```python
# Week 11 - Assessed exercises

# In these assessed exercises. We're going to perform some model comparison on a
# handwriting recognition multi-class data set. We're going to divide it up into
# training, validation and test sets. We're going to run different parameter
# values on the training and validation sets to determine the optimal parameters
# Then we're going to run the optimal values on the test set to compare models

# The models we're going to use are:
# - Random forests
# - k nearest neighbours
# - Multi-layer perceptron (a type of neural network)
# You can load in these classifiers with the following commands
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.neural_network import MLPClassifier
from sklearn import metrics

# Some other packages we may need
from sklearn import datasets
import numpy as np
import numpy.random as npr

# Load in the digits data with
digits = datasets.load_digits()
# Remember that each sklearn data set comes with a target object (the response)
# and a data object (the explanatory variables). These data concern handwriting
# recognition so the response is a digit (0 to 9) and the explanatory variables
```

```python
29 # are levels of grey on an 8 by 8 grid.
30 # You can get a plot of any row (a handwriting sample) with:
31 import matplotlib.pyplot as plt
32
33 choose_row = 100
34 plt.gray()
35 plt.matshow(digits.images[choose_row])
36 plt.title(digits.target[choose_row])
37 plt.show()
38
39
40 # Where here I've made the title the digit it's supposed to represent (4).
41 # Looking at the plot you should see that it resembles a 4.
42 # Try changing the value of choose_row to see different digits and how they've
   been
43 # drawn. Note that this data set has an extra object 'images' that contains the 8
44 # by 8 matrices containing the pixel intensities, we will ignore this object.
45
46 # Below is a function for creating training, validation and test sets for a given
47 # matrix of observations X and vector of responses y. The function also needs a
48 # seed value so that it can reproduce the same outputs. The data is split 50%,
49 # 25%, 25% between training, validation and test, respectively. We will use this
50 # function when creating our training, validation and test sets below.
51 def train_val_test_sets(X, y, s):
52     npr.seed(s)
53     inds = npr.permutation(range(len(y)))
54     n_train = int(len(y) / 2)
55     n_val = int(3 * len(y) / 4)
```

```
56      X_train = X[inds[:n_train], :]
57      y_train = y[inds[:n_train]]
58      X_val = X[inds[n_train:n_val], :]
59      y_val = y[inds[n_train:n_val]]
60      X_test = X[inds[n_val:], :]
61      y_test = y[inds[n_val:]]
62      return X_train, X_val, X_test, y_train, y_val, y_test
63
64
65  # Q1 Write a function that runs each of the three classifiers with their default
66  # parameter values. The function inputs are the training and test sets X_train,
67  # X_test, y_train, y_test and a seed value s. The seed value should be used as
68  # the random_state argument in RandomForestClassifier and MLPClassifier. The
    function
69  # should return a dict with keys 'knn', 'rf' and 'svm'. The values should be the
70  # misclassification rate for each classifier (rounded to 3dp). Remember that
71  # there are more than two categories, so your mis-classification table will have
72  # more rows and columns to interpret.
73  def exercise1(X_train, X_test, y_train, y_test, s):
74      neigh = KNeighborsClassifier()
75      neigh.fit(X_train, y_train)
76      clf = RandomForestClassifier(random_state=s)
77      clf.fit(X_train, y_train)
78      mlp = MLPClassifier(random_state=s)
79      mlp.fit(X_train, y_train)
80      mis_rate = {'knn': round(1 - metrics.accuracy_score(y_test1, neigh.predict(
    X_test1)), 3)}
81      mis_rate['mlp'] = round(1 - metrics.accuracy_score(y_test1, mlp.predict(
```

```python
81 X_test1)), 3)
82     mis_rate['rf'] = round(1 - metrics.accuracy_score(y_test1, clf.predict(
   X_test1)), 3)
83     return mis_rate
84
85
86 # Suggested test
87 X1 = digits.data
88 y1 = digits.target
89 # We can use underscores to ignore the outputs of train_val_test_sets that we don
   't need
90 [X_train1, _, X_test1, y_train1, _, y_test1] = train_val_test_sets(X1, y1, 99)
91 print(exercise1(X_train1, X_test1, y_train1, y_test1, 123))
92
93
94 # This should return
95 # {'knn': 0.024, 'mlp': 0.031, 'rf': 0.076}
96 # You can ignore the warning messages or now
97 # Again, this should return the same answer every time you run it with the inputs
98 # X2, y2 and 99. If you use a subset of X2 and y2, or change the seed value you
99 # should expect these values to change.
100
101 # Each of the above models has key parameters which we might like to estimate.
    For
102 # example, we might want to estimate the 'best' number of neighbours to use in
    kNN
103 # To do this, we fit kNN with different values of k to the training set and
    evaluate
```

```python
104  # the performance of each model using the validation set. The k value that gives
     the
105  # best performance on the validation data is chosen as the best model. We then
106  # evaluate the performance of this model on data the classifier hasn't seen
     before,
107  # the test set.
