

Classification Systems

Daniel Cerdán, Fernando Freire

March 21, 2019

Contents

1	Classification Systems	2
1.1	Methods	3
1.2	Breast cancer	5
1.3	Prostate cancer	7

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 Shrunken 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.pipeline import Pipeline
20 from sklearn.model_selection import train_test_split, RepeatedStratifiedKFold, GridSearchCV
21 from sklearn import preprocessing
22 from sklearn.metrics import accuracy_score, make_scorer, confusion_matrix
```

1.1 Methods

Script 1.1.1 (python)

```
1 def create_datasets_from_file(data_file, header, label_pos, label_value, features_ini,
  ↳ features_fin=None):
2     """
3     """
4     data = pd.read_csv(data_file, header = header)
5     if features_fin == None:
6         X = data.values[ :, features_ini:].astype(np.float)
7     else:
8         X = data.values[ :, features_ini:features_fin].astype(np.float)
9     y = (data.values[ :, label_pos ] == label_value).astype(np.int)
10    # Split dataset between training and test
11    x_train, x_test, y_train, y_test = train_test_split(X, y, test_size=1.0/3, random_state=
12    1)
13    # Data standardization
14    scaler = preprocessing.StandardScaler().fit(x_train)
15    x_train_scaled = scaler.transform(x_train)
16    x_test_scaled = scaler.transform(x_test)
17
18    # Check standardization
19    for i in range (1, np.size(x_train_scaled,1)):
20        assert round(np.var(x_train_scaled[:,0]),3) == round(np.var(x_train_scaled[:,i]),3),\
21            "Warning: revise data standardization"
22
23    return x_train_scaled, x_test_scaled, y_train, y_test
24
25 def prediction_accuracy(x_train, x_test, y_train, y_test, method_func, method_param,
  ↳ param_value):
26     """
27     """
28     if method_param != "" :
29         params = {method_param : param_value}
30     else:
31         params = {}
32     method = globals()[method_func](**params)
33
34     # Training
35     method.fit(x_train, y_train)
36
37     # Prediction
38     y_pred = method.predict(x_test)
39     conf = confusion_matrix(y_test, y_pred)
40     TN = conf[0][0]
41     TP = conf[1][1]
42     FP = conf[0][1]
43     FN = conf[1][0]
44     print(conf)
45     print('Predicion accuracy is: %f' % ((TP + TN) / (TN + TP + FP + FN)))
46     print('True postive rate is: %f' % (TP / (TP + FN)))
47     print('True negative rate is: %f\n' % (TN / (TN + FP)))
```

```

47
48 def estimate_parameter(x_train, x_test, y_train, y_test, method_func, param, param_values):
49     """
50     """
51     # Pipeline for estimate the regularization parameter
52     pipeline = Pipeline([ ('method', globals()[method_func]()) ])
53
54     # Construct the grid the hyperparameter candidate shronk theshold
55     param_grid = { 'method__' + param : param_values }
56
57     # Evaluating
58     skfold = RepeatedStratifiedKfold(n_splits=10, n_repeats=1, random_state=0)
59     gridcv = GridSearchCV(pipeline, cv=skfold, n_jobs=1, param_grid=param_grid, \
60         scoring=make_scorer(accuracy_score))
61     result = gridcv.fit(x_train, y_train)
62
63     # Accuracies
64     accuracies = gridcv.cv_results_['mean_test_score']
65     std_accuracies = gridcv.cv_results_['std_test_score']
66
67     test_accuracies = np.ones(len(param_values))
68
69     for i in range(len(param_values)):
70         method_params = {param : param_values[ i ]}
71         method = globals()[method_func](**method_params)
72         method.fit(x_train, y_train)
73         test_accuracies[ i ] = accuracy_score(method.predict(x_test), y_test)
74
75     # Plot
76     plt.figure(figsize=(15, 10))
77     line1, = plt.plot(param_values, accuracies, 'o-', color="g")
78     line2, = plt.plot(param_values, test_accuracies, 'x-', color="r")
79     plt.fill_between(param_values, accuracies - std_accuracies / np.sqrt(10), \
80         accuracies + std_accuracies / np.sqrt(10), alpha=0.1, color="g")
81     plt.grid()
82     plt.title("Different regularization parameter values for " + method_func)
83     plt.xlabel('Regularization Parameter')
84     plt.xticks(np.round(np.array(param_values), 2))
85     plt.ylabel('Classification Accuracy')
86     plt.ylim((min(accuracies) - 0.1, max(accuracies) + 0.1))
87
88     plt.xlim((min(param_values), max(param_values)))
89     legend_handles = [ mlines.Line2D([], [], color='g', marker='o', \
90         markersize=15, label='CV-estimate'), \
91         mlines.Line2D([], [], color='r', marker='x', \
92         markersize=15, label='Test set estimate')]
93     plt.legend(handles=legend_handles, loc = 3)
94     plt.show()
95
96 def learn_dataset(data_file, header, label_pos, label_value, features_ini, features_fin=None
97 ):
98     """

```

