

In [113]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.preprocessing import LabelEncoder, OneHotEncoder
from pylab import rcParams

pd.set\_option('display.max\_columns', None)
missing\_values = ["n/a", "na", "-", "", " "]
df = pd.read\_csv("Telco\_customer\_churn.csv", na\_values=missing\_values)

In this section we will visualize our dataframe, explore different features and make preliminary conclusions. We will decide the first steps of Data wrangling accordingly.

In [114]: df.head()

Out[114]:
CustomerID Count Country State City Zip Code Lat Long Latitude Gender Senior Citizen Partner Dependents T M
0 3568 QBYK 1 United States California Los Angeles 90003 33.964131 -118.272783 33.964131 -118.272783 Male No No No No
1 9237 HQTU 1 United States California Los Angeles 90005 34.059281 34.059281 -118.307420 Female No No No Yes
2 6305 QBSQ 1 United States California Los Angeles 90006 34.048013 34.048013 -118.293953 Female No No No Yes
3 7992 PQOKP 1 United States California Los Angeles 90010 34.062125 34.062125 -118.315709 Female No Yes Yes Yes
4 0280 XJGEX 1 United States California Los Angeles 90015 34.039224 34.039224 -118.266293 Male No No No Yes

In [ ]:

At first sight, we can see that all the samples share the same value of the columns 'Count', 'Country', 'State' and the only location related attributes that vary are the columns 'City' and 'Zip Code' so we will keep these ones and drop the others. Also, we assume that in our analysis we should consider different cities and not different areas in the same city, so the columns 'Lat Long', 'Latitude', 'Longitude' will be dropped aswell. Last but not least we will drop the column 'CustomerID' as it clearly has no impact on the output.

We will be keeping all the remaining variables for now, but that doesn't mean they're all significant, but we will study the importance of each one in further steps using statistical methods.

In [115]: df.drop(['CustomerID'], axis=1, inplace=True)
df.drop(['Count'], axis=1, inplace=True)
df.drop(['Country'], axis=1, inplace=True)
df.drop(['State'], axis=1, inplace=True)
df.drop(['Lat Long'], axis=1, inplace=True)
df.drop(['Latitude'], axis=1, inplace=True)
df.drop(['Longitude'], axis=1, inplace=True)
df.drop(['Churn Label'], axis=1, inplace=True)
df.drop(['CLTV'], axis=1, inplace=True)
df.drop(['Churn Score'], axis=1, inplace=True)
df.drop(['City'], axis=1, inplace=True)

In [116]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7043 entries, 0 to 7042
Data columns (total 22 columns):
# Column Non-Null Count Dtype
0 Zip Code 7043 non-null int64
1 Gender 7043 non-null object
2 Senior Citizen 7043 non-null object
3 Partner 7043 non-null object
4 Dependents 7043 non-null object
5 Tenure Months 7043 non-null int64
6 Phone Service 7043 non-null object
7 Multiple Lines 7043 non-null object
8 Internet Service 7043 non-null object
9 Online Security 7043 non-null object
10 Online Backup 7043 non-null object
11 Device Protection 7043 non-null object
12 Tech Support 7043 non-null object
13 Streaming TV 7043 non-null object
14 Streaming Movies 7043 non-null object
15 Contract 7043 non-null object
16 Paperless Billing 7043 non-null object
17 Payment Method 7043 non-null object
18 Monthly Charges 7043 non-null float64
19 Total Charges 7032 non-null float64
20 Churn Value 7043 non-null int64
21 Churn Reason 1869 non-null object
dtypes: float64(2), int64(3), object(17)
memory usage: 1.2+ MB

In [117]: df.describe()

Zip Code Tenure Months Monthly Charges Total Charges Churn Value
count 7043.000000 7043.000000 7043.000000 7032.000000 7043.000000
mean 9352.196406 32.371149 64.761692 2283.300441 0.265370
std 1865.794555 24.559481 30.090047 2266.771362 0.441561
min 90001.000000 0.000000 18.250000 18.800000 0.000000
25% 92102.000000 9.000000 35.500000 401.450000 0.000000
50% 95352.000000 29.000000 70.350000 1397.475000 0.000000
75% 95351.000000 55.000000 89.850000 3794.737000 1.000000
max 96161.000000 72.000000 118.750000 8684.800000 1.000000

In [118]: df.dtypes

Zip Code int64
Gender object
Senior Citizen object
Partner object
Dependents object
Tenure Months int64
Phone Service object
Multiple Lines object
Internet Service object
Online Security object
Online Backup object
Device Protection object
Tech Support object
Streaming TV object
Streaming Movies object
Contract object
Paperless Billing object
Payment Method object
Monthly Charges float64
Total Charges float64
Churn Value int64
Churn Reason object
dtype: object

At this stage we will look for categorical features and turn them into numerical values so we can later use them for plotting and training purposes. First we will list the variables having only 2 possible values then apply the label encoding on them.

