

MATH444 HW3

Jiasen Zhang

1 Problem 1

Test my LDA algorithm in Figure 1 using Gaussian distribution data. The group centers are:

$$\begin{bmatrix} -2.9958 \\ -3.0944 \\ -3.1437 \end{bmatrix} \quad \begin{bmatrix} 2.0308 \\ 0.1135 \\ -0.0254 \end{bmatrix} \quad \begin{bmatrix} 4.0009 \\ 4.0621 \\ -6.0079 \end{bmatrix}$$

By computing the two LDA separating directions, the 3 clusters are classified well.

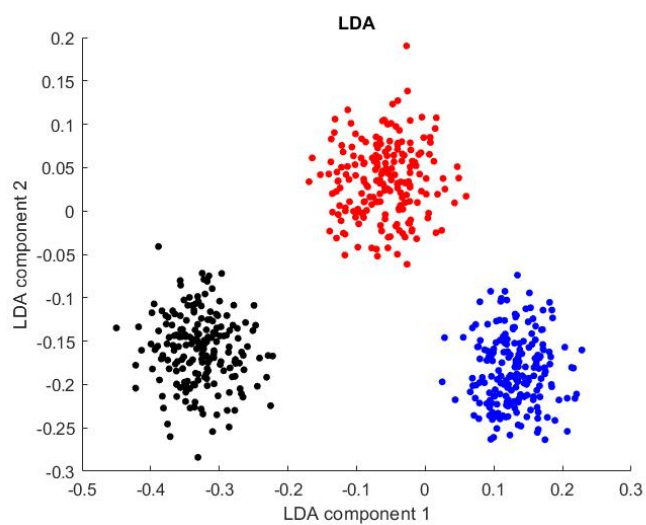


Figure 1: Problem 1

2 Problem 2

Figure 2 shows the plots of the two LDA directions(left) and the first two principal components(right).

Result of LDA is similar to that of PCA, the blue cluster is separated well, while the black and red clusters are close to each other.

However, the LDA result is better. By projecting along the maximally separated direction, the black and red clusters in LDA are more distinguishable. There are less points in the mixing area between the black and red clusters.

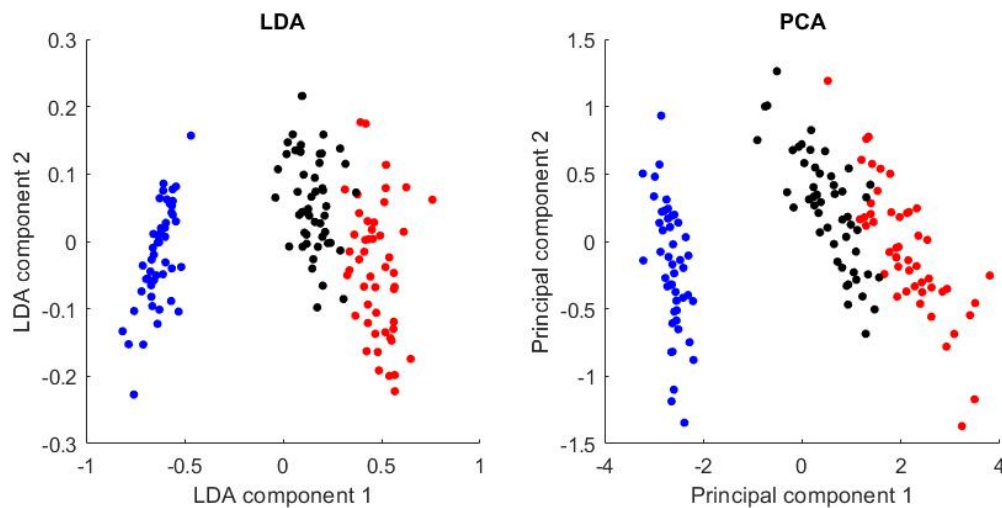


Figure 2: Problem 2: IrisData

3 Problem 3

3.1 (a)

The LDA result is shown in the left Figure 3. In my opinion, the attributes don't carry enough information to separate the two clusters, because there are still a small overlapping area between the two clusters. But I think the differences between the two clusters are characterized well, because the peaks of two clusters are not close to each other, and the overlapping area is small.

3.2 (b)

The PCA result is shown in the right Figure 3. Both of the two results have two peaks for each cluster, and have an overlapping area. However, the result of LDA is much better. For LDA result, the peaks of the two clusters are farther, and the overlapping area is smaller.