```

98      """
99
100     X_train_scaled, X_test_scaled, y_train, y_test = \
101         create_datasets_from_file(data_file, header, label_pos, label_value, features_ini,
102             ↪ features_fin=None)
103     print(X_train_scaled.shape)
104
105     # Naive Bayes accuracy
106     prediction_accuracy(X_train_scaled, X_test_scaled, y_train, y_test, "GaussianNB", "", "")
107
108     # LDA accuracy
109     prediction_accuracy(X_train_scaled, X_test_scaled, y_train, y_test,
110         ↪ "LinearDiscriminantAnalysis", "", "")
111
112     # QDA estimate reg parameter
113     param_values = np.linspace(0, 1, 10).tolist()
114     estimate_parameter(X_train_scaled, X_test_scaled, y_train, y_test, \
115         ↪ "QuadraticDiscriminantAnalysis", "reg_param", param_values)
116
117     # QDA accuracy
118     # Best parameter reg value according CV estimate
119     best_param_value = 0.444 #to compute automagically
120     prediction_accuracy(X_train_scaled, X_test_scaled, y_train, y_test, \
121         ↪ "QuadraticDiscriminantAnalysis", "reg_param", best_param_value)
122
123     # Centroids
124     # Best parameter shrink_threshold value according CV estimate
125     param_values = np.linspace(0, 8, 20).tolist()
126     estimate_parameter(X_train_scaled, X_test_scaled, y_train, y_test, \
127         ↪ "NearestCentroid", "shrink_threshold", param_values)
128
129     best_param_value = 3.4 #to compute automagically
130     # Centroids accuracy
131     prediction_accuracy(X_train_scaled, X_test_scaled, y_train, y_test, "NearestCentroid",
132         ↪ "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, "B", 2)

```

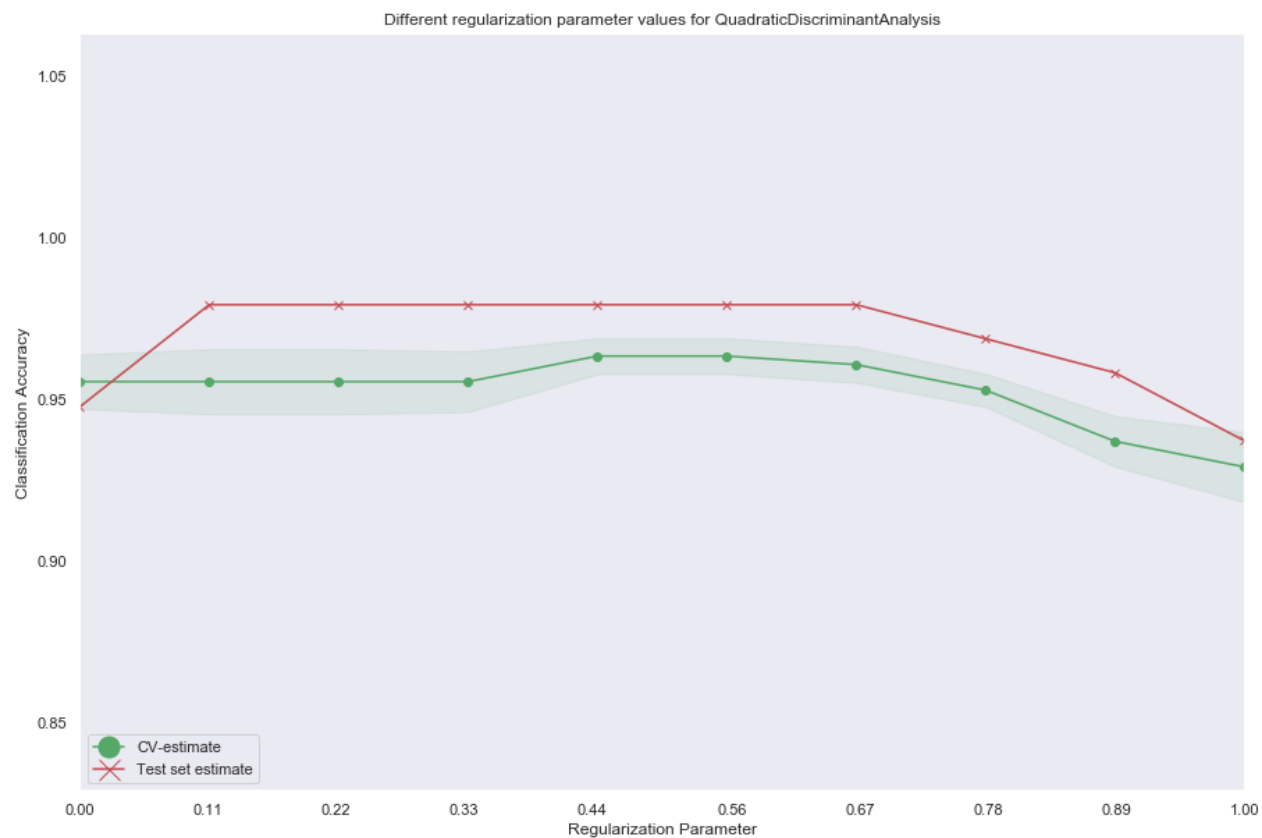
Output

