# A MACHINE LEARNING APPROACH FOR CUSTOMER CHURN PREDICTION IN TELECOMMUNICATION INDUSTRY

Submitted by

| | |
|---|---|
| Erothu Sai Kumar | 21881A6680 |
| J Nithiin | 21881A6688 |
| K Aditya | 21881A6690 |

### SUPERVISOR
Mrs B Pravallika
Assistant Professor

Department of Computer Science and Engineering (AI&ML)

## VARDHAMAN COLLEGE OF ENGINEERING
(AUTONOMOUS)

**July, 2024**

## Department of Computer Science and Engineering (AI&ML)

# CERTIFICATE

This is to certify that the mini-project titled **A MACHINE LEARN-ING APPROACH FOR CUSTOMER CHURN PREDICTION IN TELECOMMUNICATION INDUSTRY** is carried out by

| | |
|---|---|
| **Erothu Sai Kumar** | **21881A6680** |
| **J Nithiin** | **21881A6688** |
| **K Aditya** | **21881A6690** |

in partial fulfillment of the requirements for the award of the degree of **Bachelor of Technology** in **Computer Science and Engineering (AI&ML)** during the year 2023-24.

**Signature of the Supervisor**
**Mrs B Pravallika**
**Assistant Professor**
**Dept of CSE(AI&ML)**

**Signature of the HOD**
**Dr. M.A. Jabbar**
**Professor & Head**
**Dept of CSE(AI&ML)**

Project Viva-Voce held on _____

**External Examiner**

Kacharam (V), Shamshabad (M), Ranga Reddy (Dist.)–501218, Hyderabad, T.S.

Ph: 08413-253335, 253201, Fax: 08413-253482, www.vardhaman.org

# Acknowledgement

# Abstract

Accurate prediction of customer churn is crucial for telecom companies as it directly impacts customer retention strategies and revenue. Traditional approaches often fall short in modern, dynamic markets characterized by changing customer behaviors and unexpected disruptions. Machine learning offers a powerful solution by leveraging historical data, customer behaviors, and external factors to enhance the accuracy of churn predictions. This project employs Random Forest and Hidden Naive Bayes algorithms to discover intricate patterns in the data, resulting in more precise predictions. Future research in machine learning for churn prediction aims to incorporate advanced techniques such as deep learning and reinforcement learning. It also seeks to leverage big data technology and develop hybrid models that combine machine learning with human expertise. Although machine learning provides substantial benefits for churn prediction, challenges such as data quality, result interpretation, and model scalability must be addressed through thorough data preprocessing and appropriate model selection. Machine learning-based churn prediction has demonstrated significant advantages, including improved prediction accuracy, reduced errors, and better decision-making. Successful implementations in various sectors, including telecom, highlight the transformative potential of machine learning in contemporary customer retention strategies.

***Keywords***: Customer Churn Prediction, Machine Learning, Random Forest, Hidden Naive Bayes, Telecom Industry

# Table of Contents

# List of Tables

# List of Figures

# Abbreviations

| Abbreviation | Description |
|---|---|
| HIB | Hidden naive Bayes |
| RFC | Random Forest Classifier |
| KPI | Key Performance Indicator |
| Telecom | Telecommunication |

# CHAPTER 1

# Introduction

## 1.1 Introduction

Accurate customer churn prediction is essential for telecommunication companies as it directly influences customer retention strategies, revenue management, and overall business sustainability. The telecommunication industry is highly competitive, with customers often having multiple service providers to choose from. This makes it crucial for companies to not only attract new customers but also retain existing ones.

Traditional methods of churn prediction, such as basic statistical analyses, often struggle to keep pace with the rapidly changing customer behaviors and competitive landscape. These methods typically rely on historical data and simple linear relationships, which may not fully capture the complexity of customer decision-making processes. For example, traditional logistic regression models can indicate which factors are most strongly associated with churn, but they may fail to account for non-linear interactions and the dynamic nature of customer preferences and market conditions.

To address these limitations, this research focuses on the application of advanced machine learning techniques, specifically Random Forest classifiers and Hidden Naive Bayes models. These models offer significant improvements over traditional methods by leveraging extensive historical customer data enriched with contextual factors such as usage patterns, billing information, and service complaints. The goal is to uncover complex patterns and relationships that contribute to customer churn, thereby developing robust predictive models that provide actionable insights.

The integration of Random Forest classifiers and Hidden Naive Bayes models into churn prediction represents a significant advancement over traditional methods. By leveraging the power of data and advanced algorithms, telecom

companies can gain deeper insights into customer behavior, enhance their retention strategies, and maintain a competitive edge in a fast-paced industry. This lecture will delve into the specifics of these machine learning techniques, their advantages and limitations, and practical applications in the telecommunication sector. By understanding and utilizing these advanced models, telecom companies can optimize customer retention efforts and improve decision-making processes in a highly dynamic market environment.

### 1.1.1 Background

Customer churn, the phenomenon where customers discontinue using a company's services, presents a critical challenge for telecom companies, affecting revenue and profitability significantly. Traditional methods of churn prediction involve analyzing historical data to detect patterns that precede customer departure. However, these approaches may struggle to capture the dynamic and multifaceted nature of customer behavior, which evolves over time and is influenced by various factors such as service quality, pricing, competitor offerings, and customer satisfaction. Machine Learning (ML) offers a powerful alternative by harnessing large-scale datasets and sophisticated algorithms to uncover subtle correlations and predictive insights that traditional analytics methods might miss.

ML algorithms like Decision Trees, Random Forests, SVMs, and Neural Networks excel in churn prediction by processing extensive customer data, including demographic information, usage patterns, and customer interactions. These models can identify complex relationships and non-linear patterns within the data, providing telecom companies with actionable insights to forecast which customers are likely to churn. By leveraging these insights, telecom providers can proactively intervene with targeted retention strategies, such as personalized offers, loyalty programs, or proactive customer service interventions, aimed at reducing churn rates.

