# Customer Churn Analysis and Prediction

➔ **Chirag Raj (cs21b1038)**
➔ **Bhawani Shankar Dhawal (cs21b1005)**

## Abstract-

The purpose of this report is to analyze and predict customer churn using various machine learning models. Customer churn, the rate at which

customers stop doing business with an entity, is a significant metric for many businesses as it affects the revenue and long-term success directly. By leveraging historical data, this study aims to identify key factors that contribute to customer churn and to predict future trends using models such as RandomForestClassifier, LogisticRegression, GaussianNB, DecisionTreeClassifier, and XGBClassifier.This analysis seeks to provide actionable insights that can help in devising effective strategies to retain customers and reduce churn rates.

## Problem Statement -

In today's competitive market, retaining customers is just as crucial as acquiring new ones. Businesses are increasingly facing the challenge of customer churn, which can lead to significant losses in revenue and increased operational costs. Predicting customer churn allows companies to proactively implement retention strategies to prevent loss of customers. However, the significant diversity in customer behaviors and profiles makes it difficult to identify and generalize the pre-churn indicators. This report addresses the problem by applying machine learning techniques to historical data to model and predict churn probability. The objective is to determine which factors are most influential in customer turnover and to assess the effectiveness of various predictive models in forecasting churn. This investigation will enable the deployment of more personalized and timely interventions aimed at retaining customers in various business sectors.

## Model Analysis -

This section delves into the various machine learning models that were employed to predict customer churn. Each model has its unique approach to learning from data and making predictions. The performance of the

RandomForestClassifier, LogisticRegression, GaussianNB, DecisionTreeClassifier, and XGBClassifier models are evaluated based on their accuracy, precision, recall, and F1-score. This analysis includes understanding the strengths and weaknesses of each model in handling the dataset, as well as their computational efficiency and suitability for this particular application. The core aim is to determine the most effective model in terms of accuracy and reliability for predicting customer churn. Each model's methodology, the rationale for its selection, and its performance metrics are discussed in the subsequent subsections.

## 1. RandomForestClassifier-

- The RandomForestClassifier is an ensemble learning method based on the decision tree algorithm. It operates by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes of the individual trees. This model is particularly well-known for its robustness and ability to handle large datasets with higher dimensionality. It performs both classification and regression tasks with accuracy and effectively handles missing data.

- In the context of customer churn prediction, RandomForestClassifier mitigates overfitting, a common issue with decision trees, by averaging multiple deep decision trees, trained on different parts of the same training set. This improves the generalizability of the model. The model's ability to rank the importance of various features according to their impact on the outcome makes it invaluable for identifying key attributes contributing to churn.

- However, despite its advantages, the RandomForestClassifier can be computationally intensive, particularly when dealing with very large datasets and a large number of trees in the forest. This can lead to longer training times and increased complexity in tuning the model parameters like the number of trees and depth of each tree. Additionally, while the RandomForestClassifier typically provides high accuracy, its performance heavily depends on the relevance and quality of the input features.

## 2. LogisticRegression -

- Logistic Regression is a statistical model that, despite its name, is used for classification rather than regression. It predicts the probability of occurrence of an event by fitting data to a logistic function. Consequently, it is particularly suited for binary classification tasks, such as predicting whether a customer will churn or not churn.

- In terms of customer churn prediction, Logistic Regression offers a straightforward and interpretable model that does not require scaling of input features and provides coefficients that represent the importance and effect of each feature. This interpretability is crucial for business applications where understanding the influence of individual factors on churn is necessary for strategic decision-making.

- The simplicity of Logistic Regression can also be considered a limitation. It assumes a linear relationship between the independent variables and the logarithm of odds of the dependent variable, which may not always capture the complexities or interactions between features effectively. This could lead to underfitting if the actual relationships are non-linear. Moreover, it tends to perform poorly with non-linear decision boundaries compared to other algorithms. Nevertheless, its efficiency in terms of computational resources makes it a valuable first-line approach in customer churn prediction, particularly when a preliminary analysis is needed or when the dataset is not excessively large.

