



Predicting Customer Churn in the Telecommunications Industry using Machine Learning Techniques

By
Adeline Makokha
191199

**A Research Proposal Submitted in Partial Fulfillment of the Requirements for
Completion of Master of Science in Data Science and Analytics (MSc. DSA)**

iLab Africa
Strathmore Institute of Mathematical Sciences (SIMS)
Strathmore University
Nairobi, Kenya
November, 2025

DECLARATION AND APPROVAL

DECLARATION

I declare that this work has not been previously submitted and approved for the award of a degree by this or any other University. To the best of my knowledge and belief, the proposal contains no material previously published or written by another person except where due reference is made in the proposal itself.

© No part of this thesis may be reproduced without the permission of the author and Strathmore University

Student's Name: **Adeline Makokha**

Sign:



27 September 2025

APPROVAL

The proposal was reviewed and approved for defense by

Dr. Henry Muchiri

School of Computing and Engineering Sciences

Strathmore University

Sign:



13 October 2025

ABSTRACT

In the recent telecommunications industry, predicting customer churn is a recurring issue that has a big impact on long-term viability, profitability, and competitiveness. The context and research challenges are first described here, which also highlights the sector's growing rivalry and the shortcomings of conventional churn control techniques. This research emphasizes the goals of finding important churn drivers, creating a prediction model based on machine learning models, via assessing various algorithms, and creating a deployable dashboard to help in decision-making. Recent research on churn prediction is combined with theoretical framework like Customer Relationship Management (CRM) in a review of relevant literature. Previous research shows the promise of machine learning models, but it also highlights insufficient information in terms of interpretability, profit-sensitive evaluation, and a narrow concentration on telecommunication sectors; hence, the huge gap between technological customer churn prediction and business insights to inform decision-making. The methodology suggests using secondary public data from a telecommunication provider and adheres to the CRISP-DM framework. Cleaning, feature engineering, addressing class imbalance, and encoding are examples of data preprocessing procedures to be used. In order to improve performance and stakeholder trust, a number of models including traditional models and modern models will be used. These includes algorithms like logistic regression, decision trees, naive Bayes (NB), support vector machine (SVM), random forest, k-nearest neighbors (KNN), and gradient boosting (GB). Interpretability approaches like SHapley Additive exPlanations (SHAP) and Local Interpretable Model-agnostic Explanations (LIME) will also be used for analysis to get results. The creation of a reliable churn prediction system with excellent predictive accuracy, interpretability, and scalability is one of the anticipated results. From an operational point of view, the system will facilitate constant monitoring, efficient resource allocation, and proactive customer retention tactics and strategies. In addition to providing methodological insights that can be applied to other telecommunications, this research offers a deployable and replicable methodology that solves the urgent commercial demand for customer churn reduction.

Keywords: Random Forest, Customer Churn, Telecommunication, Customer Retention

ACRONYMS

1. ANNs – Artificial Neural Networks
2. API - Application Programming Interface
3. AUC - Area Under the Curve
4. CNNs – Convolutional Neural Networks
5. CRISP-DM - Cross-Industry Standard Process for Data Mining
6. CRM- Customer Relationship Management
7. DL – Deep Learning
8. DT – Decision Tree
9. ERT – Extra Random Trees
10. FT – Functional Trees
11. GB – Gradient Boosting
12. k-NN – k-Nearest Neighbors
13. LightGBM - Light Gradient Boosting Machine
14. Lime - Local Interpretable Model-agnostic Explanations
15. LMT – Logistic Model Tree
16. LR – Logistic Regression
17. ML – Machine Learning
18. NB - Naive Bayes
19. RF – Random Forest
20. ROC - Receiver Operating Characteristic
21. ROI - Return on Investment
22. SHAP - SHapley Additive exPlanations
23. SMOTE - Synthetic Minority Over-sampling Technique
24. SVM - Support Vector Machine

Contents

1	CHAPTER: INTRODUCTION	2
1.1	Background	2
1.2	Research Problem	4
1.3	Research Objectives	4
1.4	Research Questions	5
1.5	Significance and Justification	5
2	CHAPTER: LITERATURE REVIEW	7
2.1	Theoretical Framework	7
2.1.1	Customer Relationship Management	7
2.1.2	Predictive modeling using statistical and segmentation	7
2.1.3	Forest Models and Ensemble Strategies	8
2.1.4	Ensemble Learning and Explainable AI Approaches for Churn Prediction	8
2.2	Conceptual Framework	9
2.3	Empirical Studies	9
2.3.1	Handling Class Imbalance	9
2.3.2	Comparative Performance of Traditional, Ensemble, and Deep Learning Models	10
2.3.3	Knowledge-Based and Text-Driven Feature Engineering	10
2.4	Customer Churn in the Telecommunications Industry	11
2.5	Industry Relevance and Contribution	12
2.6	Gaps in Existing Literature	12
3	CHAPTER: METHODOLOGY	13
3.0.1	Proposed Workflow	13
3.1	Business Understanding	14
3.2	Data Understanding	14
3.3	Data Collection	15
3.4	Exploratory Data Analysis	15
3.5	Data Cleaning	15
3.5.1	Treating Missing Data	15
3.5.2	Treating Outliers	16
3.5.3	Data Type Conversion	16

3.5.4	Data Transformation	16
3.6	Model Selection	17
3.6.1	Logistic Regression	17
3.6.2	K-Nearest Neighbors	17
3.6.3	Decision Tree	18
3.6.4	Random Forest	18
3.6.5	Support Vector Machine	18
3.6.6	Naive Bayes	18
3.6.7	Gradient Boosting	19
3.6.8	Justification for the Selected Classifiers	19
3.6.9	Handling Class Imbalance with SMOTE	21
3.7	Model Evaluation	21
3.8	Model Deployment	22
4	CHAPTER: EXPECTED OUTCOMES	23
	References	24
5	APPENDICES	26
5.1	Project Appendix	26
5.1.1	Project Timeline: 4 Months	26
5.2	Project Budget	26
5.2.1	Estimated Project Budget (KSh)	26

1. CHAPTER: INTRODUCTION

1.1 Background

Within the worldwide telecommunications business, customer churn is defined as customers leaving their subscription or services in favor of competing offers, which represents a significant and ongoing concern in the industry. Customer churn has significantly impacted the telecommunication business's profitability, operational sustainability, and long-term market competitiveness. Annual telecommunication churn rates are above 30% Khan et al. (2024) in the global telecommunication ecosystem. This has raised the need for proactive prediction and retention measures and causing significant annual revenue losses for service providers.

