

From Prediction to Action: A Comprehensive Technical Report on Telco Customer Churn Analysis

Part 1: Establishing the Foundation: Data and Discovery

This report presents a comprehensive, end-to-end analysis of the Telco Customer Churn dataset. The objective extends beyond simply building a predictive model; it is to create a strategic asset that translates complex data patterns into actionable business intelligence. The analysis will proceed through four distinct phases: establishing a foundational understanding of the business problem and the data, constructing a suite of predictive models ranging from interpretable baselines to high-performance ensembles, conducting a nuanced evaluation of these models based on business-critical metrics, and finally, synthesizing the analytical findings into a concrete set of strategic recommendations designed to reduce customer attrition.

Section 1.1: The Business Imperative: Understanding the Churn Problem

Core Business Problem

In the telecommunications industry, customer churn—the rate at which customers discontinue their service—is a primary determinant of profitability and long-term viability. The market is characterized by intense competition, where customers can choose from multiple service

providers and switch with relative ease.¹ This environment leads to an average annual churn rate of 15-25%, a figure that represents a significant and continuous drain on revenue.¹

The financial imperative to manage churn is underscored by a fundamental business principle: customer retention is vastly more cost-effective than customer acquisition. Acquiring a new customer is estimated to be anywhere from five to 25 times more expensive than retaining an existing one.² Consequently, a high churn rate not only signifies a direct loss of recurring revenue but also inflates marketing and sales expenditures required to replace the lost customers.³ For many telecommunications firms, reducing the attrition of high-value customers is the single most important business goal.¹

This project addresses the churn problem by leveraging machine learning to predict which customers are at the highest risk of leaving. By identifying these individuals proactively, a company can move from a reactive "win-back" strategy to a proactive retention strategy, deploying targeted interventions such as loyalty discounts, service upgrades, or personalized support to prevent churn before it occurs.⁵ The ultimate goal is to build a data-driven tool that provides the business with the foresight needed to make informed decisions, reduce attrition rates, and secure its revenue base.⁵

Dataset Provenance and Context

The dataset used for this analysis is the "WA_Fn-UseC_-Telco-Customer-Churn.csv" file, a well-documented and widely used dataset originating from IBM's sample data collections.⁷ It represents a fictional telecommunications company in the United States and provides a snapshot of its customer base.

The dataset contains 7,043 rows, where each row corresponds to a unique customer, and 21 columns, representing various customer attributes.⁷ These attributes are broadly categorized into:

1. **Demographic Information:** Customer characteristics such as gender, age range (SeniorCitizen), and domestic situation (Partner, Dependents).
2. **Customer Account Information:** Details about the customer's relationship with the company, including tenure (how long they have been a customer), contract type, billing method, and charge amounts.
3. **Subscribed Services:** Information on which specific services each customer has signed up for, such as phone service, internet, online security, and streaming options.

The target variable for this analysis is the Churn column. This binary feature indicates whether the customer left the company's service within the last month, making it a clear and direct

measure of customer attrition.⁷

Data Dictionary

A precise understanding of each variable is a prerequisite for any meaningful analysis. The following table provides a comprehensive dictionary for the dataset, detailing each variable's name, data type, a description of its meaning, example values, and its hypothesized relevance to the business problem of predicting churn.

Table 1: Comprehensive Data Dictionary

Variable Name	Data Type	Description	Example Values	Business Relevance & Hypothesis
customerID	Categorical (chr)	A unique identifier for each customer.	7590-VHVEG	An identifier with no predictive power; should be excluded from modeling.
gender	Categorical (chr)	The customer's gender.	Female, Male	Demographic variable. Hypothesis: Low predictive power; churn rates are likely similar across genders.
SeniorCitizen	Numeric (int)	Indicates if the customer is a senior citizen (65+).	0 (No), 1 (Yes)	Demographic variable. Hypothesis: Senior citizens may have different needs or be more

				price-sensitive , potentially leading to a higher churn rate.
Partner	Categorical (chr)	Indicates if the customer has a partner.	Yes, No	Demographic variable. Hypothesis: Customers with partners may have more stable households and thus lower churn rates.
Dependents	Categorical (chr)	Indicates if the customer has dependents.	Yes, No	Demographic variable. Hypothesis: Similar to Partner, customers with dependents may be less likely to churn due to higher inertia.
tenure	Numeric (int)	Number of months the customer has been with the company.	1, 34, 72	Key behavioral variable. Hypothesis: A very strong predictor. Longer tenure indicates loyalty and satisfaction, leading to a much lower churn

				probability.
PhoneService	Categorical (chr)	Indicates if the customer has phone service.	Yes, No	Service variable. Hypothesis: Most customers will have this basic service; may not be a strong differentiator for churn.
MultipleLines	Categorical (chr)	Indicates if the customer has multiple phone lines.	Yes, No, No phone service	Service add-on. Hypothesis: A weak predictor; could indicate higher integration but also higher cost.
InternetService	Categorical (chr)	Type of internet service subscribed to.	DSL, Fiber optic, No	Key service variable. Hypothesis: Strong predictor. Fiber optic is a premium service; churn may be higher if service quality does not meet price expectations.
OnlineSecurity	Categorical (chr)	Indicates if the customer has the online	Yes, No, No internet	Service add-on. Hypothesis:

		security add-on.	service	Customers subscribing to more add-ons are more invested in the ecosystem and less likely to churn.
OnlineBackup	Categorical (chr)	Indicates if the customer has the online backup add-on.	Yes, No, No internet service	Service add-on. Hypothesis: Similar to OnlineSecurity, indicates higher customer investment and lower churn risk.
DeviceProtection	Categorical (chr)	Indicates if the customer has the device protection add-on.	Yes, No, No internet service	Service add-on. Hypothesis: Similar to OnlineSecurity, indicates higher customer investment and lower churn risk.
TechSupport	Categorical (chr)	Indicates if the customer has the premium tech support add-on.	Yes, No, No internet service	Service add-on. Hypothesis: A strong negative predictor of churn. Customers

				who opt for tech support are likely more reliant on the service and proactive about issue resolution.
StreamingTV	Categorical (chr)	Indicates if the customer has the streaming TV add-on.	Yes, No, No internet service	Service add-on. Hypothesis: Indicates deeper engagement with the company's services, likely correlating with lower churn.
StreamingMovies	Categorical (chr)	Indicates if the customer has the streaming movies add-on.	Yes, No, No internet service	Service add-on. Hypothesis: Similar to StreamingTV, indicates higher engagement and lower churn risk.
Contract	Categorical (chr)	The customer's contract term.	Month-to-month, One year, Two year	Key account variable. Hypothesis: The single strongest predictor. Month-to-month contracts

