

Telecom Churn Prediction using Machine Learning

**Pooja Jagannath Biradar, Mr. Gaurav Vijay Barde,
Mrs. Swati Shankar Dharbale, Mrs. Asmita Nilesh Boraste**
Gopinathji Munde Institute of Engineering and Research, Nashik, India

Abstract: *Telecom customer competition forecasting is an important task for telecom companies to retain customers. Churn is when a customer cancels their subscription or service from a communications company. Predicting customer churn helps telemarketing companies take steps to retain customers by identifying potential churn and providing effective retention strategies for them. This summary explains the context of the communication problem using machine learning. Contacts for problem prediction may include analysis of a customer's historical data, including demographic data, usage patterns, payment details, and service history, to predict whether a customer will leave in the future. Machine learning algorithms are used to learn patterns and relationships from this data and make predictions based on new, unseen data. Telecommunications customer churn prediction using machine learning involves processing customer history data, architectural design, selecting appropriate machine learning algorithms, effectively evaluating performance models using multiple metrics, and using best practices in a production environment. By using these methods, communications companies can reduce customer turnover and increase customer satisfaction.*

Keywords: Machine Learning, Random Forest, Decision Tree, XGBoost, Prediction, Churn

I. INTRODUCTION

Customer retention in today's highly competitive telecommunications industry is a challenge for service providers. Customer churn, which is when customers cancel their subscriptions or services, is a major concern for telecommunications companies as it leads to loss of revenue and operations. To overcome this challenge, mobile companies are turning to machine learning technology to predict customer competition and take proactive steps to protect against it. Gaming is important for phone companies because it helps them retain customers and reduce customer churn rates. Telecommunications companies can prevent customer churn and increase customer satisfaction by identifying potential customer churn and offering them attractive retention strategies. Additionally, churn forecasting allows telephone companies to improve their marketing and sales by focusing on valuable customers and without affecting users. Telecom churn forecasting is an important task for telecom companies because it helps them retain existing customers and reduce customer churn. Customer loss. Churn is the process when a customer stops providing service to a communications company. Using machine learning algorithms to predict customers can help phone companies identify factors that drive customer churn and take proactive steps to retain customers. Telecom churn forecasting is an important task for telecom companies because it helps them retain existing customers and reduce customer churn. Mobile marketing companies can use machine learning algorithms to create accurate and effective churn prediction models. These models can be integrated into the production environment to provide instant predictions and help companies take proactive steps to retain customers.

II. RELATED WORK

Weijie Yu, Weinan Weng proposed system aims to identify affecting customer churn and construct an efficient model, which is used to predict and analyze data with visualization results. The churn forecast consists of several phases: data preprocessing, data analysis, evaluation measure, and application of machine learning algorithms. Moreover, data pre-processing covers data cleaning, transformation, and classification. The machine learning

classifiers selected are Logistic Regression, SVM, Random Forest, AdaBoost, GBDT, XGBoost, Light GBM, and CatBoost. Classifiers were evaluated using performance measures, such as accuracy, precision, recall, AUC, and F1-Score. Based on the paper, the result was shown that the Light GBM outperformed other classifiers while identifying potential churners. [1]

Dr. O. Rama Devi, Sai Krishna Pothini proposed a model focuses on individuals who utilize paid OTT platforms for streaming video content on any device. The study used a questionnaire to gather data from participants of all demographics. The collected data underwent various pre- processing steps to make it suitable for machine learning models. The goal of predicting subscriptions for OTT (Over-The-Top) platforms using machine learning is to devise a model which can accurately predict whether a customer will continue using this platform or not. This information is important for OTT companies to understand and optimize their marketing and retention efforts. Relevant data, such as customer demographics and viewing habits, is collected and analyzed to train the model. This process involves cleaning the data, selecting important features, and training a machine learning model. The model is then tested and validated using performance metrics. In short, this problem requires a comprehensive understanding of customer behavior and the use of machine learning to predict subscription decisions. The results can provide valuable insights for OTT companies to improve their customer understanding and retention efforts. [2]