In [119]: #here we are looking for variables as described in the markdown above
def list\_2cat(df):
 cols = []
 cols\_count = 0
 for col in df.columns[1:]:
 if df[col].nunique() <= 2:
 if df[col].dtype == 'object':
 cols\_count += 1
 cols.append(col)
 print(format(cols\_count) + ' columns have less than 3 possible values:')
 return cols

6 columns have less than 3 possible values:
['Gender', 'Senior Citizen', 'Partner', 'Dependents', 'Phone Service', 'Paperless Billing']

In [120]: for col in df.columns:
 if df[col].nunique() <= 4:
 if df[col].dtype == 'object' and not(col in cols):
 print(col)
 sizes = df[col].value\_counts(sort = True)
 colors = ["orange", "purple", "green", "red"]
 rcParams['figure.figsize'] = 5,5
 plt.pie(sizes, labels=df[col].unique(), colors=colors,
 autopct='%1.1f%%', shadow=True, startangle=270,
 plt.show()
 print('\n')

Multiple Lines
No 48.1%
Yes 42.2%
No phone service 9.7%

Internet Service
DSL 44.0%
Fiber optic 34.4%
No 21.7%

Online Security
Yes 49.7%
No 29.7%
No internet service 21.7%

Online Backup
Yes 43.8%
No 34.5%
No internet service 21.7%

Device Protection
No 43.9%
Yes 34.4%
No internet service 21.7%

Tech Support
No 49.3%
Yes 29.0%
No internet service 21.7%

Streaming TV
No 39.9%
Yes 38.4%
No internet service 21.7%

Streaming Movies
No 39.5%
Yes 38.9%
No internet service 21.7%

Contract
Month-to-month 55.0%
Two year 24.1%
One year 20.9%

Payment Method
Electronic check 22.2%
Mailed check 33.6%
Bank transfer (automatic) 21.4%
Credit card (automatic) 22.8%

The PieCharts corresponding to 'Online Security', 'Online Backup', 'Device Protection', 'Tech Support', 'Streaming TV' and 'Streaming Movies', all show that 21.7% of the customers have not subscribed for internet service, and that information is already present in the 'Internet Service' column. On the other hand, if we take 'Streaming Movies' for instance, it could either be 'Yes', 'No' or 'No internet Service' but 'No' and 'No Internet Service' mean the same thing, in other words we are dealing with redundant information.

Same thing for 'Multiple Lines' feature, the category 'No phone service' means the same as 'No', so best thing to do in order to avoid redundancy is to merge categories that have the same meaning.

In [121]: df.loc[df['Multiple Lines'] == "No phone service", "Multiple Lines"] = "No"
df.loc[df['Internet Service'] == "No", "Online Security", "Online Backup", "Device Protection",
 "Tech Support", "Streaming TV", "Streaming Movies"] = "No"

In [122]: list\_cols\_2cat(df)

13 columns have less than 3 possible values:

Out[122]:
['Gender',
 'Senior Citizen',
 'Partner',
 'Dependents',
 'Phone Service',
 'Multiple Lines',
 'Online Security',
 'Online Backup',
 'Device Protection',
 'Tech Support',
 'Streaming TV',
 'Streaming Movies',
 'Paperless Billing']

Above is the final list of variables with exactly two possible values and on which we will apply the 'Label Encoding'.

