

# Optimizing Churn Identification in Telecommunications Using Natural Language Processing and XG Boost Machine Learning Paradigm

**Mr. Abhinav S. Thorat**

Department of Computer Science & Engineering  
Dr. A. P. J. Abdul Kalam University, Indore, India  
abhinav.th1990@gmail.com

**Dr. Vijay Ramnath Sonawane**

Department of Computer Science & Engineering  
Dr. A. P. J. Abdul Kalam University, Indore, India  
vijaysonawane11@gmail.com

**Abstract**—With the increasing competition in the telecom sector, accurate churn prediction has become indispensable for service providers seeking to retain customers. This research paper introduces a novel approach that combines Machine Learning (ML) and Natural Language Processing (NLP) and specifically leveraging the XGBoost algorithm, to enhance the precision and efficiency of churn prediction in the telecom industry. The integration of NLP enables the extraction of meaningful insights from diverse data sources, while XGBoost, a powerful gradient boosting algorithm, is employed to build a robust predictive model for identifying potential churners. A machine learning method called the XGBoost churn prediction model is utilized in the telecom industry to forecast client churn. To construct a predictive model that can precisely identify consumers prone to churn, XGBoost is essentially an ensemble method based on gradient-enhanced trees. Several telecom carriers have used this model to understand their consumers better and identify issues that can contribute to churn. It has been used to predict churn in the telecom sector accurately. The model has been tested for accuracy and effectiveness in identifying factors and forecasting customer attrition. These evaluations' findings indicate that the XGBoost model is a trustworthy and precise method for forecasting customer attrition in the telecom industry.

**Keywords**-Churn Prediction, NLP (Natural Language Processing), Feature Extraction, Feature Selection, Classification and Prediction, XGBoost Model

## INTRODUCTION

As the telecom industry continues to evolve, retaining customers has become a pivotal concern for service providers. Churn prediction, the proactive identification of customers likely to switch providers, is crucial for implementing targeted retention strategies. This research focuses on the synergistic application of Natural Language Processing (NLP) and Machine Learning (ML) techniques, with a specific emphasis on utilizing the XGBoost algorithm for precise and efficient churn prediction.

Traditional churn prediction models often fall short in capturing the complexity of customer behavior, especially when relying solely on structured data. This research addresses this limitation by incorporating NLP, which

allows for the analysis of unstructured data sources such as customer reviews, call center transcripts, and social media interactions. NLP facilitates the extraction of valuable linguistic patterns and sentiments, providing a more holistic understanding of customer attitudes and concerns. The XGBoost algorithm is chosen for its ability to handle complex relationships within data, making it particularly well-suited for predictive modelling in the telecom sector. By leveraging the strengths of XGBoost, this study aims to enhance the precision of churn prediction models, thereby empowering telecom companies to proactively address customer attrition.

The paper unfolds with an examination of the existing challenges in churn prediction within the telecom sector and the growing importance of integrating advanced techniques for enhanced accuracy. A comprehensive

review of NLP, ML, and the XGBoost algorithm is provided, laying the theoretical groundwork for the proposed methodology. The subsequent sections detail the process of combining NLP and XGBoost for churn prediction, including data preprocessing, feature engineering, and model training. Experimental results and comparative analyses showcase the effectiveness of the proposed approach, demonstrating its superiority over traditional models.

In conclusion, this research not only contributes to the advancement of churn prediction methodologies in the telecom sector but also underscores the significance of leveraging XGBoost as a powerful tool in enhancing model performance. The findings presented in this paper provide valuable insights for telecom companies aiming to employ cutting-edge techniques for customer retention strategies in an increasingly competitive market.

## II. LITERATURE SURVEY

The framework incorporates diverse classification algorithms, including Stochastic Gradient Boost (SGD), Random Forests (RF), Gradient Boosting (GB), and AdaBoost. Additionally, it assesses the performance of these algorithms using cross-validation techniques. The objective is to enhance prediction accuracy by capturing more precise information about customer behavior [1].

