

MATH444 Final Project

Jiasen Zhang

Problem 1 and 2

In the first 2 parts I load the data and reduce the gray scale to 6 levels. It's really enough to get good results.

Problem 3 and 4

Compute 3 GLCMs using 3 pairs of offset parameters: (2,0), (0,2) and (2,2). Stack their columns into vecotrs and put together so that the dimensionality is reduced to 108. Then compute the first 4 principal components to reduce the dimensionality to 4.

The different pairs of prinicipal components are plotted in Figure 1. The corresponding colors of the classes are shown in Table 1. We can see even with a drastic dimensionality reduction, the classes are separated well. It looks like the green points and purple points are more difficult to separate, but in fact most green points just stay around rather than mix with purple points.

Classes	1	2	3	4	5
Color	Red	Blue	Green	Purple	Black

Table 1

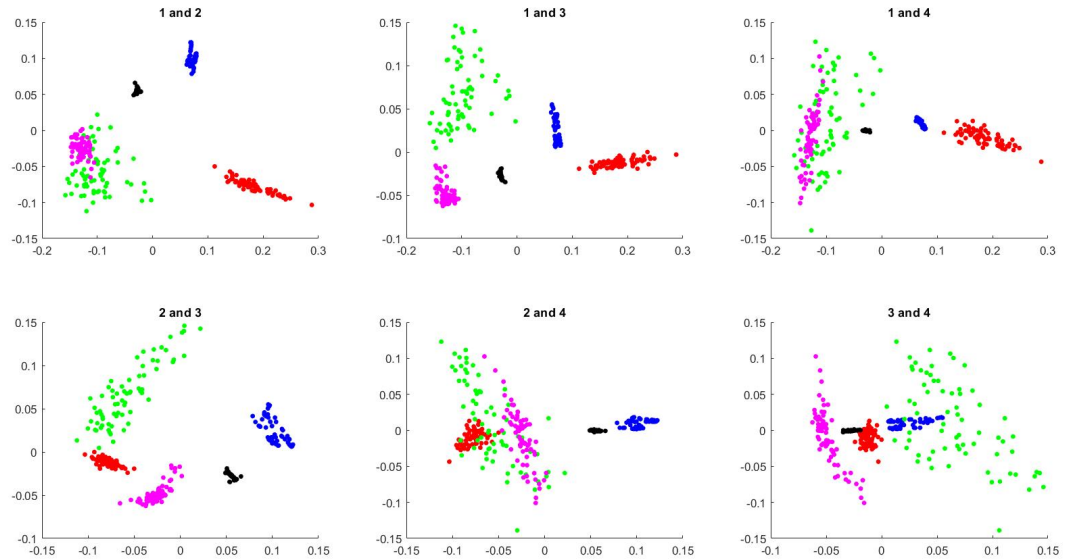


Figure 1: The first 4 principal components.

Problem 5

Using the same parameters and gray scale, compute the first 3 LDA separators and project the principal components onto these directions. As shown in Figure 2, the 5 classes are separated really well. It reflects that LDA generally performs better than PCA by making use of annotations. It also indicates that these images can be separated well with knowing the annotations, it makes sure that classification tree is suitable for these images.

