

# Customer Sales Predication Using Artificial Intelligence

M.Sheeba<sup>1</sup>, k.Srivalli<sup>2</sup>, Y.Venkataprasanna<sup>3</sup>, CH.Devaki<sup>4</sup>

<sup>1</sup>Assistant Professor, Dept of AI&DS

<sup>2,3,4</sup> Dept of AI&DS

<sup>1,2,3,4</sup> Dhanalakshmi Srinivasan University, Tamilnadu, india

**Abstract-** Customer churn prediction is an important research problem in business analytics that helps organizations identify customers who are likely to discontinue their services. This paper proposes a machine learning-based churn prediction framework using Decision Tree, Random Forest, Logistic Regression, and XGBoost algorithms. The dataset is preprocessed through missing-value handling, categorical encoding, normalization, and feature selection techniques to improve prediction accuracy. Experimental results show that the XGBoost model provides higher performance compared with other algorithms. The proposed system supports organizations in identifying high-risk customers at an early stage and enables proactive retention strategies that improve customer satisfaction and organizational profitability. This paper focuses on Customer Churn Prediction using Artificial Intelligence and Machine Learning techniques to identify customers who are likely to discontinue services in advance. The proposed system uses the Telco Customer Churn dataset and applies machine learning algorithms such as Decision Tree, Random Forest, Logistic Regression, and XGBoost to analyze customer behavior based on attributes like tenure, contract type, monthly charges, and payment methods. Data preprocessing techniques including handling missing values, encoding categorical variables, and feature selection are performed to improve prediction accuracy. Among all the algorithms, XGBoost achieved the highest performance, making it the most effective model for churn prediction. The developed system helps organizations take proactive retention strategies, reduce customer loss, and improve business decision-making through data-driven insights.

**Keywords:** Customer Churn Prediction, Machine learning, Artificial Intelligence, XGBoost, Predictive Analytics, Customer Retention

## I. INTRODUCTION

1.1 Customer retention plays a crucial role in improving the profitability and sustainability of modern organizations. In highly competitive markets such as telecommunications, banking, and subscription-based services, identifying customers who are likely to discontinue services has become an important research challenge. Traditional methods rely mainly on manual analysis and historical reporting, which are

insufficient for predicting future customer behavior accurately. Artificial Intelligence and Machine Learning techniques enable organizations to analyze customer-related data effectively and identify churn patterns at an early stage.

### 1.2 Motivation and Problem Statement

Customer churn is one of the major challenges faced by organizations because acquiring new customers is more expensive than retaining existing ones. Many companies experience significant revenue loss due to the inability to identify customers who are likely to discontinue services. Traditional analytical approaches lack predictive capability and cannot efficiently handle large volumes of customer data. Therefore, there is a strong need to develop an intelligent prediction system that can analyze customer behavior and detect potential churn in advance.

### 1.3 Project Objectives

The objectives of this research are:

#### A. Customer Churn Prediction

To predict customer departure behavior.

#### B. Algorithm Comparison

To compare machine learning models.

#### C. Data Preprocessing

To improve dataset quality.

#### D. Prediction Accuracy

To increase model performance.

#### E. Retention Support

To assist customer retention strategies.

#### F. Decision Support

To support business decisions.

## G. Accuracy Improvement

## II. LITERATURE REVIEW

Customer churn prediction has become an important research area in business analytics and machine learning because customer retention directly affects organizational profitability and long-term growth. Many researchers have proposed different statistical, data mining, and machine learning techniques to identify customers who are likely to discontinue services. Earlier studies mainly focused on traditional statistical methods, while recent research emphasizes Artificial Intelligence and Machine Learning models for improving prediction accuracy and decision-making.

Traditional churn prediction approaches used statistical analysis techniques such as mean analysis, variance analysis, regression models, and probability-based methods to analyze customer behavior patterns. These techniques helped organizations understand customer satisfaction levels and estimate the probability of churn. However, traditional methods were limited in handling large datasets and complex customer behavior patterns. They also lacked predictive intelligence and real-time analytical capability.