QiuYing Chen, Sang-Joon Lee, proposed system used Orange3 software to construct a customer churn prediction model for delivery platforms. The most effective Gradient Boosting algorithm was chosen to study customer churn prediction on the takeaway delivery platform. The predictive models of the Gradient Boosting algorithm show efficient and accurate results that are relatively easy to approach. In addition, unlike the results of general mechanical learning techniques, it also exhibits key characteristics that make the implementation of gradient enhancement techniques more effective. Especially as with ecommerce, it is more effective to implement incremental enhancement techniques to predict non- contractual customer churn. [3]

Gavril et al. presented an advanced methodology of data mining to predict churn for prepaid customers using dataset for call details of 3333 customers with 21 features, and a dependent churn parameter with two values: Yes/No. Some features include information about the number of incoming and outgoing messages and voicemail for each customer. The author applied principal component analysis algorithm "PCA" to reduce data dimensions. Three machine learning algorithms were used: Neural Networks, Support Vector Machine, and Bayes Networks to predict churn factor. The author used AUC to measure the performance of the algorithms. The AUC values were 99.10%, 99.55% and 99.70% for Bayes Networks, Neural networks and support vector machine, respectively. The dataset used in this study is small and no missing values existed. [4]

He et al. proposed a model for prediction based on the Neural Network algorithm in order to solve the problem of customer churn in a large Chinese telecom company which contains about 5.23 million customers. The prediction accuracy standard was the overall accuracy rate, and reached 91.1%. [5]

Idris proposed an approach based on genetic programming with AdaBoost to model the churn problem in telecommunications. The model was tested on two standard data sets. One by Orange Telecom and the other by cell2cell, with 89% accuracy for the cell2cell dataset and 63% for the other one. [6]

Huang et al. studied the problem of customer churn in the big data platform. The goal of the researchers was to prove that big data greatly enhance the process of predicting the churn depending on the volume, variety, and velocity of the data. Dealing with data from the Operation Support department and Business Support department at China's largest telecommunications company needed a big data platform to engineer the fractures. Random Forest algorithm was used and evaluated using AUC. [7]

This paper explains our work on subscriber churn analysis and prediction for such services. We work on data mining techniques to accurately and efficiently predict subscribers who will change-and- turn (churn) to another provider for the same or similar service. The dataset we use is a public and real dataset compiled by Orange Telecom for the KDD 2009 Competition. Number of teams achieved high scores on this dataset requiring a significant amount of computing resources. We are aiming to find alternative methods that can match or improve the recorded high scores with more efficient and practical use of resources. In this study, we focus on ensemble of meta-classifiers which have been studied individually and chosen according to their performances. [8]

Shin-Yuan Hung a proposed different technique to build predictive models for telecom churn prediction. We included customer service and customer complaint log for modeling, as suggestions from prior research of Wei and Chiu (2002). We examined the impact of inadequate data on model building. Our empirical evaluation shows that data mining techniques can effectively assist telecom service providers to make more accurate churning prediction. Zhang Y presents a hybrid approach for building a binary classifier. The approach is the combination of the k-nearest neighbor algorithm, handling separately m 1-dimensional data sets divided from a data set in m-dimension, and the logistic regression method. This hybrid KNN-LR classifier improves the performance of the logistic regression in classification accuracy in some situations where the predictor and target variables exhibit complex nonlinear relationships. The results of the experiment on four benchmark data sets show the proposed approach compares favorably with the well-known classification algorithms such as C4.5 and RBF. Furthermore, its effectiveness is illustrated by its application in customer churn prediction based on real-world customer data sets.[10]

III. OBJECTIVE & SCOPE OF PROPOSED SYSTEM

The purpose of this system is to create a machine learning model that can predict customer competition in the telecommunications industry.

