# United Way MLPP: MultinomialNB Optical Digits Example

```python
0   # -*- coding: utf-8 -*-
1
2
3   """
4   Title:                  Sample MNBayes1
5   Project Descriptor:     Machine Learning Pilot Project (United Way)
6   Project ID:             2016SoE013 (STAT_570_Consulting)
7   Record:
8   Author:                 bmarron
9   Origin Date:            28 May 2016
10  """
11
12  """
13  This script provides an example of the use of the Multinomial Bayes
14  algorithm as applied to the  Optical Recognition of Handwritten Digits
15  dataset. The original source of the dataset:
16          E. Alpaydin, C. Kaynak
17          Department of Computer Engineering
18          Bogazici University, 80815 Istanbul Turkey
19          alpaydin@boun.edu.tr
20          July 1998
21
22  The dataset is available at UC Irivine Machine Learning Repository.
23  """
24
25  """
26  Important websites:
27      #UCI Machine Learning Repository
28  http://archive.ics.uci.edu/ml/
29
30      #Python machine learning tools
31  http://scikit-learn.org/
32  """
33
34  #%% Import packages
35
36
37  import numpy as np
38  import cPickle
39  import urllib
40  import csv
41  import operator
42  import pylab as pl
43  import pandas as pd
44  from time import time
45  from sklearn import cross_validation
46  from sklearn.utils.fixes import bincount
47  from sklearn.metrics import accuracy_score
48  from sklearn.metrics import confusion_matrix
49  from sklearn.metrics import classification_report
50  from sklearn.naive_bayes import MultinomialNB
51  from sklearn.svm import SVC
52  from sklearn.svm import LinearSVC
53
54  #%%
55
56  #Download and save the raw, UCI opttical digits data and data info sheet
57
58  '''
59  http://archive.ics.uci.edu/ml/datasets/Optical+Recognition+of+Handwritten+Digits
60  http://archive.ics.uci.edu/ml/machine-learning-databases/optdigits/optdigits.tra
61  http://archive.ics.uci.edu/ml/machine-learning-databases/optdigits/optdigits.tes
62  '''
63
```

```python
80
81  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=.5,
        random_state=random_state)
82
83  #%% Confirm data content with info sheet about training dataset
84
85  #Data sheet: 0=376, 1=389, 2=380, 9=382
86
87
88  test0 = np.where(raw_optdigits_tr[:, 64] == 0)
89  test0 = list(test0[0])
90
91  test1 = np.where(raw_optdigits_tr[:, 64] == 1)
92  test1 = list(test1[0])
93
94  test2 = np.where(raw_optdigits_tr[:, 64] == 2)
95  test2 = list(test2[0])
96
97  test9 = np.where(raw_optdigits_tr[:, 64] == 9)
98  test9 = list(test9[0])
99
100 print "zeros: %s" %(len(test0))
101 print "ones: %s" %(len(test1))
102 print "twos: %s" %(len(test2))
103 print "nines: %s" %(len(test9))
104
105 #%% Split training dataset(tr_d) and test dataset(te_d) into features (X)
106
107 #and class (y)
108
109
110 tr_d_X, tr_d_y = raw_optdigits_tr[:,0:64], raw_optdigits_tr[:, 64]
111 te_d_X, te_d_y = raw_optdigits_te[:,0:64], raw_optdigits_te[:, 64]
112
113
114 #%% Save the split (processed) data
115
116
117 # Save tr_d_X
118
119 with open("tr_d_X.pkl", 'wb') as f:
120     cPickle.dump(tr_d_X, f, protocol=2)
121
122 # Save tr_d_y
123
124 with open("tr_d_y.pkl", 'wb') as f:
125     cPickle.dump(tr_d_y, f, protocol=2)
126
127 # Save te_d_X
128
129 with open("te_d_X.pkl", 'wb') as f:
130     cPickle.dump(te_d_X, f, protocol=2)
131
132 # Save te_d_y
133
134 with open("te_d_y.pkl", 'wb') as f:
135     cPickle.dump(te_d_y, f, protocol=2)
136
137 #%%
138
```

```python
139  #Reload the processed data into Python, if needed
140
141  tr_d_X = cPickle.load(open("/home/bmarron/Desktop/tr_d_X.pkl","rb"))
142
143  #Reload the processed data
144
145  tr_d_y = cPickle.load(open("/home/bmarron/Desktop/tr_d_y.pkl","rb"))
146
147  #Reload the processed data
148
149  te_d_X = cPickle.load(open("/home/bmarron/Desktop/te_d_X.pkl","rb"))
150
151  #Reload the processed data
152
153  te_d_y = cPickle.load(open("/home/bmarron/Desktop/te_d_y.pkl","rb"))
154
155
156  #%% Based on the method  at
157
158  #http://scikit-learn.org/stable/auto_examples/text/
159      mlcomp_sparse_document_classification.html#example-text-mlcomp-sparse-
160      document-classification-py
159
160  # Benchmark classifiers
161
162
163  def benchmark(clf_class, params, name):
164      print("parameters:", params)
165      t0 = time()
166      clf = clf_class(**params).fit(tr_d_X, tr_d_y)
167      print("done in %fs" % (time() - t0))
168
169      if hasattr(clf, 'coef_'):
170          print("Percentage of non zeros coef: %f"
171                % (np.mean(clf.coef_ != 0) * 100))
172          # .coef_ output is a read-only file
173
174          # Python does not define a read-write variable!
175
176          # Export the wgt vector data from .coef_
177
178          wgts = clf.coef_
179          with open('wgts_MNBayes.csv', 'w+') as f:
180              a= csv.writer(f, delimiter=',')
181              a.writerows(wgts)
182
183
184      print("Predicting the outcomes of the testing set")
185      t0 = time()
186      pred = clf.predict(te_d_X)
187      print("done in %fs" % (time() - t0))
188
189      print("Classification report on test set for classifier:")
190      print(clf)
191      print()
192      print(classification_report(te_d_y, pred))
193
194      cm = confusion_matrix(te_d_y, pred)
195      df = pd.DataFrame(cm)
196      print("Confusion  matrix:")
```

```python
197         print(cm)
198         print df.to_latex()
199
200
201         print ("Accuracy:")
202         print accuracy_score(te_d_y, pred)
203
204         # Show confusion matrix
205
206         pl.matshow(cm)
207         pl.title('Confusion matrix of the\n %s classifier \n (alpha=0.03)' % name,
                fontsize=12, y=1.10)
208         pl.colorbar()
209         pl.savefig("MNBayes1.pdf")
210
211 #%%
212
213 print("Testbenching a MultinomialNB classifier...")
214 parameters = {'alpha': 0.03}
215
216 benchmark(MultinomialNB, parameters, 'MultinomialNB')
217
218 pl.show()
219
220 #%%
221
222 # coef_ is empirical log probability of features given a class, P(x_i|y)
223
224 # coef_ in MultinomialNB mirrors the feature_log_prob_ ; useful
225
226 #for interpreting MultinomialNB as a linear model
227
228
229 # Import 'wgts_MNBayes.csv' with Spyder (green arrow); import as 'array'
230
231 # Take the absolute value of the wgts
232
233 # Define a list of feature numbers (0 to 63); actual features run i=1,2,...63
234
235
236 ###########
237
238 #import 'wgts_MNBayes.csv' with Spyder (green arrow)
239
240 ###########
241
242 abs_csv = abs(wgts_MNBayescsv)
243 features =np.arange(64)
244
245 # Create tuples (feature #, feature value)
246
247 wgts_features=[]
248 for i in range(10):
249     x=zip(features, abs_csv[i])
250     wgts_features.append(x)
251
252 # Reverse sort on the the second element in each tuple
253
254 for i in range(10):
255     wgts_features[i]=sorted(wgts_features[i], key=operator.itemgetter(1),
```

```python
          reverse=True)

256
257   df2 = pd.DataFrame(wgts_features)
258   df2.to_csv('sorted_wgts_features_MNBayes.csv')
259
260   #%% Ouput w/o panda dataframe
261
262
263   with open('sorted_wgts_features_MNBayes.csv', 'w+') as f:
264             a= csv.writer(f, delimiter=',')
265             a.writerows(wgts_features)
```