

Final Project Submission

- Student name: WAFULA SIMIYU
- Student pace: self paced / part time / full time PART-TIME
- Instructor name: SAM G/SAM KARU/WINNIE A
- Blog post URL:https://github.com/wafulaandree/Churn-in-Telecom-s-dataset.git

Customer Churn Prediction for SyriaTel

Introduction

For SyriaTel, reducing customer churn is vital for maintaining a stable customer base and maximizing revenue. Accurate predictions of churn help the company take proactive measures to improve customer satisfaction and loyalty, thus improving long-term profitability and reducing operational costs associated with acquiring new customers.

Objective:

The primary objective of this project is to predict customer churn for SyriaTel, a telecommunications company. Identifying customers likely to leave enables the company to implement retention strategies to prevent attrition and minimize revenue loss.

Data Understanding

Data Source and Properties: The dataset consists of 3,333 entries with 21 columns, covering various aspects of customer data, including demographics, service plans, usage metrics, and the target variable indicating churn.

Customer Demographics:

state: State where the customer resides. area code: Phone area code of the customer. Service Plans:

international plan: Indicates whether the customer has an international plan. voice mail plan: Indicates whether the customer has a voice mail plan. Usage Metrics:

Includes metrics such as total day minutes, total eve minutes, total night minutes, and total intl minutes, along with charges for these categories and the number of calls.

churn: Binary indicator (1 for churned, 0 for not churned). Relevance to Problem:

This data is directly related to predicting customer churn, including both features that may influence churn (e.g., service usage, plan types) and the target variable indicating whether a customer has churned. Analyzing these features helps in identifying patterns that predict customer departure.

```
In [82]:
          # Importing necessary libraries
          import pandas as pd
          import numpy as np
          import matplotlib.pyplot as plt
          import seaborn as sns
          from sklearn.model_selection import train_test_split
          from sklearn.linear_model import LogisticRegression
          from sklearn.metrics import classification_report, confusion_matrix
          from sklearn.preprocessing import StandardScaler, OneHotEncoder
          from sklearn.compose import ColumnTransformer
          from sklearn.pipeline import Pipeline
          from sklearn.impute import SimpleImputer
          from sklearn.ensemble import RandomForestClassifier
          from sklearn.model_selection import GridSearchCV
          from sklearn.linear_model import LogisticRegression
          from sklearn.metrics import classification report, confusion matrix
          from sklearn.model_selection import train_test_split
          from sklearn.compose import ColumnTransformer
          from sklearn.preprocessing import OneHotEncoder,OrdinalEncoder,StandardScaler
          from sklearn.compose import ColumnTransformer
          from sklearn.tree import DecisionTreeClassifier
          from sklearn import tree
          from IPython.display import Image
          from sklearn.tree import export_graphviz
```

Loading the dataset

```
churn_data = pd.read_csv('bigml_59c28831336c6604c800002a.csv')

# Display a few rows of the dataset
(churn_data .head())
```

Out[5]:		state	account length	area code	phone number	international plan	voice mail plan	number vmail messages	day	day	
	0	KS	128	415	382- 4657	no	yes	25	265.1	110	2
	1	ОН	107	415	371- 7191	no	yes	26	161.6	123	Ź
	2	NJ	137	415	358- 1921	no	no	0	243.4	114	2

```
3/5-
             OH
                            408
                                                  yes
                                                         no
                                                                          299.4
                                    9999
                                    330-
             OK
                      75
                           415
                                                                    0
                                                                          166.7
                                                                                        2
                                                                                 113
                                                         no
                                                  yes
                                    6626
        5 rows × 21 columns
In [6]:
         # dataset shape
         print(f"{'shape of the dataset'.title()} :- {churn_data.shape}")
       Shape Of The Dataset :- (3333, 21)
In [7]:
         # missing values
         print(f"\n {'Number of null values in every column'.title()} \n {churn_data.i
                                                                                      Number Of Null Values In Every Column
        state
                                   0
       account length
                                  0
       area code
                                  0
       phone number
                                  0
       international plan
       voice mail plan
       number vmail messages
                                  0
       total day minutes
       total day calls
       total day charge
       total eve minutes
       total eve calls
                                  0
       total eve charge
       total night minutes
                                  0
       total night calls
       total night charge
       total intl minutes
                                  0
       total intl calls
       total intl charge
                                  0
       customer service calls
                                  0
       churn
       dtype: int64
In [8]:
         # duplicate values
         print(f"\n {'number of duplicate values'.title()} :- {len(churn_data.loc[chur
        Number Of Duplicate Values :- 0
In [9]:
         # target value count
         print(f"\n {'count of each value of target column'.title()} \n {churn_data.ch
                                                                                      Count Of Each Value Of Target Column
        churn
```