With the advancement of Machine Learning, researchers introduced classification algorithms such as Decision Tree, Random Forest, Logistic Regression, Support Vector Machine, and XGBoost for customer churn prediction. These algorithms improved prediction accuracy by automatically learning customer behavior patterns from historical datasets. Machine learning models can process large amounts of customer information, identify hidden relationships among features, and provide better predictive performance compared with traditional approaches.

Decision Tree algorithms are widely used for churn classification because they create a rule-based decision structure that is easy to understand and interpret. The model divides customer data into different branches based on important attributes such as tenure, monthly charges, and service usage. Researchers found that Decision Trees provide fast classification results and help organizations understand customer churn factors clearly. However, Decision Trees may suffer from overfitting when handling complex datasets.

Random Forest was introduced to overcome the limitations of Decision Trees by combining multiple decision trees into a single predictive model. Researchers observed that

Random Forest improves classification stability, reduces overfitting problems, and provides reliable prediction accuracy. This model is highly effective for handling high-dimensional customer datasets and noisy data. Many studies reported that Random Forest achieved better performance than traditional statistical methods in churn prediction systems.

Logistic Regression is another important machine learning technique used in churn analysis. It is mainly applied for binary classification problems where customers are classified as churn or non-churn. Researchers used Logistic Regression because of its simplicity, efficiency, and probability-based prediction capability. It helps organizations estimate the likelihood of customer departure and supports business decision-making processes. Although Logistic Regression provides moderate prediction accuracy, it performs effectively on structured datasets with linear relationships.

Recent research studies focused on advanced boosting algorithms such as XGBoost because of their high prediction accuracy and efficient learning capability. XGBoost uses boosting techniques to combine multiple weak learning models into a strong predictive model. Researchers found that XGBoost handles large datasets efficiently, reduces prediction error, and minimizes overfitting problems. Many experimental studies showed that XGBoost outperformed Decision Tree, Random Forest, and Logistic Regression in customer churn prediction tasks.

Several researchers also emphasized the importance of data preprocessing and feature selection in improving model performance. Data preprocessing techniques such as missing value handling, categorical encoding, normalization, and duplicate removal help improve dataset quality and prediction accuracy. Feature selection methods identify important customer attributes such as contract type, payment method, tenure, and monthly charges, which significantly influence churn behavior.

Recent studies further explored the use of deep learning techniques, neural networks, and real-time predictive analytics for customer churn management. Researchers suggested integrating machine learning models with cloud-based business intelligence systems to support automated decision-making and proactive customer retention strategies. These advanced systems enable organizations to identify high-risk customers at an early stage and take preventive actions to reduce customer loss.

From the literature review, it is observed that machine learning techniques provide more accurate and efficient customer churn prediction compared with traditional

analytical methods. Among various algorithms, XGBoost and Random Forest achieved better prediction performance due to their ability to handle complex customer behavior patterns and large datasets effectively. Therefore, this research focuses on implementing and comparing multiple machine learning models to develop an efficient customer churn prediction system that supports business organizations in improving customer retention and profitability.

### III. RESEARCH METHODOLOGY

The research methodology describes the systematic procedure followed to develop the Customer Churn Prediction System using Artificial Intelligence and Machine Learning techniques. The proposed methodology focuses on identifying customers who are likely to discontinue services and helping organizations improve customer retention strategies. The methodology includes dataset collection, data preprocessing, feature selection, machine learning model development, prediction, and performance evaluation. The complete process is designed to improve prediction accuracy and support data-driven business decision-making.

#### 3.1 Dataset Collection

The first stage of the proposed methodology is dataset collection. The system uses the Telco Customer Churn Dataset, which contains customer-related information collected from telecommunication services. The dataset includes demographic information, customer account details, contract information, payment methods, service usage details, monthly charges, total charges, and customer churn status. These attributes help in understanding customer behavior patterns and identifying important factors responsible for customer churn.

The collected dataset provides sufficient information for training machine learning algorithms and performing predictive analysis. Proper dataset collection is important because the quality of prediction depends on the quality and completeness of customer information available in the dataset.

#### 3.2 Data Preprocessing

Data preprocessing is performed to improve dataset quality before applying machine learning algorithms. Real-world datasets often contain missing values, inconsistent records, duplicate data, and categorical information that cannot be directly processed by machine learning models. Therefore, preprocessing is necessary to convert raw customer data into a clean and structured format.

Missing values are identified and replaced using suitable techniques to avoid incorrect prediction results. Duplicate records are removed to maintain dataset consistency. Categorical attributes such as payment methods, gender, and contract types are converted into numerical format using encoding techniques. Numerical values such as monthly charges and total charges are normalized to improve model training efficiency. The preprocessing stage reduces data inconsistency, improves dataset quality, and enhances prediction accuracy.

#### 3.3 Feature Selection

Feature selection is an important step in the proposed methodology because not all customer attributes contribute equally to churn prediction. The feature selection process identifies the most relevant features influencing customer churn behavior and removes unnecessary or irrelevant data from the dataset.

Important features such as customer tenure, contract type, monthly charges, internet services, payment methods, and total charges are selected for prediction analysis. Selecting important features helps reduce computational complexity, improves model performance, and increases prediction accuracy. Feature selection also minimizes overfitting problems and supports efficient machine learning model training.

#### 3.4 Data Splitting

After preprocessing and feature selection, the dataset is divided into training and testing datasets. The training dataset is used to train machine learning algorithms, while the testing dataset is used to evaluate model performance on unseen data. Generally, eighty percent of the dataset is used for training and twenty percent is used for testing.

Data splitting is necessary to validate the effectiveness of machine learning models and prevent overfitting. It ensures that the developed system performs accurately not only on training data but also on new customer data. Proper data splitting improves prediction reliability and supports fair comparison between different algorithms.

#### 3.5 Machine Learning Model Development

The proposed system uses multiple machine learning algorithms to classify customers into churn and non-churn categories. Different algorithms are implemented and compared to identify the best predictive model for customer churn prediction.

Decision Tree is used as a classification algorithm that creates rule-based decision structures for analyzing customer behavior patterns. Random Forest improves prediction stability by combining multiple decision trees into a single ensemble model. Logistic Regression is applied for binary classification and probability-based prediction of customer churn. XGBoost is used as an advanced boosting algorithm that combines weak learning models into a strong predictive system with improved accuracy.

The use of multiple algorithms helps in analyzing prediction performance and selecting the model with the highest accuracy and efficiency.

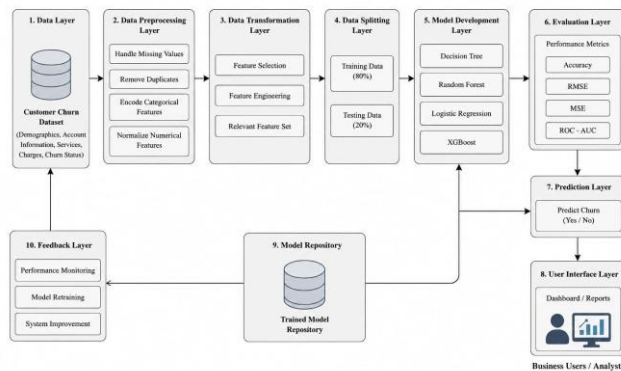


Fig. 1. Architecture of Customer Churn Prediction System

### 3.6 Model Training

During the training phase, machine learning algorithms learn customer behavior patterns from historical customer data. The training process helps the models identify relationships between customer attributes and churn status. The algorithms analyze input features, generate prediction rules, and optimize internal parameters to improve classification performance.

Proper model training is essential because the effectiveness of the churn prediction system depends on how well the algorithms learn from customer data. The trained models are capable of predicting whether a customer is likely to continue or discontinue services in the future.