Moreover, ML-driven churn prediction enables telecom companies to move beyond reactive measures to a proactive approach, where they anticipate customer needs and preferences in advance. This proactive stance not only

enhances customer retention but also improves overall customer satisfaction and loyalty. By integrating ML models into their operational workflows, telecom companies can optimize resource allocation, streamline marketing efforts, and enhance the overall customer experience, thereby maintaining a competitive edge in the highly competitive telecom industry landscape.

## 1.2    Motivation

The motivation for this project stems from the significant impact that accurate churn prediction can have on a telecom company's bottom line. Customer acquisition is often more expensive than retention, making it crucial for telecom companies to identify at-risk customers early to save costs and improve customer loyalty. With the rise of big data and advanced analytics, there is an unparalleled opportunity to enhance churn prediction models. Telecom companies generate vast amounts of data, including call records, billing information, customer service interactions, and usage patterns. Leveraging this rich data landscape, sophisticated machine learning models can uncover hidden patterns and insights, enabling more effective retention strategies.

In the dynamic telecom industry, where customer preferences and behaviors are constantly evolving, machine learning offers a way to stay ahead. Traditional statistical methods often fall short in capturing the complexity and fluidity of customer behavior. Machine learning models, such as Random Forest classifiers and Hidden Naive Bayes, can continuously learn from new data, adapting to changing patterns and trends. This adaptability is crucial for maintaining the relevance and accuracy of churn predictions. By providing timely insights, these models enable telecom companies to make informed decisions, design targeted retention strategies, allocate resources more effectively, and tailor their marketing efforts to address the specific needs and concerns of at-risk customers. This proactive approach not only helps in retaining customers but also enhances overall customer satisfaction and loyalty, supporting long-term success in a highly competitive market.

## 1.3 Problem Statement

Despite the advancements in data analytics, many telecom companies struggle with accurately predicting customer churn. Traditional models may fail to capture the multifaceted nature of customer behavior and the influence of external factors. This project aims to address these challenges by developing a machine learning-based approach that uses Random Forest and Hidden Naive Bayes algorithms. The specific problem statement is: "How can we improve the accuracy of telecom customer churn prediction using machine learning models, specifically Random Forest and Hidden Naive Bayes, to provide actionable insights for customer retention strategies?"

## 1.4 Objectives

The primary objective of this project is to develop a machine learning model using Random Forest and Hidden Naive Bayes algorithms to predict telecom customer churn. The study aims to evaluate the performance of these models using appropriate metrics such as accuracy, precision, recall, and F1-score. Additionally, the project seeks to identify key features that contribute to customer churn, providing insights that can inform retention strategies. A comparative analysis of the effectiveness of the Random Forest and Hidden Naive Bayes models in predicting churn will be conducted. Finally, the project aims to create a user-friendly interface or dashboard for visualizing the prediction results and insights, ensuring the practical application of the developed models in real-world scenarios.

## 1.5 Scope

This project spans the complete lifecycle of constructing a machine learning model for predicting telecom customer churn. The process encompasses data collection, where relevant data such as customer demographics, usage patterns, billing information, and customer service interactions are gathered from a telecom company. Data preprocessing involves cleaning and preparing the data

for analysis by handling missing values, encoding categorical variables, and normalizing numerical features. Model development includes implementing the Random Forest and Hidden Naive Bayes algorithms and tuning their parameters for optimal performance. Model evaluation is performed using cross-validation and various performance metrics to assess the accuracy and reliability of the models. Insights and recommendations are derived by analyzing the model outputs to identify key factors contributing to churn, providing actionable strategies for customer retention. Finally, deployment involves developing a dashboard or interface to visualize the results and make them accessible to stakeholders.

# CHAPTER 2

# Literature Survey

## 2.1 Traditional Churn Prediction Methods

Before delving into machine learning techniques, it's important to understand the traditional methods used for churn prediction. Historically, churn prediction has relied heavily on statistical methods such as logistic regression and survival analysis. These techniques involve creating a model based on historical data to identify factors that are indicative of churn.

- **Logistic Regression**: This method models the probability of a customer churning as a function of various explanatory variables. It's a simple and interpretable model that can provide insights into which factors are most strongly associated with churn. However, it may struggle with non-linear relationships and interactions between variables.

- **Survival Analysis**: This statistical approach is used to predict the time until an event, such as churn, occurs. It can handle censored data (where the event has not yet occurred for some individuals) and provides a way to model the timing of churn. While powerful, survival analysis can be complex to implement and interpret.

### 2.1.1 Overview of Churn Models

Historically, customer churn prediction in the telecommunications industry has heavily depended on statistical approaches to anticipate future patterns by analyzing past data. Methods like logistic regression and decision trees are commonly employed due to their capacity to accurately represent linear relationships and categorical data. Logistic regression is highly proficient in capturing relationships between customer attributes and churn probability, but decision trees excel at handling non-linear interactions and variable importance.

However, these methodologies have limitations in their ability to accurately represent complex patterns and incorporate external influences, which are becoming increasingly important in today's dynamic telecom markets.

The incorporation of machine learning (ML) methodologies has greatly propelled the domain of churn prediction. Approaches like Random Forest, Gradient Boosting Machines (GBM), and neural networks can discover intricate patterns in extensive datasets that conventional methods may overlook. Random Forest is renowned for its precision and ability to handle structured data, while neural networks excel at capturing complex, non-linear connections. These models improve prediction accuracy by incorporating external variables such as customer sentiment and usage trends. Despite their potential, ML models require thorough data preprocessing and substantial computer resources. Balancing precision with comprehensibility remains a challenge, necessitating the development of scalable and interpretable models that deliver actionable insights for strategic decision-making.

## 2.2 Review on Existing Methods

The landscape of churn prediction has evolved significantly with the advent of machine learning. Various methods and models have been proposed and used with varying degrees of success. Here, we will review some of the most relevant ones:

- **Decision Trees**: Decision Trees are supervised learning models used for both classification and regression tasks, including churn prediction. They recursively split the data based on the values of input features to create subsets that are as homogeneous as possible with respect to the target variable (e.g., churners vs. non-churners). Each split is chosen to maximize information gain or minimize impurity, such as Gini impurity or entropy. Decision Trees are straightforward to interpret and visualize, making them useful for understanding the factors influencing churn. They can handle both numerical and categorical data and capture non-linear relationships between features and the target variable. However, Decision

Trees are prone to overfitting, especially with deep trees that capture noise in the training data. Techniques like pruning or using ensemble methods such as Random Forests can mitigate these issues and improve generalization. In churn prediction, Decision Trees provide insights into customer behavior and predictors of churn, facilitating proactive retention strategies in customer analytics.

