```
03/23/18 02:15:25
/Users/caiglencross/Documents/MachineLearning/ps2/ps7/source/digits.py
```

```python
 1    """
 2    Author      : Cai Glencross & Katie Li
 3    Class       : HMC CS 158
 4    Date        : 2018 Mar 23
 5    Description : Bagging with Digits Dataset
 6                  This code was adapted from course material by Jenna
      Wiens (UMich)
 7    """
 8
 9    # python libraries
10    import collections
11
12    # numpy libraries
13    import numpy as np
14
15    # matplotlib libraries
16    import matplotlib as mpl
17    mpl.use('TkAgg')
18
19    import matplotlib.pyplot as plt
20
21    # scikit-learn libraries
22    from sklearn.datasets import load_digits
23    from sklearn.tree import DecisionTreeClassifier
24    from sklearn.ensemble import BaggingClassifier
25    from sklearn.ensemble import RandomForestClassifier
26    from sklearn.model_selection import train_test_split
27    from sklearn import metrics
28
29
30    ################################################################################
      ##
31    # bagging functions
32    ################################################################################
      ##
33
34    def bagging_ensemble(X_train, y_train, X_test, y_test,
      max_features=None, num_clf=10) :
35        """
36        Compute performance of bagging ensemble classifier.
37
38        Parameters
39        --------------------
40            X_train       -- numpy array of shape (n_train,d), training
      features
41            y_train       -- numpy array of shape (n_train,),  training
```

```
                 targets
42               X_test      -- numpy array of shape (n_test,d),  test
           features
43               y_test      -- numpy array of shape (n_test,),   test
           targets
44               max_features -- int, number of features to consider when
           looking for best split
45               num_clf      -- int, number of decision tree classifiers in
           bagging ensemble
46
47           Returns
48           -------------------
49               accuracy     -- float, accuracy of bagging ensemble
           classifier on test data
50           """
51           base_clf = DecisionTreeClassifier(criterion='entropy',
       max_features=max_features)
52           clf = BaggingClassifier(base_clf, n_estimators=num_clf)
53           clf.fit(X_train, y_train)
54           y_pred = clf.predict(X_test)
55           return metrics.accuracy_score(y_test, y_pred)
56
57
58      def random_forest(X_train, y_train, X_test, y_test, max_features,
       num_clf=10,
59                         bagging=bagging_ensemble) :
60           """
61           Wrapper around bagging_ensemble to use feature-limited decision
       trees.
62
63           Additional Parameters
64           -------------------
65               bagging      -- bagging_ensemble or bagging_ensemble2
66           """
67           return bagging(X_train, y_train, X_test, y_test,
68                          max_features=max_features, num_clf=num_clf)
69
70
71      def bagging_ensemble2(X_train, y_train, X_test, y_test,
       max_features=None, num_clf=10) :
72           """
73           Compute performance of bagging ensemble classifier.
74
75           You are allowed to use DecisionTreeClassifier but NOT
       BaggingClassifier.
76
77           Details
78           - Train num_clf base classifiers using bootstrap samples from
       X_train and y_train.
79               Use DecisionTreeClassifier with information gain as base
```

```
         classifier.
 80          Hints: Use np.random.choice(...) for bootstrap samples.
 81                Make sure to use same indices from X_train and y_train.
 82      - Predict using X_test and y_test.
 83        For each base classifier, track predictions on X_test.
 84        Make ensemble prediction using using majority vote.
 85      - Return accuracy compared to y_test.
 86
 87      Same parameters and return values as bagging_ensemble(...)
 88      """
 89
 90      n_train, d = X_train.shape
 91
 92      ### ========== TODO : START ========== ###
 93      # extra credit: implement bagging ensemble (see details above)
 94
 95      return 0.0
 96      ### ========== TODO : START ========== ###
 97
 98
 99  ######################################################################
     ##
100  # plotting functions
101  ######################################################################
     ##
102
103  def plot_scores(max_features, bagging_scores, random_forest_scores)
     :
104      """
105      Plot values in random_forest_scores and bagging_scores.
106      (The scores should use the same set of 100 different train and
     test set splits.)
