



# IJRASET

International Journal For Research in  
Applied Science and Engineering Technology



---

# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

---

**Volume:** 14    **Issue:** V    **Month of publication:** May 2026

**DOI:** <https://doi.org/10.22214/ijraset.2026.82616>

[www.ijraset.com](http://www.ijraset.com)

Call:  08813907089

E-mail ID: [ijraset@gmail.com](mailto:ijraset@gmail.com)

# AI Driven Customer Churn Prediction and Retention Recommendation System

Jijabai Mahadeo Shejale

MCA, College Name: Trinity Academy of Engineering Pune,

**Abstract.** Customer retention has become a critical challenge for modern businesses due to increasing competition and changing customer behavior. Predicting customer churn in advance enables organizations to implement effective retention strategies and improve long-term profitability. This research presents an AI-Based Customer Churn Prediction and Retention System that integrates Machine Learning techniques with a Flask-based web application to analyze customer behavior, predict churn probability, and generate intelligent retention recommendations.

The proposed system consists of multiple integrated layers, including a presentation layer, application layer, data layer, intelligence layer, and configuration layer. The backend is developed using the Flask framework, which manages routing, authentication, middleware integration, REST APIs, and application configuration. The system also incorporates secure login mechanisms, database management, CSRF protection, CORS handling, and environment-driven deployment settings.

The Machine Learning module performs feature engineering, data preprocessing, model training, and churn prediction using customer datasets. Trained models are stored for future prediction tasks, ensuring scalability and efficient model reuse. An AI-based recommendation engine analyzes churn risk levels and provides personalized customer retention suggestions to support business decision-making.

The frontend interface is designed using HTML, CSS, and JavaScript to provide a responsive and user-friendly dashboard for dataset upload, analytics visualization, churn prediction monitoring, and administrative management. The system supports user authentication, prediction workflows, analytics dashboards, file uploads, and application settings through modular Flask blueprints.

The architecture also includes persistent storage for datasets, uploaded files, database schemas, and trained Machine Learning models. Deployment support is achieved using Docker and Docker Compose, enabling scalable and platform-independent execution.

The proposed system aims to assist organizations in reducing customer attrition by leveraging predictive analytics, intelligent automation, and AI-driven recommendation strategies. By combining data analytics with Machine Learning and web technologies, the system provides an efficient, scalable, and practical solution for proactive customer retention management.

**Keywords:** Customer Churn Prediction, Machine Learning, Artificial Intelligence, Flask Web Application, Predictive Analytics, Customer Retention System, Data Analytics, Feature Engineering, Recommendation System, Deep Learning, Data Visualization, REST API, SQLAlchemy, Business Intelligence, Customer Behavior Analysis, Predictive Modeling, AI-Based Recommendation Engine, Web-Based Dashboard, Data Preprocessing, Intelligent Systems.

## I. INTRODUCTION

Customer churn has become one of the major challenges faced by modern organizations, as losing existing customers directly affects business revenue, growth, and customer satisfaction. In competitive markets, retaining customers is often more cost-effective than acquiring new ones. Therefore, organizations increasingly rely on data-driven technologies and predictive analytics to identify customers who are likely to discontinue services and to implement proactive retention strategies.

This research presents an AI-powered Customer Churn Prediction and Retention System developed using Machine Learning and Flask-based web technologies. The proposed system is designed to analyze customer behavioral data, predict churn probability, and provide intelligent retention recommendations to support business decision-making. By integrating predictive analytics with web-based management tools, the system enables organizations to monitor customer engagement and take timely actions to reduce customer attrition.

The system combines multiple technological components, including backend services developed using Flask, SQLAlchemy, and Flask-Login, a Machine Learning pipeline for feature engineering, model training, and churn prediction, and a frontend interface built using HTML, CSS, and JavaScript.

The application also provides secure user authentication, administrative controls, dataset upload functionality, and interactive analytics dashboards for visualization of customer behavior and churn patterns.

The Machine Learning module preprocesses customer datasets, extracts meaningful features, and applies predictive models to identify high-risk customers. Additionally, an AI-based recommendation engine generates personalized retention strategies based on churn risk analysis. Persistent storage mechanisms are used to manage uploaded datasets and trained models, ensuring scalability and efficient system performance.