Falca

```
True
                483
       Name: count, dtype: int64
In [10]:
         # information about dataset
         print(f"{'dataset info'.title()} \n ")
         churn data.info()
       Dataset Info
       <class 'pandas.core.frame.DataFrame'>
       RangeIndex: 3333 entries, 0 to 3332
       Data columns (total 21 columns):
            Column
                                  Non-Null Count Dtype
       ---
           -----
                                  -----
        0
            state
                                  3333 non-null
                                                 object
           account length
                                  3333 non-null int64
        1
                                                 int64
        2
           area code
                                  3333 non-null
        3
           phone number
                                  3333 non-null object
           international plan
                                 3333 non-null object
           voice mail plan
                                3333 non-null object
        6
           number vmail messages 3333 non-null int64
        7
           total day minutes
                                 3333 non-null float64
           total day calls
                                  3333 non-null int64
        8
                                  3333 non-null float64
           total day charge
        10 total eve minutes
                                 3333 non-null float64
        11 total eve calls
                                  3333 non-null int64
        12 total eve charge
                                  3333 non-null float64
        13 total night minutes
                                 3333 non-null float64
        14 total night calls
                                 3333 non-null int64
        15 total night charge
                                  3333 non-null float64
        16 total intl minutes
                                 3333 non-null float64
        17 total intl calls
                                  3333 non-null int64
                                 3333 non-null float64
        18 total intl charge
        19 customer service calls 3333 non-null
                                                 int64
                                  3333 non-null
       dtypes: bool(1), float64(8), int64(8), object(4)
       memory usage: 524.2+ KB
         df = churn_data
         df.shape
```

```
In [12]:
```

Out[12]: (3333, 21)

Voice Mail Plan and Number of Voicemail Messages:

Observation: The voice mail plan and number vmail messages columns appear to be closely related. Customers with a voice mail plan (voice mail plan = "yes") typically have a number vmail messages greater than 0. Implication: Since the presence of a voice mail plan correlates with having voicemail messages, including both features might introduce redundancy. It could be beneficial to retain only one of these columns to avoid multicollinearity and simplify the model. Day, Evening, Night, and International Charges and Minutes:

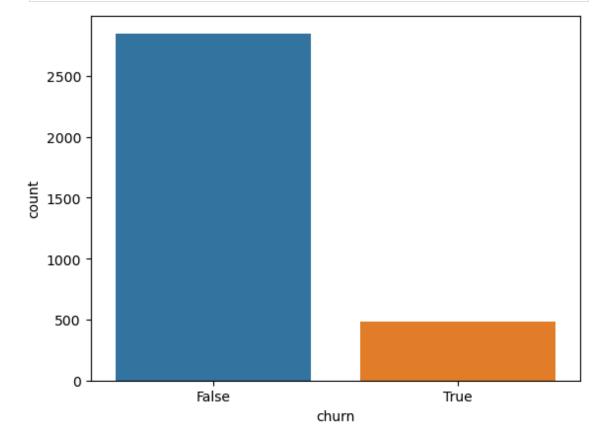
Observation: The columns total day charge, total eve charge, total night charge, and total intl charge are directly related to their respective total day minutes, total eve

minutes, total night minutes, and total intl minutes. Each charge column is calculated by multiplying the corresponding minutes column by a per-minute rate. Implication: Given this direct relationship, these pairs of columns are likely to be highly correlated. For modeling purposes, including both the minutes and charges might lead to multicollinearity. It might be more effective to use only one type of feature (e.g., minutes or charges) to prevent redundancy and improve model performance. Phone Number, State, and Area Code:

Observation: Columns like phone number, state, and area code appear to function as identifiers rather than providing substantive quantitative or qualitative information relevant to predicting churn. Phone number: Unique to each customer and does not contribute to predicting churn. State: While it might offer geographical information, it could also be a categorical variable that may not directly impact churn prediction if not processed correctly. Area code: Likely to be a categorical identifier that may not have a strong predictive power after encoding. Implication: These columns may not provide useful information for churn prediction and could be considered for removal to streamline the dataset and focus on features with predictive value.

Data Visualization

```
In [11]:
    sns.countplot(x ='churn', data = churn_data)
    plt.show()
```

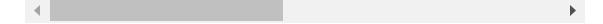


In [13]: | # New categorical feature

Out[13]:		state	account length		phone number	international plan	voice mail plan	number vmail messages	day	total day calls	c
	0	KS	128	415	382- 4657	no	yes	25	265.1	110	
	1	ОН	107	415	371- 7191	no	yes	26	161.6	123	
	2	NJ	137	415	358- 1921	no	no	0	243.4	114	
	3	ОН	84	408	375- 9999	yes	no	0	299.4	71	
	4	OK	75	415	330- 6626	yes	no	0	166.7	113	
	5	AL	118	510	391- 8027	yes	no	0	223.4	98	
	6	MA	121	510	355- 9993	no	yes	24	218.2	88	
	7	МО	147	415	329- 9001	yes	no	0	157.0	79	
	8	LA	117	408	335- 4719	no	no	0	184.5	97	
	9	WV	141	415	330- 8173	yes	yes	37	258.6	84	
	10	IN	65	415	329- 6603	no	no	0	129.1	137	
	11	RI	74	415	344- 9403	no	no	0	187.7	127	
	12	IA	168	408	363- 1107	no	no	0	128.8	96	
	13	MT	95	510	394- 8006	no	no	0	156.6	88	
	14	IA	62	415	366- 9238	no	no	0	120.7	70	
	15	NY	161	415	351- 7269	no	no	0	332.9	67	

16	ID	85	408	350- 8884	no	yes	27	196.4	139
17	VT	93	510	386- 2923	no	no	0	190.7	114
18	VA	76	510	356- 2992	no	yes	33	189.7	66
19	TX	73	415	373- 2782	no	no	0	224.4	90

20 rows × 22 columns



Total Day Minutes and Total Day Charge:

Observation: There is a perfect correlation (correlation of 1) between total day minutes and total day charge. This indicates that total day charge is directly derived from total day minutes using a fixed per-minute rate. Recommendation: To avoid redundancy, you can retain only one of these features. Either total day minutes or total day charge can be used for modeling, as they provide the same information about day usage.

Total Evening Minutes and Total Evening Charge:

Observation: Similarly, total eve minutes and total eve charge have a perfect correlation (correlation of 1). Total eve charge is calculated based on total eve minutes using a constant rate per minute. Recommendation: Since these features are perfectly correlated, you can choose to keep either total eve minutes or total eve charge, as including both does not add additional predictive value.

Total Night Minutes and Total Night Charge:

Observation: There is also a perfect correlation (correlation of 1) between total night minutes and total night charge. Total night charge is derived from total night minutes by applying a per-minute rate. Recommendation: Given this perfect correlation, you should select one of these features for inclusion in your model. Both features are redundant, so choosing either total night minutes or total night charge will suffice.

Total International Minutes and Total International Charge:

Observation: total intl minutes and total intl charge exhibit a perfect correlation (correlation of 1). This implies that total intl charge is directly calculated from total intl minutes with a fixed per-minute charge. Recommendation: As these features are perfectly correlated, you can keep either total intl minutes or total intl charge. Including both features would introduce unnecessary redundancy.

Building and evaluation of the model

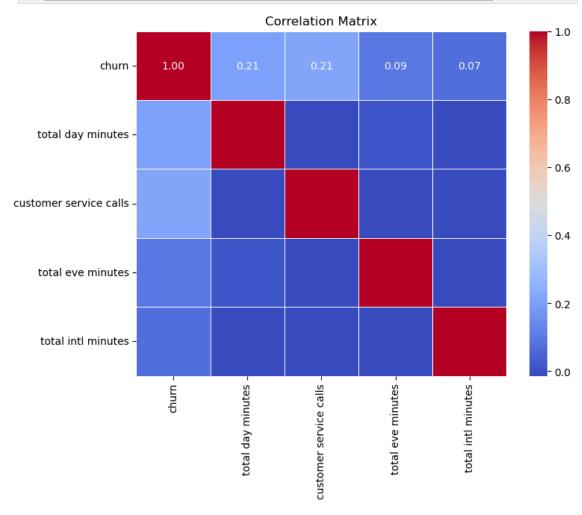
Feature Selection

ANOVA is used for comparing the distribution of a numeric variable in two or more groups II Ho = Null Hypothesis = the distribution of the variable in multiple groups is uniform Ha = Alternate Hypothesis = the distribution of the variable in multiple groups in different

we analyse the pvalue, lets say for confidence interval of 95%, significance level = 5%

if pvalue > 0.05 = accept the null hypothesis and the feature is NOT important if pvalue < 0.05 = reject the null hypothesis and the feature is important

```
In [40]:
           df.columns
Out[40]: Index(['state', 'account length', 'area code', 'phone number',
                  'international plan', 'voice mail plan', 'number vmail messages',
                  'total day minutes', 'total day calls', 'total day charge', 'total eve minutes', 'total eve calls', 'total eve charge',
                  'total night minutes', 'total night calls', 'total night charge',
                  'total intl minutes', 'total intl calls', 'total intl charge',
                  'customer service calls', 'churn', 'vmail_messages'],
                 dtype='object')
In [41]:
           numerics =['account length', 'number vmail messages',
                   'total day minutes', 'total day calls', 'total day charge',
                   'total eve minutes', 'total eve calls', 'total eve charge', 'total night minutes', 'total night calls', 'total night charge',
                   'total intl minutes', 'total intl calls', 'total intl charge',
                   'customer service calls']
           xnum = df[numerics]
           y = df['churn']
           from sklearn.feature_selection import f_classif
           fval,pval = f_classif(xnum,y)
           for i in range(len(numerics)):print(numerics[i],pval[i])
         account length 0.33976000705720666
         number vmail messages 2.1175218402696038e-07
        total day minutes 5.300278227509361e-33
        total day calls 0.28670102402211844
        total day charge 5.30060595239102e-33
        total eve minutes 8.011338561256927e-08
        total eve calls 0.5941305829720491
        total eve charge 8.036524227754477e-08
        total night minutes 0.04046648463758881
        total night calls 0.7230277872081609
        total night charge 0.040451218769160205
        total intl minutes 8.05731126549437e-05
        total intl calls 0.002274701409850077
        total intl charge 8.018753583047257e-05
         customer service calls 3.900360240185746e-34
In [84]:
           # Select features for correlation analysis (including 'churn')
```



