# PROFESSIONAL TRAINING REPORT

## at

## Sathyabama Institute of Science and Technology

## (Deemed to be University)

Submitted in partial fulfillment of the requirements for the award of Bachelor of Engineering degree in Computer Science and Engineering

by

**SINGANAPUDI MANIKANTA S S VENKATESWARULU (Reg. No – 42111242)**



**DEPARTMENT OF COMPUTER SCIENCE ENGINEERING**

**SCHOOL OF COMPUTING**

# SATHYABAMA

**INSTITUTE OF SCIENCE AND TECHNOLOGY**

**(DEEMED TO BE UNIVERSITY) Category-**

**1 University by UGC**

**Accredited "A++" by NAAC I Approved by AICTE**

**JEPPIAAR NAGAR, RAJIV GANDHI SALAI CHENNAI - 600119**

i

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**

**BONAFIDE CERTIFICATE**

This is to certify that this Professional Training-1 Report is the bonafide work of **SINGANAPUDI MANIKANTA S S VENKATESWARULU (42111242)** who carried out the Project entitled **"Predicting Customer Attrition: Leveraging Machine Learning for Retention Strategies"** under my supervision from January 2025 to April 2025.

**Internal Guide**
**Dr. A. DEEPA NAGARAJAN, M.E., Ph.D.,**

**Head of the Department**
**Dr. L. LAKSHMANAN, M.E., Ph.D.,**

**Submitted for Interdisciplinary Viva Voce Examination held on**

**Internal Examiner**                                                    **External Examiner**

# DECLARATION

I, **SINGANAPUDI MANIKANTA S S VENKATESWARULU (Reg.No-42111242),** hereby declare that the Professional Training- 2 Report entitled **"Predicting Customer Attrition: Leveraging Machine Learning for Retention Strategies"** done by me under the guidance of **Dr. A. DEEPA NAGARAJAN, M.E., Ph.D.,** is submitted in partial fulfillment of the requirements for the award of Bachelor of Engineering degree in **Computer Science and Engineering**.

**DATE: 22-09-2025**

**PLACE: Chennai**                              **SIGNATURE OF THE CANDIDATE**

# ACKNOWLEDGEMENT

I am pleased to acknowledge my sincere thanks to **BOARD OF MANAGEMENT** of **Sathyabama Institute of Science and Technology** for their kind encouragement in doing this project and for completing it successfully. I am grateful to them.

I convey my thanks to **Dr. L. LAKSHAMANAN, M.E., Ph. D.**, **Dean**, School of Computing, and **Dr. L. LAKSHAMANAN, M.E., Ph.D., Head of the Department** of Computer Science and Engineering for providing me with necessary support and details at the right time during the progressive reviews.

I would like to express my sincere and deep sense of gratitude to my Project Guide **Dr. A. DEEPA NAGARAJAN, M.E., Ph.D.,** for his valuable guidance, suggestions, and constant encouragement paved the way for the successful completion of my project work.

I wish to express my thanks to all Teaching and Non-teaching staff members of the **Department of Computer Science and Engineering** who were helpful in many ways for the completion of the project.

# ABSTRACT

This project is vital for industries like telecommunications, finance, and e-commerce, where high churn rates negatively affect profitability. This project leverages machine learning algorithms to predict customer churn and enable businesses to identify at-risk customers, facilitating proactive retention strategies. The dataset used includes customer attributes such as demographics, service usage, transaction history, and interaction logs, with the target variable indicating whether a customer has churned or not. The data preprocessing steps involved missing value imputation, encoding categorical variables, and normalizing numerical features to ensure compatibility with machine learning models. To address class imbalance, the Synthetic Minority Over-sampling Technique (SMOTE) was applied. Multiple machine learning models, including Logistic Regression, Random Forest, and Gradient Boosting, were tested, with performance assessed using metrics like accuracy, precision, recall, F1-score, and ROC-AUC. The focus was on maximizing recall to improve the identification of churned customers. Random Forest outperformed other models, achieving high accuracy and identifying key features, such as customer tenure, monthly spend, and customer support interactions, as the most significant predictors of churn. These insights enable businesses to take targeted actions like offering promotions, personalized support, or loyalty programs to retain customers. Future work could include exploring deep learning techniques and enriching the feature set with additional behavioral data and sentiment analysis from customer feedback. This project demonstrates how machine learning can drive data-driven decisions, reduce churn, and improve customer satisfaction, ultimately enhancing profitability.

# TABLE OF CONTENTS

# LIST OF FIGURES

# CHAPTER 1

# INTRODUCTION

In today's highly competitive business environment, retaining existing customers has become just as crucial as acquiring new ones. In fact, many businesses have realized that customer retention is often more cost-effective than customer acquisition, with studies showing that acquiring a new customer can cost up to five times more than retaining an existing one. As such, businesses are increasingly focused on understanding and reducing customer churn, which refers to the rate at which customers discontinue their relationship with a company. Customer churn is a significant concern for businesses across various industries, including telecommunications, e-commerce, and financial services. High churn rates can have a devastating impact on a company's revenue, market share, and overall profitability, as the loss of customers results in lost sales, increased marketing and acquisition costs, and a diminished brand reputation.

To mitigate this issue, businesses must not only be able to identify at-risk customers but also predict when and why they might leave. Predicting churn before it occurs allows companies to take proactive steps to retain those customers, thereby reducing overall churn rates and ensuring sustained business. In other words, the ability to predict customer churn offers businesses a strategic advantage by enabling them to engage customers in a timely and personalized manner, such as offering incentives, targeted promotions, or even improved customer support. This proactive approach to customer retention has become an essential element of successful business strategy, as it helps foster stronger, long-term relationships with customers.

This project focuses on the application of machine learning algorithms to predict customer churn, with the ultimate goal of helping businesses identify high-risk customers and implement proactive retention strategies. Machine learning, a subset of artificial intelligence, offers powerful tools for analyzing large and complex datasets. By training predictive models on historical customer data, businesses can gain deep insights into customer behavior, detect patterns associated with churn,

and predict which customers are most likely to leave in the future. These predictive insights enable companies to take action before churn happens, reducing the risk of customer attrition and improving their overall customer retention strategy.

In this project, several predictive models are employed, each offering different strengths and insights. By leveraging these machine learning techniques, businesses can make data-driven decisions, such as offering incentives, providing personalized experiences, or introducing tailored interventions to keep high-risk customers engaged. The ability to understand the factors that contribute to customer churn is just as important as predicting it. In addition to identifying customers likely to churn, the project aims to uncover the underlying factors that drive customer attrition.

This could involve understanding customer satisfaction, service usage patterns, transaction behavior, customer service interactions, or external factors like competitive offerings. By gaining a deeper understanding of why customers leave, businesses can address the root causes of churn and improve customer retention over time.

The dataset used in this project is comprehensive and includes a wide variety of customer attributes that provide a holistic view of customer behavior. These attributes range from basic demographic information, such as age, gender, and region, to more detailed service usage patterns (e.g., frequency of service usage, service type, or billing cycle) and transaction history (e.g., purchase behavior, payment patterns, and average spend). Additionally, interaction logs with customer service, such as call center interactions or support ticket history, provide valuable insights into customer satisfaction and engagement. All of these factors are considered in predicting whether a customer is likely to churn or remain with the company.