In [123]: encoded\_cols = []
label\_encoder = LabelEncoder()
label\_encoder.count = 0
for col in df.columns[1:]:
 if df[col].nunique() <= 2:
 if df[col].dtype == 'object':
 if (not(df[col].equals(df['Churn Reason']))):
 label\_encoder.fit(df[col])
 df[col] = label\_encoder.transform(df[col])
 label\_encoder.count += 1
 encoded\_cols.append(col)
print(format(label\_encoder.count) + ' columns were encoded')
encoded\_cols

13 columns were encoded

Out[123]:
['Gender',
 'Senior Citizen',
 'Partner',
 'Dependents',
 'Phone Service',
 'Multiple Lines',
 'Online Security',
 'Online Backup',
 'Device Protection',
 'Tech Support',
 'Streaming TV',
 'Streaming Movies',
 'Paperless Billing']

Moving on, we will try to locate missing values in every column:

In [124]: df.isna().sum()

Zip Code 0
Gender 0
Senior Citizen 0
Partner 0
Dependents 0
Tenure Months 0
Phone Service 0
Multiple Lines 0
Internet Service 0
Online Security 0
Online Backup 0
Device Protection 0
Tech Support 0
Streaming TV 0
Streaming Movies 0
Contract 0
Paperless Billing 0
Payment Method 0
Monthly Charges 0
Total Charges 11
Churn Value 0
Churn Reason 5174
dtype: int64

It's clear that 'Total Charges' contains 11 missing values and we need to take a look at their corresponding rows to decide whether to fill them or drop them.

In [125]: print(df.loc[df["Total Charges"].isnull()][["Churn Value"]].unique())

[0]

If we take a look at 'Tenure Months' and 'Churn Value' columns for this portion of data, which represents the customers with 0 Total Charges, we can see that they're all new customers that haven't yet churned and haven't spent enough periods to make their behavior significant, so this is non-significant data that better be dropped for lack of significance.

In [126]: df.dropna(subset=["Total Charges"])

In [127]: df.loc[df["Churn Reason"] == "NaN"]

Out[127]:
Zip Code Gender Senior Citizen Partner Dependents Tenure Months Phone Service Multiple Lines Internet Service Online Security Online Backup Device Protection Tech Support Streaming TV
0 0 0 0 0 0 2 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 0 0 0 1 2 0 0 0 0 0 0 0 0 0 0 0 1 70.70 1 0 0
2 0 0 0 1 8 0 0 0 1 1 1 1 1 99.65 1 0 0
3 0 1 1 28 0 1 1 1 1 1 1 1 104.80 1 0 0
4 0 0 0 1 49 0 0 0 1 1 1 1 1 103.70 1 0 0

We also have 5174 after dropping 11 lines in the previous step missing value in the 'Churn Reason' column which makes sense because that is the exact same number of customers that haven't churned yet (no churn -> no reason). This column will be dropped anyway before reasoning in dropping phases, but for now we will fill every 'NaN' with the string 'Hasn't churned yet' so we can plot graphs to see which reasons are more common etc.

In [128]: sizes = df['Churn Value'].value\_counts(sort = True)
colors = ["orange", "red"]
rcParams['figure.figsize'] = 5,5
plt.pie(sizes, labels=[sizes[0], "Yes"], colors=colors,
 autopct='%1.1f%%', shadow=True, startangle=270,
 plt.title('Churn distribution')
 plt.show()

churn distribution
[5163, no]
73.4%
[1869, yes] 26.6%

In [129]: df.loc[df["Churn Reason"].isna(), "Churn Reason"] = "Hasn't churned yet"

In [130]: from sklearn.feature\_selection import SelectKBest
from sklearn.feature\_selection import chi2

df\_new\_x = df[['Internet Service', 'Contract', 'Payment Method', 'Gender']].copy()
# df\_new\_x.drop(['Churn Value'], axis=1, inplace=True)
# df\_new\_x.drop(['Churn Reason'], axis=1, inplace=True)
# df\_new\_x.drop(['CLTV'], axis=1, inplace=True)
df\_new\_x = pd.get\_dummies(df\_new\_x, columns=['Internet Service', 'Contract', 'Payment Method'])
df\_new\_y = df['Churn Value'].copy()

In [131]: df\_new\_x.dtypes

Internet Service object
Contract object
Payment Method object
Gender int32
dtype: object

In [132]: encoded\_cols = []
label\_encoder = LabelEncoder()
label\_encoder.count = 0
for col in df\_new\_x.columns[0:]:
 if df\_new\_x[col].dtype == 'object':
 label\_encoder.fit(df[col])
 df\_new\_x[col] = label\_encoder.transform(df\_new\_x[col])
 label\_encoder.count += 1
 encoded\_cols.append(col)
print(format(label\_encoder.count) + ' columns were encoded')
encoded\_cols

