

Classification Systems

Daniel Cerdán, Fernando Freire

March 22, 2019

Contents

1	Classification Systems	2
1.1	Methods	3
1.1.1	Implementation details	3
1.2	Breast cancer	7
1.3	Prostate cancer	9
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
```

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

```

25         break
26     return component_number, cumulative_variance, explained_variance
27
28
29 def create_datasets_from_file(data_file, header, random_state, label_pos,
30                               label_value, features_ini, features_fin=None,
31                               with_dim_red=False, retained_variance=99.0):
32     """Create training and test sets from file
33
34     Args:
35         data_file (string): Name of the data file (csv) of samples a features
36         header (string): None or position of the header (pandas read_csv parameter)
37         random_state (int): Seed for the random split (as needed for sklearn
38     ↪ train_test_split)
39         label_pos (int): Column of the labels in data_file
40         label_value (int): Value of the label to assign internal '1' value
41         features_ini (int): First column of features in data_file
42     ↪ features_fin (int): Last column + 1 of features in data_file. If None, last
43     ↪ 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,
64                                                         ↪ random_state=random_state)
65
66     # Data standardization
67     scaler = preprocessing.StandardScaler().fit(x_train)
68     x_train_scaled = scaler.transform(x_train)
69     x_test_scaled = scaler.transform(x_test)
70
71     # Check standardization
72     for i in range (1, np.size(x_train_scaled,1)):
73         assert round(np.var(x_train_scaled[:,0]),3) == round(np.var(x_train_scaled[:,i]),3),\
74                 "Warning: revise data standardization")
75
76     if with_dim_red:
77         desired_variance = retained_variance
78         component_number, _, _ =\

```

```

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

```

```

125         method_func (string) : name of the learning method
126         param (string): name of learning method parameter
127         param_values (list of float): list of parameter values to try
128     Returns:
129         (float): best parameter value to use in prediction
130
131     """
132     # Pipeline for estimate the regularization parameter
133     pipeline = Pipeline([ ('method', globals()[method_func]() )])
134
135     # Construct the grid the hyperparameter candidate shronk theshold
136     param_grid = { 'method__' + param : param_values }
137
138     # Evaluating
139     skfold = RepeatedStratifiedKfold(n_splits=10, n_repeats=1, random_state=0)
140     gridcv = GridSearchCV(pipeline, cv=skfold, n_jobs=1, param_grid=param_grid,\
141         scoring=make_scorer(accuracy_score))
142     result = gridcv.fit(x_train, y_train)
143
144     # Accuracies
145     accuracies = gridcv.cv_results_['mean_test_score']
146     std_accurrencies = gridcv.cv_results_['std_test_score']
147
148     test_accurrencies = np.ones(len(param_values))
149
150     for i in range(len(param_values)):
151         method_params = {param : param_values[ i ]}
152         method = globals()[method_func](**method_params)
153         method.fit(x_train, y_train)
154         test_accurrencies[ i ] = accuracy_score(method.predict(x_test), y_test)
155
156     max_test_accuracy = max(test_accurrencies)
157
158     # Obtain best_param_value as max
159     best_param_value = 0
160     best_train_accuracy = 0
161     for i in range(len(param_values)):
162         if test_accurrencies[ i ] == max_test_accuracy:
163             if accuracies[i] > best_train_accuracy:
164                 best_train_accuracy = accuracies[i]
165                 best_param_value = param_values[i]
166
167     #best_param_value = param_values[ np.argmax(accuracies) ]
168     # Plot
169     plt.figure(figsize=(15, 10))
170     plt.grid()
171     line1, = plt.plot(param_values, accuracies, 'o-', color="g")
172     line2, = plt.plot(param_values, test_accurrencies, 'x-', color="r")
173     plt.fill_between(param_values, accuracies - std_accurrencies / np.sqrt(10), \
174         accuracies + std_accurrencies / np.sqrt(10), alpha=0.1, color="g")
175     plt.title("Different hyper-parameter " + param + " values for " + method_func)
176     plt.xlabel('Hyper-parameter')

```

```

177 plt.xticks(np.round(np.array(param_values), 2))
178 plt.ylabel('Classification Accuracy')
179 plt.ylim((min(min(accuracies), min(test_accuracies)) - 0.1,
180           min(1.02, max(max(accuracies), max(test_accuracies)) + 0.1)))
181
182 plt.xlim((min(param_values), max(param_values)))
183 legend_handles = [ mlines.Line2D([], [], color='g', marker='o', \
184                           markersize=15, label='CV-estimate'), \
185                   mlines.Line2D([], [], color='r', marker='x', \
186                           markersize=15, label='Test set estimate')]
187 plt.legend(handles=legend_handles, loc = 3)
188 plt.show()
189
190 print("Best param value:", best_param_value)
191 return best_param_value
192
193 def learn_dataset(data_file, header, random_state, label_pos,
194                  label_value, features_ini, features_fin=None,
195                  with_dim_red=False, retained_variance=99.0):
196     """Learn data sets from file, methods:
197         1. The Naive Bayes Classifier
198         2. LDA
199         3. QDA
200         4. Nearest Shrunken Centroids Classifier
201     Args:
202         data_file (string): Name of the data file (csv) of samples a features
203         header (string): None or position of the header (pandas read_csv parameter)
204         random_state (int): Seed for the random split of sets (as needed for sklearn
205     ↪ train_test_split)
206         label_pos (int): Column of the labels in data_file
207         label_value (int): Value of the label to assign internal '1' value
208         features_ini (int): First column of features in data_file
209         features_fin (int): Last column + 1 of features in data_file. If None, last
210     ↪ column of file
211         with_dim_red (bool): If True, it performs a dimensionality reduction by PCA
212         retained_variance (float): If dimensionality reduction, variance to retain
213     """
214     X_train_scaled, X_test_scaled, y_train, y_test = \
215         create_datasets_from_file(data_file, header, random_state,
216                                  label_pos, label_value, features_ini,
217                                  ↪ features_fin=features_fin,
218                                  with_dim_red=with_dim_red,
219                                  ↪ retained_variance=retained_variance)
220     print(X_train_scaled.shape)
221
222     # Naive Bayes accuracy
223     prediction_accuracy(X_train_scaled, X_test_scaled, y_train, y_test, "GaussianNB", "", "")
224
225     # LDA accuracy
226     prediction_accuracy(X_train_scaled, X_test_scaled, y_train, y_test,
227                         ↪ "LinearDiscriminantAnalysis", "", "")

```

```

224
225     # QDA estimate reg parameter
226     param_values = np.linspace(0, 1, 10).tolist()
227     best_param_value = estimate_parameter(X_train_scaled, X_test_scaled, y_train, y_test, \
228                                         "QuadraticDiscriminantAnalysis", "reg_param", param_values)
229
230     # QDA accuracy
231     # Best parameter reg value according CV estimate
232     prediction_accuracy(X_train_scaled, X_test_scaled, y_train, y_test, \
233                         "QuadraticDiscriminantAnalysis", "reg_param", best_param_value)
234
235     # Centroids
236     # Best parameter shrink_threshold value according CV estimate
237     param_values = np.linspace(0, 8, 20).tolist()
238     best_param_value = estimate_parameter(X_train_scaled, X_test_scaled, y_train, y_test, \
239                                         "NearestCentroid", "shrink_threshold", param_values)
240
241     # Centroids accuracy
242     prediction_accuracy(X_train_scaled, X_test_scaled, y_train, y_test, "NearestCentroid",
243                         ↪ "shrink_threshold", best_param_value)

```

1.2 Breast cancer

Script 1.2.1 (python)

```

1  # Breast Cancer
2  data_file = './data/wdbc.csv'
3  learn_dataset(data_file, None, 1, 1, "B", 2, features_fin = None, with_dim_red = False)

```

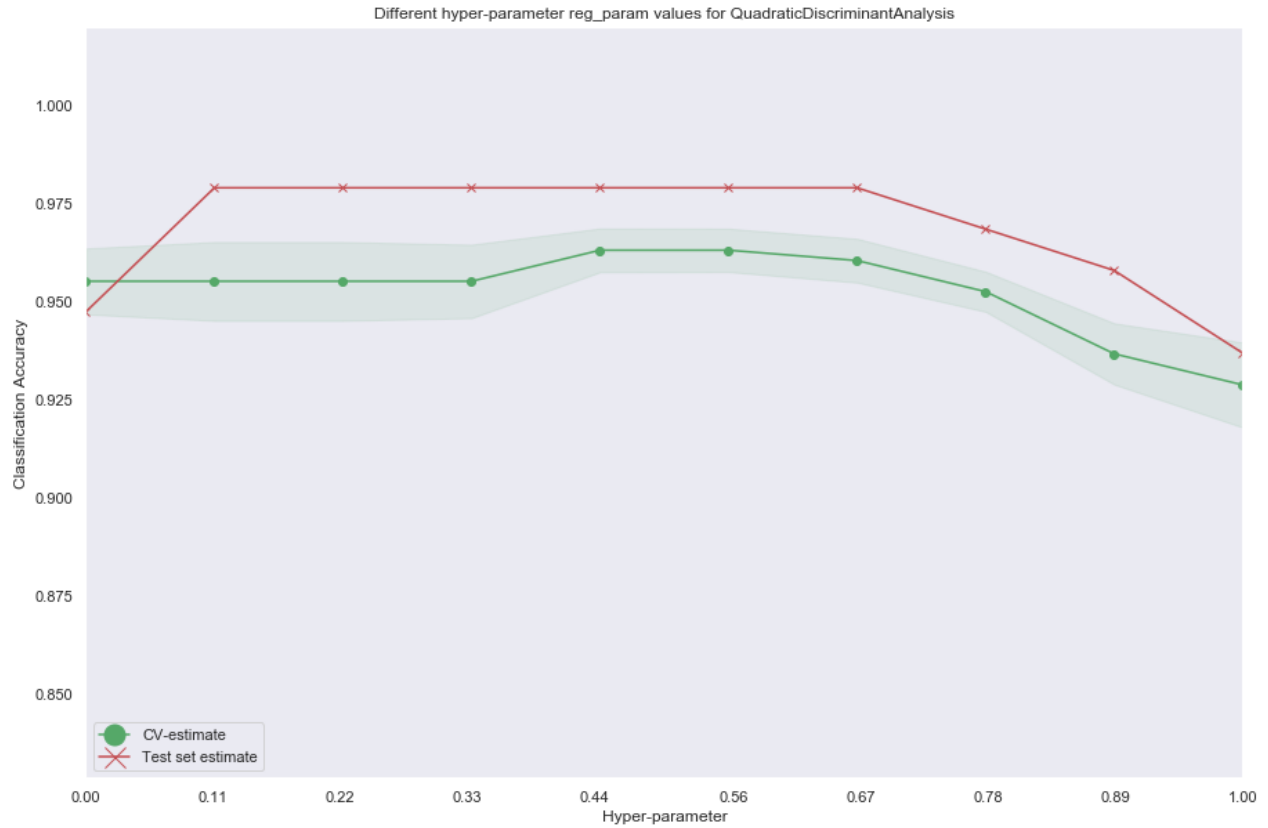
Output

```

(379, 30)
[[ 61   5]
 [  7 117]]
Prediction accuracy is: 0.936842
True postive rate is: 0.943548
True negative rate is: 0.924242

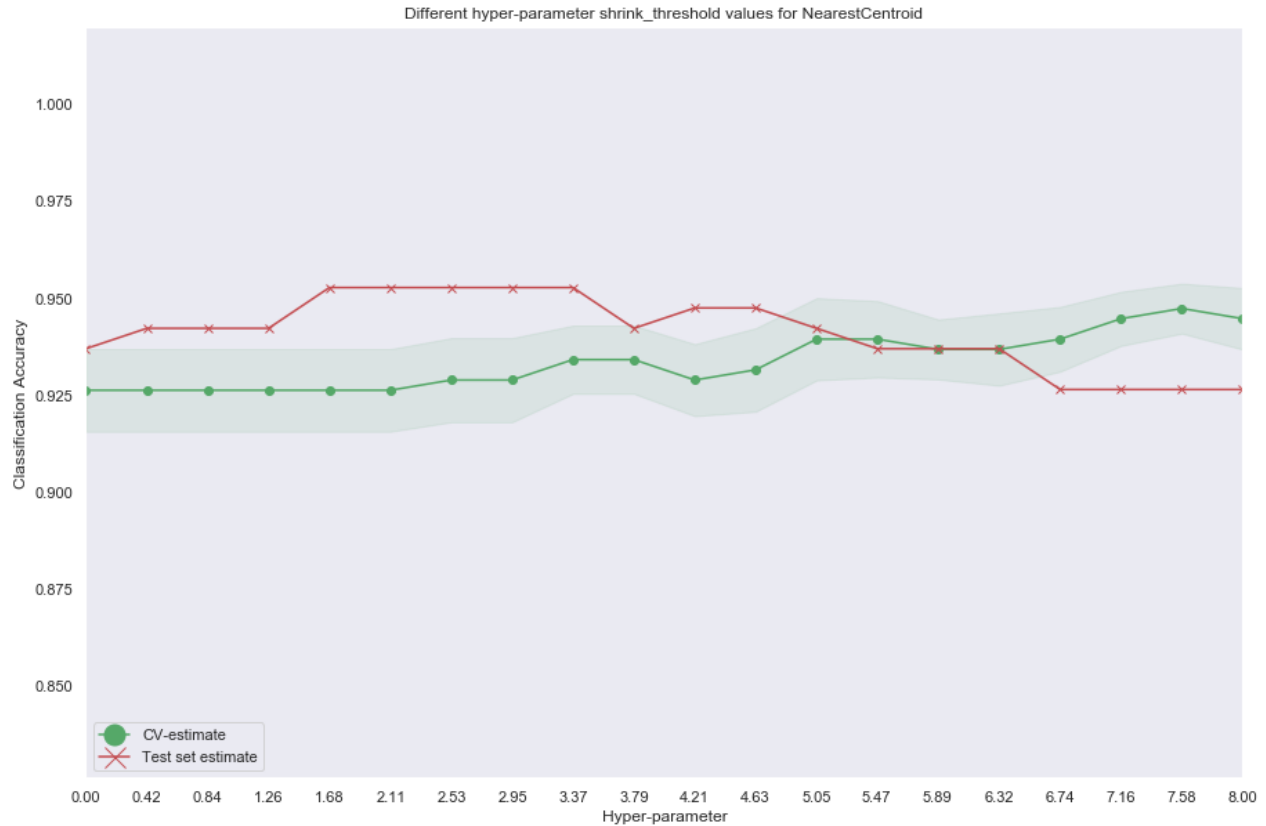
[[ 60   6]
 [  1 123]]
Prediction accuracy is: 0.963158
True postive rate is: 0.991935
True negative rate is: 0.909091

```

Output

```
Best param value: 0.4444444444444444
[[ 62  4]
 [  0 124]]
Prediction accuracy is: 0.978947
True postive rate is: 1.000000
True negative rate is: 0.939394
```



Output

```
Best param value: 3.3684210526315788
[[ 60  6]
 [ 3 121]]
Prediction accuracy is: 0.952632
True positive rate is: 0.975806
True negative rate is: 0.909091
```

1.3 Prostate cancer

Script 1.3.1 (python)

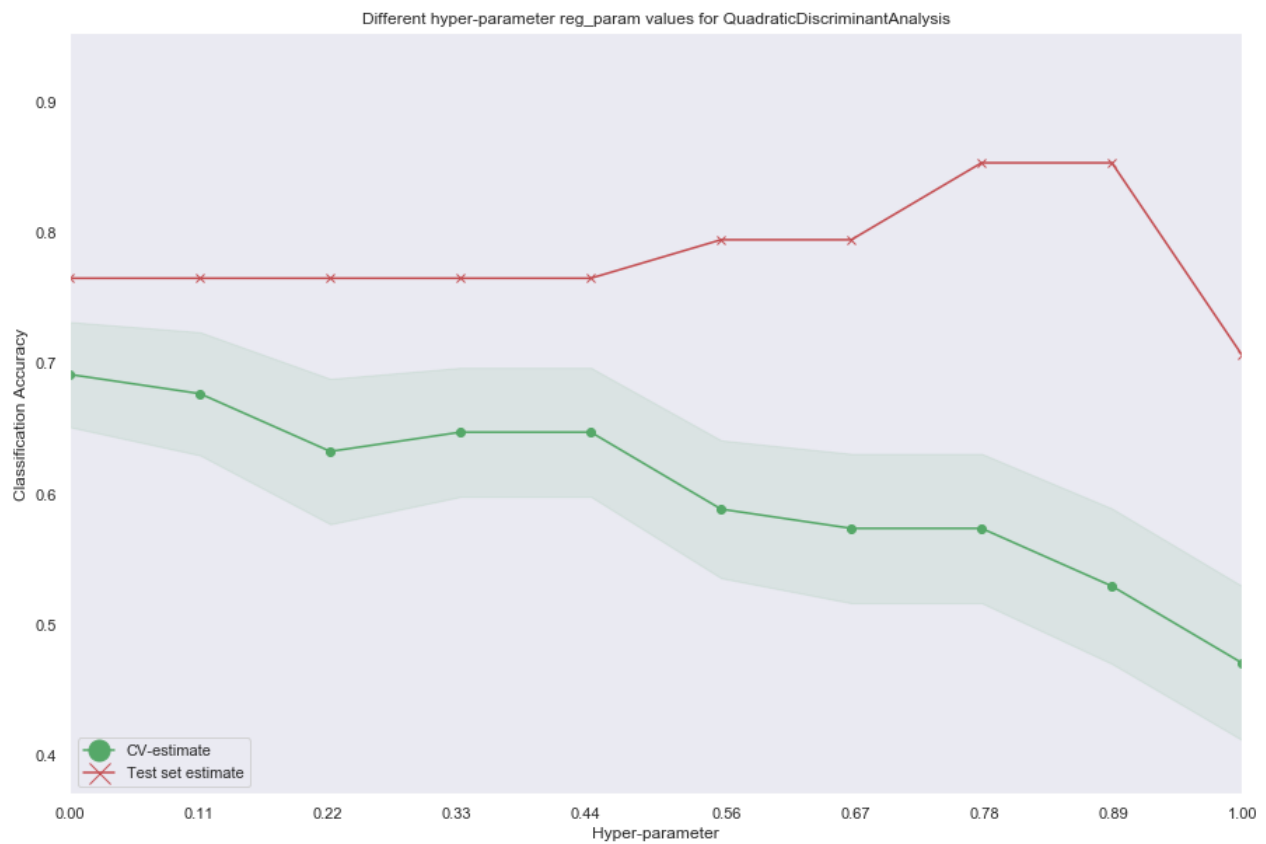
```
1 # Prostate Cancer
2 data_file = './data/prostate.csv'
3 learn_dataset(data_file, 0, 1, -1, 1, 0, -1, with_dim_red = True)
```

Output

```
Features reduced to 63
(68, 63)
```

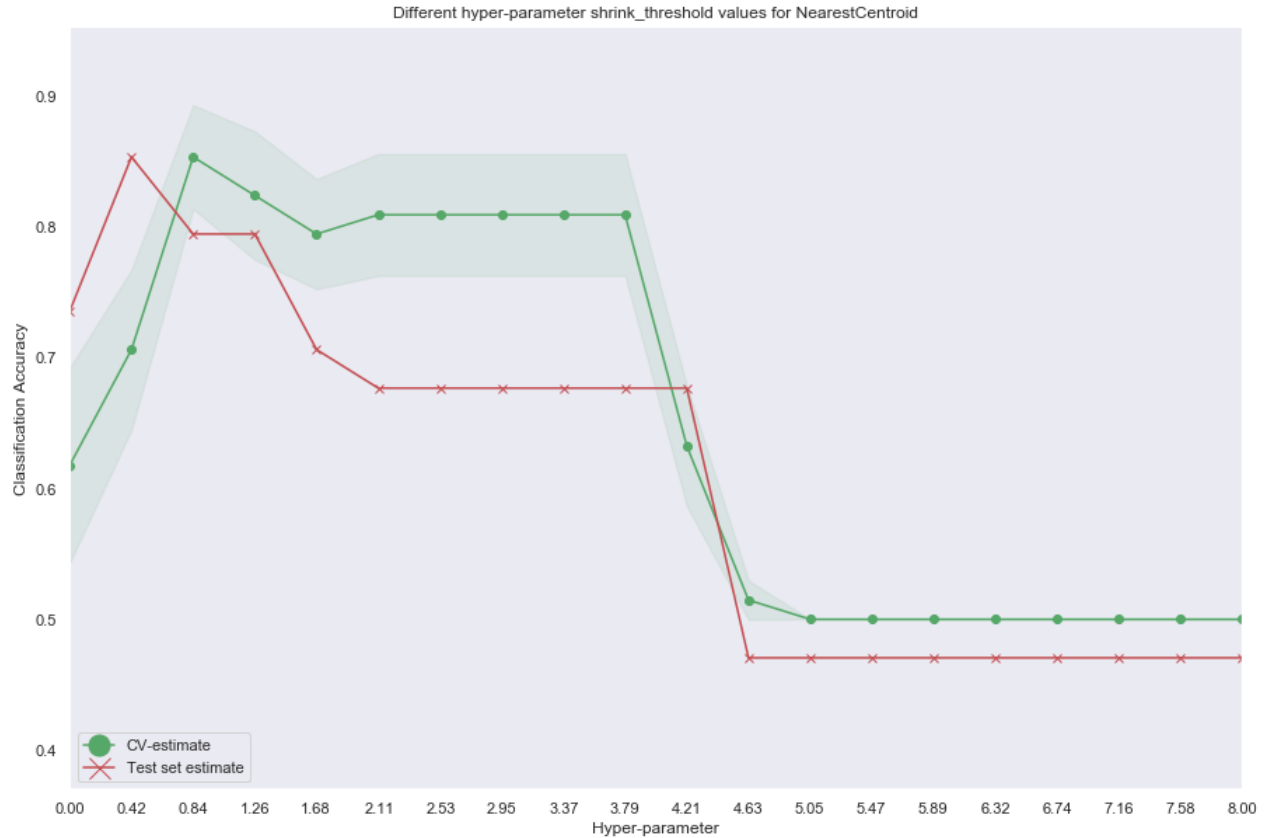
```
[[11  5]
 [ 1 17]]
Prediction accuracy is: 0.823529
True postive rate is: 0.944444
True negative rate is: 0.687500
```

```
[[16  0]
 [ 3 15]]
Prediction accuracy is: 0.911765
True postive rate is: 0.833333
True negative rate is: 1.000000
```



Output

```
Best param value: 0.7777777777777777
[[12  4]
 [ 1 17]]
Prediction accuracy is: 0.852941
True postive rate is: 0.944444
True negative rate is: 0.750000
```



Output

```
Best param value: 0.42105263157894735
[[14  2]
 [ 3 15]]
Prediction accuracy is: 0.852941
True postive rate is: 0.833333
True negative rate is: 0.875000
```

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).