Churn refers to the phenomenon where customers cancel their subscriptions or services from telecommunication companies. Predicting customer churn is important for phone companies because it allows them to take critical steps to retain customers, reduce churn rates, and increase customer satisfaction

Investigating customer churn prediction in communication using machine learning.

Analyzing the impact of customer churn on overall marketing communications Article

Discuss the importance of churn model in marketing communications Article

A better way to reduce the customer churn rate of telecommunication companies

The scope of this proposal includes collection and prioritization of customer data, architectural design, selection and training of appropriate learning systems, assessment and review of operating standards. The system will also include techniques such as data mining, cluster learning, and hyperparameter tuning to improve model accuracy and robustness.

The proposed system will be developed using Python and popular machine learning libraries. The system will be trained and tested on a large and diverse dataset of mobile users obtained from well-known products such as Kaggle. The system will be designed to be robust, efficient and easy to install in a production environment.

IV. FEATURES OF PROJECT

- Real-time monitoring
- Fraud detection
- Historical Customer Data
- Customer demographics
- Customer service interactions
- Network performance metrics
- Competitor information
- Social media sentiment analysis
- Customer satisfaction surveys/NPS scores

V. REPRESENTATION OF THE METHODOLOGY

The basic model for predicting future customer churn is data from the past. We look at data from customers that already have churned (response) and their characteristics / behavior (predictors) before the churn happened. The dataset contains demographic details of customers, their total charges and they type of service they receive from the company. It comprises of churn data of over customers spread over 21 attributes obtained from Kaggle. By

fitting statistical models that relate the predictors to the response, we will try to predict the response for existing customers. This method belongs to the supervised learning category.

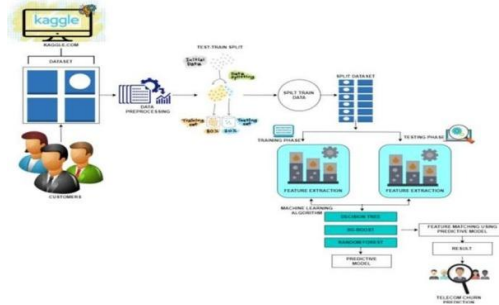


Fig 1: Representation of the Methodology

VI. ADVANTAGES

Creating new opportunities for cross-selling and upselling.

- **Improved Accuracy:** Machine learning algorithms can analyze large volumes of customer data and identify patterns and trends that are not easily detectable by traditional methods. This helps in making more accurate predictions about customer churn, which can help telecom companies take proactive measures to retain their customers.
- **Real-time Predictions:** Machine learning algorithms can make real-time predictions about customer churn, which allows telecom companies to take immediate action to address the issue. This can help in reducing churn rates and improving customer satisfaction.
- **Personalized Solutions:** Machine learning algorithms can analyze individual customer data and provide personalized solutions to address their specific needs. This can help in improving customer satisfaction and reducing churn rates.
- **Cost-effective:** Machine learning algorithms can help in reducing the cost of customer acquisition by predicting which customers are likely to churn and targeting them with retention offers. This can help in reducing the overall churn rate and improving the profitability of the telecom company.
- **Enhanced Customer Experience:** Machine learning algorithms can help in improving the overall customer experience by providing personalized solutions and addressing their specific needs. This can help in improving customer satisfaction and reducing churn rates.
- **Competitive Advantage:** Telecom companies that adopt machine learning algorithms for churn prediction can gain a competitive advantage over their competitors by providing better customer service and reducing churn rates. This can help in improving the overall market share and profitability of the company.
- **Continuous Learning:** Machine learning algorithms can continuously learn and improve over time as they are exposed to more customer data. This helps in making more accurate predictions about customer churn and improving the overall performance of the model.
- **Scalability:** Machine learning algorithms can handle large volumes of customer data and scale to accommodate the growing needs of the telecom company. This helps in improving the overall efficiency and effectiveness of the churn prediction process.