### 3.7 Customer Churn Prediction

After model training, the developed system predicts customer churn by classifying customers into churn and non-churn categories. Customers identified as high-risk are considered likely to discontinue services. These predictions help organizations take proactive retention actions such as providing discounts, improving customer support, or offering personalized services.

The prediction system enables organizations to reduce customer loss, improve customer satisfaction, and increase profitability through early identification of churn behavior. Accurate prediction also supports better business planning and customer relationship management.

### 3.8 Model Evaluation

The performance of machine learning algorithms is evaluated using different evaluation metrics such as Accuracy, Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and Confusion Matrix analysis. These evaluation techniques help measure prediction accuracy, classification efficiency, and model reliability.

The experimental results indicate that XGBoost achieved the highest prediction accuracy among all implemented algorithms. Random Forest produced stable classification results, while Logistic Regression showed moderate prediction performance. Decision Tree provided basic classification capability but lower accuracy compared with advanced algorithms.

The evaluation process helps identify the best-performing model and ensures the reliability of the proposed customer churn prediction system.

### 3.9 Result Analysis

The final stage of the methodology involves analyzing prediction results and comparing the performance of different machine learning models. The analysis shows that data preprocessing and feature selection significantly improved model performance and prediction accuracy. The developed system successfully identified customers who are likely to discontinue services and supported organizations in implementing proactive customer retention strategies.

The proposed methodology demonstrates that Artificial Intelligence and Machine Learning techniques can effectively predict customer churn and provide valuable insights for business decision-making. The system helps organizations improve customer satisfaction, reduce revenue loss, and enhance overall business performance through intelligent predictive analytics.

## IV. MATHEMATICAL AND OPTIMIZATION

### 1. Mean Squared Error (MSE)

Purpose: Measures the average squared difference between actual and predicted values.

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

## 2. Root Mean Squared Error (RMSE)

Purpose: Shows the overall prediction error of the model.

$$RMSE = \sqrt{MSE}$$

## 3. Accuracy Formula

Purpose: Calculates correct prediction percentage.

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$$

## 4. Precision Formula

Purpose: Measures how many predicted churn customers are actually churn customers.

$$Precision = \frac{TP}{TP+FP}$$

## 5. Recall Formula

Purpose: Measures how many actual churn customers are correctly identified.

$$P(Y = 1) = \frac{1}{1+e^{-(b_0+b_1x_1+b_2x_2+\dots+b_nx_n)}}$$

## 6. F1-Score Formula

Purpose: Harmonic mean of precision and recall.

$$Recall = \frac{TP}{TP+FN}$$

## 7. Logistic Regression Formula

Purpose: Used for churn probability prediction.

$$Entropy = - \sum p_i \log_2(p_i)$$

## V. RESULTAND DISCUSSION

The proposed Customer Churn Prediction System was implemented using Artificial Intelligence and Machine Learning techniques to identify customers who are likely to discontinue services. Different machine learning algorithms such as Decision Tree, Random Forest, Logistic Regression, and XGBoost were applied to the Telco Customer Churn Dataset to analyze customer behavior patterns and predict churn probability. The system was evaluated using various performance metrics to compare the efficiency and prediction capability of each algorithm.

### 5.1 Experimental Results

The experimental analysis showed that machine learning techniques can effectively predict customer churn using customer-related attributes such as tenure, contract type, payment method, internet services, monthly charges, and total charges. During the implementation process, the dataset was

preprocessed using missing value handling, categorical encoding, normalization, and feature selection techniques to improve prediction performance.

The machine learning models were trained using historical customer data and tested using unseen data to evaluate prediction accuracy. The experimental results demonstrated that all implemented algorithms successfully classified customers into churn and non-churn categories, but their prediction performance varied based on model complexity and learning capability.

Among all the algorithms, XGBoost achieved the highest prediction accuracy because of its advanced boosting mechanism and efficient handling of complex customer behavior patterns. Random Forest also produced stable and reliable classification results due to its ensemble learning capability. Logistic Regression showed moderate performance and worked efficiently for binary classification problems, while Decision Tree provided basic classification results with comparatively lower accuracy.

The proposed system successfully identified high-risk customers in advance and supported organizations in understanding customer churn behavior more effectively.

### 5.2 Performance Evaluation

The developed models were evaluated using different performance metrics such as Accuracy, Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and ROC-AUC score. These metrics helped measure prediction efficiency, classification capability, and model reliability.

Accuracy was used to determine the percentage of correctly classified churn and non-churn customers. The results indicated that XGBoost achieved the highest accuracy among all algorithms. Random Forest also produced high accuracy with stable predictions, whereas Logistic Regression and Decision Tree showed comparatively moderate performance.

MSE and RMSE values were used to measure prediction error. Lower error values indicated better model performance. XGBoost produced the lowest RMSE and MSE values, showing that the model generated more accurate churn predictions with minimal error.

The ROC-AUC score was used to analyze the classification capability of the models. Higher ROC-AUC values indicated better discrimination between churn and non-churn customers. XGBoost and Random Forest achieved

better ROC-AUC performance compared with other algorithms, demonstrating their effectiveness in customer churn classification tasks.

### 5.3 Discussion

The results of the proposed research demonstrate that Artificial Intelligence and Machine Learning techniques provide efficient solutions for customer churn prediction. The study shows that preprocessing and feature selection significantly improve model performance and prediction accuracy. Proper handling of missing values, normalization, and encoding techniques enhanced the quality of the dataset and reduced prediction errors.

The comparison between different machine learning algorithms indicates that ensemble and boosting techniques provide better predictive performance than traditional classification models. Decision Tree algorithms are simple and easy to understand, but they may suffer from overfitting problems when handling large and complex datasets. Logistic Regression provides efficient binary classification, but its prediction capability is limited when customer behavior becomes highly complex.

Random Forest improved prediction stability by combining multiple decision trees and reducing overfitting problems. However, XGBoost achieved the best overall performance because of its boosting mechanism, optimization capability, and efficient handling of large datasets. The algorithm successfully identified hidden relationships between customer attributes and churn behavior, resulting in higher prediction accuracy.

The proposed customer churn prediction system can help organizations take proactive retention actions before customers discontinue services. Businesses can use the prediction results to provide personalized offers, improve customer support services, and enhance customer satisfaction. Early identification of churn customers helps reduce revenue loss and supports better business decision-making.

The developed system also demonstrates the importance of predictive analytics in modern business environments. By integrating machine learning models into organizational systems, companies can automate customer analysis processes and improve operational efficiency. The proposed methodology can be further extended using deep learning models, real-time analytics, and cloud-based predictive systems for improved scalability and performance.

## VI. CONCLUSION

This research focused on Customer Churn Prediction using Artificial Intelligence and Machine Learning techniques to identify customers who are likely to discontinue services in advance. The proposed system used the Telco Customer Churn Dataset to analyze customer behavior based on attributes such as tenure, contract type, payment method, monthly charges, total charges, and internet services. Different machine learning algorithms including Decision Tree, Random Forest, Logistic Regression, and XGBoost were implemented to predict customer churn accurately.

The collected dataset was preprocessed using techniques such as missing value handling, duplicate removal, categorical encoding, normalization, and feature selection to improve dataset quality and prediction performance. These preprocessing techniques helped reduce inconsistencies in the dataset and enhanced the efficiency of machine learning models.

The experimental analysis showed that all implemented algorithms successfully classified customers into churn and non-churn categories. Among all the algorithms, XGBoost achieved the highest prediction accuracy and best overall performance because of its advanced boosting capability and efficient handling of complex customer behavior patterns. Random Forest also provided stable and reliable prediction results, while Logistic Regression and Decision Tree showed moderate classification performance.