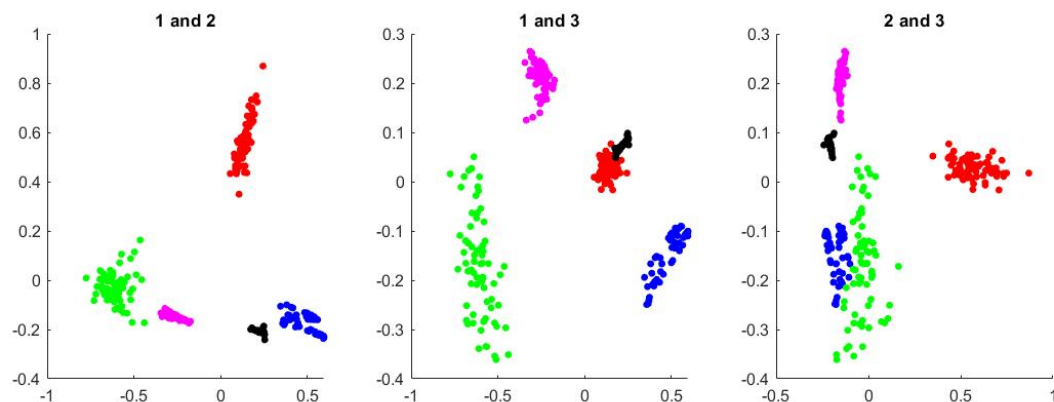


Figure 2: The first 3 LDA components (6 levels).

If I reduce the gray scale to 4 levels, the result is good too, as shown in Figure 3.

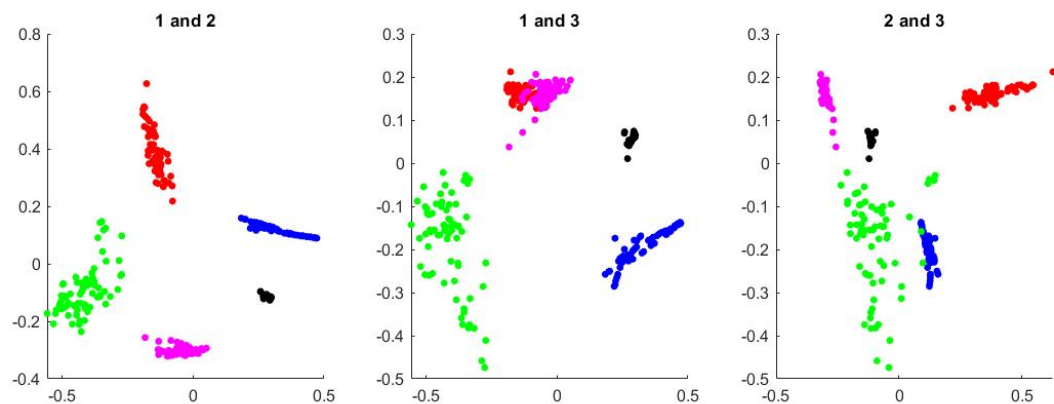


Figure 3: The first 3 LDA components (4 levels).

Problem 6

With the PCA results, select the first 200 data as the training set and generate the maximal tree, the structure is shown in Figure 4. There are 27 intervals totally with 14 pure leaves.

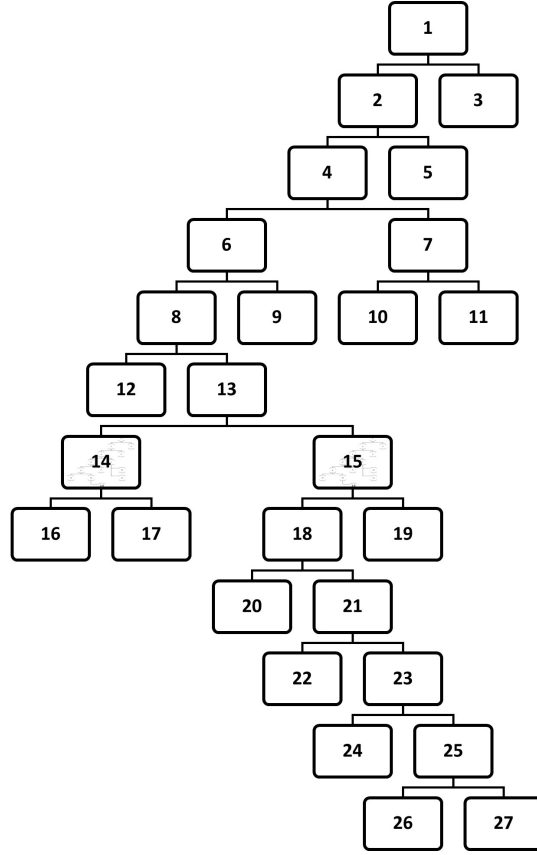


Figure 4: The maximal tree.

