

# WGU D209 TASK 2 REV 1 - MATTINSON

## Classification Tree - Telecom Churn Data

Mike Mattinson

Master of Science, Data Analytics, WGU.edu

### Task 1 - Decision Classification Tree

November 15, 2021

**Configure Notebook.** Import and configure packages. All of the code for importing and configuring is contained in a `imports .PY` file. Also, there is a second `helpers .PY` file to define a few functions used throughout this notebook.

```
In [1]: # import and configure packages
from imports import *
%matplotlib inline
warnings.filterwarnings('ignore')
```

```
P:\code\wgu\py39\Scripts\python.exe
python version: 3.9.7
pandas version: 1.3.0
numpy version: 1.19.5
scipy version: 1.7.1
sklearn version: 1.0.1
matplotlib version: 3.4.2
seaborn version: 0.11.2
graphviz version: 0.17
```

```
In [2]: from helpers import *
```

```
getFilename version: 1.0
saveTable version: 1.0
describeData version: 1.0
createScatter version: 1.0
createBarplot version: 1.1
get_unique_numbers version: 1.0
createCorrelationMatrix version: 1.0
createStackedHistogram version: 1.0
```

## Part I: Research Question

A. Describe the purpose of this data mining report by doing the following:

**A1. Propose one question relevant to a real-world organizational situation that you will answer using one of the following prediction methods: (a) decision trees, (b) random forests or (c) advanced regression (i.e., lasso or ridge regression)**

**Primary Question.** A telecomm company is concerned about customer churn. It has collected customer data on 10,000 customers and the analysis will attempt to determine if the target value 'Churn', which is a yes-no categorical variable, if 'Churn' can be predicted accurately using only a minimum number of predictor variables. The prediction analysis will use the DecisionClassificationTree to predict the target value.

**A2. Define one goal of the data analysis. Ensure that your goal is reasonable within the scope of the scenario and is represented in the available data.**

**Primary Goal.** Primary goals will look at all predictor variables and then reduce the number of predictor variables while maintaining a reasonable level of accuracy. The company provided data will be split into training and testing data, 70% training and 30% testing. I will use the confusion matrix to show the total number of correct and incorrect predictions in the training and testing datasets.

## **Part II: Method Justification**

B. Explain the reasons for your chosen prediction method from part A1 by doing the following:

**B1. Explain how the prediction method you chose analyzes the selected data set. Include expected outcomes.**

**Explain Method.** I choose to use the DecisionClassificationTree to model the data in order to predict the target value using a minimum number of features. Once the model is trained using 70% of the data, then I will calculate the accuracy of the model. If the model is doing well, anyone will be able to enter the model with unclassified data and the model will accurately predict the target value.

**B2. Summarize one assumption of the chosen prediction method.**

**One Assumption.** The model is created using the training data and the

size of the training data is sufficient large to make the best decision points throughout the tree.

**B3. List the packages or libraries you have chosen for Python or R, and justify how each item on the list supports the analysis.**

All of the Python packages required for this analysis were loaded at the very top of th notebook. The packages and version numbers are presented. Besides the normal Python packages (numpy, scipy, matplotlib, pandas, etc.) the primary package required to create and view the prediction model comes from sklearn. Also, I use two (2) .PY files instead of including all that code in this notebook, I reuse it in other notebooks. Imports.py has all of the required packages and Helpers.py has many helpful functions that allow me to standardize my tables, figures and other parts of the notebook. Both of these .PY files will be included with the notebook for reference.

**Part III: Data Preparation**

C. Perform data preparation for the chosen data set by doing the following:

**C1. Describe one data preprocessing goal relevant to the prediction method from part A1.**

**Goal.** The data provided by the company includes both numeric and categoricald data, the first primary goal will be to convert the categorical data into data that the model can understand. I plan to convert the categorical data using the sklearn.preprocessing LabelEncoder. This will happen in section 'C3' below.

**C2. Identify the initial data set variables that you will use to perform the analysis for the classification question from part A1, and classify each variable as continuous or categorical.**

TABLE 3-1.DESCRPTION OF DATA IN TELECO CHURN DATASET	

In the telecommunications industry, customers can choose from multiple service providers and actively switch from one provider to another. Customer “churn” is defined as the percentage of customers who stopped using a provider’s product or service during a certain time frame. In this highly competitive market, some telecommunications industries can experience average annual churn rates as high as 25 percent. Given that it costs 10 times more to acquire a new customer than to retain an existing one, customer retention has now become even more important than customer acquisition.

For many providers, retaining highly profitable customers is the number one business goal. To reduce customer churn, telecommunications companies need to predict which customers are at high risk of churn.