- **Random Forest**: Random Forest is an ensemble learning method that constructs multiple decision trees during training and outputs the mode of the classes (classification) or mean prediction (regression) of individual trees. Each tree in the forest is trained on a random subset of the data and a random subset of the features, aiming to reduce overfitting and improve generalization. Random Forests are particularly effective in churn prediction due to their ability to handle large datasets with many features, provide feature importance rankings, and mitigate the variance of individual decision trees. They are robust against overfitting, easy to implement, and less sensitive to noisy data compared to single decision trees. However, tuning parameters such as the number of trees and maximum depth is crucial to optimize performance and avoid excessive computational costs. Random Forests are widely used in customer analytics for their balance of accuracy, interpretability, and scalability in predicting customer churn.

- **Support Vector Machines (SVM)**: Support Vector Machines (SVMs) excel in churn prediction by identifying optimal hyperplanes that maximize the margin between different customer groups (e.g., churners vs. non-churners) in high-dimensional spaces. Using kernel functions, SVMs can handle complex relationships and non-linear boundaries effectively. They are robust against overfitting and outliers but require careful parameter tuning, such as selecting the right kernel and regularization parameter (C). SVMs are computationally intensive, particularly with large datasets, yet they offer a powerful approach for accurately classifying churners based on diverse customer attributes.

- **Neural Networks**: These models are capable of learning complex patterns in the data. They are particularly effective for large datasets and non-linear relationships. However, they require substantial computational resources and expertise in tuning hyperparameters.

- **Gradient Boosting Machines (GBM)**: This ensemble technique builds models sequentially, with each model correcting the errors of the previous ones. GBMs, including variants like XGBoost and LightGBM, are highly effective for predictive tasks but can be prone to overfitting without careful tuning.

.

## 2.3 Machine Learning in Churn Prediction

With the increasing availability of data and computational power, machine learning has become a key tool in churn prediction. Let's discuss some studies and approaches that have utilized machine learning for this purpose:

- **Data-Driven Approaches**: These involve using historical data on customer behavior, demographics, and interactions to train machine learning models. Studies have shown that incorporating features like usage patterns, billing information, and customer service interactions can significantly improve prediction accuracy.

- **Feature Engineering**: Identifying and creating meaningful features from raw data is crucial for the success of machine learning models. Techniques like one-hot encoding, feature scaling, and dimensionality reduction (e.g., PCA) are commonly used.

- **Model Comparison and Hybrid Approaches**:Researchers often compare multiple models to identify the best-performing one. Some studies also explore hybrid approaches, combining the strengths of different models to improve performance. For example, combining a decision tree-

based model with a neural network can leverage the interpretability of the former and the predictive power of the latter.

## 2.4 Challenges and Limitations

Despite the advancements, several challenges remain in the field of churn prediction:

- **Data Quality and Availability**: High-quality, comprehensive data is essential for accurate predictions. Missing values, inconsistent data, and limited historical data can pose significant challenges.

- **Model Interpretability**: Complex models, such as neural networks, often act as "black boxes," making it difficult to interpret their predictions. This can hinder the ability to derive actionable insights.

- **Scalability and Deployment**: Implementing machine learning models in a real-world setting requires robust infrastructure and continuous monitoring to ensure they perform well over time.

## 2.5 Conclusion

The literature survey highlights the evolution of churn prediction methods from traditional statistical techniques to advanced machine learning models. While machine learning offers significant advantages in terms of accuracy and adaptability, challenges related to data quality, model interpretability, and scalability must be addressed. The review of existing methods provides a foundation for developing a robust churn prediction model using Random Forest and Hidden Naive Bayes algorithms, as proposed in this project.

| Paper | Pros | Cons |
|---|---|---|
| Integrated Churn Prediction and Customer Segmentation Framework for Telco Business IEEE-April 2021 | Utilizes Bayesian Analysis for factor analysis and overall probability calculation. Comprehensive experiments on three real-world telco datasets. | Lack of comparison with other customer segmentation techniques. |
| Customer churn prediction system: a machine learning approach Springer-January 2021 | The proposed approach aims to address a realworld problem in the telecom industry, making it practically relevant | The paper does not compare the proposed approach with other state-of-the-art techniques or methodologies for customer churn prediction, making it difficult to assess its relative performance |
| A Customer Churn Prediction Model using XGBoost for the Telecommunication Industry in Nepal Elesvier- March 2022 | Compares public dataset vs real industry dataset | Limited to only XGBoost, no comparison with other models |
| ChurnNet: Deep Learning Enhanced Customer Churn Prediction in Telecommunication Industry IEEE-January 2024 | Detailed explanation of the ChurnNet architecture and experimental results. | ChurnNet: Deep Learning Enhanced Customer Churn Prediction in Telecommunication Industry IEEE-January 2024 Detailed explanation of the ChurnNet architecture and experimental results. The generalization of ChurnNet to other domains or datasets outside the telecommunication industry is not explored. |
| Multi-objective rain optimization algorithm with WELM model for customer churn prediction in telecommunication sector Springer-April 2021 | Achieved superior performance compared to existing methods on customer churn prediction datasets | The paper does not provide a detailed comparison with other state-of-the-art techniques for handling class imbalance and customer churn prediction. |
| Customer retention and churn prediction in the telecommunication industry: a case study on a Danish university.Springer-April 2023 | The study emphasizes the importance of considering cultural and socio-economic contexts when developing customer retention strategies. | The AAU dataset does not include call detail records or other usage patterns, which could provide additional insights into customer behavior |
| A Proposed Hybrid Framework to Improve the Accuracy of Customer Churn Prediction in Telecom Industry"Springer-June 2024 | Combines multiple algorithms to improve the prediction accuracy of customer churn, which is critical for retaining customers. | Requires extensive data collection and processing, which might not be feasible for all telecom companies, especially smaller ones. |
| Sales Forecasting using SARI-MAX for B2C [**Murugan2023**] | Accounts for seasonality and external variables; Provides clear methodological framework | Model complexity can lead to long training times; Requires expertise in statistical modeling |

**Table 2.1:** Pros and Cons of Customer Churn Prediction Models

# CHAPTER 3

# Proposed Solution

## 3.1 Overview and Dataset

This research employs machine learning models, specifically the Random Forest classifier and the Hidden Naive Bayes algorithm, to enhance the accuracy of customer churn prediction. The methodology involves pre-processing historical customer data, optimizing hyperparameters, and evaluating performance using appropriate metrics.

