

5-2023

OTT SUBSCRIBER CHURN PREDICTION USING MACHINE LEARNING

Needhi Devan Senthil Kumar
California State University - San Bernardino

Follow this and additional works at: <https://scholarworks.lib.csusb.edu/etd>



Part of the [Other Computer Engineering Commons](#)

Recommended Citation

Senthil Kumar, Needhi Devan, "OTT SUBSCRIBER CHURN PREDICTION USING MACHINE LEARNING" (2023). *Electronic Theses, Projects, and Dissertations*. 1660.
<https://scholarworks.lib.csusb.edu/etd/1660>

This Project is brought to you for free and open access by the Office of Graduate Studies at CSUSB ScholarWorks. It has been accepted for inclusion in Electronic Theses, Projects, and Dissertations by an authorized administrator of CSUSB ScholarWorks. For more information, please contact scholarworks@csusb.edu.

OTT SUBSCRIBER CHURN PREDICTION
USING MACHINE LEARNING

A Project
Presented to the
Faculty of
California State University,
San Bernardino

In Partial Fulfillment
of the Requirements for the Degree
Master of Science
in
Information Systems and Technology

by
Needhi Devan Senthil Kumar

May 2023

OTT SUBSCRIBER CHURN PREDICTION
USING MACHINE LEARNING

A Project
Presented to the
Faculty of
California State University,
San Bernardino

by
Needhi Devan Senthil Kumar

May 2023

Approved by:

Dr. William Butler, Member, Committee Chair

Dr. Conrad Shayo, Member, Reader & Department Chair, Information and
Decision Sciences

© 2023 Needhi Devan Senthil Kumar

ABSTRACT

Subscriber churn is a critical issue for companies that rely on recurring revenue from subscription-based services like the OTT platform. Machine Learning algorithms can be used to predict churn and develop targeted retention strategies to address the specific needs and concerns of at-risk subscribers. The research questions are 1) What Machine Learning algorithms are used to overcome subscriber churn? 2) How to predict subscribers' churn in the OTT platform using Machine Learning? 3) How to retain subscribers and improve customer targeting? The dataset was collected from the Kaggle repository and implemented it into the various prediction algorithms used in previous research. Then, evaluate the performance of each algorithm to find out the highest accuracy model. The findings and conclusion for each question are 1) Logistic regression, multi-layer perceptron, random forest, decision trees, and gradient boosting machines were identified as effective algorithms for churn prediction analysis. 2) By sending the test data to a trained model by their historical dataset, customers are likely to leave a company (i.e., churn) based on their characteristics can be predicted. 3) Personalized offers and promotions, improving customer service, developing loyalty programs, and optimizing pricing strategies were suggested strategies for retaining subscribers. The gradient boosting machine model was found to have the highest accuracy and maximum AUROC, making it a powerful tool in the fight against customer churn. Areas for further study include incorporating unstructured data sources, deep learning

techniques, and integrating real-time data sources to improve the accuracy and effectiveness of churn prediction models.

Search Term: Subscribers, Churn Prediction, OTT Platform

ACKNOWLEDGEMENTS

I would like to thank my parents and my friends for their support and encouragement throughout the process of this project.

Also, I would like to thank Dr. William Butler and Dr. Conrad Shayo for guiding me to finish this project.

TABLE OF CONTENTS

ABSTRACT	iii
ACKNOWLEDGEMENTS	v
LIST OF TABLES	viii
LIST OF FIGURES	ix
CHAPTER ONE: INTRODUCTION	
Problem Statement	1
Research Questions	2
Organization of the Project	2
CHAPTER TWO: LITERATURE REVIEW	3
CHAPTER THREE: RESEARCH METHODOLOGY	
Logistic Regression	6
Multi-Layer Perceptron.....	7
Random Forest.....	7
Decision Tree.....	8
Gradient Boosting Machine.....	8
Data Collection	9
CHAPTER FOUR: DATA DESCRIPTION AND ANALYSIS	12
System Requirements	13
Implementation	13
Experimental Results	13