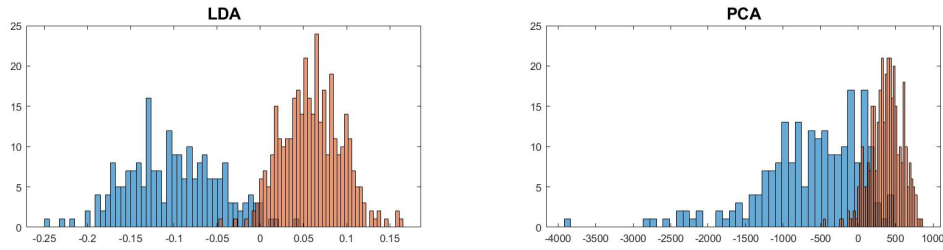


Figure 3: Problem 3: WisconsinBreastCancerData

3.3 (c)

I choose the larger components 6, 15, 17, 18, 20, 28 and 30 and remove them. The dimension is reduced from 30 to 7. As shown in Figure 4, the LDA result with reduced data still captures the differences between the clusters, but it's not better than the result with full data. It still can't separate the two clusters completely, although PCA result is better. Although some attributes play less important roles in classifying, we can't overlook their effects.

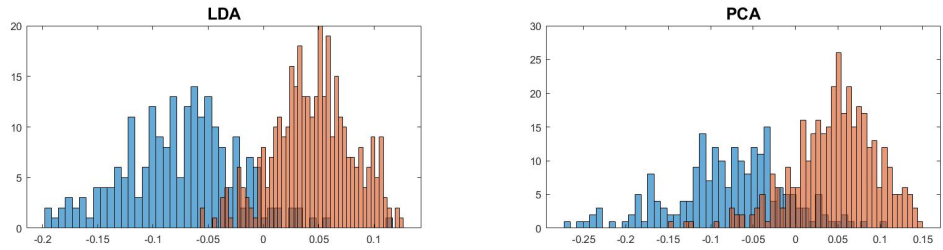


Figure 4: Problem 3: WisconsinBreastCancerData (reduced)

4 Problem 4

The result is shown in Figure 5. Based on the collected attributes, the three clusters are neatly separated.

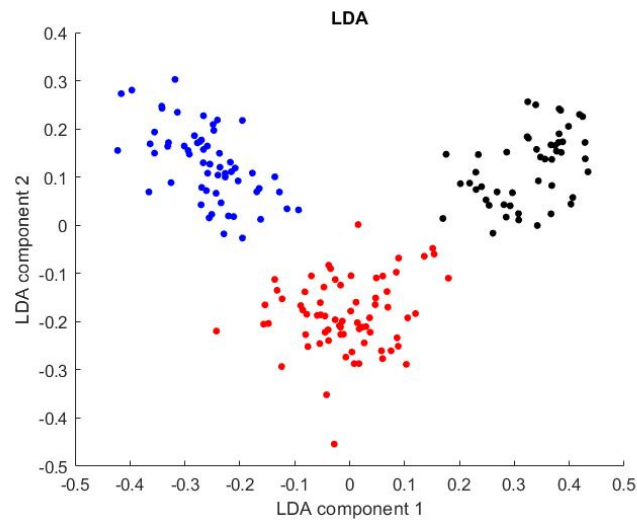


Figure 5: Problem 4: WineData

I choose the six largest components 1, 6, 7, 8, 11 and 12 and remove them. The dimension is reduced from 13 to 6. As shown in Figure 6, the three clusters are still separated but the result is worse than above.

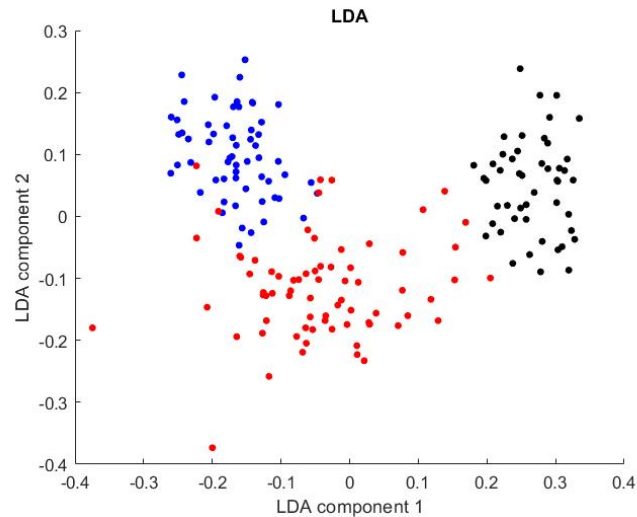


Figure 6: Problem 4: WineData (reduced)

5 Problem 5

Compute the three LDA separating vectors and plot them in Figure 7 and 8. Figure 7 shows LDA component 1,2. Figure 8 shows LDA component 1,3 and 2,3.

In Figure 7, the first two LDA directions separate the four groups well. While in Figure 8, both of the two pairs of LDA directions can't separate the groups clearly.