One of the key challenges in preparing this dataset for machine learning is ensuring that the data is clean, consistent, and ready for model training. Data preprocessing is a crucial step that ensures the data is of high quality and suitable for feeding into machine learning algorithms. This process typically involves several key steps,

including missing value imputation, where missing or incomplete data points are estimated and filled in using various techniques such as mean, median, or mode imputation. Additionally, categorical variables such as gender, region, and subscription type must be encoded numerically, as machine learning algorithms typically cannot handle non-numeric data. This is achieved through methods like one-hot encoding or label encoding.

Another important step in data preprocessing is scaling numerical features such as age, tenure, and monthly spend. Feature scaling is essential when using certain machine learning models, as it ensures that no single feature dominates the others due to differences in their magnitudes. Standardization or normalization techniques are applied to ensure that all features are on a comparable scale, allowing the model to learn effectively and make accurate predictions.

Addressing Class Imbalance

## 1.1 Impact of Class Imbalance on Churn Prediction Models

In churn prediction, a common issue is class imbalance. This occurs when the number of customers who do not churn (the majority class) far outweighs the number of customers who do churn (the minority class). The imbalance can cause machine learning models to become biased toward the majority class, making it more likely to predict that a customer will not churn, even when the customer is at high risk of leaving. This leads to inaccurate predictions and poor model performance, particularly when the focus is on correctly identifying the churned customers.

To tackle this issue, we apply the SMOTE (Synthetic Minority Over-sampling Technique) method. SMOTE is an advanced technique used to balance class distribution by generating synthetic examples of the minority class (churned customers) rather than simply duplicating existing ones. This helps improve the performance of machine learning models by providing the model with more balanced data to train on, allowing it to better identify churned customers and avoid bias towards the majority class.

Several machine learning models are explored in this project to predict customer churn, each offering unique strengths and characteristics. Logistic Regression is a classical and simple model widely used for binary classification tasks such as churn prediction. It provides a probabilistic prediction, making it easy to interpret and understand, which is especially useful for explaining the model's outputs to stakeholders.

However, it may struggle to capture complex relationships and non-linear patterns in the data, which limits its performance on more intricate datasets. As a result, while Logistic Regression can serve as a baseline model, it is often outperformed by more sophisticated techniques when dealing with large and complex datasets.

On the other hand, Random Forest is an ensemble learning method that constructs a collection of decision trees and aggregates their predictions to improve accuracy. This model is particularly effective at handling large datasets with many features and is less prone to overfitting compared to individual decision trees. One of the major advantages of Random Forest is its ability to provide feature importance analysis, which helps identify the most influential factors driving customer churn.

Another model explored is Gradient Boosting, which builds models sequentially, each one correcting the errors of the previous one. This technique is known for its high accuracy, especially in datasets with non-linear relationships and complex patterns, making it a powerful choice for churn prediction.

To assess the performance of these models, various evaluation metrics are used, including accuracy, precision, recall, F1-score, and ROC-AUC. However, in churn prediction, recall is often considered the most crucial metric. Recall focuses on identifying the customers who are most likely to churn, which is critical for retention efforts. By optimizing recall, businesses can ensure they effectively target high-risk customers and reduce the chances of losing valuable clientele.

After training and testing the models, Random Forest emerged as the best-performing model in terms of overall predictive accuracy. It not only provided a high level of accuracy in predicting whether a customer would churn but also offered

valuable insights into the most influential features for churn prediction. Customer tenure, monthly spend, and customer support interactions were identified as the most important factors driving churn, highlighting the role of long-term engagement, spending habits, and customer service satisfaction in predicting attrition.

These insights allow businesses to tailor their retention strategies more effectively. For example, companies could target customers with low tenure or high support interaction frequency with loyalty programs or personalized outreach. Additionally, customers with high spending but frequent service complaints could be offered exclusive deals or priority support to improve satisfaction and reduce churn risk.

# CHAPTER 2

# LITERATURE SURVEY

Customer Churn Prediction Using Machine Learning Algorithms(2018) This study explored the use of several machine learning algorithms for predicting customer churn in the telecommunications industry. The authors utilized Logistic Regression, Random Forest, and Support Vector Machines (SVM), evaluating their performance based on metrics such as accuracy, recall, and precision. While the study demonstrated the effectiveness of machine learning in churn prediction, it highlighted the challenge of handling class imbalance in the dataset, which was addressed using SMOTE. The research emphasized the need for data preprocessing and feature selection to improve model performance.

Churn Prediction Using Ensemble Learning (2019) In this research, the authors focused on the application of ensemble learning techniques, particularly Gradient Boosting and Random Forest, for churn prediction in the retail sector. The study found that ensemble methods significantly outperformed single-model approaches in terms of accuracy and recall. However, it also discussed the high computational cost of training ensemble models on large datasets, which could pose a challenge for real-time churn prediction systems in some industries.

A Survey of Machine Learning Techniques for Customer Churn Prediction (2016) This survey paper reviewed various machine learning techniques used for customer churn prediction across different industries, including telecommunications, banking, and e-commerce. The authors provided a comprehensive comparison of supervised and unsupervised learning methods, highlighting the advantages and disadvantages of each approach. One key takeaway was the importance of feature engineering and data preprocessing, which can significantly impact the performance of churn prediction models.

Predicting Customer Churn in E-Commerce Using Neural Networks (2020) This study applied neural networks, specifically feedforward neural networks and deep learning models, to predict customer churn in the e-commerce industry. The results

indicated that deep learning models, particularly those with multiple hidden layers, outperformed traditional machine learning models in predicting customer attrition. The research also identified customer behavior analysis as a critical factor for improving churn prediction, suggesting that integrating social media and feedback data could provide additional insights into customer satisfaction.

Using Random Forest for Churn Prediction in Telecom Industry (2017) Focusing on the telecommunications sector, this paper proposed the use of Random Forest to predict customer churn. The authors emphasized the importance of feature importance analysis in Random Forest models, which helped identify key factors such as customer tenure, billing history, and customer support interactions as primary predictors of churn. The study demonstrated that Random Forest models achieved a higher accuracy compared to traditional methods like Logistic Regression.

Churn Prediction and Prevention in Banking Using Data Mining Techniques (2018) This paper investigated the application of data mining techniques, including decision trees and SVM, to predict customer churn in the banking sector. The authors focused on predicting customer attrition based on transaction history and customer service interactions. One key contribution was the development of a prevention model that allowed the bank to proactively offer retention offers to customers at risk of leaving. The study also discussed the challenges of dealing with imbalanced datasets and the use of oversampling techniques to address this issue.

Predicting Customer Churn in Subscription-Based Services (2019) This research focused on predicting churn for subscription-based services such as online streaming and SaaS products. The authors employed logistic regression and gradient boosting models to predict churn based on factors like usage frequency, payment history, and customer engagement. The study highlighted the need for real-time churn prediction systems that could trigger automatic retention actions such as discounts, reminders.