## 3. GaussianNB -

- GaussianNB, or Gaussian Naive Bayes, is a variant of the Naive Bayes classifier that assumes the features follow a normal distribution. This assumption simplifies calculations, thereby allowing for a quick and easy way to predict the class of a dataset. It is particularly effective for large datasets and has been widely used in applications needing probabilistic classification.

- For the purpose of customer churn analysis, GaussianNB works well when the features of the dataset are continuous and normally distributed. The model calculates the probabilities of each class for a given customer and then predicts the class with the highest probability. Its strength lies in its simplicity and speed, particularly in cases where the assumption of independence among features holds reasonably true.

- However, the primary limitation of GaussianNB emerges from its foundational assumptions. In the real world, features often exhibit some level of correlation, and the assumption of normal distribution does not always hold, which can affect the model's accuracy. Additionally, the model's performance can degrade with features that are not normally distributed unless transformations are

applied. Despite these challenges, GaussianNB remains a popular choice due to its efficiency and effectiveness in predictive tasks where the conditions align closely with its assumptions.

## 4. DecisionTreeClassifier -

- The DecisionTreeClassifier is a non-parametric supervised learning method used for classification and regression tasks. It works by breaking down a dataset into smaller subsets while at the same time an associated decision tree is incrementally developed. The final result is a tree with decision nodes and leaf nodes. For customer churn prediction, a decision node represents a test on a specific feature, while a leaf node represents a class label (churn or not churn).

- One of the key advantages of using DecisionTreeClassifier is its ability to handle both numerical and categorical data. It's easy to understand and interpret, which makes it a great tool for decision making. The visual representation of a Decision Tree is very intuitive and the logic of the decisions made by the model can be easily explained to non-technical team members.

- However, Decision Trees have their limitations. They are prone to overfitting, especially with very complex trees when the model learns noise in the training data as valid signals, thus performing poorly on unseen data. This can be somewhat mitigated by pruning the tree to remove parts of the model that have little power in predicting target variables. Moreover, Decision Trees are sensitive to the specifics of the data on which they are trained. Small changes in the training data can result in very different tree structures and this instability can be a drawback in dynamic environments. Despite these issues, when used with careful tuning and in conjunction with other models, the DecisionTreeClassifier can be a very effective tool in predicting customer churn.
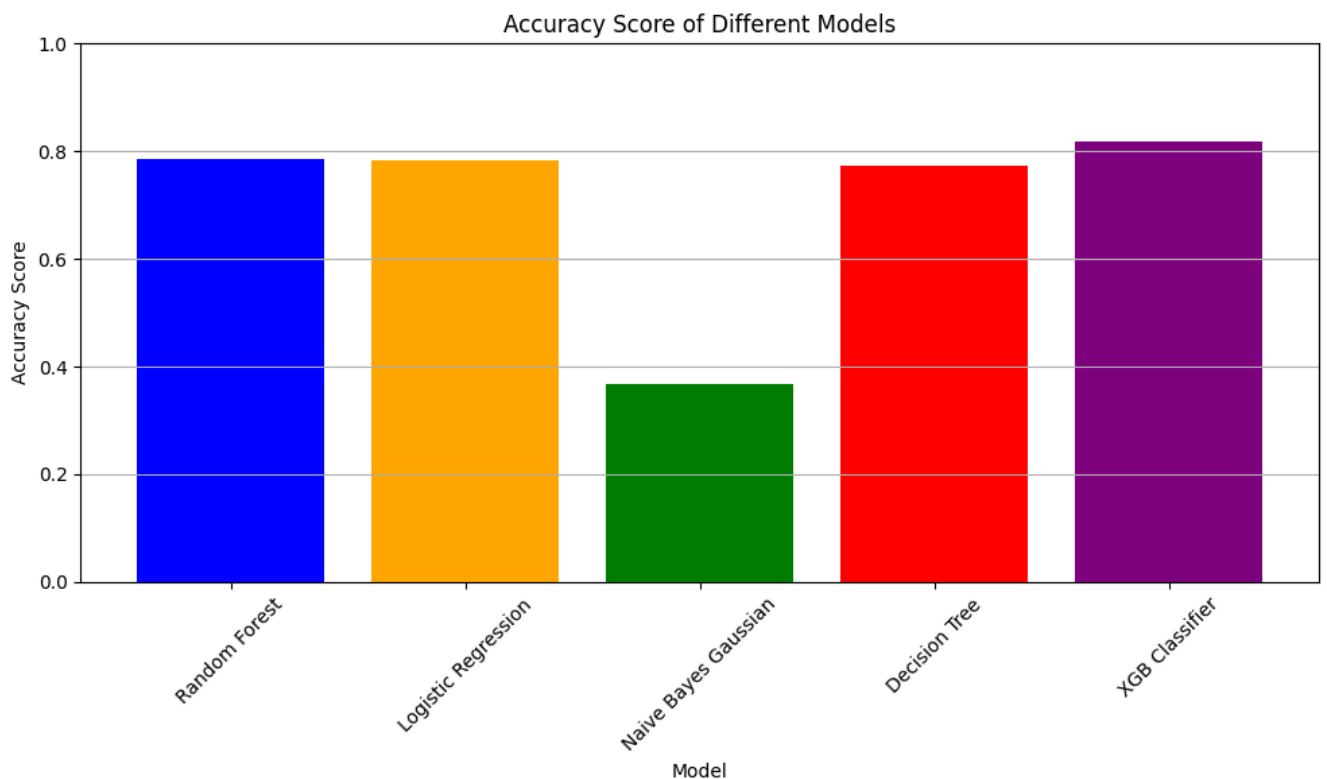
## 5. XGBClassifier -

- XGBClassifier stands for eXtreme Gradient Boosting Classifier, which is part of a library designed for high performance and high efficiency. It is an implementation of gradient boosted decision trees designed for speed and performance. In the context of customer churn prediction, the XGBClassifier is particularly notable for its ability to handle large and complex datasets with a mix of categorical and numerical features effectively.

- The major strength of the XGBClassifier lies in its model building strategy. It incorporates techniques such as tree pruning, handling missing values, and regularization which help in improving the model's accuracy and preventing overfitting. The model uses gradient boosting frameworks that are highly efficient, as they combine multiple weak models to create a strong predictive model.

- For churn prediction, XGBClassifier offers a significant advantage due to its robust handling of various types of data and its ability to model complex relationships in data. It outperforms many other classifiers in both speed and accuracy due to its sophisticated handling of underfitting and overfitting, making it especially effective for this application. The boosted trees algorithm also allows it to evaluate the importance of different features in the churn decision, providing valuable insights into which factors are most influential in customer retention strategies.

- The superior performance of XGBClassifier in predicting customer churn can be attributed to these advanced capabilities, particularly its enhancement over other ensemble techniques through regularization and tree pruning strategies, which helps in managing real-world, noisy data effectively while maintaining high computational efficiency.

## Results and Discussion -

- The analysis conducted across various models for predicting customer churn reveals a superior performance by the XGBClassifier. When compared to RandomForestClassifier, LogisticRegression, GaussianNB, and DecisionTreeClassifier, the XGBClassifier not only demonstrated higher accuracy but also showed robustness in handling the complexities of the dataset.

- While RandomForestClassifier and DecisionTreeClassifier both exhibited commendable performance, especially in capturing the non-linear relationships between features, they were both prone to overfitting. RandomForestClassifier, despite being more generalized than a single DecisionTree, could not match the precision and handling of noisy data that XGBClassifier offered.

DecisionTreeClassifier, although highly interpretable, was significantly affected by small changes in the data, leading to instabilities in the predictions.

- LogisticRegression and GaussianNB were less effective compared to tree-based methods. LogisticRegression, being a linear model, struggled with the complexity and the multi-dimensional nature of the dataset. GaussianNB, on the other hand, was fast and effective for partial datasets but its assumption of the independence of features did not hold well in the interconnected feature set of this particular churn data.

