

ML0101EN-Clas-K-Nearest-neighbors-CustCat-py-v1

October 17, 2019

K-Nearest Neighbors

In this Lab you will load a customer dataset, fit the data, and use K-Nearest Neighbors to predict a data point. But what is **K-Nearest Neighbors**?

K-Nearest Neighbors is an algorithm for supervised learning. Where the data is ‘trained’ with data points corresponding to their classification. Once a point is to be predicted, it takes into account the ‘K’ nearest points to it to determine it’s classification.

0.0.1 Here’s an visualization of the K-Nearest Neighbors algorithm.

In this case, we have data points of Class A and B. We want to predict what the star (test data point) is. If we consider a k value of 3 (3 nearest data points) we will obtain a prediction of Class B. Yet if we consider a k value of 6, we will obtain a prediction of Class A.

In this sense, it is important to consider the value of k. But hopefully from this diagram, you should get a sense of what the K-Nearest Neighbors algorithm is. It considers the ‘K’ Nearest Neighbors (points) when it predicts the classification of the test point.

Table of contents

```
<ol>
  <li><a href="#about_dataset">About the dataset</a></li>
  <li><a href="#visualization_analysis">Data Visualization and Analysis</a></li>
  <li><a href="#classification">Classification</a></li>
</ol>
```

Lets load required libraries

```
[41]: import itertools
import numpy as np
import matplotlib.pyplot as plt
from matplotlib.ticker import NullFormatter
import pandas as pd
import numpy as np
import matplotlib.ticker as ticker
from sklearn import preprocessing
%matplotlib inline
```

```
<h2>About the dataset</h2>
```

Imagine a telecommunications provider has segmented its customer base by service usage patterns, categorizing the customers into four groups. If demographic data can be used to predict group membership, the company can customize offers for individual prospective customers. It is a classification problem. That is, given the dataset, with predefined labels, we need to build a model to be used to predict class of a new or unknown case.

The example focuses on using demographic data, such as region, age, and marital, to predict usage patterns.

The target field, called **custcat**, has four possible values that correspond to the four customer groups, as follows: 1- Basic Service 2- E-Service 3- Plus Service 4- Total Service

Our objective is to build a classifier, to predict the class of unknown cases. We will use a specific type of classification called K nearest neighbour.

Lets download the dataset. To download the data, we will use !wget to download it from IBM Object Storage.

```
[5]: !wget -O teleCust1000t.csv https://s3-api.us-geo.objectstorage.softlayer.net/
    ↪ cf-courses-data/CognitiveClass/ML0101ENv3/labs/teleCust1000t.csv
```

```
--2019-10-16 22:41:24-- https://s3-api.us-geo.objectstorage.softlayer.net/cf-
courses-data/CognitiveClass/ML0101ENv3/labs/teleCust1000t.csv
Resolving s3-api.us-geo.objectstorage.softlayer.net (s3-api.us-
geo.objectstorage.softlayer.net)... 67.228.254.193
Connecting to s3-api.us-geo.objectstorage.softlayer.net (s3-api.us-
geo.objectstorage.softlayer.net)|67.228.254.193|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 37048 (36K) [text/csv]
Saving to: 'teleCust1000t.csv'
```

```
teleCust1000t.csv 100%[=====>] 36.18K --.-KB/s in 0.02s
```

```
2019-10-16 22:41:24 (1.68 MB/s) - 'teleCust1000t.csv' saved [37048/37048]
```

Did you know? When it comes to Machine Learning, you will likely be working with large datasets. As a business, where can you host your data? IBM is offering a unique opportunity for businesses, with 10 Tb of IBM Cloud Object Storage: [Sign up now for free](#)

0.0.2 Load Data From CSV File

```
[23]: df = pd.read_csv('teleCust1000t.csv')
    df.head(10)
```

```
[23]:
```

	region	tenure	age	marital	address	income	ed	employ	retire	gender	\
0	2	13	44	1	9	64.0	4	5	0.0	0	
1	3	11	33	1	7	136.0	5	5	0.0	0	
2	3	68	52	1	24	116.0	1	29	0.0	1	
3	2	33	33	0	12	33.0	2	0	0.0	1	

4	2	23	30	1	9	30.0	1	2	0.0	0
5	2	41	39	0	17	78.0	2	16	0.0	1
6	3	45	22	1	2	19.0	2	4	0.0	1
7	2	38	35	0	5	76.0	2	10	0.0	0
8	3	45	59	1	7	166.0	4	31	0.0	0
9	1	68	41	1	21	72.0	1	22	0.0	0

	reside	custcat
0	2	1
1	6	4
2	2	3
3	1	1
4	4	3
5	1	3
6	5	2
7	3	4
8	5	3
9	3	2

<h2>Data Visualization and Analysis</h2>

Let's see how many of each class is in our data set

```
[24]: df['custcat'].value_counts()
      #df['age'].value_counts()

      #df.sort_values('age',ascending=True)  # sort values of age
```

