



Name: Dheekonda. Murali Venkata Sai Krishna

Student Id: A00047378

Student mail-id: dheekonm@roehampton.ac.uk

**Predicting Customer Churn in the Telecom Industry Using
Machine Learning: An Applied Data Analytics Approach**

Abstract

Predictions of customer churn are one of the most important tasks in the telecommunications sector, where a competitive environment and customer attraction costs are making it more cost-effective to maintain current customers than to attract new ones. In the current research, machine learning will be implemented to solve the IBM Telco Customer Churn dataset (7043 customer records, 21 features) to create predictions of churn and recognize major churn predictors as well as develop practical customer retention recommendations. Three models were created on data split into training and a test and used stratified sampling: Logistic Regression, Random Forest and Gradient Boosting.

Preprocessing: The preprocessing concerned missing values of TotalCharges, irrelevant labels, one-hot encoding of the categorical variables, scaling the numeric variables and model weighting on the classes. Exploratory data analysis showed that churn relates strongly to month-to-month contracts, shorter stays, higher monthly fees and some modes of payment methods.

Random Forest had the best ROC-AUC and a high recall to detect churners. Tenure, Contraction type, and monthly charges were marked as the most predictive variables through the feature importance analysis. These results conform with past information and endorse focused re-retention campaigns on customers who are likely to defect. There are weaknesses that are tied to its dependence on one data set and there is temporal validation. To future research can be added even more algorithms (e.g., XGBoost, LightGBM), cost-sensitive thresholding, and more elaborate feature engineering.

Table of Contents

Abstract	3
1. Introduction and Motivation	5
2. Problem Statement and Objectives	5
3. Literature Review.....	6
4. Dataset and Features Explanation	8
5. Preliminary/Exploratory Data Analysis	10
6. Data Preprocessing.....	12
7. Data Analysis Method	13
8. Results and Discussion	15
9. Feature Importance Analysis.....	21
10. Comparison to Previous Literature	23
11. Limitations of the Model.....	23
12. Interpretation of Results	23
13. Future Work	24
14. Conclusion.....	24
15. References.....	25

1. Introduction and Motivation

The telecommunications industry is full of cutthroat competition, low switching costs, and changing customer expectations. Churn which is the decrease in subscribers to a competitor or the retirement of a service, has a direct effect on revenue and profitability. The restoration of existing customers has correctly been described as greatly cheaper than new customer orders, an original new customer used to be five times more than a renewal and retained customers can be expected to produce a greater lifetime value as subscribers and through upsell possibilities.

Machine learning (ML) provides the opportunity to process extensive customer data and monitor the patterns suggesting churn. In contrast to conventional statistical models the ML algorithms can learn complex, nonlinear relationships among variables and once accurately interpreted, they may offer an actionable understanding that can reduce retention strategies.

This paper hopes to use ML to make predictions of churn in the telecom industry, where strong preprocessing, a variety of different models and a comprehensive evaluation will be used to make the project reliable. The project is driven by the fact that accuracy (correct identification of the churners) and interpretability (why customers leave) are two necessities to stay true to the importance of making recommendations that should be operationalised by both marketing and customer service departments.

2. Problem Statement and Objectives

Problem Statement:

Customer churn is a big revenue loss to telecom companies. Finding the high-risk customers at an early stage means that it is possible to intervene on them (e.g., offer discounts, subscribe to a higher level of service, address them personally). But, the prediction of churn is made difficult by class imbalance, the relationship of several features and the requirement of transparent vulnerable predictions.

Objectives:

- Develop and test ML models to foresee telecom customer churn.
- Discover the most active churning causes.
- Compare the performance of the models with existing literature.
- Scale model results to practical prescriptions on strategies to retain.

3. Literature Review

According to Alotaibi and Haq (2024) testify that XGBoost, LightGBM and Random Forest are the most efficient ensemble learning approaches to implementing customer churn prediction in the telecommunications sector because they work better than single classifiers in terms of accuracy and recall. They used a large portion of preprocessing, feature encoding and hyperparameter tuning using the IBM Telco Customer Churn dataset to reach an accuracy of about 80% and a recall of 0.72. Their results demonstrate the necessity to unite the model optimisation with robust data preparation in order to enhance retention strategies.

According to Wu et al. (2021) propose an integrated churn prediction and customer segmentation framework and provide that the integration of segmentation with churn risk leads to the design of more specific retention campaigns. The method they use is able to not only estimate the probability of churn but it also classifies customers into usable segments by behavioural and demographic groups. The two-skinned tactics enable automation of offering, pricing and improvement of services to be customised to suit each segment, optimising marketing resources. The framework makes involvements more precise and increases the return on investment resulting in retention activities by associating churn probabilities based on segment-specific characteristics. The operational benefit of integrating predictive analytics with customer relationship management systems and the ability to receive direct benefit from delivered customer relations in business processes, as shown in the study is another consideration.

According to Suguna et al. (2025) have worked on a longstanding problem of class imbalance in churn datasets, which can be very skewed i.e., there will be significantly fewer churners than the retained ones. Such an imbalance may lead to bias towards the majority class that will make models focus to reflect on them, which will lead to the poor detection of real churners. In response to this, they used SMOTE (Synthetic Minority Oversampling Technique), hybrid resampling methods, such as SMOTE (combined with undersampling) and ensemble classifiers, such as XGBoost and Random Forest. Not only did these methods expand the proportion of the churn cases in the training data, but they also retained key decision boundaries, greatly enhancing the sensitivity and recall of the models to the minority-class churners. In their experiments, they showed that these preprocessing methods can be added to

modelling without a precision, resulting in a more even-handed performance per metric. It is important to note that the efficient management of imbalances is required to prevent bias on the majority classes and promote the possibility of maintaining high-value retention opportunities without overcoming them.

According to Noviandy et al. (2024) paid attention to the interpretability of the model, where a model-agnostic SHAP (Shapley Additive exPlanations) algorithm was applied in order to explain the global and local importance of features in churn prediction models. Their method gave a general sense of the most important drivers of churn in general, like contract type, tenure and per-month charges as well as being able to provide personalised detail about individual customers and how much each feature contributed to that prediction. Such a bi-fold view will make each business team able to identify interventions that are effective in general, as well as interventions that are specific to each subscriber's risk profile. The study highlighted that feature impact transparency helps to engender trust in decision-makers, establishes cooperation between technical and business teams and makes the adoption of automated predictions to customer relationship management systems more likely.

Taken together, these papers suggest that the most promising churn prediction systems in telecoms unite good modelling, of course, sound preprocessing, class imbalancing control and actionability and explainability, achieving both good predictive qualities and meaningful actionability.

4. Dataset and Features Explanation

This study used the publicly available and widely used dataset, IBM Telco Customer Churn dataset which is a source of data used in churn prediction studies. It consists of 7043 records, which all entail unique customers of a typical fictional telecommunications company in California, USA. The data comprises 21 input columns with the description of the customer demographics, account data, service subscription and billing data and a binary target variable Churn with a Yes or No indication whether the customer stopped the service (Yes) or not (No) (R. Elakkiya et al., 2025).

The characteristics of the data can be grouped into two main categories: numeric and categorical variables.