You are an analyst on a team of analysts in a popular telecommunications company, which serves customers in all regions of the United States. You have been asked to clean the raw data set in preparation to explore the data, identify trends, and compare key metrics.

```
In [4]: # Load data from .csv
raw = pd.read_csv('data/churn_clean.csv')
raw = raw.drop(columns=['CaseOrder','UID',
                        'Customer_id','Interaction',
                        'City', 'State','County','Zip',
                        'Lat','Lng','Job','TimeZone'])
```

---

## TABLE 3-2.RAW DATA

Initial state of dataset before any manipulations.

```
In [5]: saveTable(data=raw, title='RAW', sect='C2',
               course='D209', task='Task2', caption='3 2')
```

	0	1	2	3
<b>Population</b>	38	10446	3735	13863
<b>Area</b>	Urban	Urban	Urban	Suburban
<b>Children</b>	0	1	4	1
<b>Age</b>	68	27	50	48
<b>Income</b>	28561.99	21704.77	9609.57	18925.23
<b>Marital</b>	Widowed	Married	Widowed	Married
<b>Gender</b>	Male	Female	Female	Male
<b>Churn</b>	No	Yes	No	No
<b>Outage_sec_perweek</b>	7.978	11.699	10.753	14.914
<b>Email</b>	10	12	9	15
<b>Contacts</b>	0	0	0	2
<b>Yearly_equip_failure</b>	1	1	1	0
<b>Techie</b>	No	Yes	Yes	Yes
<b>Contract</b>	One year	Month-to-month	Two Year	Two Year
<b>Port_modem</b>	Yes	No	Yes	No
<b>Tablet</b>	Yes	Yes	No	No
<b>InternetService</b>	Fiber Optic	Fiber Optic	DSL	DSL
<b>Phone</b>	Yes	Yes	Yes	Yes
<b>Multiple</b>	No	Yes	Yes	No
<b>OnlineSecurity</b>	Yes	Yes	No	Yes
<b>OnlineBackup</b>	Yes	No	No	No
<b>DeviceProtection</b>	No	No	No	No
<b>TechSupport</b>	No	No	No	No
<b>StreamingTV</b>	No	Yes	No	Yes
<b>StreamingMovies</b>	Yes	Yes	Yes	No
<b>PaperlessBilling</b>	Yes	Yes	Yes	Yes
<b>PaymentMethod</b>	Credit Card (automatic)	Bank Transfer(automatic)	Credit Card (automatic)	Mailed Check
<b>Tenure</b>	6.796	1.157	15.754	17.087
<b>MonthlyCharge</b>	172.456	242.633	159.948	119.957
<b>Bandwidth_GB_Year</b>	904.536	800.983	2054.707	2164.579
<b>Item1</b>	5	3	4	4
<b>Item2</b>	5	4	4	4
<b>Item3</b>	5	3	2	4
<b>Item4</b>	3	3	4	2
<b>Item5</b>	4	4	4	5

	0	1	2	3
Item6	4	3	3	4
Item7	3	4	3	3
Item8	4	4	3	3

shape: (10000, 38)

