	Customer Churn occurs when a customer stops using a company's products or services. Customer churn affects profitability, especially in industries where revenues are heaveily dependent on subscription s( e.g. banks, telephone and internet service providers, pay-TV compa nies, insurance firsms, etc.). It is estimated that acquiring a new customer can cost up to five times more than retaining an existing one. Therefore, customer churn analysis is eesential as it can help a business:  • Identify problems in its services( e.g. poor quality product/service, poor customer support, wrong target audience, etc.)
	<ul> <li>Make correct strategic decisions that would lead to higher customer satisfaction and consequently higher customer retention.</li> <li>Main objective of the analysis</li> <li>The goal of this work is to understand and predict customer churn for a bank. Here, we are not interested in doing inferences but only</li> </ul>
	the prediction. Therefore, we will go into exploratory analysis but clean the data and use its geometry to classify the churns and non-churns using K-NN algorithm.  Data  We start our analysis with a data frame of IOK observations on the following 14 variables.
	<ul> <li>Exited: our target variable indicating whether a client exited or not the bank services (Yes/No)</li> <li>RowNumber: row indexing count</li> <li>Customerid: the unique identifier for each bank client</li> <li>Surname: client's surname</li> </ul>
	<ul> <li>Geography: the country of origin of the client. A factor with 3 levels (France, Germany and Spain)</li> <li>Gender: a factor with 2 levels (Male/Female)</li> <li>CreditScore: a numerical variable ranging from O to 850</li> </ul>
	<ul> <li>Age: a discrete variable ranging from 18 to 92</li> <li>Tenure: a factor with 11 levels corresponding to the number of months of retention</li> <li>Balance: numerical variable reporting client's balance amount</li> <li>NumOf Products: a factor with 4 levels (1, 2, 3 and 4)</li> </ul>
	<ul> <li>HasCrCard: a factor reporting Yes if client owns a credit card and No otherwise</li> <li>IsActiveMember: a factor with two levels</li> <li>EstimatedSalary: a numerical variable giving bank's est imation of t heir clients annual salary (in euro)</li> </ul>
In [1]:	Notice that we have three columns to remove since they do not provide any information. They are RowNumber, Customerid and Surname and we will have to predict Exited column where churn customers are labelled as 1 and non-churns as 0.  import os os.chdir('C:\\Users\\barla\\OneDrive\\Desktop\\Academia\\Data Science and Economics\\Datasets') import pandas as pd import numpy as np
	<pre>import matplotlib import matplotlib.pyplot as plt %matplotlib inline  # Tools for scaling data, PCA, and standard datasets from sklearn import preprocessing, decomposition, datasets</pre>
	<pre># Tools for tracking learning curves and perform cross validation from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score, validation_curve, learn # The k-NN learning algorithm from sklearn.neighbors import KNeighborsClassifier as kNN from sklearn.metrics import roc_curve, roc_auc_score</pre>
In [2]:	<pre>from sklearn.metrics import auc from sklearn.metrics import confusion_matrix</pre>
	Data columns (total 14 columns):  # Column Non-Null Count Dtype
	5 Gender 10000 non-null object 6 Age 10000 non-null int64 7 Tenure 10000 non-null int64 8 Balance 10000 non-null float64 9 NumOfProducts 10000 non-null int64 10 HasCrCard 10000 non-null int64 11 IsActiveMember 10000 non-null int64 12 EstimatedSalary 10000 non-null float64
In [3]:	<pre>13 Exited</pre>
Out[3]:	CreditScore         Geography         Gender         Age         Tenure         Balance         NumOfProducts         HasCrCard         IsActiveMember         EstimatedSalary         Exited           0         619         France         Female         42         2         0.00         1         1         1         1         101348.88         1           1         608         Spain         Female         41         1         83807.86         1         0         1         112542.58         0           2         502         France         Female         42         8         159660.80         3         1         0         113931.57         1
In [4]:	3 699 France Female 39 1 0.00 2 0 0 93826.63 0 4 850 Spain Female 43 2 125510.82 1 1 1 79084.10 0  #Create dummies  gender dummies = data.replace(to replace={'Gender': {'Female': 0, 'Male':1}})
Out[4]:	<pre>a = pd.get_dummies(data['Geography'], prefix = "Geo_dummy") frames = [gender_dummies, a] data = pd.concat(frames, axis=1) data = data.drop(columns=['Geography']) data.head()</pre>
	0       619       0       42       2       0.00       1       1       1       101348.88       1       1         1       608       0       41       1       83807.86       1       0       1       112542.58       0       0         2       502       0       42       8       159660.80       3       1       0       113931.57       1       1         3       699       0       39       1       0.00       2       0       0       93826.63       0       1
In [5]:	4 850 0 43 2 125510.82 1 1 1 1 79084.10 0 0  #Drop the label column  x=data.drop(['Exited'],axis=1) y=data['Exited']
In [6]:	We have around 8000 customers who did not churn and around 2000 customers who churned. This is an inbalanced dataset which is a problem for our analysis. We know that in K-NN the classifications are based on majority voting, however here using majority rule is not
In [7]:	a good idea since the algorithm will fail to correctly classify those who did churn due to the inbalanceness that exist in our dataset. That is why we will instead use the probabilities assigned to each data point for its classification.  X_train, X_test, y_train, y_test = train_test_split(x,y, test_size= 0.3, random_state = 42, stratify = y)  We have a relatively big data, so it makes sense to divide it as %70 training and %30 test set.
	We are now ready to train a classifier for this dataset. We create a 1-NN classifier object by invoking the function kNN(n_neighbors=1), where kNN() is the alias we created when we imported the module for KNeighborsClassifier(). The object is assigned to the variable knn.  What nearest neigbor is really doing is the following: we have 10 measurements for each observation, so a matrix with 10 dimensions.  We cannot see it neither view it but we know that it exists in the space. It is a mathematical construction. So we have this clouds of points in this 10 dimensional space and each point has a color. Color red for exited and green for non-exited, let's say. We split these
	points into subsets, training and test. Assume that we take out the test points and we are left %70 of original points in a random cloud. We give those points to the algorithm. What the algorithm does is to produce a classifier based on this %70 of original points. We then put back in those taken out points without the label and let the algorithm finds their place in this 10 dimensional space. What is the color of this point? Red or Green?  What K means does is to memorize all the training set. It takes %70 of the points and stores them in the memory. Everytime we give to
In [8]:	it a new point, it looks at the K closest point that it had in the training set. It measures the distance in this 30 dimensional space with Euclidean formula which works well for any finite dimensions.  knn = kNN (n_neighbors=1) knn.fit (X_train, y_train) knn.score (X_train, y_train), knn.score (X_test, y_test)
	Training accuracy turns out be 1 while testing is .68. We already knew that it would be %100 accurate on the train dataset since the algoritm memorizes all the train dataset. On the test data, on the other hand, the accurary is about %68 percent where every point is compared with the closes training point which is then given the label of the closest training point.
In [9]:	knn.fit(X_train, y_train) knn.score(X_train, y_train), knn.score(X_test, y_test)
	How come the training accuracy is not 1? It still memorizes the training set but it is applying the nearest neighbor rule. So suppose that we have a training point and we(algorithm) know the label of this training point because we memorized it. However, we are a three nearest neighbor classifer. We should use the three training points closest to this training point to make a classification. That's the definition of training error. The closest point to the training point will be itself but then we have two others points and they might have labels equal or disagreeing with our training point. So if we have a training point in a sea of training points of different color, it will be
	classified wrongly. Even if the algorithm knows its label, it still applies the three nearest neighbor rule.  However, we are doing better on the new data, namely, the test data! That's actually what we are interested in! We want to see how it does with the data it has never seen before.  Predictably, the training accuracy went down (by about 17%), while the test accuracy is now pretty close to the training accuracy.
