

▼ Telecom Churn: Logistic Regression with PCA

With 21 predictor variables, we need to predict whether a particular customer will switch to another telecom provider or not. In telecom terminology, customer attrition is referred to as 'churn'.

▼ Importing and Merging Data

```
# Importing Pandas and NumPy
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

# Importing all datasets
churn_data = pd.read_csv("churn_data.csv")
customer_data = pd.read_csv("customer_data.csv")
internet_data = pd.read_csv("internet_data.csv")

print(len(churn_data))
print(len(customer_data))
print(len(internet_data))

7043
7043
7043

#Merging on 'customerID'
df_1 = pd.merge(churn_data, customer_data, how='inner', on='customerID')

#Final dataframe with all predictor variables
telecom = pd.merge(df_1, internet_data, how='inner', on='customerID')
```

▼ Let's understand the structure of our dataframe

```
# Let's see the head of our master dataset
telecom.head()
```

	customerID	tenure	PhoneService	Contract	PaperlessBilling	PaymentMethod	MonthlyChurn
0	7590-VHVEG	1	No	Month-to-month	Yes	Electronic check	
1	5575-GNVDE	34	Yes	One year	No	Mailed check	
2	3668-QPYBK	2	Yes	Month-to-month	Yes	Mailed check	
3	7795-CFOCW	45	No	One year	No	Bank transfer (automatic)	
4	9237-HQITU	2	Yes	Month-to-month	Yes	Electronic check	

5 rows × 21 columns

```
telecom.describe()
```

	tenure	MonthlyCharges	SeniorCitizen
count	7043.000000	7043.000000	7043.000000
mean	32.371149	64.761692	0.162147
std	24.550121	32.200017	0.260010

▼ Data Preparation

```
25%      9 000000      35 500000      0 000000
# Converting Yes to 1 and No to 0
telecom['PhoneService'] = telecom['PhoneService'].map({'Yes': 1, 'No': 0})
telecom['PaperlessBilling'] = telecom['PaperlessBilling'].map({'Yes': 1, 'No': 0})
telecom['Churn'] = telecom['Churn'].map({'Yes': 1, 'No': 0})
telecom['Partner'] = telecom['Partner'].map({'Yes': 1, 'No': 0})
telecom['Dependents'] = telecom['Dependents'].map({'Yes': 1, 'No': 0})
```

▼ Dummy Variable Creation

```
# Creating a dummy variable for the variable 'Contract' and dropping the first one.
cont = pd.get_dummies(telecom['Contract'],prefix='Contract',drop_first=True)
#Adding the results to the master dataframe
telecom = pd.concat([telecom,cont],axis=1)

# Creating a dummy variable for the variable 'PaymentMethod' and dropping the first one.
pm = pd.get_dummies(telecom['PaymentMethod'],prefix='PaymentMethod',drop_first=True)
#Adding the results to the master dataframe
telecom = pd.concat([telecom,pm],axis=1)

# Creating a dummy variable for the variable 'gender' and dropping the first one.
gen = pd.get_dummies(telecom['gender'],prefix='gender',drop_first=True)
#Adding the results to the master dataframe
telecom = pd.concat([telecom,gen],axis=1)

# Creating a dummy variable for the variable 'MultipleLines' and dropping the first one.
ml = pd.get_dummies(telecom['MultipleLines'],prefix='MultipleLines')
# dropping MultipleLines_No phone service column
ml1 = ml.drop(['MultipleLines_No phone service'],1)
#Adding the results to the master dataframe
telecom = pd.concat([telecom,ml1],axis=1)

# Creating a dummy variable for the variable 'InternetService' and dropping the first one.
iser = pd.get_dummies(telecom['InternetService'],prefix='InternetService',drop_first=True)
#Adding the results to the master dataframe
telecom = pd.concat([telecom,iser],axis=1)

# Creating a dummy variable for the variable 'OnlineSecurity'.
os = pd.get_dummies(telecom['OnlineSecurity'],prefix='OnlineSecurity')
os1= os.drop(['OnlineSecurity_No internet service'],1)
#Adding the results to the master dataframe
telecom = pd.concat([telecom,os1],axis=1)

# Creating a dummy variable for the variable 'OnlineBackup'.
ob =pd.get_dummies(telecom['OnlineBackup'],prefix='OnlineBackup')
ob1 =ob.drop(['OnlineBackup_No internet service'],1)
#Adding the results to the master dataframe
telecom = pd.concat([telecom,ob1],axis=1)

# Creating a dummy variable for the variable 'DeviceProtection'.
dp =pd.get_dummies(telecom['DeviceProtection'],prefix='DeviceProtection')
dp1 = dp.drop(['DeviceProtection_No internet service'],1)
#Adding the results to the master dataframe
telecom = pd.concat([telecom,dp1],axis=1)

# Creating a dummy variable for the variable 'TechSupport'.
ts =pd.get_dummies(telecom['TechSupport'],prefix='TechSupport')
ts1 = ts.drop(['TechSupport_No internet service'],1)
#Adding the results to the master dataframe
telecom = pd.concat([telecom,ts1],axis=1)

# Creating a dummy variable for the variable 'StreamingTV'.
st =pd.get_dummies(telecom['StreamingTV'],prefix='StreamingTV')
st1 = st.drop(['StreamingTV_No internet service'],1)
#Adding the results to the master dataframe
telecom = pd.concat([telecom,st1],axis=1)
```

```

telecom = pd.concat([telecom,st1],axis=1)

# Creating a dummy variable for the variable 'StreamingMovies'.
sm =pd.get_dummies(telecom['StreamingMovies'],prefix='StreamingMovies')
sm1 = sm.drop(['StreamingMovies_No internet service'],1)
#Adding the results to the master dataframe
telecom = pd.concat([telecom,sm1],axis=1)