Table saved to: TABLES/D209\_TASK2\_C2\_TAB\_3\_2\_RAW.CSV

### C3. Explain each of the steps used to prepare the data for the analysis. Identify the code segment for each step.

#### Step 1.

The company data has already been cleaned and prepared. The only step required prior to creating the prediction model is to encode the categorical data. So, the raw data will be copied and then the clean data will encode all of the categorical data. The clean data will have all of the original variables but for example, all of the categorical data will have integer numbers representing the raw data, Area = 2 indicates Area=Rural in the raw data, etc.

```
In [6]: # start with a copy of raw data
clean = raw.copy()
```

```
In [7]: # define target and features
target = 'Churn'
X = clean.drop(columns=[target])
y = clean[target]
print(type(X))
print(type(y))
```

```
<class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.series.Series'>
```

```
In [8]: # define categorical features
cat_features = X.select_dtypes(include=['object']).columns
print(cat_features)
```

```
Index(['Area', 'Marital', 'Gender', 'Techie', 'Contract', 'Port_modem',
      'Tablet', 'InternetService', 'Phone', 'Multiple', 'OnlineSecurity',
      'OnlineBackup', 'DeviceProtection', 'TechSupport', 'StreamingTV',
      'StreamingMovies', 'PaperlessBilling', 'PaymentMethod'],
      dtype='object')
```

```
In [9]: # define numerical features
num_features = X.select_dtypes(exclude=['object']).columns
print(num_features)
```

```
Index(['Population', 'Children', 'Age', 'Income', 'Outage_sec_perweek',
      'Email', 'Contacts', 'Yearly_equip_failure', 'Tenure', 'MonthlyCharge',
      'Bandwidth_GB_Year', 'Item1', 'Item2', 'Item3', 'Item4', 'Item5',
      'Item6', 'Item7', 'Item8'],
      dtype='object')
```

```
In [10]: from sklearn.preprocessing import LabelEncoder
categorical_feature_names = []
label_encoders = {}
for categorical in cat_features:
    label_encoders[categorical] = LabelEncoder()
    clean[categorical] = label_encoders[categorical].fit_transform(clean[categorical])
    names = label_encoders[categorical].classes_.tolist()
    print('Label encoder %s - values: %s' % (categorical, names))
    if categorical == target:
        continue
    categorical_feature_names.extend([categorical + '_' + str(name) for name in names])
```

```
Label encoder Area - values: ['Rural', 'Suburban', 'Urban']
Label encoder Marital - values: ['Divorced', 'Married', 'Never Married', 'Separated', 'Widowed']
Label encoder Gender - values: ['Female', 'Male', 'Nonbinary']
Label encoder Techie - values: ['No', 'Yes']
Label encoder Contract - values: ['Month-to-month', 'One year', 'Two Year']
Label encoder Port_modem - values: ['No', 'Yes']
Label encoder Tablet - values: ['No', 'Yes']
Label encoder InternetService - values: ['DSL', 'Fiber Optic', 'None']
Label encoder Phone - values: ['No', 'Yes']
Label encoder Multiple - values: ['No', 'Yes']
Label encoder OnlineSecurity - values: ['No', 'Yes']
Label encoder OnlineBackup - values: ['No', 'Yes']
Label encoder DeviceProtection - values: ['No', 'Yes']
Label encoder TechSupport - values: ['No', 'Yes']
Label encoder StreamingTV - values: ['No', 'Yes']
Label encoder StreamingMovies - values: ['No', 'Yes']
Label encoder PaperlessBilling - values: ['No', 'Yes']
Label encoder PaymentMethod - values: ['Bank Transfer(automatic)', 'Credit Card (automatic)', 'Electronic Check', 'Mailed Check']
```

```
In [12]: print(categorical_feature_names)
print(label_encoders)
```

```
['Area_Rural', 'Area_Suburban', 'Area_Urban', 'Marital_Divorced', 'Marital_Married', 'Marital_Never Married', 'Marital_Separated', 'Marital_Widowed', 'Gender_Female', 'Gender_Male', 'Gender_Nonbinary', 'Techie_No', 'Techie_Yes', 'Contract_Month-to-month', 'Contract_One year', 'Contract_Two Year', 'Port_modem_No', 'Port_modem_Yes', 'Tablet_No', 'Tablet_Yes', 'InternetService_DSL', 'InternetService_Fiber Optic', 'InternetService_None', 'Phone_No', 'Phone_Yes', 'Multiple_No', 'Multiple_Yes', 'OnlineSecurity_No', 'OnlineSecurity_Yes', 'OnlineBackup_No', 'OnlineBackup_Yes', 'DeviceProtection_No', 'DeviceProtection_Yes', 'TechSupport_No', 'TechSupport_Yes', 'StreamingTV_No', 'StreamingTV_Yes', 'StreamingMovies_No', 'StreamingMovies_Yes', 'PaperlessBilling_No', 'PaperlessBilling_Yes', 'PaymentMethod_Bank Transfer(automatic)', 'PaymentMethod_Credit Card (automatic)', 'PaymentMethod_Electronic Check', 'PaymentMethod_Mailed Check']
{'Area': LabelEncoder(), 'Marital': LabelEncoder(), 'Gender': LabelEncoder(), 'Techie': LabelEncoder(), 'Contract': LabelEncoder(), 'Port_modem': LabelEncoder(), 'Tablet': LabelEncoder(), 'InternetService': LabelEncoder(), 'Phone': LabelEncoder(), 'Multiple': LabelEncoder(), 'OnlineSecurity': LabelEncoder(), 'OnlineBackup': LabelEncoder(), 'DeviceProtection': LabelEncoder(), 'TechSupport': LabelEncoder(), 'StreamingTV': LabelEncoder(), 'StreamingMovies': LabelEncoder(), 'PaperlessBilling': LabelEncoder(), 'PaymentMethod': LabelEncoder()}
```

#### C4. Provide Clean Data

Provide a copy of the cleaned data set.

---

**TABLE 3-3.CLEAN DATA**



```
In [13]: #clean = raw.copy()
saveTable(data=clean, title='CLEAN', sect='C4',
          course='D209', caption='3 2')
```

	0	1	2	3
<b>Population</b>	38	10446	3735	13863
<b>Area</b>	2	2	2	1
<b>Children</b>	0	1	4	1
<b>Age</b>	68	27	50	48
<b>Income</b>	28561.99	21704.77	9609.57	18925.23
<b>Marital</b>	4	1	4	1
<b>Gender</b>	1	0	0	1
<b>Churn</b>	No	Yes	No	No
<b>Outage_sec_perweek</b>	7.978	11.699	10.753	14.914
<b>Email</b>	10	12	9	15
<b>Contacts</b>	0	0	0	2
<b>Yearly_equip_failure</b>	1	1	1	0
<b>Techie</b>	0	1	1	1
<b>Contract</b>	1	0	2	2
<b>Port_modem</b>	1	0	1	0
<b>Tablet</b>	1	1	0	0
<b>InternetService</b>	1	1	0	0
<b>Phone</b>	1	1	1	1
<b>Multiple</b>	0	1	1	0
<b>OnlineSecurity</b>	1	1	0	1
<b>OnlineBackup</b>	1	0	0	0
<b>DeviceProtection</b>	0	0	0	0
<b>TechSupport</b>	0	0	0	0
<b>StreamingTV</b>	0	1	0	1
<b>StreamingMovies</b>	1	1	1	0
<b>PaperlessBilling</b>	1	1	1	1
<b>PaymentMethod</b>	1	0	1	3
<b>Tenure</b>	6.796	1.157	15.754	17.087
<b>MonthlyCharge</b>	172.456	242.633	159.948	119.957
<b>Bandwidth_GB_Year</b>	904.536	800.983	2054.707	2164.579
<b>Item1</b>	5	3	4	4
<b>Item2</b>	5	4	4	4
<b>Item3</b>	5	3	2	4
<b>Item4</b>	3	3	4	2

	0	1	2	3
Item5	4	4	4	5
Item6	4	3	3	4
Item7	3	4	3	3
Item8	4	4	3	3

shape: (10000, 38)

Table saved to: TABLES/D209\_TASK1\_C4\_TAB\_3\_2\_CLEAN.CSV

## Part IV Analysis

D. Perform the data analysis and report on the results by doing the following:

### D1. Split the data into training and test data sets and provide the file(s).

```
In [14]: # define primary feature and target data
target= 'Churn' # target data
X = clean.loc[:, clean.columns != target]
y = clean.loc[:, clean.columns == target]
```

```
In [15]: # train test split raw data
tts = train_test_split(X, y, test_size=0.3, random_state=13)
(X_train, X_test, y_train, y_test)=tts
print('X_train: {}'.format(X_train.shape))
print('y_train: {}'.format(y_train.shape))
print('X_test: {}'.format(X_test.shape))
print('y_test: {}'.format(y_test.shape))
```

X\_train: (7000, 37)

y\_train: (7000, 1)

X\_test: (3000, 37)

y\_test: (3000, 1)

---

### TABLE 4-4. TRAINING DATA

```
In [16]: saveTable(data=X_train.merge(y_train,
    left_index=True, right_index=True),
    title='TRAIN', sect='D1',
    course='D209', caption='4 1')
```

	4847	9992	4621	5774
<b>Population</b>	48905	4261	10218	54601
<b>Area</b>	0	1	2	2
<b>Children</b>	0	1	1	4
<b>Age</b>	63	18	27	41
<b>Income</b>	27506.01	35876.21	43148.68	35600.06
<b>Marital</b>	1	0	1	3
<b>Gender</b>	1	1	0	0
<b>Outage_sec_perweek</b>	10.455	9.213	10.526	11.035
<b>Email</b>	14	15	9	17
<b>Contacts</b>	1	2	0	0
<b>Yearly_equip_failure</b>	1	0	0	0
<b>Techie</b>	0	1	1	0
<b>Contract</b>	1	2	2	0
<b>Port_modem</b>	0	0	0	0
<b>Tablet</b>	0	1	1	0
<b>InternetService</b>	2	1	1	0
<b>Phone</b>	1	1	1	1
<b>Multiple</b>	0	0	0	0
<b>OnlineSecurity</b>	0	0	0	0
<b>OnlineBackup</b>	0	0	0	1
<b>DeviceProtection</b>	0	0	0	1
<b>TechSupport</b>	1	0	1	1
<b>StreamingTV</b>	0	1	0	0
<b>StreamingMovies</b>	0	0	0	0
<b>PaperlessBilling</b>	1	1	1	1
<b>PaymentMethod</b>	2	2	2	3
<b>Tenure</b>	9.525	56.472	2.612	58.787
<b>MonthlyCharge</b>	92.488	137.439	124.964	139.983
<b>Bandwidth_GB_Year</b>	820.042	5001.371	320.107	5562.052
<b>Item1</b>	2	3	5	4
<b>Item2</b>	3	2	5	4
<b>Item3</b>	4	4	4	3
<b>Item4</b>	4	3	1	4
<b>Item5</b>	3	4	4	3

	4847	9992	4621	5774
Item6	2	3	5	3
Item7	3	4	4	3
Item8	4	3	3	3
Churn	No	No	No	No

shape: (7000, 38)

Table saved to: TABLES/D209\_TASK1\_D1\_TAB\_4\_1\_TRAIN.CSV

---

## TABLE 4-5.TEST DATA

```
In [17]: saveTable(data=X_test.merge(y_test,
    left_index=True, right_index=True),
    title='TEST', sect='D1',
    course='D209', caption='4 2')
```

	5952	1783	4811	145
<b>Population</b>	6399	32000	18951	5157
<b>Area</b>	2	1	2	0
<b>Children</b>	4	4	1	8
<b>Age</b>	59	83	68	75
<b>Income</b>	36343.28	40031.03	39053.35	36342.31
<b>Marital</b>	2	4	2	0
<b>Gender</b>	1	0	0	1
<b>Outage_sec_perweek</b>	8.923	14.75	11.206	15.409
<b>Email</b>	13	10	10	15
<b>Contacts</b>	0	2	7	2
<b>Yearly_equip_failure</b>	0	0	0	0
<b>Techie</b>	1	0	0	0
<b>Contract</b>	2	1	2	0
<b>Port_modem</b>	0	0	1	0
<b>Tablet</b>	0	1	1	0
<b>InternetService</b>	2	2	1	1
<b>Phone</b>	1	1	1	1
<b>Multiple</b>	0	0	0	1
<b>OnlineSecurity</b>	0	1	0	0
<b>OnlineBackup</b>	1	0	1	1
<b>DeviceProtection</b>	0	0	0	1
<b>TechSupport</b>	1	0	0	1
<b>StreamingTV</b>	0	0	1	1
<b>StreamingMovies</b>	0	1	1	0
<b>PaperlessBilling</b>	1	1	1	0
<b>PaymentMethod</b>	3	2	2	0
<b>Tenure</b>	56.633	2.851	5.664	2.733
<b>MonthlyCharge</b>	114.984	117.483	230.105	217.473
<b>Bandwidth_GB_Year</b>	4910.179	479.016	982.717	882.116
<b>Item1</b>	3	2	4	5
<b>Item2</b>	4	3	4	4
<b>Item3</b>	5	3	3	4
<b>Item4</b>	2	3	3	4
<b>Item5</b>	5	3	5	3

	5952	1783	4811	145
Item6	2	2	4	5
Item7	2	4	4	5
Item8	2	4	4	3
Churn	No	No	Yes	Yes

shape: (3000, 38)

Table saved to: TABLES/D209\_TASK1\_D1\_TAB\_4\_2\_TEST.CSV

## D2. Describe the analysis technique you used to appropriately analyze the data. Include screenshots of the intermediate calculations you performed.

The data for the model is the training data created above in section 'D1'. I am selected max\_depth=2 for the model in order to create a simple model. I want to be able to create a 2D decision boundary plot using a scatter plot.

```
In [49]: # create model
dt = DecisionTreeClassifier(max_depth=2, random_state=13)
dt.fit(X_train, y_train)
```

```
Out[49]: DecisionTreeClassifier(max_depth=2, random_state=13)
```

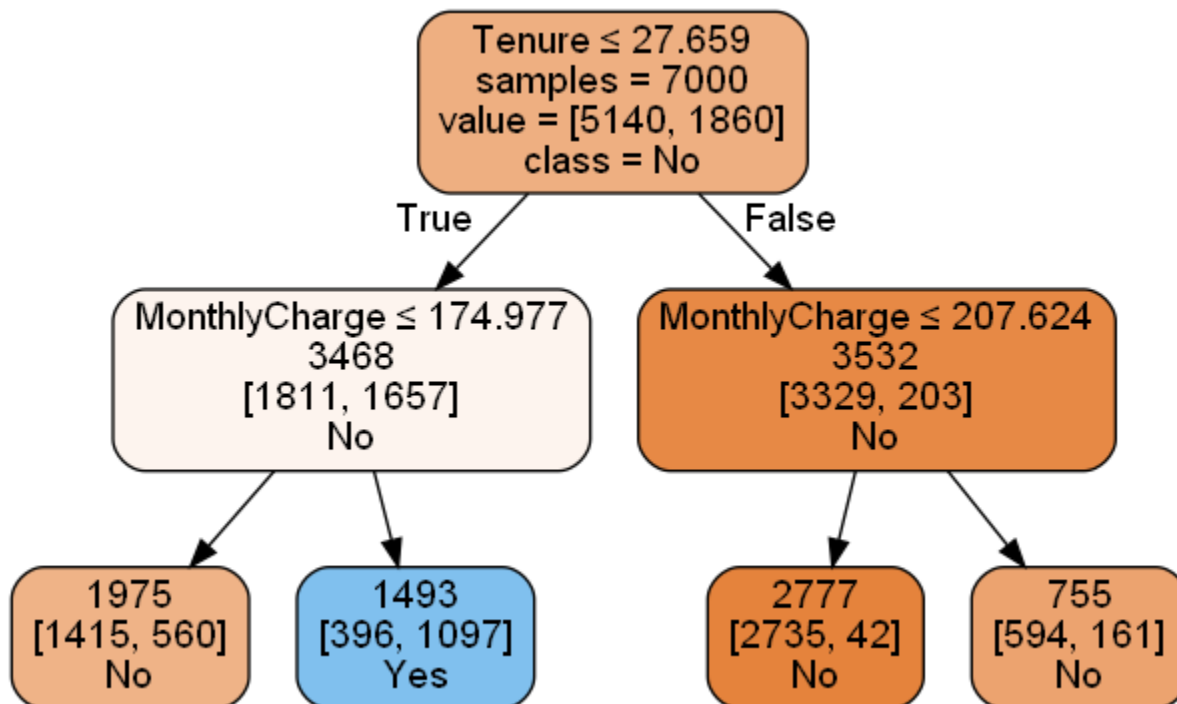
---

**FIGURE 4-1.CLASSIFICATION TREE (MAX\_DEPTH=2)**

```
In [50]: # visualize model
print('Target: [{}: {}]'.format(target, ', '.join(dt.classes_)))
plotDecisionTree(dt, feature_names=X_train.columns.to_list(),
                 class_names=dt.classes_)
```

Target: [Churn: No, Yes]

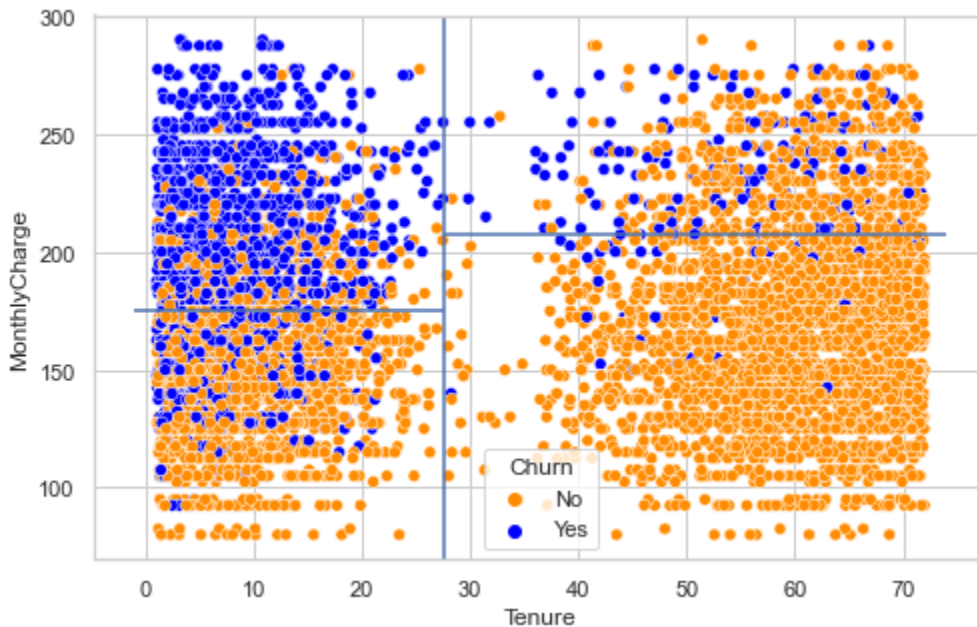
Out[50]:



**FIGURE 4-2.DECISION BOUNDARIES**

```
In [51]: # plot decision boundaries
fig, ax = plt.subplots()
fig.set_size_inches(8, 5)
sns.scatterplot(x="Tenure", y="MonthlyCharge",
                palette=['darkorange', 'blue'], hue=target,
                data=y_train.merge(X_train, left_index=True, right_index=True))
ax.axvline(x=27.659)
ax.hlines(y=174.977, xmin=-1, xmax=27.659) # horizontal line segment
ax.hlines(y=207.624, xmin=27.659, xmax=74) # horizontal line segment
```

```
Out[51]: <matplotlib.collections.LineCollection at 0x1e2acf11fd0>
```



The figure above shows the model from another slightly different vantage point. The scatter plot of the two key predictive features is broken down into four (4) section correlating to the four (4) terminal nodes of the decision tree. Each of the vertical and horizontal line segments are created using the numbers shown in the decision tree decision nodes.

### D3. Provide the code used to perform the classification analysis from part D2.

**Code.** All code and output is contained within this Jupyter notebook. The notebook file is called **D209\_2\_x.ipynb** and the associated PDF version is called **D209\_2\_x - Jupyter Notebook.pdf**.

## Part V: Data Summary and Implications



## E1. Explain the accuracy and the mean squared error (MSE) of your prediction model.

For the classification tree, the measurement of accuracy is found within the confusion matrix. Below are the confusion matrices for training data (showing 83.44% accuracy) and the test data (showing 83.57% accuracy). The accuracy is calculated by adding the number of True Positives and True Negatives on the diagonal, then divide that number by the total number of records.

```
In [45]: # training summary
classificationSummary(y_train, dt.predict(X_train))
```

Confusion Matrix (Accuracy 0.8344)

	Prediction	
Actual	0	1
0	4744	396
1	763	1097

```
In [46]: # test summary
classificationSummary(y_test, dt.predict(X_test))
```

Confusion Matrix (Accuracy 0.8357)

	Prediction	
Actual	0	1
0	2042	168
1	325	465

## E2. Discuss the results and implications of your prediction analysis.

The final prediction model uses only two (2) of the original features to predict the target class, and it does it correctly 83% of the time, pretty accurate for such a simple model. The simple model and graphical tree will be easy to explain to those individuals who will be using this prediction model.

## E3. Discuss one limitation of your data analysis.

**Limitations.** One limitation of the predictive analysis is that it correctly predicts target class only 83% of the time, that is, there is the 17% chance that the predicted class is incorrect.

## E4. Recommend a course of action for the real-world organizational situation from part A1 based on your results and implications discussed in part E2.

**Recommendations.** Recommend that company include policy to include the churn prediction for all customers at intervals during the lifespan of the customer with the company, for example, during the first month, then again a 1-3 monthly intervals. Also, the prediction can be used anytime the customer's MonthlyCharge changes. Lastly, the company can have a 'Customer Support' office reachout to customer's with Churned predicted, possibly lowering their current monthly charge by some calculated amount.

## Part VI: Demonstration

**F. Provide a Panopto video recording that includes a demonstration of the functionality of the code used for the analysis and a summary of the programming environment.**

**Video.** Panopto video was created and is located at: [https://wgu.edu\\_\(https://wgu.edu\)\\_](https://wgu.edu_(https://wgu.edu)_)

**G. Record the web sources used to acquire data or segments of third-party code to support the analysis. Ensure the web sources are reliable.**

**Configure Scrollbars.** Disable scrollbars in notebook.

```
In [ ]: %%javascript
        IPython.OutputArea.prototype._should_scroll = function(lines) {
            return false;
        }
```

**Disable Auto Scroll.** Disable automatically scroll to bottom.

```
In [ ]: %%javascript
        require("notebook/js/notebook").Notebook.prototype.scroll_to_bottom = function () {}
```

**Toggle Notebook Warnings.** Use the following code to toggle warning messages in the notebook. Another piece of code courtesy of stackoverflow (2021).

<https://stackoverflow.com/questions/9031783/hide-all-warnings-in-ipython>  
(<https://stackoverflow.com/questions/9031783/hide-all-warnings-in-ipython>)

```
In [ ]: from IPython.display import HTML
HTML('''<script>
code_show_err=false;
function code_toggle_err() {
  if (code_show_err){
    $('div.output_stderr').hide();
  } else {
    $('div.output_stderr').show();
  }
  code_show_err = !code_show_err
}
$( document ).ready(code_toggle_err);
</script>
To toggle on/off output_stderr, click <a href="javascript:code_toggle_err()">here</a>.''' )
```

**Terminal List Files.** List all of the files from the current working directory.

Ref: (1) [Fessel, K. \(2021\). How to save a matplotlib figure and fix text cutting off || Matplotlib Tips \(https://www.youtube.com/watch?v=C8MT-A7Mvk4&ab\\_channel=KimberlyFessel\)](https://www.youtube.com/watch?v=C8MT-A7Mvk4&ab_channel=KimberlyFessel) Retrieved from [https://www.youtube.com/watch?v=C8MT-A7Mvk4&ab\\_channel=KimberlyFessel](https://www.youtube.com/watch?v=C8MT-A7Mvk4&ab_channel=KimberlyFessel)

```
In [ ]: !ls
```

```
In [ ]: !du -h *.*
```

**List Installed Packages.** List of all installed PIP packages and the versions.

Ref: (1) [https://pip.pypa.io/en/stable/cli/pip\\_list/](https://pip.pypa.io/en/stable/cli/pip_list/)

```
!pip list
```

**Update Package.** Update a specific package within notebook.

Ref: (1) <https://stackoverflow.com/questions/54453219/why-can-i-see-pip-list-sklearn-but-not-in-jupyter-when-i-run-a-code>

```
In [ ]: !python -m pip install -U scikit-learn
```

**Merget Two Dataframes.** Code to merge two dataframes.

Ref: (1) <https://stackoverflow.com/questions/26265819/how-to-merge-a-series-and-dataframe>

```
In [ ]: # merge X and y back together, for example
d = X.merge(y, left_index=True, right_index=True)
display(d.head())
```

**List.index() Function.** The .index() method returns the index of the specified element in the list.

Ref: (1) <https://www.programiz.com/python-programming/methods/list/index>

```
In [ ]: animals = ['cat', 'dog', 'rabbit', 'horse']
# get the index of 'dog'
index = animals.index('dog')
print(index)
```

**Row Index Names in Pandas.** Code to get rows/index names in a Pandas dataframe.

Ref: (1) <https://www.geeksforgeeks.org/how-to-get-rows-index-names-in-pandas-dataframe/>

```
In [ ]: # making data frame
data = cleanData

# calling head() method
# storing in new variable
data_top = data.head()

# iterating the columns
for row in data_top.index:
    print(row, end = " ")
```

**Tutorial Python Subplots.** Tutorial: Python Subplots

Ref: (1) <https://www.kaggle.com/asimislam/tutorial-python-subplots>

```
In [ ]: # Categorical Data
heart_CAT = ['Churn']

# Categorical Data
a = 2 # number of rows
b = 3 # number of columns
c = 1 # initialize plot counter

fig = plt.figure(figsize=(14,10))

for i in heart_CAT:
    plt.subplot(a, b, c)
    plt.title('{}, subplot: {}'.format(i, a, b, c))
    plt.xlabel(i)
    sns.countplot(x=i, data=cleanData, palette='hls')
    c = c + 1

plt.show()
```

**PASS FIG TO CUSTOM PLOT FUNCTION.** A great way to do this is to pass a figure object to your code and have your function add an axis then return the updated figure.

Ref: (1) <https://stackoverflow.com/questions/43925337/matplotlib-returning-a-plot-object>

```

In [ ]: def plot_hist_overlay(feature, fig, p, bins=8):

    # data
    df_yes = cleanData[cleanData.Churn==True][feature]
    df_no = cleanData[cleanData.Churn==False][feature]

    # plot stacked hist
    ax = f.add_subplot() # here is where you add the subplot to f
    plt.hist([df_yes,df_no], bins=bins, stacked=True)

    # add title
    plt.title(feature + ' grouped by target', size=16)

    # tick marks
    ax.set_xticks([])
    #ax.set_yticks([]) # use default

    # add axis labels
    plt.xlabel(feature)
    plt.ylabel('# Churn')

    # add Legend
    ax.legend(['Churn - Yes', 'Churn - No'])

    return(f)

target = 'Churn'
features = ['MonthlyCharge', 'Tenure']
bins = 6
for idx,fea in enumerate(features):
    fig_size = (6,5)
    f = plt.figure(figsize=fig_size)
    f = plot_hist_overlay(fea, fig=f, p=idx+1, bins=bins)
    file = getFilename(fea, 'z1', 'fig 9 ' + str(idx+1)) # getFilename using helper
    plt.gcf().text(0.1, 0, file, fontsize=14)

    # data table
    b = pd.cut(cleanData[fea], bins=bins) # create bins (b) of numeric feature
    dt = pd.crosstab(cleanData[target], b)
    plt.gcf().text(0.1, -.4, dt.T.to_string(), fontsize=14)
    #print(dt.T)

    f.savefig(file, dpi=150, bbox_inches='tight')

#f = plot_hist_overlay('MonthlyCharge', fig=f, p=3)
#f = plot_hist_overlay('Tenure', fig=f, p=2)

```

### Enabling Jupyter Notebook extensions.

Ref: (1) <https://tjth.jupyter.org/en/latest/howto/admin/enable-extensions.html>

```

pip install jupyter_contrib_nbextensions
jupyter contrib nbextension install --sys-prefix
jupyter nbextension enable scratchpad/main --sys-prefix
jupyter nbextension list

```

## How to Use HTML to Open a Link in a New Tab.

Ref: (1)<https://www.freecodecamp.org/news/how-to-use-html-to-open-link-in-new-tab/>

```
<p>Check out <a href="https://www.freecodecamp.org/" target="_blank" rel="noopener norereferrer">freeCodeCamp</a>.</p>
```

**CSS Tutorial.** This is a great resource for CSS code with many examples.

Ref: (1)<https://www.w3schools.com/css/default.asp>

**HTML Tutorial.** This is a great resource for HTML code with many examples.

Ref: (1)<https://www.w3schools.com/html/default.asp>

**Inline Styles in HTML.** Usually, CSS is written in a separate CSS file (with file extension .css) or in a 'style' tag inside of the 'head' tag, but there is a third place which is also valid. The third place you can write CSS is inside of an HTML tag, using the style attribute. When CSS is written using the style attribute, it's called an "inline style". In general, this is not considered a best practice. However, there are times when inline styles are the right (or only) choice.

Ref: (1) <https://www.codecademy.com/articles/html-inline-styles>

## H. Acknowledge sources, using in-text citations and references, for content that is quoted, paraphrased, or summarized.

**Deitel, P. +** (2020). Intro to Python for Computer Science and Data Science: Learning to Program with AI, Big Data and the Cloud

**Geron, A.** (2019). Hands-on Machine Learning with Scikit-Learn, Keras and TensorFlow: Concepts, Tools and Techniques to Build Intelligent Systems

**Larose, C. +** (2019). Data Science Using Python and R

**Rite, S.** (2018). Demystifying 'Confusion Matrix' Confusion  
<https://towardsdatascience.com/demystifying-confusion-matrix-confusion-9e82201592fd>

**Robinson, S.** (2021). K-Nearest Neighbors Algorithm in Python and Scikit-Learn  
<https://stackabuse.com/k-nearest-neighbors-algorithm-in-python-and-scikit-learn/>

**Sharma, A.** (2021) . K-Nearest Neighbors (KNN) on Customer Churn Data  
<https://medium.com/data-science-on-customer-churn-data/k-nearest-neighbors-knn-on-customer-churn-data-40e9b2bb9266>

**Shmueli, G. +** (2020) . Data Mining for Business Analytics: Concepts, Techniques, and Application in Python

In [ ]: