

STAT 441 Project Proposal

Maisha Thasin, Joseph Wang (Group 25, undergraduate)

Winter 2023

Dataset

The Telco Customer Churn dataset, provided by IBM, contains information about a fictional telecommunications (telco) company which provided home phone and internet services to 7043 customers. The dataset provides demographic and business-related metrics for each customer, as well as identifying whether the customer switched providers (customer churn).

The dataset can be downloaded from the following link: https://accelerator.ca.analytics.ibm.com/bi/?perspective=authoring&pathRef=.public_folders%2FIBM%2BAccelerator%2BCatalog%2FContent%2FDAT00067

The dataset contains 33 variables for 7043 observations, but not all variables are fit to be predictive features. We removed columns relating to unique IDs, geographical information, and dashboarding aggregation, as well as “duplicate” columns (those that are identical to another column except for formatting), and columns related to the response (such as the churn reason and predicted lifetime value to the company).

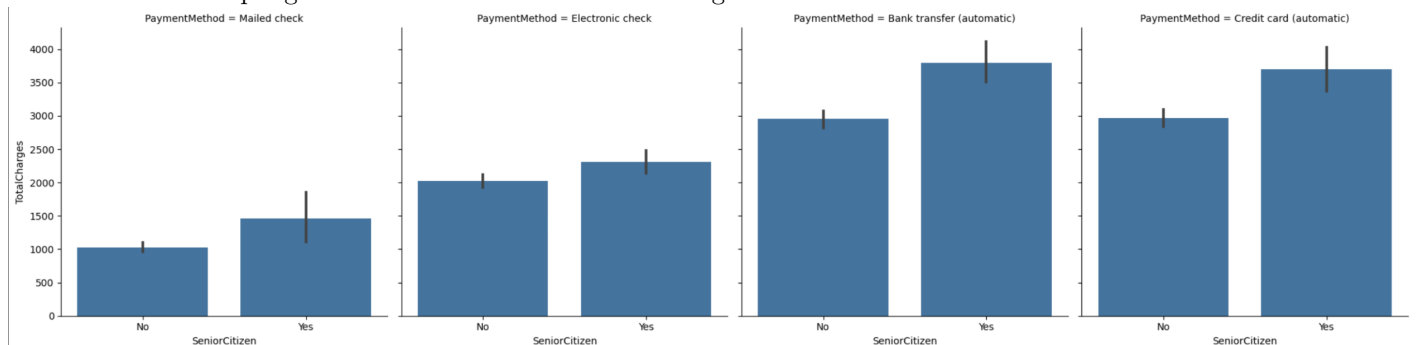
We are then left with 19 features: Gender, Senior Citizen, Partner, Dependents, Tenure Months, Phone Service, Multiple Lines, Internet Service, Online Security, Online Backup, Device Protection Plan, Tech Support, Streaming TV, Streaming Movies, Contract, Paperless Billing, Payment method, Monthly Charge, and Total Charges. Each feature is categorical except Tenure, which is (integer) numeric, and Monthly Charge/Total Charges, which are (float) numeric. The descriptions of the features can be found at <https://community.ibm.com/community/user/businessanalytics/blogs/steven-macko/2019/07/11/telco-customer-churn-1113>.

Finally, the response variable is Churn Value, binary on whether the customer left the company this quarter.

