# MIDTERM REPORT: Predicting Customer Churn in the Telecom Industry Using Data Mining

**Team Members: Tajwar-Ul Hoque, Luis Escamilla, Isaiah Malcom**

## 1. Data

Our team is using the Telco Customer Churn dataset from Kaggle. This dataset includes 7,043 instances and 21 attributes that describe customer demographics, service subscriptions, billing information, and churn status. The target variable is Churn, which shows whether a customer left the company or not. The dataset contains both numerical and categorical variables, so we prepared it carefully for modeling.

In preprocessing, we first converted the TotalCharges column to numeric and filled missing values using the median to keep the data consistent. We dropped customerID since it has no predictive value. The Churn column was encoded as a binary variable, with "Yes" mapped to 1 and "No" mapped to 0. We then applied one-hot encoding to all categorical features to convert them into a numeric format suitable for machine learning models. Finally, we split the dataset into training (80%) and testing (20%) subsets, using stratified sampling to maintain the same churn proportion in both sets. These preprocessing steps were completed using pandas and scikit-learn in Python, ensuring the data is clean and ready for modeling.
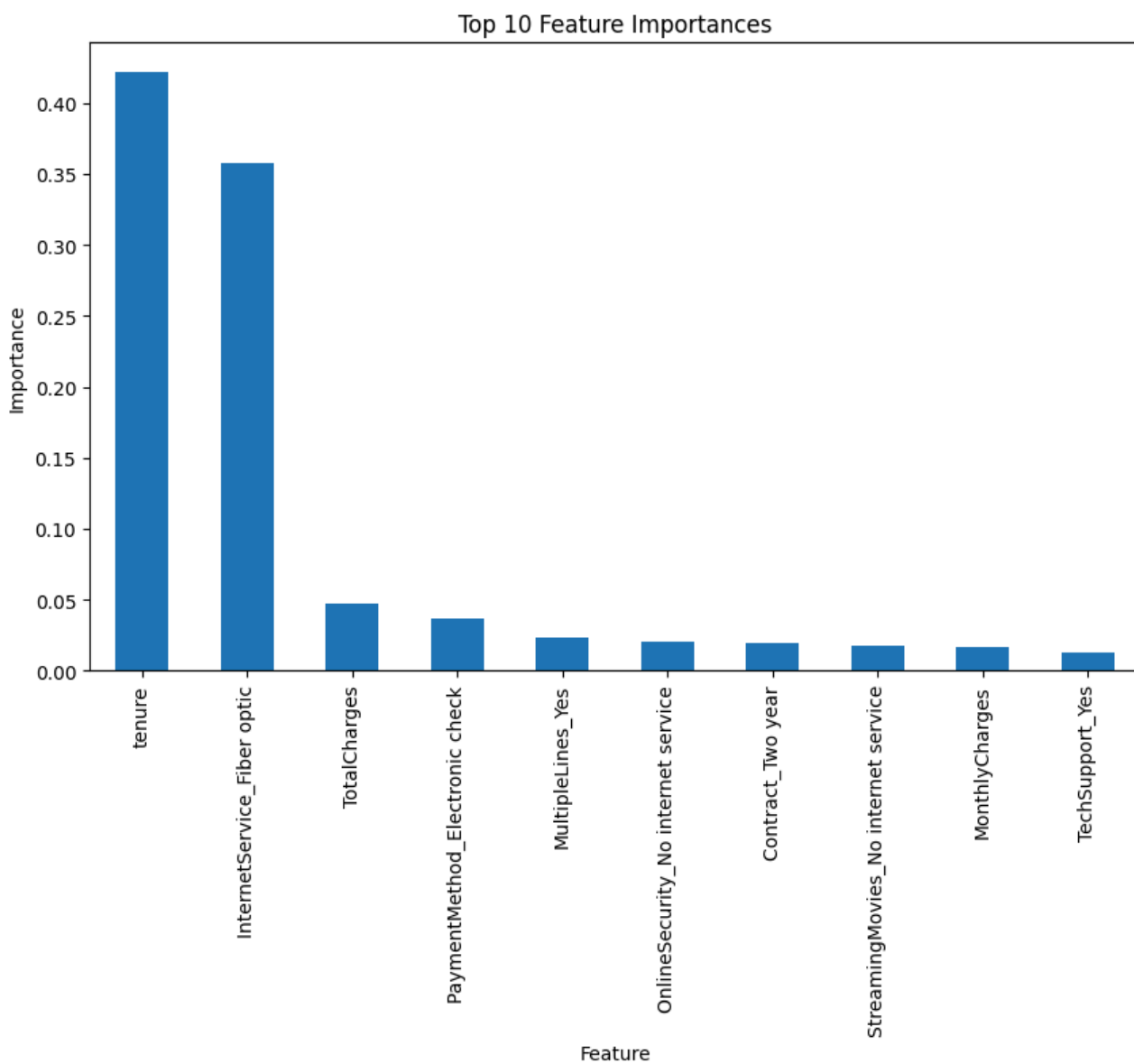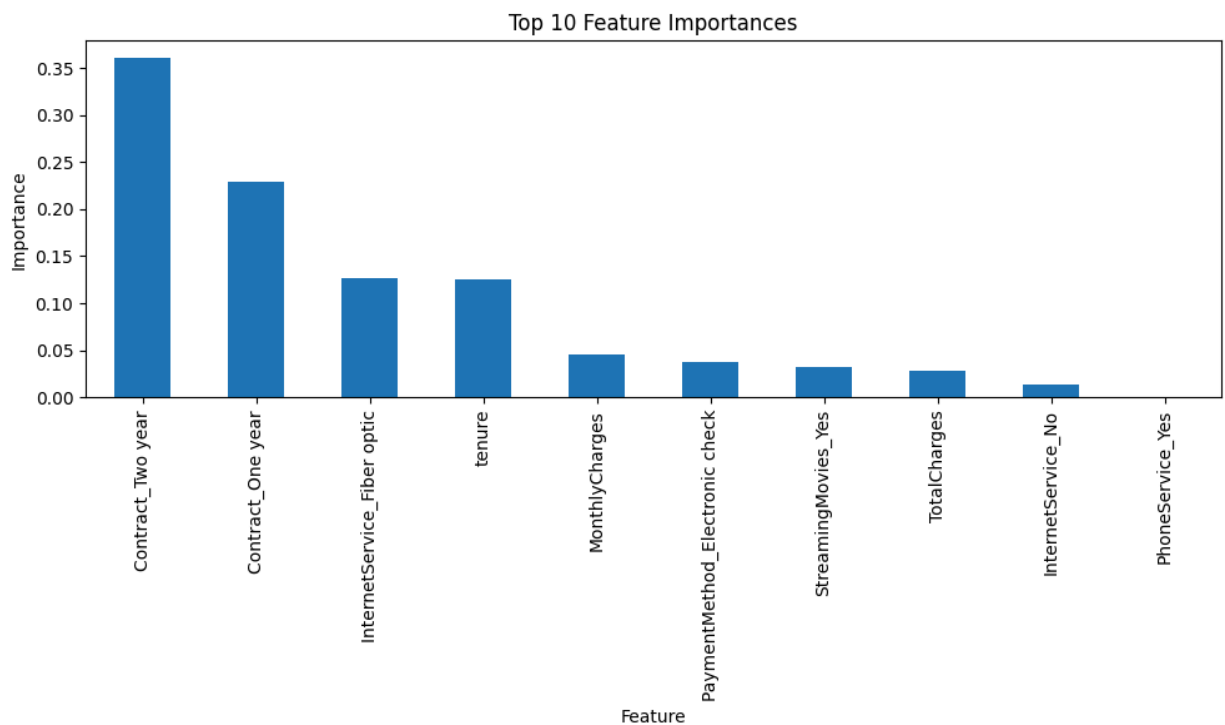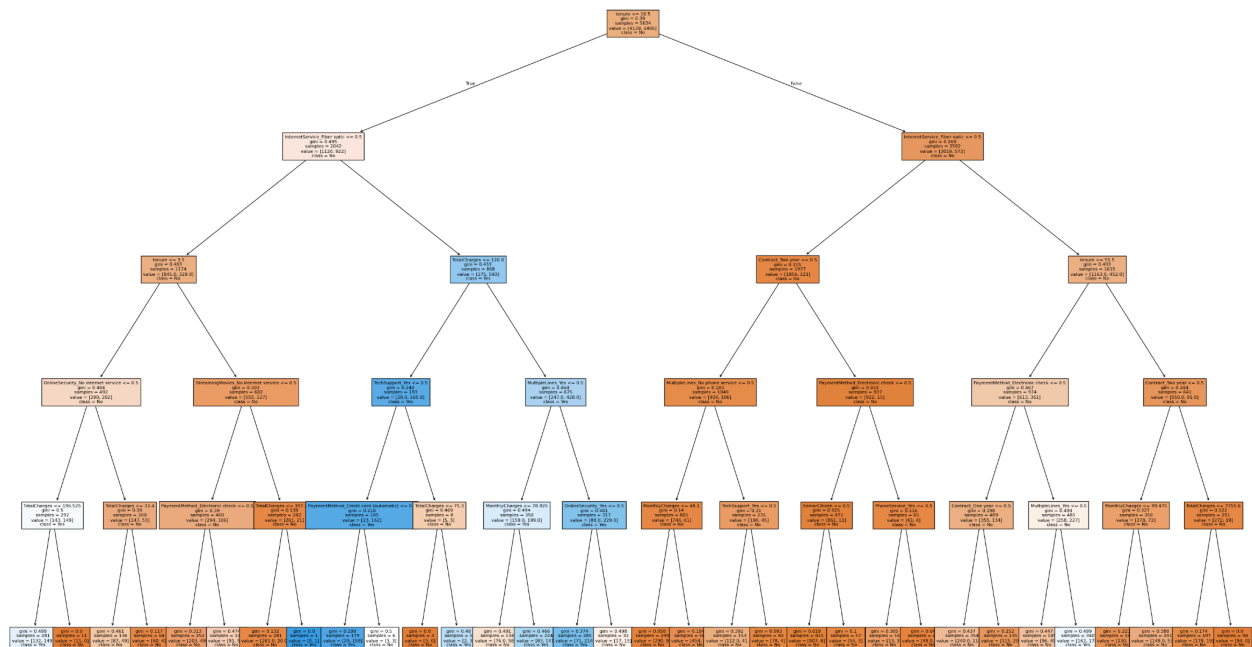
## 2. Data Mining Task

The main task for this project is classification. Our goal is to predict which customers are most likely to churn based on the available features. We are using supervised learning techniques for this purpose and plan to evaluate all models using accuracy, precision, recall, F1-score, and ROC-AUC. By building models that can accurately classify churn, we can identify at-risk customers and help telecom companies take proactive steps to retain them.

## 3. Progress

So far, our team has completed the first round of model training and evaluation using a Decision Tree Classifier as the baseline model. This model was trained using the Gini criterion with a maximum depth of 5. The Decision Tree achieved an accuracy of about 79% on the test set, showing solid baseline performance for this stage of the project. The results were further analyzed using classification reports and confusion matrices to verify consistency across precision and recall.

To better understand model behavior, we visualized the tree and saved it as tree.svg. We also generated a feature importance bar chart, saved as top_ten.svg, which highlights the top ten features contributing most to churn prediction. The most influential factors were Contract type, tenure, and MonthlyCharges, which align with expected business insights for telecom companies.

Top 10 Feature Importances

Top 10 Feature Importances



Our next steps include developing Logistic Regression, Random Forest, and XGBoost models , and potentially neural networks to compare their results against the baseline Decision Tree and determine which approach performs best.

**4. Challenges**

The main challenges we have faced involve handling class imbalance between churned and retained customers and balancing interpretability and performance. Decision Trees provide clear explanations but can struggle with overfitting if not tuned properly. The dataset is also moderately small, which limits the advantage of complex models. We plan to address these challenges by testing models with class weighting, performing cross-validation, and tuning hyperparameters for better generalization.

**5. Schedule and Next Steps**

At this point, we have completed data cleaning, preprocessing, and initial Decision Tree modeling. The upcoming phase of our project will focus on building and comparing additional models, including Logistic Regression, Random Forest, and XGBoost. We will also tune model hyperparameters using GridSearchCV to optimize performance. Once those models are evaluated, we will compare all results using multiple metrics such as accuracy, precision, recall, F1-score, and ROC-AUC.

We also plan to conduct a deeper feature importance analysis across all models and prepare visualizations to support our final report. The last stage will focus on compiling results, preparing visuals, and writing the final presentation and report.

**6. Summary**

Overall, our project is progressing according to the planned schedule. The preprocessing and baseline Decision Tree modeling steps were successful, giving us an initial accuracy of around 79%. The early results already highlight key factors affecting churn such as contract length, tenure, and monthly billing amounts. Moving forward, we aim to refine our predictive accuracy through advanced models and further analysis. This will help us provide stronger business and data-driven insights in the final stage of the project.

Link To The GitHub Repo:
https://github.com/Isaiah392/Data-Mining-Customer-Churn

**NOTE:** Please reach out to the team if you have any further question.