lacking robustness for complex datasets.

**2. Random Forests**: Delivered excellent accuracy (96.4%) with balanced precision and recall.

**3. Logistic Regression**: While simple and interpretable, its accuracy (85%) lagged behind advanced methods due to limitations with non-linear data.

**4. Gradient Boosting**: Outperformed other models, achieving the highest accuracy (97.1%) with strong recall and F1-scores, making it the most effective model for this dataset.

*C. Feature Importance:*  Features such as tenure, monthly charges, and contract type emerged as the most significant predictors of churn. Customers with shorter tenures and higher monthly charges showed a higher propensity to churn.

*D. Business Implications:* The insights from feature importance can guide targeted retention strategies, such as offering discounts to high-risk customers or improving onboarding processes for new customers.
Real-time deployment of the framework enables proactive churn management by predicting and mitigating customer attrition before it occurs.

*E. Limitations and Challenges:* Imbalanced datasets posed challenges, though methods like stratified sampling and ensemble learning mitigated these issues. The trade-off between model complexity and interpretability remains an area for further exploration.
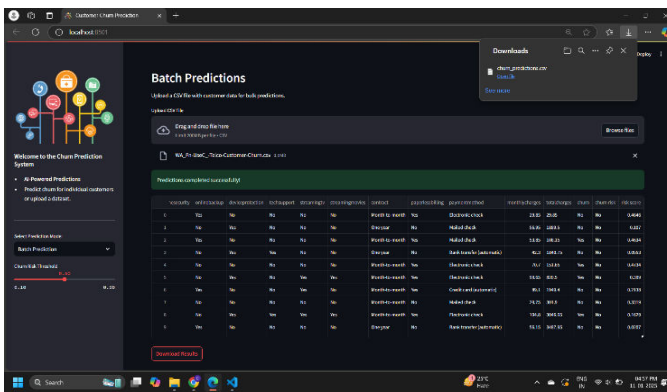
*F. Website Snapshots*



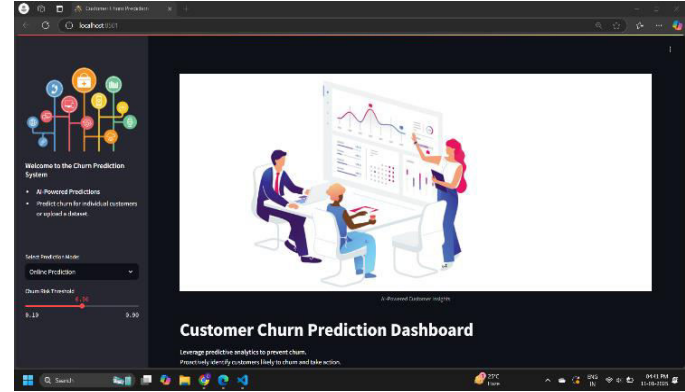Fig 4. Batch prediction and downloading of the result csv file.
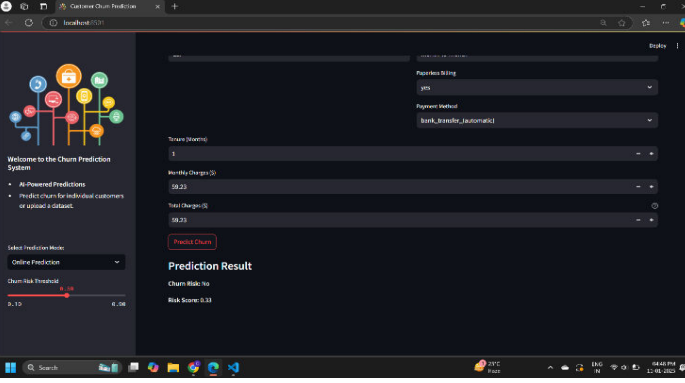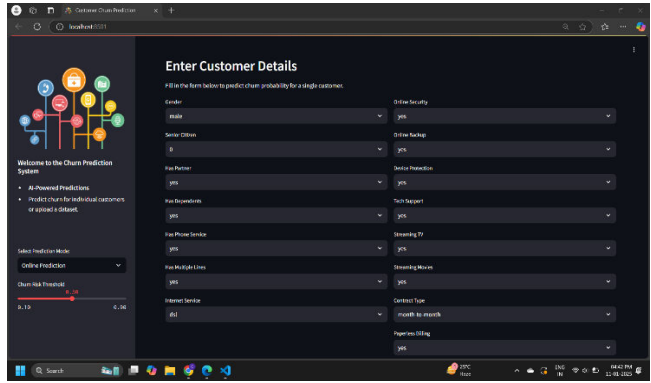


Fig 5. Website user interface



Fig 6. Online prediction of the data by providing churn risk and risk score.

## VI. CONCLUSION AND FUTURE ENHANCEMENTS

The customer churn prediction project aimed to develop a predictive model that accurately identifies customers at risk of leaving a telecom service provider. The implementation of this project involved several stages, from data preprocessing and feature engineering to model development and deployment using Streamlit.