CHAPTER FIVE: DISCUSSION, CONCLUSION, AND AREAS OF FURTHER
STUDY

Discussion	18
Conclusion	19
Area for Further Study	20
REFERENCES	21

LIST OF TABLES

Table 1. Data Description	12
---------------------------------	----

LIST OF FIGURES

Figure 1. Classification Report of Logistics Regression.....	14
Figure 2. Classification Report of Multi-Layer Perceptron	14
Figure 3. Classification Report of Random Forest Classifier	15
Figure 4. Classification Report of Decision Tree Classifier	15
Figure 5. Classification Report of Gradient Boosting Machine	16
Figure 6. Values of AUROC.....	17

CHAPTER ONE

INTRODUCTION

OTT (Over the Top) is a term used to describe the delivery of audio, video, and other media content over the Internet without the need for a traditional cable or satellite television provider. OTT services provide users with an alternative source for content that is normally only available through a cable or satellite provider. They also offer a cost-effective way to access content, often eliminating the need for expensive contracts and monthly fees. OTT services typically require users to sign up for an account, either through a website or an app, and then they can access the content they desire. OTT services can provide a range of options, from streaming live TV channels to on-demand shows and movies. These services also often provide additional features, such as the ability to save favorite shows, get recommendations, and access exclusive content (Fitzgerald, 2019). Let's start by stating the problem statement, then go on to the research questions and how this culminating experience project is structured.

Problem Statement

Subscriber churn is a major issue for online streaming services such as Over-the-top (OTT) platforms (Madden et al., 1999b). Churn is the rate at which subscribers to a particular OTT service cancel their subscriptions, and it's a major part of customer retention for OTT services. Churn is a complex issue, with potential causes ranging from rising user-acquisition costs to a lack of personalized content. It is an important metric for OTT services, as it can indicate customer satisfaction, as well as the effectiveness of marketing and promotional efforts. A high churn rate can indicate that customers are not satisfied with the service, or that the marketing efforts

are not effective. It can also indicate that the service is too expensive, or that there are better alternatives available (Kuldeep et al., n.d.).

Research Questions

1. What Machine Learning algorithms are used to overcome subscriber churn?
2. How to predict subscribers' churn in the OTT platform using Machine Learning?
3. How to retain subscribers and improve customer targeting?

Organization of Project

By building an optimized model that can be used to predict subscribers' churn, the goal is to explore how Machine Learning algorithms can be used to decrease subscribers' churn in the OTT platform. This culminating experience project is organized as follows: Chapter 2 will describe the literature review of my research on the existing churn analysis of various industries. Chapter 3 will provide research methodologies. Chapter 4 will provide the data collection and implementation of research methodologies. Based on the experimental result, Chapter 5 is composed of a discussion, conclusion, and area for further study.

CHAPTER TWO

LITERATURE REVIEW

To learn more about this issue and come up with a workable solution, research and implementations done in the past by other authors were examined. Let's provide our findings in this section related to the research question: Q1. What Machine Learning algorithms are used to overcome subscriber churn?

Agrawal et al., (2018) discussed the problem of customer churn and analyzes previous works to identify gaps in the solutions implemented. Agrawal et al., (2018) proposed a multi-layered Artificial Neural Network model to predict customer churn, resulting in an accuracy of 80.03%. Ahmad et al., (2019) aimed to build a system to predict customer churn in the Syriatel telecom company. The XGBOOST algorithm, which had an AUC value of 93.301% produced the best results. XGBOOST continued to produce the highest results with an AUC of 89% when tested on a fresh dataset pertaining to various time periods. The outcomes of predicting churn in the telecom business have been proven to be improved by the usage of Social Network Analysis elements. Their research is more focused on Machine learning algorithms used in churn prediction in various industries, which is a part of our project but their research is less focused on subscriber churn on OTT platforms.

