

```
###### STAT4080 Data Programming with Python (online) - Project  
# k nearest neighbours on the TunedIT data set
```

```
# Import packages
```

```
from pandas import Series, DataFrame  
import pandas as pd  
import numpy as np  
import numpy.random as npr
```

```
# For the project we will study the method of k nearest neighbours applied to a  
# music classification data set. These data come from the TunedIT website  
# http://tunedit.org/challenge/music-retrieval/genres  
# Each row corresponds to a different sample of music from a certain genre.  
# The original challenge was to classify the different genres (the original  
# prize for this was hard cash!). However we will just focus on a sample of the  
# data (~4000 samples) which is either rock or not. There are 191  
# characteristics (go back to the website if you want to read about these)  
# The general tasks of this exercise are to:  
# - Load the data set  
# - Standardise all the columns  
# - Divide the data set up into a training and test set  
# - Write a function which runs k nearest neighbours (kNN) on the data set.  
# (Don't worry you don't need to know anything about kNN)  
# - Check which value of k produces the smallest misclassification rate on the  
# training set  
# - Predict on the test set and see how it does
```

```
# Q1 Load in the data using the pandas read_csv function. The last variable  
# 'RockOrNot' determines whether the music genre for that sample is rock or not  
# What percentage of the songs in this data set are rock songs (to 1 d.p.)?
```

```
dataSet=pd.read_csv("tunedit_genres.csv")  
round(dataSet['RockOrNot'].value_counts(normalize=True)[1] * 100,1)
```

```
# Ans: 48.83 percentage of dataset has rock songs.
```

```
# Q2 To perform a classification algorithm, you need to define a classification  
# variable and separate it from the other variables. We will use 'RockOrNot' as  
# our classification variable. Write a piece of code to separate the data into a  
# DataFrames X and a Series y, where X contains a standardised version of  
# everything except for the classification variable ('RockOrNot'), and y contains  
# only the classification variable. To standardise the variables in X, you need  
# to subtract the mean and divide by the standard deviation
```

```
X=dataSet.drop('RockOrNot',1)  
y=dataSet['RockOrNot']  
Mean=Series(X.mean())  
standardDeviation=Series(X.std())  
X=(X-Mean)/standardDeviation
```

```
# Q3 Which variable in X has the largest correlation with y?
```

```
X.corrwith(y).idxmax()
```

```
# Ans: 'PAR_SFM_M' in X has the highest correlation with y.
```

```
# Q4 When performing a classification problem, you fit the model to a portion of
```

```
# your data, and use the remaining data to determine how good the model fit was.
# Write a piece of code to divide X and y into training and test sets, use 75%
# of the data for training and keep 25% for testing. The data should be randomly
# selected, hence, you cannot simply take the first, say, 3000 rows. If you select
# rows 1,4,7,8,13,... of X for your training set, you must also select rows
# 1,4,7,8,13,... of y for training set. Additionally, the data in the training
# set cannot appear in the test set, and vice versa, so that when recombined,
# all data is accounted for. Use the seed 123 when generating random numbers
# Note: The data may not spilt equally into 75% and 25% portions. In this
# situation you should round to the nearest integer.
```

Ans:

```
trainSety=y.sample(frac=0.75,random_state=123)
trainSetX=X.sample(frac=0.75,random_state=123)
testSety=y.drop(trainSety.index)
testSetX=X.drop(trainSetX.index)
trainSetX.index=range(len(trainSetX))
trainSety.index=range(len(trainSety))
testSetX.index=range(len(testSetX))
testSety.index=range(len(testSety))
```

```
# Q5 What is the percentage of rock songs in the training dataset and in the
# test dataset? Are they the same as the value found in Q1?
```

```
round(trainSety.value_counts(normalize=True)[1] * 100,1)
round(testSety.value_counts(normalize=True)[1] * 100,1)
```

Ans: The percentage of rock songs in training dataset is 49.4% and in the test dataset
They are nearly the same as that in Q1.

```
# Q6 Now we're going to write a function to run kNN on the data sets. kNN works
# by the following algorithm:
```

- # 1) Choose a value of k (usually odd)
- # 2) For each observation, find its k closest neighbours
- # 3) Take the majority vote (mean) of these neighbours
- # 4) Classify observation based on majority vote

```
# We're going to use standard Euclidean distance to find the distance between
# observations, defined as  $\sqrt{(x_i - x_j)^T (x_i - x_j)}$ 
# A useful short cut for this is the scipy functions pdist and squareform
```

```
# The function inputs are:
```

- # - DataFrame X of explanatory variables
- # - binary Series y of classification values
- # - value of k (you can assume this is always an odd number)

```
# The function should produce:
```

- # - Series y_star of predicted classification values

```
from scipy.spatial.distance import pdist, squareform
```

```
def kNN(X,y,k):
    # Find the number of observation
    noOfObs = len(X)
    # Set up return values
    y_star = []
    # Calculate the distance matrix for the observations in X
    dist = squareform(pdist(X))
```

```

# Make all the diagonals very large so it can't choose itself as a closest neighbour
np.fill_diagonal(dist,inf)
# Loop through each observation to create predictions
for m in range(noOfObs):
    a=dist[m].argsort()[:k]
    y_nearest=[]
    # Find the y values of the k nearest neighbours
    for noOfObs in a:
        y_nearest.append(y[noOfObs])
    # Now allocate to y_star
    y_star.append(round(Series(y_nearest).mean(),0))
return y_star

```

Q7 The misclassification rate is the percentage of times the output of a classifier doesn't match the classification value. Calculate the misclassification rate of the kNN classifier for X_train and y_train, with k=3.

```

difference=(trainSety-kNN(trainSetX,trainSety,3)).abs()
round(difference.value_counts(normalize=True)[1]*100,1)

# Ans: The misclassification rate of the kNN classifier is 4.7%

```

Q8 The best choice for k depends on the data. Write a function kNN_select that will run a kNN classification for a range of k values, and compute the misclassification rate for each.

The function inputs are:
- DataFrame X of explanatory variables
- binary Series y of classification values
- a list of k values k_vals

The function should produce:
- a Series mis_class_rates, indexed by k, with the misclassification rates for each k value in k_vals

```

def kNN_select(X,y,k_vals):
    mis_class_rates=[]
    for i in k_vals:
        diff=(y-kNN(X,y,i)).abs()
        a=round(diff.value_counts(normalize=True)[1]*100,1)
        mis_class_rates.append(a)
    mis_class_rates=Series(mis_class_rates,index=k_vals)
    return mis_class_rates

```

Q9 Run the function kNN_select on the training data for k = [1, 3, 5, 7, 9] and find the value of k with the best misclassification rate. Use the best value of k to report the mis-classification rate for the test data. What is the misclassification percentage with this k on the test set?

```

difference=Series(testSety-kNN(testSetX,testSety,kNN_select(trainSetX,trainSety,[1,3,5,7,9]))
round(difference.value_counts(normalize=True)[1]*100,1)

```

Ans: The misclassification percentage with this k on the test set is 5.0

Q10 Write a function to generalise the k nearest neighbours classification

```

# algorithm. The function should:
# - Separate out the classification variable for the other variables in the dataset,
#   i.e. create X and y.
# - Divide X and y into training and test set, where the number in each is
#   specified by 'percent_train'.
# - Run the k nearest neighbours classification on the training data, for a set
#   of k values, computing the mis-classification rate for each k
# - Find the k that gives the lowest mis-classification rate for the training data,
#   and hence, the classification with the best fit to the data.
# - Use the best k value to run the k nearest neighbours classification on the test
#   data, and calculate the mis-classification rate
# The function should return the mis-classification rate for a k nearest neighbours
# classification on the test data, using the best k value for the training data
# You can call the functions from Q6 and Q8 inside this function, provided they
# generalise, i.e. will work for any dataset, not just the TunedIT dataset.
def kNN_classification(df, class_column, seed, percent_train, k_vals):
    # df - DataFrame to
    # class_column - column of df to be used as classification variable, should
    #                specified as a string
    # seed - seed value for creating the training/test sets
    # percent_train - percentage of data to be used as training data
    # k_vals - set of k values to be tests for best classification

    # Separate X and y

    # Divide into training and test

    # Compute the mis-classification rates for each for the values in k_vals

    # Find the best k value, by finding the minimum entry of mis_class_rates

    # Run the classification on the test set to see how well the 'best fit'
    # classifier does on new data generated from the same source

    # Calculate the mis-classification rates for the test data

    y=df[class_column]
    X=df.drop(class_column,1)
    Mean=Series(X.mean())
    standardDeviation=Series(X.std())
    X=(X-Mean)/standardDeviation
    trainSety=y.sample(frac=percent_train,random_state=seed)
    trainSetX=X.sample(frac=percent_train,random_state=seed)
    testSety=y.drop(trainSety.index)
    testSetX=X.drop(trainSetX.index)
    trainSety.index=range(len(trainSety))
    trainSetX.index=range(len(trainSetX))
    testSety.index=range(len(testSety))
    testSetX.index=range(len(testSetX))
    difference=Series(testSety-kNN(testSetX,testSety,kNN_select(trainSetX,trainSety,k_val:
    mis_class_test=round(difference.value_counts(normalize=True)[1]*100,1)
    return mis_class_test

# Test your function with the TunedIT data set, with class_column = 'RockOrNot',
# seed = the value from Q4, percent_train = 0.75, and k_vals = set of k values
# from Q8, and confirm that it gives the same answer as Q9.

# Now test your function with another dataset, to ensure that your code

```

*# generalises. You can use the house_votes.csv dataset, with 'Party' as the
classifier. Select the other parameters as you wish.
This dataset contains the voting records of 435 congressman and women in the
US House of Representatives. The parties are specified as 1 for democrat and 0
for republican, and the votes are labelled as 1 for yes, -1 for no and 0 for
abstained.
Your kNN classifier should return a mis-classification for the test data (with
the best fit k value) of ~8%.*

```
kNN_classification(dataSet, 'RockOrNot', 123, 0.75, [1, 3, 5, 7, 9])
```

#kNN returns 5.0% on dataset

```
dataSet2=pd.read_csv("house_votes.csv")  
kNN_classification(dataSet2, 'Party', 123, 0.75, [1, 3, 5, 7, 9])
```

#kNN returns 8.3% on dataset2