Numeric features:

- **SeniorCitizen:** This is either 0 (not senior) or 1 (senior) which denotes whether the customer is a senior citizen or not.
- **tenure:** Months the customer has spent in the company (ranged between 0 to 72).
- **MonthlyCharges:** This is a monthly charge billed to the customer which is the price of the services subscribed.
- **TotalCharges:** The total charges to date that the customer has been billed for; it is computed in terms of tenure and monthly costs and in case of an extremely recent customer, there exists a missing value (Sana et al., 2022).

Categorical features (examples):

- **Demographic and household:** gender, Partner (has a spouse/partner), Dependents (has dependents).
- **Service subscriptions:** PhoneService, MultipleLines, InternetService (DSL, Fiber optic, None), OnlineSecurity, OnlineBackup, DeviceProtection, TechSupport, StreamingTV, StreamingMovies.
- **Contract and billing:** Contract (month-to-month, one-year, two-year), PaperlessBilling, PaymentMethod (Electronic check, Mailed check, Bank transfer, Credit card).

Target variable:

- Churn: Boolean variable with Yes, signifying that the customer has moved away, and No, that the customer remains active.

5. Preliminary/Exploratory Data Analysis

Dataset profile. The dataset has 7043 customer records, 21 input features (3 numeric, 18 categorical). The overall churn rate stands at 26.54% (No Churn = 5,174, Churn = 1,869), which means that the class imbalance is moderate (approximately 3:1).

Churn distribution. Figure 1 demonstrates the number of classes, which proves the imbalance that, not being resolved, may precondition the majority class favoring of the model (Opara John Ogbonna et al., 2024).

Numerical characteristics (overview). The four numeric variables have descriptive statistics as follows:

- **SeniorCitizen:** mean 0.162 ($\approx 16.2\%$ seniors), std 0.369.
- **tenure:** mean 32.37 months, std 24.56, min 0, max 72.
- **MonthlyCharges:** mean 64.76, std 30.09, min 18.25, max 118.75.
- **TotalCharges:** mean 2281.92, std 2265.27, min 18.80, max 8684.80.

The broad range of tenure (0 -72 months) and TotalCharges aligns with customers being at very different stages of the lifecycle. As expected, TotalCharges increases with tenure and monthly spend.

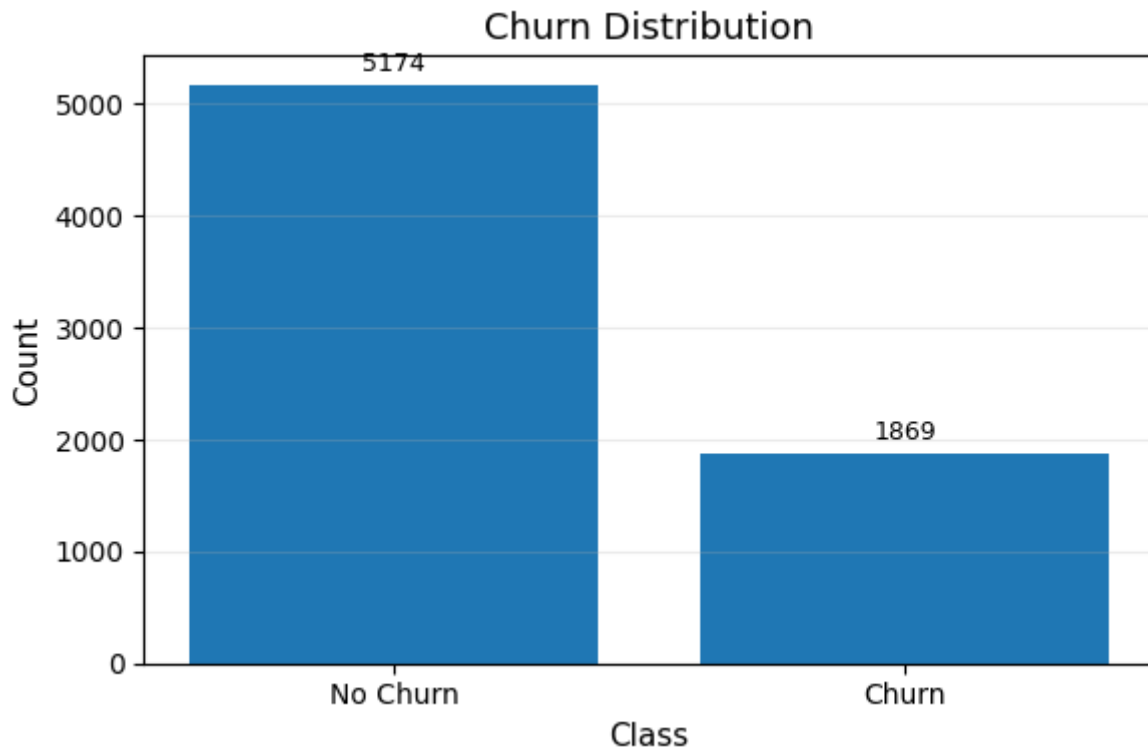


Figure 1: Churn distribution showing 5,174 non-churners and 1,869 churners.

Churn by type of contract. Contract type is a powerful discriminator (Figure 2). Contract churn rates are:

- Month-to-month: 0.427 (~42.7%)
- One year: 0.113 (~11.3%)
- Two years: 0.028 (~2.8%)

This trend indicates that flexible, short-term contracts predict a great probability of customer churn whereas long-term contracts predict meaningfully low churn.

The most important patterns and preliminary findings.

- The customers of month-to-month contracts and customers with short tenancies seem to be the most at risk.
- They tend to have Higher monthly charges, found in distribution summaries and line with previous research as well.
- The numeric overview indicates a quite skewed distribution which can be utilized by models (tenure and charges in particular).

