
Model Selection for Customer Churn Prediction in the Telecommunications Industry

Adam Cooke, Anthony Gheen, Nitin Shirsat
School of Data Science and Society
University of North Carolina
Chapel Hill, NC 27514

Abstract

This report aims to identify an improved methodology for implementing predictive models to identify customers at risk of churn. We perform an extensive comparison of six models used heavily in churn prediction. For each model, we identify the sampling method and scaling method that provided the highest F_1 score. Once the optimal sampling and scaling methods were identified, we further increase F_1 scores by optimizing hyperparameters for each of the models. Finally, we train and test model performance on three new datasets to determine if the identified methods can be used as a baseline for predicting customer churn in similar or different industries.

1 Introduction

Customer churn, defined as the termination of customer subscriptions to services, represents one of the critical challenges faced by businesses operating subscription-based models, especially in industries like telecommunications, banking, and online services. Churn directly impacts revenue streams, customer lifetime value, and overall market competitiveness. For telecommunications companies, particularly, customer acquisition costs (CAC) are substantially high, driven by significant investments in network infrastructure, marketing campaigns, and regulatory compliance. Thus, retention strategies driven by predictive analytics offer high returns on investment.

Effective churn prediction enables organizations not only to proactively engage and retain customers at risk but also to optimize marketing strategies by identifying and investing resources in high-value customers. Moreover, accurate churn prediction models help reduce unnecessary customer retention expenditure on customers unlikely to churn, improving overall operational efficiency.

Previous research in churn prediction spans a variety of analytical techniques, including logistic regression, decision trees, ensemble methods, and neural networks. Each technique provides unique advantages regarding interpretability, computational efficiency, and predictive accuracy. However, real-world datasets in churn scenarios are typically characterized by severe class imbalance, where the number of churn events is significantly smaller than non-churn events. This imbalance makes traditional accuracy metrics inadequate and demands specialized sampling strategies (undersampling, oversampling, synthetic generation) and robust performance metrics such as the F1-score and ROC-AUC score.

In our work, we rigorously explore multiple state-of-the-art sampling and scaling methodologies combined with leading classification algorithms to identify an optimal approach for churn prediction in telecommunications. Unlike prior works that often rely solely on accuracy, our evaluation emphasizes business-critical metrics such as F1-score and ROC-AUC, providing balanced insights into the practical utility of models. Furthermore, we perform comprehensive hyperparameter optimization and propose future experiments involving model explainability techniques (e.g., SHAP analysis),

feature engineering enhancements, and cross-industry validation to increase the robustness and generalizability of our findings.

2 Related Work

J. Burez et. al. examined the effects of undersampling and oversampling when training classification models to account for scenarios where the predicted event is rare. If the true proportion of customers who churned was 5% in the original data set, under-sampling was used to artificially inflate the proportion of churned customers to up to 50%. They determined that this under-sampling technique increased the AUC for a logistic regression model but had no significant impact when using Random Forests.

O. Adwan et. al. used a Multilayer Perceptron Neural Network (MLP) to predict whether a customer was at risk of churn. They used change-on-error (CoE) analysis to identify which of the features had a large impact on the probability of a customer churning, and looked at the weights of each layer of the neural network to identify which features had the most impact.

Florian and Damien used recurrent neural networks to predict customer churn in the online gambling industry. This method uses a time series element to predict churn, and specifically predicted whether or not a player would churn within the next 30 days of a current time t . Predicting churn in this way is needed in the online gambling industry because there is not a simple way of identifying whether or not a customer has churned, compared to telecommunications where a customer has to manually opt out or cancel their services. Using an RNN could be useful in predicting churn in the telecommunications industry but seems unnecessarily complex for this task.

3 Methods

3.1 Data

Table 1: Features and Descriptions

Feature	Description
Gender	Gender: The customer's gender: Male, Female
Age	Age: The customer's current age, in years, at the time the fiscal quarter ended.
Under 30	Under 30: Indicates if the customer is under 30 years old: Yes, No
Senior Citizen	Senior Citizen: Indicates if the customer is 65 or older: Yes, No
Partner	Married: Indicates if the customer is married: Yes, No
Dependents	Dependents: Indicates if the customer lives with any dependents: Yes, No. Dependents could be children, parents, grandparents, etc.
Number of Dependents	Number of Dependents: Indicates the number of dependents that live with the customer.
Referred a Friend	Referred a Friend: Indicates if the customer has ever referred a friend or family member to this company: Yes, No
Number of Referrals	Number of Referrals: Indicates the number of referrals to date that the customer has made.
Tenure	Tenure in Months: Indicates the total amount of months that the customer has been with the company by the end of the quarter specified above.
Phone Service	Phone Service: Indicates if the customer subscribes to home phone service with the company: Yes, No
Avg Monthly Long Distance Charges	Avg Monthly Long Distance Charges: Indicates the customer's average long distance charges, calculated to the end of the quarter specified above.
Multiple Lines	Multiple Lines: Indicates if the customer subscribes to multiple telephone lines with the company: Yes, No
Internet Service	Internet Service: Indicates if the customer subscribes to home internet service with the company: Yes, No
Internet Type	Internet Type: Indicates what type of Internet service the customer subscribes to with the company: No, DSL, Fiber Optic, Cable.
Avg Monthly GB Download	Avg Monthly GB Download: Indicates the customer's average download volume in gigabytes, calculated to the end of the quarter specified above.
Online Security	Online Security: Indicates if the customer subscribes to an additional online security service provided by the company: Yes, No
Online Backup	Online Backup: Indicates if the customer subscribes to an additional online backup service provided by the company: Yes, No
Device Protection Plan	Device Protection Plan: Indicates if the customer subscribes to an additional device protection plan for their Internet equipment provided by the company: Yes, No
Premium Tech Support	Premium Tech Support: Indicates if the customer subscribes to an additional technical support plan from the company with reduced wait times: Yes, No
Streaming TV	Streaming TV: Indicates if the customer uses their Internet service to stream television programming from a third party provider: Yes, No. The company does not charge an additional fee for this service.
Streaming Movies	Streaming Movies: Indicates if the customer uses their Internet service to stream movies from a third party provider: Yes, No. The company does not charge an additional fee for this service.
Streaming Music	Streaming Music: Indicates if the customer uses their Internet service to stream music from a third party provider: Yes, No. The company does not charge an additional fee for this service.
Unlimited Data	Unlimited Data: Indicates if the customer has paid an additional monthly fee to have unlimited data downloads/uploads: Yes, No
Contract	Contract: Indicates the customer's current contract type: Month-to-Month, One Year, Two Year.
Paperless Billing	Paperless Billing: Indicates if the customer has chosen paperless billing: Yes, No
Payment Method	Payment Method: Indicates how the customer pays their bill: Bank Withdrawal, Credit Card, Mailed Check.
Monthly Charges	Monthly Charge: Indicates the customer's current total monthly charge for all their services from the company.
Total Charges	Total Charges: Indicates the customer's total charges, calculated to the end of the quarter specified above.
Total Refunds	Total Refunds: Indicates the customer's total refunds, calculated to the end of the quarter specified above.
Total Extra Data Charges	Total Extra Data Charges: Indicates the customer's total charges for extra data downloads above those specified in their plan, by the end of the quarter specified above.
Total Long Distance Charges	Total Long Distance Charges: Indicates the customer's total charges for long distance above those specified in their plan, by the end of the quarter specified above.
Total Revenue	Total Revenue: Indicates the customer's total revenue generated, calculated to the end of the quarter specified above.
Churn	Churn Label: Yes = the customer left the company this quarter. No = the customer remained with the company. Directly related to Churn Value.

The data used in to train the below models is the 'Telco' dataset published by IBM. This data set consists of simulated data in a fictional telecommunication company. The full list of features and their descriptions are in Table 1. All categorical variables were one-hot encoded and numerical variables will be scaled using various scaling techniques detailed below. The full dataset has a total of 7032 rows, of which 1869 (27%) are customers who have churned. The distributions of the numerical and categorical features are listed in Figures 1 and 2.



Figure 1: Distribution of Numerical Features

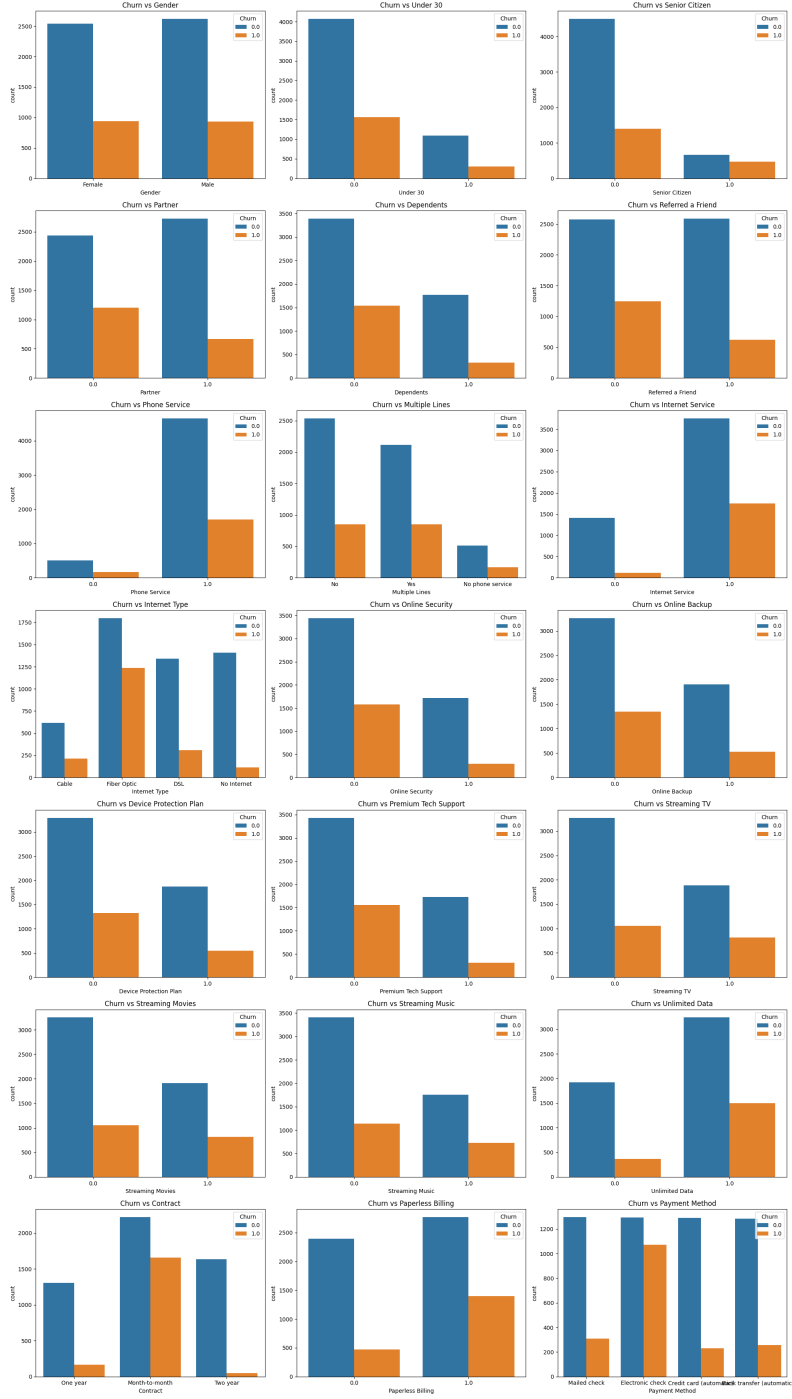


Figure 2: Distribution of Categorical Features

To identify the optimal method of training a classification model to predict customer churn, multiple sampling methods and scaling methods will be employed. For each model we train, we will build train datasets using the following steps:

1. Split the original dataframe into training dataset and testing dataset using a 5-fold randomly shuffled cross validation
2. Determine scaling technique to be used and scale continuous variables
3. Determine sampling method to be used and resample from training dataset

4. Fit model using the training set and score based on test dataset

The combination of model, scaling method, and sampling method will be reported on by taking the average of the evaluation metrics across all 5 folds of the cross validation splits.

3.2 Sampling Methods

J. Burez et. al described the potential need to use different sampling methods when predicting customer churn, due to the relative rarity of a customer disconnecting services. In this report, we will train each model using four sampling techniques: no sampling at all, random under sampling, random over-sampling, and Synthetic Minority Over-sampling Technique (SMOTE). We are including a technique that does not use any sampling method at all to establish a baseline of model performance. Random under-sampling and random over-sampling both artificially increase the proportion of churners in the data, but they are not the same. Under-sampling removes instances from the majority class until the desired proportion is met, and oversampling will duplicate instances of the minority class to increase the proportion of churners to non-churners. Lastly, we will be using SMOTE. SMOTE is another oversampling technique, but differs from the oversampling technique above in that it generates synthetic examples to balance out the minority class. SMOTE generates new synthetic instances by taking one real instance x_i of the minority class, taking several nearest neighbors of the same class, and performing random interpolation to obtain new instances.

3.3 Scaling Methods

For scaling, we will use 4 approaches: we will train the models without any scaling of variables, using ScikitLearn’s StandardScaler, ScikitLearn’s MinMax Scaler and ScikitLearn’s RobustScaler. The only columns that we are scaling are the columns that contain continuous numerical values. Binary indicator feature values and encoded values will not be scaled.

Scikit Learn’s StandardScaler scales each observation by subtracting the mean of that feature and dividing by the standard deviation. RobustScaler scales each observation by subtracting the median from each observation and dividing it by the interquartile range, which is the 75th percentile minus the 25th percentile. MinMax scaler scales each observation using the below formula:

$$x_{scaled} = \frac{x_i - x_{min}}{x_{max} - x_{min}} \quad (1)$$

3.4 Models Used

We are comparing six different models and comparing their ability to predict churn. The models we will use are the following.

1. Decision Tree Classifier
2. Logistic Regression
3. Gaussian Naive Bayes Classifier
4. Random Forest Classifier
5. Support Vector Machine
6. XGBoost Classifier

3.5 Validation Metrics

There are many metrics that can be used to validate a classification model, but we will primarily be interested in the F_1 score and area under the receiver operator characteristic curve (AUC-ROC). F_1 score is comprised of two other performance metrics, precision and recall.

$$Precision = \frac{TP}{TP + FP} \quad (2)$$

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

Where TP denotes 'true positive' (customers who had churned and the model predicted churn), FP denotes 'false positive' (customers who did not churn but model predicted churn), and FN denotes 'false negative' (customers who churned but model predicted not churn). F_1 score is defined as the harmonic mean of precision and recall:

$$F_1 \text{ score} = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}} \quad (4)$$

F_1 score is one of the best ways we can evaluate these models' ability to predict churn, because this value will be low if either precision or recall is low. In customer churn, incorrectly classifying a customer as at risk of churn results in lowering prices for customers who were never at risk of leaving. Incorrect labeling of a customer as not being at risk of churn likely results in losing that customer. Both of these scenarios are costly to the business, and F_1 score accounts for both scenarios.

We will also use AUC-ROC, which is the area under the curve plotted with the true positive rate on the vertical axis and the true negative rate on the horizontal axis. Simply put, the closer this value is to 1, the better the performance of the model.