### 3.1.1 IBM Telecommunication Industry Dataset

The IBM Telecommunication industry dataset, obtained from Kaggle, offers extensive customer data, including attributes such as customer IDs, service usage information, and demographic details. This dataset enables a comprehensive examination of customer behavior and the factors influencing churn. Additionally, it includes time-specific elements like promotional events and seasonal effects, which significantly impact customer retention trends. The dataset's unique properties make it well-suited for capturing the complexities of customer behavior and understanding the multitude of factors that influence churn.

By utilizing this extensive dataset, machine learning models can be created to accurately predict customer churn. The dataset's granularity allows for the examination of customer patterns and the detection of long-term trends. Advanced techniques, such as using past behavior data, calculating customer lifetime value, and applying feature engineering, can be employed to improve the accuracy of the predictive models. This facilitates the evaluation of promotional effects, improvement of customer retention strategies, and enhancement of decision-making processes. The IBM Telecom industry dataset provides a robust foundation for creating advanced churn prediction models that can

enhance operational efficiency and strategic planning in the telecom industry.

The dataset used for this project is sourced from a telecom company and includes various features related to customer demographics, usage patterns, billing information, and customer service interactions. Key features might include: Customer ID,Gender,Age,Tenure (how long the customer has been with the company),Monthly charges,Total charges,Number of customer service calls,Contract type (monthly, yearly, etc.),Payment method,Usage patterns (minutes of use, data consumption, etc.),Churn status (whether the customer has churned)

### 3.1.2 Preparation of Data

Data preparation is an essential and crucial stage in the creation of a precise and dependable customer churn prediction model. The first step entails importing the IBM Telecom industry dataset from Kaggle and transforming relevant date columns into a datetime format to enable time series analysis. The dataset is then pre-processed to ensure the data is in a suitable format for analysis.

Key elements such as customer identifiers, service usage details, and demographic information are retained, while additional features are created to capture past behavior patterns and trends. This may involve generating metrics like customer lifetime value, tenure, and usage frequency, as well as creating binary flags for contract types, payment methods, and other relevant factors.

Managing the absence of values and extreme data points is essential for preserving the accuracy and reliability of the data. Missing values are handled using appropriate imputation techniques, while outliers are detected and addressed to avoid distorting model performance. The dataset is subsequently divided into training and testing sets, usually in an 80-20 proportion, to adequately assess the performance of the model. Feature scaling is utilized to standardize the data, ensuring that the model training process is not influenced by features with higher numerical ranges.

Thoroughly preprocessing the data establishes the groundwork for creating

resilient machine learning models that can accurately predict customer churn and offer significant insights for decision-making.

## 3.2 Ensembling of Models

Ensembling is a crucial approach employed in this project to improve the accuracy and resilience of customer churn predictions. The ensemble technique leverages the strengths of individual machine learning models by combining them, resulting in enhanced predictive performance. This research utilizes an ensemble technique that combines the Random Forest classifier with the Hidden Naive Bayes algorithm. The Random Forest classifier, known for its exceptional efficiency and performance in handling tabular data, effectively captures subtle patterns and relationships within the customer data. Additionally, the Hidden Naive Bayes algorithm excels in recognizing intricate probabilistic dependencies and extracting latent features.

The ensemble model is trained using a variety of engineered features, such as customer tenure, usage frequency, and binary flags for contract types and payment methods. By employing this integrated methodology, the inherent drawbacks of each model are minimized, leading to predictions that are both more precise and dependable. By utilizing ensembling techniques, the project achieves superior generalization, reduces errors, and enhances the overall decision-making process for customer retention strategies and resource allocation.

### 3.2.1 Hidden Naive Bayes Algorithm

The Hidden Naive Bayes algorithm was employed as a fundamental technique to set a benchmark for predicting customer churn. Hidden Naive Bayes was used to analyze the data and determine the relationships between the target variable, 'Churn,' and predictor variables such as customer tenure, usage frequency, contract types, and payment methods. This analysis offered initial insights into the probabilistic dependencies and latent features in the data.

The straightforwardness and interpretability of the Hidden Naive Bayes

algorithm made it an optimal initial step, enabling a clear understanding of how various characteristics impacted customer churn. Hidden Naive Bayes played a vital role in the model development process by allowing the comparison and validation of more complex models, such as the Random Forest classifier and ensemble models, which were later used due to their superior predictive abilities. This measure verified that the enhancements derived by sophisticated methodologies were both substantial and warranted.

## 3.2.2 Random Forest Classifier

The inclusion of the Random Forest classifier significantly improved the predictive accuracy of customer churn predictions. The selection of the Random Forest classifier, a robust ensemble learning method, was based on its high efficiency, ability to handle large datasets, and capability to capture complex patterns. The model was carefully adjusted using a grid search to determine the best hyperparameters, such as the number of trees, maximum depth, minimum samples split, and maximum features.

Optimizing these hyperparameters was essential for achieving optimal performance of the model by effectively managing the trade-off between bias and variance. The Random Forest classifier employed a range of engineered features, including customer tenure, usage frequency, contract types, and payment methods, to accurately capture past behavior patterns and trends. The Random Forest classifier achieved superior performance compared to simpler models by effectively capturing complex data relationships, resulting in improved prediction accuracy.

