

```

1 # Week 11 - Assessed exercises
2
3 # In these assessed exercises. We're going to perform some model comparison on a
4 # handwriting recognition multi-class data set. We're going to divide it up into
5 # training, validation and test sets. We're going to run different parameter
6 # values on the training and validation sets to determine the optimal parameters
7 # Then we're going to run the optimal values on the test set to compare models
8
9 # The models we're going to use are:
10 # - Random forests
11 # - k nearest neighbours
12 # - Multi-layer perceptron (a type of neural network)
13 # You can load in these classifiers with the following commands
14 import matplotlib.pyplot as plt
15 from sklearn.neighbors import KNeighborsClassifier
16 from sklearn.ensemble import RandomForestClassifier
17 from sklearn.neural_network import MLPClassifier
18
19 # Some other packages we may need
20 from sklearn import datasets
21 import numpy as np
22 import numpy.random as npr
23 import pandas as pd
24 from sklearn.metrics import accuracy_score
25
26 # Load in the digits data with
27 digits = datasets.load_digits()
28 # Remember that each sklearn data set comes with a target object (the response)
29 # and a data object (the explanatory variables). These data concern handwriting
30 # recognition so the response is a digit (0 to 9) and the explanatory variables
31 # are levels of grey on an 8 by 8 grid.
32 # You can get a plot of any row (a handwriting sample) with:
33 choose_row = 100
34 plt.gray()
35 plt.matshow(digits.images[choose_row])
36 plt.title(digits.target[choose_row])
37 # Where here I've made the title the digit it's supposed to represent (4).
38 # Looking at the plot you should see that it resembles a 4.
39 # Try changing the value of choose_row to see different digits and how they've been
40 # drawn. Note that this data set has an extra object 'images' that contains the 8
41 # by 8 matrices containing the pixel intensities, we will ignore this object.
42
43 # Below is a function for creating training, validation and test sets for a given
44 # matrix of observations X and vector of responses y. The function also needs a
45 # seed value so that it can reproduce the same outputs. The data is split 50%,
46 # 25%, 25% between training, validation and test, respectively. We will use this
47 # function when creating our training, validation and test sets below.
48
49
50 def train_val_test_sets(X, y, s):
51     npr.seed(s)
52     inds = npr.permutation(range(len(y)))
53     n_train = int(len(y)/2)
54     n_val = int(3*len(y)/4)
55     X_train = X[inds[:n_train], :]
56     y_train = y[inds[:n_train]]
57     X_val = X[inds[n_train:n_val], :]
58     y_val = y[inds[n_train:n_val]]
59     X_test = X[inds[n_val:], :]
60     y_test = y[inds[n_val:]]
61     return X_train, X_val, X_test, y_train, y_val, y_test

```

```

62
63 # Q1 Write a function that runs each of the three classifiers with their default
64 # parameter values. The function inputs are the training and test sets X_train,
65 # X_test, y_train, y_test and a seed value s. The seed value should be used as
66 # the random_state argument in RandomForestClassifier and MLPClassifier. The
  function
67 # should return a dict with keys 'knn', 'rf' and 'svm'. The values should be the
68 # misclassification rate for each classifier (rounded to 3dp). Remember that
69 # there are more than two categories, so your mis-classification table will have
70 # more rows and columns to interpret.
71
72
73 def exercise1(X_train, X_test, y_train, y_test, s):
74     # create each model based on each Classification algorithm
75     # mis_rate list is used to store every misclassification rate for each model
76     knn = KNeighborsClassifier()
77     rf = RandomForestClassifier(random_state=s)
78     mlp = MLPClassifier(random_state=s)
79     model = [knn, mlp, rf]
80     mis_rate = []
81     for m in model:
82         mm = m.fit(X_train, y_train)
83         y_pred = mm.predict(X_test)
84         temp = 1-accuracy_score(y_test, y_pred)
85         temp = round(temp, 3)
86         mis_rate.append(temp)
87     dic = {'knn': mis_rate[0], 'mlp': mis_rate[1], 'rf': mis_rate[2]}
88     return dic
89
90
91 # Suggested test
92 X1 = digits.data
93 y1 = digits.target
94 # We can use underscores to ignore the outputs of train_val_test_sets that we don't
  need
95 [X_train1, __, X_test1, y_train1, __, y_test1] = train_val_test_sets(X1, y1, 99)
96 print(exercise1(X_train1, X_test1, y_train1, y_test1, 123))
97 # This should return
98 # {'knn': 0.024, 'mlp': 0.031, 'rf': 0.076}
99 # You can ignore the warning messages or now
100 # Again, this should return the same answer every time you run it with the inputs
101 # X2, y2 and 99. If you use a subset of X2 and y2, or change the seed value you
102 # should expect these values to change.
103
104 # Each of the above models has key parameters which we might like to estimate. For
105 # example, we might want to estimate the 'best' number of neighbours to use in kNN
106 # To do this, we fit kNN with different values of k to the training set and evaluate
107 # the performance of each model using the validation set. The k value that gives the
108 # best performance on the validation data is chosen as the best model. We then
109 # evaluate the performance of this model on data the classifier hasn't seen before,
110 # the test set.
111
112 # Q2 Write a function that determines the 'best' number of neighbours k to use in
113 # the kNN classifier and evaluates the performance of the best model on the test
114 # set. The function inputs are the training, validation and test sets and a list of
115 # values of k to try. The function should return a dict with the best k value (key:
116 # 'k') and the misclassification rate for the test set (key: 'MR') (rounded to 3dp).
117 # Ensure that you use these exact keys.
118
119
120 def exercise2(X_train, X_val, X_test, y_train, y_val, y_test, kvals):