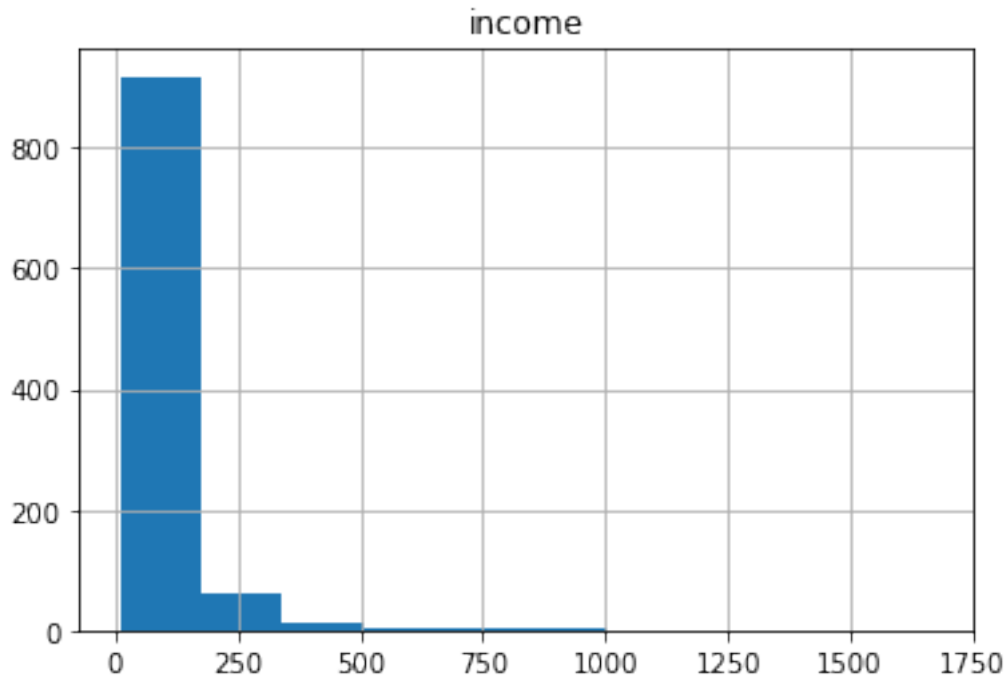
```
[24]: 3    281
      1    266
      4    236
      2    217
      Name: custcat, dtype: int64
```

281 Plus Service, 266 Basic-service, 236 Total Service, and 217 E-Service customers

You can easily explore your data using visualization techniques:

```
[29]: df.hist(column='income', bins=10)
```

```
[29]: array([[<matplotlib.axes._subplots.AxesSubplot object at 0x7f7f3017e0f0>]],
      dtype=object)
```



0.0.3 Feature set

Lets define feature sets, X:

```
[30]: df.columns
```

```
[30]: Index(['region', 'tenure', 'age', 'marital', 'address', 'income', 'ed',
            'employ', 'retire', 'gender', 'reside', 'custcat'],
            dtype='object')
```

To use scikit-learn library, we have to convert the Pandas data frame to a Numpy array:

```
[43]: X = df[['region', 'tenure', 'age', 'marital', 'address', 'income', 'ed',
            → 'employ', 'retire', 'gender', 'reside']].values #.astype(float)
X[0:5]
```

```
[43]: array([[ 2., 13., 44., 1., 9., 64., 4., 5., 0., 0., 2.],
            [ 3., 11., 33., 1., 7., 136., 5., 5., 0., 0., 6.],
            [ 3., 68., 52., 1., 24., 116., 1., 29., 0., 1., 2.],
            [ 2., 33., 33., 0., 12., 33., 2., 0., 0., 1., 1.],
            [ 2., 23., 30., 1., 9., 30., 1., 2., 0., 0., 4.]])
```

What are our labels?

```
[40]: y = df['custcat'].values
      y[0:5]
```

```
[40]: array([1, 4, 3, 1, 3])
```