3 columns were encoded

Out[132]:
['Internet Service', 'Contract', 'Payment Method']

In [133]: bestFeatures = SelectKBest(chi2, k=4)
fit = bestFeatures.fit(df\_new\_x, df\_new\_y)
df\_scores = pd.DataFrame(fit.pvalvals\_)
dfcolumns = pd.DataFrame(df\_new\_x.columns)
print(df\_scores.alargest(4, 'Score'))

Variable Score
3 Gender 6.140659e+01
0 Internet Service 1.827433e+03
2 Payment Method 1.395318e+14
1 Contract 5.186154e+244

The results above show that among the 4 variables, only 'Payment Method' and 'Contract' are in play when it comes to customers' tendency to churn so based on this we will drop the other two.

In [134]: df.drop(['Gender'], axis=1, inplace=True)
df.drop(['Internet Service'], axis=1, inplace=True)

As we can see thanks to this plot, the distribution of 'Churn\_Label' in terms of 'Zip\_Code' is the same for samples with churn\_value=0 and churn\_value=1, so the variable 'Zip\_Code' has no impact on the 'Churn\_Value' so it is to be Dropped.

In [135]: df.drop(['Zip Code'], axis=1, inplace=True)

Moving forward, we will use the 'Correlation Matrix' and 'Correlation Vector' for feature selection.

In [136]: plt.figure(figsize=(25,20))
sns.heatmap(df.corr(), cmap="RdBu\_r", center=0.0, annot=True)
plt.title('Correlation Matrix', fontsize=16)

Out[136]: Text(0.5, 1.0, 'Correlation Matrix')

Correlation Matrix
Senior Citizen 1.0000 -0.0017 -0.137 0.028 0.0084 0.18 0.007 0.007 0.06 -0.061 0.11 0.12 0.16 0.12 0.1 0.13
Partner 0.011 1.000 0.36 0.06 0.038 0.184 0.14 0.14 0.15 0.12 0.12 0.12 -0.014 0.098 0.12 0.15
Dependents 0.11 0.34 1.000 0.11 -0.0818 0.038 0.01 0.038 0.031 0.029 0.037 0.032 0.12 0.14 0.032 0.26
Tenure Months 0.016 0.033 0.023 1.000 0.1 0.0075 0.138 0.002 0.16 0.13 0.16 0.16 0.26 0.26 0.047 0.36
Phone Service 0.008 0.014 0.0038 0.0071 1.000 0.126 0.432 0.002 0.07 0.005 0.001 0.003 0.017 0.15 0.11 0.012
Multiple Lines 0.14 0.14 0.028 0.1 0.29 1.000 0.2 0.2 0.1 0.25 0.26 0.16 0.11 0.47 0.04
Online Security -0.019 0.14 0.011 0.12 0.002 0.099 1.000 0.1 0.28 0.27 0.25 0.18 0.13 0.0041 0.1 0.17
Online Backup 0.007 0.14 0.008 0.16 0.002 0.1 0.27 1.000 0.1 0.1 0.19 0.18 0.27 0.13 0.44 0.12 0.002
Device Protection 0.26 0.15 0.021 0.16 0.01 0.1 0.27 0.1 1.000 0.19 0.4 0.1 0.1 0.1 0.1 0.1 0.06
Tech Support -0.041 0.12 0.029 0.12 0.001 0.1 0.25 0.19 0.1 0.1 1.000 0.28 0.08 0.14 0.43 0.14
Streaming TV 0.11 0.12 0.057 0.126 0.001 0.26 0.18 0.26 0.26 0.1 0.12 1.000 0.12 0.12 0.13 0.052 0.083
Streaming Movies 0.122 0.122 0.057 0.126 0.001 0.26 0.13 0.17 0.14 0.12 0.1 0.12 1.000 0.12 0.13 0.11 0.052 0.083
Paperless Billing 0.16 0.014 0.12 0.008 0.017 0.16 0.09 0.13 0.1 0.108 0.12 0.11 0.1 0.105 0.19 0.19
Monthly Charges 0.22 0.088 0.14 0.25 0.25 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.19
Total Charges 0.1 0.132 0.021 0.01 0.17 0.47 0.41 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.2
Churn Value 0.15 0.15 0.15 0.15 0.15 0.012 0.06 0.17 0.002 0.006 0.15 0.001 0.001 0.15 0.15 0.1 1
Churn Reason 0.001 0.001 0.001 0.001 0.001 0.001 0.001 0.001 0.001 0.001 0.001 0.001 0.001 0.001 0.001 0.001 0.001

This Matrix shows the Correlation index between all variables two by two but we will only focus on the correlation of each variable with the target so that we will use correlation vector.

In [137]: plt.figure(figsize=(20,15))
sns.heatmap(df.corr()[['Churn Value']], cmap="RdBu\_r", center=0.0, annot=True)
plt.title('Correlation Vector', fontsize=16)

Out[137]: Text(0.5, 1.0, 'Correlation Vector')

Correlation Vector
Senior Citizen 0.15
Partner -0.15
Dependents -0.25
Tenure Months -0.012
Phone Service 0.002
Multiple Lines 0.04
Online Security -0.17
Online Backup -0.002
Device Protection -0.06
Tech Support -0.16
Streaming TV 0.063
Streaming Movies 0.061
Paperless Billing 0.13
Monthly Charges 0.19
Total Charges -0.2
Churn Value 1

We already know that 'Total Charges' is a calculated column (Total Charges = Tenure Months \* Monthly charges) in addition to that we can see through the Correlation Matrix that it is highly correlated with 'Tenure Months' so it will not impact our model if we dropped it.

If we take a look at the Correlation Vector, we can see that the correlation index of 'Device Protection', 'Multiple Lines', 'Online Backup', 'Phone Service' with the target variable is very low, so also these are to be taken away.

We will also be dropping the 'Churn Reason' column, because our purpose is to predict a customer's tendency to churn, and the churn reason appears in only a particular case where the predicted outcome will churn (according to the model).

In [138]: df.drop(['Total Charges'], axis=1, inplace=True) #high correlation with another var and it is calculated
df.drop(['Device Protection'], axis=1, inplace=True) #low correlation correlation
df.drop(['Multiple Lines'], axis=1, inplace=True) #low correlation correlation
df.drop(['Online Backup'], axis=1, inplace=True) #low correlation correlation
df.drop(['Phone Service'], axis=1, inplace=True) #low correlation correlation
df.drop(['Churn Reason'], axis=1, inplace=True) #useless in predictive modelling

So we're left with only two categorical variables that have more than two categories, and we will proceed with applying the OneHotEncoding on both of them:

In [139]: df = pd.get\_dummies(data=df, columns=['Payment Method', drop\_first=True])
df = pd.get\_dummies(data=df, columns=['Contract'], drop\_first=True)

In [140]: df.head()

Out[140]:
Senior Citizen Partner Dependents Tenure Months Online Security Tech Support Streaming TV Streaming Movies Paperless Billing Monthly Charges Churn Value Payment Method\_Credit Card (automatic) Method\_E
0 0 0 0 0 2 0 0 0 0 0 0 0 0 0 1 70.70 1 0
1 0 0 0 1 2 0 0 0 0 0 0 0 0 0 1 99.65 1 0
2 0 0 1 28 0 1 1 1 1 1 1 1 104.80 1 0
4 0 0 0 1 49 0 0 0 1 1 1 1 1 103.70 1 0

Now that we have our data cleaned, prepared and optimized, we will move on to the modeling part, we're dealing with a classification problem and Logistic regression is quite performant in our case, and to get our model ready we need to go through a few steps:

-Split our data into train and test data
-Standardize our data
-Initialize the logistic regression model with default parameters
-Find the best combination of hyper-parameters that will allow our model to give us the best score
-create the final model
-finally evaluate our model's performance

In [141]: #splitting data into train and test data
from sklearn.model\_selection import train\_test\_split
y=df['Churn Value']
x=df.drop(['Churn Value'], axis=1)
x\_train,x\_test,y\_train,y\_test=train\_test\_split(x,y,test\_size=0.3,random\_state=3)

In [142]: #data standardization
from sklearn.preprocessing import MinMaxScaler
sc=MinMaxScaler()
x\_train=sc.fit\_transform(x\_train)
x\_test=sc.transform(x\_test)

In [143]: #model initialization
from sklearn.linear\_model import LogisticRegression
logReg = LogisticRegression()
logReg.fit(x\_train, y\_train)
print('train score:', logReg.score(x\_train, y\_train))
print('test score:', logReg.score(x\_test, y\_test))
train score: 0.8110524177163755
test score: 0.796208308056872

In [144]: #finding most optimal values for hyper-parameters
param = {'penalty': ['l1', 'l2', 'elasticnet'], 'none'},
 'C': np.linspace(0.10, 100),
 'solver': ['newton-cg', 'lbfgs', 'lbfgs', 'sag', 'saga'],
 'max\_iter': [100, 1000, 2000, 3000]}

from sklearn.model\_selection import GridSearchCV
grid = GridSearchCV(logReg, param\_grid=param, cv=5, n\_jobs=-1)
grid.fit(x\_train, y\_train)
grid.best\_params\_

Out[144]: {'C': 0.30303030303030304, 'max\_iter': 100, 'penalty': 'l1', 'solver': 'liblinear'}

In [145]: grid.best\_score\_

Out[145]: 0.809024829383847

In [146]: #creating the final model with best hyper-parameters
logReg\_best\_params = LogisticRegression(C = 0.7070707070707071, max\_iter=100, penalty="none", solver = "sag")
logReg\_best\_params.fit(x\_train, y\_train)
print('train score:', logReg\_best\_params.score(x\_train, y\_train))
print('test score:', logReg\_best\_params.score(x\_test, y\_test))

train score: 0.8110524177163755
test score: 0.795260663507109

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\linear\_model\\_logistic.py:1320: UserWarning: Setting penalty='none' will ignore the C and l1\_ratio parameters
warnings.warn(

In [147]: #confusion matrix
from sklearn.metrics import plot\_confusion\_matrix, classification\_report
plot\_confusion\_matrix(logReg\_best\_params, x\_test, y\_test)
plt.show()

True \ Predicted 0 1
0 1380 168
1 264 258
Predicted label 0 1
Color scale 200 1200



```
In [148]: y_pred = logReg_best_params.predict(x_test)
print(classification_report(y_test, y_pred, digits=6))

              precision    recall  f1-score   support

    0       0.839416   0.891473   0.864662       1548
    1       0.639485   0.530249   0.579767       562

 accuracy          0.739451   0.710961   0.722214       2110
 macro avg         0.739451   0.710961   0.722214       2110
 weighted avg      0.786164   0.795261   0.788780       2110
```

Our Current model has a score of 70.52% and According to the confusion matrix above, the precision for the first category (0.839) is significantly higher than the precision for the second category (0.639), that means that the error is more likely to occur when the prediction result is "1" (the customer has churned).

Now that we're done with the linear regression model, we will try to elaborate a Decision Tree Model for our problem and for that we will go through almost the steps as before, then compare it to the previous model and see what we can get for a result.

```
In [149]: #splitting data into train and test data
from sklearn.tree import DecisionTreeClassifier
y=df['Churn Value']
x=df.drop(['Churn Value'], axis=1)
X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.3,random_state=3)

In [151]: #data standardization
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)

In [152]: #model initialization
decTree = DecisionTreeClassifier(random_state=0)
decTree.fit(X_train, y_train)
print('Le train score est :', decTree.score(X_train, y_train))
print('Le test score est :', decTree.score(X_test, y_test))

Le train score est : 0.9971556277935798
Le test score est : 0.716587677251184
```

```
In [153]: #finding most optimal values for hyper-parameters
param = ('criterion': ['gini', 'entropy'],
        'max_depth': np.arange(1,10),)

grid = GridSearchCV(decTree, param_grid=param, cv =5, n_jobs=-1)
grid.fit(X_train, y_train)

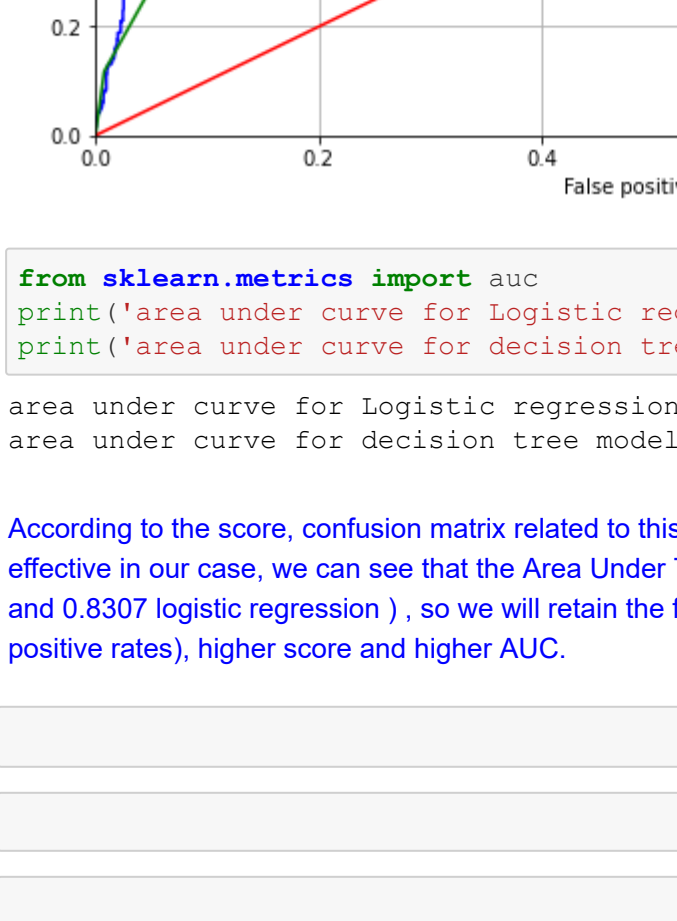
grid.best_params_

Out[153]: {'criterion': 'gini', 'max_depth': 4}
```

```
In [154]: #creating the final model with best hyper-parameters
decTree_best_params = DecisionTreeClassifier(random_state=0, criterion="gini", max_depth=4)
decTree_best_params.fit(X_train, y_train)
print('Le train score est :', decTree_best_params.score(X_train, y_train))
print('Le test score est :', decTree_best_params.score(X_test, y_test))

Le train score est : 0.8039414872003251
Le test score est : 0.7729857813905214
```

```
In [155]: #confusion matrix
from sklearn.metrics import plot_confusion_matrix, classification_report
plt.confusion_matrix(decTree_best_params, X_test, y_test)
plt.show()
```



```
In [156]: y_pred = decTree_best_params.predict(X_test)
print(classification_report(y_test, y_pred, digits=6))

              precision    recall  f1-score   support

    0       0.819968   0.885013   0.851196       1548
    1       0.594533   0.464413   0.521479       562

 accuracy          0.707201   0.674713   0.686337       2110
 macro avg         0.709850   0.772986   0.763376       2110
 weighted avg      0.739850   0.772986   0.763376       2110
```

```
In [157]: #models comparison in terms of AUC indicator
from sklearn.metrics import roc_curve
y_scores_logReg = logReg_best_params.predict_proba(X_test)
y_scores_decTree = decTree_best_params.predict_proba(X_test)

fprLr, tprLr, thresholdsLr = roc_curve(y_test, y_scores_logReg[:, 1])
fprDt, tprDt, thresholdsDt = roc_curve(y_test, y_scores_decTree[:, 1])

plt.figure(figsize=(10,5))
plt.plot(fprLr, tprLr, 'blue',label="Logistic regression" )
plt.plot(fprDt, tprDt, 'green',label="Decision Tree")
plt.xlim([0,1])
plt.ylim([0,1])
plt.xlabel('False positive rate')
plt.ylabel('True positive rate')
plt.plot([0,1], [0,1], 'r')
plt.grid(True)
plt.legend()
plt.show()
```



```
In [158]: from sklearn.metrics import auc
print('area under curve for Logistic regression model: '+str(auc(fprLr,tprLr)))
print('area under curve for Decision tree model: '+str(auc(fprDt, tprDt)))

area under curve for Logistic regression model: 0.8307246406797428
area under curve for Decision tree model: 0.7912224015375119
```

According to the score, confusion matrix related to this model and the AUC plot above we can say that the logistic regression model is more effective in our case, we can see that the Area Under The curve for the decision tree model is higher than the other (0.7912 decision tree and 0.8307 logistic regression) , so we will retain the first model as most performant for its higher precision (lower false negative and false positive rates), higher score and higher AUC.

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]: