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

# Import packages
from scipy.spatial.distance import pdist, squareform
from pandas import Series, DataFrame
import pandas as pd
import numpy as np
import numpy.random as npr
from scipy.stats import pearsonr


# 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.)?
song = pd.read_csv("./tunedit_genres.csv")
rock_ratio = song['RockOrNot'].mean()
# Ans:
print("About {:.4f} of songs are rock songs.".format(rock_ratio))


# 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
y = song.pop("RockOrNot").values
X = song.values
labels = song.columns
X_col = X.shape[1]
for i in range(X_col):
    mean_val = X[:, i].mean()
    std_val = X[:, i].std()
    X[:, i] = (X[:, i] - mean_val)/std_val


# Q3 Which variable in X has the largest correlation with y?
corr_list = list()
# Using Pearson correlation coefficient to determine which values of a column in X
# have the largest correlation with y
```

```python
62 for i in range(X_col):
63     corr_list.append(pearsonr(X[:, i], y)[0])
64
65 # find which values of a column in X have the largest negative correlation or
66 # positive correlation
67 index_var = 0
68 if max(corr_list) > abs(min(corr_list)):
69     index_var = np.argmax(corr_list)
70 else:
71     index_var = np.argmin(corr_list)
72
73 # Ans:
74 print("The variable in X has the largest correlation with y is {:s}".format(
75     labels[index_var]))
76
77
78 # Q4 When performing a classification problem, you fit the model to a portion of
79 # your data, and use the remaining data to determine how good the model fit was.
80 # Write a piece of code to divide X and y into training and test sets, use 75%
81 # of the data for training and keep 25% for testing. The data should be randomly
82 # selected, hence, you cannot simply take the first, say, 3000 rows. If you select
83 # rows 1,4,7,8,13,... of X for your training set, you must also select rows
84 # 1,4,7,8,13,... of y for training set. Additionally, the data in the training
85 # set cannot appear in the test set, and vice versa, so that when recombined,
86 # all data is accounted for. Use the seed 123 when generating random numbers
87 # Note: The data may not spilt equally into 75% and 25% portions. In this
88 # situation you should round to the nearest integer.
89
90 # Ans:
91 song = pd.read_csv("./tunedit_genres.csv")
92 # create a spliting function to split the original data set into 4 part(i.e.
   X_train, X_test, y_train, y_test)
93 # input: DataFrame, the size of training data set, random seed, the name of the
   classification label.
94 # output: X_train, X_test, y_train, y_test
95
96
97 def split_df(df, train_size, seed, target):
98     np.random.seed(seed)
99     # randomly shuffle the data set
100    rand_index = np.random.permutation(df.index)
101    rand_df = df.reindex(rand_index)
102    # set the last training row of a data frame according to the train_size
103    last_train_row = int(len(rand_df)*train_size)
104    # get the training data set and testing data set
105    train_df = rand_df.iloc[0:last_train_row]
106    test_df = rand_df.iloc[last_train_row:]
107    # spliting the data frame into X_train, X_test, y_train, y_test
108    y_train = train_df.pop(target).values
109    X_train = train_df.values
110    y_test = test_df.pop(target).values
111    X_test = test_df.values
112
113    return X_train, X_test, y_train, y_test
114
115
116 X_train, X_test, y_train, y_test = split_df(song, 0.75, 123, "RockOrNot")
117
118 # Q5 What is the percentage of rock songs in the training dataset and in the
119 # test dataset? Are they the same as the value found in Q1?
120
```

```python
121  # Ans:
122  print("About {:.4f} of the songs in training dataset are rock song.".format(
123      y_train.mean()))
124  print("About {:.4f} of the songs in testing dataset are rock song.".format(
125      y_test.mean()))
126  # They are not as same as the value found in Q1.
127
128
129  # Q6 Now we're going to write a function to run kNN on the data sets. kNN works
130  # by the following algorithm:
131  # 1) Choose a value of k (usually odd)
132  # 2) For each observation, find its k closest neighbours
133  # 3) Take the majority vote (mean) of these neighbours
134  # 4) Classify observation based on majority vote
135
136  # We're going to use standard Euclidean distance to find the distance between
137  # observations, defined as sqrt( (xi - xj)^T (xi-xj) )
138  # A useful short cut for this is the scipy functions pdist and squareform
139
140  # The function inputs are:
141  # - DataFrame X of explanatory variables
142  # - binary Series y of classification values
143  # - value of k (you can assume this is always an odd number)
144
145  # The function should produce:
146  # - Series y_star of predicted classification values
147
148
149  def kNN(X, y, k):
150      # Find the number of obsvervation
151      n = X.shape[0]
152      # Set up return values
153      y_star = list()
154      # Calculate the distance matrix for the observations in X
155      dist = squareform(pdist(X, 'euclidean'))
156      # Loop through each observation to create predictions
157      for i in range(n):
158          d = sorted(list(dist[i]))
159          sum = 0
160          # Find the y values of the k nearest neighbours
161          for j in range(1, k+1):
162              y_index = list(dist[i]).index(d[j])
163              sum += y[y_index]
164          # Now allocate to y_star
165          if (sum/k) > 0.5:
166              y_star.append(1)
167          else:
168              y_star.append(0)
169      return y_star
170
171
172  # Q7 The misclassification rate is the percentage of times the output of a
173  # classifier doesn't match the classification value. Calculate the
174  # misclassification rate of the kNN classifier for X_train and y_train, with k=3.
175  test1 = kNN(X_train, y_train, 3)
176  mis_count = 0
177  for i in range(len(test1)):
178      if (test1[i] != y_train[i]):
179          mis_count += 1
180  # Ans:
181  print("The misclassification rate is about {:.4f} ".format(
```