Mobile communication continues to dominate global connectivity, with mobile devices, both smartphones and feature phones, accounting for the majority of voice and data traffic worldwide. Despite this saturation, the mobile segment remains one of the rapidly growing areas within the broader telecommunications company. As competition intensifies, the strategic focus has shifted from acquiring new subscribers to retaining existing ones, a pattern seen across many mature markets. In this context, churn represents the proportion of customers who discontinue a service in favor of alternative providers as referenced by Khan et al. (2024). For telecommunication operators, tracking churn has become a core performance indicator, similar to how long-term service industries monitor client loyalty. The study further suggests that annual customer losses in the post-COVID period range between 30% and 35% due to hardship circumstances. These rates may escalate further as new entrants and larger competitors reshape the market, where some telecommunications companies which decided to increase their prices while others have reduced their prices, leading to a flexible shift of customers preferring other networks. At the core of the issue lies the financial burden of customer acquisition, as gaining a new subscriber typically requires significantly more investment than keeping an existing one, making effective churn management an essential business priority for telecommunication providers.

From a financial standpoint, keeping existing subscribers is far less costly than attracting new ones. Research indicates that the expense of acquiring new users may range between five to seven times higher than that of maintaining loyal customers Usman-Hamza et al. (2022). An increase in the number of churners can cause a reduction in revenue stability in the market, constraining the organizational competitiveness. The rapid growth and technological evolution of the telecommunications industry have broadened the number of companies operating in the sector, intensifying competitive

pressure. As a result, customer churn has become a persistent challenge, driven by saturated markets, aggressive competitive offerings, and frequent release of attractive service bundles. In such a flexible environment, telecommunication operators must continually find ways to protect and grow their revenue streams while also considering the dynamic factors in the industry. Strategies commonly recommended include attracting new users, increasing the value of existing customers, and prolonging subscriber loyalty. This makes churn prediction an essential component of modern telecommunication operations.

A study by Usman-Hamza et al. (2022) shows that Customer Churn Prediction (CCP) provides organizations with the ability to anticipate customer loss and design targeted retention initiatives that strengthen revenue performance and competitive position. Telecommunication firms now possess extensive datasets on subscriber behavior ranging from voice, data, and messaging records to demographic profiles and billing histories. These rich information sources are valuable for identifying patterns that signal a likelihood of churn. The challenge for operators is to use these insights proactively, detecting early signs of disengagement or inactivity before a customer decides to leave.

Globally, service-oriented industries particularly telecommunications continue to grapple with customer churn as digitalization accelerates and competition intensifies. International studies show that telecommunication markets generate significant revenue, especially in developing regions, yet face persistent challenges as operators introduce new technologies and service offerings to retain increasingly mobile customers. With monthly churn rates reported in several markets and evidence that retaining a subscriber is far more cost-effective than acquiring a new one, customer loyalty has become a strategic priority worldwide. Over the years, machine learning research has played a critical role in supporting these efforts by enabling firms to identify subscribers likely to leave and by revealing behavioral patterns embedded in customer databases. While global literature has extensively explored churn prediction through techniques such as feature extraction, ensemble learning, and gradient boosting, a growing number of local studies highlight the need to move beyond simple binary classification and incorporate contextual understanding of why customers churn. Recent contributions, such as Rahman et al. (2024), demonstrate that models like Logistic Regression, Random Forests, LightGBM, and ensemble workflows can handle diverse datasets and accurately flag high-risk customers before they exit. However, for regions, where telecommunication services underpin financial transactions, education, and everyday communication, combining predictive accuracy with interpretable insights is essential for designing effective retention strategies that reflect local user behavior and market realities

The study highlighted by Mwaura (2021) demonstrates that customer experience management plays a central role in reducing churn within Kenya's telecommunications industry. Using survey

data from 195 telecommunication staff, the research found that social factors, quality of service interfaces, and the atmosphere of retail outlets all contribute positively to customer retention, while pricing pressures increase the likelihood of churn. These findings emphasize that customer loyalty in Kenya is shaped not only by service performance but also by relational, experiential, and affordability dimensions. As the sector becomes more competitive and essential to daily digital activity, telecommunication operators must prioritize customer-centered strategies leveraging social engagement, enhancing user interactions, improving retail environments, and offering price-sensitive plans to strengthen retention and maintain market stability.

The rapid digital adoption in the market and socio-economic impact on mobile services, including calls and data usage, has really led to high competition in the telecommunications market. By combining demographic information, revenue information, and service usage, the customer can be used to intervene in the customer retention strategies, which will save on marketing expenditure and improve customer loyalty to one network.

1.2 Research Problem

Industries have really intensified the increase in customer expectations, and customers are preferring affordable and reliable services. Most of the customers have left the network and prefer other telecommunications service providers that offer good and affordable service qualities. Most of the previous churn rate predictions remain reactive, that is, relying on static indicators and manual analysis, which can be inefficient, prone to human bias, and also time-consuming to come up with business decisions. This existing prediction also has limitations when it comes to bridging the gap between the technological concept of customer churn and the business insights. These methods have struggled to predict customer churn and have resulted in revenue leakage due to the high customer churn rate in the industry. Mwaura (2021).

1.3 Research Objectives

The purpose of this study is to address the existing gaps in customer churn prediction within the telecommunications industry through the following objectives:

Main Objective

To design and evaluate a machine learning framework capable of accurately predicting customer churn in the telecommunications sector.

Specific Objectives

1. To analyze and document the methodological and conceptual challenges present in current churn prediction research within the telecommunications field.
2. To construct a predictive churn model based on machine learning algorithms using telecommunication customer data.
3. To assess and compare the predictive performance of selected machine learning models using appropriate evaluation metrics.
4. To implement the best-performing churn prediction model in an interactive and interpretable dashboard for practical use by decision-makers.

1.4 Research Questions

In alignment with the stated objectives, this study will be guided by the following research questions:

1. What methodological and data-related limitations characterize existing churn prediction models in the telecommunications industry?
2. Which customer attributes, such as demographic, contractual, financial, or behavioral factors, have the strongest influence on churn likelihood?
3. How can machine learning algorithms be developed, optimized, and validated to enhance the accuracy of churn prediction?
4. What differences exist in the predictive performance of various machine learning approaches (e.g., logistic regression, random forest)?
5. In what ways can a predictive churn model be integrated into an interactive dashboard to support strategic decision-making and customer retention initiatives?

1.5 Significance and Justification

Globally, effective management of churn is mandatory in order to ensure sustainability, profitability, and competitiveness in the telecommunications industry. The strategies for customer retention have really reduced the high cost that is put in place to acquire new customers in the network. This has also stabilized the telecommunications revenue stream and has safeguarded the market share as well as enhanced a long-term customer value relationship with the industry.

The implications of customer churn have extended beyond the corporate world and are now also impacting the emerging markets. This is evident in our day-to-day lives through e-learning systems,

e-health systems, and small businesses' operations, where telecommunication services have a major impact and are one of the drivers of the economy. There is slow progress towards digital inclusion and also economic development as a result of high churn rates.

This research is justified on three key fronts:

- **Economic efficiency:** By identifying at-risk customers before they leave, providers can target retention efforts more precisely, thereby reducing costs and maximizing return on marketing investment.
- **Operational effectiveness:** The integration of the predictive model into CRM systems ensures actionable insights are delivered to decision-makers, enabling timely and effective interventions.
- **Scalability and replicability:** While tailored to socio-economic and competitive environments, the proposed framework offers a methodological blueprint that can be adapted for other telecommunication markets and industries facing similar churn-related challenges.

This research seeks to advance both practical and theoretical understanding of customer churn within the telecommunications sector. It aims to establish a sustainable, data-informed framework that integrates machine learning techniques into day-to-day business operations to enhance customer retention.

2. CHAPTER: LITERATURE REVIEW

There has been a lot of shift over the past years in the digital evolution of technology; there are a lot of markets and startups that have opened up, and as a result, there is customer demand and more personalized quality services in the telecommunications industry. Most of these changes have really widened access to communication and led to the need for improvement. With this highly competitive ecosystem of customer churn, keeping the existing customers has been more critical than acquiring new customers into the network. One of the main challenges is the subscriber or customer making the decision to stop using the network, for one reason or another. Many telecommunications operators have lost a third of their customer base to other operators, making customer churn a critical global issue. These losses have significantly translated to revenue erosion in the industry; hence the need to develop and adopt a data-driven approach that can curb customer churn Nagarkar (2022).

2.1 Theoretical Framework

2.1.1 Customer Relationship Management

Shahabikargar et al. (2025) provides an important contribution to Customer Relationship Management (CRM) by demonstrating how deeper insights into customer cognition and emotional signals can significantly strengthen churn prevention strategies. Traditional CRM systems often rely on demographic attributes, usage behavior, or billing history; however, this study shows that much richer information lies within customer-company interactions, particularly unstructured text such as emails or service requests. Importantly, the study highlights that this approach not only enhances customer retention but can also extend to broader CRM applications such as personalized engagement, early detection of customer dissatisfaction, and proactive service intervention, positioning ChurnKB as a powerful tool for building stronger, more responsive customer relationships.

2.1.2 Predictive modeling using statistical and segmentation

Zhang et al. (2022) study provides an important theoretical contribution to understanding churn behavior by integrating customer segmentation with statistical prediction techniques. Using data from three major Chinese telecommunication operators, the research employed Fisher discriminant analysis and logistic regression to construct predictive models capable of identifying customers at risk of leaving. The findings indicated that logistic regression delivered superior performance, achieving a prediction accuracy of 93.94%, outperforming the discriminant approach. This demonstrates the ef-

fectiveness of regression-based models in capturing behavioral differences across customer segments and translating them into actionable churn insights. Zhang et al. (2022) work reinforces the theoretical premise that combining segmentation strategies with interpretable statistical models enhances the ability of telecommunication firms to anticipate churn and implement targeted retention interventions, ultimately supporting more profitable and customer-centric decision-making.

2.1.3 Forest Models and Ensemble Strategies

Usman-Hamza et al. (2022) offers a comprehensive theoretical perspective on how intelligent decision forest algorithms can improve customer churn prediction within the telecommunications domain, where customer acquisition costs are high and competitive pressure is intense. Their work highlights that conventional rule-based systems are limited in scalability and that standard machine learning models often underperform due to the significant class imbalance between churners and non-churners. To overcome these challenges, the researchers implemented and compared several decision forest variants such as logistic model trees, random forests, and functional trees while also developing advanced ensemble methods incorporating weighted soft voting and stacking to enhance predictive performance.

Using publicly available benchmark datasets, the study demonstrated that these decision forest models consistently outperformed baseline machine learning methods and showed strong resilience when dealing with imbalanced data. The findings highlight the theoretical and practical value of ensemble-based decision forests in producing more stable and accurate churn predictions, supporting their adoption for telecommunications CCP and broader machine learning applications.

2.1.4 Ensemble Learning and Explainable AI Approaches for Churn Prediction

Recent work examining churn dynamics in the telecommunications sector, an industry known to record annual churn levels exceeding 30%, has highlighted the value of ensemble learning in developing accurate and interpretable prediction models.

The research explored several ensemble and decision tree-based algorithms, such as boosted trees, random forests, and standard decision trees, to assess their effectiveness in predicting customer attrition using extensive telecommunications datasets. Among all the models examined, the Random Forest algorithm produced the most reliable outcomes, attaining an accuracy of 91.66%, a precision of 82.2%, and a recall of 81.8%. These results indicate its strong capability to detect customers who are likely to discontinue their subscriptions.

Moreover, the incorporation of explainable artificial intelligence (XAI) tools, specifically SHAP

and LIME, enhanced the interpretability of the model by revealing how individual features contributed to churn risk. Overall, the study highlights that integrating ensemble methods with explainable AI frameworks not only elevates predictive performance but also fosters model transparency, a crucial factor for telecom firms aiming to build data-driven and trustworthy retention systems.

2.2 Conceptual Framework

The conceptual models guiding this research are mainly integrated in four components that include:

- **Behavioral and Demographic Profiling:** This mainly entails customer characteristics like service usage, billing routine and revenue impact just to establish most of the key predictors in customer churn.
- **Predictive Modeling:** The study will implement and compare multiple machine learning techniques, namely Logistic Regression, Decision Trees, Random Forest, LightGBM, and Ensemble Learning, to assess their predictive reliability and consistency in identifying potential churners.
- **Retention Insights:** The predicted churn probabilities will be transformed into practical retention strategies, including customized promotions, focused customer engagement, and segmentation-based interventions designed to meet specific client needs.
- **Model Evaluation:** Model performance will be analyzed using key metrics such as precision, recall, F1-score, and ROC-AUC, providing a comprehensive view of how effectively each algorithm distinguishes between churn and non-churn cases.

2.3 Empirical Studies

2.3.1 Handling Class Imbalance

Research in telecommunications churn analysis emphasizes that revenue loss stemming from customer exit remains one of the sector's most pressing challenges. As retention strategies increasingly outweigh acquisition-focused approaches in their cost-effectiveness, predictive modeling has become central to identifying customers at risk. However, many traditional machine learning models struggle with the highly imbalanced structure of telecommunication datasets, where churners represent only a small fraction of the customer base. Addressing this limitation, Nagarkar (2022) evaluated a large, operator-specific dataset from Nepal comprising 52,332 subscribers and demonstrated that XGBoost could effectively capture churn behavior despite the skewed class distribution. The model achieved

strong performance, recording 97% accuracy and an 88% F1-score on the native dataset, which included 6,128 churners and 46,204 non-churners. When applied to a smaller publicly available dataset of 3,333 records for comparison, XGBoost also delivered improved results, achieving 96.25% accuracy and an 86.34% F1-score. These findings highlight the robustness of gradient-boosting algorithms in handling real-world telecommunication data and reinforce the importance of using representative datasets to validate churn prediction models.