	Next, we use the function learning_curve() to inspect the evolution of training and test performance of 7-NN for increasing sizes of the training set.  For each value of the training set size, a 5-fold stratified cross-validation is performed to estimate the risk.  sizes= np.arange(0.1, 0.9, 0.1)
In [11]:	<pre>train_size, train_score, val_score = learning_curve(kNN(n_neighbors = 7), x, y, train_sizes= sizes, cv=5)</pre>
	plt.ylabel("Accuracy") plt.show()  K-NN vs training size
	0.795 -  Validation Accuracy  Training Accuracy
	0.775  1000 2000 3000 4000 5000 6000  Training size  For this particular dataset, we conclude that the training size does not change the result a lot. This is probably the inbalancesness that exist in our dataset. The relative sizes of churns and non-churns do not really change a lot with respect to the training size provided to
In [12]:	the algorithm. We can still do more in developing our model. We could measure the behaviour of the training and test set as a function of the number of the nearest neighbor, K. So, we basically want to know what value of parameter to use, what value of K maximizes our accuracy on the test data?  neighbors = range(1, 30, 1)
Out [12] •	<pre>train_score, val_score = validation_curve(kNN(), x, y, 'n_neighbors', neighbors, cv=5) train_score, val_score  C:\Users\barla\anaconda3\lib\site-packages\sklearn\utils\validation.py:67: FutureWarning: Pass param_name=n_ neighbors, param_range=range(1, 30) as keyword args. From version 0.25 passing these as positional arguments will result in an error    warnings.warn("Pass {} as keyword args. From version 0.25 " (array([[1.</pre>
Out[12]:	[0.8375 , 0.84425 , 0.839 , 0.839875, 0.840125], [0.841875, 0.84375 , 0.843 , 0.839625, 0.838875], [0.817 , 0.812875, 0.815875, 0.814875, 0.81575 ], [0.816 , 0.814 , 0.8155 , 0.818375, 0.816125], [0.80625 , 0.80625 , 0.807125, 0.807125, 0.806625], [0.805875, 0.807875, 0.807125, 0.807125, 0.8045 ], [0.8005 , 0.80175 , 0.8025 , 0.802125, 0.8025 ], [0.801875, 0.801375, 0.800625, 0.80225 , 0.802375],
	[0.79875 , 0.8
	[0.7965 , 0.7965 , 0.796 , 0.797875, 0.798 ], [0.79675 , 0.798  , 0.796625, 0.798625, 0.796875], [0.796625 , 0.797875 , 0.79625 , 0.797625 , 0.797 ], [0.79775 , 0.7985  , 0.796  , 0.7975  , 0.797375], [0.796875 , 0.797875 , 0.796375 , 0.796875 , 0.7978 ], [0.796375 , 0.799   , 0.79675  , 0.79725  , 0.797875], [0.796375 , 0.7975   , 0.796625  , 0.797
	[0.796375, 0.796875, 0.796 , 0.796375, 0.796875], [0.796875, 0.797375, 0.796125, 0.7965 , 0.79725 ], [0.7965 , 0.796625, 0.796375, 0.796375], [0.7965 , 0.79675 , 0.796375, 0.7965 , 0.7965 ]]),  array([[0.7 , 0.677 , 0.6745, 0.692 , 0.68 ], [0.772 , 0.769 , 0.769 , 0.7775, 0.7745], [0.7315, 0.74 , 0.7295, 0.74 , 0.7365], [0.7855, 0.7815, 0.7815, 0.7815, 0.781 ], [0.767 , 0.7655, 0.767 , 0.7645, 0.766 ],
	[0.7895, 0.787 , 0.79 , 0.788 , 0.792 ], [0.7785, 0.781 , 0.781 , 0.7785, 0.7815], [0.786 , 0.7945, 0.793 , 0.786 , 0.7915], [0.782 , 0.7885, 0.7855, 0.7795, 0.783 ], [0.7905, 0.7935, 0.7955, 0.7905, 0.79 ], [0.787 , 0.7865, 0.79 , 0.7855, 0.7855], [0.791 , 0.7925, 0.795 , 0.791 , 0.7915], [0.79 , 0.7875, 0.793 , 0.792 , 0.7875],
	[0.794 , 0.7945, 0.796 , 0.7955, 0.7915], [0.7925, 0.7935, 0.7945, 0.7945, 0.789 ], [0.7945, 0.7935, 0.7965, 0.7965, 0.794 ], [0.795 , 0.7935, 0.796 , 0.796 , 0.793 ], [0.7955, 0.7955, 0.796 , 0.7965, 0.7945], [0.7965 , 0.7965 , 0.796 , 0.7965 , 0.7945], [0.795 , 0.797 , 0.796 , 0.7965 , 0.7945], [0.795 , 0.794 , 0.7955 , 0.796 , 0.793 ],