```python
182        mis_count/len(test1)))
183
184 # Q8 The best choice for k depends on the data. Write a function kNN_select that
185 # will run a kNN classification for a range of k values, and compute the
186 # misclassification rate for each.
187
188 # The function inputs are:
189 # - DataFrame X of explanatory variables
190 # - binary Series y of classification values
191 # - a list of k values k_vals
192
193 # The function should produce:
194 # - a Series mis_class_rates, indexed by k, with the misclassification rates for
195 # each k value in k_vals
196
197
198 def kNN_select(X, y, k_vals):
199     mis_class_rates = list()
200     for i in range(len(k_vals)):
201         test1 = kNN(X, y, k_vals[i])
202         mis_count = 0
203         for j in range(len(test1)):
204             if (test1[j] != y[j]):
205                 mis_count += 1
206         rate = mis_count/len(test1)
207         mis_class_rates.append(rate)
208     mis_class_rates = pd.Series(mis_class_rates, index=k_vals)
209     return mis_class_rates
210
211
212 # Q9 Run the function kNN_select on the training data for k = [1, 3, 5, 7, 9]
213 # and find the value of k with the best misclassification rate. Use the best
214 # value of k to report the mis-classification rate for the test data. What is
215 # the misclassification percentage with this k on the test set?
216 k = [1, 3, 5, 7, 9]
217 temp = kNN_select(X_train, y_train, k)
218 test2 = kNN(X_test, y_test, temp.idxmin())
219 mis_count = 0
220 for i in range(len(test2)):
221     if (test2[i] != y_test[i]):
222         mis_count += 1
223 # Ans:
224 print("The best k value for test data is {0:d} with about {1:.4f} misclassification
    rate.".format(
225     temp.idxmin(),  mis_count/len(test2)))
226
227 # Q10 Write a function to generalise the k nearest neighbours classification
228 # algorithm. The function should:
229 # - Separate out the classification variable for the other variables in the dataset,
230 #   i.e. create X and y.
231 # - Divide X and y into training and test set, where the number in each is
232 #   specified by 'percent_train'.
233 # - Run the k nearest neighbours classification on the training data, for a set
234 #   of k values, computing the mis-classification rate for each k
235 # - Find the k that gives the lowest mis-classification rate for the training data,
236 #   and hence, the classification with the best fit to the data.
237 # - Use the best k value to run the k nearest neighbours classification on the test
238 #   data, and calculate the mis-classification rate
239 # The function should return the mis-classification rate for a k nearest neighbours
240 # classification on the test data, using the best k value for the training data
241 # You can call the functions from Q6 and Q8 inside this function, provided they
```

```python
242  # generalise, i.e. will work for any dataset, not just the TunedIT dataset.
243
244
245  def kNN_classification(df, class_column, seed, percent_train, k_vals):
246      # df            - DataFrame to
247      # class_column  - column of df to be used as classification variable, should
248      #                   specified as a string
249      # seed          - seed value for creating the training/test sets
250      # percent_train - percentage of data to be used as training data
251      # k_vals        - set of k values to be tests for best classification
252      # Divide into training and test
253      X_train, X_test, y_train, y_test = split_df(
254          df, percent_train, seed, class_column)
255      # Compute the mis-classification rates for each for the values in k_vals
256      temp = kNN_select(X_train, y_train, k_vals)
257      # Find the best k value, by finding the minimum entry of mis_class_rates
258      k = temp.idxmin()
259      # Run the classification on the test set to see how well the 'best fit'
260      # classifier does on new data generated from the same source
261      y_star = kNN(X_test, y_test, k)
262      # Calculate the mis-classification rates for the test data
263      mis_count = 0
264      for i in range(len(y_star)):
265          if (y_star[i] != y_test[i]):
266              mis_count += 1
267      mis_class_test = mis_count/len(y_star)
268      return mis_class_test
269
270
271  # Test your function with the TunedIT data set, with class_column = 'RockOrNot',
272  # seed = the value from Q4, percent_train = 0.75, and k_vals = set of k values
273  # from Q8, and confirm that it gives the same answer as Q9.
274
275  k = [1, 3, 5, 7, 9]
276  song = pd.read_csv("./tunedit_genres.csv")
277  kNN_classification(song, "RockOrNot", 123, 0.75, k)
278  # The mis-classification rate for test set is 0.395, which is as same as Q9's.
279
280  # Now test your function with another dataset, to ensure that your code
281  # generalises. You can use the house_votes.csv dataset, with 'Party' as the
282  # classifier. Select the other parameters as you wish.
283  # This dataset contains the voting records of 435 congressman and women in the
284  # US House of Representatives. The parties are specified as 1 for democrat and 0
285  # for republican, and the votes are labelled as 1 for yes, -1 for no and 0 for
286  # abstained.
287  # Your kNN classifier should return a mis-classification for the test data (with
288  # the best fit k value) of ~8%.
289  df = pd.read_csv("./house_votes.csv")
290  mis_rate = []
291  seed = 123
292  for i in range(10):
293      rate = kNN_classification(df, "Party", seed, 0.75, k)
294      seed -= 1
295      mis_rate.append(rate)
296  pd.Series(mis_rate).mean()
297  # Because the misclassification rate is significantly influenced by the seed value
       and the way to split the data set,
298  # Then try 10 times with different seed values to get an average misclassification,
       which is about 8%
299
```