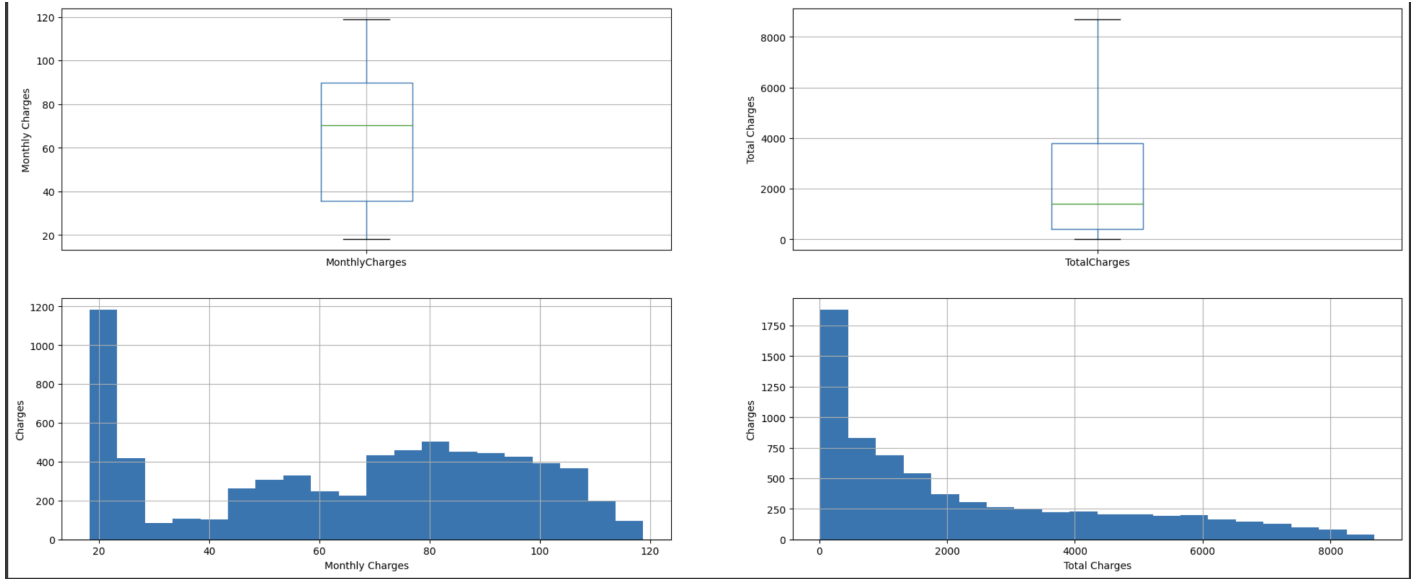
Exploratory Data Analysis

We started out with some EDA to understand the data before applying any machine learning techniques.

We noticed that the data is imbalanced; 73% of the data is labelled as No whereas 27% is labelled Yes, so we might have to choose our resampling methods to ensure that our training data is not biased.



Notice from the graph that there is a clear relationship between Senior citizens = Yes and the total charges they have. Senior citizens seem to be paying extra than non senior citizens, and customers that tend to pay by mailed check pay less than those who do bank transfers.



Furthermore, from our Boxplots and distribution we notice that the underlying charges are not distributed uniformly. And the total charges, although calculated with the monthly charges, do not have exactly the same distributions and are more left skewed than the Monthly charges. This may be the result of new customers who have not yet had multiple monthly charges for their total (quarterly) charge.

Classification Methods

We would like to compare the predictive accuracy of the more complex/computationally intensive models with the predictive accuracy of the more straightforward/computationally simple models. We plan to set aside some of the dataset as “new data” for this purpose; all of the models will train on the same subset of the dataset (called the visible data) and never see the “new data”. This way, we have a fair comparison of predictive power across all of the different models.

1. Logistic Regression: This is a widely used linear regression model for binary classification problems and it is less prone to overfitting. We will use scikit-learn to build our implementation. To calibrate the model, we can tune the regularization parameter using cross-validation.
2. Decision Trees: A decision tree algorithm can be used to classify customers as churners or non-churners based on a set of features. We will use scikit-learn to build our implementation. We can tune parameters like the maximum depth of the tree and the minimum number of samples required to split a node, as well as the cost complexity pruning parameter alpha and the number of folds used for cross validation. This algorithm is prone to overfitting if there is too much depth.
3. Random Forests: This is an ensemble learning algorithm that creates multiple decision trees (as above) and aggregates their predictions. We will use scikit-learn to build our implementation. We can tune parameters like the number of trees in the forest, the maximum depth of the trees, and the number of features considered for each split.
4. Support Vector Machines (SVMs): SVMs are particularly useful when the number of features is high relative to the number of observations, as they can handle high-dimensional data well. They try to find a hyperplane that separates the churners from the non-churners. SVMs are implemented in libraries such as scikit-learn in Python with the SVC class. We can tune parameters like the kernel type, the regularization parameter, and the kernel coefficient. With our EDA we can figure out if there is a higher number of outliers and use hard or soft margins, since we have a lot of features SVM is a good choice. The SVC class takes several parameters, including the type of kernel to use (linear, polynomial), the regularization parameter (C), and the kernel coefficient (gamma for RBF kernel, degree for polynomial kernel).
5. K-Nearest Neighbors (KNN): This is a non-parametric algorithm that classifies a new data point based on the k-nearest points in the training data. KNN is implemented in libraries such as scikit-learn in Python and we can use k-fold cross validation to calculate the best ‘k’. We can tune parameters like the number of neighbors (k) and the distance metric.
6. Neural Networks: Neural networks are a popular algorithm for classification problems for non-linear data. We will implement this using TensorFlow in Python. Since the dataset has a large number of training data this algorithm is a good choice. We can tune parameters like the number of layers, the number of neurons in each layer, and the activation functions.

After testing multiple models for accuracy, the next step is to select the best model based on its performance on the validation set. The validation set is used to estimate the performance of each model on unseen data, which is a better estimate of the model's true performance than the training set.

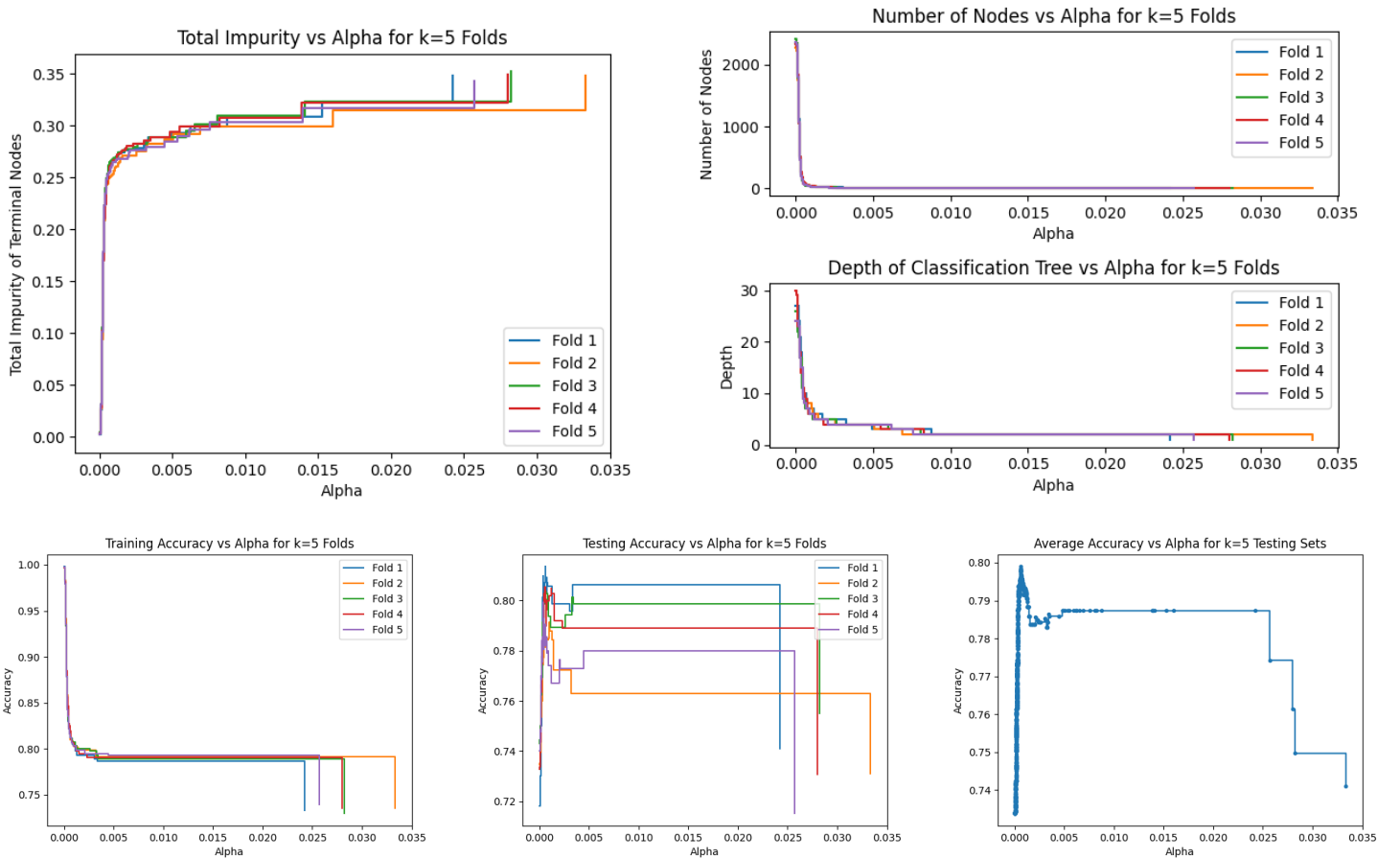
Once the best model has been selected, we will retrain the model on the entire (visible) dataset to get maximum predictive performance.

Preliminary Results

For our preliminary results, we fit two different models which we will use two of the computationally simple models to compare against.

Our first model is a Decision Tree (aka Classification Tree), built using Recursive Binary Splitting on the entire dataset. Then we used Cost Complexity Pruning and k-fold cross-validation (with $k=5$) to find the optimal value for alpha which defines a subtree.

We can produce 5 metrics for each of the folds as a function of the cost complexity parameter alpha: the total impurity of leaf/terminal nodes, the number of nodes in the tree, the depth of the tree, the training accuracy, and the testing accuracy. Finally, we can also produce the average testing accuracy across the 5 folds.



Our second model is a simple SVM that uses 90% of the pre-processed dataset to train the SVM classifier from scikit-learn. We fit the model on 4 different hyperparameters; default hyperparameters, Regularization parameter = 100, Regularization parameter = 100 and Kernel = polynomial, and another model with Regularization parameter = 100 and Kernel = sigmoid.

We then used this model to predict on our test dataset and predicted a model accuracy to 0.76, 0.80, 0.76 and 0.63 respectively.

Appendix