```

```

121 # mis_rate list is used to store every misclassification rate for each model
122 mis_rate = []
123 for k in kvals:
124     knn = KNeighborsClassifier(n_neighbors=k)
125     mm = knn.fit(X_train, y_train)
126     y_pred = mm.predict(X_val)
127     temp = 1-accuracy_score(y_val, y_pred)
128     temp = round(temp, 3)
129     mis_rate.append(temp)
130 mis_rate = pd.Series(mis_rate, index=kvals)
131 # choose the best k value to test it on the test data set
132 k_best = mis_rate.idxmin()
133 knn = KNeighborsClassifier(n_neighbors=k_best)
134 mm = knn.fit(X_train, y_train)
135 y_pred = mm.predict(X_test)
136 temp = 1-accuracy_score(y_test, y_pred)
137 temp = round(temp, 3)
138
139 dic = {'k': k_best, 'MR': temp}
140 return dic
141
142
143 # Suggestes test
144 print(exercise2(*train_val_test_sets(X1, y1, 199), range(1, 22)))
145 # This should return {'k': 2, 'MR': 0.031}
146 # If you change the seed value for creating your training, validation and test sets
147 # you can expect to get different values for k and the missclassification rate.
148
149 # Q3 Write a function that determines the 'best' number of trees (n_estimators) to
150 # use in the random forest classifier and evaluates the performance of the best
151 # model
152 # on the test set. The function inputs are the training, validation and test sets,
153 # a list of values of n_estimators to try and a seed value s to use as the
154 # random_state
155 # for the classifier. The function should return a dict with the best number of
156 # trees
157 # (key: 'Trees') and the misclassification rate for the test set (key: 'MR')
158 # (rounded to 3dp).
159 # Ensure that you use these exact keys.
160
161 def exercise3(X_train, X_val, X_test, y_train, y_val, y_test, tree_vals, s):
162     # mis_rate list is used to store every misclassification rate for each model
163     mis_rate = []
164     for t in tree_vals:
165         rf = RandomForestClassifier(random_state=s, n_estimators=t)
166         mm = rf.fit(X_train, y_train)
167         y_pred = mm.predict(X_val)
168         temp = 1-accuracy_score(y_val, y_pred)
169         temp = round(temp, 3)
170         mis_rate.append(temp)
171     mis_rate = pd.Series(mis_rate, index=tree_vals)
172     # choose the best value of trees and test it on test data set
173     t_best = mis_rate.idxmin()
174     rf = RandomForestClassifier(n_estimators=t_best, random_state=s)
175     mm = rf.fit(X_train, y_train)
176     y_pred = mm.predict(X_test)
177     temp = 1-accuracy_score(y_test, y_pred)
178     temp = round(temp, 3)
179
180     dic = {'Trees': t_best, 'MR': temp}

```

```

178     return dic
179
180
181 # Suggestes test
182 print(exercise3(*train_val_test_sets(X1, y1, 99), range(5, 101, 5), 23))
183 # This should return {'Trees': 55, 'MR': 0.038}
184 # Again, changing the seed value for creating your training, validation and test
    sets
185 # will change the number of trees and the missclassification rate. As will changing
186 # the seed value for the random state of the classifier
187
188 # Q4 The parameter we wish to estimate for the multi-layer perceptron classifier is
189 # the number of neurons in the hidden layers of the neural network. To change this
190 # parameter include hidden_layer_sizes=num_neurons as an input to the MLPClassifier
191 # function. Write a function that determines the 'best' number of neurons in the
192 # multi-layer perceptron classifier and evaluates the performance of the best model
193 # on the test set. The function inputs are the training, validation and test sets,
194 # a list of values of hidden_layer_sizes to try and a seed value s to use as the
195 # random_state for the classifier. The function should return a dict with the best
196 # number of neurons (key: 'Neurons') and the misclassification rate for the test
197 # set (key: 'MR') (rounded to 3dp).
198 # Ensure that you use these exact keys.
199
200
201 def exercise4(X_train, X_val, X_test, y_train, y_val, y_test, layer_vals, s):
202     # mis_rate list is used to store every misclassification rate for each model
203     mis_rate = []
204     for l in layer_vals:
205         mlp = MLPClassifier(random_state=s, hidden_layer_sizes=l)
206         mm = mlp.fit(X_train, y_train)
207         y_pred = mm.predict(X_val)
208         temp = 1-accuracy_score(y_val, y_pred)
209         temp = round(temp, 3)
210         mis_rate.append(temp)
211     mis_rate = pd.Series(mis_rate, index=layer_vals)
212     # choose the best number of hidden layer sizes and test it on the test data set
213     l_best = mis_rate.idxmin()
214     mlp = MLPClassifier(random_state=s, hidden_layer_sizes=l_best)
215     mm = mlp.fit(X_train, y_train)
216     y_pred = mm.predict(X_test)
217     temp = 1-accuracy_score(y_test, y_pred)
218     temp = round(temp, 3)
219
220     dic = {'Neurons': l_best, 'MR': temp}
221     return dic
222
223 # Suggested test
224 print(exercise4(*train_val_test_sets(X1, y1, 175), range(50, 1551, 100), 45))
225 # This should return {'Neurons': 550, 'MR': 0.033}
226 # As before, changing either seed value will change the number of neurons and the
227 # missclassification rate.
228 # Note that this function will take ~20s to run
229

```