To answer the research questions: Q2. How to predict subscribers' churn in the OTT platform using Machine Learning?

De Caigny et al., (2018) aimed to explore the viability of using the Logistic Linear Model (LLM) as a classification technique in telecom customer churn prediction. The results showed that LLM performed better than using Logistic Regression (LR) and Decision Tree (DT) as standalone techniques and at least equally well as two homogeneous ensemble methods, Random Forest (RF) and

Logistic Model Tree (LMT). The LLM provides a comprehensible method with acting ability, and it can enrich both DT and LR by adding the coefficients of logistic regression to the leaves in a decision tree and by fitting several logistic regressions to consider specific group characteristics. Huang et al., (2012) introduced a brand-new feature set for predicting customer turnover in the telecom sector, which includes call details, account information, bill information, and other forms of data. However, the research examined how customer behavioral characteristic data is used to predict churn of the customers in financial industries, it lacks the methods in predicting subscriber churn in the OTT platform.

To answer the research questions: Q3. How to retain subscribers and improve customer targeting?

Ullah et al., (2019) aimed to build a churn prediction model for a telecom company to improve its CRM and retain valuable customers. They used machine learning techniques to analyze customer data and identify the main factors contributing to churn. The results showed that the proposed model performed better than other techniques and produced a better F-measure result of 88% using Random Forest and J48. Sung Won Kim et al., (2019) also performed cluster profiling to better understand the risk of churn for different groups of customers and provided guidelines for customer retention. Their research is more focused on the customer retention churn prediction model in telecom, but their research is less focused on subscriber churn on OTT platforms.

Therefore, the above reviews discuss different approaches to predicting customer churn. Agrawal et al., (2018) proposed a multi-layered Artificial Neural Network model that identified attributes related to churn rate and achieved an accuracy of 80.03%. Ahmad et al., (2019) applied feature engineering,

transformation, and selection to prepare data for four tree-based algorithms and found that the xgboost algorithm produced the best results. De Caigny et al., (2018) explored the viability of using the Logistic Linear Model as a classification technique and found that it performed better than other classifiers. Huang et al., (2012) presented a new feature set for churn prediction and evaluated seven modeling techniques, finding that the Logistic Regression and DT/SVM were suitable for predicting true churn rate and false churn rate. Finally, Ullah et al., (2019) used machine learning techniques to build a churn prediction model and achieved an F-measure result of 88% using Random Forest and J48. These studies highlight the importance of feature selection and modeling technique selection for effective churn prediction in the telecom, financial and social network industries and let's implement these methods into the OTT platform industry.

CHAPTER THREE

RESEARCH METHODOLOGY

Here, the research methodologies used to answer the research questions have been explained. To answer the research question: Q1. What Machine Learning algorithms are used to overcome subscriber churn? From the literature review the research come to the solution that Machine learning algorithms used to overcome churn in other industries are Logistics Regression, Multi-Layer Perceptron, Random Forest, Decision Trees, and Gradient Boosting Machines, let's implement these methods into the OTT platform industry to find subscriber churn. Let's introduce the machine learning algorithms that are used to predict subscriber churn analysis and list some of their benefits and drawbacks.

Logistic Regression

A supervised learning method used in machine learning for binary classification issues is logistic regression. The sigmoid function is used to predict the likelihood of a binary outcome by estimating the coefficients of the input variables. The implementation of logistic regression is straightforward, computationally effective, and appropriate for small datasets. To comprehend the underlying causes of the result, the coefficients can be understood. However, it is limited to binary classification issues, presumes the independence of the data, and presupposes a linear relationship between the input variables and the result. As a foundational model for forecasting binary outcomes, logistic regression is frequently employed in the fields of healthcare, finance, and marketing (Huang et al., 2012).