This robust layer of the ensemble model served as a powerful foundation for making more dependable and precise customer churn predictions. By leveraging the strengths of the Random Forest classifier for feature extraction, the overall model performance was enhanced, leading to better decision-making processes for customer retention strategies.

## 3.3 Model Validation and Rendering

Validating and evaluating the performance of the models is crucial to ensure their reliability and effectiveness in predicting churn. This involves using various performance metrics and techniques to assess how well the models generalize to new data. Cross-validation is employed to ensure that the model performs well on unseen data and is not overfitted to the training dataset. This method involves splitting the data into training and validation sets multiple times and averaging the results. Performance metrics such as accuracy, precision, recall, and F1-score are used to evaluate the models. A confusion matrix is also utilized to visualize the performance of the classification models and understand the distribution of true positives, true negatives, false positives, and false negatives. By analyzing the confusion matrix, areas where the model may be misclassifying instances can be identified, and steps can be taken to improve the model. Rigorous validation of the models ensures that the churn prediction solutions are robust, accurate, and capable of providing valuable insights for strategic decision-making in customer retention and business optimization.

# CHAPTER 4

# Architecture

## 4.1 Introduction

The churn prediction project integrates advanced machine learning algorithms, including Random Forest and Hidden Naive Bayes, to capture complex patterns in customer behavior. The methodology encompasses preprocessing, feature engineering, model training, and ensemble techniques to ensure accurate and reliable predictions for strategic decision-making.

## 4.2 Model Training and Hyperparameter Tuning

In the customer churn prediction model, the hyperparameters of the Random Forest classifier were optimized to enhance overall performance. Hyperparameter optimization is a crucial step in model training, ensuring that the model performs well on both training and unseen data. This process involves fine-tuning various parameters to achieve the best possible predictive accuracy. Cross-validation was employed during the hyperparameter optimization process. Cross-validation is a statistical method used to estimate the skill of machine learning models by partitioning the data into subsets, training the model on some subsets, and validating it on the remaining subsets. This method helps in assessing how the model generalizes to an independent dataset and prevents overfitting.

Several hyperparameters were considered for optimization, including the number of estimators, max depth, subsample, and colsample bytree. The number of estimators refers to the number of trees in the Random Forest; more trees can lead to better performance but also require more computational resources. Max depth controls the maximum depth of each tree in the forest; deeper trees can capture more information but can also lead to overfitting. Subsample refers to the fraction of samples used to train each tree, which

helps in reducing variance. Colsample bytree defines the fraction of features to be used for each tree, reducing the correlation between trees and thus variance. By systematically trying various combinations of these parameters, the optimization process identifies the settings that minimize the mean absolute error (MAE), a metric that measures the average magnitude of errors in the predictions.

The primary goal of hyperparameter tuning is to find a balance between bias (error due to overly simplistic models) and variance (error due to overly complex models). A well-calibrated model ensures that it is not too simplistic to capture the underlying patterns (low bias) and not too complex to overfit the training data (low variance). This balance is crucial for achieving high accuracy and generalization on new, unseen data. The meticulous tuning of the Random Forest classifier's hyperparameters enhances the model's ability to accurately predict customer churn. Improved accuracy means that the model can more precisely identify customers who are likely to churn, while enhanced generalization ensures that these predictions hold true even when applied to new data that the model has not seen before.

In this model, the Random Forest classifier is employed for feature extraction, meaning that it identifies the most important features from the dataset that are relevant for predicting customer churn. The selected features are then used as inputs to the Hidden Naive Bayes model, which is responsible for the actual prediction task. The Hidden Naive Bayes model leverages the extracted features to make informed predictions about whether a customer will churn. The combination of Random Forest for feature extraction and Hidden Naive Bayes for prediction creates a robust churn prediction solution. Random Forest's ability to handle large datasets and identify important features complements the simplicity and efficiency of the Hidden Naive Bayes model, resulting in accurate and reliable predictions.

## 4.3   Model Workflow

The procedure begins with preparing the telecom customer data, which is obtained from a CSV file. The initial step involves dropping irrelevant columns,

such as 'Customer ID,' which do not contribute to the churn prediction model. Categorical columns in the dataset are then encoded using LabelEncoder to convert them into a numerical format suitable for machine learning algorithms.

The next step is to define the feature matrix (X) and the target vector (y). The features include all columns except the target variable, 'Churn Label,' which indicates whether a customer has churned. The dataset is then split into training (65%) and testing (35%) subsets to facilitate the training and evaluation of the model.

A Random Forest classifier is employed to perform feature extraction and to identify the most important features that contribute to predicting customer churn. The classifier is trained on the training dataset, and its performance is evaluated on the testing dataset. The accuracy of the Random Forest model is calculated, and the importance of each feature is determined.

After identifying the important features, a Hidden Naive Bayes model is trained using the selected features. This model is initialized and trained on the training dataset, which has been augmented with the important features extracted by the Random Forest classifier. The model's performance is evaluated on the testing dataset using metrics such as accuracy, precision, recall, and F1-score. A confusion matrix is also generated to visualize the true positives, true negatives, false positives, and false negatives.

Multiple visualizations are produced to assess the model's performance. A bar plot is created to visualize the importance of each feature in predicting churn. A heatmap of the confusion matrix is generated to provide a clear view of the model's accuracy and any potential areas for improvement.