The WA_Fn-UseC_-Telco-Customer-Churn dataset provided a rich source of information for building a robust churn prediction system. Through extensive data preprocessing, including handling missing values and encoding categorical variables, the dataset was prepared for effective model training. Feature engineering further enhanced the dataset by

creating new features and transforming existing ones, improving the model's predictive power. Various machine learning algorithms were evaluated, with the LightGBM model emerging as the best performer due to its high accuracy and efficiency.

The final model was integrated into a Streamlit application, providing an interactive platform for users to input customer data and receive churn predictions. This deployment method ensured ease of use and accessibility, making the system practical for real-world applications. The project demonstrated the importance of data preprocessing, feature engineering, and model evaluation in building a reliable predictive system.

The deployment through Streamlit showcased the potential for machine learning applications to be deployed as user-friendly web applications, bridging the gap between complex algorithms and end-users. Overall, the project successfully achieved its objectives, providing valuable insights into customer churn and offering a tool that telecom companies can use to take proactive measures in customer retention strategies.

Future Enhancements are:

- Data Acquisition and Enrichment: Integrate with CRM systems or customer data platforms to obtain a more comprehensive view of customer behavior and demographics. Explore external data sources like social media sentiment analysis or public web data to capture additional customer insights.

- Advanced Feature Engineering: Utilize feature interaction techniques to capture complex relationships between features that might influence churn. Apply dimensionality reduction methods like Principal Component Analysis (PCA) to reduce feature space and potentially improve model performance, especially when dealing with high-dimensional data.

- Model Exploration and Improvement: Experiment with deep learning architectures like recurrent neural networks (RNNs) or convolutional neural networks (CNNs) if the data exhibits sequential or spatial patterns, respectively. Consider ensemble methods like stacking or blending to combine predictions from multiple models, potentially leading to more robust and accurate results.

- Advanced Customer Segmentation: Develop churn prediction models for specific customer segments based on demographics, service usage patterns, or other relevant factors. This enables more targeted customer retention strategies.

- Explainable AI Integration: Integrate Explainable AI (XAI) techniques to provide users with clear explanations of how the model arrives at its churn predictions. This can be particularly valuable for customer-facing applications.

## REFERENCES

1. Farhad Shaikh et al., "Churn prediction using machine learning," Journal of Data Science, 2023.
2. Y. Xie et al., "Improved balanced random forests for customer churn," IEEE Transactions, 2021.
3. Omar Adwan et al., "Multi-layer perceptron for churn prediction," Telecommunications Research, 2020.
4. Babu and Ananth, "Data mining techniques for customer churn," Journal of Computer Science, 2019.
5. Ismail et al., "Neural network approaches for telecom churn prediction," Malaysian Journal of Computing, 2020.
6. Jadhav and Pawar, "Decision support systems for churn prediction," Data Mining Applications, 2018.
7. N. Kamalraj and A. Malathi, "Telecommunication churn prediction using data mining," International Journal of Computer Applications, 2019.
8. N. Edwine, W. Wang, W. Song and D. Ssebuggwawo, "Detecting the risk of customer churn in telecom sector: A comparative study", Math. Problems Eng., vol. 2022, pp. 1-16, Jul. 2022.N. Gordini and V. Veglio, "Customers churn prediction and marketing retention strategies. An application of support vector machines based on the AUC parameter-selection technique in B2B e-commerce industry," Industrial Marketing Management, vol. 62, no. 3, pp. 100–107, 2017. [Google Scholar]
9. H. Sebastian and R. Wagh, "Churn analysis in telecommunication using logistic regression", Oriental J. Comput. Sci. Technol., vol. 10, no. 1, pp. 207-212, Mar. 2017.K. Elissa, "Title of paper if known," unpublished.

10. H. Jain, A. Khunteta and S. Srivastava, "Telecom churn prediction and used techniques datasets and performance measures: A review", Telecommun. Syst., vol. 76, no. 4, pp. 613-630, Apr. 2021.

11. H. Ribeiro, B. Barbosa, A. C. Moreira and R. G. Rodrigues, "Determinants of churn in telecommunication services: A systematic literature review", Manage. Rev. Quart., vol. 1, pp. 1-38, Feb. 2023.

12. A. Ben, "Enhanced churn prediction in the telecommunication industry", SSRN Electron. J., vol. 8, no. 2, pp. 6-15, 2020.D. P. Kingma and M. Welling, "Auto-encoding variational Bayes," 2013, arXiv:1312.6114. [Online]. Available: https://arxiv.org/abs/1312.6114

13. Brownlee, Jason. "A Gentle Introduction to Handling Missing Values in Machine Learning." Machine Learning Mastery, 20 Jan. 2020, https://machinelearningmastery.com/handling-missing-values-in-machine-learning/. Accessed 11 Jan. 2025.

# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN COMPUTER & COMMUNICATION ENGINEERING

9940 572 462  6381 907 438  ijircce@gmail.com

Scan to save the contact details