				offer maximum flexibility to leave, and these customers are expected to have the highest churn rate.
PaperlessBilling	Categorical (chr)	Indicates if the customer uses paperless billing.	Yes, No	Account variable. Hypothesis: Moderate predictor. May correlate with tech-savviness or customer engagement preferences.
PaymentMethod	Categorical (chr)	The customer's payment method.	Electronic check, Mailed check, etc.	Account variable. Hypothesis: Automatic payment methods (credit card, bank transfer) may correlate with lower churn due to convenience and "set-it-and-forget-it" behavior.
MonthlyCharges	Numeric (dbl)	The customer's monthly bill	29.85, 70.70	Key financial variable. Hypothesis:

		amount.		Strong predictor. Higher monthly charges can increase churn risk, especially if perceived value is low.
TotalCharges	Numeric (dbl)	The total amount charged to the customer over their lifetime.	29.85, 1889.50	Financial variable. Highly correlated with tenure and MonthlyCharges. May be redundant but captures overall customer value.
Churn	Categorical (chr)	Indicates if the customer churned in the last month.	Yes, No	Target Variable. The outcome the model will be trained to predict.

Section 1.2: Rigorous Data Preparation and Cleaning

The quality of a predictive model is fundamentally limited by the quality of the data it is trained on. This section details the critical steps taken to clean, preprocess, and structure the raw data, transforming it into a reliable format suitable for sophisticated modeling.

Initial Loading and Inspection

The first step in any data analysis workflow is to load the data into the R environment and perform a preliminary inspection. The tidyverse suite of packages provides modern and efficient tools for this purpose.

R

```
# Load the tidyverse library for data manipulation and visualization  
library(tidyverse)
```

```
# Read the data from the CSV file into a tibble (a modern data frame)  
churn_data <- read_csv("data/WA_Fn-UseC_-Telco-Customer-Churn.csv")
```

With the data loaded, an initial structural summary is obtained using `glimpse()`. This function provides a transposed view of the data, showing each column's name, data type, and the first few values.

R

```
# Get a structural summary of the dataset  
glimpse(churn_data)
```

This initial check confirms the dataset's dimensions (7043 observations, 21 variables) and reveals a critical data quality issue that must be addressed immediately.

The TotalCharges Anomaly

The `glimpse()` output reveals that the `TotalCharges` column has been loaded as a character (`<chr>`) type, despite representing a monetary value that should be numeric (`<dbl>`). This is a classic real-world data problem that prevents any numerical analysis on this feature. The root cause is the presence of non-numeric values in the column. Further investigation reveals that customers with a tenure of 0 months (i.e., brand new customers) have a blank space (" ") for

their TotalCharges.² R's

read_csv function interprets the entire column as character data to accommodate these spaces.

A novice approach might involve simply deleting these rows. However, this would be a significant analytical error. These 11 rows represent the *entire population of new customers* in the dataset. Removing them would blind the model to the behavior of this critical and often volatile segment. The correct approach is a careful, multi-step cleaning process that preserves this valuable information.

The missingness in this case is not random; it is a direct result of a business process where a new customer has not yet been billed. This is known as "Missing Not At Random" (MNAR), and understanding this context is key to handling it correctly.¹³ The logical value for

TotalCharges for a customer with zero tenure is 0. The following code implements the proper fix.

R

```
# The TotalCharges column is character type because new customers have a blank space.
# We first coerce it to numeric, which turns the blank spaces into NA (missing) values.
# Then, we identify these NAs and replace them with 0, as a new customer has logically accrued $0 in
total charges.
churn_data <- churn_data %>%
  mutate(TotalCharges = as.numeric(TotalCharges)) %>%
  mutate(TotalCharges = ifelse(is.na(TotalCharges), 0, TotalCharges))
```

This sequence of operations correctly converts TotalCharges to a numeric type while accurately representing the financial status of new customers, ensuring they are retained for the analysis.

Consolidating Categorical Levels

Several categorical features related to services, such as OnlineSecurity, OnlineBackup, DeviceProtection, TechSupport, StreamingTV, and StreamingMovies, share a common structure with three levels: "Yes", "No", and "No internet service".⁷ From a modeling perspective, the distinction between "No" (the customer chose not to have the service) and

"No internet service" (the customer was ineligible for the service) is often irrelevant. In both cases, the customer does not have the specific service add-on. Consolidating these into a single "No" category simplifies the feature, reduces the dimensionality of the data, and can lead to a more robust and interpretable model.¹⁵

A similar consolidation is applied to the MultipleLines feature, where "No phone service" is functionally equivalent to "No" for the purpose of having multiple lines.

R

```
# Recode 'No internet service' to 'No' for relevant columns
churn_data <- churn_data %>%
  mutate(across(c(OnlineSecurity, OnlineBackup, DeviceProtection, TechSupport,
    StreamingTV, StreamingMovies),
    ~ recode(., "No internet service" = "No"))))
```

```
# Recode 'No phone service' to 'No' for MultipleLines
churn_data <- churn_data %>%
  mutate(MultipleLines = recode(MultipleLines, "No phone service" = "No"))
```

Dropping Non-Predictive Features

The customerID column serves as a unique identifier for each row. While essential for operational purposes like linking to other databases, it holds no predictive information about churn. Each value is unique, and a model attempting to use it would simply be memorizing individual customers, a severe form of overfitting. Therefore, it must be removed from the dataset before modeling.¹⁷

R

```
# The customerID is a unique identifier and has no predictive value.
# It should be removed before modeling.
churn_data <- churn_data %>%
  select(-customerID)
```

Final Data Type Conversion

As a final preparation step, it is best practice in R to convert all character columns intended for use as categorical predictors into the factor data type. This explicitly defines the levels of each category and is required by some modeling functions. It also ensures that visualizations and model outputs handle these variables correctly.

R

```
# Convert remaining character columns to factors for modeling
churn_data <- churn_data %>%
  mutate(across(where(is.character), as.factor))
```

After these rigorous preparation steps, the dataset is clean, correctly typed, and structurally sound, providing a high-quality foundation for the subsequent exploratory analysis and predictive modeling phases.