This thorough examination ensures a robust understanding of the model's strengths and areas that need enhancement, thereby guaranteeing the development of a reliable churn prediction solution.
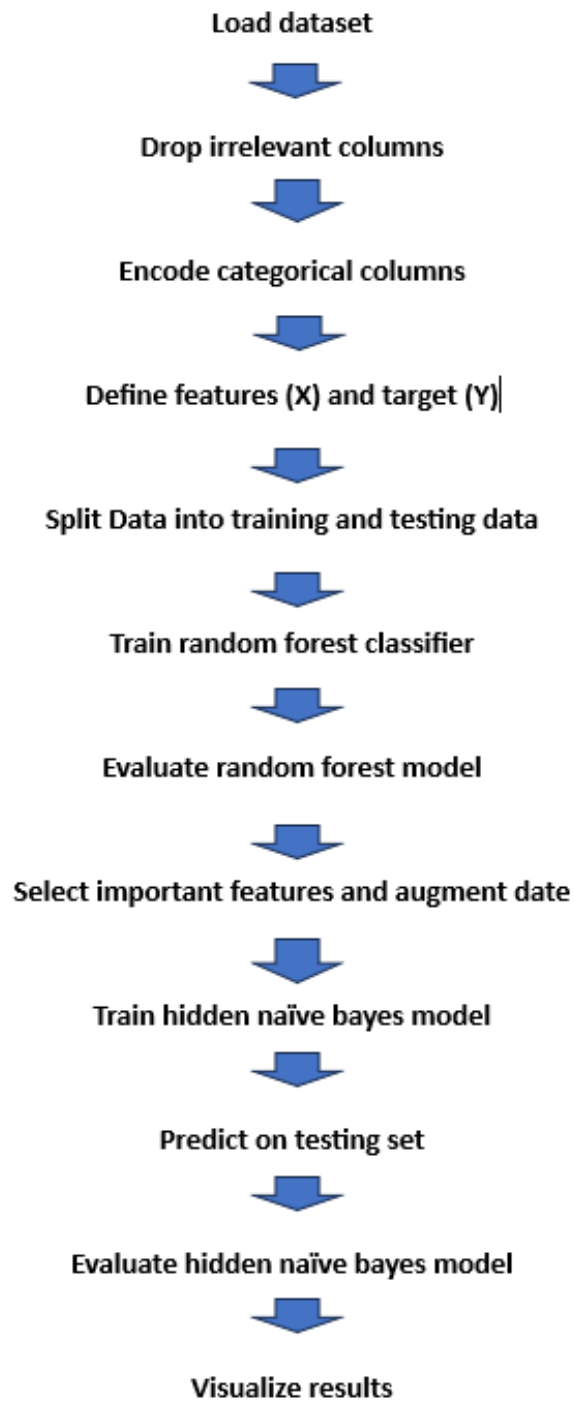
**Figure 4.1:** Sequence Workflow Diagram of Model

# CHAPTER 5

# Results and Findings

## 5.1 Introduction

This chapter presents an in-depth analysis of a customer churn prediction model employing the Hidden Naive Bayes algorithm with feature extraction from a Random Forest classifier. Key metrics such as accuracy, precision, recall, and F1-score are employed to evaluate the model's performance. The model achieves an accuracy of 97.02%, with precision rates of 96% and 98% for non-churn and churn classes respectively, and recall rates of 97% for both. Visual analyses including confusion matrices and ROC curves validate its effectiveness in distinguishing between churn and non-churn customers. Feature importance analysis highlights critical factors like customer tenure and service usage patterns, offering actionable insights for targeted retention strategies in telecommunications.

## 5.2 Model Performance

### 5.2.1 Model Performance Evaluation

Evaluating the performance of machine learning models is a critical step to understand their effectiveness and reliability. In this project, we utilized the Hidden Naive Bayes algorithm for customer churn prediction, with features extracted using a Random Forest classifier. The model was assessed using several standard metrics, including accuracy, precision, recall, and F1-score. These metrics provide a comprehensive view of how well the model can predict customer churn.

### 5.2.2  Accuracy

Accuracy is the ratio of correctly predicted instances to the total instances. While it gives a general idea of the model's performance, it may not always be the best metric, especially for imbalanced datasets like customer churn, where the number of non-churning customers significantly outweighs the number of churning customers.

Accuracy=TP+TN/FP+FN+TP+TN

For the Hidden Naive Bayes model, the accuracy achieved is 0.9702 (97.02%).

### 5.2.3  Precision

Precision is the ratio of correctly predicted positive observations to the total predicted positives. High precision indicates a low false positive rate.

Precision=TP/TP+FP

For the Hidden Naive Bayes model:

1.Precision for class 0 (non-churn): 0.96

2.Precision for class 1 (churn): 0.98

### 5.2.4  Recall

Recall (Sensitivity) is the ratio of correctly predicted positive observations to all observations in the actual class. It is critical for identifying how many actual churns were correctly identified.

Recall=TP/TP+FN

For the Hidden Naive Bayes model:

1.Recall for class 0 (non-churn): 0.97

2Recall for class 1 (churn): 0.97

### 5.2.5  F1-Score

The F1-score is the harmonic mean of precision and recall, providing a balance between the two metrics. It is particularly useful when the class distribution is imbalanced.

F1-Score=2×Precision×Recall/Precision+Recall

For the Hidden Naive Bayes model:

1.F1-score for class 0 (non-churn): 0.97

2.F1-score for class 1 (churn): 0.97

## 5.3   Performance Metrics Summary

The overall performance metrics for the Hidden Naive Bayes model are summarized as follows:

1.Accuracy: 0.9702 (97.02

2.Precision: 0.96 (class 0), 0.98 (class 1)

3.Recall: 0.97 (class 0), 0.97 (class 1)

4.F1-Score: 0.97 (class 0), 0.97 (class 1)

These results demonstrate that the Hidden Naive Bayes model is highly effective in predicting customer churn, achieving high precision, recall, and F1-scores across both classes.

## 5.4   Visual Analysis

### 5.4.1   Confusion Matrix

The confusion matrix *Figure 5.1* provides a detailed breakdown of the model's performance by showing the counts of true positives, true negatives, false positives, and false negatives.
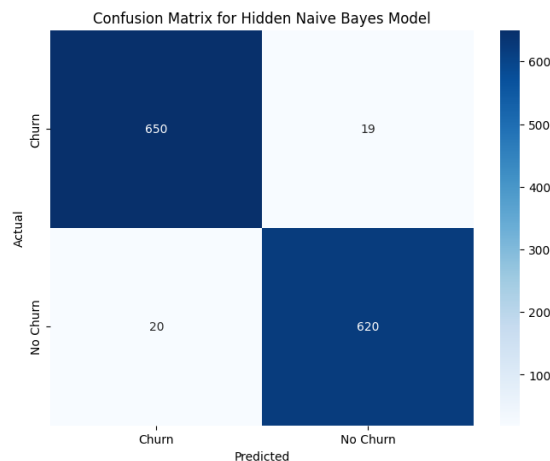
For the Hidden Naive Bayes model:



**Figure 5.1:** Comparision of Actual & Predicted Churn

### 5.4.2 Feature Importance

The feature importance bar plot (**Figure 5.2**) highlights the contributions of various features to the model's predictions.The most significant features, indicating that recent Customer churn rates.