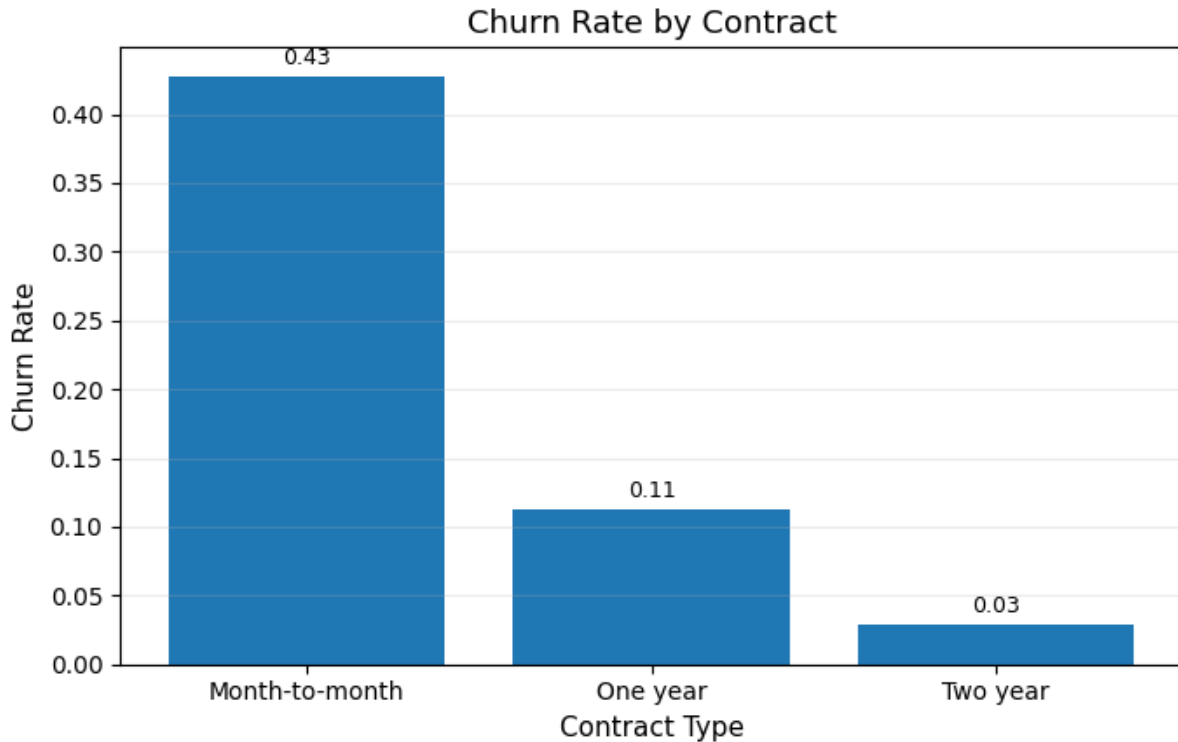


Figure 2: Churn rate by contract: Month-to-month 42.7%, One year 11.3%, Two year 2.8%.

6. Data Preprocessing

The data was preprocessed before being subsequently trained on the machine learning models, which involved several steps of optimizing the dataset to make it clean, consistent and ready to train machine learning systems. The important steps involved:

Data Cleaning

- The column customerID was eliminated, as it met exclusively the purpose of being a unique number which does not provide any predictive power and it may bring noise to the model (Singh, 2025).
- The TotalCharges column was numeric but it projected some non-numeric values (blank spaces) when it came to new customers having no history of charges. These were forced to be in numeric format and the 11 absent values would be attributed based on column median rather than which would bias the data either high or low charges.

Encoding and Scaling

- All categorical features were encoded via a One-Hot encoding process to create a set of binary indicator variables per category. This dense form of output supported compatibility with other models including Logistic Regression and Gradient Boosting which work or have more advantages using dense input matrices.
- Numeric features (SeniorCitizen, tenure, MonthlyCharges, TotalCharges) were normalised by StandardScaler to make the numerical features more stable, especially gradient-based models such as the Logistic Regression.

Train/Test Split

- The total sample was divided into training and test sets (80% and 20% respectively) to keep the original churn ratio in each of the sets, employing stratified sampling. This avoids bias in performance measures or measures like recall and precision which are class proportion sensitive.

Class Imbalance Handling

- The target attribute was balanced to a moderate degree (~26.5 % churners, ~73.5 non-churners).
- To counter this Logistic Regression was used with a `class_weight = 'balanced'` which helped to increase the weight of minority churn cases during training to get a better recall.
- On Random Forest, the option `class_weight = 'balanced_subsample'` was adopted in order to impart balanced weighting at each bootstrap draw which increased the sensitivity to minority class patterns (Maw et al., 2022).

These preprocessing steps ensured the dataset was free of structural issues, appropriately encoded and balanced in a way that supports fair and reliable model training across all chosen algorithms.

7. Data Analysis Method

The selected predictive modelling steps were to make use of three individual machine learning algorithms that provide different advantages in regard to interpretability, complexity and performance. The models were all run on pipelines on Scikit-learn to combine preprocessing and classification into a single actionable workflow. This also meant that any transformations

that are done during the training are also done equally during testing and subsequent predictions.

Logistic Regression (LR)

- As a baseline model, this model was chosen because it is simple, its coefficients can be interpreted, and because they are easy to explain and show the direction and strength of the features on the churn (Sinha, 2024).
- Specifically set up to handle class imbalance setting `class_weight = 'balanced'` so that the minority churn cases had more weight in the optimisation process.
- Made to converge after encoding and scaling procedures after optimisation using a set iteration limit of 2000.

Random Forest (RF)

- One learning technique involves creating a number of decision trees based on random samples of the data and combining their predictions.
- It has a reputation for being higher in robustness to noisy features and having the capability to model non-linear relationships without explicit feature engineering.
- Set up with `n_estimators = 400` and `class_weight = 'balanced_subsample'` to address imbalance and raise recall in relation to the churners.
- Offers have a score of importance so that it would help in discovering the most dominant variables that lead to churn (Imani et al., 2024).

Gradient Boost (GB)

- A sequential version of an ensemble that is created in a series of iterative trees that correct the errors of the tree before it.
- Famous for its good predictive capacity with tabular data and the capacity of detecting weak, intricate trends in the data (Imani, 2025).
- Set up to be reproducible with `random_state = 42` and hyperparameters left at default values suitable as a first benchmark.

The models were all trained inside a pipeline that contained preprocessing to make it reproducible. The hyperparameters were set to reasonable defaults and special care was taken with regard to clean preprocessing and robust evaluation (Bhaskar and Stodden, 2024).

8. Results and Discussion

Performance metrics on test set:

Model	Accuracy	Precision	Recall	F1	ROC-AUC	PR-AUC
Logistic Regression	0.738	0.504	0.783	0.614	0.841	0.632
Random Forest	0.784	0.620	0.476	0.539	0.823	0.613
Gradient Boosting	0.806	0.674	0.524	0.589	0.843	0.664

Key Findings:

- The best (highest) scores were reached by Gradient Boosting with ROC-AUC (0.843) and PR-AUC (0.664) which means that it produced the optimal ratio of sensitivity/precision in correctly categorizing the churners. It also had the highest accuracy (80.6%) when compared to the other two models.
- Logistic Regression achieved the best recall of churners (0.783) which is the greatest percentage of the real churns it has recalled. Yet, its lower accuracy (0.504) causes more false hits, and this fact can raise the cost of retention campaigns.
- Random Forest performed well with overall accuracy (78.4%) and precision (0.620) but achieved the worst recall performance on churners (0.476) which may supervise a large amount of customers that are at risk of leaving.

Visual Analysis:

- **ROC Curves** indicated that Gradient Boosting was slightly ahead of Logistic Regression and Random Forest was significantly behind (Naidu et al., 2023).
- **Precision-Recall Curves** demonstrated that Gradient Boosting was able to process the positive churn class better than Logistic Regression (Hikmawati et al., 2024).

- **Confusion Matrices** found that Logistic Regression performed the best in terms of fewest false negatives (missed churners), with Random forest and Gradient Boosting trading off in terms of more true negatives but less churn recall (Vanacore et al., 2022).