Multi-Layer Perceptron

Machine learning uses the Multi-Layer Perceptron (MLP), an artificial neural network, for a variety of tasks like classification, regression, and prediction. It has many layers of neurons, which combine the inputs in a weighted way and then pass them via an activation function. Using non-linear activation functions, MLP can manage relationships between inputs and outputs that are not linear. It is an effective tool for difficult classification and regression problems, but it can be computationally expensive and prone to overfitting for big datasets. MLP is widely utilized in numerous applications, including financial forecasting, speech recognition, image recognition, and natural language processing (Agrawal et al., 2018).

Random Forest

It is an ensemble learning technique that mixes the outputs of numerous decision trees to increase accuracy and lessen the chance of overfitting. A group of decision trees are trained in Random Forest utilizing random subsets of characteristics and data. The algorithm chooses the appropriate characteristic to divide on at each node depending on a parameter like information gain or Gini index. The outcome is then determined by which decision tree received the most votes (De Caigny et al., 2018).

In order to obtain high accuracy and lower the risk of overfitting, it integrates numerous decision trees. It is helpful for churn prediction studies where data quality may be an issue because it can manage missing data and is robust to outliers. Yet, compared to decision trees, Random Forests might be more computationally expensive and challenging to understand. Moreover, it can favor categorical variables with more levels and not be appropriate for small datasets (Ullah et al., 2019).

Decision Tree

A tree-like structure known as a decision tree algorithm was developed to depict decisions and their potential outcomes based on the characteristics of the data. The tree structure is produced by recursively dividing the data into subsets depending on the most significant predictors, and the algorithm selects the optimal feature to split on based on a criterion like information gain or Gini index at each node (Ahmad et al., 2019).

Categorical variables, non-linear correlations, and big datasets with numerous variables can all be handled via decision trees. They can be unstable with slight changes in data, are biased toward variables with more levels, and are prone to overfitting. They are also only able to predict binary outcomes, which may not be useful for doing so in cases when there are several alternative outcomes, such as forecasting churn (Ullah et al., 2019).

Gradient Boosting Machine

Gradient Boosting Machine (GBM) is a powerful and popular machine learning algorithm that can be used for regression, classification, and ranking tasks. GBM is capable of handling high-dimensional data and effectively reducing bias and variance in complex models. It works by building an ensemble of weak learners, such as decision trees, in a stepwise manner. The algorithm starts with a simple model and then iteratively improves upon it by fitting new models to the residual errors of the previous models. GBM optimizes a loss function using a gradient descent approach to minimize the discrepancy between predicted and actual values by adjusting the weights of the weak learners. The process is repeated until convergence or a pre-specified number of iterations is reached (Ahmad et al., 2019).

To answer the research questions: Q2. How to predict subscribers' churn in the OTT platform using Machine Learning?

The data related to subscribers' churn, customer behavior, and demographics will be collected from the OTT platform. The data may also include customer feedback, preferences, and engagement with the platform. The collected data will be preprocessed to remove any inconsistencies, missing values, and outliers. Feature engineering techniques will also be applied to create new features that can be useful in predicting churn and retention.

Data Collection

To do churn there are several characteristics that are important for subscriber churn prediction analysis in OTT (over-the-top) services, and they are:

1. User demographics: Age, gender, location, and other demographic data can be significant predictors of churn.
2. Usage patterns: Information on how frequently and how long a user accesses the service, which features are used, and when they are used can provide insights into churn likelihood.
3. Payment information: Whether a user is using a free trial, or a paid subscription, their billing history, and their payment method can be used to predict churn.
4. Content preferences: The types of content a user consumes, the ratings and reviews they provide, and their engagement with content can be used to predict churn.
5. Customer support interactions: Data on customer support tickets, complaints, and other interactions can indicate whether a user is at risk of churn.

6. Technical data: Information on the user's device type, network quality, and other technical data can also be useful in predicting churn.

