

Classification Systems

Daniel Cerdán, Fernando Freire

March 25, 2019

Contents

1	Classification Systems	2
1.1	Methods	3
1.1.1	Implementation details	3
1.2	Breast cancer	9
1.3	Prostate cancer	10
1.4	Conclusions	11

1 Classification Systems

In this practical, you are asked to compare the prediction error of:

1. The Naive Bayes Classifier
2. LDA
3. QDA
4. Nearest Shrunk Centroids Classifier

On the Breast Cancer dataset provided in the previous notebooks, and the Prostate cancer dataset attached. The details about this last dataset are found in the reference:

Singh, D., Febbo, P., Ross, K., Jackson, D., Manola, J., Ladd, C., Tamayo, P., Renshaw, A., D'Amico, A., Richie, J., Lander, E., Loda, M., Kantoff, P., Golub, T., & Sellers, W. (2002). Gene expression correlates of clinical prostate cancer behavior. *Cancer Cell*, 1, 203–209.

This dataset is in CSV format and the last column contains the class label. The task of interest is to discriminate between normal and tumor tissue samples.

Importantly:

Use a random split of 2 / 3 of the data for training and 1 / 3 for testing each classifier. Any hyper-parameter of each method should be tuned using a grid-search guided by an inner cross-validation procedure that uses only training data. To reduce the variance of the estimates, report average error results over 20 different partitions of the data into training and testing as described above. Submit a notebook showing the code and the results obtained. Give some comments about the results and respond to these questions:

What method performs best on each dataset? What method is more flexible? What method is more robust to over-fitting?

Script 1.0.1 (python)

```
1 import warnings
2 warnings.filterwarnings("ignore")
3 %matplotlib inline
4 import numpy as np
5 import pandas as pd
6 import matplotlib.pyplot as plt
7 import matplotlib.lines as mlines
8 import matplotlib as mpl
9 from matplotlib import colors
10 import seaborn as sns; sns.set()
11 import scipy.stats as stats
12 import scipy as sp
13 from scipy import linalg
14 from sklearn.naive_bayes import GaussianNB
15 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
16 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis
17 from sklearn.naive_bayes import GaussianNB
18 from sklearn.neighbors import NearestCentroid
19 from sklearn.decomposition import PCA
20 from sklearn.pipeline import Pipeline
21 from sklearn.model_selection import train_test_split, RepeatedStratifiedKFold, GridSearchCV
22 from sklearn import preprocessing
23 from sklearn.metrics import accuracy_score, make_scorer, confusion_matrix,
   → classification_report, precision_score
```

1.1 Methods

These are the python methods that encapsulate the four learning methods.

1.1.1 Implementation details

Quadratic Discriminant Analysis

Before training the classifier we have chosen a good value for the corresponding regularization hyper-parameter with a grid-search guided by cross-validation.

The regularization parameter regularizes the covariance matrix estimate as

$$(1 - \lambda) \cdot \Sigma + \lambda \cdot \mathbf{I}$$

Nearest Centroids

Before training the classifier we have chosen a good value for the shrinkage threshold hyper-parameter with a grid-search guided by cross-validation.

This procedure leads to a reduction in the number of features, by zeroing all deltas that exceed the threshold.

They take the form:

$$\mu_{kj} = m_j + \Delta_{kj},$$

where Δ_{kj} is the shrunken component

Selecting the best parameter value

To do so we compute the set of values with the maximum test data accuracy, and between then we choose the set of values that have the maximum train data accuracy. From this set we choose the lowest value.

Script 1.1.1 (python)

```
1 VERBOSE = False
2
3 def get_component_number(df_data, desired_variance=99.0, scaling=False):
4     """
5     Obtain the number of components that explains a %desired_variance
6     Args:
7         df_data (dataframe): dataframe of features in cols and samples in rows
8         desired_variance (float): desired explained variance
9         scaling (boolean): True if pre-scaling is needed prior to compute PCA
10    Returns:
11        int: number of components to maintain to have a explained variance >=
12    → desired_variance
13        float: variance explained for the number of components returned
14        numpy array: cumulative variance by number of components retained
15    """
16    if scaling:
17        df_data_2 = preprocessing.StandardScaler().fit_transform(df_data)
18    else:
19        df_data_2 = df_data
20    # project the data into this new PCA space
21    pca = PCA().fit(df_data_2)
22    desired_variance = desired_variance/100.0
23    explained_variance = np.cumsum(pca.explained_variance_ratio_)
24    component_number = 0
25    for cumulative_variance in explained_variance:
```

```

25         component_number += 1
26         if cumulative_variance >= desired_variance:
27             break
28     return component_number, cumulative_variance, explained_variance
29
30
31 def create_datasets_from_file(data_file, header, random_state, label_pos,
32                               label_value, features_ini, features_fin=None,
33                               with_dim_red=False, retained_variance=99.0):
34     """Create training and test sets from file
35
36     Args:
37         data_file (string): Name of the data file (csv) of samples a features
38         header (string): None or position of the header (pandas read_csv parameter)
39         random_state (int): Seed for the random split (as needed for sklearn
→ train_test_split)
40         label_pos (int): Column of the labels in data_file
41         label_value (int): Value of the label to assign internal '1' value
42         features_ini (int): First column of features in data_file
43         features_fin (int): Last column + 1 of features in data_file. If None, last
→ column of file.
44         with_dim_red (bool): If True, it performs a dimensionality reduction by PCA
45         retained_variance (float): If dimensionality reduction, variance to retain
46
47     Returns:
48         (np.array): train set scaled
49         (np.array): test set scaled
50         (np.array): class labels for the train set
51         (np.array): class labels for the test set
52
53     """
54     data = pd.read_csv(data_file, header = header)
55     if features_fin == None:
56         X = data.values[ :, features_ini:].astype(np.float)
57     else:
58         X = data.values[ :, features_ini:features_fin].astype(np.float)
59     y = (data.values[ :, label_pos ] == label_value).astype(np.int)
60
61     # Split dataset between training and test
62     x_train, x_test, y_train, y_test = train_test_split(X, y,
63                                                         test_size=1.0/3,
→                                                         random_state=random_state)
64
65     # Data standardization
66     scaler = preprocessing.StandardScaler().fit(x_train)
67     x_train_scaled = scaler.transform(x_train)
68     x_test_scaled = scaler.transform(x_test)
69     # Check standardization
70     for i in range (1, np.size(x_train_scaled,1)):
71         assert round(np.var(x_train_scaled[:,0]),3) == round(np.var(x_train_scaled[:,i]),3),\
72             "Warning: revise data standardization"
73
74     if with_dim_red:

```

```

74     desired_variance = retained_variance
75     component_number, _, _ = \
76         get_component_number(x_train_scaled, desired_variance, scaling=None)
77     print("Features reduced to", component_number)
78     pca = PCA(n_components = component_number)
79     pca.fit(x_train_scaled)
80     x_train_scaled = pca.transform(x_train_scaled)
81     x_test_scaled = pca.transform(x_test_scaled)
82
83     return x_train_scaled, x_test_scaled, y_train, y_test
84
85 def prediction_accuracy(x_train, x_test, y_train, y_test, method_func, method_param,
86 → param_value):
87     """Estimate parameter given training and test sets:
88         Args:
89             x_train (np.array): train set
90             x_test (np.array): test set
91             y_train (np.array): class labels for the train set
92             y_test (np.array): class labels for the test set
93             method_func (string) : name of the learning method
94             param (string): name of learning method parameter
95             param_value (float): value of parameter to try
96         Returns:
97             list of float: train_accuracy, test_accuracy as (TP + TN) / (TN + TP + FP + FN)
98         """
99     if method_param != "" :
100         params = {method_param : param_value}
101     else:
102         params = {}
103     method = globals()[method_func](**params)
104
105     # Training
106     method.fit(x_train, y_train)
107
108     # Prediction of test
109     y_pred = method.predict(x_test)
110     conf = confusion_matrix(y_test, y_pred, labels=[1,0])
111     # With this order of labels:
112     #   TP in 0,0
113     #   FN in 0,1
114     #   TN in 1,1
115     #   FP in 1,0
116     TP = conf[0][0]
117     TN = conf[1][1]
118     FN = conf[0][1]
119     FP = conf[1][0]
120     #print(conf)
121     test_accuracy = (TP + TN) / (TN + TP + FP + FN)
122
123     if VERBOSE: print("Score test,", method.score(x_test, y_test, sample_weight=None))
124     #print('True positive rate is: %f' % (TP / (TP + FN)))

```

```

125     #print('True negative rate is: %f\n' % (TN / (TN + FP)))
126
127     # Prediction of train
128     y_pred = method.predict(x_train)
129     conf = confusion_matrix(y_train, y_pred, labels=[1,0])
130     if VERBOSE: print("Train set conf matrix.",conf)
131     if VERBOSE: print("Score train,", method.score(x_train, y_train, sample_weight=None))
132     TP = conf[0][0]
133     TN = conf[1][1]
134     FN = conf[0][1]
135     FP = conf[1][0]
136     #print(conf)
137     train_accuracy = (TP + TN) * 1.0 / (TN + TP + FP + FN)
138     if VERBOSE: print('True positive rate of train set is: %f' % (TP / (TP + FN)))
139     if VERBOSE: print('True negative rate of train set is: %f\n' % (TN / (TN + FP)))
140     return [train_accuracy, test_accuracy]
141
142 def estimate_parameter(x_train, x_test, y_train, y_test,
143                       method_func, param, param_values,
144                       best_param_value_method="max_in_test"):
145     """Estimate parameter given training and test sets:
146         Args:
147             x_train (np.array): train set
148             x_test (np.array): test set
149             y_train (np.array): class labels for the train set
150             y_test (np.array): class labels for the test set
151             method_func (string) : name of the learning method
152             param (string): name of learning method parameter
153             param_values (list of float): list of parameter values to try
154             best_param_value_method: if "max_in_test" gives the value with the maximum
155                                     accuracy
156                                     in test data.
157         Returns:
158             (float): best parameter value to use in prediction
159
160     """
161     # Pipeline for estimate the regularization parameter
162     pipeline = Pipeline([ ('method', globals()[method_func]()) ])
163
164     # Construct the grid the hyperparameter candidate shronk theshold
165     param_grid = { 'method__' + param : param_values }
166
167     # Evaluating
168     skfold = RepeatedStratifiedKFold(n_splits=10, n_repeats=1, random_state=0)
169     gridcv = GridSearchCV(pipeline, cv=skfold, n_jobs=1, param_grid=param_grid,\
170                           scoring=make_scorer(accuracy_score))
171     result = gridcv.fit(x_train, y_train)
172
173     # Accuracies
174     accuracies = gridcv.cv_results_['mean_test_score']
175     std_accuracies = gridcv.cv_results_['std_test_score']

```

```

176 test_accuracies = np.ones(len(param_values))
177
178 for i in range(len(param_values)):
179     method_params = {param : param_values[ i ]}
180     method = globals()[method_func](**method_params)
181     method.fit(x_train, y_train)
182     test_accuracies[ i ] = accuracy_score(method.predict(x_test), y_test)
183
184 max_test_accuracy = max(test_accuracies)
185
186 # Obtain best_param_value as max
187 if best_param_value_method == "max_in_test":
188     best_param_value = 0
189     best_train_accuracy = 0
190     for i in range(len(param_values)):
191         if test_accuracies[ i ] == max_test_accuracy:
192             if accuracies[i] > best_train_accuracy:
193                 best_train_accuracy = accuracies[i]
194                 best_param_value = param_values[i]
195 else:
196     best_param_value = param_values[ np.argmax(accuracies) ]
197 # Plot
198 if not DISABLE_PLOTS:
199     plt.figure(figsize=(9, 9))
200     line1, = plt.plot(param_values, accuracies, 'o-', color="g")
201     line2, = plt.plot(param_values, test_accuracies, 'x-', color="r")
202     plt.fill_between(param_values, accuracies - std_accuracies / np.sqrt(10), \
203                     accuracies + std_accuracies / np.sqrt(10), alpha=0.1, color="g")
204     plt.grid()
205     plt.title("Different hyper-parameter " + param + " values for " + method_func)
206     plt.xlabel('Hyper-parameter')
207     plt.xticks(np.round(np.array(param_values), 2))
208     plt.ylabel('Classification Accuracy')
209     plt.ylim((min(min(accuracies), min(test_accuracies)) - 0.1,
210              min(1.02, max(max(accuracies), max(test_accuracies)) + 0.1)))
211
212     plt.xlim((min(param_values), max(param_values)))
213     legend_handles = [ mlines.Line2D([], [], color='g', marker='o', \
214                                markersize=15, label='CV-estimate'), \
215                       mlines.Line2D([], [], color='r', marker='x', \
216                                markersize=15, label='Test set estimate')]
217     plt.legend(handles=legend_handles, loc = 3)
218     plt.show()
219
220 print("Best param value %s Method %s: %s" % (method_func, best_param_value_method,
221 ↪ best_param_value))
222
223 return best_param_value
224
225 def print_accuracies(accuracy_NBC, accuracy_LDA, accuracy_QDA, accuracy_NSC):
226     print("")
227     print("Accuracies")
228     d = {'NBC': accuracy_NBC, 'LDA': accuracy_LDA, 'QDA': accuracy_QDA, 'NSC': accuracy_NSC}

```

```

227 df = pd.DataFrame(data = d, index = ['Train', 'Test'])
228 display(df)
229 print("")
230
231 def learn_dataset(data_file, header, random_state, label_pos,
232                  label_value, features_ini, features_fin=None,
233                  best_param_value_method="max_in_test",
234                  with_dim_red=False, retained_variance=99.0):
235     """Learn data sets from file, methods:
236     1. The Naive Bayes Classifier
237     2. LDA
238     3. QDA
239     4. Nearest Shrunken Centroids Classifier
240     Args:
241     data_file (string): Name of the data file (csv) of samples a features
242     header (string): None or position of the header (pandas read_csv parameter)
243     random_state (int): Seed for the random split of sets (as needed for sklearn
244     ↪ train_test_split)
245     label_pos (int): Column of the labels in data_file
246     ↪ label_value (int): Value of the label to assign internal '1' value. We consider
247     ↪ this label as
248     ↪ the positive label in prediction validation. We assign malign or cancer status to
249     ↪ this label.
250     ↪ features_ini (int): First column of features in data_file
251     ↪ features_fin (int): Last column + 1 of features in data_file. If None, last
252     ↪ column of file
253     ↪ best_param_value_method (str): if "max_in_test" gives the value with the maximum
254     ↪ accuracy
255     ↪ in test data
256     ↪ with_dim_red (bool): If True, it performs a dimensionality reduction by PCA
257     ↪ retained_variance (float): If dimensionality reduction, variance to retain
258     """
259     X_train_scaled, X_test_scaled, y_train, y_test = \
260         create_datasets_from_file(data_file, header, random_state,
261                                  label_pos, label_value, features_ini,
262                                  ↪ features_fin=features_fin,
263                                  with_dim_red=with_dim_red,
264                                  ↪ retained_variance=retained_variance)
265     if VERBOSE: print(X_train_scaled.shape)
266
267     if VERBOSE: print("NBC")
268     # Naive Bayes accuracy
269     accuracy_NBC = prediction_accuracy(X_train_scaled, X_test_scaled, y_train, y_test,
270     ↪ "GaussianNB", "", "")
271
272     # LDA accuracy
273     if VERBOSE: print("LDA")
274     accuracy_LDA = prediction_accuracy(X_train_scaled, X_test_scaled, y_train, y_test,
275     ↪ "LinearDiscriminantAnalysis", "", "")
276
277     # QDA estimate reg parameter

```



```

270 if VERBOSE: print("QDA")
271 param_values = np.linspace(0, 1, 10).tolist()
272 best_param_value = estimate_parameter(X_train_scaled, X_test_scaled, y_train, y_test, \
273                                     "QuadraticDiscriminantAnalysis", "reg_param", param_values, \
274                                     best_param_value_method)
275 # QDA accuracy
276 # Best parameter reg value according CV estimate
277 accuracy_QDA = prediction_accuracy(X_train_scaled, X_test_scaled, y_train, y_test, \
278                                   "QuadraticDiscriminantAnalysis", "reg_param", best_param_value)
279
280 # Centroids
281 if VERBOSE: print("NSC")
282 # Best parameter shrink_threshold value according CV estimate
283 param_values = np.linspace(0, 8, 20).tolist()
284 best_param_value = estimate_parameter(X_train_scaled, X_test_scaled, y_train, y_test, \
285                                     "NearestCentroid", "shrink_threshold", param_values, \
286                                     best_param_value_method)
287 # Centroids accuracy
288 accuracy_NSC = prediction_accuracy(X_train_scaled, X_test_scaled, y_train, y_test, \
289                                   "NearestCentroid", "shrink_threshold",
290                                   → best_param_value)
291 print_accuracies(accuracy_NBC, accuracy_LDA, accuracy_QDA, accuracy_NSC)

```

1.2 Breast cancer

Script 1.2.1 (python)

```

1 # Breast Cancer
2 DISABLE_PLOTS = True
3 VERBOSE = False
4 DIM_RED = False
5 #learn_dataset(data_file, None, 1, 1, "M", 2, features_fin = None, with_dim_red = False)
6 learn_dataset(data_file = './data/wdbc.csv', header = None, random_state=1,
7               label_pos=1, label_value="M", features_ini = 2, features_fin = None,
8               best_param_value_method="max_in_test",
9               with_dim_red = DIM_RED, retained_variance = 99.0)
10
11 learn_dataset(data_file = './data/wdbc.csv', header = None, random_state=1,
12               label_pos=1, label_value="M", features_ini = 2, features_fin = None,
13               best_param_value_method="max_in_cv",
14               with_dim_red = DIM_RED, retained_variance = 99.0)

```

Output

Best param value QuadraticDiscriminantAnalysis Method max_in_test: 0.5555555555555556
 Best param value NearestCentroid Method max_in_test: 3.3684210526315788

Accuracies

NBC LDA QDA NSC

Train	0.936675	0.973615	0.970976	0.931398
Test	0.936842	0.963158	0.978947	0.952632

Output

Best param value QuadraticDiscriminantAnalysis Method max_in_cv: 0.5555555555555556
 Best param value NearestCentroid Method max_in_cv: 7.578947368421052

Accuracies

	NBC	LDA	QDA	NSC
Train	0.936675	0.973615	0.970976	0.944591
Test	0.936842	0.963158	0.978947	0.926316

Output

1.3 Prostate cancer

Script 1.3.1 (python)

```

1 # Prostate Cancer
2 #learn_dataset(data_file, 0, 1, -1, 1, 0, -1, with_dim_red = True)
3 DISABLE_PLOTS = True
4 VERBOSE = False
5 DIM_RED = True
6 learn_dataset(data_file = './data/prostate.csv', header = 0, random_state = 1,
7               label_pos = -1, label_value = 1, features_ini = 0, features_fin = -1,
8               best_param_value_method = "max_in_test",
9               with_dim_red = DIM_RED, retained_variance = 99.0)
10
11 learn_dataset(data_file = './data/prostate.csv', header = 0, random_state = 1,
12               label_pos = -1, label_value = 1, features_ini = 0, features_fin = -1,
13               best_param_value_method = "max_in_cv",
14               with_dim_red = DIM_RED, retained_variance = 99.0)

```

Output

Features reduced to 63
 Best param value QuadraticDiscriminantAnalysis Method max_in_test: 0.7777777777777777
 Best param value NearestCentroid Method max_in_test: 0.42105263157894735

Accuracies

	NBC	LDA	QDA	NSC
Train	0.941176	1.000000	0.838235	0.882353
Test	0.823529	0.911765	0.852941	0.852941

Output

```

Features reduced to 63
Best param value QuadraticDiscriminantAnalysis Method max_in_cv: 0.0
Best param value NearestCentroid Method max_in_cv: 0.8421052631578947

Accuracies

```

	NBC	LDA	QDA	NSC
Train	0.941176	1.000000	1.000000	0.852941
Test	0.823529	0.911765	0.764706	0.794118

Output

1.4 Conclusions

We observe that **QDA** performs very poorly in the prostate dataset, given the high dimensionality of this dataset, which do not ease the accurate computation of the covariance matrices. Perhaps if we perform previously a dimensionality reduction by PCA, we'll improve this result.

NSC performs in this case much better due to the reduced number of parameters and the feature selection properties of this classifier and more consistently between both cases (prostate and breast).