107
108      Parameters
109      --------------------
110          max_features        -- list, number of features considered
     when looking for best split
111          bagging_scores       -- list, accuracies for bagging
     ensemble classifier using DTs
112          random_forest_scores -- list, accuracies for random forest
     classifier
113      """
114
115      plt.figure()
116      plt.plot(max_features, bagging_scores, '--', label='bagging')
117      plt.plot(max_features, random_forest_scores, '--', label='random
     forest')
118      plt.xlabel('max features considered per split')
119      plt.ylabel('accuracy')
120      plt.legend(loc='upper right')
```

```python
121        plt.show()
122
123
124    def plot_histograms(bagging_scores, random_forest_scores):
125        """
126        Plot histograms of values in random_forest_scores and
       bagging_scores.
127        (The scores should use the same set of 100 different train and
       test set splits.)
128
129        Parameters
130        --------------------
131            bagging_scores        -- list, accuracies for bagging
       ensemble classifier using DTs
132            random_forest_scores -- list, accuracies for random forest
       classifier
133        """
134
135        bins = np.linspace(0.8, 1.0, 100)
136        plt.figure()
137        plt.hist(bagging_scores, bins, alpha=0.5, label='bagging')
138        plt.hist(random_forest_scores, bins, alpha=0.5, label='random
       forest')
139        plt.xlabel('accuracy')
140        plt.ylabel('frequency')
141        plt.legend(loc='upper left')
142        plt.show()
143
144
145    ###########################################################################
       ##
146    # main
147    ###########################################################################
       ##
148
149    def main():
150        np.random.seed(1234)
151
152        # load digits dataset
153        digits = load_digits(4)
154        X = digits.data
155        y = digits.target
156
157        # evaluation parameters
158        num_trials = 100
159
160        # sklearn or home-grown bagging ensemble
161        bagging = bagging_ensemble
162
163        #========================================
```

```
164        # vary number of features
165
166        # calculate accuracy of bagging ensemble and random forest
167        #    for 100 random training and test set splits
168        # make sure to use same splits to enable proper comparison
169        max_features_vector = range(1,65,2)
170        bagging_scores = []
171        random_forest_scores = collections.defaultdict(list)
172        for i in range(num_trials):
173            print i
174            X_train, X_test, y_train, y_test = train_test_split(X, y,
       test_size=0.2)
175            bagging_scores.append(bagging(X_train, y_train, X_test,
       y_test))
176            for m in max_features_vector :
177                random_forest_scores[m].append(random_forest(X_train,
       y_train, X_test, y_test, m,

       bagging=bagging))
179
180        # analyze how performance of bagging and random forest changes
       with m
181        bagging_results = []
182        random_forest_results = []
183        for m in max_features_vector :
184            bagging_results.append(np.median(np.array(bagging_scores)))
185    #        print m, np.median(np.array(random_forest_scores[m]))
186
       random_forest_results.append(np.median(np.array(random_forest_scores
       [m])))
187        plot_scores(max_features_vector, bagging_results,
       random_forest_results)
188
189        #========================================
190        # plot histograms of performances for max_features=8
191        bagging_scores = []
192        random_forest_scores = []
193        for i in range(num_trials) :
194            X_train, X_test, y_train, y_test = train_test_split(X, y,
       test_size=0.2)
195            bagging_scores.append(bagging(X_train, y_train, X_test,
       y_test))
196            random_forest_scores.append(random_forest(X_train, y_train,
       X_test, y_test, 8,
197                                        bagging=bagging))
198        plot_histograms(bagging_scores, random_forest_scores)
199
200        ### ========== TODO : START ========== ###
201        # part d: determine pixel importance
202
```

```
203        print(X)
204
205        randomForest = RandomForestClassifier()
206        randomForest.fit(X, y)
207        importances = randomForest.feature_importances_
208        importances = importances.reshape(digits.images[0].shape)
209
210        # Plot pixel importances
211        plt.matshow(importances, cmap=plt.cm.hot)
212        plt.title("Pixel importances with forests of trees")
213        plt.colorbar()
214        plt.show()
215
216
217
218
219
220        ### ========== TODO : END ========== ###
221
222
223    if __name__ == "__main__" :
224        main()
```