Preprocesing

```
In [ ]:
       def col unique values(col name):
         ## input : category variables
          print(f"Unique Values :- \n {churn_data[col_name].unique()}")
          print(f"Number of Unique values :- {churn_data[col_name].nunique()}\n\n")
        total_col_names = churn_data.columns
In [ ]:
        #columns
        total_col_names = churn_data.columns
In [ ]:
       # find numeric columns (int & float, bool)
        num_cols = churn_data._get_numeric_data().columns
        # getting category columns
        cat_col_names = list(set(total_col_names) - set(num_cols))
        for col_name in cat_col_names:
          ## check unique values of every category column
          col_unique_values(col_name)
      ******* Col Name : international plan *********
      Unique Values :-
       ['no' 'yes']
      Number of Unique values :- 2
      ********* Col Name : vmail messages **********
      Unique Values :-
       ['Normal Users', 'No VM plan', 'High Frequency users']
      Categories (3, object): ['No VM plan' < 'Normal Users' < 'High Frequency user
      Number of Unique values :- 3
      ******** Col Name : state *********
      Unique Values :-
       ['KS' 'OH' 'NJ' 'OK' 'AL' 'MO' 'LA' 'WV' 'IN' 'RI' 'IA' 'MY' 'NY'
       'ID' 'VT' 'VA' 'TX' 'FL' 'CO' 'AZ' 'SC' 'NE' 'WY' 'HI' 'IL' 'NH' 'GA'
```

```
In [ ]: # Apply label encoding operation on category columns
    def label_encoding(col_name):
        le = LabelEncoder()
        churn_data[col_name] = le.fit_transform(churn_data[col_name])
    for col_name in cat_col_names:
        label_encoding(col_name)
```

	state	account length		international plan	voice mail plan	number vmail messages	total day minutes	day	total day charge	to € minu
0	KS	128	415	no	yes	25	265.1	110	45.07	19
1	ОН	107	415	no	yes	26	161.6	123	27.47	19
2	NJ	137	415	no	no	0	243.4	114	41.38	12
3	ОН	84	408	yes	no	0	299.4	71	50.90	6
4	OK	75	415	yes	no	0	166.7	113	28.34	14

5 rows × 21 columns

```
def model_building(model_name):
    model = model_name
    model.fit(X_train, y_train)
    y_prediction = model.predict(X_test)
    print(classification_report(y_test, y_prediction))
```

['account length', 'number vmail messages', 'total day minutes', 'total day cal ls', 'total day charge', 'total eve minutes', 'total eve calls', 'total eve cha rge', 'total night minutes', 'total night charge', 'total intl minutes', 'total intl calls', 'total intl charge', 'customer service call s', 'state_AK', 'state_AL', 'state_AR', 'state_AZ', 'state_CA', 'state_CO', 'st ate_CT', 'state_DC', 'state_DE', 'state_FL', 'state_GA', 'state_HI', 'state_I A', 'state_ID', 'state_IL', 'state_IN', 'state_KS', 'state_KY', 'state_LA', 'st ate_MA', 'state_MD', 'state_ME', 'state_MI', 'state_MN', 'state_MO', 'state_M S', 'state_MT', 'state_NC', 'state_ND', 'state_NE', 'state_NH', 'state_NJ', 'st ate_NM', 'state_NV', 'state_NY', 'state_OH', 'state_OK', 'state_OR', 'state_P A', 'state_RI', 'state_SC', 'state_SD', 'state_TN', 'state_TX', 'state_UT', 'st ate_VA', 'state_VT', 'state_WA', 'state_WI', 'state_WV', 'international plan_no', 'international plan_no', 'voice mail plan_no', 'voice mail plan_yes']

```
In [44]: x.head()
```

Out[44]:		international plan	vmail_messages	total day minutes	total eve minutes	total night minutes	total intl minutes	customer service calls	
	0	no	Normal Users	265.1	197.4	244.7	10.0	1	
	1	no	Normal Users	161.6	195.5	254.4	13.7	1	
	2	no	No VM plan	243.4	121.2	162.6	12.2	0	
	3	yes	No VM plan	299.4	61.9	196.9	6.6	2	
	4	yes	No VM plan	166.7	148.3	186.9	10.1	3	