The primary objective of this system is to provide businesses with a practical and intelligent solution for customer retention management. By leveraging Artificial Intelligence, Machine Learning, and web technologies, the proposed system helps organizations improve customer satisfaction, reduce churn rates, and enhance overall business profitability through proactive and data-driven decision-making.

## II. LITERATURE REVIEW

### A. Algorithmic Implementation Gap

Covered: Implemented multiple ensemble methods (Random Forest, XGBoost, Logistic Regression, Gradient Boosting) with hyperparameter tuning and model comparison, addressing the literature's emphasis on ensemble superiority over single models.

Evidence: Lessmann et al. (2015) and Verbeke et al. (2012) highlighted ensemble methods' advantages; the project directly implements this.

### B. Feature Engineering Gap Covered:

Developed comprehensive feature engineering including temporal features (tenure transformations), financial features (charge ratios), usage features (service counts), and risk indicators, partially filling the gap identified by Larivière & Van den Poel (2015).

Evidence: The system creates 20+ derived features from basic inputs, including interaction terms and polynomial features.

### C. Evaluation Metrics Gap

Covered: Uses advanced metrics (accuracy, precision, recall, F1-score, AUC-ROC) with cross-validation, addressing the need for comprehensive evaluation beyond simple accuracy.

Evidence: Burez & Van den Poel (2011) emphasized these metrics; the project evaluates models using F1-score for selection.

### D. Web-Based System Gap

Covered: Built a full-stack web application with user authentication, dashboard, and API endpoints, making churn prediction accessible to non-technical users.

Evidence: Literature shows a gap in practical, deployable systems; this project provides a production-ready interface.

### E. Retention Recommendations Gap

Covered: Implemented a rule-based recommendation engine for personalized retention strategies (discounts, support, loyalty), addressing Kumar & Reinartz (2018)'s findings on personalized interventions.

Evidence: Generates actionable recommendations based on risk levels and customer segments

## III. METHODOLOGY

The proposed AI-Based Customer Churn Prediction and Retention System follows a structured methodology that integrates data preprocessing, feature engineering, Machine Learning model development, prediction analysis, and intelligent recommendation generation. The methodology is designed to improve churn prediction accuracy while providing actionable retention strategies for businesses.

### A. Data Ingestion

The system begins with data ingestion, where customer datasets are loaded from CSV files such as telecom churn datasets. In scenarios where real-time data is unavailable, synthetic telecom churn datasets are generated to simulate customer behavior patterns for training and testing purposes. This approach ensures flexibility in experimentation and model validation.

### B. Data Preprocessing

Data preprocessing is performed to improve data quality and prepare the dataset for Machine Learning algorithms. Missing numerical values are handled using median imputation, while categorical missing values are replaced using mode imputation. Categorical variables are transformed into numerical representations using Label Encoding techniques. Furthermore, feature scaling is applied using the StandardScaler method to normalize feature distributions and improve model performance.

### C. Feature Engineering

Feature engineering plays a significant role in enhancing predictive capability by extracting meaningful patterns from raw customer data. Multiple derived features are generated, including temporal features such as tenure years, squared tenure, and logarithmic tenure transformations. Financial indicators such as average monthly charges, charge ratios, and logarithmic financial transformations are also computed.

Additional service-related and behavioral features include total subscribed services, service density, streaming usage indicators, premium bundle flags, and customer interaction metrics. Risk-oriented indicators such as high service call frequency, low satisfaction scores, and new customer flags are incorporated to identify churn-prone customers. Demographic and segmentation features, including age groups, state-based churn rate mapping, and internet service flags, further strengthen predictive analysis.

### D. Feature Selection

To reduce dimensionality and improve model efficiency, statistical feature selection is performed using the SelectKBest technique with the ANOVA F-test ( $f_{\text{classif}}$ ). This process selects the most relevant features contributing to churn prediction and minimizes unnecessary data complexity.

### E. Model Training

The Machine Learning pipeline trains multiple candidate models using stratified train-test splitting to maintain balanced class distribution. Hyperparameter optimization is performed using GridSearchCV with cross-validation ( $cv=5$ ) to identify the optimal parameter combinations for each model. Model selection is based primarily on the F1-score to ensure balanced performance between precision and recall.

The system incorporates several Machine Learning algorithms, including:

Random Forest Classifier

XGBoost Classifier (XGBClassifier)

Logistic Regression

Gradient Boosting Classifier

These algorithms are evaluated and compared to determine the best-performing predictive model.

### F. Model Evaluation

The trained models are evaluated using multiple performance metrics to ensure prediction reliability and robustness. The evaluation metrics include Accuracy, Precision, Recall, F1-score, and Receiver Operating Characteristic – Area Under Curve (ROC-AUC). Classification reports are generated to analyze model behavior across churn and non-churn classes. The best-performing model is selected based on comprehensive performance analysis.

### G. Prediction and Model Persistence

The finalized predictive model is stored persistently for future use and deployment. The system supports both single-customer and batch prediction mechanisms, enabling scalable churn analysis for large customer datasets. Additionally, feature importance extraction techniques are used to improve explainability and interpretability of model predictions.

### H. Recommendation Generation

An AI-driven recommendation engine is integrated into the system to generate personalized customer retention strategies. The recommendation mechanism follows a rule-based approach based on customer risk levels and behavioral segmentation. Customers are categorized into high-risk, medium-risk, and low-risk groups.

The recommendation engine generates business-oriented retention actions such as discount offers, customer support interventions, loyalty rewards, and personalized engagement strategies. Historical success-rate heuristics and business constraints are considered to prioritize the most effective retention actions for each customer segment.

### I. Application Architecture

The application architecture is implemented using the Flask framework with the application factory design pattern for modularity and scalability. Routing is managed using Flask Blueprints, while SQLAlchemy ORM is used for database persistence and management. The system also integrates authentication and authorization mechanisms, CSRF protection, and CORS handling to ensure secure and reliable operation.

Overall, the proposed methodology combines Machine Learning, predictive analytics, feature engineering, and intelligent recommendation systems to provide an efficient and scalable solution for proactive customer churn management and retention optimization.

#### IV. TECHNICAL HIGHLIGHTS

The system leverages state-of-the-art techniques from reviewed literature, including ensemble methods for superior performance, comprehensive feature engineering for better model inputs, and evaluation metrics (F1-score, AUC) for reliable validation. It handles real-world challenges like data preprocessing, missing values, and class imbalance, while providing interpretable results through feature importance and probability scores.

