

03/23/18 02:15:59

/Users/caiglencross/Documents/MachineLearning/ps2/ps7/source/soybean.py

```
1  """
2  Author      : Cai Glencross & Katie Li
3  Class       : HMC CS 158
4  Date        : 2018 Mar 23
5  Description : Multiclass Classification on Soybean Dataset
6              This code was adapted from course material by Tommi
   Jaakola (MIT)
7  """
8
9  # utilities
10 from util import *
11
12 # scikit-learn libraries
13 from sklearn.svm import SVC
14 from sklearn import metrics
15 import math
16
17
18 #####
19 ##
20 # output code functions
21 #####
22 ##
23
24 def generate_output_codes(num_classes, code_type) :
25     """
26     Generate output codes for multiclass classification.
27
28     For one-versus-all
29         num_classifiers = num_classes
30         Each binary task sets one class to +1 and the rest to -1.
31         R is ordered so that the positive class is along the
32         diagonal.
33
34     For one-versus-one
35         num_classifiers = num_classes choose 2
36         Each binary task sets one class to +1, another class to -1,
37         and the rest to 0.
38         R is ordered so that
39             the first class is positive and each following class is
40             successively negative
41             the second class is positive and each following class is
42             successively negative
43             etc
44
45     Parameters
```

```

40  -----
41      num_classes    -- int, number of classes
42      code_type      -- string, type of output code
43                      allowable: 'ova', 'ovo'
44
45  Returns
46  -----
47      R              -- numpy array of shape (num_classes,
num_classifiers),
48                      output code
49      """
50
51  ### ===== TODO : START ===== ###
52  # part a: generate output codes
53  # hint : initialize with np.ones(...) and np.zeros(...)
54  if (code_type == 'ova'):
55      R = -1* np.ones((num_classes, num_classes)) +
2*np.identity(num_classes)
56  if (code_type == 'ovo'):
57      #generate the correct number of rows
58      f = math.factorial
59      n = num_classes
60      #find n choose 2 classes
61      nC2_classes = f(n) / f(2) / f(n-2)
62      R = np.zeros((num_classes, nC2_classes))
63      current_col = 0
64      for i in range(num_classes):
65          for j in range(i+1,num_classes):
66              R[i, current_col] = 1
67              R[j, current_col] = -1
68              current_col = current_col + 1
69  ### ===== TODO : END ===== ###
70  return R
71
72
73  def load_code(filename) :
74      """
75      Load code from file.
76
77      Parameters
78      -----
79      filename -- string, filename
80      """
81
82      # determine filename
83      import util
84      dir = os.path.dirname(util.__file__)
85      f = os.path.join(dir, '..', 'data', filename)
86
87      # load data

```

```

88     with open(f, 'r') as fid :
89         data = np.loadtxt(fid, delimiter=",")
90
91     return data
92
93
94 def test_output_codes():
95     R_act = generate_output_codes(3, 'ova')
96     R_exp = np.array([[ 1, -1, -1],
97                      [-1, 1, -1],
98                      [-1, -1, 1]])
99     assert (R_exp == R_act).all(), "'ova' incorrect"
100
101     R_act = generate_output_codes(3, 'ovo')
102     R_exp = np.array([[ 1, 1, 0],
103                      [-1, 0, 1],
104                      [ 0, -1, -1]])
105     assert (R_exp == R_act).all(), "'ovo' incorrect"
106
107
108 #####
109 ##
110 # loss functions
111 #####
112 ##
113
114 def compute_losses(loss_type, R, discrim_func, alpha=2) :
115     """
116     Given output code and distances (for one example), compute
117     losses (for each class).
118
119     hamming : Loss = (1 - sign(z)) / 2
120     sigmoid : Loss = 1 / (1 + exp(alpha * z))
121     logistic : Loss = log(1 + exp(-alpha * z))
122
123     Parameters
124     -----
125     loss_type      -- string, loss function
126                      allowable: 'hamming', 'sigmoid', 'logistic'
127     R              -- numpy array of shape (num_classes,
128 num_classifiers)
129                      output code
130     discrim_func   -- numpy array of shape (num_classifiers,)
131                      distance of sample to hyperplanes, one per
132 classifier
133     alpha          -- float, parameter for sigmoid and logistic
134 functions
135
136     Returns
137     -----

```

```
132         losses      -- numpy array of shape (num_classes,), losses
133         """
134
135         # element-wise multiplication of matrices of shape (num_classes,
num_classifiers)
136         # tiled matrix created from (vertically) repeating discrim_func
num_classes times
137         z = R * np.tile(discrim_func, (R.shape[0],1))      # element-wise
138
139         # compute losses in matrix form
140         if loss_type == 'hamming' :
141             losses = np.abs(1 - np.sign(z)) * 0.5
142
143         elif loss_type == 'sigmoid' :
144             losses = 1./(1 + np.exp(alpha * z))
145
146         elif loss_type == 'logistic' :
147             # compute in this way to avoid numerical issues
148             # log(1 + exp(-alpha * z)) = -log(1 / (1 + exp(-alpha * z)))
149             eps = np.spacing(1) # numpy spacing(1) = matlab eps
150             val = 1./(1 + np.exp(-alpha * z))
151             losses = -np.log(val + eps)
152
153         else :
154             raise Exception("Error! Unknown loss function!")
155
156         # sum over losses of binary classifiers to determine loss for
each class
157         losses = np.sum(losses, 1) # sum over each row
158
159         return losses
160
161
162 def hamming_losses(R, discrim_func) :
163     """
164     Wrapper around compute_losses for hamming loss function.
165     """
166     return compute_losses('hamming', R, discrim_func)
167
168
169 def sigmoid_losses(R, discrim_func, alpha=2) :
170     """
171     Wrapper around compute_losses for sigmoid loss function.
172     """
173     return compute_losses('sigmoid', R, discrim_func, alpha)
174
175
176 def logistic_losses(R, discrim_func, alpha=2) :
177     """
178     Wrapper around compute_losses for logistic loss function.
```

```

179         """
180         return compute_losses('logistic', R, discrim_func, alpha)
181
182
183 #####
184 ##
185 # classes
186 #####
187 ##
188
189 class MulticlassSVM :
190     """
191     Multiclass SVM.
192
193     Attributes
194     -----
195         R          -- numpy array of shape (num_classes,
num_classifiers)
196         output code
197         svms       -- list of length num_classifiers
198                     binary classifiers, one for each column of R
199         classes    -- numpy array of shape (num_classes,) classes
200
201     Parameters
202     -----
203         R          -- numpy array of shape (num_classes,
num_classifiers)
204         output code
205         C          -- numpy array of shape (num_classifiers,1) or
float
206                     penalty parameter C of the error term
207         kernel     -- string, kernel type
208                     see SVC documentation
209         kwargs     -- additional named arguments to SVC
210     """
211
212     num_classes, num_classifiers = R.shape
213
214     # store output code
215     self.R = R
216
217     # use first value of C if dimension mismatch
218     try :
219         if len(C) != num_classifiers :
220             raise Warning("dimension mismatch between R and C ")
+
221                                     "=> using first value in C")
222         C = np.ones((num_classifiers,)) * C[0]

```

```

223     except :
224         C = np.ones((num_classifiers,)) * C
225
226     # set up and store classifier corresponding to jth column of
R
227     self.svms = [None for _ in xrange(num_classifiers)]
228     for j in xrange(num_classifiers) :
229         svm = SVC(kernel=kernel, C=C[j], **kwargs)
230         self.svms[j] = svm
231
232
233     def fit(self, X, y) :
234         """
235         Learn the multiclass classifier (based on SVMs).
236
237         Parameters
238         -----
239             X      -- numpy array of shape (n,d), features
240             y      -- numpy array of shape (n,), targets
241
242         Returns
243         -----
244             self -- an instance of self
245         """
246
247         classes = np.unique(y)
248         num_classes, num_classifiers = self.R.shape
249         if len(classes) != num_classes :
250             raise Exception('num_classes mismatched between R and
data')
251         self.classes = classes      # keep track for prediction
252
253         ### ===== TODO : START ===== ###
254         # part c: train binary classifiers
255
256         # HERE IS ONE WAY (THERE MAY BE OTHER APPROACHES)
257         #
258         # keep two lists, pos_ndx and neg_ndx, that store indices
259         #   of examples to classify as pos / neg for current binary
task
260         #
261         # for each class C
262         # a) find indices for which examples have class equal to C
263         #     [use np.nonzero(CONDITION)[0]]
264         # b) update pos_ndx and neg_ndx based on output code R[i,j]
265         #     where i = class index, j = classifier index
266         #
267         # set X_train using X with pos_ndx and neg_ndx
268         # set y_train using y with pos_ndx and neg_ndx
269         #     y_train should contain only {+1,-1}

```

```

270 #
271 # train the binary classifier
272 n,d = X.shape
273 R = self.R
274
275 for i in range(0, len(self.svms)):
276     pos_ndx = []
277     neg_ndx = []
278     for ndx in range(0, n):
279         class_index = list(classes).index(y[ndx])
280         if R[class_index, i] == 1:
281             pos_ndx.append(ndx)
282         elif R[class_index, i] == -1:
283             neg_ndx.append(ndx)
284
285     new_X_pos = X[pos_ndx,:]
286     new_X_neg = X[neg_ndx,:]
287
288     #get the new_X with labels
289     new_X = np.vstack((new_X_pos, new_X_neg))
290     new_Y = np.append(
291         np.ones(len(pos_ndx)),
292         (-1*np.ones(len(neg_ndx))) )
293
294     self.svms[i].fit(new_X, new_Y)
295     return self
296     ### ===== TODO : END ===== ###
297
298 def predict(self, X, loss_func=hamming_losses) :
299     """
300     Predict the optimal class.
301
302     Parameters
303     -----
304     X          -- numpy array of shape (n,d), features
305     loss_func  -- loss function
306                  allowable: hamming_losses, logistic_losses,
sigmoid_losses
307
308     Returns
309     -----
310     y          -- numpy array of shape (n,), predictions
311     """
312
313     n,d = X.shape
314     num_classes, num_classifiers = self.R.shape
315
316     # setup predictions
317     y = np.zeros(n)

```

```

318
319     ### ===== TODO : START ===== ###
320     # part d: predict multiclass class
321     #
322     # HERE IS ONE WAY (THERE MAY BE OTHER APPROACHES)
323     #
324     # for each example
325     #     predict distances to hyperplanes using
SVC.decision_function(...)
326     #     find class with minimum loss (be sure to look up in
self.classes)
327     #
328     # if you have a choice between multiple occurrences of the
minimum values,
329     # use the index corresponding to the first occurrence
330     for i in range(n):
331         dists = np.empty(num_classifiers)
332         for j in range(num_classifiers):
333             X_i_resaped = X[i,:].reshape(1,-1)
334             dist = self.svms[j].decision_function(X_i_resaped)
335
336             dists[j] = dist
337
338
339             losses = loss_func(self.R, dists)
340             best_class_index = np.argmin(losses)
341             y[i] = self.classes[best_class_index]
342
343     ### ===== TODO : END ===== ###
344
345     return y
346
347
348     #####
349     ##
350     #####
351     ##
352     def main() :
353         # load data
354         converters = {35: ord} # label (column 35) is a character
355         train_data = load_data("soybean_train.csv", converters)
356         test_data = load_data("soybean_test.csv", converters)
357         num_classes = 15
358
359         # part b : generate output codes
360         test_output_codes()
361
362         # plot loss functions

```



```

363 z = np.arange(-2, 3, 0.01)
364 hamming = map(lambda u: (1 - np.sign(u))/2, z)
365 sigmoid1 = map(lambda u: (1/ (1+np.exp(u))), z)
366 sigmoid2 = map(lambda u: (1/ (1+np.exp(2*u))), z)
367 logistic1 = map(lambda u: (math.log(1 + np.exp(-u))), z)
368 logistic2 = map(lambda u: (math.log(1 + np.exp(-2*u))), z)
369
370 plt.plot(z, hamming, label="Hamming")
371 plt.plot(z, sigmoid1, label="Sigmoid1")
372 plt.plot(z, sigmoid2, label="Sigmoid2")
373 plt.plot(z, logistic1, label="Logistic1")
374 plt.plot(z, logistic2, label="Logistic2")
375 plt.legend()
376 #plt.show()
377
378
379 ### ===== TODO : START ===== ###
380 # parts c-e : train component classifiers, make predictions,
381 #             compare output codes and loss functions
382 #
383 # use generate_output_codes(...) to generate OVA and OVO codes
384 # use load_code(...) to load random codes
385 #
386 # for each output code and loss function
387 #   train a multiclass SVM on training data and evaluate on test
data
388 #   setup the binary classifiers using the specified parameters
from the handout
389 #
390 # if you implemented MulticlassSVM.fit(...) correctly,
391 #   using OVA, your first trained binary classifier
392 #   should have the following indices for support vectors
393 #   array([ 12,  22,  29,  37,  41,  44,  49,  55,  76, 134,
394 #           157, 161, 167, 168,    0,    3,    7])
395 #
396 # if you implemented MulticlassSVM.predict(...) correctly,
397 #   using OVA and Hamming loss, you should find 54 errors
398
399 ova = generate_output_codes(len(np.unique(train_data.y)), 'ova')
400
401 multiclass = MulticlassSVM(ova, C=10, kernel='poly', gamma=1,
degree=4, coef0=1)
402 multiclass = multiclass.fit(train_data.X, train_data.y)
403
404 predictions = multiclass.predict(test_data.X,
loss_func=hamming_losses)
405
406 errors = metrics.zero_one_loss(predictions, test_data.y,
normalize=False)
407 print(errors)

```

```
408
409     print(multiclass.svms[0].support_)
410     ### ===== TODO : END ===== ###
411
412 if __name__ == "__main__" :
413     main()
```