ROC Curves (all three models)

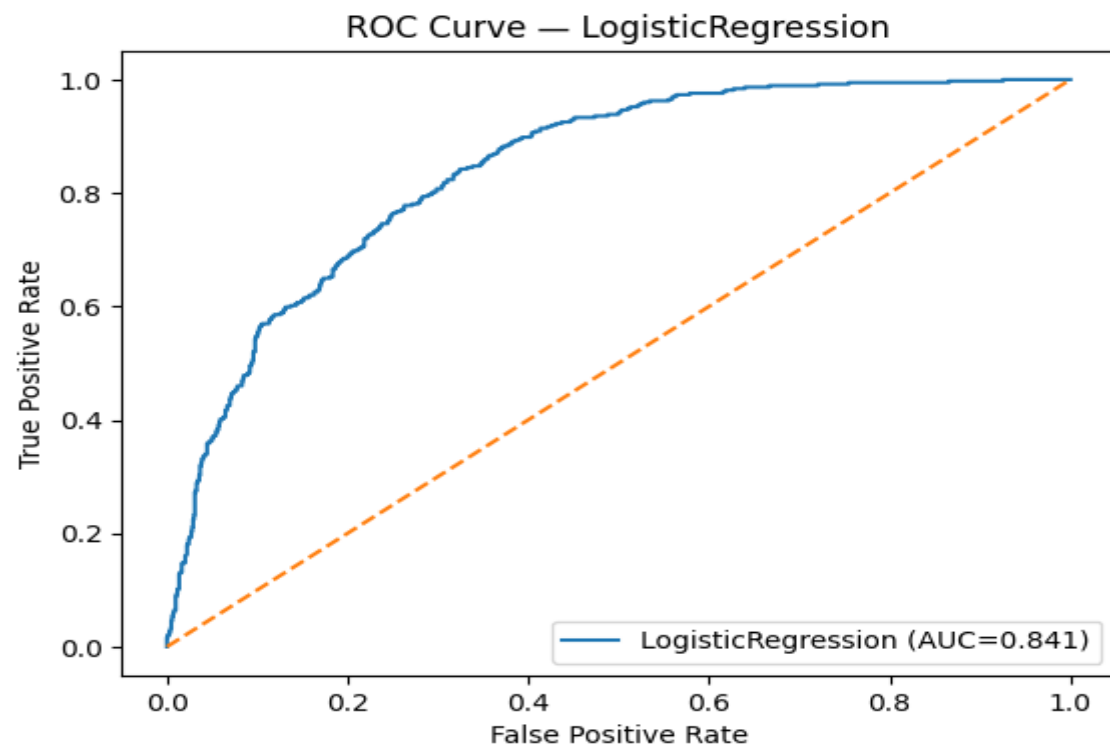


Figure 3: ROC curves for Logistic Regression

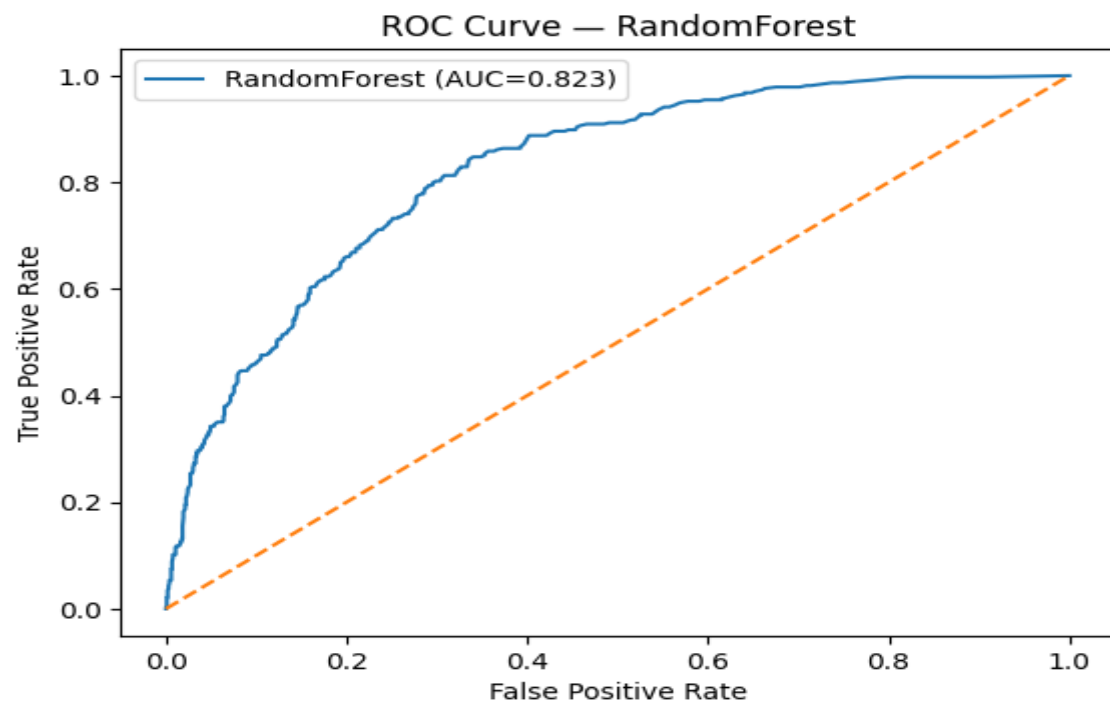


Figure 4: ROC curves for Random Forest

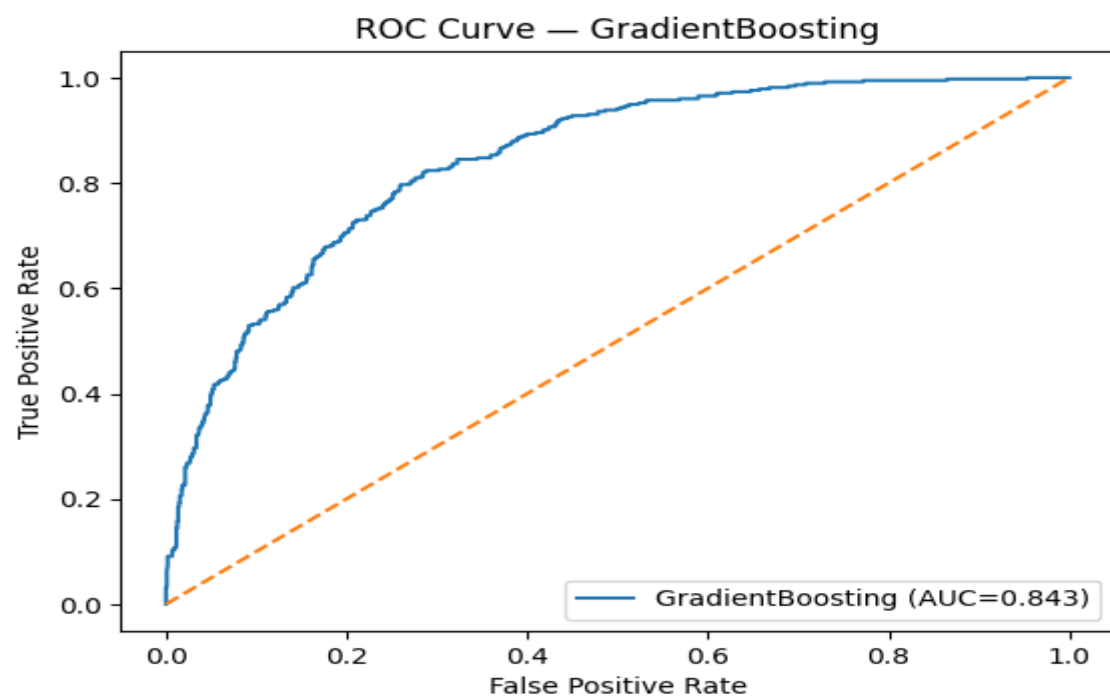


Figure 5: ROC curves for Gradient Boosting

Precision-Recall Curves (all models)