0.1 Normalize Data

Data Standardization give data zero mean and unit variance, it is good practice, especially for algorithms such as KNN which is based on distance of cases:

```
[47]: X = preprocessing.StandardScaler().fit(X).transform(X.astype(float))
      X[0:5]
```

```
[47]: array([[ -0.02696767, -1.055125  ,  0.18450456,  1.0100505 , -0.25303431,
          -0.12650641,  1.0877526 , -0.5941226 , -0.22207644, -1.03459817,
          -0.23065004],
        [ 1.19883553, -1.14880563, -0.69181243,  1.0100505 , -0.4514148 ,
          0.54644972,  1.9062271 , -0.5941226 , -0.22207644, -1.03459817,
          2.55666158],
        [ 1.19883553,  1.52109247,  0.82182601,  1.0100505 ,  1.23481934,
          0.35951747, -1.36767088,  1.78752803, -0.22207644,  0.96655883,
          -0.23065004],
        [-0.02696767, -0.11831864, -0.69181243, -0.9900495 ,  0.04453642,
          -0.41625141, -0.54919639, -1.09029981, -0.22207644,  0.96655883,
          -0.92747794],
        [-0.02696767, -0.58672182, -0.93080797,  1.0100505 , -0.25303431,
          -0.44429125, -1.36767088, -0.89182893, -0.22207644, -1.03459817,
          1.16300577]])
```

0.1.1 Train Test Split

Out of Sample Accuracy is the percentage of correct predictions that the model makes on data that that the model has NOT been trained on. Doing a train and test on the same dataset will most likely have low out-of-sample accuracy, due to the likelihood of being over-fit.

It is important that our models have a high, out-of-sample accuracy, because the purpose of any model, of course, is to make correct predictions on unknown data. So how can we improve out-of-sample accuracy? One way is to use an evaluation approach called Train/Test Split. Train/Test Split involves splitting the dataset into training and testing sets respectively, which are mutually exclusive. After which, you train with the training set and test with the testing set.

This will provide a more accurate evaluation on out-of-sample accuracy because the testing dataset is not part of the dataset that have been used to train the data. It is more realistic for real world problems.

```
[55]: from sklearn.model_selection import train_test_split
      X_train, X_test, y_train, y_test = train_test_split( X, y, test_size=0.2,
      ↪random_state=4)
```

```
print ('Train set:', X_train.shape, y_train.shape)
print ('Test set:', X_test.shape, y_test.shape)
```

Train set: (800, 11) (800,)

Test set: (200, 11) (200,)

<h2>Classification</h2>

K nearest neighbor (KNN)

Import library Classifier implementing the k-nearest neighbors vote.

```
[56]: from sklearn.neighbors import KNeighborsClassifier
```

0.1.2 Training

Lets start the algorithm with k=4 for now:

```
[59]: k = 4
      #Train Model and Predict
      neigh = KNeighborsClassifier(n_neighbors = k).fit(X_train,y_train)
      neigh
```

```
[59]: KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
      metric_params=None, n_jobs=None, n_neighbors=4, p=2,
      weights='uniform')
```

0.1.3 Predicting

we can use the model to predict the test set:

```
[60]: yhat = neigh.predict(X_test)
      yhat[0:5]
```

```
[60]: array([1, 1, 3, 2, 4])
```

0.1.4 Accuracy evaluation

In multilabel classification, **accuracy classification score** is a function that computes subset accuracy. This function is equal to the jaccard_similarity_score function. Essentially, it calculates how closely the actual labels and predicted labels are matched in the test set.

```
[61]: from sklearn import metrics
      print("Train set Accuracy: ", metrics.accuracy_score(y_train, neigh.
      ↪predict(X_train)))
      print("Test set Accuracy: ", metrics.accuracy_score(y_test, yhat))
```

Train set Accuracy: 0.5475

Test set Accuracy: 0.32

0.2 Practice

Can you build the model again, but this time with k=6?

```
[68]: k = 6
      #Train Model and Predict
      neighs = KNeighborsClassifier(n_neighbors = k).fit(X_train,y_train)
      neighs

      yhat = neighs.predict(X_test)
      from sklearn import metrics
      print("Train set Accuracy: ", metrics.accuracy_score(y_train, neighs.
        ↪predict(X_train)))
      print("Test set Accuracy: ", metrics.accuracy_score(y_test, yhat))
```

Train set Accuracy: 0.51625

Test set Accuracy: 0.31

Double-click [here](#) for the solution.

What about other K? K in KNN, is the number of nearest neighbors to examine. It is supposed to be specified by the User. So, how can we choose right value for K? The general solution is to reserve a part of your data for testing the accuracy of the model. Then chose k =1, use the training part for modeling, and calculate the accuracy of prediction using all samples in your test set. Repeat this process, increasing the k, and see which k is the best for your model.

We can calculate the accuracy of KNN for different Ks.

```
[106]: print('KNN Begin for best value of K')
      Ks = 10
      mean_acc = np.zeros((Ks-1))
      #print('mean_acc begin',mean_acc)
      std_acc = np.zeros((Ks-1))
      #print('std_acc begin',std_acc)
      ConfusionMx = [];
      for n in range(1,Ks):

          #Train Model and Predict
          neigh = KNeighborsClassifier(n_neighbors = n).fit(X_train,y_train)
          yhat=neigh.predict(X_test)
          mean_acc[n-1] = metrics.accuracy_score(y_test, yhat)
          # print('Mean acc',mean_acc)

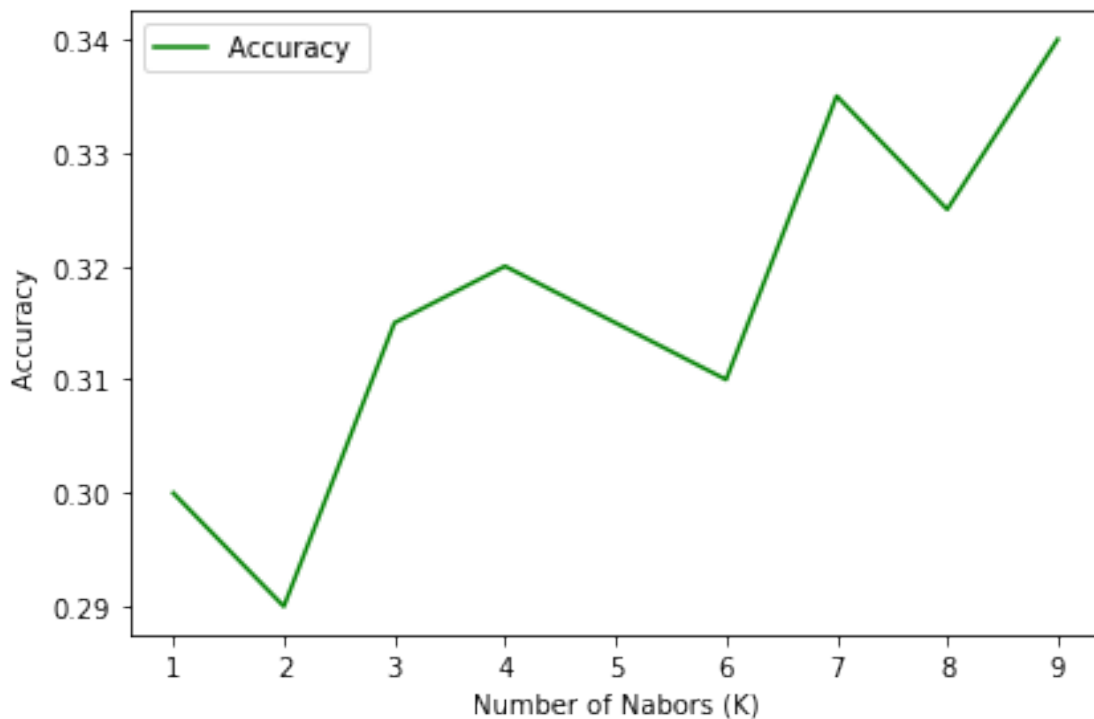

          std_acc[n-1]=np.std(yhat==y_test)/np.sqrt(yhat.shape[0])
          # print('Std acc',std_acc)
```

```
mean_acc
#std_acc
print('The Best Accuracy = ',mean_acc.max(),'with k=',mean_acc.argmax()+1)
```

KNN Begin for best value of K
The Best Accuracy = 0.34 with k= 9

Plot model accuracy for Different number of Neighbors

```
[104]: plt.plot(range(1,Ks),mean_acc,'g')
#plt.fill_between(range(1,Ks),mean_acc - 1 * std_acc,mean_acc + 1 * std_acc,
↳alpha=0.10)
plt.legend(('Accuracy ', '+/- 3xstd'))
plt.ylabel('Accuracy ')
plt.xlabel('Number of Nabors (K)')
plt.tight_layout()
plt.show()
```



```
[105]: print( "The best accuracy was with", mean_acc.max(), "with k=", mean_acc.
↳argmax()+1)
```

The best accuracy was with 0.34 with k= 9

Want to learn more?

IBM SPSS Modeler is a comprehensive analytics platform that has many machine learning algorithms. It has been designed to bring predictive intelligence to decisions made by individuals, by groups, by systems – by your enterprise as a whole. A free trial is available through this course, available here: [SPSS Modeler](#)

Also, you can use Watson Studio to run these notebooks faster with bigger datasets. Watson Studio is IBM's leading cloud solution for data scientists, built by data scientists. With Jupyter notebooks, RStudio, Apache Spark and popular libraries pre-packaged in the cloud, Watson Studio enables data scientists to collaborate on their projects without having to install anything. Join the fast-growing community of Watson Studio users today with a free account at [Watson Studio](#)

Thanks for completing this lesson!

Author: Saeed Aghabozorgi

Saeed Aghabozorgi, PhD is a Data Scientist in IBM with a track record of developing enterprise level applications that substantially increases clients' ability to turn data into actionable knowledge. He is a researcher in data mining field and expert in developing advanced analytic methods like machine learning and statistical modelling on large datasets.

Copyright © 2018 Cognitive Class. This notebook and its source code are released under the terms of the MIT License.