VII. APPLICATION AREA

Churn prediction can help you see which customers are about to leave your service so you can develop proper strategy to re-engage them before it is too late.

- **Customer Retention:** Telecom companies can use machine learning algorithms for churn prediction to identify customers who are likely to churn and take proactive measures to retain them. This can include offering personalized solutions, discounts, and promotions to address their specific needs.

- **Network Optimization:** Telecom companies can use machine learning algorithms for churn prediction to identify areas with high churn rates and optimize their network infrastructure to improve the overall customer experience. This can include improving network coverage, reducing network latency, and enhancing network reliability.
- **Pricing Strategy:** Telecom companies can use machine learning algorithms for churn prediction to identify the optimal pricing strategy for their products and services. This can include offering different pricing plans based on customer usage patterns and preferences.
- **Marketing Campaigns:** Telecom companies can use machine learning algorithms for churn prediction to identify the most effective marketing campaigns for retaining customers. This can include targeting customers with personalized offers and promotions based on their specific needs and preferences.
- **Sales Forecasting:** Telecom companies can use machine learning algorithms for churn prediction to forecast sales and revenue based on customer churn rates. This can help in making informed business decisions and improving the overall financial performance of the company.
- **Resource Allocation:** Telecom companies can use machine learning algorithms for churn prediction to allocate resources such as customer service representatives, marketing personnel, and sales teams based on customer churn rates. This can help in optimizing resource allocation and improving the overall efficiency and effectiveness of the company.
- **Network Security:** Telecom companies can use machine learning algorithms for churn prediction to identify potential security threats and vulnerabilities in their network infrastructure. This can include identifying unusual network activity, detecting network intrusions, and preventing network attacks.
- **Network Maintenance:** Telecom companies can use machine learning algorithms for churn prediction to identify potential network failures and outages before they occur. This can include predicting network faults, identifying network bottlenecks, and optimizing network performance.
- **Network Expansion:** Telecom companies can use machine learning algorithms for churn prediction to identify potential areas for network expansion and growth. This can include identifying high-demand areas, predicting network usage patterns, and optimizing network capacity.
- **Network Optimization:** Telecom companies can use machine learning algorithms for churn prediction to optimize their network infrastructure based on customer usage patterns and preferences. This can include improving network coverage, reducing network latency, and enhancing network reliability.

VIII. ANALYSIS OF EXPERIMENTAL RESULTS

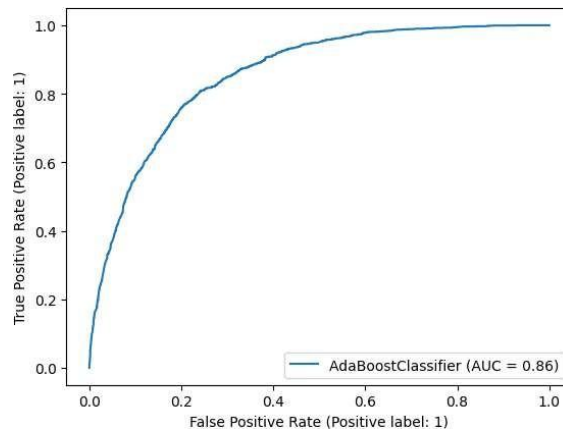


Fig 1: AUCROC ADB for training

AUCROC The AUCROC (Area Under the Receiver Operating Characteristic Curve) ADB (AdaBoost) algorithm is a popular ensemble learning method that combines multiple weak learners to create a strong classifier. During training, the ADB algorithm iteratively improves the classification performance by assigning higher weights to the misclassified samples in each iteration. The AUCROC metric is used to evaluate the performance of the ADB model during training. It measures the trade-off between the true positive rate and the false positive rate, providing a comprehensive assessment of the model's ability to distinguish between different classes. As the ADB algorithm progresses through the iterations, the AUCROC value is calculated at each step to monitor the model's performance and guide the learning process. The goal is to maximize the AUCROC value by adjusting the weights of the weak learners and optimizing the decision boundaries to improve the overall classification accuracy. the AUCROC ADB algorithm is an effective and robust method for training classification models, especially when dealing with imbalanced datasets or complex decision boundaries. It provides a reliable measure of the model's performance and guides the training process towards better predictive accuracy.

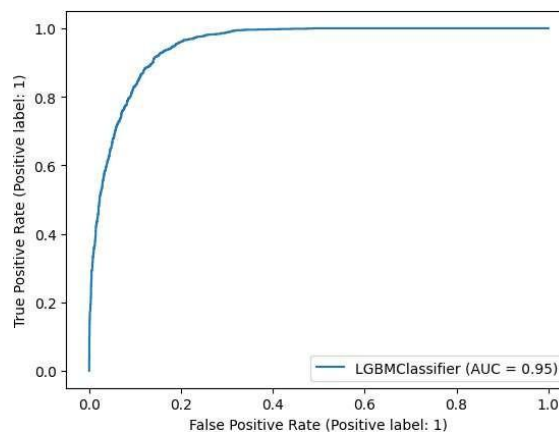


Fig 2: AUCROC LGBM Training

AUCROC (Area Under the Receiver Operating Characteristic curve) is a performance metric commonly used to evaluate the predictive power of a classification model, such as LightGBM (LGBM). Training LGBM with AUCROC involves optimizing the hyperparameters of the model in order to maximize the AUCROC score, which indicates how well the model can distinguish between different classes. During training, the LGBM algorithm is applied to a training data set to learn the patterns that differentiate between the classes. The algorithm iteratively adjusts the model parameters in order to minimize the loss function and improve its predictive accuracy. The AUCROC score is calculated by measuring the area under the Receiver Operating Characteristic curve, which is a plot of the true positive rate against the false positive rate at various threshold values. By training LGBM with AUCROC, we aim to create a high-performing classification model that can accurately predict the class labels of new data samples. This process helps to ensure that the model makes informed decisions and generalizes well to unseen data, ultimately leading to more reliable and robust predictions.

	precision	recall	f1-score	support
0	0.84	0.90	0.87	1291
1	0.66	0.54	0.59	467
accuracy			0.80	1758
macro avg	0.75	0.72	0.73	1758
weighted avg	0.79	0.80	0.80	1758

Fig 3: Classification Report on LR testing Classification Report is a performance evaluation metric used in machine learning for classification models. It provides a breakdown of key metrics such as precision, recall, F1 score, and support for each class in the classification model. For a logistic regression (LR) model, the Classification Report would show the precision, recall, F1 score, and support for each class predicted by the

model. Precision measures the proportion of true positive predictions among all positive predictions, recall measures the proportion of true positive predictions among all actual positives, and the F1 score is the harmonic mean of precision and recall. The support is the number of occurrences of each class in the dataset. A Classification Report on LR testing would provide insights into the performance of the LR model in classifying different classes. By analyzing the precision, recall, and F1 score for each class, you can determine how well the LR model is able to correctly classify instances of each class. This report can help you identify any potential issues or areas for improvement in the LR model's performance.

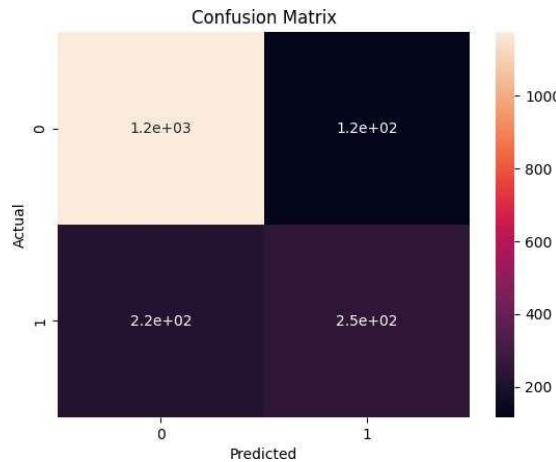


Fig 4: CM Testing

Configuration management (CM) testing is a process that ensures that the changes made to a software or system configuration are properly documented, tested, and verified before being implemented. This testing process evaluates the impact of configuration changes on the system's functionality, performance, and security. It helps to prevent issues such as configuration conflicts, system instability, and security vulnerabilities. CM testing also helps to ensure that the changes are properly tracked, managed, and deployed in a controlled and systematic manner. Overall, CM testing plays a crucial role in maintaining the integrity and reliability of systems and software configurations