```

▼ Dropping the repeated variables

```

# We have created dummies for the below variables, so we can drop them
telecom = telecom.drop(['Contract','PaymentMethod','gender','MultipleLines','InternetService', 'OnlineSecurity', 'OnlineBackup', 'DeviceProtection','TechSupport', 'StreamingTV', 'StreamingMovies'], 1)

#The variable was imported as a string we need to convert it to float
telecom['TotalCharges'] =telecom['TotalCharges'].convert_objects(convert_numeric=True)
#telecom['tenure'] = telecom['tenure'].astype(int).astype(float)

```

C:\ProgramData\Anaconda3\lib\site-packages\ipykernel_launcher.py:2: FutureWarning: convert_objects is deprecated. Use the data-type specific methods instead.

```
telecom.info()
```

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 7043 entries, 0 to 7042
Data columns (total 32 columns):
customerID           7043 non-null object
tenure                7043 non-null int64
PhoneService          7043 non-null int64
PaperlessBilling      7043 non-null int64
MonthlyCharges        7043 non-null float64
TotalCharges          7032 non-null float64
Churn                 7043 non-null int64
SeniorCitizen         7043 non-null int64
Partner               7043 non-null int64
Dependents            7043 non-null int64
Contract_One year    7043 non-null uint8
Contract_Two year    7043 non-null uint8
PaymentMethod_Credit card (automatic) 7043 non-null uint8
PaymentMethod_Electronic check       7043 non-null uint8
PaymentMethod_Mailed check          7043 non-null uint8
gender_Male             7043 non-null uint8
MultipleLines_No        7043 non-null uint8
MultipleLines_Yes       7043 non-null uint8
InternetService_Fiber optic     7043 non-null uint8
InternetService_No       7043 non-null uint8
OnlineSecurity_No       7043 non-null uint8
OnlineSecurity_Yes      7043 non-null uint8
OnlineBackup_No          7043 non-null uint8
OnlineBackup_Yes         7043 non-null uint8
DeviceProtection_No     7043 non-null uint8
DeviceProtection_Yes    7043 non-null uint8
TechSupport_No          7043 non-null uint8
TechSupport_Yes         7043 non-null uint8
StreamingTV_No          7043 non-null uint8
StreamingTV_Yes         7043 non-null uint8
StreamingMovies_No       7043 non-null uint8
StreamingMovies_Yes     7043 non-null uint8
dtypes: float64(2), int64(7), object(1), uint8(22)
memory usage: 756.6+ KB

```

Now we can see we have all variables as integer.

▼ Checking for Outliers

```

# Checking for outliers in the continuous variables
num_telecom = telecom[['tenure','MonthlyCharges','SeniorCitizen','TotalCharges']]

# Checking outliers at 25%,50%,75%,90%,95% and 99%
num_telecom.describe(percentiles=[.25,.5,.75,.90,.95,.99])

```

	tenure	MonthlyCharges	SeniorCitizen	TotalCharges
count	7043.000000	7043.000000	7043.000000	7032.000000
mean	32.371149	64.761692	0.162147	2283.300441
std	24.559481	30.090047	0.368612	2266.771362
min	0.000000	18.250000	0.000000	18.800000
25%	9.000000	35.500000	0.000000	401.450000
50%	29.000000	70.350000	0.000000	1397.475000
75%	55.000000	89.850000	0.000000	3794.737500
90%	69.000000	102.600000	1.000000	5976.640000
95%	72.000000	107.400000	1.000000	6923.590000
99%	72.000000	114.729000	1.000000	8039.883000
max	72.000000	118.750000	1.000000	8684.800000

From the distribution shown above, you can see that there no outlier in your data. The numbers are gradually increasing.

▼ Checking for Missing Values and Inputing Them

It means that $11/7043 = 0.001561834$ i.e 0.1%, best is to remove these observations from the analysis

```
# Checking the percentage of missing values
round(100*(telecom.isnull().sum()/len(telecom.index)), 2)

customerID          0.00
tenure              0.00
PhoneService        0.00
PaperlessBilling    0.00
MonthlyCharges      0.00
TotalCharges        0.16
Churn               0.00
SeniorCitizen       0.00
Partner             0.00
Dependents          0.00
Contract_One year   0.00
Contract_Two year   0.00
PaymentMethod_Credit card (automatic) 0.00
PaymentMethod_Electronic check        0.00
PaymentMethod_Mailed check           0.00
gender_Male          0.00
MultipleLines_No     0.00
MultipleLines_Yes    0.00
InternetService_Fiber optic         0.00
InternetService_No          0.00
OnlineSecurity_No        0.00
OnlineSecurity_Yes       0.00
OnlineBackup_No          0.00
OnlineBackup_Yes         0.00
DeviceProtection_No      0.00
DeviceProtection_Yes     0.00
TechSupport_No          0.00
TechSupport_Yes         0.00
StreamingTV_No          0.00
StreamingTV_Yes         0.00
StreamingMovies_No      0.00
StreamingMovies_Yes     0.00
dtype: float64
```

```
# Removing NaN TotalCharges rows
telecom = telecom[~np.isnan(telecom['TotalCharges'])]
```

```
# Checking percentage of missing values after removing the missing values
round(100*(telecom.isnull().sum()/len(telecom.index)), 2)
```

```
customerID          0.0
tenure              0.0
PhoneService        0.0
PaperlessBilling    0.0
MonthlyCharges      0.0
TotalCharges        0.0
```