2.3.2 Comparative Performance of Traditional, Ensemble, and Deep Learning Models

Saha et al. (2023) carried out an extensive comparative investigation into churn prediction by evaluating a broad spectrum of machine learning and deep learning algorithms. Their work incorporated both conventional classifiers, such as logistic regression, decision trees, and k-nearest neighbors, and ensemble approaches, including AdaBoost, Random Forest, Gradient Boosting, Extreme Randomized Trees, and various bagging and stacking combinations. In addition, Artificial Neural Networks (ANNs) and Convolutional Neural Networks (CNNs) were implemented to explore the performance differences between traditional algorithms and deep learning models. Using publicly available datasets from the Southeast Asian and U.S. telecommunications markets, the researchers reported that deep learning architectures outperformed other models. Specifically, CNNs reached accuracy levels of 99% and 98% on the two datasets, while ANNs achieved comparable results with 98% and 99% accuracy. These results highlight that deep learning models are more adept at identifying intricate behavioral trends in customer data, making them particularly valuable for telecommunication companies striving for precise churn prediction in competitive environments.

2.3.3 Knowledge-Based and Text-Driven Feature Engineering

Shahabikargar et al. (2025) Recent work offers an important advancement in churn prediction by shifting the focus from traditional demographic or usage-based variables to richer, knowledge-driven features extracted from customer interactions. The study emphasizes that understanding customer cognition, emotions, and behavioral signals is essential for anticipating churn more accurately than relying solely on structured CRM records. To address this gap, the researchers introduced the Customer Churn-related Knowledge Base (ChurnKB), a feature engineering framework that incorporates domain expertise, textual data mining techniques such as TF-IDF, cosine similarity, tokenization, and stemming, as well as generative AI models capable of interpreting unstructured text like emails. By integrating these knowledge-based feature into machine learning models, including Random Forests, Logistic Regression, Multilayer Perceptrons, and XGBoost, the study demonstrated substantial performance gains. Notably, the F1-score of XGBoost improved from 0.5752 to 0.7891, illustrating

the value of incorporating cognitive and behavioral indicators alongside conventional features. Beyond telecommunication churn management, the authors highlight that the same knowledge-based feature extraction approach can support broader applications such as personalized marketing, online harm detection, and mental health monitoring, signaling its potential as a versatile tool for business intelligence and digital safety

The last decade has seen remarkable progress in how customer churn is modeled, particularly with the rise of machine learning and deep learning approaches. Usman-Hamza et al. (2022) showed that decision-forest ensembles equipped with weighted voting can handle imbalanced telecommunication datasets more effectively than baseline classifiers, demonstrating great improvements in predictive stability. In another strand of research, they Saha et al. (2023) experimented with convolutional neural networks on mixed datasets from India and the United States, concluding that financial metrics such as return on investment should carry equal importance to raw predictive performance when evaluating model usefulness.

Geographical context continues to shape model behavior. In Nepal, Nagarkar (2022) used XG-Boost on more than fifty thousand mobile subscribers and observed that recharge regularity and prolonged inactivity were the strongest signals of customer retention, enabling more targeted interventions. In the Syrian context, the introduction of social network analysis features pushed AUC values from 0.84 to 0.933, illustrating the influence of relational information on churn outcomes. More recently, Poudel and Sharma (2024) proposed models that integrate explainability tools directly into the prediction workflow, noting that transparency is essential for organizational acceptance of churn analytics.

A broader review Shahabikargar et al. (2025) emphasized the rapid increase in external data integration and the need for interpretability alongside accuracy. Their synthesis underscores a growing shift toward models that not only perform well but can also justify their predictions to decision-makers.

2.4 Customer Churn in the Telecommunications Industry

Churn continues to pose a significant threat to revenue certainty and long-term profitability in the telecommunications sector. The industry statistics show that annual churn rates in some regions exceed 30%, even in highly developed markets Chang et al. (2024). The challenge is intensified in Sub-Saharan Africa, where mobile networks support essential services such as digital services like payments, remote education, remote jobs, and online sales, Mwaura (2021), where most of the communications rely on online platforms; hence, the need to buy resources like data bundles for

communication. Despite the growing body of churn research, very few studies directly address market conditions. This lack of contextual evidence highlights an important gap that limits the applicability of existing predictive approaches of machine learning models across the region.

2.5 Industry Relevance and Contribution

The review of prior research points to several consistent themes. Advanced machine learning techniques, especially ensemble and deep learning methods, tend to outperform traditional statistical models like survival analysis, which mainly focus on predicting the timing of churn rather than just whether churn happens. This analysis can be more complex to implement and also interpret compared to the model models like Random Forest or gradient boosting when it comes to churn prediction. At the same time, organizations increasingly require models that balance predictive power with financial clarity, enabling business teams and stakeholders to interpret outputs and integrate them into decision processes as per as revenue is concerned. Another recurring insight is the disproportionate focus on other markets, leaving emerging economies underrepresented in current literature.

The framework proposed in this study responds directly to these gaps by prioritizing scalable modeling approaches, interpretability, and practical alignment with the operational realities of telecommunication providers.

2.6 Gaps in Existing Literature

Limited Feature Diversity: Many datasets rely heavily on transactional CRM variables such as usage volume, payment history, and subscription tenure while overlooking high-value indicators like financial components, Mwaura (2021). This restricts the richness of models and can weaken financial prediction quality.

Narrow Evaluation Practices: There is a strong preference for accuracy-based metrics, yet these do not always reflect the financial consequences of churn. Return on investment-linked indicators, including lifetime value measures, are rarely incorporated despite their relevance for managerial decision-making Shahabikargar et al. (2025).

Interpretability Limitations: Although advanced algorithms perform well, their outputs are not always intuitive for business users. The limited transparency often reduces trust in business stakeholders and slows the adoption of churn systems as practical decision-support tools Chang et al. (2024).

3. CHAPTER: METHODOLOGY

This study adopts the Cross-Industry Standard Process for Data Mining (CRISP-DM) methodology to ensure a structured, repeatable, and scalable approach to customer churn prediction in the telecommunications sector, as illustrated in Figure 3.1. This reproducible approach shows that it can be plugged into any other telecommunications industry.