Hotencoding categorical features

```
In [47]:
          preprocessor = ColumnTransformer([('ohe',OneHotEncoder(),[1]),
                                            ('ode',OrdinalEncoder(),[0]),
                                             ('sc',StandardScaler(),[2,3,4,5,6])],remaind
In [48]:
          x_new = preprocessor.fit_transform(x)
          pd.DataFrame(x_new).head()
Out[48]:
                  1
                       2
                           3
                                     4
                                               5
                                                         6
                                                                    7
                                                                              8
```

```
0 0.0 0.0 1.0
               0.0
                    1.566767 -0.070610 0.866743 -0.085008 -0.427932
  0.0
       0.0
          1.0
               0.0
                    -0.333738 -0.108080
                                        1.058571
                                                   1.240482 -0.427932
      1.0 0.0
               0.0
                     1.168304 -1.573383 -0.756869
 0.0
                                                   0.703121 -1.188218
 0.0
      1.0 0.0
               1.0
                     2.196596 -2.742865 -0.078551 -1.303026
                                                              0.332354
4 0.0 1.0 0.0 1.0 -0.240090 -1.038932 -0.276311 -0.049184
                                                              1.092641
```

```
In [49]:
# train test split
xtrain,xtest,ytrain,ytest = train_test_split(x_new,y,test_size=0.2,random_staprint(x.shape)
print(xtrain.shape)
print(ytest.shape)
print(ytrain.shape)
print(ytest.shape)

# (3333, 7)
(2666, 9)
(667, 9)
(3333,)
(2666,)
(667,)
```

Logistic regression

```
In [51]: model = LogisticRegression(class_weight='balanced')
    model.fit(xtrain,ytrain)
```

Out[51]: LogisticRegression(class_weight='balanced')

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [52]: #Performance Analysis
    from sklearn import metrics
    ypred = model.predict(xtest)
    print("Accuracy : ",metrics.accuracy_score(ytest,ypred))
    print("Recall : ",metrics.recall_score(ytest,ypred))
    print("F1 score : ",metrics.f1_score(ytest,ypred))
    print("Precision : ",metrics.precision_score(ytest,ypred))

Accuracy : 0.775112443778111
    Recall : 0.7934782608695652
    F1 score : 0.4932432432432
    Precision : 0.35784313725490197
In [53]: # performance analysis on train data
    ypred2 = model.predict(xtrain)
```

```
print("Accuracy : ",metrics.accuracy_score(ytrain,ypred2))
print("Recall : ",metrics.recall_score(ytrain,ypred2))
print("F1 score : ",metrics.f1_score(ytrain,ypred2))
print("Precision : ",metrics.precision_score(ytrain,ypred2))
```

Accuracy: 0.764066016504126
Recall: 0.7570332480818415
F1 score: 0.4848484848484849
Precision: 0.3566265060240964

Classifiers

Decision Tree

```
Out[54]:
             0.0 0.0
                       1.0
                            0.0 265.1 197.4 244.7 10.0 1.0
               0.0 0.0 1.0 0.0 161.6 195.5 254.4 13.7 1.0
                   1.0 0.0
                            0.0
                                243.4
                                       121.2
                                             162.6 12.2
                                299.4
                                             196.9
               0.0 1.0 0.0
                            1.0
                                        61.9
                                                    6.6
                                                         2.0
                0.0 1.0 0.0
                            1.0 166.7
                                       148.3
                                             186.9
                                                   10.1
               0.0
                   0.0 1.0 0.0 156.2 215.5 279.1
                                                         2.0
          3329 0.0 1.0 0.0 0.0 231.1 153.4 191.3
                                                    9.6
                                                        3.0
          3330 0.0 1.0 0.0 0.0 180.8 288.8 191.9 14.1 2.0
          3331 0.0 1.0 0.0 1.0 213.8 159.6 139.2
                                                    5.0 2.0
          3332 0.0 0.0 1.0 0.0 234.4 265.9 241.4 13.7 0.0
```

3333 rows × 9 columns

```
# train test split
from sklearn.model_selection import train_test_split
xtrain,xtest,ytrain,ytest = train_test_split(x_new,y,test_size=0.2,random_sta
print(x.shape)
print(xtrain.shape)
print(xtest.shape)
```

```
print(y.snape)
print(ytrain.shape)
print(ytest.shape)