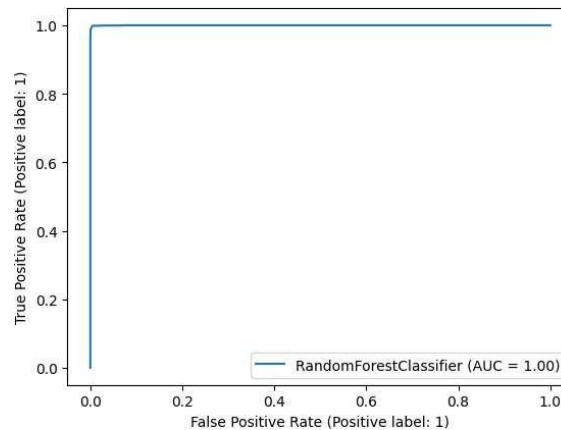


Fig 5: AUCROC RF Training

AUCROC RF Training is a machine learning training process that utilizes the Random Forest algorithm to evaluate the performance of a model through the area under the receiver operating characteristic curve (AUCROC). This training process involves building an ensemble of decision trees to classify data and calculate the AUCROC score, which is a measure of the model's ability to distinguish between classes. By training the model with the Random Forest algorithm, it can effectively handle large datasets with high-dimensional features and provide

accurate predictions for classification tasks. This training process helps improve the model's performance and enhance its ability to make informed decisions based on the input data.

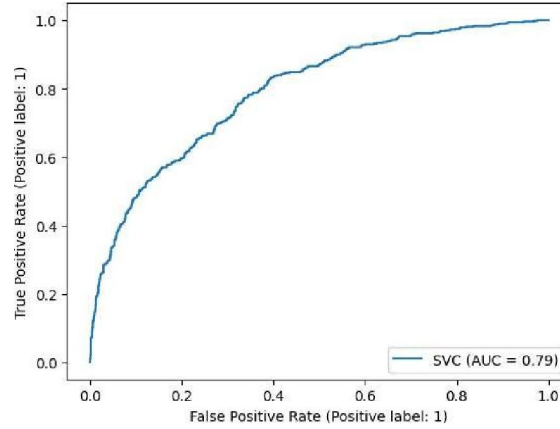
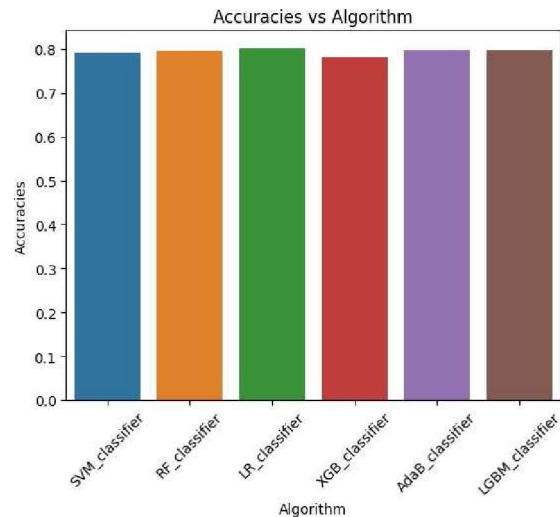