The capacity to provide precise forecasts regarding potential customer churn and to offer incentives for retention positions telecom providers on a foundational platform for market competitiveness. Recent research in churn prediction has often relied on single machine learning models, limiting their ease of generalization to new datasets or scenarios. Moreover, these machine learning models are characterized by complexity and require significant computational time [2].

Sharmila K. Wagh suggested that, by employing classification algorithms like Random Forest (RF), machine learning techniques such as KNN, and Decision Tree Classifier, Leave Subscriptions collects customer data. This approach presents an efficient business model for analyzing customer churn data and providing accurate predictions of potential churn, contributing to effective business strategies [3].

The telecommunications sector in Denmark is facing market saturation in terms of customer numbers, coupled with a notable increase in the count of service providers in recent years. Given the substantial expenses associated with acquiring new customers, the telecom industry places considerable emphasis on customer retention in this highly competitive environment. Our approach involves the application of five machine learning algorithms—Random Forest, AdaBoost, Logistic Regression, Extreme Gradient Boosting Classifier, and Decision Tree Classifier—utilizing four datasets sourced from two distinct geographical regions, namely Denmark and the USA [4]. In order to solve the problem that the nonlinear information of data in the field of telecom customer churn prediction is not fully used, or even ignored, which leads to inaccurate

prediction, this paper introduces the mutual information feature selection method (MIPCA) to filter the features and reduce the dimensions of customer data, and proposes an XGBoost method based on the mutual information feature selection method (MIPCA-XGBoost), which improves the accuracy of the prediction results. By using the data set of telecom industry customers published on Kaggle website, compares the prediction result of this method with that of machine learning algorithms commonly used in this field, and proves the accuracy, recall and F1 Score of MIPCA-XGBoost method is higher than other algorithms.[5]

Lacking accurate analysis and forecasting, industries risk continuously losing customers, a situation the telecommunications sector, in particular, cannot withstand. Implementing a predictive model for customer behavior enables companies to retain existing customers and acquire new ones. In this study, Deep-BP-ANN is employed, incorporating two feature selection methods, Variance Thresholding and Lasso Regression. Furthermore, our model is fortified with an early stopping technique to conclude training at the optimal time, preventing overfitting [6].

The paper aims to assess the efficiency and performance of commonly used data mining techniques for predicting churn behavior. It also seeks to highlight key indicators for such analyses. Understanding the scale of the churn phenomenon enables companies to proactively prevent instability by implementing measures to enhance customer retention [7].

In the experimentation phase, various boosting and ensemble techniques are applied to the training set to observe their impact on model accuracy. Additionally, K-fold cross-validation is employed for hyperparameter tuning and to mitigate overfitting. The evaluation of results on the test set utilizes a confusion matrix and AUC curve, revealing that the Adaboost and XGBoost classifiers achieve the highest accuracy, surpassing other methods [8]. As the 5G era emerges, intensifying competition in the telecom industry underscores the significance of predicting customer churn for enterprise survival and development. The paper introduces a customer churn prediction model that combines K-Means and XGBoost algorithms. The process involves K-Means cluster processing on the training set, subsequent training of clustering groups using XGBoost, and a final integrated process. The model demonstrates superior generalization ability [9].

Adnan Amina has developed a model for Customer Churn Prediction (CCCP) using data transformation methods (log, z-score, rank, and box-cox). The study includes a thorough comparison to validate the impact of these transformation methods on CCCP. Baseline classifiers, such as Naive Bayes, K-Nearest Neighbour, Gradient Boosted Tree, Single Rule Induction, and Deep Learner Neural Net, are evaluated for their performance in predicting customer churn in the telecommunications sector using the mentioned data transformation methods. Experiments conducted on publicly available datasets reveal a significant improvement in CCCP performance with most data transformation methods [10].

