

# ECEN 649 Pattern Recognition - Spring 2017

## Computer Project 2

Geethik Narayana Kamineni  
UIN - 825008990

---

### Introduction

The project involved determining the best gene feature sets using techniques of error estimation, feature selection and classifier design which best discriminate the two prognosis classes in the breast tumor biopsies data from the following cancer study:

van de Vijver, M.J., He, Y.D., van't Veer, L.J., et al. (2002), "A gene-expression signature as a predictor of survival in breast cancer." *New Eng. J. Med.*, 347, 1999-2009.

### Procedure

The process to obtain the best gene features consisted of implementing the two classifier techniques *DLDA* and *3NN* along with the feature selection methods of *Exhaustive Search* and *Sequential Forward Search* on the data considering the determining error estimate to be the resubstitution error.

### DLDA

Diagonal linear discriminant analysis (*DLDA*) is a special case of the *LDA* wherein the sample covariance matrix is constrained to be only along the diagonal. This is implemented in *Matlab* using the *fitcdiscr* function using the following convention:

$$mdl = \text{fitcdiscr}(\text{trainX}, \text{trainY}', \text{DiscrimType}', \text{diaglinear}')$$

### 3NN

$k$  – Nearest Neighbor rule (*KNN*) where in  $k = 3$  is implemented as one another strategy in finding the best gene sets. Its implementation in *Matlab* is through the function *fitcknn* as follows:

$$mdl = \text{fitcknn}(\text{trainX}, \text{trainY}', \text{NumNeighbors}', 3)$$

## Error criterion

The resubstitution error which is the apparent error of the designed classifier on the training data is given by:

$$\varepsilon_n^r = \frac{1}{n} \sum_{i=0}^n |Y_i - \psi(X_i)|$$

is implemented in *Matlab* on the predictions made by the *predict* function as follows:

$$\begin{aligned} label_{pred} &= predict mdl, trainX \\ \varepsilon_n^r &= sum(abs(label_{pred}-trainY))/size(trainY,1) \end{aligned}$$

The test set error estimate on the test data  $(x_i^{(m)}, y_i^{(m)})$  which is given by:

$$\varepsilon_{n,m} = \frac{1}{m} \sum_{i=0}^m |Y_i^{(m)} - \psi(X_i^{(m)})|$$

is even implemented in a similar fashion to the resubstitution error but instead on the test data.

## Feature Selection

For the case of determining the feature sets with size from 1 to 3 using *Exhaustive Search*, all the combinations are being generated using the *combnk* combinatorial function in *Matlab* and looped over each time to find the best discriminating gene set for each classifier.

In the case of *Sequential Forward Search* to determine the top eight features among the 70-gene sets *sequentialfs* function in *Matlab* is being used whose implementation is as follows:

```
[fs, history] = sequentialfs(criterion, trainX, trainY, 'cv', 'resubstitution', 'nfeatures', 8);
```