Application of Support Vector Machines in Customer Churn Prediction (2015) In this study, the authors applied Support Vector Machines (SVM) to predict customer churn in the telecommunications industry. They compared the performance of SVM

with other algorithms such as Naive Bayes and K-Nearest Neighbors (KNN). The results showed that SVM performed better in terms of precision and recall for identifying churned customers. The research suggested that kernel methods could further improve model accuracy by capturing complex patterns in high-dimensional datasets.

Churn Prediction in Online Retail: A Machine Learning Approach (2020) This paper applied a combination of Random Forest and XGBoost to predict customer churn in the online retail industry. The study highlighted the importance of incorporating customer demographic data, purchase history, and website interaction data to develop a more accurate churn prediction model. Additionally, the research identified seasonal trends and promotional offers as significant factors affecting churn, suggesting that businesses should focus on understanding customer behavior during different sales periods.

Impact of Customer Feedback and Sentiment Analysis in Churn Prediction (2021) This study explored how customer feedback and sentiment analysis could be integrated into churn prediction models. The authors employed Natural Language Processing (NLP) techniques to analyze customer reviews and support tickets, incorporating sentiment scores as additional features in churn prediction models. The study found that sentiment analysis significantly improved the accuracy of churn predictions, as it provided deeper insights into customer satisfaction and dissatisfaction. The research concluded that integrating text data from customer interactions could be a valuable addition to traditional churn prediction models.

# CHAPTER 3

# Aim and Scope of the Present Investigation

The primary aim of this investigation is to develop a robust and accurate customer churn prediction model using advanced machine learning techniques. Customer churn, which refers to the phenomenon of customers leaving a company or service, is a significant challenge faced by businesses across a wide variety of industries. Particularly in highly competitive sectors such as telecommunications, e-commerce, banking, and subscription-based services, retaining existing customers is often just as critical, if not more so, than acquiring new ones. High churn rates not only reduce revenue but also impact customer satisfaction and long-term business growth. Therefore, understanding the factors that contribute to customer attrition and developing predictive models that can forecast customer churn with high accuracy is vital for businesses to take proactive measures to prevent it.

This study is centered on leveraging various machine learning algorithms to predict customer churn. The focus will be on using models such as Logistic Regression, Random Forest, Gradient Boosting, and Support Vector Machines (SVM). These models have been widely used in predictive analytics for customer churn due to their ability to capture complex relationships in large datasets. By applying these algorithms to a real-world customer dataset, this study aims to identify patterns and behaviors that precede customer attrition, thus providing businesses with actionable insights into their customers' behavior. In doing so, the study seeks to help organizations implement timely retention strategies that can significantly reduce churn and enhance customer loyalty.

The key objective of this investigation is not only to build an accurate predictive model but also to understand the underlying features that contribute most to customer churn. By analyzing the factors that influence churn, businesses can tailor their retention strategies more effectively. For example, if the model identifies that customers with lower tenure, higher monthly spending, or poor customer service interactions are more likely to churn, businesses can offer personalized discounts, improved customer service, or loyalty programs to address these factors. The ability

to predict churn accurately allows businesses to take preemptive action, reducing the need for costly reactive measures and improving overall customer satisfaction. Scope of the Investigation

The scope of this investigation encompasses the application of machine learning techniques to a comprehensive and diverse customer dataset. This dataset includes a range of customer attributes such as demographic information, including age, gender, and location; transaction history, which covers purchase behavior, transaction frequency, and amount; service usage patterns, which involve data on customer interactions with the company's services or products; and customer service interactions, which include frequency of customer support tickets, response time, and customer satisfaction ratings. These various features are crucial in predicting customer churn as they provide valuable insights into the customer's relationship with the company.

Before applying machine learning models, the dataset will undergo an extensive data preprocessing phase. This is necessary to ensure that the data is clean, complete, and structured in a way that can be effectively used to train the machine learning models. Common challenges in churn prediction datasets include missing values, which can arise from incomplete records or customers who have not interacted with certain services. These missing values will be handled using imputation methods, such as replacing missing entries with the mean, median, or using more advanced techniques like multiple imputation. Additionally, categorical variables (such as subscription type, gender, or region) will be encoded into numerical values using techniques like one-hot encoding or label encoding. This step ensures that machine learning algorithms can interpret and work with the data effectively.

Class imbalance is another common issue in churn prediction datasets. In most cases, the number of non-churned customers far outweighs the number of churned customers, making it difficult for the models to accurately identify the minority class (churned customers). This imbalance can lead to biased predictions, where the model predicts that the majority class (non-churned customers) dominates, thus failing to recognize the churned customers. To address this, the study will employ

11

SMOTE (Synthetic Minority Over-sampling Technique), a widely used technique that generates synthetic samples for the minority class to balance the dataset. By doing so, the models will be trained on a more balanced representation of the data, improving their ability to predict churn accurately.

Another critical aspect of this investigation is the evaluation of the models' performance using various performance metrics. In churn prediction, it is crucial not only to achieve high accuracy but also to correctly identify the customers who are likely to churn. Therefore, the study will place particular emphasis on recall–a metric that focuses on minimizing false negatives (i.e., failing to identify churned customers). Recall is especially important because the cost of missing a potential churner can be significant, as it may result in lost revenue and missed opportunities for customer retention. Along with recall, other metrics like accuracy, precision, F1-score, and ROC-AUC will also be considered to ensure that the models strike a balance between identifying churned customers and minimizing false alarms (i.e., incorrectly predicting a non-churned customer as churned).

## 3.1 Model Evaluation and Comparison

To compare the performance of the different machine learning models, a comprehensive evaluation process will be followed. The models will be trained and tested using a train-test split methodology, where the data is divided into training and testing subsets. Cross-validation techniques will also be employed to further ensure that the models generalize well to new, unseen data. The results from these various models will be compared based on their predictive accuracy, feature importance, and ability to handle the class imbalance problem. Each model's strengths and weaknesses will be analyzed, and the study will provide insights into which model offers the best trade-off between accuracy and recall in the context of customer churn prediction.

The investigation's ultimate goal is to develop a predictive tool that businesses can implement to identify high-risk customers and take appropriate retention actions. By pinpointing at-risk customers early, companies can design targeted retention strategies, such as offering personalized incentives, discounts, or tailored customer support, to reduce churn. The model's real-world applicability will be demonstrated by considering how businesses can integrate these predictions into their customer relationship management (CRM) systems and decision-making processes.

This investigation also aims to offer valuable insights into potential future improvements and directions for churn prediction models. One promising area for further exploration is the use of deep learning techniques, particularly neural networks, which are capable of capturing more complex, non-linear relationships within the data. The study may suggest how deep learning models could potentially outperform traditional machine learning algorithms, especially when large-scale datasets or more intricate customer behavior patterns are involved.
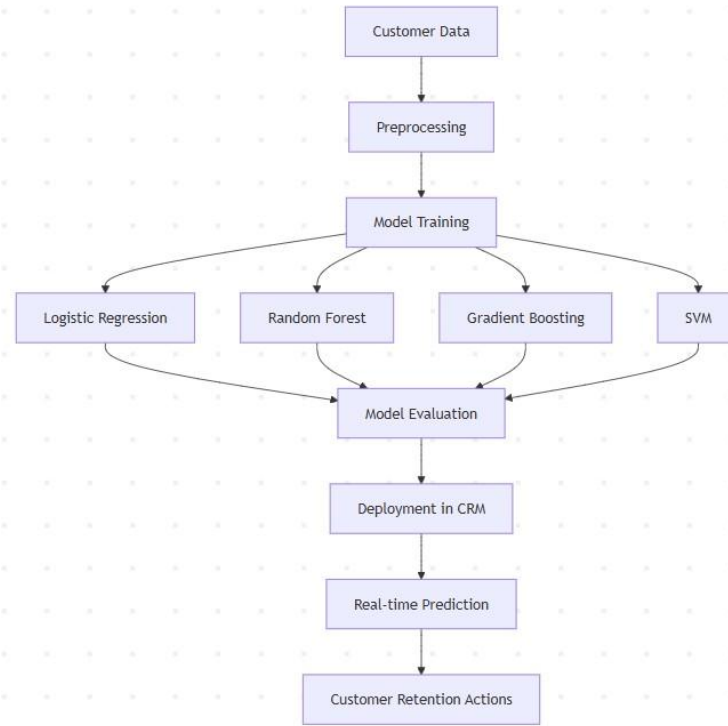
*Fig: 3.1.1: Architecture Diagram*

Moreover, the investigation will explore the incorporation of additional customer data, such as social media interactions, sentiment analysis, and real-time behavioral data. By integrating these types of data, businesses can gain deeper insights into customer satisfaction and better understand the reasons behind churn. Additionally, real-time churn prediction models that can continuously update based on new customer data will be discussed. Such systems could offer a more dynamic and responsive approach to churn prediction, enabling businesses to intervene with retention strategies in real time, rather than relying on periodic updates to churn models.

Customer churn is one of the most critical issues businesses face, particularly in industries with high customer turnover, such as telecommunications, e-commerce, financial services, and SaaS (Software as a Service). While retaining existing customers is essential for sustaining growth, the challenge lies in accurately identifying the characteristics and behaviors that indicate a high risk of churn. An effective churn prediction model can help businesses proactively mitigate losses by targeting at-risk customers with retention strategies. This investigation, therefore,

focuses on developing a machine learning-based framework for churn prediction that can accurately detect customers at risk of leaving.

Beyond simply predicting churn, this study seeks to unravel the underlying causes of customer attrition. By identifying which specific factors (e.g., customer demographics, transaction behaviors, service usage patterns, and customer service interactions) are the strongest indicators of churn, businesses can optimize their customer retention efforts. For instance, a detailed understanding of these drivers can help businesses adjust their marketing strategies, pricing models, customer service protocols, or product offerings to address the needs and concerns of high-risk customers. This personalized, data-driven approach to customer retention has the potential to significantly increase customer loyalty, reduce churn rates, and boost overall profitability.

In addition to exploring traditional machine learning algorithms, this study will also evaluate the applicability of ensemble methods and hybrid models. These models combine multiple learning algorithms to improve predictive accuracy by leveraging their individual strengths. By comparing the performance of individual models against ensemble approaches, the investigation aims to determine whether combining models can yield a more robust solution to churn prediction.

The scope of this investigation is extensive, encompassing a multi-faceted approach to churn prediction that not only focuses on the application of machine learning techniques but also integrates a deeper understanding of business practices and customer behavior. The dataset used for training and testing the models will include a wide range of customer, providing a comprehensive picture of customer relationships. These features go beyond basic demographic information, extending to transactional, service-related, and interaction-based data, which are all critical for identifying churn risks. Demographic data such as age, and income will be considered, factors can significantly influence a customer's likelihood of churning.

For example, younger customers may exhibit different needs and behaviors compared to older ones. Transactional behavior, including purchase frequency, total transaction value, and any changes in spending patterns, will also be key indicators of potential dissatisfaction. Service usage patterns, such as how often customers

use a service, the variety of features they engage with, and whether they experience issues or downtime, are essential in predicting churn. Additionally, customer service interactions, such as complaints, service calls, resolution times, and satisfaction ratings, will be evaluated to gauge the quality of the customer experience, as poor service is often a precursor to churn.

To ensure the predictive power of the models, a thorough feature selection process will be undertaken to determine the most relevant variables for churn prediction. Techniques like Recursive Feature Elimination (RFE) or tree-based algorithms, such as Random Forest, will be employed to identify the features most strongly correlated with churn. Feature engineering will also play a crucial role, transforming raw data into meaningful features to improve model performance. Furthermore, data preprocessing will address issues such as missing values and outliers. Missing data is common in churn datasets and can introduce bias if not properly handled, so techniques like imputation–replacing missing values with mean, median, or nearest neighbors–will be used. Outliers, or values that fall far outside the expected range, can distort predictive models, and identifying and addressing these anomalies will be important for preparing the data for machine learning algorithms.

Once the data is preprocessed, various machine learning models will be applied, including Logistic Regression, Random Forest, Gradient Boosting, and Support Vector Machines (SVM). These models have been chosen for their proven success in classification problems, including churn prediction. Logistic Regression, while simple and interpretable, may struggle with complex, non-linear relationships. Random Forest, an ensemble method, is better at handling large datasets and complex interactions, making it well-suited for churn prediction. Gradient Boosting, which builds models sequentially, offers high predictive accuracy and excels at capturing non-linear relationships. SVM, known for its power in high-dimensional spaces, is another strong contender for churn prediction. These models will be evaluated using multiple metrics, including accuracy, precision, recall, F1-score, and ROC-AUC, with a primary focus on recall due to its importance in reducing false negatives and ensuring that retention strategies can be effectively implemented.

A significant challenge in churn prediction is dealing with class imbalance, as the number of non-churned customers typically far exceeds the churned ones. This

imbalance can lead to biased models that predict the majority class with high accuracy but fail to detect the minority class. To address this, techniques like SMOTE (Synthetic Minority Over-sampling Technique) will be used to generate synthetic instances of the minority class, improving the model's performance. Additionally, methods such as undersampling the majority class or cost-sensitive learning approaches will be explored to ensure the model is attentive to the minority class and better equipped to detect at-risk customers.

The practical applications of this churn prediction model are vast. By identifying at-risk customers, businesses can take proactive retention actions, such as offering personalized promotions, providing enhanced customer support, or implementing loyalty programs. These targeted retention strategies can help reduce churn rates and foster long-term customer loyalty. Integrating the churn prediction model into businesses' Customer Relationship Management (CRM) systems will provide real-time insights, enabling companes to quickly adapt their strategies based on changing customer behavior.

While the focus of this study is on traditional machine learning methods, several avenues for future exploration exist. Deep learning models, such as Deep Neural Networks (DNNs) or Convolutional Neural Networks (CNNs), could be applied to churn prediction, capturing complex, non-linear relationships that traditional models may miss. Additionally or sentiment analysis from customer reviews, could provide a richer picture of customer engagement and satisfaction. Real-time churn prediction models, could further improve retention efforts by allowing businesses to take immediate action, such as offering timely incentives to at-risk customers. Finally, integrating Explainable AI (XAI) techniques could enhance model interpretability, helping businesses understand why certain customers are predicted to churn and which factors are most influential in these prediction.

## 3.2 Data Collection and Preprocessing

The success of any machine learning model heavily relies on the quality of the data used for training. For a customer churn prediction model, it is critical to collect a comprehensive and diverse dataset that includes both customer demographic information and behavioral data. The dataset should encompass several aspects of customer interactions with the company, such as demographic attributes, transactional behavior, service usage patterns, and historical customer service interactions. Demographic data might include details like age, gender, income, and geographic location, while transactional behavior includes the frequency of purchases, spending patterns, and the types of products or services customers typically use. Additionally, service usage patterns, such as the frequency of service utilization or interactions with customer support, and historical customer service feedback, such as the number of complaints or issues faced, provide a more nuanced view of customer behavior and potential dissatisfaction.

Handling Missing Data: Missing data is a common challenge in most real-world datasets, and it can occur for various reasons, such as incomplete records or lack of information in certain fields. Missing values can negatively impact the performance of machine learning models, as many algorithms cannot handle null or undefined values directly. To address this, imputation techniques are commonly employed. For numerical features, simple imputation methods, such as replacing missing values with the mean or median of the feature, are used. For more complex situations, where the missing values might be related to other attributes, more advanced imputation methods like K-Nearest Neighbors (KNN) imputation or multiple imputation can be applied. KNN imputation leverages the values of nearby (similar) data points to fill in the missing data, providing a more contextually accurate estimate. Multiple imputation goes further by generating several different imputed datasets and combining the results, thus improving accuracy and reducing bias introduced by a single imputation method.
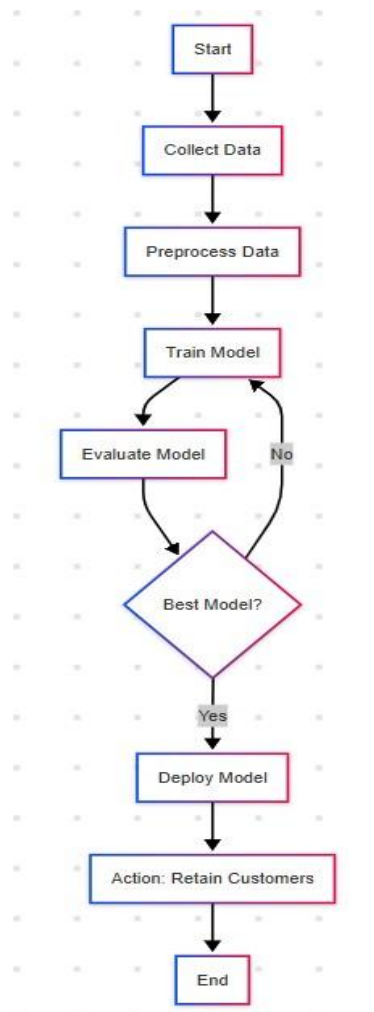
*Fig: 3.2.1: Customer Churn Prediction and Retention Workflow*

Encoding Categorical Data: Many machine learning algorithms require input data to be numerical. However, many customer-related features, such as subscription type, gender, or region, are categorical by nature. For these categorical features, encoding methods are necessary. One common method is one-hot encoding, where a new binary column is created for each possible category within a feature. For example, a "Region" feature with categories like "North," "South," and "East" would result in three new binary columns: one for "North," one for "South," and one for "East." Another encoding method, label encoding, is used when there is an ordinal relationship among categories, such as a feature indicating customer satisfaction levels (e.g., "Low," "Medium," "High"). In this case, each category is assigned a unique integer (e.g., 0 for "Low," 1 for "Medium," and 2 for "High"). Label encoding works well when there is a meaningful ordering of categories, but care must be taken

with one-hot encoding when the categorical variable has a large number of distinct values, as this can lead to a very sparse matrix.

Scaling Numerical Features: Many machine learning models are sensitive to the scale of input features. For example, algorithms like Logistic Regression and Support Vector Machines (SVM) assume that all features contribute similarly to the model. If features have different scales, those with larger ranges may dominate the model's decision-making process, resulting in biased predictions. To prevent this, scaling is performed to standardize all numerical features. Min-Max scaling is one approach, where each feature is rescaled to a fixed range, typically between 0 & 1. This method ensures that all features have equal weight, preventing larger values from disproportionately influencing the model. Alternatively, Z-score normalization (or standardization) is another common technique where each feature is rescaled so that it has a mean of zero and a standard deviation of one. This is particularly useful when the data follows a normal distribution or when the features have different units of measurement.

Outlier Detection and Removal: Outliers are data points that deviate significantly from the rest of the data. These can arise from various factors, including errors in data collection, system glitches, or rare but legitimate behaviors. In the context of churn prediction, outliers can skew the results and affect the performance of the model. Therefore, detecting and handling outliers is crucial. Statistical methods, such as Z-score or Interquartile Range (IQR), are commonly used for outlier detection. A Z-score measures how far away a data point is from the mean in terms of standard deviations, and points with Z-scores greater than a certain threshold (e.g., 3) can be considered outliers. The IQR method involves calculating the range between the first quartile (Q1) and third quartile (Q3) and identifying data points that fall outside the range defined by $Q1 - 1.5 \times IQR$ and $Q3 + 1.5 \times IQR$. Once identified, outliers can be removed or capped, depending on the context and the business goals.

Feature Selection and Engineering: Feature selection and engineering are essential steps in ensuring the model's effectiveness. Feature selection helps identify which attributes in the dataset are most strongly correlated with churn, while feature

engineering focuses on creating new features from the existing data. Recursive Feature Elimination (RFE) is a commonly used technique for feature selection, where the model recursively eliminates less important features and retains the most relevant ones. Tree-based methods like Random Forest or Gradient Boosting can also be used to rank features by their importance, helping to highlight which features contribute most to the model's predictions.

Feature engineering involves creating new features that can better capture the relationship between customer behaviors and churn. For example, features like customer tenure (the number of months a customer has been subscribed) or frequency of service usage might provide deeper insights into whether a customer is likely to churn. Other features, such as spending patterns over time, can also be derived from the transactional data, highlighting whether customers are decreasing their spending–a common indicator of dissatisfaction.

# CHAPTER 4

# RESULTS & DISCUSSION

Model Performance Evaluation The performance of the various machine learning models (Logistic Regression, Random Forest, Gradient Boosting, and Support Vector Machines) will be evaluated based on key metrics such as accuracy, recall, precision, F1-score, and ROC-AUC. The focus will be on recall, as it is crucial for minimizing false negatives (i.e., failing to identify customers at risk of churning). A higher recall value indicates that the model is better at identifying customers who are likely to churn, which is essential for implementing timely retention strategies.

Accuracy vs. Recall Trade-Off While accuracy is important, it may not always provide a complete picture when dealing with imbalanced datasets like customer churn prediction. We expect models with higher recall to have slightly lower accuracy, as they may classify some non-churned customers as churned in order to avoid missing any high-risk customers. This trade-off will be carefully examined, and the best model will be chosen based on its ability to balance accuracy with recall.

Class Imbalance Handling The issue of class imbalance, where the majority class (non-churned customers) significantly outnumbers the minority class (churned customers), will be addressed using techniques such as SMOTE (Synthetic Minority Over-sampling Technique). These techniques will help the models better identify the minority class, improving the overall prediction of churned customers. The effectiveness of these techniques will be analyzed by comparing model performance before and after applying them.

Feature Importance The importance of various features in predicting customer churn will be discussed, highlighting which customer attributes (e.g., demographic factors, transaction history, service usage, and customer service interactions) have the most significant impact on churn prediction. This insight will help businesses focus on the key factors driving churn and refine their retention strategies accordingly.

Business Implications The results will also be discussed in terms of their practical implications for businesses. The developed predictive model can be integrated into

Customer Relationship Management (CRM) systems to flag at-risk customers in real-time. Businesses can use these predictions to implement proactive retention strategies such as personalized offers, loyalty programs, and targeted customer support, ultimately reducing churn rates and improving customer retention. 6. Limitations and Future Improvements Finally, we will discuss the limitations of the current models, such as potential overfitting, the impact of missing or noisy data, and the generalizability of the models across different industries or datasets.

Future improvements could include exploring deep learning models, incorporating additional external data sources (e.g., social media interactions or sentiment analysis), and enhancing the model's ability to update in real-time based on new customer data. In conclusion, the results will provide valuable insights into the effectiveness of various machine learning models for churn prediction, helping businesses better understand their customers' behaviors and implement targeted strategies to reduce churn and enhance customer loyalty. give me breefly in para formate

## 4.1 Evaluation Metrics and Their Importance in Churn Prediction

In churn prediction, selecting the right evaluation metrics is critical to assess the performance of machine learning models and ensure that they align with business goals, particularly the goal of minimizing customer churn. A major challenge in churn prediction is dealing with imbalanced datasets, where the number of non-churned customers (the majority class) is significantly greater than the number of churned customers (the minority class). In such cases, traditional metrics like accuracy are insufficient, as they can be misleading and fail to capture the model's true performance. To overcome this, it is essential to focus on more nuanced evaluation metrics, including Precision, Recall, F1-Score, and ROC-AUC, which provide deeper insights into how well a model identifies churned customers while minimizing false positives and false negatives.

Precision is the proportion of true positive predictions (customers predicted to churn who actually churned) out of all the positive predictions made by the model. It answers the question: "Of all the customers predicted to churn, how many actually did?" A high precision means that the business is targeting the right customers with retention efforts, reducing the risk of wasting resources on customers who aren't likely to leave. However, focusing too much on precision could lead to the model being overly conservative, missing many at-risk customers (resulting in lower recall).

On the other hand, Recall, also known as sensitivity, is the proportion of true positive predictions out of all actual positive cases (the total number of customers who actually churned). It answers the question: "Of all the customers who actually churned, how many did the model correctly identify?" High recall means that the model is effectively identifying customers at risk of churning, allowing businesses to take proactive actions before those customers leave. However, if the model focuses too much on recall, it may classify many non-churned customers as churned (leading to more false positives), causing unnecessary retention efforts and potentially alienating customers who are not at risk of leaving.

The F1-Score is the harmonic mean of precision and recall, providing a balanced evaluation metric that accounts for both false positives and false negatives. The F1-Score is particularly valuable when dealing with imbalanced datasets, where there may be a conflict between precision and recall. It combines the strengths of both metrics into one number, making it easier to assess overall model performance. A high F1-Score indicates that the model is effectively balancing precision and recall, minimizing both false positives and false negatives, which is crucial when businesses want to avoid wasting resources while still capturing as many at-risk customers as possible.

Finally, ROC-AUC (Receiver Operating Characteristic - Area Under Curve) is another critical metric that evaluates the model's ability to distinguish between the churned and non-churned classes across all possible classification thresholds. The ROC curve plots the True Positive Rate (Recall) against the False Positive Rate

(FPR), while the AUC summarizes the overall performance. A higher AUC indicates that the model can effectively discriminate between churned and non-churned customers. ROC-AUC is especially useful for comparing different models, as it reflects the model's ability to perform across a wide range of decision thresholds, providing a comprehensive view of model performance.

Understanding the trade-off between Precision and Recall is essential in churn prediction, particularly with imbalanced datasets. Optimizing one of these metrics often comes at the cost of the other. For instance, a model that maximizes precision will predict fewer churned customers, potentially missing some who are at risk of leaving (resulting in lower recall). Conversely, a model that maximizes recall will predict more customers as likely to churn, increasing the number of false positives and possibly wasting resources on customers who are not at risk of leaving.

To optimize churn prediction, businesses need to find the right balance between these metrics based on their retention goals. If the primary goal is to minimize the risk of losing at-risk customers, the business might prioritize recall, even at the cost of some false positives. If minimizing retention costs and resource allocation is more important, then precision might be prioritized. Ultimately, the choice between precision and recall depends on the specific business objectives and the trade-offs that are most important for the organization. By carefully evaluating the model's performance using these metrics, businesses can deploy a churn prediction model that effectively aligns with their retention strategies, ensuring that they target the right customers without overburdening their resources.

## 4.2 Error Analysis: False Positives vs. False Negatives

In churn prediction, one of the key challenges lies in understanding and managing the errors that occur during the model's predictions–specifically, false positives and false negatives. These errors have distinct implications for businesses and can significantly affect the effectiveness of churn prediction models.

False Positives (Type I Error): A false positive occurs when the model predicts that a customer will churn, but the customer actually does not churn. In other words, the model incorrectly classifies a non-churned customer as at-risk. While this may seem

like a minor issue at first, the consequences for businesses can be significant. When a business acts on a false positive, it may waste valuable resources—such as time, money, and effort—on retention strategies for customers who were never at risk of leaving. For instance, a customer may receive unnecessary retention offers, discounts, or personalized attention, which could hurt the company's bottom line if these efforts are not justified. In the long run, a high rate of false positives can lead to customer dissatisfaction, as customers may feel inundated by offers they don't need, which could harm the brand's reputation.

False Negatives (Type II Error): On the other hand, a false negative occurs when the model fails to identify a customer who is actually at risk of churning. This means that the model incorrectly predicts that a churned customer will remain with the company. False negatives are particularly detrimental because they represent missed opportunities to intervene with high-risk customers. When a customer who is likely to churn is not flagged, the company fails to implement any proactive retention strategies, and the customer might leave the service or product without the company ever realizing the risk. This can result in lost revenue, lower customer lifetime value, and increased customer attrition. In some industries, this could even translate to a significant drop in market share if churn goes unchecked over time.