- The XGBClassifier outperformed other models due to its advanced handling of underfitting and overfitting through gradient boosting and regularization. These methods not only help in optimizing the performance but also in achieving better generalization on unseen data. It effectively addressed issues like missing values and varied feature scales which often negatively impact other models. Moreover, its capability to sequentially build upon the errors of the previous trees and to focus on difficult cases gave it a higher predictive accuracy.

- Therefore, based on the comparative results and the characteristics of the analyzed models, the XGBClassifier stands out as the most suitable model for the task of predicting customer churn. Its ability to deliver high performance while managing common data issues effectively makes it the preferred choice in the context of this analysis.



Accuracy Score of Different Models

## Business Impact -

- In the context of a telecommunications company, the predictive model that categorizes customers into three categories—stayed, churned, and joined—has specific implications for business strategies aimed at improving service delivery and customer loyalty.

- For customers identified as 'stayed' (output 0), the company can focus on upselling and cross-selling more services, such as upgraded data plans or additional lines, based on the customer's usage patterns and satisfaction levels. Regular personalized communication can be employed to ensure these customers are aware of the latest technologies, offers, and benefits, thereby reinforcing their decision to stay.

- For those labeled as 'churned' (output 1), the company can implement immediate recovery strategies. This may involve contacting these customers through call centers to offer specially tailored plans, discounts, or bonuses as incentives to reconsider their decision. Understanding the reasons for their dissatisfaction through direct feedback is crucial and can lead to service improvements that prevent future churn.

- Lastly, for new customers or those predicted as 'joined' (output 2), the company should ensure a seamless onboarding process. This includes clear guidance on plan options, network usage, and troubleshooting, along with introductory offers that enhance their initial experience and foster long-term loyalty. Early engagement through welcome calls or emails can also help in gathering initial feedback and establishing a positive relationship.

- Through these focused strategies, the telecommunications company can effectively manage different customer segments, thus improving retention rates, reducing churn, and ultimately maximizing profitability and market share.

# Features and Limitations -

- The project incorporated a range of features to analyze and predict customer churn effectively. These features included customer demographic information, account details, service usage patterns, and historical transaction data. The diversity of data types, from categorical to continuous, allowed for a comprehensive analysis of factors influencing churn.

- However, the project also encountered certain limitations that could impact the overall outcome and general applicability of the results. One significant limitation was the potential bias in the dataset due to the method of data collection, which may not perfectly represent the entire customer base. This could lead to skewed results and models that perform well on the dataset but less effectively on new, real-world data.

- Another limitation was the handling and interpretation of complex relationships and interactions between various features. Despite the advanced capabilities of models like the XGBClassifier, simplifying these relationships might lead to the loss of critical information. Additionally, the reliance on historical data assumes that past patterns will continue to apply in the future, which may not always be the case due to changing market dynamics and customer behaviors.

- Moreover, the project did not incorporate real-time data, which can be crucial for understanding and responding to immediate churn risks. The models used also required substantial computational resources for training, especially when handling large datasets, potentially limiting the scalability of the solution in a resource-constrained environment.

- Overall, while the project leveraged powerful analytical tools and diverse data, these limitations highlight the importance of continuous model evaluation and updates, considering dynamic customer behaviors and technological advancements.

# Conclusion -

This report has comprehensively analyzed the phenomenon of customer churn through the application of various machine learning models, with a clear demonstration that the XGBClassifier provides the highest precision in predicting churn. The study illuminated the influence of diverse features on customer retention and underscored the critical role of sophisticated predictive analytics in formulating effective customer retention strategies.

Despite the promising results, the analysis acknowledged inherent limitations related to data bias, resource demands, and the static nature of used historical data. These factors emphasize the necessity for ongoing refinement of the models and adaptation to evolving market conditions and technologies.

Ultimately, the insights garnered from this study should guide further research and practical implementations aimed at reducing customer churn. By continuously enhancing the predictive capabilities and addressing the identified limitations, businesses can better satisfy their customer base, improve retention, and sustain competitive advantage in the market.