3.6 Hyper Parameter Tuning

After identifying the optimal combinations of scaling and sampling techniques for each classification model based on average F_1 scores (as shown in Table 2), we performed hyperparameter tuning to further enhance model performance. We used GridSearchCV, a technique that uses all possible combinations of specified hyper-parameters using cross validation to find the best combination of hyper-parameters. We predefined values for some of the more crucial hyper-parameters for each model, and only defined a select few values for each parameter to save on computational resources. The GridSearchCV was completed using 5-fold cross validation and the models were evaluated across multiple metrics including accuracy, precision, recall, F_1 score, and AUC scores. However, the F_1 score was the primary metric used to select the optimal hyper-parameters for each classification model. The following parameters are those for which we optimized for each classification model type.

Logistic Regression:

- Regularization strength (C): [0.1, 1.0, 10] - The smaller the value the stronger the regularization (Default: 1.0)
- Penalty: [L2] - Type of regularization applied to the model (Default: L2). L1 is excluded due to lbfgs incompatibility
- Solver: [liblinear, lbfgs] - Algorithm used to optimize the parameters (Default: lbfgs)
- Class weight: [None, balanced] - Balanced automatically adjust the weights inversely proportional to class frequencies (Default: None)

Support Vector Classifier:

- Regularization strength (C): [0.1, 1.0, 10] - The smaller the value the stronger the regularization (Default: 1.0)
- Kernel: [Linear, RBF, Sigmoid] - Kernel type for used in algorithm (Default: RBF)
- Gamma: [scale, auto, 0.1, 1.0] - Kernel coefficient (Default: Scale)
- Class weight: [None, balanced] - Balanced automatically adjust the weights inversely proportional to class frequencies (Default: None)

Decision Tree Classifier:

- Criterion: [gini, entropy, log_loss] - The function used to measure the quality of a split (Default: gini)
- Max depth: [None, 5, 10, 15] - Maximum depth of the tree (Default: None)
- Min samples split: [2, 10, 20] - The minimum number of samples required to split an internal node (Default: 2)

- Class weight: [None, balanced] - Balanced automatically adjust the weights inversely proportional to class frequencies (Default: None)

Random Forest Classifier:

- N estimators: [100, 200, 300] - Number of trees in the forest (Default: 100)
- Criterion: [gini, entropy, log_loss] - The function used to measure the quality of a split (Default: gini)
- Max depth: [None, 5, 10, 15] - Maximum depth of the tree (Default: None)
- Min samples split: [2, 10, 20] - The minimum number of samples required to split an internal node (Default: 2)
- Max features: [None, sqrt, log2] - The number of features to consider when looking for the best split (Default Sqrt)
- Class weight: [None, balanced] – class balancing mechanism

Gaussian Naive Bayes Classifier:

- Var smoothing: [1×10^{-9} , 1×10^{-7} , 1×10^{-5}] - Artificially adding a value to the variance of each feature, widening the distribution and accounting for more samples further from the mean (Default: 1×10^{-9})

XGBoost Classifier:

- N estimators: [100, 200, 300] - Number of trees in the forest (Default: 100)
- Learning rate: [0.1, 0.2, 0.3] - Step size shrinkage (Default: 0.3)
- Max depth: [3, 6] - Maximum depth of the trees (Default: 6)
- Subsample: [0.8, 1.0] - The fraction of samples to be randomly sampled to build each tree (Default: 1.0)
- Colsample bytree: [0.8, 1.0] - The fraction of features to be randomly sampled for building each tree (Default: 1.0)

3.7 SHAP Analysis

In order to effectively communicate model results to stakeholders, Shapley additive explanations (SHAP) will be employed. SHAP analysis provides a method for measuring each feature’s contribution to the predicted value for a given sample. cite shap. SHAP Analysis is rooted in game theory, and builds upon concepts used to determine a fair payout in a collaborative game where individual contributions are not equal. While we will not dive deep into the mathematical foundations of SHAP Analysis, the below formula illustrates the concept at a high level.

$$f_i = \phi_0 + \sum_{features} \phi_{i,j} \quad (5)$$

Where f_i indicates the prediction of the model for sample i , ϕ_0 indicates the average prediction of the model, and $\phi_{i,j}$ indicates the SHAP values for each sample i . Each individual SHAP value represents the impact that feature had on the model’s prediction, relative to the average prediction of the model.

4 Results

While attempting to find the optimal combination of scaling and sampling techniques for each type of classification model, 96 models were trained. Tables 2 and 3 contain the combination of scaling and sampling techniques that lead to the highest F_1 and AUC scores, respectively. Interesting to note, the RobustScaler was not among any of the combinations of high performing models, 2 of the 12 listed in tables 2 and 3 scored better with no scaling, and another 2 of the 12 listed in tables 2 and 3 scored better with no sampling techniques applied. Most of the models performed best when using the StandardScaler and random oversampling techniques applied. XGBoost performed the best for

maximizing both the F_1 and AUC scores. The combination of no scaling and random oversampling were optimal for F_1 score, and the combination of using the StandardScaler and SMOTE were optimal for AUC scores. Other models like Logistic Regression and Random Forest Classifiers Support Vector Classifiers performed comparatively to the F_1 scores of XGBoost having strong recall, but poor precision causes their F_1 scores to drop.

There was more variance in sampling techniques. When maximizing F_1 score, the Decision Tree and Random Forest classifiers performed best with undersampling, the Support Vector Machine, XGBoost and Logistic Regression performed better with oversampling, and only the Naive Bayes classifier did best with SMOTE.

Figure 3 contains the AUC curves for all combinations of the models trained for each type of classification model. In the decision tree classifier, logistic regression, and support vector machine ROC curves, there is a significant variance between the ROC curves depending on the combination of scaling technique and sampling method. For naive bayes, random forest, and XGBoost, there does not appear to be much change to the ROC curves when the sampling and scaling methods change.

We then took the 6 combinations of scaling and sampling techniques and utilized GridSearchCV to perform hyper parameter tuning, aiming to see if we could further enhance model performance. The metrics shown in table 4 show a substantial improvement in model performance across all models. Most notably, the Support Vector Classifier shown the most significant improvement in performance with the F_1 score increasing by .258 from 0.674 to 0.931, resulting in the highest F_1 score of the 6 models. XGBoost also improved similarly and held the highest AUC score. All models improved with hyper parameter tuning. These results show that initial scaling and sampling technique combination choices are critical, however, hyperparameter optimization is crucial for unlocking the potential performance of each model.

Table 2: Combination of Scaling and Sampling that maximizes F_1 Score

classification_model	scaling_technique	sampling_technique	accuracy	precision	recall	f1_score	log_loss	roc_auc
xgboost_classifier	no_scaling	random_oversampling	0.817974	0.642409	0.711859	0.675115	0.406087	0.886566
random_forest_classifier	standard_scaler	random_undersampling	0.787257	0.568583	0.826544	0.673595	0.458929	0.875245
support_vector_classifier	standard_scaler	random_oversampling	0.787684	0.56976	0.824689	0.673592	0.422451	0.88003
logistic_regression	standard_scaler	random_oversampling	0.780573	0.558455	0.834342	0.668893	0.431611	0.883752
gaussian_naive_bayes_classifier	standard_scaler	smote	0.769764	0.547335	0.780207	0.642949	2.185817	0.849863
decision_tree_classifier	standard_scaler	random_undersampling	0.730519	0.49491	0.725554	0.588345	9.713094	0.728828

Table 3: Combination of Scaling and Sampling that maximizes AUC

classification_model	scaling_technique	sampling_technique	accuracy	precision	recall	f1_score	log_loss	roc_auc
xgboost_classifier	standard_scaler	smote	0.829492	0.68649	0.661427	0.673444	0.38513	0.890087
support_vector_classifier	standard_scaler	no_sampling	0.827218	0.690538	0.63485	0.661133	0.366044	0.885327
logistic_regression	standard_scaler	random_oversampling	0.787684	0.56976	0.824689	0.673592	0.422451	0.88003
random_forest_classifier	standard_scaler	smote	0.822099	0.670228	0.650816	0.659916	0.416036	0.879037
gaussian_naive_bayes_classifier	no_scaling	no_sampling	0.766352	0.541913	0.786551	0.641472	1.763892	0.85059
decision_tree_classifier	standard_scaler	random_undersampling	0.730519	0.49491	0.725554	0.588345	9.713094	0.728828

Table 4: Optimized hyper parameters selected based on the best average F_1 Score

classification_model	scaling_technique	sampling_technique	best_hyperparameters	accuracy	precision	recall	f1_score	roc_auc
support_vector_classifier	standard_scaler	random_oversampling	{'C': 10, 'class_weight': 'None', 'gamma': 1.0, 'kernel': 'rbf'}	0.933082	0.960004	0.903858	0.931037	0.967825
xgboost_classifier	no_scaling	random_oversampling	{'colsample_bytree': 0.8, 'learning_rate': 0.2, 'max_depth': 6, 'n_estimators': 300, 'subsample': 0.8}	0.912841	0.874553	0.96793	0.918355	0.971081
random_forest_classifier	standard_scaler	random_undersampling	{'class_weight': 'None', 'criterion': 'entropy', 'max_depth': 10, 'max_features': 'None', 'min_samples_split': 2, 'n_estimators': 300}	0.969596	0.967085	0.854196	0.807968	0.878474
logistic_regression	standard_scaler	random_oversampling	{'C': 1.0, 'class_weight': 'balanced', 'penalty': 'l2', 'solver': 'lbfgs'}	0.99632	0.781263	0.832183	0.805904	0.887956
gaussian_naive_bayes_classifier	standard_scaler	smote	{'var_smoothing': 1e-09}	0.787236	0.775824	0.808266	0.791667	0.8697
decision_tree_classifier	standard_scaler	random_undersampling	{'class_weight': 'None', 'criterion': 'gini', 'max_depth': 5, 'min_samples_split': 2}	0.776082	0.744953	0.842526	0.789896	0.848311

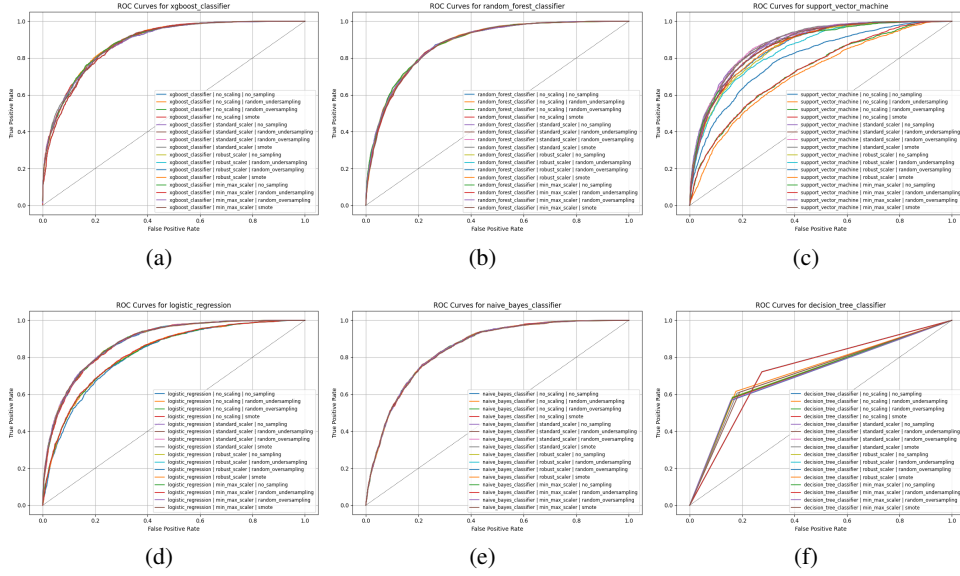


Figure 3: AUC curves for each type of classification model

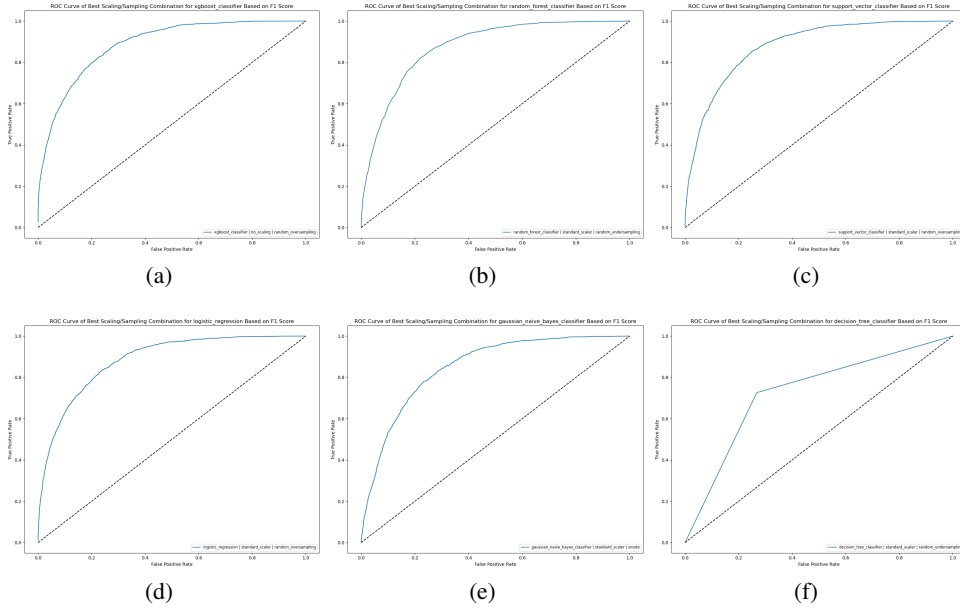


Figure 4: AUC curves for the optimal scaling/sampling combinations for each type of classification model

5 Applications beyond initial dataset

Now that we have identified the optimal scaling and sampling technique for each model and identified the hyperparameters that led to the highest F1 scores using the original dataset, we want to apply those techniques and hyperparameters to different datasets to test how applicable these methods are. Before moving to a completely new dataset, we performed feature engineering on the original dataset before training the models in an attempt to further increase the F1 scores and AUC. The feature engineering consisted of the following:

- Binning age, monthly charges, customer tenure, and average monthly GB download into 10 distinct buckets each
- Created a "Service Penetration Rate" column defined as the total number of services a customer is subscribed to divided by the number of services available
- A binary indicator to identify if a customer has any family
- A "Family Size" column indicating combining the partner indicator and the sum of all dependents
- A "Referral Rate" defined as the total number of referrals for a customer divided by the customer's tenure
- A "Refund Rate" defined as the total refund amount divided by the total charges
- A "Monthly Cost per GB" column defined as the monthly charges divided by the average monthly download in GB
- An "Extra Charges" ratio defined as the ratio of total extra data charges over the total charges

We then dropped the original features used for binning, one-hot encoded the categorical features, and trained the optimal model combinations on this new dataset. With respect to F_1 score the results were similar to the results while using the original dataset, except for the Support Vector Classifier. The recall in the Support Vector Classifier dropped below 0.02, which severely impacted the F_1 score. The full results of this model are listed in Table 5.

Table 5: Model Performance using Engineered/Derived Features

Classification Model	Scaling Technique	Sampling Technique	Accuracy	Precision	Recall	F1 Score	Log Loss	ROC AUC
Random Forest Classifier	Standard Scaler	Random Undersampling	0.777302	0.552262	0.860853	0.672624	0.434098	0.884096
Logistic Regression	Standard Scaler	Random Oversampling	0.777588	0.555098	0.821599	0.662332	0.433754	0.881479
XGBoost Classifier	No Scaling	Random Oversampling	0.817972	0.655546	0.667395	0.661047	0.453056	0.882279
Decision Tree Classifier	Standard Scaler	Random Undersampling	0.747298	0.516246	0.846182	0.640395	0.610188	0.839351
Gaussian Naive Bayes Classifier	Standard Scaler	SMOTE	0.778297	0.571152	0.687393	0.622214	2.851317	0.838566
Support Vector Classifier	Standard Scaler	Random Oversampling	0.733077	0.459865	0.019299	0.036977	0.753589	0.764371

We identified a new dataset also from the telecommunications industry containing 64374 records with 8 features and indicator of whether the customer churned or not. Three of these features were categorical and one hot encoded.

We trained and tested each model with their respective combination of scaling, sampling, and hyperparameters. Table 6 contains the average the model performance metrics across a 5-fold cross validation split.

Table 6: Model Performance using second Telecommunications Dataset

Classification Model	Scaling Technique	Sampling Technique	Accuracy	Precision	Recall	F1 Score	Log Loss	ROC AUC
XGBoost Classifier	No Scaling	Random Oversampling	0.999922	0.999869	0.999968	0.999918	0.000582	0.999992
Random Forest Classifier	Standard Scaler	Random Undersampling	0.998742	0.998261	0.999082	0.998671	0.004244	0.999978
Decision Tree Classifier	Standard Scaler	Random Undersampling	0.956613	0.930023	0.982346	0.955465	0.088584	0.99441
Support Vector Classifier	Standard Scaler	Random Oversampling	0.942228	0.934457	0.944284	0.939341	0.143388	0.987544
Gaussian Naive Bayes Classifier	Standard Scaler	SMOTE	0.835151	0.800956	0.867575	0.832931	0.403325	0.908518
Logistic Regression	Standard Scaler	Random Oversampling	0.825504	0.800325	0.84158	0.820427	0.393823	0.904097

Similar to the original dataset, the Support Vector Classifier and XGBoost model had the highest F1 score and AUC. All models improved significantly as compared to the original dataset, with all models except the Naive Bayes Classifier and Logistic Regression having an F1 score of over 0.939. This appears to suggest that these model parameters and techniques could be implemented to predict churn in any industry with decent success.

We finally applied the same methodologies to a third dataset, this time in the banking industry. The models performed significantly worse when trained and tested in this dataset, with the highest F1 score slightly below 0.6. XGBoost was the best performing model in terms of F1 score, which was the case with both the original dataset and the second telecommunications dataset. The significantly lower F1 scores across the board could indicate that these combinations are not necessarily applicable to datasets or industries beyond the initial dataset.

Table 7: Model Performance using Banking Dataset

Classification Model	Scaling Technique	Sampling Technique	Accuracy	Precision	Recall	F1 Score	Log Loss	ROC AUC
XG Boost Classifier	No Scaling	Random Oversampling	0.8384	0.607759	0.584399	0.595701	0.422642	0.836589
Random Forest Classifier	Standard Scaler	Random Undersampling	0.7810	0.476573	0.756371	0.584504	0.474847	0.856955
Decision Tree Classifier	Standard Scaler	Random Undersampling	0.7449	0.430041	0.761398	0.548784	0.618796	0.830766
Gaussian Naive Bayes Classifier	Standard Scaler	SMOTE	0.7361	0.408822	0.661405	0.505227	0.549287	0.781643
Logistic Regression	Standard Scaler	Random Undersampling	0.7117	0.383714	0.687141	0.492409	0.577906	0.768511
Support Vector Classifier	Standard Scaler	Random Undersampling	0.7991	0.509257	0.361364	0.422674	0.698014	0.764247

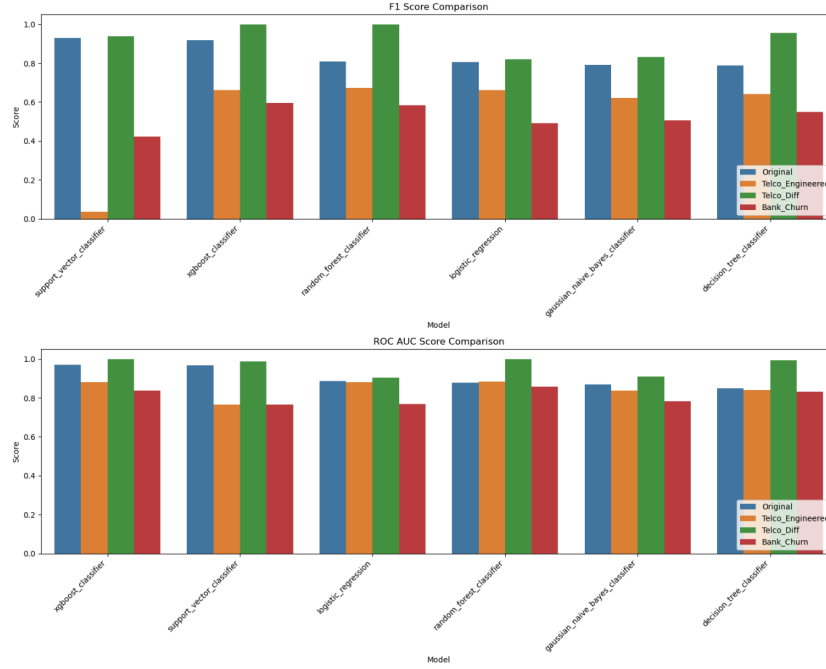


Figure 5: F1 score and ROC AUC score Bar plots for each classification model on each dataset

6 Conclusion

Some of the key takeaways are that XGBoost, Random Forest, and Naive Bayes classifiers all had little variance between combinations of sampling and scaling techniques. Support Vector Classifier, Logistic Regression, and the Decision Tree Classifier all had a significantly higher variance in the ROC Curves given different combinations.

It appears that the model, sampling, and scaling combinations identified in this report may be applicable to any instance where predicting customer churn is of benefit. The model seemed to maintain high AUC values and F1 scores beyond the original dataset where optimal scaling/sampling strategies and hyperparameters were identified. This could be used as a framework or starting point for any industry where minimizing churn is important.

To continue in this research, it would be critical to identify other verified, high quality datasets to apply the combinations of scaling and sampling techniques with the optimized hyperparameters. With more applications of these combinations, we would be able to better assess whether this methodology is applicable to churn prediction as a whole. It would also be interesting to assess the variances between each combination for each of the chosen models to see if the pattern of low variance for XGBoost, Random Forest and Naive Bayes and high variance for Support Vector Classifier, Logistic Regression, and Decision Trees follows. Regardless of the performance of these specific model combinations for predicting churn in future datasets or industries, this iterative process serves as a good framework in identifying the ideal data pre-processing techniques and hyperparameters that will apply to any dataset.

ALL BELOW IS NOT PART OF FINAL PAPER, JUST DIDNT WANT TO DELETE PREVIOUS "FUTURE PLANS" FROM MIDWAY REPORT

7 Future Plans

In the following weeks, we plan on continuing to further explore the current methods as well as look into other promising directions for future exploration. We recently were able to perform hyperparameter tuning on all the best sampling and scaling combinations for each type of classification model based on F_1 score. We plan on taking those optimized models and performing feature importance analysis using techniques like SHAP (shapely additive explanations) (Adam, ETA: 07/13/2025). This will help us understand which features influence the outcome of a customer churning the most. This is a crucial step that will allow business stakeholders to gain actionable insights from our data.

Another future exploration that we plan to investigate is the development of a neural network tailored for churn prediction. We plan on utilizing different structures, activation functions, loss functions, and optimization algorithms to evaluate model performance and potentially capture more complex patterns in our data (Adam, ETA: 07/13/2025).

We also plan on testing our models on a different churn dataset possibly from a different industry (one with similar class imbalances) and comparing results (Anthony, ETA: 07/13/2025). This will help us understand whether certain classification models and combinations of scaling and sampling techniques can generalize well and be transferred across domains and different datasets. Another avenue of testing on different data that we may also pursue would be to test on our current dataset but with new engineered features and compare model performance that way (Nitin, ETA: 07/13/2025).

We are also open to feedback and guidance to see what direction we should be going moving forward with current methods or new methods. As we work on finalizing and completing the remainder of the project, we will need to complete a few deliverables. The 8 page paper NeurIPS paper (Team, ETA: 07/20/2025), a slide deck (PowerPoint) with our results, completed Github repository (Adam, ETA: 07/25/2025), and rehearsal/practice for the spotlight presentation will need to be completed before the final class meeting on 07/29/2025 (Team, ETA: 07/26-27/2025).

Everything below this is from the NeurIPS Template and should be disregarded

References

References follow the acknowledgments in the camera-ready paper. Use unnumbered first-level heading for the references. Any choice of citation style is acceptable as long as you are consistent. It is permissible to reduce the font size to small (9 point) when listing the references. Note that the Reference section does not count towards the page limit.

[1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In G. Tesauro, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp. 609–616. Cambridge, MA: MIT Press.

[2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the GEneral NEural Simulation System*. New York: TELOS/Springer-Verlag.

[3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249-5262.

A Appendix / supplemental material

Optionally include supplemental material (complete proofs, additional experiments and plots) in appendix. All such materials **SHOULD be included in the main submission**.

NeurIPS Paper Checklist

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found.

IMPORTANT, please:

- **Delete this instruction block, but keep the section heading “NeurIPS paper checklist”,**
- **Keep the checklist subsection headings, questions/answers and guidelines below.**
- **Do not modify the questions and only use the provided macros for your answers.**

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope?

Answer: [TODO]

Justification: [TODO]

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [TODO]

Justification: [TODO]

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that the paper does not include experiments.

- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).

- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.

- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that the paper poses no such risks.

- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.