Next prune the tree. The times of pruning is 4. The pruned nodes and α are shown in Table 2. We can see α is an increasing sequence. The numbers of nodes are shown in Figure 4. In the 4th row, I test the trees with testing data and compute the rates of success. From the table we can see the maximal tree and that pruned for one time have the highest rate of success 0.89. Because the pruned tree has lower complexity, it can be the optimal tree (the tree pruned for one time).

Pruning steps	0	1	2	3	4
Pruned nodes	N/A	7, 23	18	6	1
α	0	0.005	0.0125	0.015	0.0183
Rate of success	0.89	0.89	0.88	0.76	0.17

Table 2

Problem 7

The structure of the optimal tree is shown in Figure 5. Use it to classify the rest of 50 images, the rate of success is 0.88. The confusion matrix is:

$$\begin{bmatrix} 9 & 0 & 0 & 0 & 0 \\ 0 & 11 & 0 & 0 & 0 \\ 0 & 0 & 8 & 0 & 0 \\ 0 & 0 & 4 & 6 & 0 \\ 0 & 0 & 2 & 0 & 10 \end{bmatrix}$$

In the matrix we can see the class 3 and 4 are difficult to separate, this result corresponds to Figure 1, in which the class 3 and 4 stay really close.

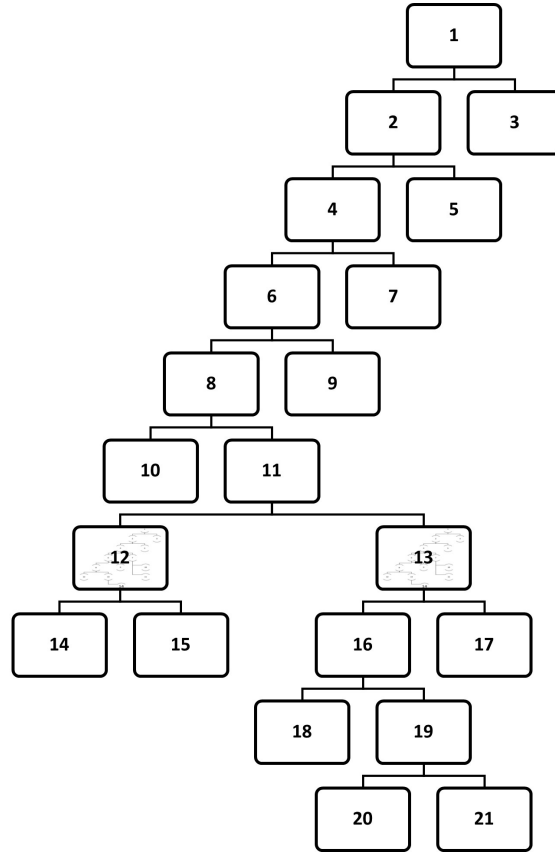


Figure 5: The optimal tree.

Matlab Code

Problem 1-7

The code has 7 parts and two subroutines.

```
1 clear all;clc;
2 % Math444 Final Project
3 % Jiasen Zhang
4
5 load TestImages.mat;
6 p = 350;
7
8
9 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
10 % 2: Reduce the gray scales of the images
11 k = 6; % number of gray scales
12 unit = 1/k;
13 X = ceil(X/unit);
14 X(X==0)=1;
15 %imagesc(reshape(X(:,2),128,128));colormap(gray);
16
17
18 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
19 % 3: Three GLCMs
20 GLCM = zeros(3*k^2,p);
21 mu=2;
22 nu=0;
23 for n=1:p
24     C = reshape(X(:,n),128,128);
25     C_shift = NaN(128,128);
26     C_shift(1:128-mu,1:128-nu) = C(1+mu:128,1+nu:128);
27     G = zeros(k,k);
28     for i=1:k
29         Ii = (C==i);
30         for j=1:k
31             Ij = (C_shift==j);
32             G(i,j) = sum(sum(Ii.*Ij));
33         end
34     end
35     G = G/(sum(sum(G))); % normalize G
36     GLCM(1:k^2,n) = reshape(G/3,k^2,1);
37 end
38 mu=0;
39 nu=2;
40 for n=1:p
41     C = reshape(X(:,n),128,128);
42     C_shift = NaN(128,128);
43     C_shift(1:128-mu,1:128-nu) = C(1+mu:128,1+nu:128);
44     G = zeros(k,k);
45     for i=1:k
46         Ii = (C==i);
47         for j=1:k
48             Ij = (C_shift==j);
49             G(i,j) = sum(sum(Ii.*Ij));
50         end
51     end
52     G = G/(sum(sum(G))); % normalize G
53     GLCM(k*k+1:2*k^2,n) = reshape(G/3,k^2,1);
54 end
55 mu=2;
56 nu=2;
57 for n=1:p
58     C = reshape(X(:,n),128,128);
59     C_shift = NaN(128,128);
60     C_shift(1:128-mu,1:128-nu) = C(1+mu:128,1+nu:128);
61     G = zeros(k,k);
```

```

62     for i=1:k
63         Ii = (C==i);
64         for j=1:k
65             Ij = (C_shift==j);
66             G(i,j) = sum(sum(Ii.*Ij));
67         end
68     end
69     G = G/(sum(sum(G))); % normalize G
70     GLCM(2*k*k+1:3*k^2,n) = reshape(G/3,k^2,1);
71 end
72 % Compute the first 4 principal components
73 GLCMbar = mean(GLCM,2);
74 GLCMc = GLCM - GLCMbar; % centered data
75 [u,~,~] = svd(GLCMc);
76 z = u(:,1:4)'*GLCMc;
77
78
79 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
80 % 4: plot using different colors or the same color
81 nplot = 1;
82 figure(1);
83 for i=1:3
84     for j=i+1:4
85         subplot(2,3,nplot);
86         scatter(z(i,I==1),z(j,I==1),'r','k.','SizeData',200); hold on;
87         scatter(z(i,I==2),z(j,I==2),'b','k.','SizeData',200);
88         scatter(z(i,I==3),z(j,I==3),'g','k.','SizeData',200);
89         scatter(z(i,I==4),z(j,I==4),'m','k.','SizeData',200);
90         scatter(z(i,I==5),z(j,I==5),'k','k.','SizeData',200);
91         title(strcat(num2str(i),{32},'and',{32},num2str(j)));
92         nplot = nplot+1;
93     end
94 end
95
96
97 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
98 % 5: Compute the first 3 LDA separators
99 [zLDA,~] = lda(5, z, I);
100 nplot=1;
101 figure(2);
102 for i=1:2
103     for j=i+1:3
104         subplot(1,3,nplot);
105         scatter(zLDA(i,I==1),zLDA(j,I==1),'r','k.','SizeData',200); hold on;
106         scatter(zLDA(i,I==2),zLDA(j,I==2),'b','k.','SizeData',200);
107         scatter(zLDA(i,I==3),zLDA(j,I==3),'g','k.','SizeData',200);
108         scatter(zLDA(i,I==4),zLDA(j,I==4),'m','k.','SizeData',200);
109         scatter(zLDA(i,I==5),zLDA(j,I==5),'k','k.','SizeData',200);
110         title(strcat(num2str(i),{32},'and',{32},num2str(j)));
111         nplot = nplot+1;
112     end
113 end
114
115
116 % %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
117 % Training set and testing set
118 I_train = I(1:200);
119 I_test = I(201:300);
120 I0 = I(301:end);
121 z_train = z(:,1:200);
122 z_test = z(:,201:300);
123 z0 = z(:,301:end);
124
125
126 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
127 % 6: Classification tree: using z_train
128 % initialize the tree
129 % structure R

```

```

130 R(1).idx = 1:200; % index of data in the interval
131 R(1).data = z_train;
132 R(1).n = 1; % number of R(1)
133 R(1).split_idx = [];
134 R(1).split_value = [];
135 R(1).leftR = []; % left child
136 R(1).rightR = []; % right child
137 [R(1).g, R(1).c] = gini(I_train(R(1).idx)); % gini index and classification
138 Lpure = [];
139 Lmixed = [1];
140 m = 1;
141
142 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
143 % 6.1: Build a maximal tree
144 while(1)
145 % pick one index from Lmixed randomly:
146 pick = Lmixed(1);
147 Lmixed = Lmixed(2:end);
148
149 % find split index and value of R(pick)
150 [j, s] = optimal(R(pick).idx, I_train, z_train);
151 R(pick).split_idx = j;
152 R(pick).split_value = s;
153
154 % build idxleft and idxright, generate m+1 and m+2
155 R(pick).leftR = m+1;
156 R(pick).rightR = m+2;
157 dataj = R(pick).data(j,:);
158 idx_left = find(dataj<=s);
159 idx_right = find(dataj>s);
160 %
161 R(m+1).idx = R(pick).idx(idx_left);
162 R(m+1).data = R(pick).data(:,idx_left);
163 R(m+1).n = m+1;
164 R(m+2).idx = R(pick).idx(idx_right);
165 R(m+2).data = R(pick).data(:,idx_right);
166 R(m+2).n = m+2;
167
168 % decide if new leaves are pure or mixed
169 [R(m+1).g, R(m+1).c] = gini(I_train(R(m+1).idx));
170 if R(m+1).g==0
171     Lpure = [Lpure,m+1];
172 else
173     Lmixed = [Lmixed,m+1];
174 end
175 [R(m+2).g, R(m+2).c] = gini(I_train(R(m+2).idx));
176 if R(m+2).g==0
177     Lpure = [Lpure,m+2];
178 else
179     Lmixed = [Lmixed,m+2];
180 end
181 m=m+2;
182
183 % If all leaves are pure, stop
184 if isempty(Lmixed)==1
185     break;
186 end
187
188 end
189 Rmax = R;
190
191 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
192 % 6.2: Prune the tree
193 % initial sucesor matrix A
194 A = zeros(m,m);
195 for i=1:m
196     if ~isempty(R(i).leftR)
197         A(i,R(i).leftR)=1;

```

```

198     A(i,R(i).rightR)=1;
199     end
200 end
201 alpha0 = [0]; % save alpha
202 T(1,:) = R; % put the maximal tree into T
203 nNodes = m;
204 numLeaf = []; % number of leaves of each pruned tree
205
206 while(m>1)
207     %for nn=1:4
208     % update Lpure and Lmixed
209     Lpure = [];
210     Lmixed = [];
211     for i=1:m
212         if isempty(R(i).leftR)==1 % if it's a leaf
213             if R(i).g==0
214                 Lpure = [Lpure, i];
215             else
216                 Lmixed = [Lmixed, i];
217             end
218         end
219     end
220     % genealogy matrix G
221     AA = A;
222     G = zeros(m,m);
223     while(sum(sum(AA))>0)
224         G = G+AA;
225         AA = A*AA;
226     end
227     % find non-leaf intervals for the tree
228     Jleaf = [];
229     Jnonleaf = [];
230     for i=1:m
231         if sum(G(i,:))==0
232             Jleaf = [Jleaf,i];
233         else
234             Jnonleaf = [Jnonleaf,i];
235         end
236     end
237     numLeaf = [numLeaf,length(Jleaf)];
238     % compute alpha for each nonleaf interval
239     alpha = zeros(1,length(Jnonleaf));
240     for i=1:length(Jnonleaf)
241         j0 = Jnonleaf(i);
242         % frequency multiply misclassification rate: v*r
243         if R(j0).c == R(R(j0).leftR).c
244             vr = length(R(R(j0).rightR).idx)/length(R(1).idx);
245         elseif R(j0).c == R(R(j0).rightR).c
246             vr = length(R(R(j0).leftR).idx)/length(R(1).idx);
247         end
248         % compute # of pure leaves for each subtree
249         successor = find(G(j0,:));
250         %nL = sum(ismember([R(successor).n],Lpure))+sum(ismember([R(successor).n],Lmixed));
251         nL = sum(ismember(successor,Lpure))+sum(ismember(successor,Lmixed));
252         % v and r of leaves
253         nominator = vr;
254         for j=1:length(successor)
255             if isempty(R(successor(j)).leftR) % If it's a leaf
256                 vrL = sum(I-train(R(successor(j)).idx)~=R(successor(j)).c)/length(R(1).idx);
257                 nominator = nominator - vrL;
258             end
259         end
260         % compute alpha
261         alpha(i) = nominator/(nL-1);
262     end
263     % choose smallest alpha and prune the corresponding subtree
264     alpha0 = [alpha0,min(alpha)];
265     nPrune = Jnonleaf(alpha==min(alpha))

```



```

266 nKeep = find(sum(G(nPrune,:),1)==0);
267 R = R(nKeep);
268 A = A(nKeep,nKeep);
269 % m, and save the pruned tree into T
270 m = length(nKeep);
271 nNodes = [nNodes,m];
272 for ell=1:m
273     if sum(A(ell,:))==0
274         R(ell).split_idx = [];
275         R(ell).split_value = [];
276         R(ell).leftR = [];
277         R(ell).rightR = [];
278     else
279         i_1 = find(A(ell,:)==1,1,'first');
280         i_2 = find(A(ell,:)==1,1,'last');
281         R(ell).leftR = i_1;
282         R(ell).rightR = i_2;
283     end
284 end
285 T(end+1,1:m) = R;
286 end
287 numLeaf = [numLeaf,1];
288 alpha0
289
290
291
292 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
293 % 6.3: choose the tree with testing data
294 [numTree,~] = size(T);
295 rate = zeros(1,numTree);
296 for n=1:numTree % for each tree
297     % classification of testing data
298     R = T(n,1:nNodes(n));
299     class = zeros(1,length(I_test)); % result of classification
300     for k=1:length(I_test) % for each data vector in the testing data
301         i=1;
302         while(isempty(R(i).leftR)==0) % if it's not a leaf
303             j = R(i).split_idx;
304             s = R(i).split_value;
305             if z_test(j,k)<=s
306                 i = R(i).leftR;
307             elseif z_test(j,k)>s
308                 i = R(i).rightR;
309             end
310         end
311         class(k) = R(i).c;
312     end
313     rate(n) = sum(class==I_test)/length(I_test);
314     % confusionM = zeros(5,5);
315     % for k=1:length(I_test)
316     %     confusionM(class(k),I_test(k)) = confusionM(class(k),I_test(k))+1;
317     % end
318
319 end
320 rate
321
322
323
324 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
325 % 7: Test the tree using the rest of data
326 n=find(rate==max(rate));
327 n=n(end)
328 R = T(n,1:nNodes(n));
329 class = zeros(1,length(I0)); % result of classification
330 for k=1:length(I0) % for each data vector in the testing data
331     i=1;
332     while(isempty(R(i).leftR)==0) % if it's not a leaf
333         j = R(i).split_idx;