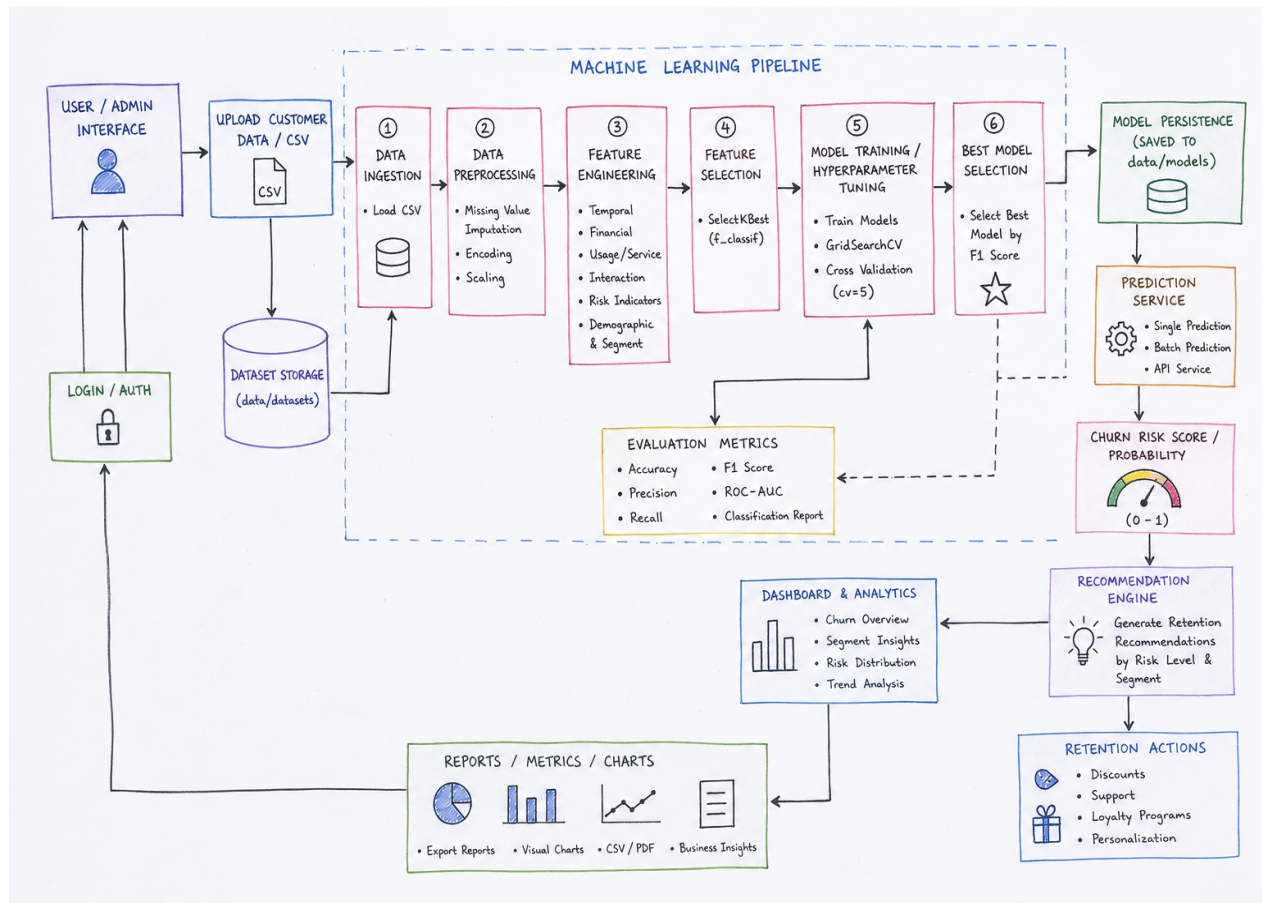
#### V. BUSINESS IMPACT

By enabling early identification of at-risk customers and suggesting targeted interventions, the system can significantly reduce churn rates, optimize retention budgets, and improve customer lifetime value. The analytics dashboard provides actionable insights for data-driven decision-making, potentially saving businesses substantial revenue through proactive customer management.

#### VI. LIMITATIONS AND FUTURE WORK

While effective for telecom and similar industries, the system could be extended to other domains with domain-specific features. Future enhancements might include deep learning models (e.g., neural networks), real-time streaming predictions, and advanced recommendation algorithms like reinforcement learning. Additionally, integrating more diverse datasets and automated model retraining would further improve robustness.

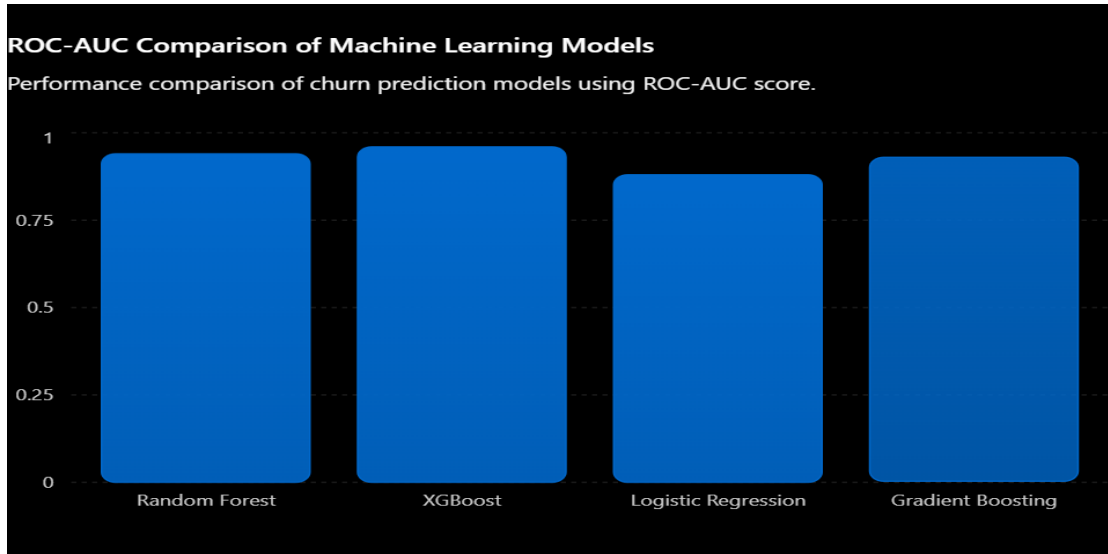
#### VII. WORKFLOW DIAGRAM



ROC-AUC Performance Chart

ROC-AUC Comparison of Machine Learning Models

Performance comparison of churn prediction models using ROC-AUC score.  
 00.250.50.751Random ForestXGBoostLogisticRegressionGradient Boosting