	[0.7945, 0.7965, 0.796 , 0.797 , 0.796 ], [0.7945, 0.7955, 0.7965, 0.7965, 0.7935], [0.7965, 0.797 , 0.796 , 0.7975, 0.795 ], [0.7965, 0.7965, 0.7955, 0.7975, 0.793 ], [0.796 , 0.796 , 0.796 , 0.7975, 0.7955], [0.796 , 0.797 , 0.796 , 0.7975, 0.7955], [0.795 , 0.7965, 0.7965, 0.797 , 0.7965], [0.795 , 0.7965, 0.7965, 0.7965, 0.796 ]]))
In [13]:	<pre>plt.title('K-NN vs number of neighbors') plt.plot(neighbors, np.mean(val_score,1), label='Testing Accuracy') plt.plot(neighbors, np.mean(train_score,1), label="Training Accuracy") plt.legend() plt.xlabel('Number of neighbors') plt.ylabel('Accuracy') plt.show()</pre>
	K-NN vs number of neighbors  1.00 - Testing Accuracy Training Accuracy  0.90 - Training Accuracy
	0.85 0.75 0.70
	The left side of the plot, where the testing accuracy and the training accuracy is very apart from each other is the region where there is overfitting. The training accuracy is really high because it is fitted so well on it but it is low in the testing accuracy. Afte a certain number of neighbors which is around 18, our accuracy becomes constant. Probably because, the particular distribution of the data points on the 10 dimensional space leads to this result. After around K=18, the number of neighbors provided do not change the accuracy.
In [14]: Out[14]:	<pre>learner = GridSearchCV(estimator=kNN(), param_grid = k_grid, cv=5, return_train_score=True) learner.fit(x,y) learner.best_params_, learner.best_score_</pre>
	So now, let's assume that we do not care about the testing. What we care about is only finding the right parameter. So, we take all the data we have and just want to determine the best value of the parameter for this data. So what do we do? We first of all decide which values to explore and then create a grid as above which is a dictionary. we then train the algorithm for each one of the values in the grid and we repeat the operation for each of those values. So we pick a value, we run 5 cross-validations on it: meaning that we train the algorithm 5 times with the same value of K each time taking out %20 of the data and train on the remaning %80 and then compute the mean. And then we take another value of K in the grid and repeat this operation for each value in the grid and we are looking for the
	best cross-validated performance: the best value of K that is associated with the highest average.  With K=41, we get the highest testing accuracy. However, this is a very high value for K and it is actually against the idea of similarity that is behind the algorithm. Also, we know that the accuracy becomes stable around 18 so we do not need to make it more complex with 41 parameters.
In [15]: Out[15]:	scores = cross_val_score (model, x, y, cv=5) scores.mean()  0.7963  We saw that cross-validation allows us to use the data for choosing a good value of the parameter. However, we are still left with the
	problem of estimating the risk of the classifier generated by the algorithm. Nested cross-validation provides a way of estimating the risk of a classifier generated by an algorithm whose parameters are tuned using cross-validation on the training set.  In the following example, we:  Run 5-fold cross-validation on the entire dataset.
In [16]:	
Out[16]:	learner = GridSearchCV(estimator=kNN(), param_grid=k_grid, cv=5) # internal C-V scores = cross_val_score(learner, x, y, cv=5) # external C-V scores, scores.mean()  (array([0.796 , 0.796 , 0.7965, 0.7965]), 0.7963)  Note that the nested cross-validated estimate turned out be the same as CV.