```

```

334         s = R(i).split_value;
335         if z0(j,k) ≤ s
336             i = R(i).leftR;
337         elseif z0(j,k) > s
338             i = R(i).rightR;
339         end
340     end
341     class(k) = R(i).c;
342
343 end
344 confusionM = zeros(5,5);
345 for k=1:length(I0)
346     confusionM(class(k),I0(k)) = confusionM(class(k),I0(k))+1;
347 end
348 rate = sum(class==I0)/length(I0)
349 confusionM
350
351 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
352 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
353 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
354 % function of optimal splitting
355 function [split_idx, split_value] = optimal(I, C, X)
356 % optimal splitting
357 % I = R(k).idx
358 % C = full annotation
359 % X = full data matrix
360
361 q = length(I);
362 RX = X(:,I); % data of I
363 RC = C(I); % annotation of I
364 giniRC = gini(RC);
365 %split_idx = zeros(size(X,1),1);
366 split_value = zeros(size(X,1),1);
367 dg = zeros(size(X,1),1);
368 dgn = zeros(1,q-1);
369 sigma=zeros(1,q-1);
370 % for each attribute
371 for n=1:size(X,1) % n=1:4
372     % sort RX with attribute n
373     [sortRX, idx] = sort(RX(n,:));
374     RC = RC(idx);
375     % for each data of attribute n
376     for j=1:q-1
377         % compute q-1 split values
378         sigma(j) = (sortRX(j)+sortRX(j+1))/2;
379         % compute Rr and Rl from RC
380         idx1=(sortRX ≤ sigma(j));
381         idx2=(sortRX > sigma(j));
382         Rl = RC(idx1);
383         Rr = RC(idx2);
384         % compute q-1 impurity change
385         dgn(j) = giniRC - length(Rl)*gini(Rl)/q - length(Rr)*gini(Rr)/q;
386     end
387     temp = sigma(dgn==max(dgn));
388     % split value and gini index for attribute n=1:4
389     split_value(n) = (temp(1)+temp(end))/2;
390     temp = dgn(dgn==max(dgn));
391     dg(n) = (temp(1)+temp(end))/2;
392 end
393
394 % find the largest impurity change
395 temp = find(dg==max(dg));
396 temp = temp(1);
397 split_idx = temp;
398 split_value = split_value(temp);
399
400 end
401

```

```

402
403 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
404 % function of gini index
405 function [g,c] = gini(annotation)
406 % g = gini index
407 % c = classification by majority vote
408 P = zeros(5,1);
409 g = 1;
410 for j=1:5
411     P(j) = sum(annotation==j)/length(annotation); % frequency
412     g = g - P(j)^2;
413 end
414 c = find(P == max(P)); % classification by majority vote
415 c = c(unidrnd(length(c)));
416
417 end

```

LDA function

```

1 function [z,Q] = lda(k, X, I)
2 % Jiasen Zhang LDA
3
4 % input
5 % k = number of clusters
6 % X = data
7 % I = partition
8
9 % output
10 % z = LDA reduced data
11 % Q = seperating direction
12
13 [n,p]=size(X);
14
15 c=mean(X,2); % global mean
16 ck = zeros(n,k); % cluster means
17 S_w = zeros(n,n);
18 for L=1:k
19     % cluster means
20     ck(:,L) = mean(X(:,I==L),2);
21     % compute S_L and S_w
22     X_LC = X(:,I==L) - ck(:,L);
23     % S_L = X_LC*X_LC';
24     S_w = S_w + X_LC*X_LC';
25 end
26
27 % compute S_b
28 X_bar = zeros(n,p);
29 for j=1:p
30     L = I(j); % cluster
31     X_bar(:,j) = ck(:,L);
32 end
33 S_b = (X_bar-c)*(X_bar-c)';
34
35 % get matrix A
36 d1 = max(eig(S_w));
37 tau = 1e-16; % if not positive definite, change to 1e-16
38 S_we = S_w + tau*d1*d1*eye(n);
39 K = chol(S_we);
40 A = inv(K)'*S_b/K;
41
42 % nQ largest eigenvectors of A and solve Q
43 nQ = 4;
44 Q = zeros(n,nQ);
45 [v,d]=eig(A);
46 d = diag(d)/sum(diag(d)); % proportion of trace

```

```

47 fprintf('Proportion of trace:\n');
48 for j=1:nQ
49     maxindex = find(d==max(d)); % find the maximum eigenvalue
50     Q(:,j) = K\ v(:,maxindex);
51     fprintf('%0.8f\n',d(maxindex)); % show proportions of first k-1 eigenvalues
52     d(maxindex)=0;
53 end
54
55 % LDA reduced data
56 Z = Q'*X;
57 end

```