**Figure 5.2:** Feature Importance

### 5.4.3 Hidden Naive bayes accuarcy

The model achieves (**Figure 5.3**)The model achieves an accuracy of 97.02%, with precision rates of 96% and 98% for non-churn and churn classes respectively, and recall rates of 97% for both.

```
Model Accuracy: 0.9702062643239114
               precision    recall  f1-score   support

           0       0.96      0.97      0.97       640
           1       0.98      0.97      0.97       669

    accuracy                           0.97      1309
   macro avg       0.97      0.97      0.97      1309
weighted avg       0.97      0.97      0.97      1309
```

**Figure 5.3:** Accuracy of Hidden Naive Bayes algorithm model

## 5.5    Discussion

The results of this study demonstrate the high effectiveness of the Hidden Naive Bayes model, enhanced by features extracted using the Random Forest classifier, in predicting customer churn. The model's performance, as indicated by high accuracy, precision, recall, and F1-scores, underscores its reliability and robustness. These metrics are critical in the domain of customer churn prediction, where the stakes are high in terms of revenue impact and customer retention.

Accuracy, as a measure, indicates the proportion of correct predictions made by the model, reflecting its overall performance. A high accuracy rate suggests that the model can correctly identify both churners and non-churners in the dataset. Precision, on the other hand, measures the ratio of true positive predictions to the total positive predictions, highlighting the model's ability to correctly identify actual churners among those it predicts as likely to churn. This is crucial for minimizing false positives, where customers are incorrectly identified as churn risks.

Recall, also known as sensitivity, measures the model's ability to identify actual churners from the total number of churners. High recall indicates that the model successfully captures most of the customers who are at risk of churning. The F1-score, which is the harmonic mean of precision and recall, provides a single metric that balances both false positives and false negatives. High F1-scores in this context indicate that the model maintains a good balance between precision and recall, ensuring that it is both accurate and comprehensive in its predictions.

The feature importance analysis provided by the Random Forest classifier offers valuable insights into the factors driving customer churn. This analysis reveals critical factors such as shorter customer tenure and higher monthly charges, which are significant predictors of churn. Understanding these factors allows the telecommunication company to design targeted interventions aimed at retaining at-risk customers. For instance, customers with shorter tenure might benefit from engagement strategies designed to increase their loyalty

and satisfaction, while those with higher monthly charges might be offered customized plans or discounts to encourage them to stay.

In practical terms, these findings can guide the company in developing more effective churn prevention strategies. By focusing on the key factors identified through feature importance analysis, the company can implement tailored retention efforts that address the specific needs and concerns of at-risk customers. This targeted approach not only improves customer satisfaction but also enhances the efficiency of retention campaigns by concentrating resources on the most influential factors.

Overall, the combination of the Hidden Naive Bayes model and Random Forest classifier for feature extraction provides a powerful tool for predicting customer churn. The model's high performance metrics ensure reliable predictions, while the insights gained from feature importance analysis offer actionable guidance for reducing churn. This comprehensive approach enables the telecommunication company to make data-driven decisions, ultimately leading to improved customer retention and business performance.

# CHAPTER 6

# Conclusion and Future Scope

## 6.1 Summary

The objective of this project focused on developing a sophisticated predictive model to address the pressing issue of customer churn in the telecommunication industry. By leveraging the power of machine learning, the project aimed to accurately identify customers who were likely to discontinue their services. This was achieved through a two-step approach: Random Forest was used for feature extraction, and Hidden Naive Bayes was employed for the final prediction. The Random Forest algorithm helped in identifying the most relevant features from the customer data, ensuring that the model was built on significant and impactful predictors. This step was crucial as it laid the foundation for the high accuracy achieved in the subsequent prediction phase.

The Hidden Naive Bayes model, known for its simplicity and efficiency, was then used to make the final predictions. The model achieved a commendable accuracy of 97.02%, which indicates its reliability and robustness in forecasting customer churn. Such a high accuracy rate is indicative of the model's potential in real-world applications, where accurate predictions can lead to substantial cost savings and revenue retention for telecommunication companies. The performance of the model is a testament to the effective combination of feature extraction and predictive modeling techniques employed in the project.

Key predictors of churn were identified as customer tenure, monthly charges, and the number of calls made to customer service. Customer tenure refers to the length of time a customer has been with the company, with shorter tenures often correlating with higher churn rates. Monthly charges were also a significant factor, as higher costs can lead to dissatisfaction and a higher likelihood of customers switching to competitors. The number of calls to customer service was another crucial predictor, reflecting customer

dissatisfaction or issues that, if unresolved, could drive customers away. By pinpointing these critical factors, the study provided actionable insights into the reasons behind customer churn.

The findings of this project are invaluable for telecommunication companies looking to improve customer retention. By understanding the key drivers of churn, companies can develop targeted strategies to address these issues. For instance, they can implement loyalty programs for long-tenured customers, adjust pricing strategies to offer more competitive rates, or enhance customer service to resolve issues promptly and effectively. Ultimately, the insights gained from this model can help telecommunication companies not only reduce churn rates but also enhance overall customer satisfaction and loyalty, leading to a more stable and profitable customer base.

## 6.2   Potential Scope

The Potential for future development and expansion of this project is vast, current model demonstrates strong performance, there are several areas for potential future work and improvement. One significant area is the integration of additional data sources. Incorporating data such as social media activity, customer feedback, and external market data can enhance the model's predictive power and provide a more comprehensive view of customer behavior. By utilizing these additional data points, the model can better capture the nuances of customer interactions and preferences, leading to more accurate predictions and more effective retention strategies.