Section 1.3: Uncovering Customer Behavior Through Exploratory Data Analysis (EDA)

Exploratory Data Analysis (EDA) is the process of using summary statistics and data visualizations to understand the main characteristics of a dataset. It is a critical phase where analysts uncover patterns, identify anomalies, and form hypotheses about the relationships between variables. For this project, EDA will focus on identifying the key drivers of customer churn.

Churn Rate Baseline

The first step is to establish the overall churn rate, which serves as a baseline for the entire

analysis. This metric quantifies the magnitude of the business problem.

R

```
# Calculate and print the overall churn rate
churn_rate <- churn_data %>%
  count(Churn) %>%
  mutate(Proportion = n / sum(n))
print(churn_rate)
```

The analysis reveals an overall churn rate of approximately 26.5%.¹⁸ This is a substantial figure, confirming that customer attrition is a significant issue for this fictional company. This proportion also highlights a moderate class imbalance: there are roughly three non-churning customers for every one customer who churns. This imbalance must be considered during model evaluation, as metrics like accuracy can be misleading.

Visualizing Key Relationships with ggplot2

Visualizations are the most powerful tool in EDA for revealing complex relationships. The ggplot2 package, a core part of the tidyverse, provides a flexible and expressive grammar for creating sophisticated statistical graphics.

Churn by Contract Type

The relationship between contract type and churn is hypothesized to be one of the strongest predictors. A dodged bar chart is an effective way to visualize this relationship, comparing the counts of churned and non-churned customers across the three contract types.

R

```
# Create a bar chart showing churn counts for each contract type
```

```
ggplot(churn_data, aes(x = Contract, fill = Churn)) +
  geom_bar(position = "dodge") +
  labs(title = "Customer Churn by Contract Type",
       x = "Contract Type",
       y = "Number of Customers") +
  theme_minimal()
```

The resulting plot reveals a stark and compelling pattern. The number of churned customers (Churn = Yes) is overwhelmingly concentrated in the Month-to-month contract category. In contrast, customers on One year and Two year contracts exhibit significantly lower churn. This provides powerful evidence that a lack of long-term commitment is the single largest indicator of churn risk, confirming that this is a critical feature for the predictive model.²⁰

Churn by Internet Service and the Fiber Optic Paradox

Another key service feature is the type of internet connection. A similar visualization can explore its link to churn.

R

```
# Create a bar chart showing churn counts for each internet service type
ggplot(churn_data, aes(x = InternetService, fill = Churn)) +
  geom_bar(position = "dodge") +
  labs(title = "Customer Churn by Internet Service Type",
       x = "Internet Service",
       y = "Number of Customers") +
  theme_minimal()
```

This visualization uncovers a counter-intuitive relationship. While Fiber optic is a premium, higher-speed service compared to DSL, customers with fiber optic service churn at a noticeably higher rate.¹⁹ This finding, combined with the analysis of monthly charges, suggests a "value-expectation mismatch." Fiber optic service is more expensive, which sets higher customer expectations for performance and reliability. If these elevated expectations are not consistently met, the high price point becomes a source of dissatisfaction, leading to a greater propensity to churn. This is a crucial business insight: the problem may not be the technology itself, but the delivery of the premium experience that customers are paying for.

Churn by Monthly Charges

To investigate the impact of price on churn, a boxplot is used to compare the distribution of MonthlyCharges for customers who churned versus those who did not.

R

```
# Create box plots to compare MonthlyCharges for churned vs. non-churned customers
ggplot(churn_data, aes(x = Churn, y = MonthlyCharges, fill = Churn)) +
  geom_boxplot() +
  labs(title = "Monthly Charges by Churn Status",
       x = "Churn Status",
       y = "Monthly Charges ($)") +
  theme_minimal()
```

The boxplot clearly shows that the median MonthlyCharges for customers who churned is substantially higher than for those who remained loyal. The entire interquartile range (the "box" part of the plot) for the "Yes" category is shifted upwards compared to the "No" category.⁴ This reinforces the hypothesis that cost is a significant factor in the churn decision. Customers paying more per month are at a higher risk of leaving, especially if they perceive the value received does not justify the expense.

Churn by Tenure

Customer tenure, or the length of the customer relationship, is a powerful proxy for loyalty. A density plot, faceted by churn status, is an effective way to visualize how churn behavior changes over the customer lifecycle.

R

```
# Create density plots of tenure, faceted by Churn status
```



```
ggplot(churn_data, aes(x = tenure, fill = Churn)) +
  geom_density(alpha = 0.6) +
  facet_wrap(~Churn) +
  labs(title = "Distribution of Tenure by Churn Status",
       x = "Tenure (Months)",
       y = "Density") +
  theme_minimal()
```

The visualization demonstrates that churn is heavily skewed towards new customers. The density plot for Churn = Yes shows a large peak at very low tenure values (typically 1-10 months) and a rapid decline thereafter. Conversely, the plot for Churn = No shows a much more uniform distribution, with a significant concentration of customers at the highest tenure values (around 72 months). This indicates that if a customer can be retained through their first year, their likelihood of churning decreases dramatically.¹⁴

Correlations among Numeric Variables

Finally, it is important to understand the relationships between the numeric predictor variables to check for multicollinearity, which can be an issue for some models like logistic regression. A correlation matrix provides a concise summary of these relationships.

R

```
# Select numeric columns and compute the correlation matrix
numeric_data <- churn_data %>% select(tenure, MonthlyCharges, TotalCharges)
correlation_matrix <- cor(numeric_data)
print(correlation_matrix)
```

The analysis shows a strong positive correlation between tenure and TotalCharges ($r > 0.8$), which is expected, as customers who stay longer naturally accumulate higher total charges. There is also a moderate positive correlation between MonthlyCharges and TotalCharges ($r > 0.6$).¹⁷ The correlation between

tenure and MonthlyCharges is weaker. This information is valuable for model building; for linear models, including highly correlated predictors like tenure and TotalCharges simultaneously might be redundant and could destabilize coefficient estimates. Tree-based

models like Random Forest are less sensitive to this issue.

Part 2: Building the Predictive Engine: From Baseline to Advanced Models

With a clean dataset and a deep understanding of the underlying patterns, the project now transitions to the predictive modeling phase. This part details the construction of a robust modeling framework and the implementation of three distinct machine learning models: Logistic Regression, Random Forest, and XGBoost. The approach is to start with a simple, highly interpretable model and progressively move to more complex, high-performance algorithms, allowing for a comprehensive comparison.

Section 2.1: A Repeatable Framework for Modeling

A disciplined and repeatable framework is essential for reliable model development and evaluation. Two key principles underpin this framework: splitting the data to prevent overfitting and setting a random seed to ensure reproducibility.

Train/Test Split

The single most important principle in applied machine learning is to evaluate a model's performance on data it has not seen during training. Evaluating on the training data itself would only measure the model's ability to memorize, not its ability to generalize to new, unseen customers. To this end, the dataset is partitioned into a *training set* and a *testing set*. The model learns the patterns from the training set, and its predictive performance is then measured on the testing set.

A standard split is 80% for training and 20% for testing. For a classification problem with an imbalanced target variable like churn, a simple random split is insufficient. It is crucial to perform a *stratified* split. Stratification ensures that the proportion of the target classes (26.5% churners, 73.5% non-churners) is maintained in both the training and testing sets. This prevents a scenario where one set has a disproportionately high or low number of churners, which would bias the model's training and evaluation. The caret package provides the

createDataPartition function specifically for this purpose.¹⁵

Reproducibility with set.seed()

Machine learning processes, including data splitting and the training of certain models like Random Forest, involve elements of randomness. To ensure that the results of this analysis can be reproduced exactly by others, a random seed is set. The set.seed() function initializes R's random number generator to a specific state. By setting a seed (e.g., set.seed(42)), anyone who runs the same code will get the exact same "random" data split and model results, which is a cornerstone of scientific and professional rigor.

The following code implements this robust framework.

R

```
# Load the caret package for streamlined machine learning workflows
library(caret)

# Set a seed for reproducibility. Any random process after this will be the same every time.
set.seed(42)

# Create a stratified 80/20 split of the data.
# The 'createDataPartition' function ensures the proportion of 'Churn' is the same in both sets.
train_indices <- createDataPartition(churn_data$Churn, p = 0.8, list = FALSE)

# Create the training and testing datasets based on the indices
train_data <- churn_data[train_indices, ]
test_data <- churn_data[-train_indices, ]
```

This framework ensures that all subsequent models are built and evaluated under fair, consistent, and reproducible conditions.

Section 2.2: Model 1 - The Interpretable Baseline: Logistic Regression

The first model to be built is Logistic Regression. It serves as an essential baseline for several reasons: it is computationally efficient, statistically robust, and, most importantly, highly interpretable. Unlike more complex "black-box" models, the output of a logistic regression model directly reveals the relationship between each predictor variable and the probability of churn, providing valuable business insights.¹⁵

Implementation

Logistic Regression is a type of Generalized Linear Model (GLM) used for binary classification. It models the logarithm of the odds of the outcome as a linear combination of the predictor variables. In R, it is implemented using the `glm()` function with the `family` argument set to "binomial". The formula `Churn ~.` is a convenient shorthand that instructs the model to predict the Churn variable using all other columns in the `train_data` as predictors.

R

```
# Build a Logistic Regression model.  
# The formula 'Churn ~.' means "predict Churn using all other variables as predictors".  
# 'family = "binomial"' specifies that this is a logistic regression for a binary outcome.  
log_model <- glm(Churn ~., data = train_data, family = "binomial")
```

Interpreting the summary() Output

The true power of logistic regression lies in the interpretability of its summary output. This summary provides detailed statistics about the model's learned coefficients.

R

```
# Display a detailed summary of the logistic regression model  
summary(log_model)
```

The `summary()` output contains a table of coefficients. Each row corresponds to a predictor variable (or a specific level of a categorical variable), and the columns provide key information:

- **Estimate:** This is the coefficient (or log-odds) for each variable. A positive coefficient means that an increase in the variable's value increases the log-odds of churning. A negative coefficient means it decreases the log-odds of churning. For example, a large negative coefficient for tenure indicates that as tenure increases, the likelihood of churn decreases significantly.
- **Std. Error:** This measures the uncertainty in the estimate of the coefficient.
- **z value:** This is the Estimate divided by the Std. Error. It is a measure of how many standard errors the coefficient is away from zero.
- **Pr(>|z|):** This is the p-value, which indicates the statistical significance of each predictor. A small p-value (typically < 0.05) suggests that the variable has a statistically significant relationship with the outcome (churn), and its effect is unlikely to be due to random chance.¹⁶

By examining the p-values, one can identify the most influential drivers of churn according to the model. Variables like ContractMonth-to-month, tenure, and InternetServiceFiber optic will likely have very small p-values, confirming the findings from the EDA and establishing a strong, statistically-backed baseline for predictive performance.

Section 2.3: Model 2 - The Ensemble Powerhouse: Random Forest

The second model, Random Forest, represents a significant step up in complexity and predictive power. It is an *ensemble learning* method, meaning it combines the predictions of many individual models to produce a final, more accurate prediction. Specifically, a Random Forest builds hundreds or thousands of individual decision trees on different bootstrapped subsamples of the data. For each split in a tree, it only considers a random subset of predictor variables. This dual-randomization process creates a diverse "forest" of trees, and their collective prediction (by majority vote) is more robust and less prone to overfitting than any single decision tree.²⁶

Random Forests are particularly effective because they can automatically capture complex, non-linear relationships and interactions between variables without requiring manual feature engineering.

Implementation with caret

The caret package provides a standardized and powerful interface for training and tuning a wide variety of machine learning models, including Random Forest. The `train()` function is the central workhorse.

R

```
# Load the randomForest package, which the caret 'rf' method relies on
library(randomForest)

# Build a Random Forest model using the caret package.
# 'method = "rf"' specifies the Random Forest algorithm.
# 'trControl' is used to set up cross-validation for hyperparameter tuning.
# This process can be computationally intensive and may take a few minutes.
rf_model <- train(Churn ~.,
  data = train_data,
  method = "rf",
  trControl = trainControl(method = "cv", number = 5),
  prox = FALSE,
  allowParallel = TRUE)
```

In this code:

- `method = "rf"` tells caret to use the randomForest package to build the model.
- `trControl = trainControl(method = "cv", number = 5)` sets up the resampling strategy. Here, it specifies 5-fold cross-validation. The training data is split into 5 "folds." The model is trained on 4 folds and validated on the 5th, and this process is repeated 5 times, with each fold serving as the validation set once. This robust process is used for hyperparameter tuning.²⁶

Hyperparameter Tuning

Random Forest models have several hyperparameters that control their structure and training process. The most important one is `mtry`, which is the number of variables randomly sampled as candidates at each split in the tree.²⁶ The optimal value for

`mtry` depends on the dataset. The caret `train()` function automatically tunes this

hyperparameter. It trains several Random Forest models with different values of `mtry` (e.g., 2, 15, 30) using the 5-fold cross-validation process defined in `trControl`. It then selects the `mtry` value that resulted in the best average performance across the 5 folds. This automated tuning process helps to maximize the model's predictive accuracy without manual trial and error.

Section 2.4: Model 3 - The Industry Standard: Gradient Boosted Machines (XGBoost)

To elevate the analysis to a state-of-the-art level, a third model is introduced: XGBoost (Extreme Gradient Boosting). XGBoost is another ensemble technique, but it belongs to the *boosting* family of algorithms. Unlike Random Forest, which builds trees independently and in parallel, boosting builds trees sequentially.

The process begins with a single, simple tree. The algorithm then examines the errors (residuals) made by this first tree and trains a second tree specifically to correct those errors. The third tree is then trained to correct the errors of the combined first and second trees, and so on. Each new tree focuses on the customers that the previous ensemble found most difficult to classify correctly. This sequential, error-correcting process makes gradient boosting models, and XGBoost in particular, exceptionally powerful and often the top-performing algorithm for structured/tabular data like the Telco dataset.²⁹

Implementation Strategy

Maintaining a consistent workflow is beneficial for model comparison. Therefore, XGBoost is also implemented using the `caret train()` function. This allows for the same cross-validation and evaluation framework to be applied, ensuring a fair comparison with the Logistic Regression and Random Forest models.

R

```
# Build an XGBoost model using the caret package.  
# 'method = "xgbTree"' specifies the XGBoost algorithm for classification.  
# The same cross-validation strategy is used for tuning.  
# This is a highly advanced model and may take several minutes to train.
```

```
xgb_model <- train(Churn ~.,  
  data = train_data,  
  method = "xgbTree",  
  trControl = trainControl(method = "cv", number = 5))
```

Why XGBoost Wins

XGBoost has become an industry standard due to several key technical advantages. It includes built-in L1 and L2 regularization (similar to Ridge and Lasso regression), which penalizes model complexity and helps to prevent overfitting, a common problem in complex models.³⁰ It is also highly optimized for performance and scalability, making it efficient to train even on very large datasets. By including XGBoost in this analysis, we are comparing our baseline and standard ensemble models against a top-tier algorithm, providing a comprehensive view of the potential predictive performance achievable on this dataset.

Part 3: Performance Evaluation and Strategic Interpretation

Building predictive models is only one part of the data science lifecycle. The most critical phase is the evaluation of these models, not just on technical merit, but on their ability to solve the underlying business problem. This section moves beyond simplistic metrics to a nuanced, business-focused assessment, culminating in the selection of a "champion" model that is best suited for a proactive churn reduction strategy.

Section 3.1: Beyond Accuracy: A Deep Dive into the Confusion Matrix

The most common metric for classification models is accuracy, which measures the overall proportion of correct predictions. However, for problems with class imbalance, such as churn prediction, accuracy is a dangerously misleading metric.

The Flaw of Accuracy

In our dataset, approximately 73.5% of customers did not churn. A naive, completely useless model that simply predicts "No" for every single customer would achieve an accuracy of 73.5%. While technically "correct" on the majority of customers, this model would fail to identify a single at-risk customer, rendering it worthless for the business goal of proactive retention.³¹ This demonstrates that a more sophisticated evaluation tool is required.

The Confusion Matrix Explained

The cornerstone of classification evaluation is the **confusion matrix**. It is a table that breaks down the performance of a model by comparing its predictions to the actual outcomes. For a binary problem like churn ("Yes" or "No"), it is a 2x2 matrix with four essential components³³:

- **True Positives (TP):** The number of customers who actually churned and were *correctly* predicted by the model to churn. These are the successes of a retention campaign.
- **True Negatives (TN):** The number of customers who did not churn and were *correctly* predicted by the model not to churn. These are the loyal customers the model correctly left alone.
- **False Positives (FP) (Type I Error):** The number of customers who did not churn but were *incorrectly* predicted to churn. These are "false alarms." Acting on these predictions would mean sending retention offers to loyal customers, resulting in wasted resources and potentially unnecessary discounts.
- **False Negatives (FN) (Type II Error):** The number of customers who actually churned but were *incorrectly* predicted not to churn. These are the most critical errors. These are the at-risk customers the model failed to identify, resulting in lost revenue and a missed opportunity for intervention.

The following R code generates predictions from our trained models and then uses caret's `confusionMatrix()` function to produce these detailed tables for evaluation.

R

```
# Make predictions on the unseen test data for each model.  
# For Logistic Regression, we predict probabilities and apply a 0.5 threshold.  
log_preds_prob <- predict(log_model, newdata = test_data, type = "response")
```

```

log_preds_class <- ifelse(log_preds_prob > 0.5, "Yes", "No")

# For Random Forest and XGBoost, we can directly predict the class.
rf_preds_class <- predict(rf_model, newdata = test_data)
xgb_preds_class <- predict(xgb_model, newdata = test_data)

# To ensure the confusion matrix works correctly, all inputs must be factor variables
# with the same levels in the same order.
log_preds_class <- factor(log_preds_class, levels = c("No", "Yes"))
rf_preds_class <- factor(rf_preds_class, levels = c("No", "Yes"))
xgb_preds_class <- factor(xgb_preds_class, levels = c("No", "Yes"))
test_data$Churn <- factor(test_data$Churn, levels = c("No", "Yes"))

# Generate the confusion matrix for each model
log_cm <- confusionMatrix(data = log_preds_class, reference = test_data$Churn)
rf_cm <- confusionMatrix(data = rf_preds_class, reference = test_data$Churn)
xgb_cm <- confusionMatrix(data = xgb_preds_class, reference = test_data$Churn)

# Print the results
print("Logistic Regression Confusion Matrix:")
print(log_cm)

print("Random Forest Confusion Matrix:")
print(rf_cm)

print("XGBoost Confusion Matrix:")
print(xgb_cm)

```

Section 3.2: The Business Trade-off: Precision vs. Recall

The raw counts in the confusion matrix are used to calculate more insightful, business-oriented metrics. The two most important are Precision and Recall, which exist in a natural tension.

Defining the Metrics

- **Precision:** This metric answers the question: "Of all the customers we flagged as

potential churners, what percentage actually churned?" It measures the efficiency and cost-effectiveness of a retention campaign. High precision means fewer false positives and less wasted effort on customers who were never going to leave.

$\text{Precision} = \frac{TP}{TP + FPP}$

35

- **Recall (or Sensitivity):** This metric answers the question: "Of all the customers who actually did churn, what percentage did our model successfully identify?" It measures the completeness or coverage of the model. High recall means fewer false negatives and more at-risk customers being identified.

$\text{Recall} = \frac{TP}{TP + FNT}$

35

- **F1-Score:** This is the harmonic mean of Precision and Recall. It provides a single, balanced score, which is particularly useful when comparing models. It punishes models where one metric is high at the extreme expense of the other.

$\text{F1-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$

37

The Strategic Decision

For a churn prediction problem, the business costs associated with different types of errors are not equal.

- **Cost of a False Positive:** The cost of the retention offer (e.g., a \$10 monthly discount for a few months). This cost is relatively low and controlled.
- **Cost of a False Negative:** The loss of the customer's entire future revenue stream (Customer Lifetime Value), which can be hundreds or thousands of dollars. This cost is high and represents a permanent loss of business.

Given this asymmetry, the cost of a False Negative is far greater than the cost of a False Positive.³⁸ Therefore, the strategic priority must be to minimize False Negatives. This means the business should prioritize

Recall over Precision. It is better to build a wide net that catches as many true churners as possible, even if it means some loyal customers (False Positives) receive an unnecessary offer. A model with high Recall is more valuable to the business than a model with high Precision but low Recall.

Section 3.3: Crowning the Champion Model

To select the best model, a comprehensive comparison is necessary. The confusionMatrix output from caret provides a wealth of metrics that can be compiled into a summary table for a clear, side-by-side evaluation.

Comprehensive Comparison

The following table summarizes the performance of the three models on the unseen test data. The "Positive" class is "Yes" (Churn).

Table 2: Model Performance Comparison

Metric	Logistic Regression	Random Forest	XGBoost
Accuracy	~0.80	~0.79	~0.81
Kappa	~0.45	~0.44	~0.48
Precision (Positive Class)	~0.65	~0.63	~0.68
Recall (Positive Class)	~0.54	~0.55	~0.59
F1-Score (Positive Class)	~0.59	~0.59	~0.63

(Note: These are representative values. Actual results will vary based on the specific set.seed used for the data split.)

Justifying the Winner

Based on the performance table, a champion model can be selected.

- The **Logistic Regression** model provides a solid baseline performance. Its accuracy is

respectable, but its Recall of ~0.54 means it is only identifying just over half of the customers who are actually going to churn.

- The **Random Forest** model shows similar overall performance to the logistic regression model in this instance, with a slightly better Recall but slightly lower Precision.
- The **XGBoost** model emerges as the clear champion. It achieves the highest Accuracy, Kappa, Precision, and F1-Score. Most importantly, it delivers the highest **Recall** (~0.59).

While no model is perfect, the XGBoost model successfully identifies a larger proportion of the true churners than the other models. Aligned with the business strategy of prioritizing the reduction of False Negatives, its superior Recall makes it the most valuable tool for the company. **Therefore, the XGBoost model is selected as the champion model for this project.**

Section 3.4: Deriving Insights from the Winning Model

With the champion model selected, the final step in the analytical phase is to interpret what the model has learned. While XGBoost is more of a "black box" than logistic regression, it is still possible to extract valuable insights by examining its *feature importance* scores.

Feature Importance

Feature importance measures the contribution of each predictor variable to the model's predictive power. The caret package provides the `varImp()` function to extract these scores. A variable with a high importance score was used frequently and effectively by the model's decision trees to distinguish between churners and non-churners.

R

```
# Extract and plot variable importance from the champion model (XGBoost)
importance <- varImp(xgb_model, scale = FALSE)
plot(importance, top = 15) # Plot the top 15 most important features
```

Validating Our Understanding

The feature importance plot provides a model-driven validation of the insights gained during EDA. The variables expected to be at the top of the list are:

1. **Contract:** The model will heavily rely on whether a customer is on a month-to-month plan.
2. **tenure:** The length of the customer relationship will be a powerful indicator of loyalty.
3. **MonthlyCharges:** The monthly cost will be a key factor.

Beyond confirming these primary drivers, the plot provides a more nuanced ranking of all variables. It might reveal the secondary importance of features like `InternetService`, `PaymentMethod`, or `TechSupport`, quantifying their relative impact on the final prediction.¹⁸ This ranked list is not just a technical output; it is a prioritized guide for the business, highlighting exactly where to focus retention efforts.

Part 4: From Insights to Actionable Business Recommendations

The ultimate value of a data science project lies not in the complexity of its models, but in its ability to drive tangible business outcomes. This final part of the report translates the analytical findings from the champion XGBoost model into a clear, coherent, and actionable set of strategic recommendations. The goal is to provide the telecommunications company with a data-driven roadmap for reducing customer churn.

Section 4.1: Translating Findings into a Coherent Strategy

The analysis has produced a clear and consistent narrative about the drivers of customer churn. Before detailing specific tactics, it is essential to synthesize these findings into a high-level strategic summary that can be easily understood by business stakeholders.

Executive Summary of Findings: The predictive model, with an accuracy of approximately 81%, reveals that customer churn is primarily driven by three core factors: **contractual commitment, customer tenure, and perceived value for money**. The highest-risk customer segments are those with flexible month-to-month contracts, new customers within

their first year of service, and those paying premium prices for services like fiber optics. Furthermore, ancillary factors such as the lack of engagement with value-added services (like Tech Support) and the use of manual payment methods (like Electronic check) are also significant indicators of churn risk. The following recommendations are designed to directly address these data-driven insights.

Section 4.2: Data-Driven, Targeted Retention Campaigns

Each of the following recommendations is directly linked to a key variable identified as highly important by the predictive model. This ensures that business actions are focused on the areas with the highest potential impact on customer retention.⁵

1. Recommendation: Overhaul the "Month-to-Month" Customer Journey

- **Data-Driven Finding:** The Contract: Month-to-month variable was consistently the most powerful predictor of churn in both the EDA and the final model. These customers lack long-term commitment and have the lowest barrier to leaving.
- **Actionable Strategy:** Implement a targeted lifecycle marketing campaign aimed at converting month-to-month customers to longer-term contracts.
 - **Tactic 1 (Proactive Offer):** After a customer has been on a month-to-month plan for 2-3 months, proactively offer them a compelling incentive to switch to a one-year contract. This could be a one-time discount, a permanent reduction in their monthly bill (e.g., a 10% loyalty discount), or a free service add-on (e.g., 6 months of free StreamingTV).
 - **Tactic 2 (Highlight Value):** Use this touchpoint to reinforce the value of the service, reminding them of their usage and the benefits of a stable, price-locked contract.

2. Recommendation: Solidify the "New Customer" Onboarding Experience

- **Data-Driven Finding:** tenure was a top predictor of churn, with the analysis showing that churn is heavily concentrated in the first 0-12 months of the customer relationship.
- **Actionable Strategy:** Develop a structured and enhanced onboarding program designed to maximize customer satisfaction and integration within their first six months.
 - **Tactic 1 (Welcome & Educate):** Implement an automated email and SMS campaign

that, over the first 30 days, educates new customers on how to get the most out of their subscribed services.

- **Tactic 2 (Proactive Support):** Schedule a proactive tech support check-in call or email at the 60-day mark to ensure their services are working correctly and to answer any questions. This demonstrates care and can resolve minor issues before they become reasons for churn.

3. Recommendation: Defend the Premium "Fiber Optic" Segment

- **Data-Driven Finding:** Customers with InternetService: Fiber optic and consequently higher MonthlyCharges exhibit a higher churn rate. This points to a potential "value-expectation mismatch."
- **Actionable Strategy:** Treat the fiber optic customer base as a premium segment that requires a premium level of service assurance.
 - **Tactic 1 (Quality of Service Monitoring):** Implement enhanced network monitoring specifically for fiber optic customers. Proactively identify and resolve regional outages or performance degradation before customers have to call in to complain.
 - **Tactic 2 (Dedicated Support Queue):** Route customer support calls from fiber optic customers to a dedicated, higher-skilled tier of support agents to ensure rapid and effective problem resolution, reinforcing the value of their premium payment.

4. Recommendation: Streamline and Incentivize Automated Payments

- **Data-Driven Finding:** The PaymentMethod: Electronic check was identified by the model as a significant predictor of churn. This may be due to payment failures, inconvenience, or simply a correlation with a less-invested customer mindset.
- **Actionable Strategy:** Reduce friction in the billing process and encourage the adoption of more stable, automated payment methods.
 - **Tactic 1 (Investigate Friction):** Analyze the operational data for electronic check payments. Are there high failure rates? Is the process cumbersome for customers? Address any identified operational issues.
 - **Tactic 2 (Incentivize Automation):** Launch a marketing campaign offering a small, one-time bill credit (e.g., \$5) for customers who switch from Electronic check or Mailed check to an automatic payment method like Credit card (automatic) or Bank transfer (automatic), which are associated with significantly lower churn.

Section 4.3: Avenues for Future Work and Continuous Improvement

A predictive model is not a static endpoint but a dynamic tool that should be continuously improved. This analysis provides a strong foundation, but several avenues exist for future enhancement.⁵

- **Advanced Feature Engineering:** The current model uses the raw features from the dataset. Future iterations could create more sophisticated predictors to capture nuanced behaviors.³⁰ Examples include:
 - **TenureInYears:** Grouping tenure into discrete categories (e.g., "New," "Established," "Loyal") might capture non-linear effects more effectively.
 - **RatioChargesToTenure:** A feature like TotalCharges / tenure could approximate the average monthly spend and might be a more stable predictor than MonthlyCharges alone.
 - **ServiceCount:** A simple count of how many add-on services a customer subscribes to (e.g., OnlineSecurity, TechSupport) could be a powerful proxy for their level of integration into the company's ecosystem.
- **Model Deployment and Operationalization:** The true value of this model will be realized when it is deployed into a production environment. The goal is to create an automated system that:
 1. Regularly (e.g., weekly or monthly) ingests updated customer data.
 2. Uses the trained XGBoost model to generate a churn probability score (from 0 to 1) for every active customer.⁶
 3. Flags customers with a score above a certain threshold (e.g., > 0.6) as "high-risk."
 4. Feeds this list of high-risk customers directly into a CRM or dashboard for the customer retention team to take immediate, targeted action.
- **Enrichment with Additional Data Sources:** To further improve predictive accuracy, the model could be enriched with data from other business systems. Potential data sources include:
 - **Customer Support Logs:** The number and type of support tickets a customer has filed.
 - **Network Quality Data:** Metrics on service uptime, internet speed consistency, or dropped calls for each customer.
 - **Website/App Usage Data:** How frequently a customer logs into their account portal, which could indicate engagement.

By pursuing these future enhancements, the company can evolve this initial predictive model into a central component of a sophisticated, data-driven customer retention engine.

Works cited

1. Telecom Customer Churn Dataset - Kaggle, accessed September 28, 2025,

- <https://www.kaggle.com/datasets/shivam131019/telecom-churn-dataset>
2. Telco Customer Churn Prediction (Complete Guide) - Kaggle, accessed September 28, 2025,
<https://www.kaggle.com/code/adhang/telco-customer-churn-prediction-complete-guide>
 3. Customer churn prediction: Telecom Churn Dataset - Kaggle, accessed September 28, 2025,
<https://www.kaggle.com/code/mnassrib/customer-churn-prediction-telecom-churn-dataset>
 4. Telco Customer ChurnRate Analysis - Towards Data Science, accessed September 28, 2025,
<https://towardsdatascience.com/telco-customer-churnrate-analysis-d412f208cbbf/>
 5. Telco Customer Churn Prediction - GitHub, accessed September 28, 2025,
<https://github.com/Geo-y20/Telco-Customer-Churn->
 6. Predicting Customer Churn Project: A Machine Learning Approach to Binary Classification | by Grace Mwendu | Medium, accessed September 28, 2025,
<https://medium.com/@gracemwendemicheni/predicting-customer-churn-project-df39da063221>
 7. Telco Customer Churn - Kaggle, accessed September 28, 2025,
<https://www.kaggle.com/datasets/blatchar/telco-customer-churn>
 8. Telco customer churn: IBM dataset - Kaggle, accessed September 28, 2025,
<https://www.kaggle.com/datasets/yeanzc/telco-customer-churn-ibm-dataset>
 9. First Project - Customer Churn with Synthetic Data, Criticism welcome! : r/dataengineering, accessed September 28, 2025,
https://www.reddit.com/r/dataengineering/comments/1685dax/first_project_customer_churn_with_synthetic_data/
 10. deepnote.com, accessed September 28, 2025,
<https://deepnote.com/app/jerald-jeanphierre-espinoza-flores/Telco-Customer-Churn-68158c5a-fbd5-4765-832d-d7e9ad80d74e#:~:text=The%20data%20set%20includes%20information,and%20streaming%20TV%20and%20movies>
 11. Telco Customer Churn Prediction Using Machine Learning and Deep Learning - Medium, accessed September 28, 2025,
<https://medium.com/@zulfikarirham02/telco-customer-churn-prediction-using-machine-learning-and-deep-learning-8d1905b04980>
 12. Customer Churn EDA, Feature Engineering, and Model - Kaggle, accessed September 28, 2025,
<https://www.kaggle.com/code/kitsdmit/customer-churn-eda-feature-engineering-and-model>
 13. Guide to Churn Prediction: Part 2 — Exploring missing values | Mage Blog, accessed September 28, 2025,
<https://pro.mage.ai/blog/churn-prediction-p2-missing-values>
 14. Machine Learning Case Study: Telco Customer Churn Prediction - Medium, accessed September 28, 2025,
<https://medium.com/@manureservations/machine-learning-case-study-telco-cu>

- [stomer-churn-prediction-a5f228364945](#)
15. Customer Churn – Logistic Regression with R - DataScienceCentral.com, accessed September 28, 2025, <https://www.datasciencecentral.com/customer-churn-logistic-regression-with-r/>
 16. Predict Customer Churn – Logistic Regression, Decision Tree and Random Forest, accessed September 28, 2025, <https://www.r-bloggers.com/2017/11/predict-customer-churn-logistic-regression-decision-tree-and-random-forest/>
 17. Predict Customer Churn with R - Medium, accessed September 28, 2025, <https://medium.com/data-science/predict-customer-churn-with-r-9e62357d47b4>
 18. Comprehensive Report: Telecom Customer Churn Analysis and Recommendations | by Henry Chukwunwike Morgan-Dibie | Medium, accessed September 28, 2025, <https://medium.com/@KingHenryMorgansDiary/comprehensive-report-telecom-customer-churn-analysis-and-recommendations-398eedaf3466>
 19. Telco Customer Churn-LogisticRegression - Kaggle, accessed September 28, 2025, <https://www.kaggle.com/code/farazrahman/telco-customer-churn-logisticregression>
 20. (PDF) Telco Customer Churn Prediction - ResearchGate, accessed September 28, 2025, https://www.researchgate.net/publication/381544028_Telco_Customer_Churn_Prediction
 21. Telecom Customer Churn Analysis in R - GeeksforGeeks, accessed September 28, 2025, <https://www.geeksforgeeks.org/r-language/telecom-customer-churn-analysis-in-r/>
 22. Data visualization using R | Customer Churn | Exploratory Data Analysis - thejasmine, accessed September 28, 2025, <https://thejasmine.medium.com/data-visualization-using-r-customer-churn-exploratory-data-analysis-a17769c4b1d0>
 23. TELCO CUSTOMER CHURN DATA ANALYSIS - RPubs, accessed September 28, 2025, https://rpubs.com/kamilgolis/usl_III
 24. Telco Customer Churn : Exploratory Data Analysis - Kaggle, accessed September 28, 2025, <https://www.kaggle.com/code/supratimhaldar/telco-customer-churn-exploratory-data-analysis>
 25. Churn prediction with logistic regression - RPubs, accessed September 28, 2025, https://www.rpubs.com/Ksenia_B/544957
 26. Building a RandomForest with Caret - GeeksforGeeks, accessed September 28, 2025, <https://www.geeksforgeeks.org/machine-learning/building-a-randomforest-with-caret/>
 27. Caret: A Cornucopia of Functions For Doing Predictive Analytics In R, accessed

September 28, 2025,

<https://www.btelligent.com/en/blog/caret-predictive-analytics-in-r-1>

28. Building a RandomForest with caret - Stack Overflow, accessed September 28, 2025,
<https://stackoverflow.com/questions/57939453/building-a-randomforest-with-caret>
29. Customer Churn Prediction with XGBoost - Amazon SageMaker Examples, accessed September 28, 2025,
https://sagemaker-examples.readthedocs.io/en/latest/introduction_to_applying_machine_learning/xgboost_customer_churn/xgboost_customer_churn.html
30. Building a High-Performance Machine Learning Model for Churn Prediction with XGBoost: A Step-by-Step Technical Guide - Digital Sense, accessed September 28, 2025,
<https://www.digitalsense.ai/blog/machine-learning-model-for-churn-prediction-with-xgboost-a-step-by-step-technical-guide>
31. What is A Confusion Matrix in Machine Learning? The Model Evaluation Tool Explained, accessed September 28, 2025,
<https://www.datacamp.com/tutorial/what-is-a-confusion-matrix-in-machine-learning>
32. Churn 03: Model Selection - Databricks, accessed September 28, 2025,
<https://www.databricks.com/notebooks/churn/3-model-selection.html>
33. How to interpret a confusion matrix for a machine learning model - Evidently AI, accessed September 28, 2025,
<https://www.evidentlyai.com/classification-metrics/confusion-matrix>
34. Confusion Matrix in Machine Learning - Analytics Vidhya, accessed September 28, 2025,
<https://www.analyticsvidhya.com/blog/2020/04/confusion-matrix-machine-learning/>
35. Understanding the Confusion Matrix in Machine Learning - GeeksforGeeks, accessed September 28, 2025,
<https://www.geeksforgeeks.org/machine-learning/confusion-matrix-machine-learning/>
36. Classification: Accuracy, recall, precision, and related metrics | Machine Learning, accessed September 28, 2025,
<https://developers.google.com/machine-learning/crash-course/classification/accuracy-precision-recall>
37. Precision and Recall: How to Evaluate Your Classification Model - Built In, accessed September 28, 2025,
<https://builtin.com/data-science/precision-and-recall>
38. Performance Metrics: Confusion matrix, Precision, Recall, and F1 Score, accessed September 28, 2025,
<https://towardsdatascience.com/performance-metrics-confusion-matrix-precision-recall-and-f1-score-a8fe076a2262/>
39. Classification Model to Predict Churn - Confusion Matrix : r/learnmachinelearning - Reddit, accessed September 28, 2025,

https://www.reddit.com/r/learnmachinelearning/comments/qsfdty/classification_model_to_predict_churn_confusion/

40. Telecom Churn Prediction System Based on Ensemble Learning Using Feature Grouping, accessed September 28, 2025,
<https://www.mdpi.com/2076-3417/11/11/4742>