where the *fs* stores the logical vector of the gene sets obtained after the final iteration whereas *history* holds the logical vectors in each iteration and also the criterion values of each iteration. The criterion function is being based on the *DLDA* and *3NN* classifiers. The code for the analysis implementation is as follows:

```
1 % Computer Project 2 Code
2 %Get data
3
```

```

4 trainObj = importdata("train.txt",'t');
5 train_x = trainObj.data;
6 train_y = train_x(:,72);
7 train_x = train_x(:,1:71);
8 trainHeader = trainObj.colheaders;
9
10 testObj = importdata("test.txt",'t');
11 test_x = testObj.data;
12 test_y = test_x(:,72);
13 test_x = test_x(:,1:71);
14 testHeader = testObj.colheaders;
15
16 %Exhaustive search
17 fprintf("*****Exhaustive for size-1*****\n");
18
19 %size-1 feature set
20 %3NN
21 for i=2:1:size(train_x,2)
22     new_train = [train_x(:,i)];
23     mdl_knn = fitcknn(new_train,train_y,'NumNeighbors',3);
24     new_train_label = mdl_knn.predict(new_train);
25
26     if i>2
27         if res_error > sum(abs(new_train_label - train_y))/size(train_x,1)
28             res_error = sum(abs(new_train_label - train_y))/size(train_x,1);
29             best_knn_1 = i;
30         end
31     else
32         res_error = sum(abs(new_train_label - train_y))/size(train_x,1);
33         best_knn_1 = 2;
34     end
35 end
36
37 best1 = char(trainHeader(best_knn_1));
38 fprintf('Best set of size-1 for 3NN is {%s} and the resubstitution error is %f\n',best1,res_error)
39
40 %DLDA
41 for i=2:1:size(train_x,2)
42     new_train = [train_x(:,i)];
43     mdl_dlda = fitcdiscr(new_train,train_y,'DiscrimType','diaglinear');
44     new_train_label = predict(mdl_dlda, new_train);
45
46     if i>2
47         if res_error_dlda > sum(abs(new_train_label - train_y))/size(train_x,1)
48             res_error_dlda = sum(abs(new_train_label - train_y))/size(train_x,1);
49             best_dlda_1 = i;
50         end
51     else
52         res_error_dlda = sum(abs(new_train_label - train_y))/size(train_x,1);
53         best_dlda_1 = 2;
54     end
55 end
56
57 best1_dlda = char(trainHeader(best_dlda_1));
58 fprintf('Best set of size-1 for DLDA is {%s} and the resubstitution error is %f\n',best1_dlda,res_error_dlda)
59
60 %Test Set error calculation
61 best_train_knn = [train_x(:,best_knn_1)];
62 best_mdl_knn = fitcknn(best_train_knn, train_y,'NumNeighbors',3);
63 best_test_knn = [test_x(:,best_knn_1)];
64 test_label_knn = best_mdl_knn.predict(best_test_knn);
65 test_error_knn = sum(abs(test_label_knn - test_y))/size(test_x,1);
66 fprintf('Test set error for 3NN on best set {%s} of size-1 is %f\n',best1,test_error_knn);
67
68 best_train_dlda = [train_x(:,best_dlda_1)];
69 best_mdl_dlda = fitcdiscr(best_train_dlda, train_y,'DiscrimType','diaglinear');
70 best_test_dlda = [test_x(:,best_dlda_1)];
71 test_label_dlda = predict(best_mdl_dlda, best_test_dlda);
72 test_error_dlda = sum(abs(test_label_dlda - test_y))/size(test_x,1);
73 fprintf('Test set error for DLDA on best set {%s} of size-1 is %f\n',best1_dlda,test_error_dlda);
74
75 fprintf("*****\n");
76 fprintf("*****Exhaustive for size-2*****\n");
77
78 %-----$
79 %size-2 feature set
80 C = combnk(2:71,2);
81
82 %3NN
83 for i=1:1:size(C,1)