The proposed customer churn prediction system helps organizations identify high-risk customers at an early stage and supports proactive customer retention strategies. By using predictive analytics, organizations can improve customer satisfaction, reduce customer loss, and increase profitability through data-driven decision-making.

The research demonstrates that Artificial Intelligence and Machine Learning techniques provide effective solutions for customer churn prediction and business analytics. In future work, the proposed system can be further improved using deep learning models, real-time predictive analytics, and cloud-based systems to enhance scalability, automation, and prediction accuracy.


## VII. ACKNOWLEDGMENT

We express our sincere gratitude to our guide and faculty members for their valuable guidance, continuous support, and encouragement throughout the completion of this research work on “Customer Churn Prediction Using Artificial

Intelligence and Machine Learning Techniques.” Their suggestions and motivation greatly helped us in completing this project successfully.

We also thank our institution and department for providing the necessary facilities and academic support required for this research. Finally, we express our heartfelt thanks to our parents and friends for their constant encouragement and support during the completion of this work.

## REFERENCES

- [1] I. Idris, A. Khan, and Y. S. Lee, “Intelligent Churn Prediction in Telecom Using Machine Learning Methods,” *Expert Systems with Applications*, 2012.
- [2] T. Hadden, A. Tiwari, R. Roy, and D. Ruta, “Computer Assisted Customer Churn Management: State-of-the-Art and Future Trends,” *Computers & Operations Research*, 2007.
- [3] J. Han, M. Kamber, and J. Pei, *Data Mining: Concepts and Techniques*, Morgan Kaufmann Publishers, 2011.
- [4] T. Chen and C. Guestrin, “XGBoost: A Scalable Tree Boosting System,” *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2016.
- [5] L. Breiman, “Random Forests,” *Machine Learning Journal*, 2001.
- [6] D. W. Hosmer and S. Lemeshow, *Applied Logistic Regression*, Wiley Publications, 2013.
- [7] S. Amin, A. Anwar, A. Adnan, M. Nawaz, A. Alawfi, A. Hussain, and K. Huang, “Customer Churn Prediction in Telecommunication Industry Using Data Mining Techniques,” *IEEE Access*, 2019.
- [8] E. Verbeke, D. Martens, and B. Baesens, “Social Network Analysis for Customer Churn Prediction,” *Applied Soft Computing*, 2014.
- [9] K. Coussement and D. Van den Poel, “Churn Prediction in Subscription Services Using Support Vector Machines,” *Expert Systems with Applications*, 2008.
- [10] V. R. Prykhodko and A. A. Pakhomov, “Machine Learning Approaches for Customer Churn Prediction,” *Procedia Computer Science*, 2020.
- [11] F. Provost and T. Fawcett, *Data Science for Business: What You Need to Know About Data Mining and Data-Analytic Thinking*, O’Reilly Media, 2013.
- [12] G. James, D. Witten, T. Hastie, and R. Tibshirani, *An Introduction to Statistical Learning*, Springer, 2013.
- [13] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Springer, 2009.
- [14] P. C. Pendharkar, “Genetic Algorithm Based Neural Network Approaches for Predicting Customer Churn,” *Expert Systems with Applications*, 2009.
- [15] A. Ng, *Machine Learning Yearning*, DeepLearning.ai, 2018.
- [16] IBM Watson Analytics, “Telco Customer Churn Dataset Documentation,” IBM Corporation.
- [17] Scikit-learn Machine Learning Library Documentation, <https://scikit-learn.org> 
- [18] I. Witten, E. Frank, and M. Hall, *Data Mining: Practical Machine Learning Tools and Techniques*, Morgan Kaufmann Publishers, 2011.
- [19] C. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006.
- [20] S. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach*, Pearson Education, 2010.
- [21] R. Kohavi and F. Provost, “Glossary of Terms,” *Machine Learning Journal*, 1998.
- [22] Y. LeCun, Y. Bengio, and G. Hinton, “Deep Learning,” *Nature Journal*, 2015.