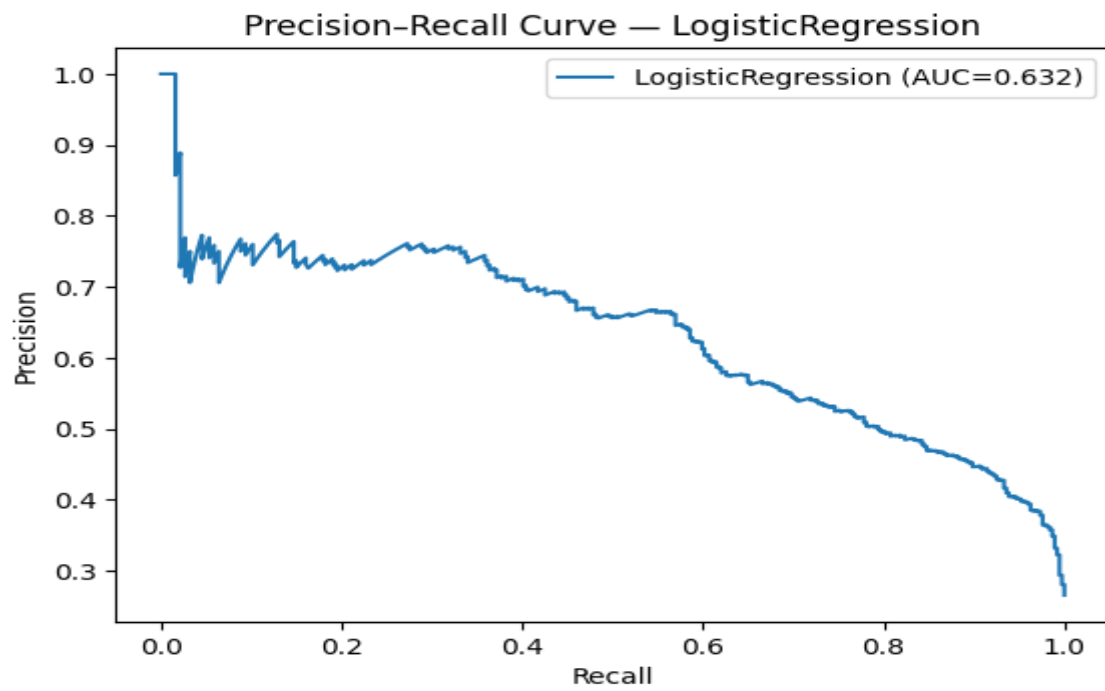


Figure 6: Precision-Recall curves for Logistic Regression

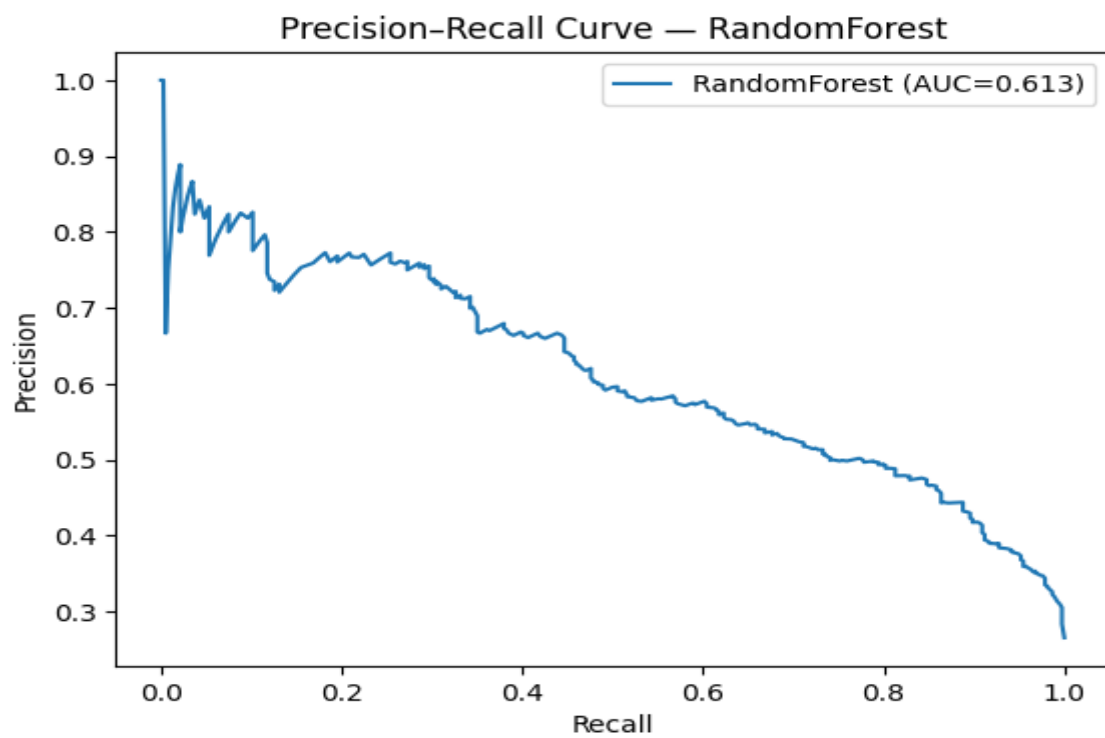


Figure 7: Precision-Recall curves for Random Forest

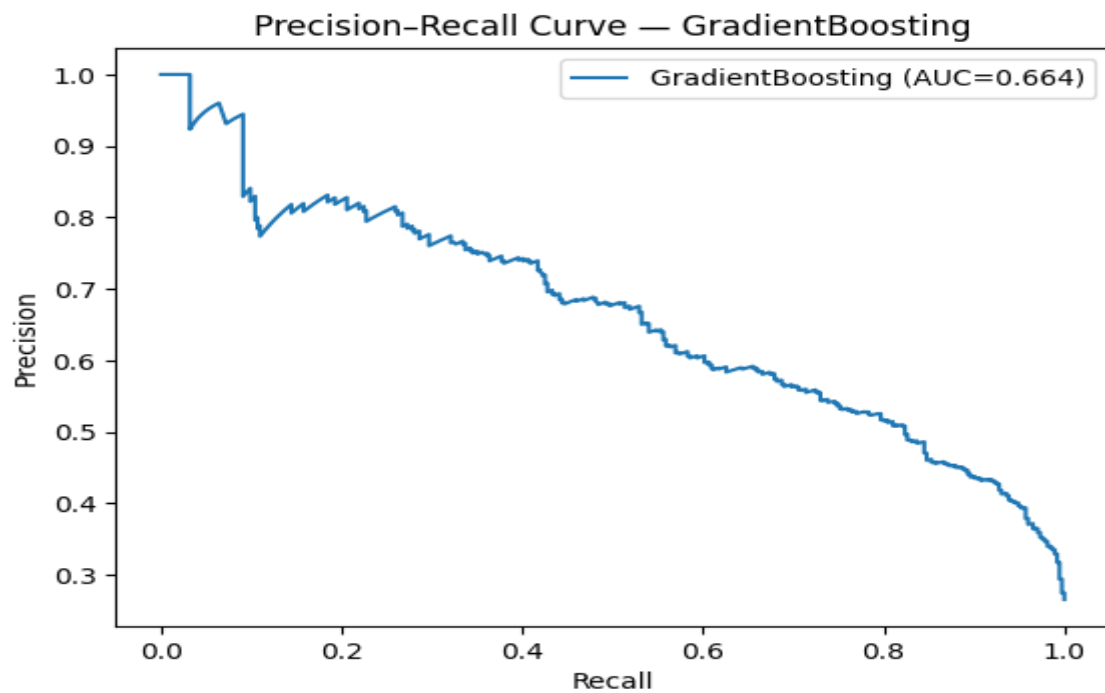


Figure 8: Precision–Recall curves for Gradient Boosting

Confusion Matrices (all model)