```

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

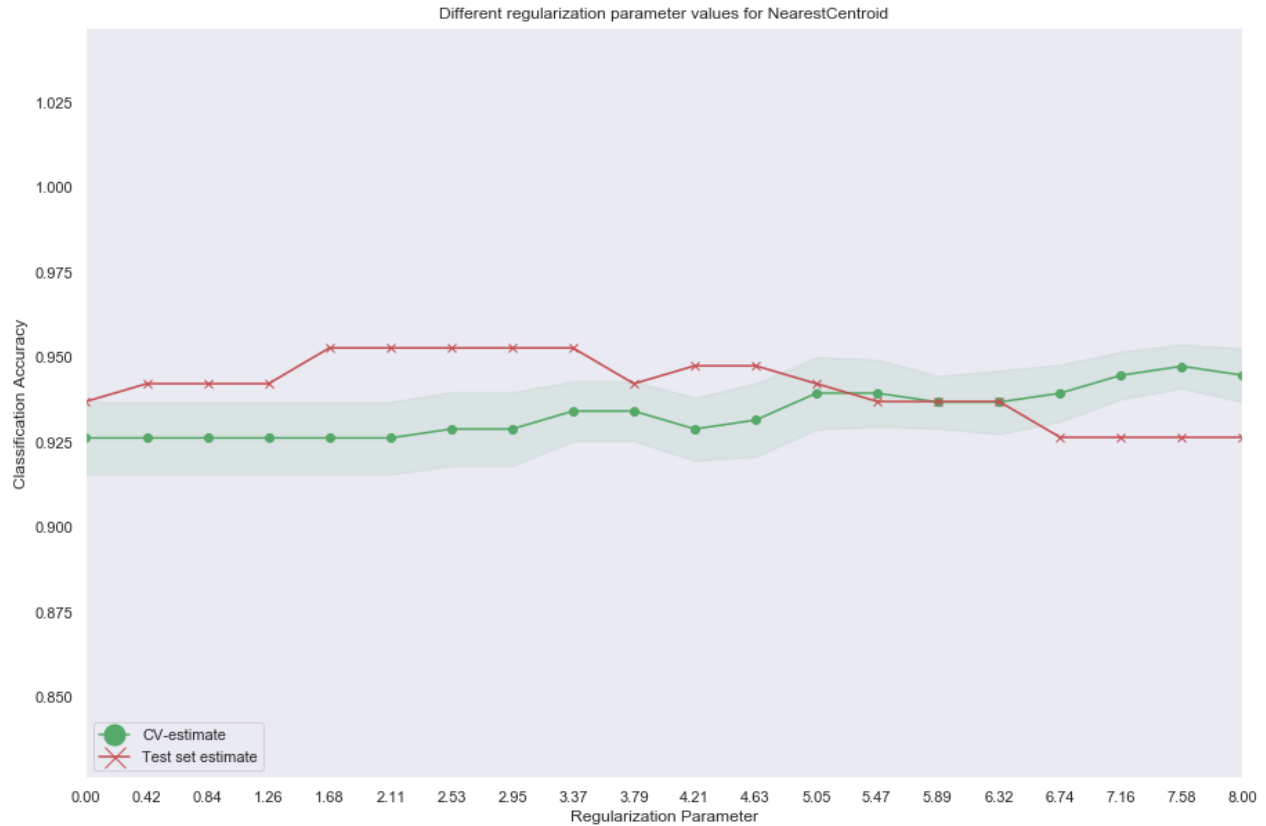
```

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



Output

```
[[ 62  4]
 [  0 124]]
Prediction accuracy is: 0.978947
True positive rate is: 1.000000
True negative rate is: 0.939394
```



Output

```
[[ 60  6]
 [ 3 121]]
Prediction accuracy is: 0.952632
True postive rate is: 0.975806
True negative rate is: 0.909091
```

1.3 Prostate cancer

Script 1.3.1 (python)

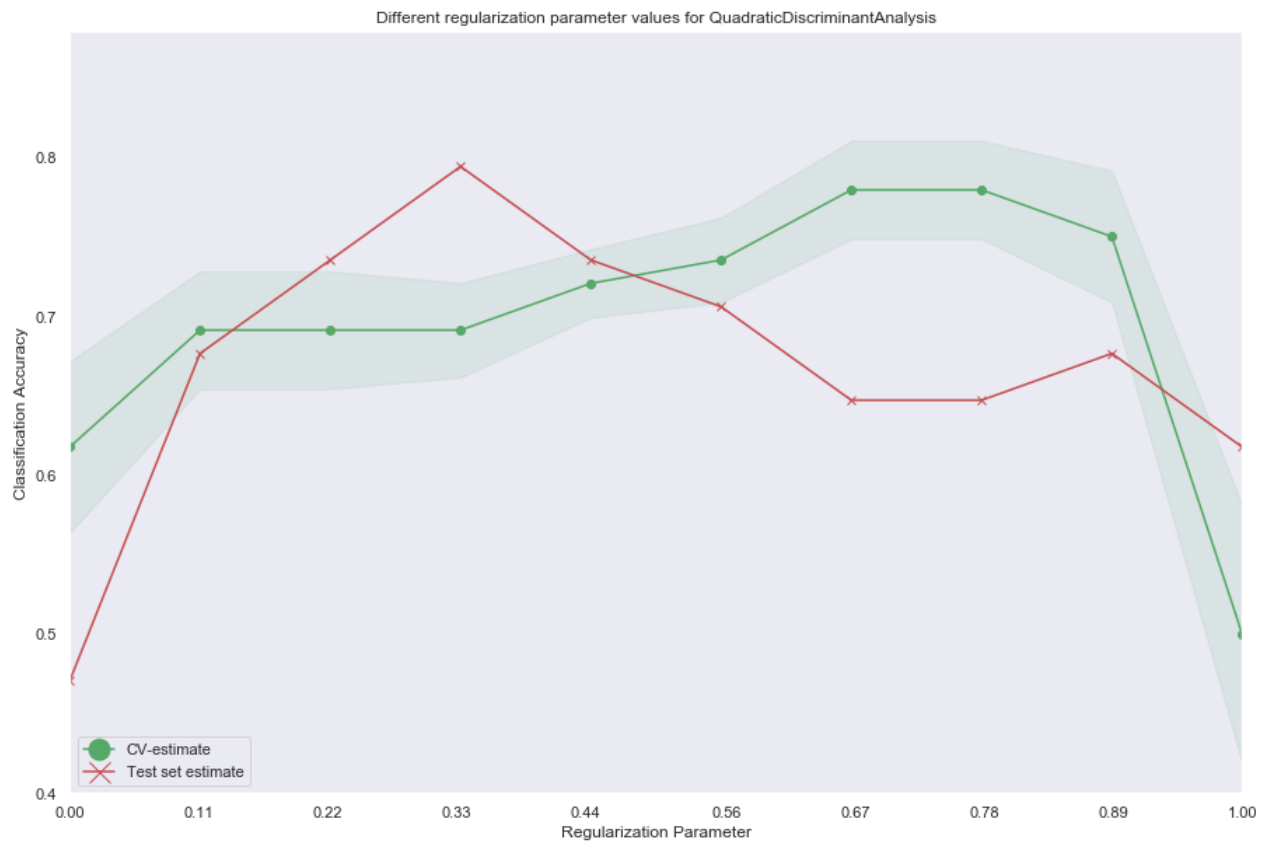
```
1 # Prostate Cancer
2 data_file = './data/prostate.csv'
3 learn_dataset(data_file, 0, -1, 1, 0, -1)
```

Output

```
(68, 12626)
[[16  0]
 [ 0 18]]
```

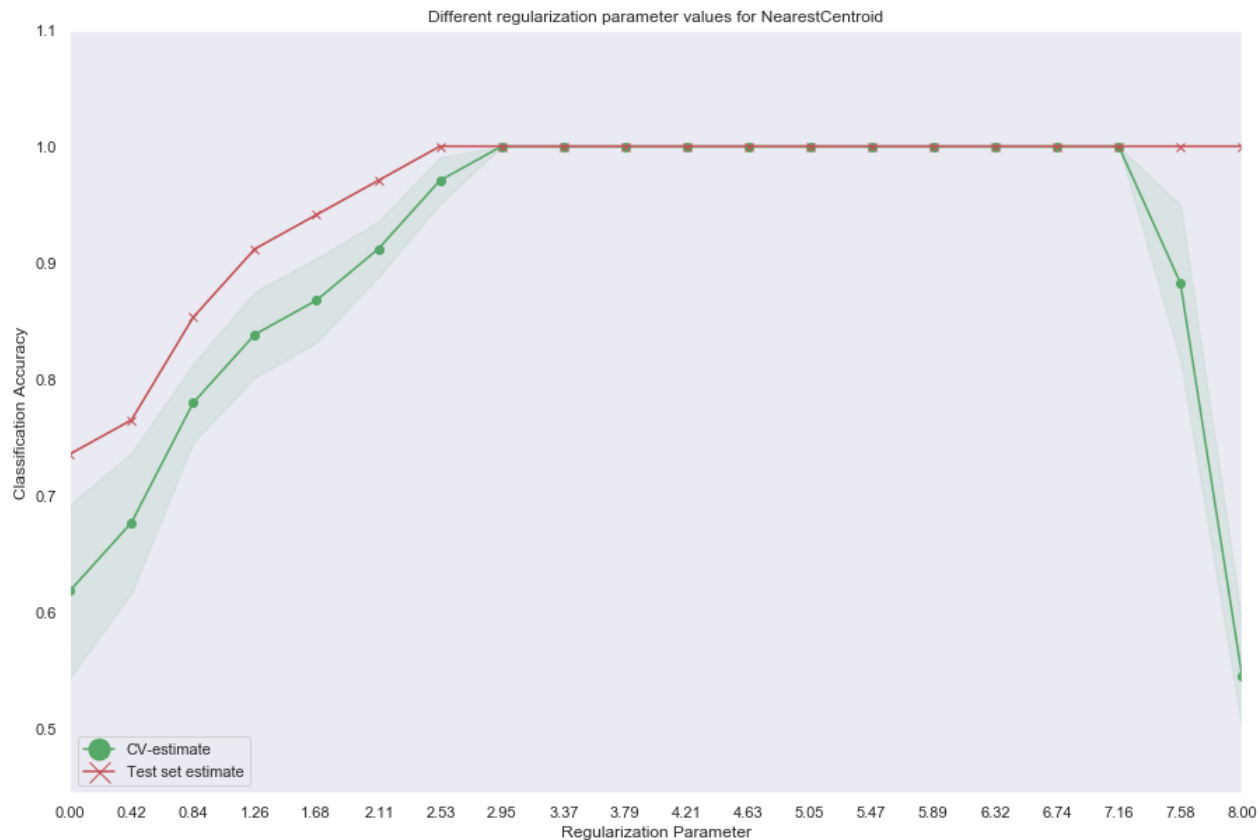
```
Prediction accuracy is: 1.000000  
True positive rate is: 1.000000  
True negative rate is: 1.000000
```

```
[[15  1]  
 [ 4 14]]  
Prediction accuracy is: 0.852941  
True positive rate is: 0.777778  
True negative rate is: 0.937500
```



Output

```
[[13  3]  
 [ 6 12]]  
Prediction accuracy is: 0.735294  
True positive rate is: 0.666667  
True negative rate is: 0.812500
```

Output

```
[[16  0]  
 [ 0 18]]
```

Prediction accuracy is: 1.000000

True positive rate is: 1.000000

True negative rate is: 1.000000