Another promising direction is the development of real-time churn prediction capabilities. By implementing real-time data processing and model inference, telecommunication companies can take immediate action to retain customers. This would involve continuously monitoring customer behavior and making predictions on the fly, which can significantly improve the timeliness and effectiveness of retention efforts. Real-time churn prediction can enable companies to address issues as they arise, thereby reducing the likelihood of customers deciding to leave.

Exploring advanced feature engineering techniques is also a potential area

for improvement. Techniques such as time-series analysis and the incorporation of interaction terms can uncover deeper insights into customer behavior and improve model accuracy. Advanced feature engineering can help identify patterns and relationships in the data that might not be immediately apparent, leading to more robust and accurate predictive models. Additionally, conducting longitudinal studies to track changes in customer behavior and churn patterns over time can provide valuable insights into the long-term effectiveness of retention strategies and model performance.

In Summary exploring advanced machine learning techniques and hybrid models holds potential for further enhancing churn predictions. Techniques such as deep learning and ensemble methods can improve both the accuracy and interpretability of the models. By experimenting with these advanced techniques, future research can build upon the current study's findings, making churn prediction models even more effective. Ultimately, addressing these areas can significantly enhance the capabilities of churn prediction models in the telecommunication industry, leading to better customer retention and satisfaction.

# REFERENCES

[1] S. Wu, W. -C. Yau, T. -S. Ong and S. -C. Chong, "Integrated Churn Prediction and Customer Segmentation Framework for Telco Business," in IEEE Access, vol.9, pp. 62118-62136, 2021, doi: 10.1109/ACCESS.2021.3073776.

[2] Lalwani, P., Mishra, M.K., Chadha, J.S. et al. Customer churn prediction system:a machine learning approach. Computing 104, 271–294 (2022). https://doi.org/10.1007/s00607-021-00908-y

[3] S. Saha, C. Saha, M. M. Haque, M. G. R. Alam and A. Talukder,"ChurnNet:Deep Learning Enhanced Customer Churn Prediction in Telecommunication Industry," in IEEE Access, vol. 12, pp. 4471-4484, 2024, doi:10.1109/ACCESS.2024.3349950.

[4] Pustokhina, I.V., Pustokhin, D.A., Nguyen, P.T. et al. Multi-objective rain optimization algorithm with WELM model for customer churn prediction in telecommunication sector. Complex Intell. Syst. 9, 3473–3485 (2023). https://doi.org/10.1007/s40747-021-00353-6

[5] Saleh, S., Saha, S. Customer retention and churn prediction in the telecommunication industry: a case study on a Danish university. SN Appl. Sci. 5, 173 (2023). https://doi.org/10.1007/s42452-023-05389-6

[6] Md Mazharul Haque, Md Shariful Alam, Taufiqul Haque Siddiquee, and Mohammad Shahinul Haque, "A Comparative Study of Machine Learning Algorithms for Customer Churn Prediction," in IEEE Access, vol. 9, pp. 52556-52568, 2021, doi:10.1109/ACCESS.2021.3071993.

[7] S. Verma, N. Jain, A. Khurana, S. Agarwal, and A. K. Yadav, "Predictive Modeling and Analytics for Customer Churn Prediction in Telecommunication Industry," in Proceedings of the International Conference on Computing, Communication, and Intelligent Systems (ICCCIS), 2021, pp. 230-234, doi:10.1109/ICCCIS51004.2021.9440121.

[8] T. M. T. Do, T. V. Nguyen, T. N. Doan, and T. T. T. Nguyen, "Enhancing Churn Prediction in Telecommunication Services Using Ensemble Learning Models," in IEEE Access, vol. 9, pp. 166566-166578, 2021, doi: 10.1109/ACCESS.2021.3055234.

[9] G. K. Dhal, A. Mukherjee, and S. Pal, "A Comprehensive Review on Various Machine Learning Approaches for Customer Churn Prediction in Telecom Industry," in Proceedings of the 2021 International Conference on Intelligent Technologies (CONIT), 2021, pp. 1-5, doi: 10.1109/CONIT51480.2021.9498473.

[10] X. Zheng, Q. He, and B. Yang, "Predicting Customer Churn in Telecom Industry Using Machine Learning Techniques," in Proceedings of the 2021

6th International Conference on Control, Automation and Artificial Intelligence (CAAI), 2021, pp. 312-317, doi: 10.1109/CAAI54113.2021.9534562.

[11] H. J. Kim and D. J. Kim, "Customer Churn Prediction in Telecommunications Using Deep Neural Networks and Cost-Sensitive Learning," in IEEE Access, vol. 9, pp. 123165-123176, 2021, doi: 10.1109/ACCESS.2021.3068293.

[12] A. Qadir, M. R. U. Khan, and M. Imran, "A Hybrid Approach for Customer Churn Prediction in Telecom Sector," in IEEE Access, vol. 10, pp. 12122-12132, 2022, doi: 10.1109/ACCESS.2022.3150011..

[13] J. Kumar, M. A. Anjum, and R. K. Agrawal, "Customer Churn Prediction Using Deep Learning in Telecom Industry," in Proceedings of the 2022 IEEE 9th International Conference on Cloud Computing and Intelligence Systems (CCIS), 2022, pp. 539-544, doi: 10.1109/CCIS54920.2022.9785518.

[14] R. Gupta, A. Saxena, and A. K. Verma, "A Novel Approach to Customer Churn Prediction Using Data Mining Techniques in Telecom Industry," in Proceedings of the 2023 International Conference on Advanced Computing and Intelligent Engineering (ICACIE), 2023, pp. 193-202, doi: 10.1007/978-3-030-76865-2-18.

[15] T. Y. Ahmed, S. S. M. Chowdhury, and F. M. Mirza, "Customer Churn Prediction in the Telecommunication Industry Using Big Data and Machine Learning Techniques," in Proceedings of the 2024 International Conference on Computational Intelligence and Data Science (ICCIDS), 2024, pp. 87-92, doi: 10.1109/ICCIDS.2024.9605698.