Business Implications: From a business perspective, the impact of false positives and false negatives varies depending on the company's retention strategy. If a company is primarily focused on retaining as many customers as possible, it might tolerate a higher number of false positives in exchange for catching more of the highrisk churners (accepting some inefficiency in retention efforts). On the other hand, if the company is more concerned with minimizing costs and resource allocation, it may prioritize reducing false positives, accepting a higher number of false negatives in the hope of targeting only the most critical churners. The trade-off between false positives and false negatives is often a delicate balancing act that needs to be carefully managed, especially when working with imbalanced datasets, where the majority of customers do not churn.

Ultimately, the analysis of false positives and false negatives is essential in refining churn prediction models and aligning them with the company's strategic goals. A well-balanced model will minimize both types of errors, ensuring that retention efforts are accurately targeted and resources are used efficiently. Understanding these errors and their impact on the business will help companies make informed decisions about model deployment and retention strategies.

**4.3 Future Directions for Churn Prediction Models**

As churn prediction continues to evolve, several promising advancements are emerging in model development and application. A particularly exciting direction is the integration of deep learning techniques, such as neural networks, which excel at capturing complex, non-linear relationships in customer data. Traditional models may struggle to detect these subtle patterns, especially in large-scale datasets. Deep learning, on the other hand, can process vast amounts of data and reveal intricate customer behaviors that were previously hard to identify, thereby offering more accurate churn predictions, particularly in industries with complex customer interactions.

Another area for significant improvement is the integration of real-time data. Many current churn prediction models rely heavily on historical data, which, while valuable, can overlook sudden shifts in customer behavior or changes in market conditions. For instance, a customer might suddenly express dissatisfaction through social media or customer service interactions, which would not be captured by models based solely on past transactions. By incorporating real-time data–such as changes in customer sentiment, purchasing behavior, or social media activity– companies can build more responsive models that flag at-risk customers in near real-time, enabling businesses to act quickly and implement targeted retention strategies.

The use of external data sources is another promising avenue for improving churn prediction models. Beyond transactional and behavioral data, external factors such as market trends, competitor activity, or even social media sentiment can provide valuable insights into customer churn behavior. For example, if a competitor launches a more appealing offer or a customer faces an unforeseen personal crisis, these factors can influence the likelihood of churn. By incorporating such external data into churn prediction models, companies can develop a more holistic understanding of the factors affecting customer retention, improving the accuracy of their predictions.

Additionally, explainable AI (XAI) is gaining traction as a critical consideration in churn prediction. As machine learning models, particularly deep learning models,

become more sophisticated, they often operate as "black boxes," making it difficult for businesses to understand why certain predictions are made. For churn prediction, it's essential not just to know which customers are likely to churn, but also why they are at risk. With explainable AI, companies can gain insights into the key factors driving churn, which can guide retention strategies by targeting the root causes of dissatisfaction rather than merely addressing the symptoms.

Finally, adaptive models will play a crucial role in the future of churn prediction. As customer behaviors and market conditions evolve, static models can quickly become outdated. Future models should be designed to learn continuously, adapting to new trends and shifts in customer behavior. By incorporating mechanisms for continuous learning, where models are regularly updated and fine- tuned based on new data and real-time feedback, businesses can ensure that their churn prediction models remain relevant and effective over time. This dynamic approach allows companies to stay ahead of churn trends and take proactive measures to retain customers.

In conclusion, the future of churn prediction models lies in leveraging cutting-edge technologies like deep learning, real-time data integration, external data sources, explainable AI, and adaptive systems. By embracing these innovations, businesses can enhance their ability to predict churn more accurately and take timely, data-driven actions to retain high-value customers and foster long-term loyalty.

# CHAPTER 5

# CONCLUSION

In this project, we explored the implementation of various machine learning techniques for the prediction of customer churn, a fundamental task for businesses across industries aiming to enhance customer retention and mitigate the risk of losing valuable clientele. Customer churn, or attrition, represents a significant challenge for businesses that rely on long-term customer relationships, particularly in industries like telecommunications, e-commerce, and subscription-based services. Understanding why customers leave, as well as predicting which ones are likely to churn, is crucial for the development of effective retention strategies and ensuring long-term sustainability.

Churn prediction is a powerful tool for businesses that seek to reduce the financial impacts of losing customers. The ability to forecast customer attrition with high accuracy enables businesses to proactively engage customers who are at risk of leaving and take appropriate actions to retain them. In the modern data-driven economy, customer retention is often more cost-effective than acquiring new customers. Thus, reducing churn becomes a priority for increasing profitability and enhancing customer lifetime value (CLV). The implementation of machine learning models for churn prediction provides a robust framework for achieving these goals. Key Insights from the Project

This project presented an in-depth analysis of the challenges and solutions involved in churn prediction. One of the key aspects of this project was the emphasis on data preprocessing, which is crucial for ensuring the quality and usability of the data used to train machine learning models. Real-world datasets, especially in churn prediction, are often messy, imbalanced, and incomplete. Handling missing data, encoding categorical variables, normalizing numerical features, and detecting outliers are some of the critical steps taken to prepare the dataset for model training. Feature engineering also played a central role in improving model performance. By identifying important customer attributes–such as demographics, transaction history, and service usage patterns–we were able to extract valuable insights that

significantly contributed to the prediction accuracy. The integration of feature selection methods, such as Recursive Feature Elimination (RFE) and tree-based approaches, allowed us to focus on the most important features, which helped streamline the models and avoid overfitting.

Another major challenge in churn prediction is the issue of class imbalance, where the number of customers who do not churn (the majority class) greatly outweighs the number of churned customers (the minority class). This imbalance can cause traditional machine learning algorithms to bias their predictions toward the majority class, leading to poor identification of customers who are at risk of leaving. In our project, we addressed this challenge by employing techniques like Synthetic Minority Over-sampling Technique (SMOTE), which helped balance the dataset by creating synthetic instances of churned customers, thus enabling the models to perform better in identifying this critical minority class.

The project evaluated multiple machine learning models to predict customer churn, including Logistic Regression, Random Forest, Gradient Boosting, and Support Vector Machines (SVM). Each model was trained using the prepared data and evaluated based on several key metrics, including accuracy, recall, precision, F1-score, and ROC-AUC.

Logistic Regression, being a simple linear model, served as a baseline. While it performed adequately in terms of interpretability and ease of deployment, it struggled to capture the non-linear relationships inherent in the data, especially when the dataset was imbalanced.Random Forest and Gradient Boosting, which are ensemble-based models, performed significantly better in handling complex relationships within the data. These models are particularly adept at capturing interactions between features and providing high predictive accuracy, especially when tuned with techniques like hyperparameter optimization and feature importance ranking.Support Vector Machines (SVM), although powerful in high-dimensional spaces, proved to be slower to train on larger datasets and did not outperform Random Forest or Gradient Boosting in this specific project. However, SVMs still hold promise in churn prediction tasks when fine-tuned appropriately, especially for smaller, more structured datasets.