```

```

84     new_train = [train_x(:,C(i,:))];
85     mdl_knn = fitcknn(new_train,train_y,'NumNeighbors',3);
86     new_train_label = mdl_knn.predict(new_train);
87
88     if i>1
89         if res_error > sum(abs(new_train_label - train_y))/size(train_x,1)
90             res_error = sum(abs(new_train_label - train_y))/size(train_x,1);
91             best_knn_2 = i;
92         end
93     else
94         res_error = sum(abs(new_train_label - train_y))/size(train_x,1);
95         best_knn_2 = 1;
96     end
97 end
98
99 best1 = char(trainHeader(C(best_knn_2,1)));
100 best2 = char(trainHeader(C(best_knn_2,2)));
101 fprintf('Best set of size-2 for 3NN is {%s, %s} and the resubstitution error is %f\n',best1,best2,res_error)
102
103 %DLDA
104 for i=1:1:size(C,1)
105     new_train = [train_x(:,C(i,:))];
106     mdl_dlda = fitcdiscr(new_train,train_y,'DiscrimType','diagLinear');
107     new_train_label = predict(mdl_dlda, new_train);
108
109     if i>1
110         if res_error_dlda > sum(abs(new_train_label - train_y))/size(train_x,1)
111             res_error_dlda = sum(abs(new_train_label - train_y))/size(train_x,1);
112             best_dlda_2 = i;
113         end
114     else
115         res_error_dlda = sum(abs(new_train_label - train_y))/size(train_x,1);
116         best_dlda_2 = 1;
117     end
118 end
119
120 best1_dlda = char(trainHeader(C(best_dlda_2,1)));
121 best2_dlda = char(trainHeader(C(best_dlda_2,2)));
122 fprintf('Best set of size-2 for DLDA is {%s, %s} and the resubstitution error is %f\n',best1_dlda,best2_dlda,res_error_dlda)
123
124 %Test Set error calculation
125 best_train_knn = [train_x(:,C(best_knn_2,:))];
126 best_mdl_knn = fitcknn(best_train_knn, train_y,'NumNeighbors',3);
127 best_test_knn = [test_x(:,C(best_knn_2,:))];
128 test_label_knn = best_mdl_knn.predict(best_test_knn);
129 test_error_knn = sum(abs(test_label_knn - test_y))/size(test_x,1);
130 fprintf('Test set error for 3NN on best set {%s, %s} of size-2 is %f\n',best1,best2,test_error_knn);
131
132 best_train_dlda = [train_x(:,C(best_dlda_2,:))];
133 best_mdl_dlda = fitcdiscr(best_train_dlda, train_y,'DiscrimType','diaglinear');
134 best_test_dlda = [test_x(:,C(best_dlda_2,:))];
135 test_label_dlda = predict(best_mdl_dlda, best_test_dlda);
136 test_error_dlda = sum(abs(test_label_dlda - test_y))/size(test_x,1);
137 fprintf('Test set error for DLDA on best set {%s, %s} of size-2 is %f\n',best1_dlda,best2_dlda,test_error_dlda);
138
139 fprintf("*****\n");
140 fprintf("*****Exhaustive for size-3*****\n");
141
142 %-----§
143 %size-3 feature set
144 C = combnk(2:71,3);
145
146 %3NN
147 for i=1:1:size(C,1)
148     new_train = [train_x(:,C(i,:))];
149     mdl_knn = fitcknn(new_train,train_y,'NumNeighbors',3);
150     new_train_label = mdl_knn.predict(new_train);
151
152     if i>1
153         if res_error > sum(abs(new_train_label - train_y))/size(train_x,1)
154             res_error = sum(abs(new_train_label - train_y))/size(train_x,1);
155             best_knn_3 = i;
156         end
157     else
158         res_error = sum(abs(new_train_label - train_y))/size(train_x,1);
159         best_knn_3 = 1;
160     end
161 end
162
163 best1 = char(trainHeader(C(best_knn_3,1)));

```

```

164 best2 = char(trainHeader(C(best_knn_3,2)));
165 best3 = char(trainHeader(C(best_knn_3,3)));
166 fprintf('Best set of size-3 for 3NN is {%s, %s, %s} and the resubstitution error is %f\n',best1,best2,best3,res_error)
167
168 %DLDA
169 for i=1:size(C,1)
170     new_train = [train_x(:,C(i,:))];
171     mdl_dlda = fitcdiscr(new_train,train_y,'DiscrimType','diaglinear');
172     new_train_label = predict(mdl_dlda, new_train);
173
174     if i>1
175         if res_error_dlda > sum(abs(new_train_label - train_y))/size(train_x,1)
176             res_error_dlda = sum(abs(new_train_label - train_y))/size(train_x,1);
177             best_dlda_3 = i;
178         end
179     else
180         res_error_dlda = sum(abs(new_train_label - train_y))/size(train_x,1);
181         best_dlda_3 = 1;
182     end
183 end
184
185 best1_dlda = char(trainHeader(C(best_dlda_3,1)));
186 best2_dlda = char(trainHeader(C(best_dlda_3,2)));
187 best3_dlda = char(trainHeader(C(best_dlda_3,3)));
188 fprintf('Best set of size-3 for DLDA is {%s, %s, %s} and the resubstitution error is %f\n',best1_dlda,best2_dlda,best3_dlda,
    res_error_dlda)
189
190 %Test Set error calculation
191 best_train_knn = [train_x(:,C(best_knn_3,:))];
192 best_mdl_knn = fitcknn(best_train_knn, train_y,'NumNeighbors',3);
193 best_test_knn = [test_x(:,C(best_knn_3,:))];
194 test_label_knn = best_mdl_knn.predict(best_test_knn);
195 test_error_knn = sum(abs(test_label_knn - test_y))/size(test_x,1);
196 fprintf('Test set error for 3NN on best set {%s, %s, %s} of size-3 is %f\n',best1,best2,best3,test_error_knn);
197
198 best_train_dlda = [train_x(:,C(best_dlda_3,:))];
199 best_mdl_dlda = fitcdiscr(best_train_dlda, train_y,'DiscrimType','diaglinear');
200 best_test_dlda = [test_x(:,C(best_dlda_3,:))];
201 test_label_dlda = predict(best_mdl_dlda, best_test_dlda);
202 test_error_dlda = sum(abs(test_label_dlda - test_y))/size(test_x,1);
203 fprintf('Test set error for DLDA on best set {%s, %s, %s} of size-3 is %f\n',best1_dlda,best2_dlda,best3_dlda,test_error_dlda);
204
205 fprintf("*****\n");
206
207 %-----§
208 %Sequential forward Search
209 opts = statset('display','iter');
210
211 fprintf("*****Iteration for 3NN*****\n");
212 [fs,history] = sequentialfs(@my_crit_knn,train_x,train_y,'cv','resubstitution','nfeatures',8,'options',opts,'direction','
    forward');
213
214 fprintf("*****\n"); fprintf("*****Iteration for DLDA*****\n");
215 [fs1,history1] = sequentialfs(@my_crit_dlda,train_x,train_y,'cv','resubstitution','nfeatures',8,'options',opts,'direction','
    forward');
216
217 fprintf("*****\n");
218
219 %Test set estimate
220 %3NN
221 for i=1:size(history.In,1)
222     new_train = [train_x(:,history.In(i,:)=1)];
223     fprintf("(")
224     features_selected = trainHeader(history.In(i,:)=1);
225     for j=1:i-1
226         fprintf("%s, ",char(features_selected(j)));
227     end
228     fprintf("%s} and the resub-error is %f\n",char(features_selected(i)),history.Crit(i));
229     mdl_knn_fs = fitcknn(new_train,train_y,'NumNeighbors',3);
230
231     test_knn_fs = [test_x(:,history.In(i,:)=1)];
232     test_error_knn_fs = sum(abs(mdl_knn_fs.predict(test_knn_fs)- test_y))/size(test_x,1);
233     fprintf("Test error of KNN with set size-%d is %f\n",i,test_error_knn_fs);
234 end
235
236 fprintf("*****\n");
237
238 %Test Set estimate
239 %DLDA
240 for i=1:size(history1.In,1)

```

```

241 new_train = [train_x(:,history1.In(i,:)==1)];
242 fprintf("(")
243 features_selected_dlda = trainHeader(history1.In(i,:)==1);
244 for j=1:1:i-1
245     fprintf("%s, ",char(features_selected_dlda(j)));
246 end
247 fprintf("%s} and the resub-error is %f\n",char(features_selected_dlda(i)),history1.Crit(i));
248 mdl_dlda_fs = fitcdiscr(new_train,train_y,'DiscrimType','diaglinear');
249
250 test_dlda_fs = [test_x(:,history1.In(i,:)==1)];
251 test_error_dlda_fs = sum(abs(predict(mdl_dlda_fs,test_dlda_fs)- test_y))/size(test_x,1);
252 fprintf("Test error of DLDA with set size-%d is %f\n\n",i,test_error_dlda_fs);
253 end
254
255 function val_knn = my_crit_knn(xT,yT,xt,yt)
256     mdl_knn = fitcknn(xT,yT,'NumNeighbors',3);
257     val_knn = sum(abs(predict(mdl_knn,xt) - yt)) ;
258 end
259
260 function val_dlda = my_crit_dlda(xT,yT,xt,yt)
261     mdl_dlda = fitcdiscr(xT,yT,'DiscrimType','diaglinear');
262     val_dlda = sum(abs(predict(mdl_dlda,xt) - yt)) ;
263 end

```

The results obtained from the program above when tabulated are as follows for  $3NN$

$3NN$ gene sets	Resubstitution error	True error
{CENPA}	0.100000	0.225532
{G5, CENPA}	0.033333	0.212766
{G1, FGF18, ORC6L}	0.016667	0.191489
{CENPA}	0.100000	0.225532
{G5, CENPA}	0.033333	0.212766
{G5, PECO1.1, CENPA}	0.016667	0.221277
{G5, G15, PECO1.1, CENPA}	0.016667	0.246809
{G5, G15, COL4A2, PECO1.1, CENPA}	0.016667	0.238298
{G5, G8, G15, COL4A2, PECO1.1, CENPA}	0.016667	0.238298
{G5, G8, DC13, G15, COL4A2, PECO1.1, CENPA}	0.016667	0.221277
{G4, G5, G8, DC13, G15, COL4A2, PECO1.1, CENPA}	0.033333	0.204255

The results for the *DLDA* obtained are:

<i>DLDA</i> gene sets	Resubstitution error	True error
{ORC6L}	0.133333	0.170213
{G3, G5}	0.100000	0.204255
{MMP9, IGFBP5.1, CENPA}	0.050000	0.195745
{ORC6L}	0.133333	0.170213
{G4, ORC6L}	0.116667	0.170213

{G4, FLT1, ORC6L}	0.116667	0.178723
{G4, ALDH4, FLT1, ORC6L}	0.116667	0.174468
{G4, ALDH4, KIAA1442, FLT1, ORC6L}	0.116667	0.148936
{G4, ALDH4, KIAA1442, FLT1, ORC6L, IGFBP5.1}	0.100000	0.136170
{G4, ALDH4, KIAA1442, FLT1, MMP9, ORC6L, IGFBP5.1}	0.083333	0.127660
{G2, G4, ALDH4, KIAA1442, FLT1, MMP9, ORC6L, IGFBP5.1}	0.083333	0.127660

From the results we can infer that,

- the resubstitution error produced in the case of gene sets obtained from the exhaustive search are less than those obtained from the sequential forward search. Also, we can observe that the apparent error converges to zero early in exhaustive search method when compared to sequential forward search. However, in terms of computational efficiency sequential feature selection is more preferred due to its practical applicability as the maximum number of computations done is  $O(n^2)$  whereas its counterpart takes  $O(n!)$  time.

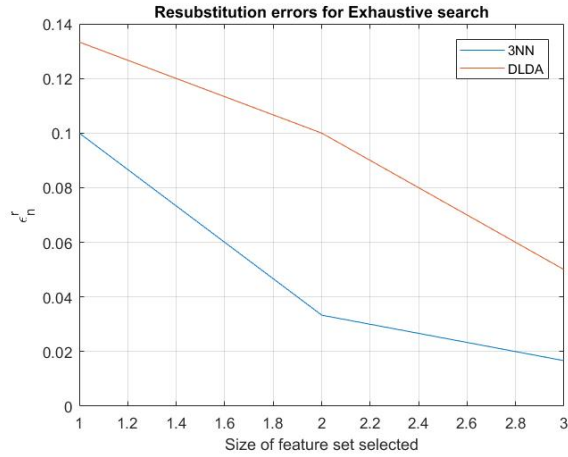


Figure 1: Plot of  $\varepsilon_n^r$  for Exhaustive search for 3NN and DLDA

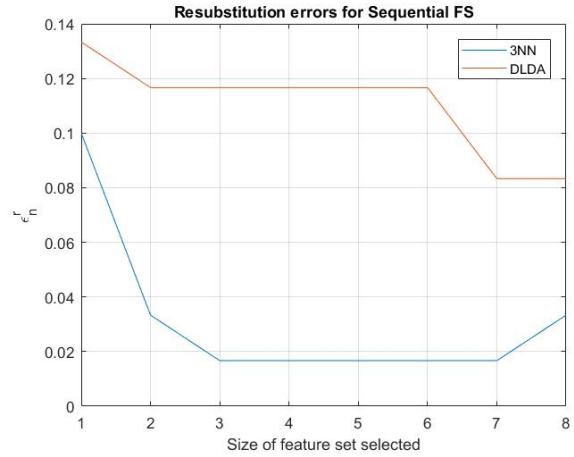


Figure 2: Plot of  $\varepsilon_n^r$  for Sequential forward search for 3NN and DLDA

- From the plot in *Figure – 1*, above, we can also observe the training error is monotonically decreasing in the exhaustive case while this is not represented in the case of sequential feature selection, this can be attributed due to the fact that since once a feature is selected, it gets

frozen even though a better combination exists. However, this monotonic behavior is not observed for the true error which in the case of exhaustive search, which might imply that the exhaustive search was over-fitting the model using the training data.

- From both the figures above, the resubstitution error in the case of  $3NN$  are less than those obtained from  $DLDA$ , whereas the test set estimates are lower for  $DLDA$ . This can be explained from the fact that  $3NN$  might be over-fitting a complex model which  $DLDA$  fails to do irrespective of the data. So, we can conclude that,  $DLDA$  appears to be a better classifier in this case of the data compared to  $3NN$ .
- In order to build better models, we can implement strategies like Plus- $l$  take- $r$  search which might perform well as there is no freezing of a feature as in the case of sequential search. Also implementing a complex classification algorithm like neural networks and support vector machines we can perform well along with increasing the number of samples considered in the training set in designing the classifier.