It is important to note that not all these characteristics may be relevant for every OTT service or every subscriber churn analysis. The specific characteristics that are most important may depend on the type of service, the target audience, and the specific factors that drive churn for that service. Since the research papers discussed in the literature review were from various industries like telecom, financial, and social networks, it is not appropriate to use their dataset. So, after searching OTT datasets in different data repository websites like UCI, Kaggle, Data.gov, google dataset search, Data.world, Reddit, and OpenML got the churn modeling dataset from the Kaggle repository which satisfies the requirement to be used to develop the churn prediction analysis machine learning model.

To answer the research questions: Q3. How to retain subscribers and improve customer targeting?

Based on the insights gained from the models, customer retention strategies such as personalized recommendations, targeted marketing, and loyalty programs will be developed to improve customer engagement and retention. The developed models and retention strategies will be implemented in the OTT platform and monitored for effectiveness. Any necessary modifications will be made to further improve the performance of the models and strategies. The effectiveness of the implemented models and retention strategies will be evaluated using appropriate metrics such as customer churn rate, customer lifetime value, and revenue. The evaluation will help in identifying the strengths and weaknesses of the implemented models and strategies and making necessary improvements.

So, the summary is choosing the best algorithm for churn analysis depending on various factors such as data size and quality, problem complexity, desired accuracy, and interpretability. It is recommended to start with a simple model like logistic regression, which is computationally efficient and easy to interpret. For larger datasets, random forests and gradient boosting can handle missing data and outliers but may be slower. Performance metrics such as accuracy, precision, recall, and F1 score should be evaluated on a holdout set of data, and interpretability should be considered. Ensemble methods like random forests and gradient boosting can be used if no single algorithm performs well on the data. Ultimately, the best algorithm is chosen by comparing the performance of multiple algorithms on the specific characteristics of the problem and the data. The following chapter will begin with an introduction to our dataset, followed by a thorough analysis and results.

CHAPTER FOUR

DATA DESCRIPTION AND ANALYSIS

In this chapter, let's see the detailed data description, and the system requirements which need to run the machine learning program. Finally, the experimental results of the models have been discussed below. The dataset, which contains 16 columns and 2000 entries, is described in the table below.

Table 1. Data Description

S.No	Attribute Name	Description
1	Year	year of the subscription
2	customer_id	id of the subscriber
3	phone_no	phone number of the subscriber
4	Gender	gender of the subscriber
5	Age	age of the subscriber
6	no_of_days_subscribed	the number of days since the subscription
7	multi_screen	single/multiple-screen subscription
8	mail_subscribed	subscribers receive mails
9	weekly_mins_watched	number of minutes watched weekly
10	minimum_daily_mins	minimum minutes watched
11	maximum_daily_mins	maximum minutes watched
12	weekly_max_night_mins	number of minutes watched at nighttime
13	videos_watched	total number of videos watched
14	maximum_days_inactive	days since inactive
15	customer_support_calls	number of subscriber support calls
16	Churn	1-yes, 0-No

System Requirements

Preprocessing the dataset, implementing the solution, and training the models were all carried out in a Jupyter notebook. Version 6.1.4 of the Jupyter notebook is utilized. Python 3.8 is the language used for programming. Data science and machine learning both employ Python, a well-liked high-level programming language. For activities like data preprocessing, model training, and evaluation, it has a sizable and vibrant community that offers considerable help and resources for machine learning projects. NumPy, Pandas, Matplotlib, Scikit-Learn, TensorFlow, PyTorch, and Keras are a few of these libraries.

Implementation

To build a Machine Learning algorithm below steps should be followed:

- 1) Load the churned data.
- 2) Preprocess the data ready for use.
- 3) Split the preprocessed data into training and testing sets.
- 4) Fit and train the model of the chosen algorithm.
- 5) Evaluate the performance of the model on the testing set.
- 6) Save and load the trained model.

Experimental Results

In this project, generated two variables, "X" and "Y," after uploading our dataset to the environment, and allocated all of the features of the dataset to "X," with the exception of the column "Target," which was assigned to "Y." After dividing the datasets, the appropriate algorithm was then employed, in this instance, logistic regression. sklearn.linear model library's LogisticRegression package has been imported. Utilizing the testing dataset, a prediction was made after the trained

datasets had been fitted into the model. Got the following Classification Report in Fig.1 and an accuracy of 86.50%. The AUROC value of the model is 0.539 as shown in Fig.6.

Classification Report				
	precision	recall	f1-score	support
0.0	0.88	0.98	0.93	347
1.0	0.45	0.09	0.16	53
accuracy			0.86	400
macro avg	0.67	0.54	0.54	400
weighted avg	0.82	0.86	0.82	400

Figure 1. Classification Report of Logistics Regression

Secondly, the Multi-Layer Perceptron algorithm is used by importing the MLPClassifier package from sklearn.neural_network library. Set the hyperparameter hidden_layer, activation, solver, alpha, and max_iter with the following values 100, 'relu', 'adam', 0.0001, and 1000 respectively. After fitting the trained datasets into the model, predicted the outcome using the testing dataset. Got the following Classification Report in Fig.2 and an accuracy of 87.50%. The AUROC value of the model is 0.552 as shown in Fig.6.

Classification Report				
	precision	recall	f1-score	support
0.0	0.88	0.99	0.93	347
1.0	0.67	0.11	0.19	53
accuracy			0.88	400
macro avg	0.77	0.55	0.56	400
weighted avg	0.85	0.88	0.83	400

Figure 2. Classification Report of Multi-Layer Perceptron

Then used the Random Forest algorithm by importing the RandomForestClassifier package from sklearn.ensemble library. After fitting the

trained datasets into the model, predicted the outcome using the testing dataset. Got the following Classification Report in Fig.3 and an accuracy of 91%. The AUROC value of the model is 0.676 as shown in Fig.6.

Classification Report				
	precision	recall	f1-score	support
0.0	0.91	0.99	0.95	347
1.0	0.90	0.36	0.51	53
accuracy			0.91	400
macro avg	0.91	0.68	0.73	400
weighted avg	0.91	0.91	0.89	400

Figure 3. Classification Report of Random Forest Classifier

Then used the Decision Tree algorithm by importing the DecisionTreeClassifier package from sklearn.tree library. After fitting the trained datasets into the model, predicted the outcome using the testing dataset. Got the following Classification Report in Fig.4 and an accuracy of 92.25%. The AUROC value of the model is 0.779 as shown in Fig.6.

Classification Report				
	precision	recall	f1-score	support
0.0	0.94	0.97	0.96	347
1.0	0.78	0.58	0.67	53
accuracy			0.92	400
macro avg	0.86	0.78	0.81	400
weighted avg	0.92	0.92	0.92	400

Figure 4. Classification Report of Decision Tree Classifier

Finally, used the Gradient Boosting Machine algorithm by importing the GradientBoostingClassifier package from sklearn.ensemble library. After fitting the trained datasets into the model, predicted the outcome using the testing dataset. Got

the following Classification Report in Fig.5 and an accuracy of 93.25%. The AUROC value of the model is 0.793 as shown in Fig.6.

Classification Report				
	precision	recall	f1-score	support
0.0	0.94	0.98	0.96	347
1.0	0.84	0.60	0.70	53
accuracy			0.93	400
macro avg	0.89	0.79	0.83	400
weighted avg	0.93	0.93	0.93	400

Figure 5. Classification Report of Gradient Boosting Machine

A common performance metric in binary classification issues is AUROC (Area Under the Receiver Operating Characteristic Curve). At various categorization criteria, the true positive rate (TPR) is plotted versus the false positive rate (FPR). The TPR, sometimes referred to as sensitivity, is calculated as the percentage of true positives (positive samples that were properly predicted to be positive) among all positive samples. The FPR is defined as the percentage of false positives (positive samples that were mistakenly predicted to be positive) among all negative samples.

The classifier's prediction scores are arranged in descending order to determine AUROC, and the threshold is steadily raised. TPR and FPR are calculated for each threshold and plotted on the ROC curve. The AUROC score, which goes from 0 to 1, is the region under this curve. A score of 1 represents flawless classification, whereas a score of 0.5 shows that the classifier is only slightly more accurate than random guessing. Unlike other measures like accuracy, precision, and recall, the AUROC is not affected by class imbalance and threshold selection. This makes it a valuable metric. It is especially helpful since the threshold can be

changed to maximize the statistic that matters most in the particular context when the cost of false positives and false negatives differ. The AUROC of the used model is shown below in Fig.6.

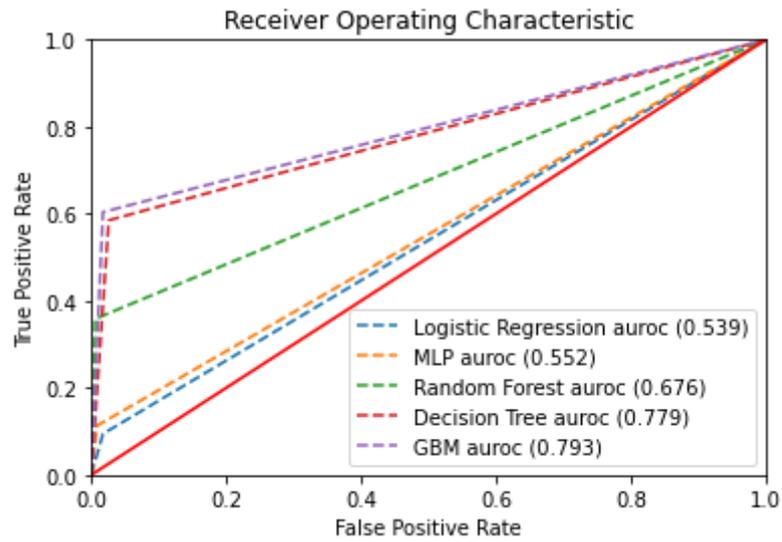


Figure 6. Values of AUROC

CHAPTER FIVE

DISCUSSION, CONCLUSION, AND AREAS OF FURTHER STUDY

Discussion

The research questions were:

- 1) What Machine Learning algorithms are used to overcome subscriber churn?
- 2) How to predict subscribers' churn in the OTT platform using Machine Learning?
- 3) How to retain subscribers and improve customer targeting?

After discussing the results and drawing a conclusion, follows suggestion for areas for further study. Here are some findings and conclusions from previous questions. In reference to the 1st question, from the literature review and experimental results, the research come to the solution that Machine learning algorithms used to overcome subscriber churn are Logistics Regression, Multi-Layer Perceptron, Random Forest, Decision Trees, and Gradient Boosting Machines. Moving on to the 2nd question, by sending the test data to a trained model by their historical dataset, customers are likely to leave a company (i.e., churn) based on their historical behavior and characteristics can be predicted.

Lastly referring to the 3rd question, here are some strategies that can be used based on the insights obtained from churn prediction analysis: Identify at-risk subscribers, personalized offers and promotions, improve customer service, develop loyalty programs, and optimize pricing strategies.

Talking about Identifying at-risk subscribers' strategy, churn prediction models can help identify subscribers who are at risk of churning. Once these

subscribers are identified, targeted retention strategies can be developed to address their specific needs and concerns. Moving on to personalized offers and promotions strategy, by analyzing subscriber behavior and preferences, churn prediction models can help companies develop personalized offers, discounted subscription plans, free trials, or exclusive content to retain subscribers.

Churn prediction models can help identify common issues that lead to subscriber churns, such as poor customer service or long wait times. Companies can use this information to improve their customer service and address subscriber concerns in a timely manner. By rewarding loyal subscribers with exclusive offers and discounts, companies can increase subscriber retention and reduce churn. The models can help identify which subscribers are most likely to respond to loyalty programs and what types of rewards are most effective. Analyzing subscriber behavior and subscription history can help companies optimize their pricing strategies.

Conclusion

In this project, five machine learning algorithms have been analyzed and got the highest accuracy and maximum AUROC of 93.25% and 0.793 respectively in the Gradient Boosting Machine Model. And so, in churn prediction, GBM can be used to identify the most significant predictors of churn and build a model that predicts which subscribers are at a high risk of leaving. GBM can handle many features, including both numerical and categorical variables, making it suitable for churn prediction analysis. The algorithm can also deal with imbalanced datasets and handle missing values, which are common in real-world churn prediction scenarios. GBM can provide accurate and robust predictions, making it a powerful tool in the fight against customer churn.

Area for Further Study

Churn prediction is a complex and dynamic field, and there are several areas of further study that could help improve the accuracy and effectiveness of churn prediction models. One area of study is the incorporation of unstructured data sources, such as social media and customer reviews, into churn prediction models. Another area of study is the use of deep learning techniques, such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs), to improve the accuracy of churn prediction models. Finally, the integration of real-time data sources, such as clickstream data and mobile app usage data, could enable more accurate and timely identification of customers at risk of churn. Overall, there are many exciting areas of further study in churn prediction, and continued research and innovation in this field have the potential to greatly improve customer retention strategies across a variety of industries.

REFERENCES

- Agrawal, S., Das, A., Gaikwad, A., & Dhage, S. (2018). Customer Churn Prediction Modelling Based on Behavioural Patterns Analysis using Deep Learning. *2018 International Conference on Smart Computing and Electronic Enterprise, ICSCEE 2018*. <https://doi.org/10.1109/ICSCEE.2018.8538420>
- Ahmad, A. K., Jafar, A., & Aljoumaa, K. (2019). Customer churn prediction in telecom using machine learning in big data platform. *Journal of Big Data*, 6(1). <https://doi.org/10.1186/s40537-019-0191-6>
- De Caigny, A., Coussement, K., & De Bock, K. W. (2018). A new hybrid classification algorithm for customer churn prediction based on logistic regression and decision trees. *European Journal of Operational Research*, 269(2). <https://doi.org/10.1016/j.ejor.2018.02.009>
- Fitzgerald, S. (2019). Over-the-Top Video Services in India: Media Imperialism after Globalization. *Media Industries*, 6.
- Huang, B., Kechadi, M. T., & Buckley, B. (2012). Customer churn prediction in telecommunications. *Expert Systems with Applications*, 39(1). <https://doi.org/10.1016/j.eswa.2011.08.024>
- Kuldeep, C., Rojhe, V., Singh, N., & Rao, A. (n.d.). *{YICCISS-2021} HJ Emerging Trends in Management Sciences*. Retrieved March 17, 2023, from www.houseofjournals.com
- Madden, G., Savage, S. J., & Coble-Neal, G. (1999a). Subscriber churn in the Australian ISP market. *Information Economics and Policy*, 11(2), 195–207. [https://doi.org/10.1016/S0167-6245\(99\)00015-3](https://doi.org/10.1016/S0167-6245(99)00015-3)

Madden, G., Savage, S. J., & Coble-Neal, G. (1999b). Subscriber churn in the Australian ISP market. *Information Economics and Policy*, 11(2), 195–207.

[https://doi.org/10.1016/S0167-6245\(99\)00015-3](https://doi.org/10.1016/S0167-6245(99)00015-3)

Ullah, I., Raza, B., Malik, A. K., Imran, M., Islam, S. U., & Kim, S. W. (2019). A Churn Prediction Model Using Random Forest: Analysis of Machine Learning Techniques for Churn Prediction and Factor Identification in Telecom Sector.

IEEE Access, 7, 60134–60149. <https://doi.org/10.1109/ACCESS.2019.2914999>