```

Churn          0.0
SeniorCitizen 0.0
Partner        0.0
Dependents     0.0
Contract_One year 0.0
Contract_Two year 0.0
PaymentMethod_Credit card (automatic) 0.0
PaymentMethod_Electronic check 0.0
PaymentMethod_Mailed check 0.0
gender_Male    0.0
MultipleLines_No 0.0
MultipleLines_Yes 0.0
InternetService_Fiber optic 0.0
InternetService_No 0.0
OnlineSecurity_No 0.0
OnlineSecurity_Yes 0.0
OnlineBackup_No 0.0
OnlineBackup_Yes 0.0
DeviceProtection_No 0.0
DeviceProtection_Yes 0.0
TechSupport_No 0.0
TechSupport_Yes 0.0
StreamingTV_No 0.0
StreamingTV_Yes 0.0
StreamingMovies_No 0.0
StreamingMovies_Yes 0.0
dtype: float64

```

Now we don't have any missing values

▼ Feature Standardisation

```

# Normalising continuous features
df = telecom[['tenure','MonthlyCharges','TotalCharges']]

normalized_df=(df-df.mean())/df.std()
telecom = telecom.drop(['tenure','MonthlyCharges','TotalCharges'], 1)
telecom = pd.concat([telecom,normalized_df],axis=1)
telecom.head()

```

	customerID	PhoneService	PaperlessBilling	Churn	SeniorCitizen	Partner	Dependents	Contract_One year	Contract_Two year
0	7590-VHVEG	0	1	0	0	1	0	0	0
1	5575-GNVDE	1	0	0	0	0	0	0	1
2	3668-QPYBK	1	1	1	0	0	0	0	0
3	7795-CFOCW	0	0	0	0	0	0	0	1
4	9237-HQITU	1	1	1	0	0	0	0	0

5 rows × 32 columns

▼ Checking the Churn Rate

```

churn = (sum(telecom['Churn'])/len(telecom['Churn'].index))*100
churn

```

26.578498293515356

We have almost 27% churn rate

▼ Model Building

Let's start by splitting our data into a training set and a test set.

▼ Splitting Data into Training and Test Sets

```
from sklearn.model_selection import train_test_split

# Putting feature variable to X
X = telecom.drop(['Churn','customerID'],axis=1)

# Putting response variable to y
y = telecom['Churn']

y.head()

0    0
1    0
2    1
3    0
4    1
Name: Churn, dtype: int64

# Splitting the data into train and test
X_train, X_test, y_train, y_test = train_test_split(X,y, train_size=0.7,test_size=0.3,random_state=100)
```

▼ Running Your First Training Model

```
import statsmodels.api as sm

# Logistic regression model
logm1 = sm.GLM(y_train,(sm.add_constant(X_train)), family = sm.families.Binomial())
logm1.fit().summary()
```

Generalized Linear Model Regression Results

Dep. Variable: Churn **No. Observations:** 4922
Model: GLM **Df Residuals:** 4898
Model Family: Binomial **Df Model:** 23

▼ Correlation Matrix

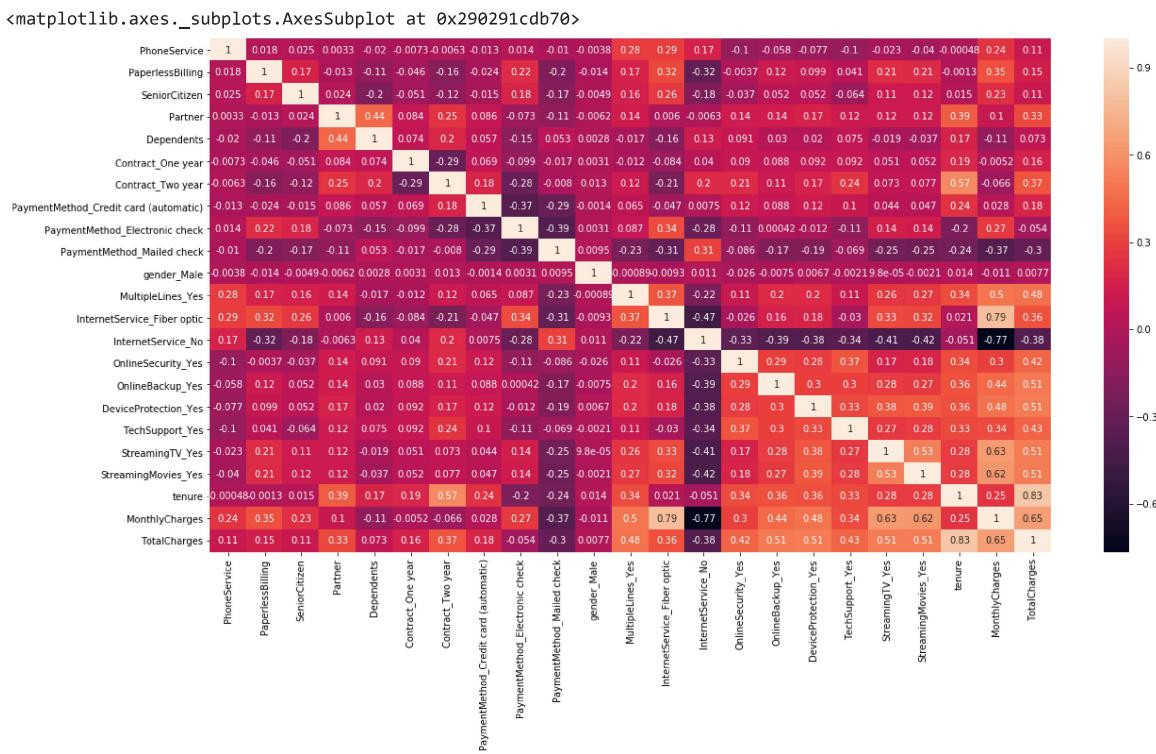
- ▼ Dropping highly correlated variables.