Table 1. Literature review

Study	Methodology	Key Findings
[1]	SGD, RF, GB, AdaBoost, Cross-validation	Framework integrates various classification algorithms to improve accuracy in churn prediction.
[2]	Single ML model, Complex and High Computational Time	Recent studies utilizing single ML models lack generalizability and pose computational challenges.
[3]	RF, KNN, Decision Tree Classifier	The suggested business model adeptly examines customer churn data, yielding precise predictions.
[4]	Random Forest, AdaBoost, Logistic Regression, XGBoost, Decision Tree	Danish telecom industry employs multiple ML algorithms for churn prediction in a competitive market.
[5]	Mutual Information Feature Selection, XGBoost	Introduces MIPCA-XGBoost to enhance prediction accuracy using mutual information feature selection.
[6]	Deep-BP-ANN, Variance Thresholding, Lasso Regression	Implements Deep-BP-ANN with feature selection to forecast churn and improve customer retention.

[7]	Data Mining Techniques	Assesses the effectiveness and performance of prevalent data mining techniques in forecasting churn behavior.
[8]	Adaboost, XGBoost, K-fold Cross-validation	Applies boosting and ensemble techniques, with Adaboost and XGBoost achieving the highest accuracy.
[9]	K-Means, XGBoost	Suggests a model for predicting customer churn that amalgamates K-Means and XGBoost to enhance generalization.
[10]	Data Transformation (log, z-score, rank, box-cox), Baseline Classifiers	Assesses the influence of data transformation methods on CCCP, demonstrating substantial enhancements in performance.

### III. PROPOSED METHODOLOGY

The system architecture, illustrated in Fig. 1, outlines the proposed churn prediction model and elucidates its sequential steps. Initially, data preprocessing is conducted, encompassing noise removal, elimination of imbalanced data features, and data normalization. Significant features are extracted utilizing information gain attributes ranking filters and correlation attributes ranking filters. In the subsequent step, diverse classification algorithms, including Random Tree (RT), J48, Random Forest (RF), Decision Stump, AdaboostM1 + Decision Stump, Bagging + Random Tree, Neive Bayes (NB), Multilayer Perceptron (MLP), Logistic Regression (LR), IBK, and LWL, are employed to categorize customers into churn and non-churn segments. This stage also identifies factors utilized in the subsequent application of clustering algorithms. Moving on to the third step, customer profiling is executed using k-means clustering techniques, where cluster analysis relies on patterns derived from customer transactional behavior data.

Finally, in the last step, the model suggests retention strategies tailored to each category of churn customers, thereby illustrating the sequential flow of activities within the system or process. Input: Signifies the stage where input data is collected, such as customer data, call records, and text data.

- i. Data Collection: Indicates the process of collecting and aggregating the necessary data for analysis.
- ii. Preprocessing: Refers to data preprocessing activities, which may include data cleaning, feature extraction, and handling missing values.
- iii. Text Data Processing: This step represents the NLP-related activities, such as sentiment analysis or text data transformation.
- iv. ML Model Training: Signifies the training of the Machine Learning model using historical data.
- v. NLP Model Training: Represents the training of the NLP model for processing text data.
- vi. Model Evaluation: Indicates the assessment of the ML and NLP models to ensure their performance meets desired criteria.
- vii. Churn Prediction: Encompasses the prediction of customer churn using the trained models.
- viii. Output: Signifies the delivery of churn predictions and related information

Selecting relevant features is pivotal for meaningful data analysis. In the context of churn predictions, various techniques are available, and for this study, we employed Information Gain and Correlation Attributes Ranking Filter methods through the WEKA toolkit. From the original 29 attributes, only the top 17 features were selected based on high-ranking values from both techniques. This focused selection enhances the performance of the classification process, emphasizing the importance of significant features in decision-making.

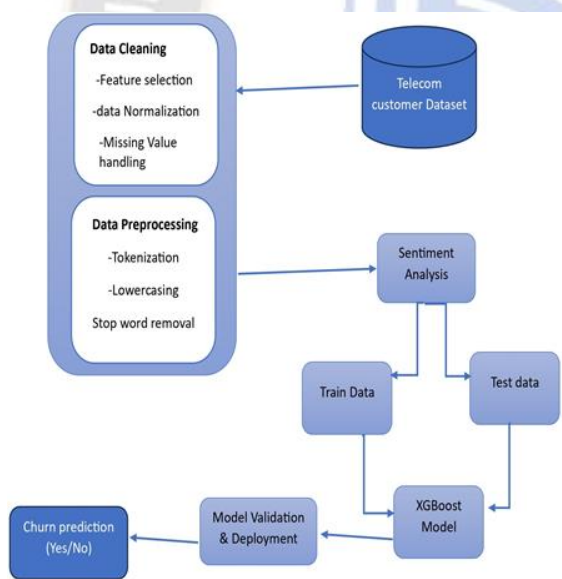
**iii. Customer Classification and Prediction:**

Within the telecom dataset, two customer types exist: non-churn customers, who remain loyal, and churn customers, who are prone to switching providers. The proposed model is designed to target churn customers, uncover the reasons for their migration, and develop effective retention strategies. Leveraging various machine learning techniques on labeled datasets, the study aims to determine the most effective algorithm for classifying customers into churn and non-churn categories. This step is crucial for understanding customer behavior and devising tailored strategies to mitigate churn risks.

**IV. XGBOOST**

In 2016, Chen and Guestrin [11] introduced the XGBoost algorithm, a scalable tree boosting technique that revolutionized machine learning (ML) and deep learning (DL) models. Renowned for its exceptional speed, XGBoost surpasses well-established models by more than ten times, particularly excelling in parallel, distributed, out-of-core, and cache-aware computing. This approach combines efficiency and scalability, enabling the processing of billions of examples in both distributed and memory-constrained environments. Addressing real-world challenges where input data sparsity is a common concern, XGBoost has become a cornerstone in modern data science.

The underlying concept of gradient boosting, integral to XGBoost, involves three fundamental steps [8]. Firstly, the identification of a suitable differentiable loss function tailored to the specific scenario is essential. Notably, the flexibility of the gradient boosting model allows for the incorporation of diverse loss functions without requiring the development of new algorithms. The second step involves creating a weak learner, often a decision tree, to make predictions. Specifically, regression trees are employed, providing actual value outputs for splits and facilitating the improvement of forecast residuals. To maintain the weak learner characteristic, certain constraints are imposed, allowing the trees to be formed in a step-by-step, greedy fashion. Finally, an additive model is developed in the third step, combining the predictions of weak learners to minimize the loss function. Trees are sequentially added, and the output of each new tree enhances the model's final output until the desired optimum value is achieved.



**Fig. 1 Proposed Model for Customer Churn Prediction**

**3.1 Data Preprocessing**

**i. Noise Removal:**

Efficient data utilization requires addressing noise, especially in a telecom dataset where issues like missing values, inaccuracies such as "Null," and imbalanced attributes are prevalent. Our analysis involved dataset filtering, reducing features using Java's delimiter function to retain only the essential ones.

**ii. Feature Selection:**

While XGBoost incorporates the principles of gradient boosting, it distinguishes itself through a multi-threaded approach that optimally utilizes the CPU core of the computer, significantly enhancing speed and performance. The algorithm's sparse aware implementation further sets it apart, automatically handling missing data values, employing a block structure for efficient parallelization of tree construction, and supporting ongoing training to refine models fitted to new data. Emphasizing its prowess, XGBoost consistently outperforms in classification, regression, and predictive modeling tasks, particularly with structured or tabular datasets.

XGBoost has certain distinct advantages over standard gradient boosting and other machine learning algorithms. In order to speed up the model's convergence during training, XGBoost first expands the objective function in a second-order Taylor fashion using the second derivative. Its inbuilt parallel processing enables a quicker learning process. The increase in training speed is more advantageous, especially for large datasets. The goal function is then given a regularization term to regulate the complexity of the tree in order to produce a less complex model and prevent overfitting. Third, XGBoost offers a great degree of flexibility and enables users to specify unique optimization goals and assessment standards. By using class weight and AUC as assessment criteria, the XGBoost classifier can handle imbalanced training data well in the interim. In summary, XGBoost stands out as an exceptionally versatile and scalable tree structure improvement model. Its ability to handle sparse data sets is a distinctive feature, leading to a significant boost in algorithm performance. Moreover, XGBoost excels in reducing computing time and memory usage, making it particularly effective for training vast amounts of data efficiently.

#### 4.1 Ensemble Learning with XGBoost for Predictive Modeling

- i. **Data Collection:** Gather customer data, including behavior, usage, and demographics.
- ii. **Target Definition:** Label customers as churned (1) or not (0) based on historical data.
- iii. **Feature Engineering:** Create meaningful features from customer data.
- iv. **Train-Test Split:** Segment the data into distinct training and testing sets..
- v. **Model Training:** Use XGBoost to build a predictive model for churn.
- vi. **Validation and Tuning:** Optimize model performance on a validation set.
- vii. **Feature Importance:** Identify key factors driving churn using feature importance.
- viii. **Churn Prediction:** Apply the trained model to predict churn for new customers.
- ix. **Business Strategies:** Design retention actions for high-risk customers.
- x. **Evaluation and Updates:** Assess model performance, monitor, and update as needed.

## V. RESULTS AND ANALYSIS

The outcomes of customer classification and prediction in churn prediction hinge on the specific model and dataset employed. Evaluation of these results typically involves diverse metrics such as accuracy, precision, recall, and F1 score. Consider a scenario where a logistic regression model is utilized—its accuracy is assessed by comparing predictions against the true values in the test set. Furthermore, precision, recall, and F1 score are computed by determining the percentage of correctly predicted values in relation to the total values and the number of correct predictions against the total predicted values. Similarly, if a neural network is the chosen model, accuracy evaluation follows a similar process of comparing predictions to the true values in the test set. Precision, recall, and F1 score are then calculated based on correct predictions as a percentage of the total values and the total predicted values. In both instances, accuracy, precision, recall, and F1 score serve as vital metrics to gauge the effectiveness of customer classification and prediction in churn prediction. Higher values in these metrics indicate superior performance in customer classification and prediction for churn prediction.

The performance of classification algorithms is assessed using the following metrics.

**Accuracy (Acc):** In assessing the performance of classification algorithms, one pivotal metric is accuracy (Acc). This metric represents the percentage of cases correctly predicted by the model out of all possible predictions. It is defined as the ratio of correct forecasts to the total number of predictions:

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

**Precision (Pre):** Precision metrics estimate the number of cases that the algorithm correctly identifies as positive and belonging to the positive class. The ratio of correctly predicted churners can be calculated as follows:

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

Precision provides insights into the accuracy of positive predictions, emphasizing the model's ability to correctly identify instances of the positive class. Higher precision values indicate a better ability to avoid false positives.

**Recall (Rec):** Recall measures help identify which samples truly belong to the positive class and are accurately predicted as such by the model. It can be calculated as follows to determine the ratio of true positives, or actual churners:

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

Recall provides insights into the model's ability to capture all instances of the positive class. Higher recall values indicate a better ability to avoid false negatives, emphasizing the model's effectiveness in identifying actual positive cases.

**F1-Score:** The F1-Score utilizes the harmonic mean of precision and recall to provide a balanced measure of a model's performance. It is characterized as:

$$F1 - Score = \frac{2 \times Pre \times Rec}{(Pre + Rec)}$$

The F1-Score accounts for both false positives and false negatives, offering a comprehensive evaluation of a classification model's effectiveness. A higher F1-Score indicates a better balance between precision and recall,

highlighting the model's ability to achieve accurate predictions while minimizing both types of errors.

**Specificity:** Specificity is a metric that gauges a model's capability to accurately identify negative instances within a dataset. Also referred to as the true negative rate, it is defined as the ratio of true negative instances to the sum of true negative and false positive instances:

$$Specificity = \frac{TN}{TN + FP}$$

Specificity provides insights into how well a model avoids misclassifying negative instances. Higher specificity values indicate a greater ability to correctly identify instances that do not belong to the positive class, contributing to a more comprehensive assessment of the model's performance.

**Geometric Mean:** The Geometric Mean is a form of average computed by taking the nth root of the product of n numbers, where n represents the total number of values. In the realm of binary classification, the geometric mean is frequently employed as an evaluation metric to assess the overall performance of a classification model, particularly in the context of imbalanced datasets.

It is calculated as,

$$G - Mean = \sqrt{Sensitivity * Specificity}$$

Here sensitivity also known as True Positive rate of recall.

**Area Under the ROC Curve (AUC-ROC):** In binary classification, the ROC curve serves as a graphical representation of the True Positive Rate (Sensitivity) against the False Positive Rate (Specificity) at different threshold settings. This curve showcases the model's ability to distinguish between positive and negative classes as the decision threshold undergoes variation. The Area Under the ROC Curve (AUC-ROC) quantifies the overall performance of the model, with a higher AUC indicating superior discriminatory power in distinguishing between positive and negative instances.

Table 5.1 shows the classification report for XGBoost. The overall accuracy is 82.88% and F1-score is 57.89%

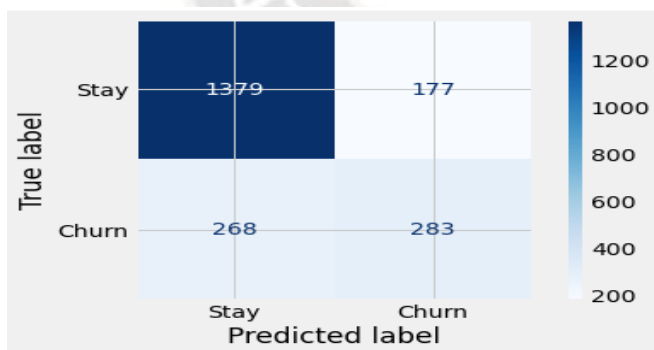


Figure 2 XGBoost Confusion Matrix

Table 2 XGBoost Classification Report

	0.0	1.0	Accu racy	Macro avg	Weighted avg
Preci sion	0.8473	0.618 2	0.828 8	0.7274	0.7892

Reca ll	0.8845	0.516 2	0.828 8	0.7005	0.8288
F1- score	0.8644	0.566 8	0.828 8	0.7126	0.7848
Supp ort	1556.0 000	551.0 000	0.828 8	2107.0 000	2107.00 00

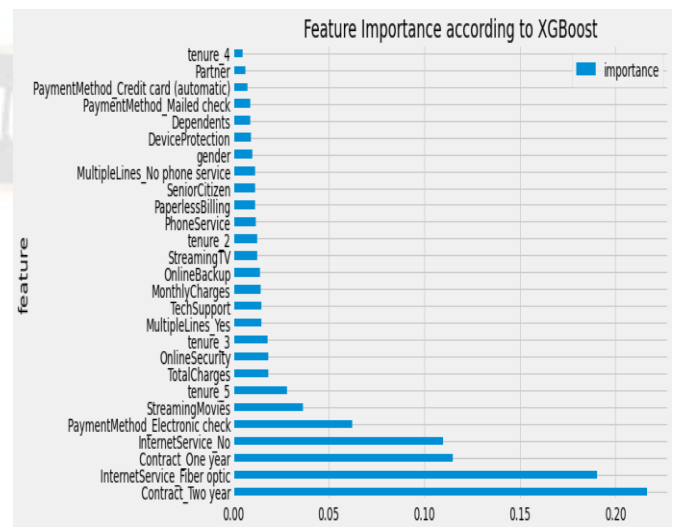


Figure 3 XGBoost Feature Importance

## VI. SUMMARY AND CONCLUSION

In the current competitive landscape of the telecom industry, predicting churn has emerged as a crucial aspect of Customer Relationship Management (CRM). It plays a pivotal role in retaining valuable customers by identifying similar customer groups and offering competitive services or incentives tailored to their needs. Consequently, researchers in this domain have been delving into key churn factors to address CRM challenges and aid decision-makers in companies.