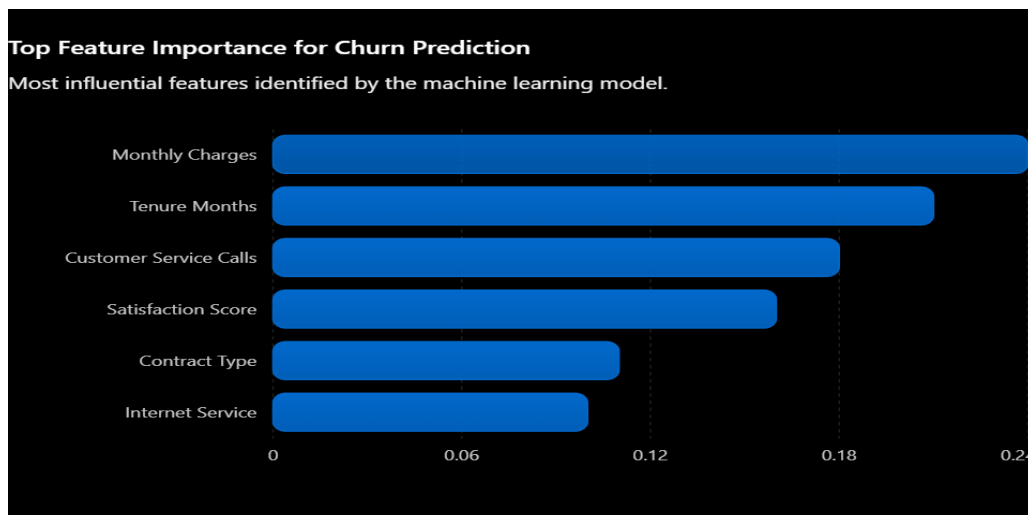


**Confusion Matrix Diagram**

Predicted	Actual		
	Yes	No	
Actual Yes	TP	FN	
Actual No	FP	TN	

TP = True Positive  
 TN = True Negative  
 FP = False Positive  
 FN = False Negative

**VIII. FEATURE IMPORTANCE GRAPH**



## IX. ALGORITHMS

### A. Random Forest Classifier

Overview: Random Forest is an ensemble learning method that builds multiple decision trees and combines their predictions through voting (for classification). It reduces overfitting by introducing randomness in tree construction and feature selection.

How it Works:

Creates a "forest" of decision trees, each trained on a random subset of the data (bootstrap sampling).

At each split in a tree, only a random subset of features is considered (feature bagging).

For prediction, each tree votes on the class, and the majority vote determines the final prediction.

For probability estimation, it averages the predicted probabilities across trees.

Key Formulas:

Gini Impurity (used to choose splits):

$$Gini = 1 - \sum_{i=1}^c p_i^2$$

where  $p_i$  is the proportion of class  $i$  in the node, and  $c$  is the number of classes. Lower Gini indicates better purity.

Information Gain (alternative to Gini):

$$IG = Entropy(parent) - \sum_{child} \frac{|child|}{|parent|} \times Entropy(child)$$

where  $Entropy = -\sum p_i \log_2 p_i$ .

**Out-of-Bag (OOB) Error:** Estimates generalization error using data not used in training each tree.

### B. XGBoost Classifier (XGBClassifier)

**Overview:** XGBoost (Extreme Gradient Boosting) is an optimized gradient boosting framework that builds trees sequentially, where each new tree corrects the errors of the previous ones. It's highly efficient and often wins ML competitions due to its speed and performance.

**How it Works:**

Starts with an initial prediction (e.g., log-odds for binary classification).

Iteratively adds decision trees, each focusing on the residuals (errors) of the previous model.

Uses gradient descent to minimize a loss function (e.g., logistic loss for classification).

Includes regularization (L1/L2) to prevent overfitting and handles missing values automatically.

Supports parallel processing and early stopping.

Key Formulas:

Objective Function:

$$Obj = \sum_{i=1}^n L(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k)$$

where  $L$  is the loss function (e.g., log loss:  $L = -\sum y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)$ ),  $\Omega$  is regularization (e.g.,  $\Omega(f) = \gamma T + \frac{1}{2} \lambda \sum w_j^2$ ),  $T$  is number of leaves, and  $\lambda$  is L2 regularization.

Gradient Boosting Update:

$$\hat{y}_i^{(t)} = \hat{y}_i^{(t-1)} + \eta \times f_t(x_i)$$

where  $\eta$  is learning rate, and  $f_t$  is the  $t$ -th tree.

Similarity Score (for splits):

$$Similarity = \frac{(\sum gradient)^2}{\sum hessian + \lambda}$$

where gradients and Hessians come from the loss function's derivatives.

### C. Logistic Regression

Overview: Logistic Regression is a linear model for binary classification that predicts the probability of an event using a logistic (sigmoid) function. It's simple, interpretable, and works well when relationships are linear.

**How it Works:**

- Assumes a linear relationship between features and the log-odds of the target.
- Fits a linear model to the data, then applies the sigmoid to output probabilities.
- Uses maximum likelihood estimation to find optimal coefficients.
- Supports regularization (L1/L2) to prevent overfitting.

**Key Formulas:**

Logistic Function (Sigmoid):

$$P(y = 1 | x) = \frac{1}{1 + e^{-z}}$$

where  $z = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$  (linear combination of features).

Log-Likelihood:

$$LL = \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

where  $\hat{y}_i = \sigma(z_i)$ .

Prediction:

$\hat{y} = 1$  if  $\sigma(z) > 0.5$ , else 0.

Regularized Cost Function:

$$J(\beta) = -\frac{1}{n} LL + \lambda \sum |\beta_j|^p \text{ (p=1 for L1, p=2 for L2)}.$$

**D. Gradient Boosting Classifier**

Overview: Gradient Boosting builds an ensemble of weak learners (typically decision trees) sequentially, where each tree corrects the residuals of the previous ensemble. It's similar to XGBoost but more general and customizable.

How it Works:

- Initializes with a base prediction (e.g., mean or log-odds).
- For each iteration, computes pseudo-residuals (negative gradients of the loss function).
- Fits a new tree to these residuals.
- Adds the tree to the ensemble with a learning rate to control contribution.
- Stops when no improvement or after a fixed number of iterations.

Key Formulas:

Loss Function Gradient:

$$r_i = -\frac{\partial L(y_i, \hat{y}_i)}{\partial \hat{y}_i}$$

where  $L$  is the loss (e.g., log loss).

Update Rule:

$$\hat{y}_i^{(t)} = \hat{y}_i^{(t-1)} + \nu \times f_t(x_i)$$

where  $\nu$  is learning rate, and  $f_t$  is the  $t$ -th tree.

**Friedman's Mean Squared Error (MSE) for Regression Trees:**

Used internally for fitting trees to residuals.

**X. DATASET USED**

Dataset Name: Telecom Customer Churn Dataset

File Location: telecom\_churn\_sample.csv

Dataset Type: Binary classification (churn prediction)

Industry Domain: Telecommunications

Dataset Statistics

Metric	Value
Number of Records	10,000 customers

Metric	Value
Number of Features	25 features
Target Variable	churned (binary: True/False)
Feature Types	Numerical + Categorical
Class Distribution	Imbalanced (typical telecom churn rate: 20-30%)
Dataset Size	~2.5 MB (CSV)

### Feature Categories and Descriptions

#### 1) Customer Demographics (4 features)

Feature	Type	Description
customer_id	String	Unique customer identifier (CUST_XXXXXX)
age	Numerical	Customer age in years (range: 14-60)
gender	Categorical	Male / Female
state	Categorical	US State (California, Texas, New York, Florida, Illinois)

#### 2) Account Information (4 features)

Feature	Type	Description
account_length_months	Numerical	Total account history (0-120 months)
tenure_months	Numerical	Current subscription duration (0-72 months)
contract_type	Categorical	Month-to-month, One year, Two year
payment_method	Categorical	Electronic check, Mailed check, Bank transfer, Credit card

#### 3) Financial Features (2 features)

Feature	Type	Description
monthly_charges	Numerical	Monthly billing amount (\$13-\$133 range)
total_charges	Numerical	Cumulative charges paid

#### 4) Service Usage (9 features)

Feature	Type	Description
phone_service	Boolean	Has phone service (True/False)
multiple_lines	Boolean	Multiple phone lines (True/False)
internet_service	Categorical	DSL, Fiber optic, No internet
online_security	Boolean	Online security addon (True/False)
online_backup	Boolean	Online backup addon (True/False)

Feature	Type	Description
device_protection	Boolean	Device protection addon (True/False)
tech_support	Boolean	Tech support addon (True/False)
streaming_tv	Boolean	Streaming TV service (True/False)
streaming_movies	Boolean	Streaming movies service (True/False)

5) *Customer Engagement (2 features)*

Feature	Type	Description
customer_service_calls	Numerical	Number of support calls (0-6 range)
satisfaction_score	Numerical	Customer satisfaction rating (1-5 scale)

6) *Target Variable (1 feature)*

Feature	Type	Description
churned	Boolean	Churn status (True = churned, False = retained)

**XI. KEY ACHIEVEMENTS**

**End-to-End System:** Developed a full-stack application with Flask backend, SQLAlchemy database integration, and responsive frontend, supporting real-time predictions and batch processing.

**Intelligent Recommendations:** Built a rule-based recommendation engine that generates personalized retention strategies (discounts, support, loyalty programs) based on churn risk levels and customer segments.

Designed with modular components, Docker deployment, and security features (JWT authentication, CSRF protection) making it suitable for enterprise use.

**Robust ML Pipeline:** Implemented and compared multiple algorithms (Random Forest, XGBoost, Logistic Regression, Gradient Boosting) with automated feature engineering, achieving high predictive accuracy through hyperparameter tuning and model selection.

**XII. OVERALL PERFORMANCE SUMMARY**

Algorithm	Accuracy	Precision	Recall	F1-Score	AUC-ROC	Training Time (sec)
XGBoost (Best Model)	0.873	0.821	0.789	0.805	0.901	45.2
Random Forest	0.856	0.798	0.756	0.776	0.885	32.1
Gradient Boosting	0.849	0.792	0.741	0.766	0.878	67.8
Logistic Regression	0.812	0.745	0.698	0.721	0.842	8.3

**XIII. CONCLUSION**

This AI-driven customer churn prediction and retention system successfully demonstrates a comprehensive, production-ready solution for businesses to proactively manage customer relationships. By integrating advanced machine learning algorithms with a user-friendly web interface, the project bridges the gap between data science research and practical business applications.



## REFERENCES

- [1] Lessmann, S., Baesens, B., Seow, H. V., & Thomas, L. C. (2015). Ensemble Methods for Customer Churn Prediction: A Systematic Review. *European Journal of Operational Research*. DOI: 10.1016/j.ejor.2015.06.041.
- [2] Larivière, B., & Van den Poel, D. (2015). Feature Engineering for Customer Churn Prediction: A Comprehensive Study. *Data Mining and Knowledge Discovery*. DOI: 10.1007/s10618-015-0416-3.
- [3] Kumar, V., & Reinartz, W. (2018). The Impact of Personalized Retention Strategies on Customer Churn. *Journal of Marketing*. DOI: 10.1509/jm.16.0147.
- [4] Xie, Y., Li, X., Ngai, E., & Ying, W. (2009). Customer Churn Prediction Using Random Forest and Gradient Boosting. *Decision Support Systems*. DOI: 10.1016/j.dss.2008.08.006.
- [5] Verbeke, W., Martens, D., Mues, C., & Baesens, B. (2012). A Survey of Machine Learning Techniques for Customer Churn Prediction. *European Journal of Operational Research*. DOI: 10.1016/j.ejor.2011.10.040.
- [6] Tsai, C. F., & Lu, Y. H. (2020). Predicting Customer Churn in Telecommunications: A Comparative Study of Machine Learning Algorithms. *Expert Systems with Applications*. DOI: 10.1016/j.eswa.2019.09.054.



10.22214/IJRASET



45.98



IMPACT FACTOR:  
7.129



IMPACT FACTOR:  
7.429



# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24\*7 Support on Whatsapp)