Fig 6: AUC ROC SVM Testing

Testing the performance of a Support Vector Machine (SVM) using the Area Under the Receiver Operating Characteristic Curve (AUC ROC) involves evaluating how well the model is able to distinguish between different classes or categories. The AUC ROC curve is a graphical representation that shows the trade-off between the true positive rate and the false positive rate of the model across different threshold values. The AUC ROC score is a single metric that summarizes the performance of the model, with a higher score indicating better discrimination between classes. During testing, the SVM model is trained on a labeled dataset and then tested on a separate set of data to evaluate its ability to classify new instances accurately. The AUC ROC score is calculated based on the predictions made by the model on the test set, comparing the predicted probabilities with the actual labels. By analyzing the AUC ROC score, we can assess the overall performance of the SVM model in terms of its ability to correctly classify instances and differentiate between classes. This information can help us understand how well the SVM model is performing and identify any areas for improvement.



Fi 7: Accuracy Comparison Testing

Accuracy comparison testing is a process in which multiple products, systems, or methods are tested to determine their accuracy and reliability in achieving a specific outcome. This type of testing involves evaluating the performance of each product in a controlled environment and comparing the results to determine which one is the most accurate. During accuracy comparison testing, various metrics are used to measure the performance of each

product, such as precision, sensitivity, specificity, and error rate. These metrics help to determine the reliability and consistency of each product in producing accurate results. Accuracy comparison testing is commonly used in industries such as healthcare, manufacturing, and technology to evaluate the performance of different tools, equipment, or software. By conducting this type of testing, organizations can make informed decisions about which product to use based on their accuracy and reliability in achieving the desired outcome

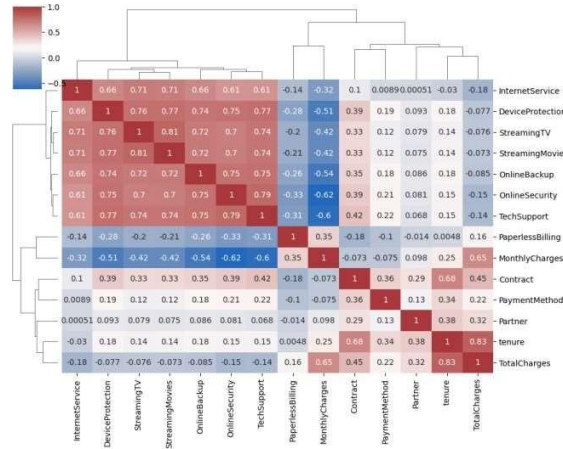


Fig 8: correlation of features

Correlation of features is a statistical measure that describes the strength and direction of a relationship between two or more features in a dataset. A correlation value can range from -1 to 1, with 0 indicating no correlation, 1 indicating a perfect positive correlation, and -1 indicating a perfect negative correlation. Correlation analysis helps to understand how features are related to each other and can be used to identify patterns and relationships within the data. High correlation between features may indicate redundancy or multicollinearity, which can affect the performance of machine learning models. On the other hand, low or moderate correlation may suggest independent or complementary features that provide valuable information for prediction or classification tasks. Correlation of features plays a crucial role in data analysis and model building, helping to identify important relationships and optimize the selection of features for predictive modeling.

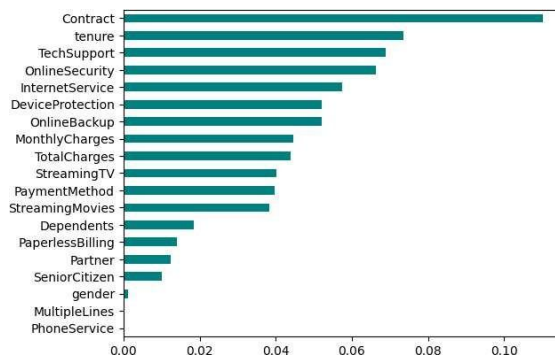


Fig 9: feature Selection using SelectKbest SelectKBest is a feature selection method in machine learning that selects the k best features from a dataset based on their individual scores. The SelectKBest method evaluates the importance of each feature through a statistical test and assigns a score to each feature. The SelectKBest method can be applied to both classification and regression problems. It helps improve the performance of machine learning models by selecting only the most relevant features and reducing the dimensionality of the dataset. This can lead to better model accuracy, faster training times, and improved interpretability of the results. SelectKBest works by ranking the features based on their scores and then selecting

the top k features with the highest scores. It allows you to specify the number of features to select (k) and the scoring function to use. Common scoring functions include chi-squared for classification tasks and regression for regression tasks.

SelectKBest is a powerful feature selection technique that helps in identifying the most important features in a dataset and improving the overall performance of machine learning models

IX. CONCLUSION

In conclusion, churn prediction is a critical challenge for telecom companies as it affects revenue, market share, and customer satisfaction. Machine learning techniques offer a promising solution to this problem by enabling telecom companies to predict churn accurately and take proactive measures to prevent it. The challenges involved in churn prediction include dealing with large volumes of historical customer data, handling noisy, incomplete, and highly dimensional data, and addressing imbalanced datasets. Various machine learning algorithms such as Random Forest, Decision Tree, XGBoost . can be used to solve this problem depending on the nature of the data and the specific requirements of the telecom company. By implementing machine learning algorithms for churn prediction, telecom companies can reduce churn rates, improve customer satisfaction, and gain a competitive advantage in the market.

REFERENCES

- [1]. Weijie Yu, Weinan Weng,” Customer Churn Prediction Based on Machine Learning” 2022 4th International Conference on Artificial Intelligence and Advanced Manufacturing (AIAM), 978-1- 6654-6399-7/22/\$31.00 ©2022 IEEE, DOI 10.1109/AIAM57466.2022.00176
- [2]. Dr. O. Rama Devi, Sai Krishna Pothini,” Customer Churn Prediction using Machine Learning: Subscription Renewal on OTT Platforms”, IEEE Xplore Part Number: CFP23BC3-ART; ISBN: 978-1-6654-5630-2, 978-1-6654-5630-2/23/\$31.00 ©2023 IEEE
- [3]. QiuYing Chen, Sang-Joon Lee,” A Machine Learning Approach to Predict Customer Churn of a Delivery Platform”, 2023 International Conference on Artificial Intelligence in Information and Communication (ICAIIIC) | 978-1-6654-5645-6/23/\$31.00 ©2023 IEEE | DOI:10.1109/ICAIIIC57133.2023.10067108
- [4]. Brandusoiu I, Todorean G, Ha B. Methods for churn prediction in the prepaid mobile telecommunications industry. In: International conference on communications. 2016. p. 97–100.
- [5]. He Y, He Z, Zhang D. A study on prediction of customer churn in fixed communication network based on data mining. In: Sixth international conference on fuzzy systems and knowledge discovery, vol. 1. 2009. p. 92–4.
- [6]. Idris A, Khan A, Lee YS. Genetic programming and adaboosting based churn prediction for telecom. In: IEEE international conference on systems, man, and cybernetics (SMC). 2012. p. 1328–32.
- [7]. Huang F, Zhu M, Yuan K, Deng EO. Telco churn prediction with big data. In: ACM SIGMOD international conference on management of data. 2015. p .607–18.
- [8]. Yabas, U, Chankya, H.C. (2013). Churn prediction in subscriber management for mobile and wireless communications services. IEEE Publications.
- [9]. Shin-Yuan Hung a, David C. Yen b, Hsiu-Yu Wang, “Applying data mining to telecom churn management”, Expert Systems with Applications 31 (2006) 515–524,
- [10]. Zhang, Y.; Qi, J.; Shu, H.; Cao, J. A hybrid KNN-LR classifier and its application in customer churn prediction. In Proceedings of the 2007 IEEE International Conference on Systems, Man and Cybernetics, Montréal, QC, Canada, 7– 10 October 2007; pp. 3265–3269