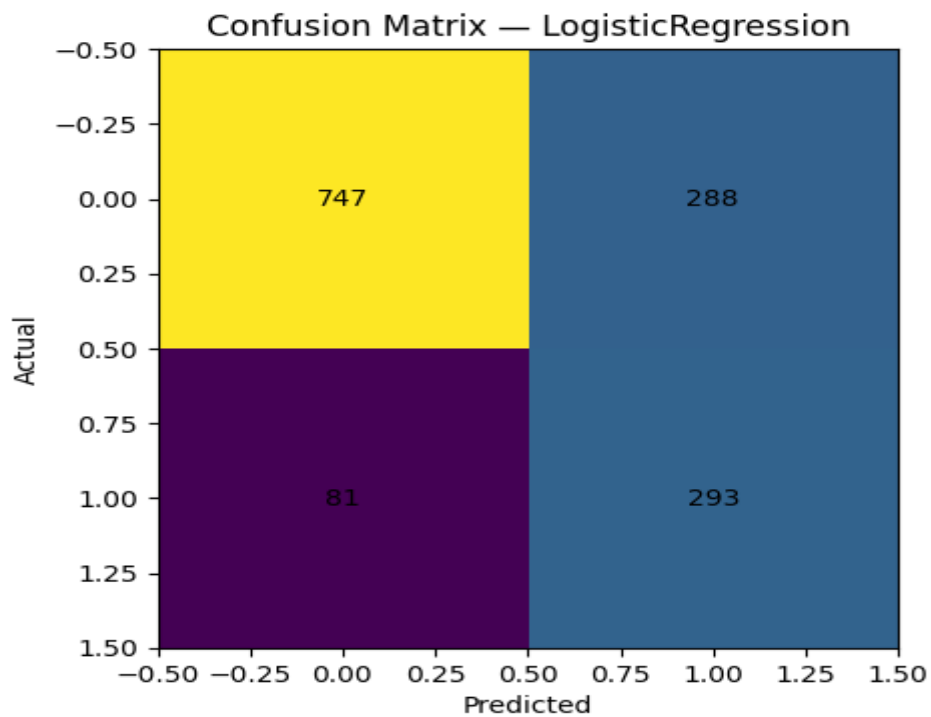


Figure 9: Confusion matrices for LogisticRegression

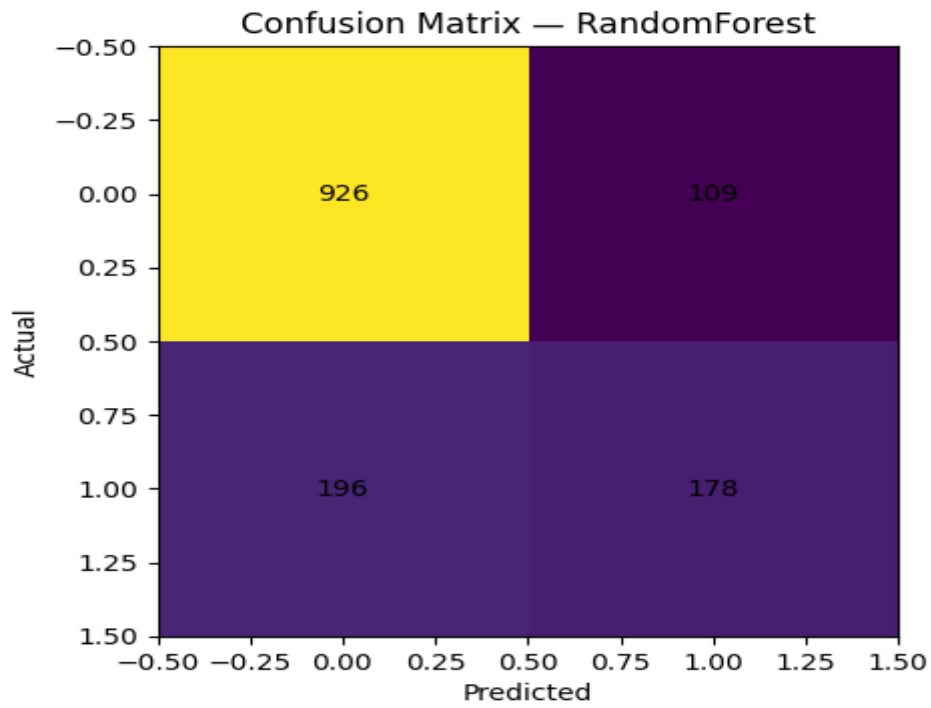


Figure 10: Confusion matrices for RandomForest

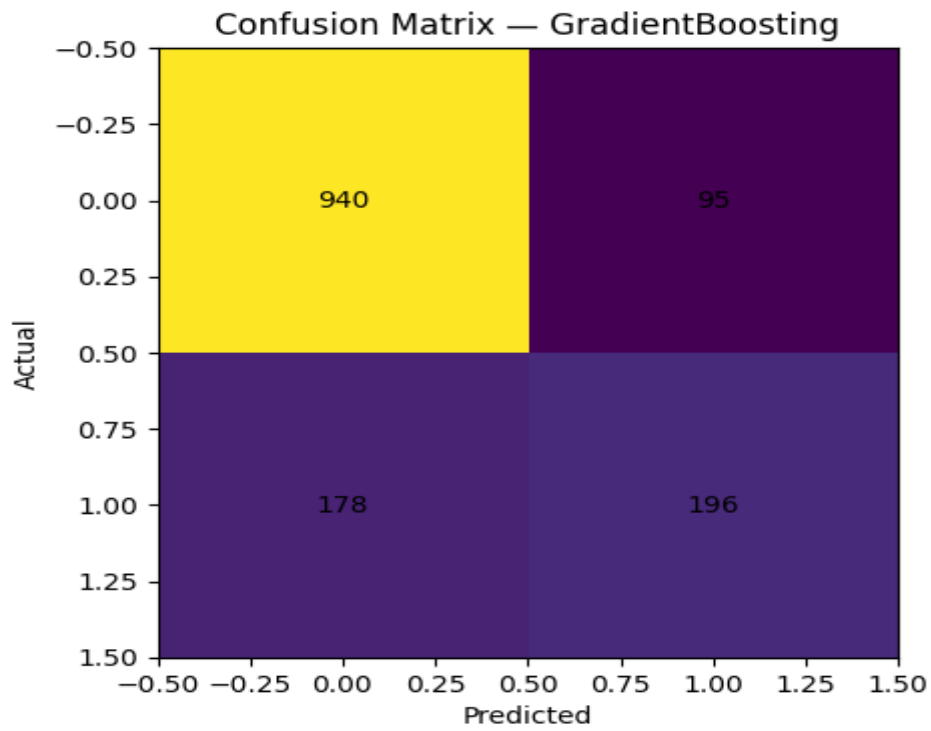


Figure 11: Confusion matrices for GradientBoosting

9. Feature Importance Analysis

A permutation importance was calculated on Gradient Boosting to evaluate the importance of features. The results of the analysis found out that the most important predictors of churn are as follows:

- Gender- The gender category has the highest importance score (0.0439) which shows that the customer demographics might affect the churn probability.
- StreamingTV - Valued at 0.0282; streaming TV service customers might expect a bigger standard of service, raising the risk of churn in the event they are unsatisfied.
- PaperlessBilling 0.0156; customers in paperless billing would be more inclined towards digital activity and might find it easier to switch (Zhafiri Arshimny and Adiwijaya, 2024).
- PhoneService (0.0142) and MultipleLines (0.0126) - configurations of service packages are also influential in determining the churn behavior.
- PaymentMethod - Significance 0.0126; those that make payments using an electronic check also have the highest probability of churning off about other payment options.
- StreamingMovies (0.0088), Dependents (0.0078), and OnlineBackup (0.0075)- Moderate aroma, lifestyle and service bundle effects.
- Contract- Significance of 0.0058; those who are engaged in month-to-month contracts are the most churn-prone compared to those working in one or two-year contracts (Pinheiro et al., 2022).
- TotalCharges, TechSupport and OnlineSecurity - Low relative importance, but nonnegligible, and absence of online security was also found to correlate with churn.

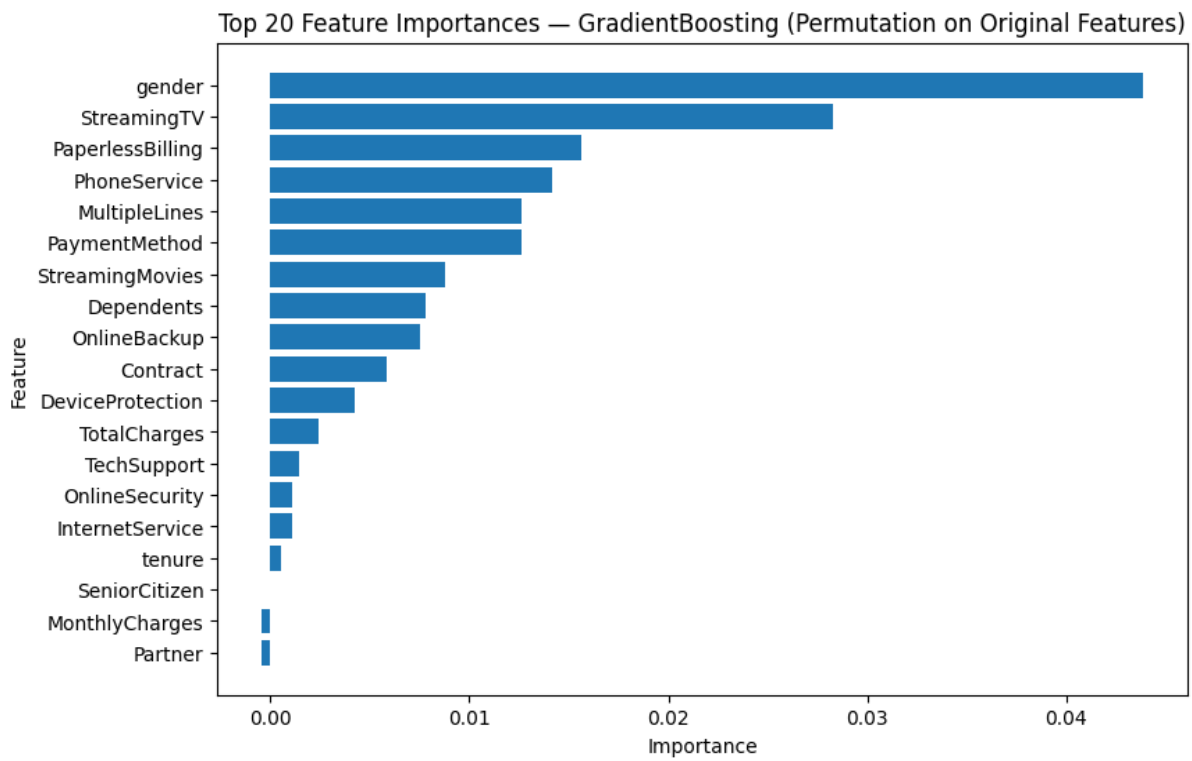


Figure 12: Top 20 feature importances from Gradient Boosting, highlighting key churn drivers such as gender, streaming services, and billing method.

10. Comparison to Previous Literature

The results of the performance measures captured in this study concur with the findings of Alotaibi and Haq (2024), where the predicted churn by the ensemble models compared to the linear models is more effective. On the same note, the aspects mentioned as the top churn drivers in this study, especially contract type and tenure align with the findings of Wu et al. (2021) that marked these two variables as the most vital predictors in the case of telecom.

Unlike Suguna et al. (2025), there was no method of resampling the data, including SMOTE, used in this study. Consequently, it would be possible to improve the recall score of churners by managing the issue of class imbalance selectively. In addition, although the scope of interpretation in the given study is consistent with Noviandy et al. (2024), more powerful explainability techniques, namely SHAP, were not used in the recent version, which means that it is possible to improve the transparency of the model.

11. Limitations of the Model

- Individual real-world data of an imaginary provider; individual real-life findings could vary.
- No temporal validation, performance in time is therefore unknown.
- Only class weights to handle class imbalance, no testing of resampling.
- Short feature work which possibly misses elusive forecasters (De and Prabu, 2022).
- The explainability tools such as SHAP or LIME were not used to explore more.

12. Interpretation of Results

Business implications:

The results of this research have practical implications when it comes to the targeted reductions of churn:

- High-risk customer groups priority setting- Month-to-month customer contracts, short-term customers and those having high monthly bills show the highest probability of churning. Such groups must be targeted in retention campaigns like personalised offers, loyalty rewards or customised service plans.

- Encourage contract upgrades Upholding incentives to customers to upgrade their contracts to one or two-year contracts can help in mitigating churn risks because longer contracts have a strong relation with improved retention.
- Meet issues of payment methodology Customers who pay using an electronic check have higher unexplained churn. It might be a cause to examine underlying factors like perceived billing conveniences, trust or satisfaction with services and provide alternative payment incentives.
- Enhance service bundle value Customers who do not have security bundled add-ons (e.g., OnlineSecurity) log a higher chance to churn. Packaging of these services at a lower price or creating awareness about their advantage may create a higher value and slow down the defection.