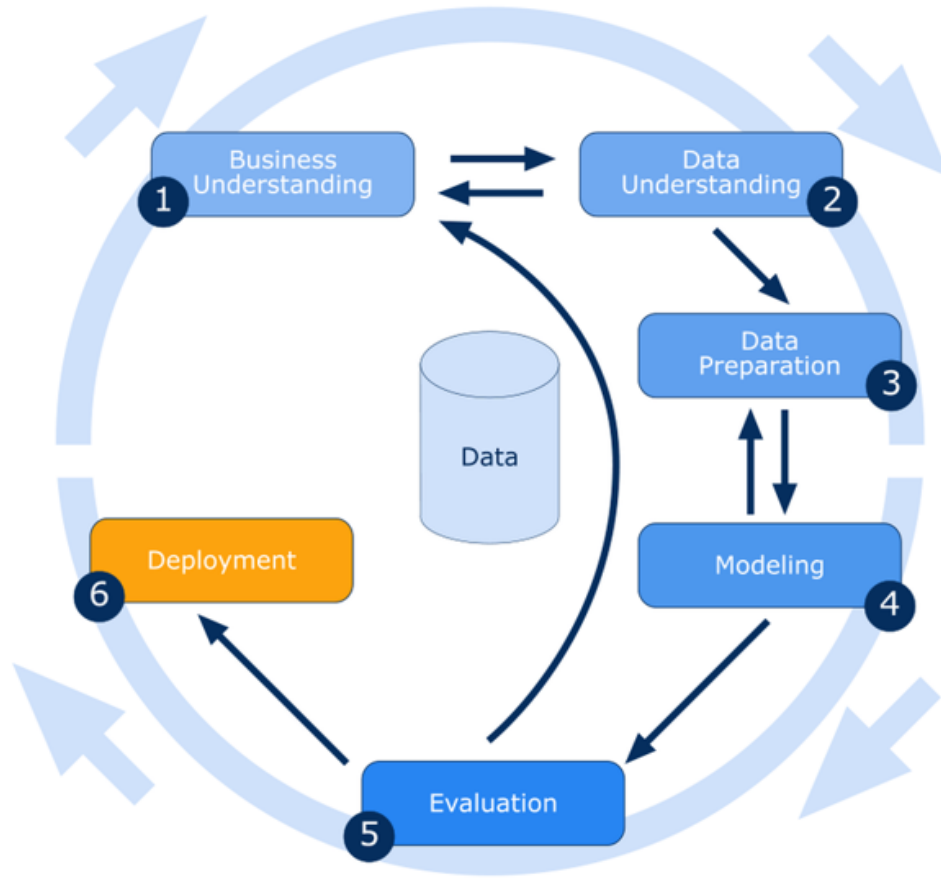


Figure 3.1: Cross-Industry Standard Process for Data Mining (CRISP-DM)

3.0.1 Proposed Workflow

The proposed workflow for this study is illustrated in Figure 3.2. It outlines the systematic process of building and evaluating machine learning models for customer churn prediction.

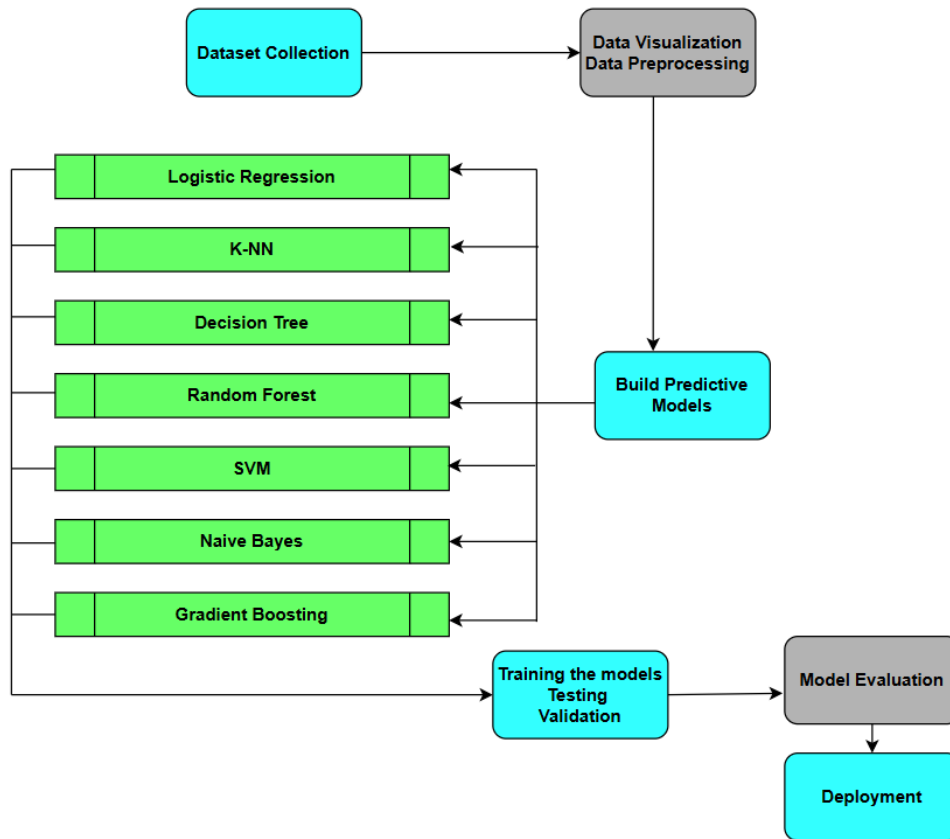


Figure 3.2: Proposed Workflow

3.1 Business Understanding

Telecommunications companies operate in a market where competition is intense, and customer expectations shift rapidly according to their preferences. As subscription growth slows and price competition increases, retaining existing users becomes a more important and cheaper method than acquiring new ones. For many operators, the financial burden of attracting new subscribers is far higher than maintaining current customers, making churn reduction central to business sustainability and operations. Despite continuous investment in marketing, loyalty programs, and service improvement initiatives, many operators still experience high churn levels that negatively affect revenue reliability and weaken long-term competitiveness. Understanding the drivers of churn and building a predictive model that can anticipate customer exit has therefore become a strategic necessity in this industry.

3.2 Data Understanding

This phase will mainly concentrate on identifying, collecting, and analyzing datasets essential for achieving project goals. Tasks include collecting initial data, describing data, exploring data, and

verifying data quality.

3.3 Data Collection

This study will use secondary data sourced from a Bulgarian telecommunication operator, which is publicly available on the Mendeley data website Tokmakov (2024). This dataset comprises customers' service usage, demographic information, and spending, among others, which will provide key indicators of customer churn, allowing us to anticipate the behaviors that contribute to customer retention and predict the main behavior that will help us retain customers in the network. The study will encourage iteration in model building and assessment until a good enough model is achieved.

3.4 Exploratory Data Analysis

Exploratory Data Analysis (EDA) will be applied to understand variable distributions, identify class imbalance cases, and diagnose missing or inconsistent values. Graphical tools, including box plots, histograms, and correlation matrices will be used to examine relationships between various attributes. These insights will inform preprocessing decisions, feature engineering strategies, and the selection of modeling techniques to ensure a reliable and interpretable predictive framework that will be used.

3.5 Data Cleaning

Data cleaning is a fundamental step that ensures the dataset is consistent, accurate, and suitable for modeling. Given that the information is extracted from telecommunication systems, it may contain missing entries, duplicated records, and inconsistent categorical labels. The cleaning process will follow a structured approach: identifying missing values, applying appropriate imputation strategies, identifying duplicates, standardizing categorical fields, and reviewing outliers in the dataset.

A structured data cleaning workflow will be followed to ensure quality and consistency. Missing values in key numerical fields will be assessed for extent and pattern. For variables with minimal missingness, mean or median imputation will be applied, while variables with substantial or non-random missingness will undergo domain-informed estimation or exclusion. Duplicate records in the dataset based on unique identifiers will be identified and removed to prevent double-counting.

3.5.1 Treating Missing Data

Missing values will be handled based on their type and extent. Numerical attributes will be imputed using mean or median statistics, while categorical fields will use mode-based imputation. Variables

with excessive or patterned missingness may undergo regression-based imputation or exclusion if they contribute minimal analytical value.

3.5.2 Treating Outliers

Outliers will be identified using statistical techniques such as z-scores and IQR-based thresholds. Depending on their origin, extreme values will be capped, transformed, or retained where they reflect genuine business behavior.

3.5.3 Data Type Conversion

All variables will be converted to appropriate data formats to ensure compatibility with analysis tools. Categorical variables will be encoded numerically, enabling their use in machine learning algorithms.

3.5.4 Data Transformation

Feature Engineering

Raw dataset will be refined to enhance the performance of machine learning models. This process will involve the creation of new features or modification of existing ones to extract meaningful information and patterns. Different techniques will contribute to the generation of features that better capture the underlying complexities of the data.

Feature Scaling

Numerical variables will undergo standardization or min-max scaling to ensure that models relying on distance metrics or gradient-based optimization behave consistently. The standardization methods, like z-score normalization and min-max scaling, will bring features within a comparable scale, preventing dominant features from influencing the model outcome.

Feature Selection

This step will involve correlation patterns, mutual information scores, and model-driven importance rankings, which will be used to remove redundant or weak predictors. Tree-based algorithms such as random forests and gradient boosting will guide feature prioritization.

Data Encoding

Categorical attributes will be encoded using either one-hot encoding, label encoding, or binary encoding techniques, depending on their hierarchy. This process will involve converting categorical data

into a numerical format, allowing algorithms to interpret and utilize this information effectively.

3.6 Model Selection

Customer churn prediction represents a supervised learning challenge where the goal is to classify customers into specific categories, typically churners or non-churners, based on a set of descriptive features. In this research, seven modern classification algorithms are examined: Logistic Regression (LR), K-Nearest Neighbors (KNN), Decision Tree (DT), Random Forest (RF), Support Vector Machine (SVM), Naive Bayes (NB), and Gradient Boosting (GB). The comparative analysis aims to determine the most accurate and interpretable approach suitable for telecommunication customer data. After model evaluation, the best-performing algorithm will be selected for final churn prediction.

3.6.1 Logistic Regression

Logistic regression estimates the likelihood that a customer will discontinue service (churn) based on their feature vector $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{ip}]$. The model computes this probability using the logistic function:

$$P(y_i = 1 | \mathbf{x}_i) = \frac{1}{1 + e^{-(\mathbf{w}^\top \mathbf{x}_i + b)}} \quad (3.1)$$

Here, \mathbf{w} represents the model coefficients and b is the bias term, both estimated through maximum likelihood optimization. Logistic regression is preferred for its interpretability, computational efficiency, and ability to quantify how individual customer attributes influence churn probability.

3.6.2 K-Nearest Neighbors

The KNN algorithm is a simple yet effective non-parametric method that classifies a new data instance based on the majority class among its k closest points in the feature space. The proximity between data points is typically measured using a distance metric such as Euclidean or Manhattan distance. KNN is advantageous for its intuitive design and adaptability to non-linear decision boundaries, though it can be computationally demanding for large-scale telecommunication datasets.

The predicted label is

$$\hat{y} = \text{mode}\{y_j \mid \mathbf{x}_j \in \mathcal{N}_k(\mathbf{x}_i)\}, \quad (3.2)$$

where $\mathcal{N}_k(\mathbf{x}_i)$ denotes the k closest neighbors according to Euclidean distance. KNN is robust to noisy data and effective when class boundaries are well separated.

3.6.3 Decision Tree

Decision trees split the dataset recursively using an impurity metric such as the Gini index:

$$\text{Gini}(t) = 1 - \sum_{c=1}^C p_c^2, \quad (3.3)$$

where p_c is the fraction of samples of class c at node t . DTs offer visual transparency and help managers trace the rules leading to churn predictions.

3.6.4 Random Forest

The Random Forest algorithm constructs an ensemble of B independent decision trees, each trained on a randomly drawn subset (bootstrap sample) of the original dataset. For any given observation, the collective prediction is determined through a majority voting process among all individual trees:

$$\hat{y} = \text{mode} h_b(\mathbf{x})_{b=1}^B \quad (3.4)$$

This ensemble approach minimizes overfitting, enhances generalization, and effectively manages high-dimensional and heterogeneous data. Moreover, Random Forest inherently provides measures of feature importance, offering valuable insight into which attributes most influence customer churn.

3.6.5 Support Vector Machine

Support vector machines classify data by identifying an optimal separating hyperplane that maximizes the distance (margin) between different classes. The optimization problem can be expressed as

$$\min_{\mathbf{w}, b} \frac{1}{2} |\mathbf{w}|^2 \quad \text{subject to} \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \quad \forall i \quad (3.5)$$

SVMs are capable of modeling complex, non-linear decision boundaries by applying kernel functions such as polynomial, radial basis function (RBF), or sigmoid kernels. This makes them suitable for telecommunication churn prediction, where behavioral relationships between variables are often non-linear.

3.6.6 Naive Bayes

Naive Bayes is a probabilistic model that leverages Bayes' theorem under the simplifying assumption that input features are conditionally independent given the class label. The class posterior probability

is computed as

$$P(y|\mathbf{x}) = \frac{P(y) \prod_{j=1}^p P(x_j|y)}{P(\mathbf{x})} \quad (3.6)$$

Although this independence assumption is rarely exact, the method performs remarkably well on large-scale, high-dimensional datasets. Its computational efficiency and simplicity make Naive Bayes particularly effective for rapid churn classification in telecommunication databases.

3.6.7 Gradient Boosting

Gradient boosting combines multiple weak learners, typically shallow decision trees, into a strong predictive model through a sequential learning process. At each stage m , the algorithm fits a new learner $h_m(\mathbf{x})$ to the pseudo-residuals of the previous model:

$$r_{im} = -\frac{\partial \mathcal{L}(y_i, F_{m-1}(\mathbf{x}_i))}{\partial F_{m-1}(\mathbf{x}_i)} \quad (3.7)$$

The ensemble model is then updated according to

$$F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \eta, h_m(\mathbf{x}) \quad (3.8)$$

where η represents the learning rate, controlling how much each new tree contributes to the model. Gradient boosting is widely recognized for its strong predictive accuracy, robustness to noise, and ability to uncover complex, non-linear dependencies prevalent in customer churn data.

3.6.8 Justification for the Selected Classifiers

The decision to employ five classification algorithms, Logistic Regression (LR), K-Nearest Neighbors (KNN), Decision Tree (DT), Random Forest (RF), Support Vector Machine (SVM), Naive Bayes (NB), and Gradient Boosting (GB) was made to ensure a balanced assessment of both interpretable and high-performing predictive models. Each algorithm contributes a unique analytical value and addresses different aspects of the customer churn phenomenon, which is inherently multidimensional in nature.

Logistic regression is adopted as the baseline model because of its simplicity, transparency, and effectiveness in binary classification problems. It estimates the likelihood of a customer discontinuing service based on explanatory variables. LR provides interpretable coefficient weights that help

managers understand the marginal influence of each variable on churn probability, making it highly suitable for initial benchmarking.

KNN offers an instance-based, non-parametric approach where each new observation is classified by the majority label among its nearest neighbors in feature space. This method captures local behavioral similarities among customers; for example, subscribers with comparable usage or spending patterns tend to exhibit similar churn tendencies. Although computationally intensive on large datasets, KNN's flexibility in modeling non-linear decision boundaries adds valuable comparative insight.

Decision trees are selected for their interpretability and intuitive structure. The model recursively partitions data using feature-based thresholds, creating a visual hierarchy of decision rules. In a churn context, a DT might reveal which customers who are most at risk. Its rule-based nature enables straightforward explanation to non-technical stakeholders and supports managerial decision-making.

The Random Forest algorithm, an ensemble of multiple decision trees, is incorporated to improve predictive accuracy and reduce overfitting. By aggregating the outputs of many trees trained on random subsets of data and features, RF captures complex, non-linear relationships while maintaining robustness to noise. Furthermore, its feature importance scores provide quantitative measures of the most influential churn drivers.

SVM is included for its strong generalization capacity, particularly in high-dimensional feature spaces. By constructing an optimal hyperplane that maximizes the margin between classes, SVM performs well even with overlapping class boundaries. Kernel functions further enhance its ability to model non-linear relationships in customer behavior, offering an advanced benchmark against tree-based models.

Naive Bayes contributes a probabilistic perspective grounded in Bayes' theorem. Despite its simplifying assumption of feature independence, it is computationally efficient and performs remarkably well with large-scale categorical data. Its ability to quickly estimate churn probabilities makes it ideal for initial screening of customers at risk, particularly when rapid, low-cost predictions are desirable.

Gradient boosting is chosen for its capacity to deliver state-of-the-art predictive accuracy by combining multiple weak learners in a sequential, error-correcting manner. Each new tree in the sequence focuses on the residual errors of the previous ensemble, allowing the model to learn complex patterns. Although parameter tuning is essential to prevent overfitting, GB's flexibility and performance make it indispensable in churn analytics. Moreover, its compatibility with interpretability frameworks such as SHAP enables a clear understanding of feature contributions.

These seven algorithms collectively span the major paradigms of supervised learning, linear,

probabilistic, instance-based, tree-based, and boosting models, allowing for a comprehensive comparison of predictive behavior. The inclusion of both interpretable (LR, DT) and high-performing (RF, GB, SVM) algorithms ensures a balanced evaluation of accuracy, computational efficiency, and managerial usability. Through this comparative analysis, the study aims to select the model that best captures customer behavior dynamics while maintaining practical relevance for operational decision-making in the telecommunications sector.

3.6.9 Handling Class Imbalance with SMOTE

In most telecommunications datasets, the number of customers who churn is typically much smaller than those who remain subscribed, resulting in a pronounced class imbalance. To mitigate this issue, the Synthetic Minority Over-sampling Technique (SMOTE) will be employed to enhance the representation of the minority (churn) class within the training data. SMOTE works by creating synthetic samples for the minority class rather than merely duplicating existing records. For each minority data point x_i , a new synthetic sample x_{new} is generated using the relationship:

$$x_{\text{new}} = x_i + \delta \times (x_{nn} - x_i), \quad \delta \sim U(0, 1) \quad (3.9)$$

Here, x_{nn} denotes one of the k nearest neighbors of x_i , and δ is a random value drawn uniformly between 0 and 1. By interpolating between existing minority samples and their neighbors, SMOTE produces realistic synthetic observations, leading to a more balanced training distribution. This balancing improves the model's sensitivity to churners and prevents it from being overly biased toward the dominant non-churn class.

3.7 Model Evaluation

The performance of the predictive models will be assessed through an 80/20 train–test split, with stratified sampling applied to maintain the natural class ratio between churners and non-churners. This ensures that both subsets reflect the true population distribution, improving the reliability of model validation.

Model effectiveness will be quantified using key performance indicators, namely precision, recall, F1-score, and the area under the Receiver Operating Characteristic Curve (AUC–ROC), expressed as follows:

$$\text{Precision} = \frac{TP}{TP + FP}, \quad \text{Recall} = \frac{TP}{TP + FN}, \quad \text{F1-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (3.10)$$

In this formulation, TP , FP , and FN represent true positives, false positives, and false negatives, respectively. Precision measures the proportion of correctly identified churners among all predicted churners, whereas recall evaluates the proportion of actual churners successfully captured by the model. The F1-score serves as a harmonic mean of precision and recall, balancing the trade-off between the two, while the AUC–ROC assesses the model’s capability to differentiate between churners and retained customers.

In the business context, false negatives, cases where actual churners are misclassified as non-churners, carry a higher financial risk than false positives. Therefore, recall will be emphasized during model selection to ensure that high-risk customers are correctly detected. The final model will be chosen based on an optimal balance between recall and precision, complemented by interpretability analysis using SHAP (SHapley Additive exPlanations) to clarify the relative influence of input variables on churn predictions.

3.8 Model Deployment

The best-performing model will be deployed as an API-driven service using FastAPI. This setup enables high-speed inference and seamless integration with systems. When customer data is submitted, the API will return a churn probability and accompanying interpretability insights derived from SHAP values.

These predictions will feed into a dashboard that segments customers by risk level and highlights key churn drivers. The modular design supports regular model retraining to reflect emerging behavioral trends, ensuring that the churn detection system remains adaptive and actionable for business teams.

The API-driven architecture ensures modularity, scalability, and continuous learning, allowing for retraining with new customer data as behavioral trends evolve. This approach bridges predictive analytics and business decision-making, enabling proactive interventions that enhance customer retention and revenue stability in the telecommunications sector.

4. CHAPTER: EXPECTED OUTCOMES

The proposed study seeks to develop a robust, machine learning–driven framework for predicting customer churn using the Bulgarian telecommunication dataset. By leveraging real-world customer data, the study will generate actionable insights that can be adapted to any telecommunications market. The outcomes will specifically address key research gaps related to feature diversity, model interpretability, and business integration of predictive analytics.

In relation to the first objective, which involves identifying challenges and limitations in existing churn prediction literature, the study will provide a comprehensive analysis of prior approaches, highlighting deficiencies such as inadequate handling of class imbalance, low interpretability of black-box models, and limited generalization across market contexts. This review will establish the foundation for a more transparent and transferable modeling framework.

Aligned with the second objective, the research will develop predictive models capable of capturing customer churn behavior through advanced machine learning algorithms. The models will integrate variables to create a multidimensional representation of churn risk. Feature engineering and class rebalancing using SMOTE will further enhance the dataset’s predictive strength and fairness.

In fulfilling the third objective, the study will rigorously evaluate and compare multiple machine learning algorithms based on performance metrics such as recall, precision, F1-score, and ROC-AUC. Emphasis will be placed on maximizing recall to ensure that high-risk churners are correctly identified, while maintaining interpretability through model explainability tools such as SHAP (SHapley Additive exPlanations).

The best-performing model will be deployed through a FastAPI-based predictive dashboard. The dashboard will display churn probabilities, key explanatory features, and visual analytics for decision support, enabling managers to design proactive retention campaigns grounded in data-driven insights.

Overall, the study is expected to produce a scalable, interpretable, and business-oriented churn prediction framework that enhances strategic decision-making and customer retention in the telecommunications sector. By combining predictive accuracy with explainable AI and seamless system integration, the research will minimize revenue leakage, improve marketing efficiency, and strengthen customer relationships. Beyond Bulgaria, the framework offers a replicable model for other telecommunication markets seeking to transition toward evidence-based, customer-centric operations.

References

- Chang, V., Hall, K., Xu, Q. A., Amao, F. O., Ganatra, M. A., & Benson, V. (2024). Prediction of customer churn behavior in the telecommunication industry using machine learning models. *Algorithms*, 17(6), 231. Retrieved from <https://research.tees.ac.uk/files/86885815/algorithms-17-00231.pdf> doi: 10.3390/a17060231
- Khan, M. T., Hall, K., Xu, Q. A., Amao, F. O., Ganatra, M. A., & Benson, V. (2024). Customer churn behaviour in the telecommunication industry. *Algorithms*, 17(6), 231. Retrieved from <https://doi.org/10.3390/a17060231> doi: 10.3390/a17060231
- Mwaura, E. T. (2021). *Effects of customer experience management on customer churn in the telco industry in kenya*. Project Report. Retrieved from <https://afribary.com/works/effects-of-customer-experience-management-on-customer-churn-in-the-telco-industry-in-kenya> (Accessed: 2025-09-27)
- Nagarkar, P. (2022). A customer churn prediction model using xgboost for the telecommunication industry in nepal. *Procedia Computer Science*, 215, 652–661. Retrieved from <https://www.sciencedirect.com/science/article/pii/S187705092202138X> doi: 10.1016/j.procs.2022.12.083
- Poudel, S., & Sharma, S. (2024). Explaining customer churn prediction in telecom industry using interpretable machine learning techniques. *Information Sciences Letters*, 36, 100043. Retrieved from <https://www.sciencedirect.com/science/article/pii/S2666827024000434> doi: 10.1016/j.isl.2024.100043
- Rahman, M., Bista, R., & Poudel, K. (2024). Explaining customer churn prediction in the telecom industry using interpretable machine learning techniques. *Information Sciences Letters*, 36, 100043. Retrieved from <https://www.sciencedirect.com/science/article/pii/S2666827024000434> doi: 10.1016/j.isl.2024.100043
- Saha, L., Tripathy, H. K., Gaber, T., El-Gohary, H., & El-Kenawy, E. S. M. (2023). Deep churn prediction method for telecommunication industry. *Sustainability*, 15(5), 4543. Retrieved from <https://www.mdpi.com/2071-1050/15/5/4543> doi: 10

.3390/su15054543

- Shahabikargar, M., Beheshti, A., Mansoor, W., Zhang, X., Foo, E. J., Jolfaei, A., ... Shabani, N. (2025). Churnkb: A generative ai-enriched knowledge base for customer churn feature engineering. *Algorithms*, 18(4), 238. Retrieved from <https://www.mdpi.com/1999-4893/18/4/238> doi: 10.3390/a18040238
- Tokmakov, D. (2024). *Customer churn dataset containing real records from a leading bulgarian telecom operator, specifically for business customers (version 1)*. Mendeley Data. Retrieved from <https://data.mendeley.com/datasets/nrb55gr66h/1> (Licensed under CC BY 4.0) doi: 10.17632/nrb55gr66h.1
- Usman-Hamza, F. E., Balogun, A. O., Capretz, L. F., Mojeed, H. A., Mahamad, S., Salihu, S. A., ... Salahdeen, N. K. (2022). Intelligent decision-forest models for customer churn prediction. *Applied Sciences*, 12(16), 8270. Retrieved from <https://www.mdpi.com/2076-3417/12/16/8270> doi:10.3390/app12168270
- Zhang, T., Moro, S., & Ramos, R. F. (2022). A data-driven approach to improve customer churn prediction based on telecom customer segmentation. *Future Internet*, 14(3), 94. Retrieved from <https://www.mdpi.com/1999-5903/14/3/94> doi: 10.3390/fi14030094

5. APPENDICES

5.1 Project Appendix

This appendix provides essential details regarding the project's condensed execution schedule and the financial resources required for its successful completion.

5.1.1 Project Timeline: 4 Months

The proposed research will follow an intensive, structured four-month schedule, emphasizing rapid progression from methodology design to the final dissertation submission.

P1: Phase I: Design and Literature Review

P2: Phase II: Data Preprocessing and Feature Extraction

P3: Phase III: Model Development and Training

P4: Phase IV: Evaluation, Analysis, and Final Drafting

5.2 Project Budget

The project budget below outlines the estimated costs required for the development, analysis, and dissemination of the machine learning based telecommunication customer churn prediction framework. All estimates are presented in Kenya Shillings (KSh) and are based on current institutional and market rates.

5.2.1 Estimated Project Budget (KSh)

Category	Estimate	Estimated Cost (KSh)
A. Computational Resources		
Data Storage and Access	250 GB cloud storage	5,000
Software Licenses	Annual license for data visualization	10,000
Subtotal A		15,000
B. Administration		
Ethical and Institutional Review	Standard university ethics	10,000
Internet and Communication	high-speed internet @ KSh 4,000/month	20,000
Overleaf/LaTeX Subscription	Overleaf Premium plan	6,000
Subtotal B		36,000
C. Data and Dissemination		
Data Access	Use of public dataset	0
Publication Fees	Open-access journal submission fees	40,000
Subtotal C		40,000
D. Equipment and Materials		
External SSD Backup Drive	1 TB external drive	10,000
Stationery and Printing	Report printing	8,000
Subtotal D		18,000
Subtotal (A + B + C + D)		109,000
Miscellaneous / Contingency (10)	unforeseen costs	30,000
Total Estimated Budget		139,000

Table 5.1: Proposal budget