```
X_test2 = X_test.drop(['MultipleLines_No','OnlineSecurity_No','OnlineBackup_No','DeviceProtection_No','TechSupport_No','StreamingTV_No','StreetParking_No'])  
X_train2 = X_train.drop(['MultipleLines_No','OnlineSecurity_No','OnlineBackup_No','DeviceProtection_No','TechSupport_No','StreamingTV_No','StreetParking_No'])
```

▼ Checking the Correlation Matrix

After dropping highly correlated variables now let's check the correlation matrix again.

```
plt.figure(figsize = (20,10))
sns.heatmap(X_train2.corr(), annot = True)
```



▼ Re-Running the Model

Now let's run our model again after dropping highly correlated variables

```
logm2 = sm.GLM(y_train,(sm.add_constant(X_train2)), family = sm.families.Binomial())
logm2.fit().summary()
```

```

Generalized Linear Model Regression Results
Dep. Variable: Churn          No. Observations: 4922
Model: GLM                   Df Residuals: 4898
Model Family: Binomial        Df Model: 23
Link Function: logit         Scale: 1.0
Method: IRLS                 Log-Likelihood: -2004.7
Date: Wed, 25 Apr 2018       Deviance: 4009.4
Time: 20:12:54               Pearson chi2: 6.07e+03
No. Iterations: 7

                                                coef  std err      z   P>|z| [0.025 0.975]
const                                         -3.9338 1.545  -2.545 0.011 -6.963 -0.905
PhoneService                                    0.9507 0.789  1.205 0.228 -0.595 2.497
PaperlessBilling                                0.3254 0.090  3.614 0.000 0.149 0.502
SeniorCitizen                                    0.3984 0.102  3.924 0.000 0.199 0.597
Partner                                         0.0274 0.004  0.300 0.600 0.146 0.221

```

▼ Feature Selection Using RFE

```

Contract Two year          -1 2455 0 212 -5 874 0 000 -1 661 -0 830

from sklearn.linear_model import LogisticRegression
logreg = LogisticRegression()
from sklearn.feature_selection import RFE
rfe = RFE(logreg, 13)           # running RFE with 13 variables as output
rfe = rfe.fit(X,y)
print(rfe.support_)            # Printing the boolean results
print(rfe.ranking_)            # Printing the ranking

[ True  True False False False  True  True False  True False False  True
 False  True  True False  True False False False False  True False  True False
 False  True False  True False  True]
[ 1   1   2 18   6   1   1 11   1 12 14   1   8   1   1   4   1 15   5 13 10   7   1   3 16
 1 17   1   9   1]
   -
# Variables selected by RFE
col = ['PhoneService', 'PaperlessBilling', 'Contract_One year', 'Contract_Two year',
       'PaymentMethod_Electronic check', 'MultipleLines_No', 'InternetService_Fiber optic', 'InternetService_No',
       'OnlineSecurity_Yes', 'TechSupport_Yes', 'StreamingMovies_No', 'tenure', 'TotalCharges']

# Let's run the model using the selected variables
from sklearn.linear_model import LogisticRegression
from sklearn import metrics
logsk = LogisticRegression(C=1e9)
#logsk.fit(X_train[col], y_train)
logsk.fit(X_train, y_train)

LogisticRegression(C=1000000000.0, class_weight=None, dual=False,
                   fit_intercept=True, intercept_scaling=1, max_iter=100,
                   multi_class='ovr', n_jobs=1, penalty='l2', random_state=None,
                   solver='liblinear', tol=0.0001, verbose=0, warm_start=False)

#Comparing the model with StatsModels
#logm4 = sm.GLM(y_train,(sm.add_constant(X_train[col])), family = sm.families.Binomial())
logm4 = sm.GLM(y_train,(sm.add_constant(X_train)), family = sm.families.Binomial())
modres = logm4.fit()
logm4.fit().summary()

```

```

Generalized Linear Model Regression Results
Dep. Variable: Churn          No. Observations: 4922
Model: GLM                    Df Residuals: 4898
Model Family: Binomial        Df Model: 23
Link Function: logit          Scale: 1.0
Method: IRLS                  Log-Likelihood: -2004.7
Date: Wed, 25 Apr 2018        Deviance: 4009.4
Time: 20:12:55                Pearson chi2: 6.07e+03
No. Iterations: 7

            coef  std err   z   P>|z| [0.025 0.975]
const      -3.2783 1.187  -2.762 0.006 -5.605 -0.952
PhoneService 0.8213 0.588  1.396 0.163 -0.332 1.974
PaperlessBilling 0.3254 0.090  3.614 0.000 0.149 0.502
SeniorCitizen 0.3984 0.102  3.924 0.000 0.199 0.597
Partner       0.0374 0.094  0.399 0.690 -0.146 0.221
Dependents    -0.1430 0.107 -1.332 0.183 -0.353 0.067
Contract_One year -0.6578 0.129 -5.106 0.000 -0.910 -0.405
Contract_Two year -1.2455 0.212 -5.874 0.000 -1.661 -0.830
PaymentMethod_Credit card (automatic) -0.2577 0.137 -1.883 0.060 -0.526 0.011
PaymentMethod_Electronic check 0.1615 0.113  1.434 0.152 -0.059 0.382
PaymentMethod_Mailed check -0.2536 0.137 -1.845 0.065 -0.523 0.016
gender_Male     -0.0346 0.078 -0.442 0.658 -0.188 0.119
MultipleLines_No 0.1295 0.205  0.632 0.527 -0.272 0.531
MultipleLines_Yes 0.6918 0.392  1.763 0.078 -0.077 1.461

```

```

X_test[col].shape
#res = modres.predict(X_test[col])

(2110, 13)

```

▼ Making Predictions

```

# Predicted probabilities
y_pred = logsk.predict_proba(X_test)
# Converting y_pred to a dataframe which is an array
y_pred_df = pd.DataFrame(y_pred)
# Converting to column dataframe
y_pred_1 = y_pred_df.iloc[:,[1]]
# Let's see the head
y_pred_1.head()

```

	1
0	0.437703
1	0.326774
2	0.004873
3	0.576256
4	0.007747

```

# Converting y_test to dataframe
y_test_df = pd.DataFrame(y_test)
y_test_df.head()

```

	Churn
942	0
3730	1
1761	0
2283	1
1872	0

```

# Putting CustID to index
y_test_df['CustID'] = y_test_df.index
# Removing index for both dataframes to append them side by side
y_pred_1.reset_index(drop=True, inplace=True)

```

```
y_test_df.reset_index(drop=True, inplace=True)
# Appending y_test_df and y_pred_1
y_pred_final = pd.concat([y_test_df,y_pred_1],axis=1)
# Renaming the column
y_pred_final= y_pred_final.rename(columns={ 1 : 'Churn_Prob'})
# Rearranging the columns
y_pred_final = y_pred_final.reindex_axis(['CustID','Churn','Churn_Prob'], axis=1)
# Let's see the head of y_pred_final
y_pred_final.head()
```

	CustID	Churn	Churn_Prob
0	942	0	0.437703
1	3730	1	0.326774
2	1761	0	0.004873
3	2283	1	0.576256
4	1872	0	0.007747

```
# Creating new column 'predicted' with 1 if Churn_Prob>0.5 else 0
y_pred_final['predicted'] = y_pred_final.Churn_Prob.map( lambda x: 1 if x > 0.5 else 0)
# Let's see the head
y_pred_final.head()
```

	CustID	Churn	Churn_Prob	predicted
0	942	0	0.437703	0
1	3730	1	0.326774	0
2	1761	0	0.004873	0
3	2283	1	0.576256	1
4	1872	0	0.007747	0

▼ Model Evaluation

```
from sklearn import metrics

# Confusion matrix
confusion = metrics.confusion_matrix( y_pred_final.Churn, y_pred_final.predicted )
confusion

array([[1354, 174],
       [248, 334]], dtype=int64)

# Predicted      Churn  not_churn  __all__
# Actual
# Churn          1359    169     1528
# not_churn      256    326     582
# __all__        1615    751     2110

#Let's check the overall accuracy.
metrics.accuracy_score(y_pred_final.Churn, y_pred_final.predicted)

0.8000000000000004

def draw_roc( actual, probs ):
    fpr, tpr, thresholds = metrics.roc_curve( actual, probs,
                                              drop_intermediate = False )
    auc_score = metrics.roc_auc_score( actual, probs )
    plt.figure(figsize=(6, 6))
    plt.plot( fpr, tpr, label='ROC curve (area = %0.2f)' % auc_score )
    plt.plot([0, 1], [0, 1], 'k--')
    plt.xlim([0.0, 1.0])
    plt.ylim([0.0, 1.05])
    plt.xlabel('False Positive Rate or [1 - True Negative Rate]')
    plt.ylabel('True Positive Rate')
    plt.title('Receiver operating characteristic example')
    plt.legend(loc="lower right")
```

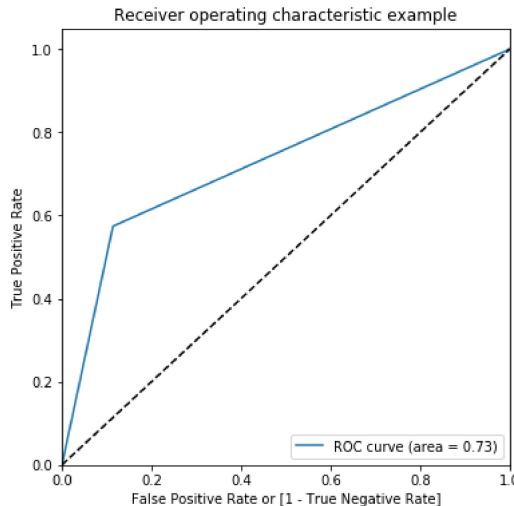
```

plt.show()

return fpr, tpr, thresholds

draw_roc(y_pred_final.Churn, y_pred_final.predicted)

```



```

(array([ 0.          ,  0.11387435,  1.        ]),
array([ 0.          ,  0.57388316,  1.        ]),
array([2, 1, 0], dtype=int64))

#draw_roc(y_pred_final.Churn, y_pred_final.predicted)
"{:2.2f}".format(metrics.roc_auc_score(y_pred_final.Churn, y_pred_final.Churn_Prob))

'0.83'

```

We see an overall AUC score of 0.83 looks like we did a decent job.

- But we did spend a lot of effort on the features and their selection.
- Can PCA help reduce our effort?

▼ PCA on the data

▼ Note -

- While computing the principal components, we must not include the entire dataset. Model building is all about doing well on the data we haven't seen yet!
- So we'll calculate the PCs using the train data, and apply them later on the test data

```

X_train.shape
# We have 30 variables after creating our dummy variables for our categories

(4922, 30)

#Importing the PCA module
from sklearn.decomposition import PCA
pca = PCA(svd_solver='randomized', random_state=42)

#Doing the PCA on the train data
pca.fit(X_train)

PCA(copy=True, iterated_power='auto', n_components=None, random_state=42,
    svd_solver='randomized', tol=0.0, whiten=False)

```

▼ Let's plot the principal components and try to make sense of them

- We'll plot original features on the first 2 principal components as axes

```
pca.components_
```

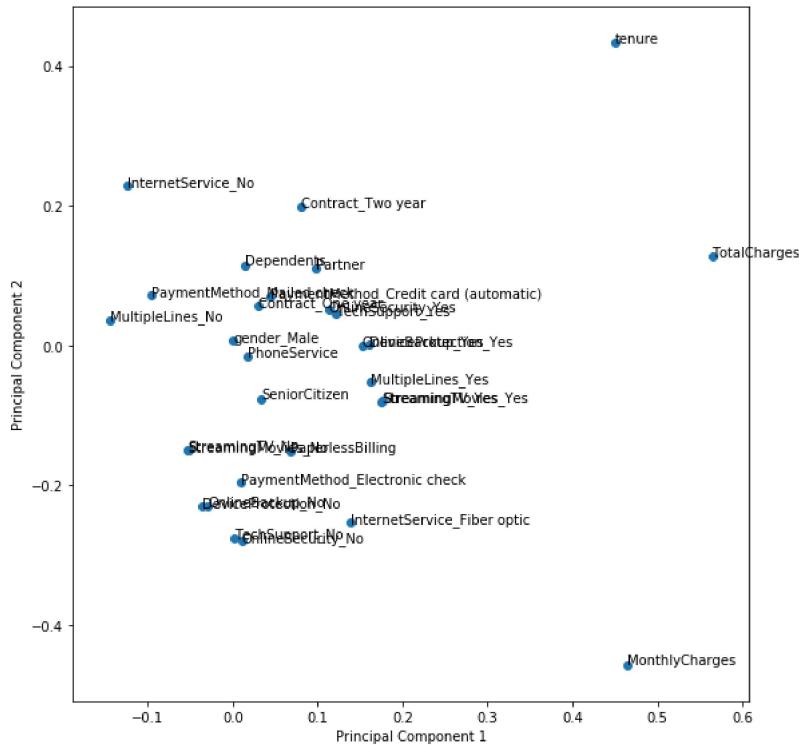
```
array([[ 1.76113503e-02,  6.73236050e-02,  3.33558780e-02,
       9.81025382e-02,  1.45144894e-02,  2.99711219e-02,
      8.03474721e-02,  4.45137169e-02,  1.03838839e-02,
     -9.64778216e-02,  5.46455027e-04,  -1.44992416e-01,
      1.62603767e-01,  1.38408762e-01,  -1.23559445e-01,
     1.07552300e-02,  1.12804215e-01,  -2.93144081e-02,
     1.52873853e-01,  -3.66586107e-02,  1.60218056e-01,
     2.26069096e-03,  1.21298754e-01,  -5.19779044e-02,
     1.75537350e-01,  -5.32208106e-02,  1.76780256e-01,
     4.50041974e-01,  4.64768058e-01,  5.64467383e-01],
   [-1.58225213e-02,  -1.51382114e-01,  -7.60391859e-02,
    1.10808074e-01,  1.15262476e-01,  5.70934546e-02,
    1.99196237e-01,  6.98211144e-02,  -1.94953309e-01,
    7.22696236e-02,  6.84716202e-03,  3.64391587e-02,
   -5.22616800e-02,  -2.53405153e-01,  2.28888352e-01,
   -2.80064773e-01,  5.11764215e-02,  -2.29155531e-01,
   2.67178910e-04,  -2.30417558e-01,  1.52920583e-03,
   -2.74825798e-01,  4.59374458e-02,  -1.48754387e-01,
   -8.01339654e-02,  -1.50215959e-01,  -7.86723933e-02,
   4.33270268e-01,  -4.56258849e-01,  1.27410027e-01],
   [-9.12291058e-02,  -2.84369927e-02,  -1.24001243e-02,
    5.63239688e-02,  4.78714602e-02,  -7.62225596e-04,
    2.57781910e-02,  5.81002273e-02,  -7.61267399e-02,
   -2.15070492e-02,  -8.51174807e-03,  -8.76630128e-02,
   -3.56609297e-03,  -9.44615765e-02,  -1.92982686e-01,
   -2.70477190e-02,  2.20030405e-01,  7.09799469e-02,
   1.22002739e-01,  2.79335424e-01,  -8.63527377e-02,
   1.14699466e-01,  7.82832203e-02,  4.96232463e-01,
   -3.03249777e-01,  4.89684449e-01,  -2.96701763e-01,
   2.43036444e-01,  -1.22710705e-01,  8.38171907e-02],
   [-1.52582548e-02,  -8.53395795e-02,  -1.24482551e-01,
    -1.15607472e-01,  1.35706616e-02,  6.18957244e-02,
   -2.55122762e-02,  1.08439476e-02,  -1.54943819e-01,
   1.50132655e-01,  -3.27194885e-02,  3.35259920e-01,
   -3.50518175e-01,  -9.95677606e-02,  -1.49234454e-01,
   -1.94633301e-01,  3.43867755e-01,  1.37892828e-02,
   1.35445172e-01,  -7.12015820e-02,  2.20436036e-01,
   -2.81687409e-01,  4.30921863e-01,  6.78808741e-02,
   8.13535803e-02,  6.63634814e-02,  8.28709730e-02,
   -2.84331268e-01,  2.16875383e-01,  -5.64112933e-03],
   [-1.14905477e-01,  -3.91805880e-02,  -2.43453437e-02,
    9.22236587e-02,  8.22694189e-02,  -7.39508732e-02,
   1.11130926e-01,  5.20967416e-02,  -4.54904532e-02,
   -1.28556843e-02,  2.36253723e-02,  -2.38918946e-01,
   1.24013469e-01,  -1.13486010e-01,  -7.56624541e-02,
   -2.67674094e-03,  7.83391950e-02,  6.40827206e-01,
   -5.65164752e-01,  -1.26242983e-01,  2.01905438e-01,
   -9.84497332e-02,  1.74112187e-01,  4.01689953e-02,
   3.54934588e-02,  -2.25276324e-02,  9.81900865e-02,
   3.70542100e-02,  -6.07852630e-02,  -6.31576961e-02],
   [ 3.64322765e-02,  2.12864636e-02,  4.16857747e-03,
    3.12327904e-01,  2.07866485e-01,  8.80938644e-02,
   -2.10506892e-02,  1.98347068e-02,  8.59334323e-02,
   -1.24909455e-01,  4.63903329e-02,  5.41058303e-01,
   -5.04626027e-01,  7.84512080e-02,  -3.51118051e-02,
   1.62939067e-01,  -1.27827262e-01,  1.91948712e-01,
   -1.56836907e-01,  -3.60056963e-02,  7.11175015e-02,
   2.33630834e-01,  -1.98519029e-01,  -9.70663787e-03,
```

```
colnames = list(X_train.columns)
pcs_df = pd.DataFrame({'PC1':pca.components_[0], 'PC2':pca.components_[1], 'Feature':colnames})
pcs_df.head()
```

	Feature	PC1	PC2
0	PhoneService	0.017611	-0.015823
1	PaperlessBilling	0.067324	-0.151382
2	SeniorCitizen	0.033356	-0.076039
3	Partner	0.098103	0.110808
4	Dependents	0.014514	0.115262

```
%matplotlib inline
fig = plt.figure(figsize = (8,8))
plt.scatter(pcs_df.PC1, pcs_df.PC2)
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
```

```
for i, txt in enumerate(pcs_df.Feature):
    plt.annotate(txt, (pcs_df.PC1[i],pcs_df.PC2[i]))
plt.tight_layout()
plt.show()
```



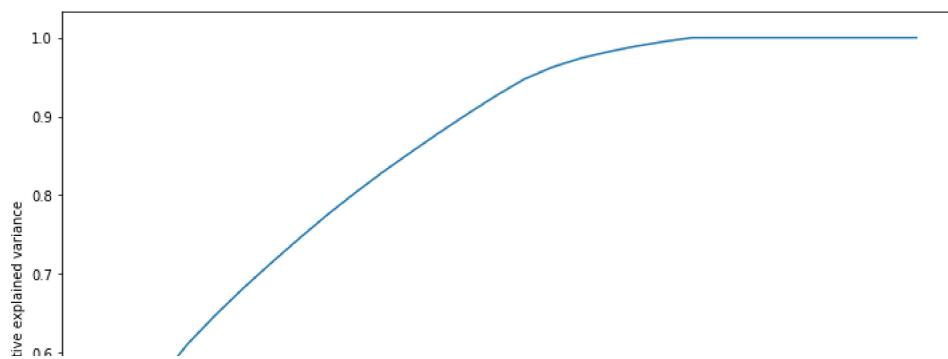
We see that the fist component is in the direction where the 'charges' variables are heavy

- These 3 components also have the highest of the loadings
- ▼ Looking at the screeplot to assess the number of needed principal components

```
pca.explained_variance_ratio_
```

```
array([ 3.27847827e-01,  1.75765007e-01,  6.03450801e-02,
       4.55648825e-02,  3.71944481e-02,  3.44939196e-02,
       3.21953215e-02,  3.12085042e-02,  3.00721929e-02,
       2.81455178e-02,  2.67941203e-02,  2.51671298e-02,
       2.45802775e-02,  2.36931522e-02,  2.27788300e-02,
       2.10982912e-02,  1.53183204e-02,  1.12792045e-02,
       8.45904722e-03,  7.12080502e-03,  5.92788522e-03,
       4.90401247e-03,  4.62235942e-05,  2.05310042e-33,
       2.05310042e-33,  2.05310042e-33,  2.18456080e-35])
```

```
#Making the screeplot - plotting the cumulative variance against the number of components
%matplotlib inline
fig = plt.figure(figsize = (12,8))
plt.plot(np.cumsum(pca.explained_variance_ratio_))
plt.xlabel('number of components')
plt.ylabel('cumulative explained variance')
plt.show()
```



- ▼ Looks like 16 components are enough to describe 95% of the variance in the dataset

- We'll choose 16 components for our modeling

```
#Using incremental PCA for efficiency - saves a lot of time on larger datasets
from sklearn.decomposition import IncrementalPCA
pca_final = IncrementalPCA(n_components=16)
```

- ▼ Basis transformation - getting the data onto our PCs

```
df_train_pca = pca_final.fit_transform(X_train)
df_train_pca.shape
```

```
(4922, 16)
```

- ▼ Creating correlation matrix for the principal components - we expect little to no correlation

```
#creating correlation matrix for the principal components
corrmat = np.corrcoef(df_train_pca.transpose())

#plotting the correlation matrix
%matplotlib inline
plt.figure(figsize = (20,10))
sns.heatmap(corrmat, annot = True)
```



```
# 1s -> 0s in diagonals
corrmat_nodiag = corrmat - np.diagflat(corrmat.diagonal())
print("max corr:",corrmat_nodiag.max(), ", min corr: ", corrmat_nodiag.min(),)
# we see that correlations are indeed very close to 0

max corr: 0.00181652621921 , min corr: -0.00107516919037
```

▼ Indeed - there is no correlation between any two components! Good job, PCA!

- We effectively have removed multicollinearity from our situation, and our models will be much more stable

```
#Applying selected components to the test data - 16 components
df_test_pca = pca_final.transform(X_test)
df_test_pca.shape

(2110, 16)
```

▼ Applying a logistic regression on our Principal Components

- We expect to get similar model performance with significantly lower features
- If we can do so, we would have done effective dimensionality reduction without losing any import information

```
#Training the model on the train data
from sklearn.linear_model import LogisticRegression
from sklearn import metrics

learner_pca = LogisticRegression()
model_pca = learner_pca.fit(df_train_pca,y_train)
```

Note

Note that we are fitting the original variable y with the transformed variables (principal components). This is not a problem because the transformation done in PCA is *linear*, which implies that you've only changed the way the new x variables are represented, though the nature of relationship between X and Y is still linear.

```
#Making prediction on the test data
pred_probs_test = model_pca.predict_proba(df_test_pca)[:,1]
"{:2.2}".format(metrics.roc_auc_score(y_test, pred_probs_test))

'0.83'
```

▼ Impressive! The same result, without all the hard work on feature selection!

Why not take it a step further and get a little more 'unsupervised' in our approach? This time, we'll let PCA select the number of components based on a variance cutoff we provide

```
pca_again = PCA(0.90)

df_train_pca2 = pca_again.fit_transform(X_train)
df_train_pca2.shape
# we see that PCA selected 14 components

(4922, 14)

#training the regression model
learner_pca2 = LogisticRegression()
model_pca2 = learner_pca2.fit(df_train_pca2,y_train)

df_test_pca2 = pca_again.transform(X_test)
df_test_pca2.shape

(2110, 14)

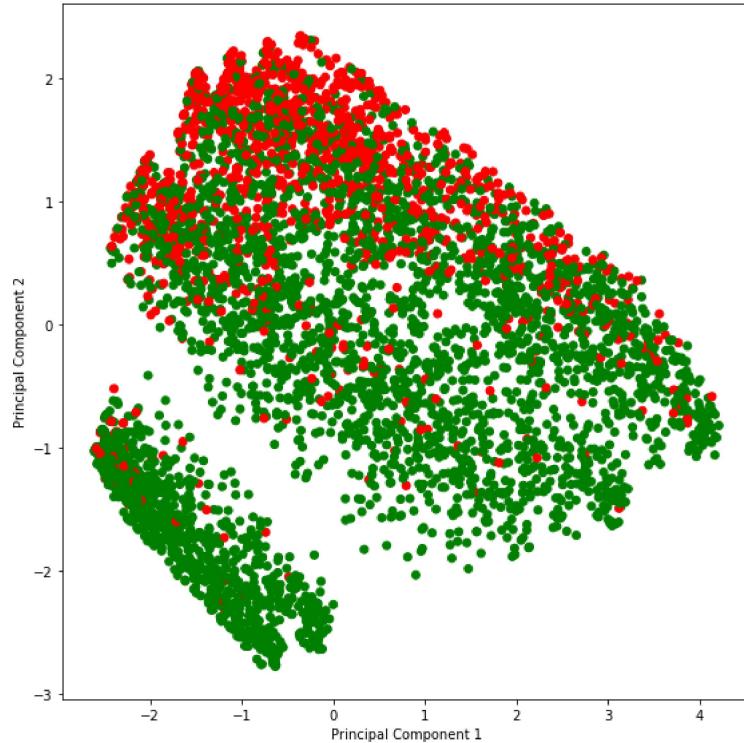
#Making prediction on the test data
pred_probs_test2 = model_pca2.predict_proba(df_test_pca2)[:,1]
```

```
"{:2.2f}".format(metrics.roc_auc_score(y_test, pred_probs_test2))
'0.83'
```

So there it is - a very similar result, without all the hassles. We have not only achieved dimensionality reduction, but also saved a lot of effort on feature selection.

- ▼ Before closing, let's also visualize the data to see if we can spot any patterns

```
%matplotlib inline
fig = plt.figure(figsize = (8,8))
plt.scatter(df_train_pca[:,0], df_train_pca[:,1], c = y_train.map({0:'green',1:'red'}))
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.tight_layout()
plt.show()
```



Looks like there is a good amount of separation in 2D, but probably not enough

Let's look at it in 3D, and we expect spread to be better (dimensions of variance, remember?)

```
%matplotlib notebook
from mpl_toolkits.mplot3d import Axes3D
fig = plt.figure(figsize=(8,8))
ax = Axes3D(fig)
# ax = plt.axes(projection='3d')
ax.scatter(df_train_pca[:,2], df_train_pca[:,0], df_train_pca[:,1], c=y_train.map({0:'green',1:'red'}))
```

```
<mpl_toolkits.mplot3d.art3d.Path3DCollection at 0x29028f9a320>
```

- ▼ So let's try building the model with just 3 principal components!

```
pca_last = PCA(n_components=3)
df_train_pca3 = pca_last.fit_transform(X_train)
df_test_pca3 = pca_last.transform(X_test)
df_test_pca3.shape

(2110, 3)

#training the regression model
learner_pca3 = LogisticRegression()
model_pca3 = learner_pca3.fit(df_train_pca3,y_train)
#Making prediction on the test data
pred_probs_test3 = model_pca3.predict_proba(df_test_pca3)[:,1]
"{:2.2f}".format(metrics.roc_auc_score(y_test, pred_probs_test3))

'0.82'
```

0.82! Isn't that just amazing!