108
109  # Q2 Write a function that determines the 'best' number of neighbours k to use in

110  # the kNN classifier and evaluates the performance of the best model on the test
111  # set. The function inputs are the training, validation and test sets and a list
     of
112  # values of k to try. The function should return a dict with the best k value (
     key:
113  # 'k') and the misclassification rate for the test set (key: 'MR') (rounded to
     3dp).
114  # Ensure that you use these exact keys.
115  def exercise2(X_train, X_val, X_test, y_train, y_val, y_test, kvals):
116      acc = []
117      for k in kvals:
118          neigh = KNeighborsClassifier(n_neighbors=k)
119          neigh.fit(X_train, y_train)
120          acc.append(round(metrics.accuracy_score(y_val, neigh.predict(X_val)), 3))
121      k_index = np.argmax(acc)
122      neigh = KNeighborsClassifier(n_neighbors=kvals[k_index])
123      neigh.fit(X_train, y_train)
124      mis_rate = round(1 - metrics.accuracy_score(y_test, neigh.predict(X_test)), 3
     )
```

```python
125         return {'k': kvals[k_index], 'MR': mis_rate}
126
127
128 # Suggestes test
129 print(exercise2(*train_val_test_sets(X1, y1, 199), range(1, 22)))
130
131
132 # This should return {'k': 2, 'MR': 0.031}
133 # If you change the seed value for creating your training, validation and test
    sets
134 # you can expect to get different values for k and the missclassification rate.
135
136 # Q3 Write a function that determines the 'best' number of trees (n_estimators)
    to
137 # use in the random forest classifier and evaluates the performance of the best
    model
138 # on the test set. The function inputs are the training, validation and test sets
    ,
139 # a list of values of n_estimators to try and a seed value s to use as the
    random_state
140 # for the classifier. The function should return a dict with the best number of
    trees
141 # (key: 'Trees') and the misclassification rate for the test set (key: 'MR') (
    rounded to 3dp).
142 # Ensure that you use these exact keys.
143 def exercise3(X_train, X_val, X_test, y_train, y_val, y_test, tree_vals, s):
144     acc = []
145     for tree in tree_vals:
```

```python
146        clf = RandomForestClassifier(n_estimators=tree, random_state=s)
147        clf.fit(X_train, y_train)
148        acc.append(round(metrics.accuracy_score(y_val, clf.predict(X_val)), 3))
149    tree_index = np.argmax(acc)
150    clf = RandomForestClassifier(n_estimators=tree_vals[tree_index], random_state
   =s)
151    clf.fit(X_train, y_train)
152    mis_rate = round(1 - metrics.accuracy_score(y_test, clf.predict(X_test)), 3)
153    return {'Trees': tree_vals[tree_index], 'MR': mis_rate}
154
155
156 # Suggestes test
157 print(exercise3(*train_val_test_sets(X1, y1, 99), range(5, 101, 5), 23))
158
159
160 # This should return {'Trees': 55, 'MR': 0.038}
161 # Again, changing the seed value for creating your training, validation and test
   sets
162 # will change the number of trees and the missclassification rate. As will
   changing
163 # the seed value for the random state of the classifier
164
165 # Q4 The parameter we wish to estimate for the multi-layer perceptron classifier
   is
166 # the number of neurons in the hidden layers of the neural network. To change
   this
167 # parameter include hidden_layer_sizes=num_neurons as an input to the
   MLPClassifier
```

```python
168 # function. Write a function that determines the 'best' number of neurons in the
169 # multi-layer perceptron classifier and evaluates the performance of the best
    model
170 # on the test set. The function inputs are the training, validation and test sets
    ,
171 # a list of values of hidden_layer_sizes to try and a seed value s to use as the
172 # random_state for the classifier. The function should return a dict with the
    best
173 # number of neurons (key: 'Neurons') and the  misclassification rate for the test

174 # set (key: 'MR') (rounded to 3dp).
175 # Ensure that you use these exact keys.
176 def exercise4(X_train, X_val, X_test, y_train, y_val, y_test, layer_vals, s):
177     acc = []
178     for layer in layer_vals:
179         mlp = MLPClassifier(hidden_layer_sizes=layer, random_state=s)
180         mlp.fit(X_train, y_train)
181         acc.append(round(metrics.accuracy_score(y_val, mlp.predict(X_val)), 3))
182     layer_index = np.argmax(acc)
183     mlp = MLPClassifier(hidden_layer_sizes=layer_vals[layer_index], random_state=
    s)
184     mlp.fit(X_train, y_train)
185     mis_rate = round(1 - metrics.accuracy_score(y_test, mlp.predict(X_test)), 3)
186     return {'Neurons': layer_vals[layer_index], 'MR': mis_rate}
187
188
189 # Suggested test
190 print(exercise4(*train_val_test_sets(X1, y1, 175), range(50, 1551, 100), 45))
```

```
191 # This should return {'Neurons': 550, 'MR': 0.033}
192 # As before, changing either seed value will change the number of neurons and the

193 # missclassification rate.
194 # Note that this function will take ~20s to run
195
```