(3333, 7)
(2666, 9)
(667, 9)
(3333,)
(2666,)
(2666,)
(667,)

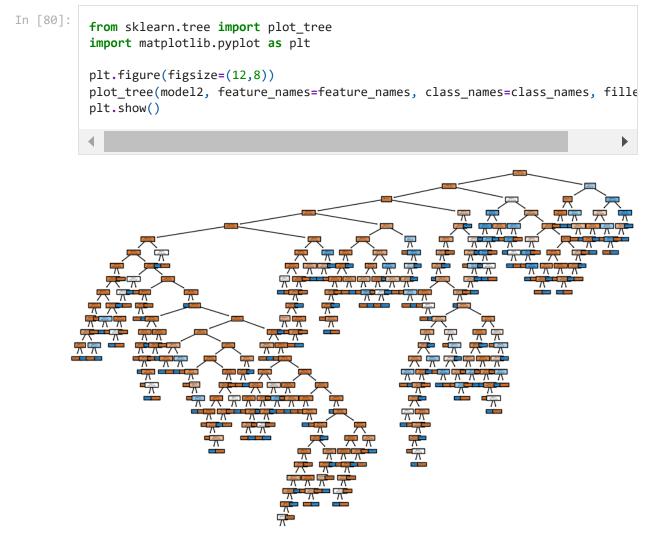
In [72]: # decision tree
model2 = DecisionTreeClassifier(random_state=5, class_weight={0:0.5,1:0.5})
model2.fit(xtrain,ytrain)
```

Out[72]: DecisionTreeClassifier(class_weight={0: 0.5, 1: 0.5}, random_state=5)
In a Jupyter environment, please rerun this cell to show the HTML representation

or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

Visualizing the decision Tree Model



```
In [81]: # performance analysis
    ypred2 = model2.predict(xtest)
    print("Accuracy : ",metrics.accuracy_score(ytest,ypred2))
    print("Recall : ",metrics.recall_score(ytest,ypred2))
    print("F1 score : ",metrics.f1_score(ytest,ypred2))
    print("Precision : ",metrics.precision_score(ytest,ypred2))
```

Accuracy: 0.8995502248875562
Recall: 0.7391304347826086
F1 score: 0.6699507389162561
Precision: 0.6126126126126126

Reccommendations

1. Enhance Customer Experience

Improve Customer Service: Since high numbers of customer service calls are a feature in your dataset, consider investing in training and resources to improve the quality of customer service. Ensure that customer service representatives are well-equipped to resolve issues efficiently and effectively.

Personalize Interactions: Use customer data to personalize interactions and offers. Personalized experiences can improve customer satisfaction and reduce churn. For example, tailor communication and promotions based on customer usage patterns and preferences.

2. Targeted Retention Campaigns

Identify At-Risk Customers: Use the churn prediction model to identify customers who are at high risk of churning. Develop targeted retention campaigns specifically designed for these customers.

Special Offers and Discounts: Offer incentives such as discounts on service plans, special promotions, or loyalty rewards to at-risk customers. These incentives can encourage them to stay with SyriaTel.

Enhanced Service Plans: For customers who frequently use customer service or have high usage, offer service plan upgrades or additional features that enhance their experience.

3. Improve Service Quality and Offerings

Analyze Usage Patterns: Examine the usage patterns of customers who churn and compare them with those who stay. This analysis can help identify gaps in service quality or features that are missing or underperforming.

Upgrade Service Plans: Regularly review and upgrade service plans based on customer

Churn-in-Telecom-s-dataset/index.ipynb at master \cdot wafulaandree/Churn-in-Telecom-s-dataset reedback and usage trends. Introduce new reatures or improve existing ones to keep customers engaged.

Reduce Service Disruptions: Ensure that there are minimal service disruptions and high-quality network coverage. Frequent disruptions can lead to customer dissatisfaction and churn.

4. Proactive Engagement

Regular Check-ins: Implement regular check-ins with customers to understand their