In [17]:	We can now do better by scaling down our dataset.  X_train, X_test, y_train, y_test = train_test_split(x,y, test_size= 0.3, random_state = 42, stratify = y) standard_scaler = preprocessing.StandardScaler()  X_train_standard = standard_scaler.fit_transform(X_train) #computes the mean and Sd for each feature  X_test_standard = standard_scaler.transform(X_test) #transforms the data according to the stats computed about the state of t
In [18]:	<pre>test_scores = [] test_scores_standard =[]  for k in neighbors:    knn = kNN(n_neighbors=k)    knn.fit(X_train, y_train)</pre>
In [19]:	<pre>test_scores.append(knn.score(X_test,y_test)) knn.fit(X_train_standard,y_train) test_scores_standard.append(knn.score(X_test_standard,y_test))  plt.title('k-NN vs. number of neighbors') plt.plot(neighbors, test_scores, label='Testing accuracy') plt.plot(neighbors, test_scores_standard,label='Testing accuracy(scaled)') plt.legend()</pre>
	plt.xlabel('Number of neighbors') plt.ylabel('Accuracy') plt.show()  k-NN vs. number of neighbors  0.84  0.82
	0.82 - 0.80 - 0.78 - 0.76 - 0.74 -
	0.72 - 0.70 - Testing accuracy Testing accuracy(scaled)  0.68 - Testing accuracy(scaled)  0.68 - Testing accuracy(scaled)  Number of neighbors
In [20]:	knn_final.fit(X_train_standard,y_train) knn_final.score(X_test_standard,y_test)
Out[20]:  In [21]:  In [22]:	Our testing accuracy is %83 which is not bad. However, we expect this model to classify badly the ones who did not churn.  predict = knn_final.predict(X_test_standard)  a,b,c,d=confusion_matrix(predict,y_test)[0,0],confusion_matrix(predict,y_test)[0,1],confusion_matrix(predict)
Out[22]:	sensitivity = (a/(a+c))*100 specificity =(d/(d+b))*100 sensitivity, specificity  (97.57220594390958, 29.296235679214405)  As we can see, we are doing pretty good in predictings those who did not churn but doing really bad in predicting those who churned.
In [23]:	This is a general problem in inbalanced datasets. We are able to classify correctly only %30 of the churned customers.
In [24]:	<pre>grid= [0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9] res = np.zeros((knn_probs.shape[0],len(grid))) row=0 while row &lt; len(knn_probs)-1:     for threshold in grid:         row=0         for probs in knn_probs[:,1]:             if(probs&gt;threshold).all():</pre>
In [25]:	<pre>if (probs&gt;threshold) .all():     res[row, grid.index(threshold)]=1     row+=1  else:     res[row, grid.index(threshold)]=0     row+=1  specificity=[0 for i in range(0, res.shape[1])]</pre>
	<pre>sensitivity=[0 for i in range(0,res.shape[1])]  for i in range(0,res.shape[1]):     a,b,c,d=confusion_matrix(res[:,i],y_test)[0,0],confusion_matrix(res[:,i],y_test)[0,1],confusion_matrix(res[inityity[i] = (a/(a+c))*100     specificity[i] = (d/(d+b))*100</pre>
In [26]:	<pre>plt.title('Sensitivity-Specificity vs. threshold') plt.plot(grid, sensitivity, label='Sensitivity') plt.plot(grid, specificity, label='Specificity') plt.legend() plt.xlabel('Threshold') plt.ylabel('Sensitivity-Specificity') plt.show()</pre>
	Sensitivity-Specificity vs. threshold
	Sensitivity Specificity
	As we can see, when the probabilities threshold used to classify the points increases, so when we approach the majority voting system which is the default in Knn, we have a really high sensitivity but low specificity. We cannot correctly classify the ones who did not churn. Around threshold 0.2, we get a nice trade off where we correctly classify around 73% of both classes which is the maximum we can get for this dataset with K-NN.
In [27]:	
In [28]:	<pre>a,b,c,d=confusion_matrix(res,y_test)[0,0],confusion_matrix(res,y_test)[0,1],confusion_matrix(res,y_test)[1,0] sensitivity = (a/(a+c))*100 specificity =(d/(d+b))*100 sensitivity, specificity</pre>
Out[28]:	(73.41984093763081, 75.77741407528642) With the threshold 0.2, we are able to classify %73 of churns and %75 of non-churns correctly.

Introduction