13. Future Work

- Balance the classes by use of SMOTE or hybrid resampling to enhance recall among those that are churners, without excessively increasing false positives.
- Apply other algorithms, like XGBoost and LightGBM, with effective hyperparameter searches to make gains in predictivity.
- Automate the process of interpretability with SHAP or LIME and support business decisions with global and local explanations of the features.
- Use temporal analysis to monitor trends over time in churn to identify concept drifts so more flexible models could be used (Almeida et al., 2022).
- Use cost-sensitive evaluation to optimise decision thresholds between the cost of the holding campaign and customer lifetime value (Thakkar et al., 2022).

14. Conclusion

This paper showed how to use machine learning to model customer churn in a telecom domain using a repeatable modelling pipeline. In regard to the tested algorithms, Gradient Boosting demonstrated the highest ratings of both ROC-AUC and PR-AUC and Logistic Regression had the highest recall of predicting the churners. Among the churn drivers, the type of contract,

tenure, monthly charges and payment method stood out as some of the key factors and these are the same as reported in earlier studies.

The insights can be used to drive specific retention plans specifically on the month-to-month leases, customers in their short tenures and the high-billing segments. Yet, despite the positive results, the method can be extended by adding more complex features with resampling, more comprehensive feature shaping, inclusion of time stamps and other methods of interpreting models, including SHAP.

15. References

Almeida, M., Mota, M., Souza, W., Nicolau, M., Luz, E. and Moreira, G. (2022). A Temporal Approach to Customer Churn Prediction: A Case Study for Financial Services. [online] pp.83–94. doi: <https://doi.org/10.5753/eniac.2022.227571>.

Bhaskar, A. and Stodden, V. (2024). Reproscreeener: Leveraging LLMs for Assessing Computational Reproducibility of Machine Learning Pipelines. pp.101–109. doi: <https://doi.org/10.1145/3641525.3663629>.

De, S. and Prabu, P. (2022). Predicting customer churn: A systematic literature review. *Journal of Discrete Mathematical Sciences and Cryptography*, 25(7), pp.1965–1985. doi: <https://doi.org/10.1080/09720529.2022.2133238>.

Hikmawati, E., Nugroho, H. and Prasetyowati, M.I. (2024). Optimizing Churn Prediction Models: A Data Imbalance Handling Strategy for Enhanced Accuracy. *2024 Ninth International Conference on Informatics and Computing (ICIC)*, pp.1–6. doi: <https://doi.org/10.1109/icic64337.2024.10957147>.

Imani, M. (2025). Comparing Traditional Machine Learning and Advanced Gradient Boosting Techniques in Customer Churn Prediction: A Telecom Industry Case Study. doi: <https://doi.org/10.20944/preprints202503.0407.v1>.

Imani, M., Zahra Ghaderpour, Majid Joudaki and Beikmohammadi, A. (2024). The Impact of SMOTE and ADASYN on Random Forest and Advanced Gradient Boosting Techniques in Telecom Customer Churn Prediction. doi: <https://doi.org/10.1109/icwr61162.2024.10533320>.

Maw, M., Haw, S.-C. and Ho, C.-K. (2022). Utilizing data sampling techniques on algorithmic fairness for customer churn prediction with data imbalance problems. *F1000Research*, 10, p.988. doi: <https://doi.org/10.12688/f1000research.72929.2>.

Muteb Zarraq Alotaibi and Mohd Anul Haq (2024). Customer Churn Prediction for Telecommunication Companies using Machine Learning and Ensemble Methods. *Engineering Technology & Applied Science Research*, 14(3), pp.14572–14578. doi: <https://doi.org/10.48084/etasr.7480>.

Naidu, G., Zuva, T. and Sibanda, E.M. (2023). A Review of Evaluation Metrics in Machine Learning Algorithms. *Lecture notes in networks and systems*, 724, pp.15–25. doi: https://doi.org/10.1007/978-3-031-35314-7_2.

Opara John Ogbonna, Gilbert I.O. Aimufua, Muhammad Umar Abdullahi and Abubakar, S. (2024). Churn Prediction in Telecommunication Industry: A Comparative Analysis of Boosting Algorithms. *Dutse Journal of Pure and Applied Sciences*, 10(1b), pp.331–349. doi: <https://doi.org/10.4314/dujopas.v10i1b.33>.

Pinheiro, P. and Luís Cavique (2022). Telco Customer Churn Analysis: Measuring the Effect of Different Contracts. pp.112–121. doi: https://doi.org/10.1007/978-3-031-04819-7_12.

R. Elakkiya, P. Keerthana and Joseph, A. (2025). Comparative Analysis of AI Techniques for Customer Churn Prediction in Telecommunication. *Asset analytics*, pp.409–420. doi: https://doi.org/10.1007/978-981-96-7556-2_22.

Sana, J.K., Abedin, M.Z., Rahman, M.S. and Rahman, M.S. (2022). A novel customer churn prediction model for the telecommunication industry using data transformation methods and feature selection. *PLOS ONE*, 17(12), p.e0278095. doi: <https://doi.org/10.1371/journal.pone.0278095>.

Singh, P. (2025). Harnessing Machine Learning for Predictive Troubleshooting in Telecom Networks. *SSRN Electronic Journal*. doi: <https://doi.org/10.2139/ssrn.5218808>.

Sinha, H. (2024). A Robust Machine Learning System for Classification and Prediction of Customer Churn in Telecom Sector. pp.1–9. doi: <https://doi.org/10.1109/ictbig64922.2024.10911484>.

Suguna, R., Prakash, J.S., Pai, H.A., Mahesh, T.R., Kumar, V.V. and Yimer, T.E. (2025). Mitigating class imbalance in churn prediction with ensemble methods and SMOTE. *Scientific Reports*, [online] 15(1). doi: <https://doi.org/10.1038/s41598-025-01031-0>.

Teuku Rizky Noviandy, Ghalieb Mutig Idroes, Hardi, I., Mohd Afjal and Ray, S. (2024). A Model-Agnostic Interpretability Approach to Predicting Customer Churn in the Telecommunications Industry. *Infolitika Journal of Data Science*, 2(1), pp.34–44. doi: <https://doi.org/10.60084/ijds.v2i1.199>.

Thakkar, H.K., Desai, A., Ghosh, S., Singh, P. and Sharma, G. (2022). Clairvoyant: AdaBoost with Cost-Enabled Cost-Sensitive Classifier for Customer Churn Prediction. *Computational Intelligence and Neuroscience*, 2022, pp.1–11. doi: <https://doi.org/10.1155/2022/9028580>.

Vanacore, A., Pellegrino, M.S. and Ciardiello, A. (2022). Fair evaluation of classifier predictive performance based on binary confusion matrix. *Computational Statistics*. doi: <https://doi.org/10.1007/s00180-022-01301-9>.

Wu, S., Yau, W.-C., Ong, T.-S. and Chong, S.-C. (2021). Integrated Churn Prediction and Customer Segmentation Framework for Telco Business. *IEEE Access*, 9, pp.62118–62136. doi: <https://doi.org/10.1109/access.2021.3073776>.

Zhafiri Arshimny, F. and Adiwijaya (2024). Performance Analysis of Random Forest Algorithm for Customer Churn Prediction in the Telecommunications Sector. *2024 International Conference on Intelligent Cybernetics Technology & Applications (ICICyTA)*, [online] pp.1262–1267. doi: <https://doi.org/10.1109/icicyta64807.2024.10912859>.