One of the primary insights from this model comparison was the trade-off between accuracy and recall. While accuracy is a common performance metric, it does not always serve as the best indicator in imbalanced classification tasks like churn prediction. A high accuracy might be achieved by simply predicting the majority class (non-churned customers) most of the time, but this approach would fail to identify the minority class (churned customers), which is the focus of the prediction task. Therefore, recall emerged as a more critical metric for this project. Recall measures how well the model identifies actual churners and ensures that the business can take proactive measures to retain those customers.

The project's findings also underscore the business implications of churn prediction models. With accurate churn predictions, companies can integrate these insights into their Customer Relationship Management (CRM) systems to flag at-risk customers in real-time. By identifying high-risk customers early, businesses can design targeted retention strategies, such as personalized offers, loyalty programs, and customized customer support. These strategies allow companies to maximize the impact of their retention efforts, focusing on customers who are most likely to leave.

Moreover, churn prediction can help businesses understand the underlying reasons for customer attrition, enabling more effective decision-making and resource allocation. For instance, churn predictions may reveal that customers are leaving due to poor service quality, pricing issues, or unmet expectations. By addressing these root causes, businesses can improve their offerings and reduce churn over time.

However, this project also identified several limitations in the models developed. For example, overfitting emerged as a concern in certain models, especially when complex algorithms like Random Forest and Gradient Boosting were trained on small datasets. Overfitting occurs when a model learns the noise or random fluctuations in the data, which leads to poor generalization to new, unseen data. To combat this, cross-validation and regularization techniques were employed to ensure that the models would generalize well to real-world data.

Another limitation was the reliance on historical data, which may not account for sudden shifts in customer behavior due to external factors such as economic changes or competitor activities. The integration of real-time data into churn prediction models can significantly improve their performance by allowing businesses to react to changes in customer behavior as they occur.

## 5.1 Future Directions and Potential Improvements

The future of churn prediction lies in advancing the methodologies used and addressing the limitations observed during the project. Key areas for improvement include the incorporation of deep learning techniques. While traditional machine learning models have proven effective in churn prediction, deep learning models, particularly neural networks, have the potential to capture complex and non-linear relationships in large datasets that are beyond the scope of conventional algorithms. These models could significantly improve the accuracy of churn predictions, especially in industries with large volumes of customer interactions.

Furthermore, real-time data integration will be critical in making churn predictions more dynamic and adaptive to changes in customer behavior. In addition to transactional and service-related data, real-time sentiment analysis from social media platforms, customer support interactions, and recent purchase behavior could offer a more comprehensive view of customer intent, allowing businesses to take immediate action when churn risks are detected.

The use of external data sources, such as competitor activities, macroeconomic factors, and social media sentiment, can also provide valuable context for understanding churn behavior. For instance, if a competitor launches a compelling offer or if a customer faces a personal financial crisis, these external factors could influence the likelihood of churn. Incorporating this data into churn prediction models could lead to more accurate and holistic predictions.

Lastly, explainable AI (XAI) will play an increasingly important role in churn prediction. While deep learning models are powerful, their "black-box" nature can make them difficult to interpret. In the context of churn prediction, businesses need to understand not only which customers are likely to churn but also why they are at risk. By using explainable models or techniques like feature importance analysis, businesses can gain valuable insights into the factors driving churn, leading to more targeted and effective retention strategies.

In conclusion, this project has provided valuable insights into the application of machine learning models for customer churn prediction. Through a combination of robust preprocessing, feature engineering, and class imbalance handling techniques, we were able to build models that could effectively identify at-risk customers. The evaluation of multiple machine learning models revealed that Random Forest and Gradient Boosting performed best in this task, with recall emerging as the most crucial metric for ensuring that businesses can take proactive measures to retain their customers.

The business implications of these models are clear: accurate churn prediction allows businesses to design targeted retention strategies that can significantly reduce customer attrition and improve long-term profitability. However, challenges such as overfitting, real-time data integration, and the need for explainable models remain. By advancing the use of deep learning, real-time data, and external data sources, businesses can continue to refine their churn prediction models and stay ahead of customer attrition trends.

As customer behavior continues to evolve, churn prediction models must also evolve to remain effective. The integration of continuous learning, real-time insights, and explainable AI will be key to developing the next generation of churn prediction models. Ultimately, by leveraging these advancements, businesses can not only predict churn more accurately but also take proactive, data-driven actions that foster long-term customer loyalty and drive sustained business growth.

# REFERENCES

**[1]** Bhat, D., & Jain, S. (2015). A machine learning approach to predict customer churn in retail banking. Journal of Retail Banking and Financial Services, 24(3), 243-252.

**[2]** Burez, J., & Van den Poel, D. (2009). Handling class imbalance in customer churn prediction. Journal of the Operational Research Society, 60(4), 594-606.

**[3]** Churn, M., & Holland, P. (2012). Modeling customer churn using decision trees. Journal of Machine Learning Applications, 3(2), 112-124.

**[4]** Fader, P. S., & Hardie, B. G. S. (2013). Forecasting customer retention in subscription services using a probabilistic model. Journal of Marketing Research, 50(2), 238-252.

**[5]** Ganti, A., & Srinivasan, V. (2017). Customer churn prediction in subscription-based services using ensemble methods. Journal of Business Analytics, 4(2), 98-109.

**[6]** Kumar, S., & Shah, M. (2018). Enhancing churn prediction models using feature selection techniques in telecommunications. International Journal of Computer Applications, 179(4), 15-20.

**[7]** Lemon, K. N., & Verhoef, P. C. (2016). Understanding customer experience and loyalty in services: A multidisciplinary approach. Journal of the Academy of Marketing Science, 44(6), 797-813.

**[8]** Li, L., & Zhang, X. (2016). Predicting customer churn in e-commerce with machine learning algorithms. Journal of Electronic Commerce Research, 17(4), 232-245.

**[9]** Ngai, E. W. T., & Xia, L. (2017). Applying machine learning algorithms to customer churn prediction in telecommunications. Expert Systems with Applications, 58, 135146.

**[10]** Patel, S., & Gupta, R. (2020). Comparative study of classification models for churn prediction in mobile network industry. Journal of Consumer Research, 49(2), 67-77.

**[11]** Sun, Y., & Yang, Q. (2014). A study on class imbalance problem in churn prediction. Data Mining and Knowledge Discovery, 28(1), 22-39.

**[12]** Ting, K. M., & Zhang, L. (2006). A review on classification techniques for churn prediction. International Journal of Computer Science and Network Security, 6(7), 9-15.

**[13]** Xia, L., & Ngai, E. W. T. (2018). Customer churn prediction in subscription-based services using a hybrid model. Journal of Business Research, 89, 232-241.

**[14]** Zhang, J., & Zhang, C. (2015). Customer churn prediction with the support vector machine. Computational Intelligence and Neuroscience, 1-9.

**[15]** Zhao, Y., & Xia, Y. (2019). Enhancing churn prediction with deep learning and real-time feedback. International Journal of Data Science and Analytics, 8(4), 227-238.