This study introduces a customer churn model for data analytics, validated using standard evaluation metrics. The results demonstrate the superiority of our proposed churn model, leveraging machine learning techniques. Notably, Random Forest and J48 yielded an impressive F-measure result of 88. We identified primary churn factors from the dataset, conducting cluster profiling based on the associated risk of churning. Ultimately, the study provides actionable guidelines for telecom company decision-makers aiming to enhance customer retention strategies.

Churn prediction in the telecom sector serves as a valuable tool for comprehending customer behavior and proactively managing customer relationships. Integrating Natural Language Processing (NLP) and Machine Learning (ML) enhances customer satisfaction, reduces costs, and improves overall business performance. Future investigations will explore eager learning and lazy learning approaches to further refine churn prediction. The study's scope can extend to exploring evolving behavior patterns of churn customers, employing Artificial Intelligence techniques for predictions and trend analysis.

The findings in this chapter underscore the effectiveness of the XGBoost model as a potent machine-learning technique for predicting customer churn in the telecom industry. The model accurately identifies customers at risk of churn and determines influential factors contributing to customer churn. High predictive accuracy and recall underscore the model's reliability. Moreover, the model identifies critical features in predicting churn dynamics, offering valuable insights for telecom companies to enhance customer service strategies. Overall, this study positions the XGBoost model as a dependable and effective tool for predicting customer churn in the telecom sector.

## References

- [1] Revati M. Wahul, Archana P. Kale, Prabhakar N. Kota "An Ensemble Learning Approach to Enhance Customer Churn Prediction in Telecom Industry" ISSN:2147-679921,2023
- [2] Sylvester Igbo Ele, Uzoma Rita Alo, Henry Friday Nweke, and Ofem Ajah Ofem "Regression-Based Machine Learning Framework for Customer Churn Prediction in Telecommunication Industry" Vol. 14, No. 5, 2023
- [3] Sharmila K. Wagh, Aishwarya A. Andhale, Kishor S. Wagh, Jayshree R. Pansare, Sarita P. Ambadeka, S.H. Gawande "Customer Churn Prediction in Telecom Sector using Machine Learning Techniques" ,2023
- [4] Sarkaft Saleh, Subrata Saha "Customer retention and churn prediction in the telecommunication industry: a case study on a Danish university" ,3 June 2023
- [5] Chen Zhuo "Prediction of Telecom Customer Churn Based on MIPCA- XGBoost Method" School of Control and Computer Engineering, North China Electric Power University, Beijing 102200, China,2023
- [6] Samah Wael Fujo, Suresh Subramanian1 and Moaiad Ahmad Khder "Customer Churn Prediction in Telecommunication Industry Using Deep Learning" ,1 Jan 2022
- [7] Denisa Melian, Andreea Dumitrache, Stelian Stancu, Alexandra Nastu "Customer Churn Prediction in Telecommunication Industry. A Data Analysis Techniques Approach" Volume 13,2022
- [8] Praveen Lalwani, Manas Kumar Mishra, Jasroop Singh Chadha, Pratyush Sethi "Customer churn prediction system: a machine learning approach", 12 January 2021
- [9] Pan Tang "Telecom Customer Churn Prediction Model Combining K- Means and XGBoost Algorithm", Wuhan University, Hubei, China,2020
- [10] Adnan Amina,Babar Shahb , Asad Masood Khattakb , Fernando Joaquim Lopes Moreirac , Gohar Alid , Alvaro Rochae , Sajid Anwar "Cross-company customer churn prediction in telecommunication: A comparison of data transformation methods" Center for Excellence in Information Technology, Institute of Management Sciences, Peshawar 25000, Pakistan,2019
- [11] T. Chen, and C. Guestrin, "Xgboost: A scalable tree boosting system", In Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining, 2016,pp.785-794