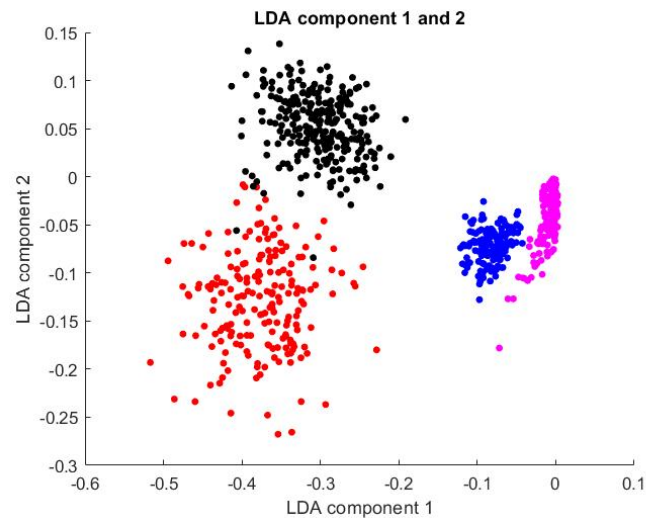


Figure 7: Problem 5: ForestSpectra

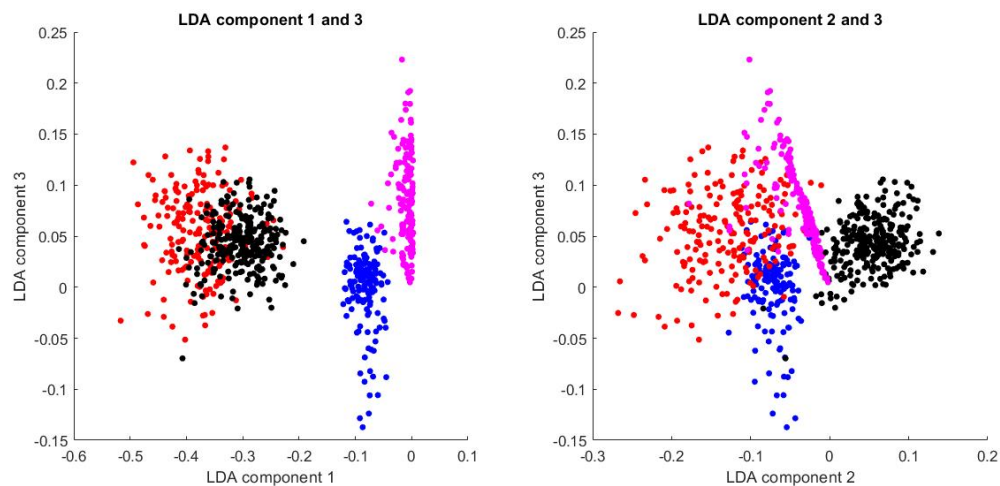


Figure 8: Problem 5: ForestSpectra

6 Matlab Code

6.1 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 = 0; % 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 % k-1 largest eigenvectors of A and solve Q
43 Q = zeros(n,k-1);
44 [v,d]=eig(A);
45 d = diag(d)/sum(diag(d)) % proportion of trace
46 fprintf('Proportion of trace:\n');
47 for j=1:k-1
48     maxindex = find(d==max(d)); % find the maximum eigenvalue
49     Q(:,j) = K\v(:,maxindex);
50     fprintf('%0.8f\n',d(maxindex)); % show proportions of first k-1 eigenvalues
51     d(maxindex)=0;
52 end
53
54 % LDA reduced data
55 z = Q'*X;
56 end
```

6.2 Problem 1

```

1 clear all;clc;
2
3 % three clusters
4 I = ones(600,1);
5 I(201:400)=2;
6 I(401:600)=3;
7
8 % generate data X
9 X = zeros(3,600);
10 X(1,1:200)=normrnd(-3,1,1,200);
11 X(2,1:200)=normrnd(-3,1,1,200);
12 X(3,1:200)=normrnd(-3,1,1,200);
13 X(1,201:400)=normrnd(2,1,1,200);
14 X(2,201:400)=normrnd(0,1,1,200);
15 X(3,201:400)=normrnd(0,1,1,200);
16 X(1,401:600)=normrnd(4,1,1,200);
17 X(2,401:600)=normrnd(4,1,1,200);
18 X(3,401:600)=normrnd(-6,1,1,200);
19 cov(X(:,1:200)')
20 cov(X(:,201:400)')
21 cov(X(:,401:600)')
22 Xbar = mean(X,2);
23 Xc = X - Xbar;
24
25 k = 3;
26 [z,Q] = lda(k, X, I);
27
28 % plot
29 scatter(z(1,I==1),z(2,I==1),'b','k.','SizeData',200);hold on;
30 scatter(z(1,I==2),z(2,I==2),'r','k.','SizeData',200);
31 scatter(z(1,I==3),z(2,I==3),'k','k.','SizeData',200);
32 xlabel('LDA component 1');
33 ylabel('LDA component 2');
34 title('LDA');

```

6.3 Problem 2

```

1 clear all;clc;
2
3 % Given data and partition
4 load IrisData
5 [n,p]=size(X);
6 I = zeros(1,p);
7 I(1:50)=1;
8 I(51:100)=2;
9 I(101:end)=3;
10 Xbar = mean(X,2);
11 Xc = X - Xbar;
12
13 k = 3;
14 [z,Q] = lda(k, Xc, I);
15
16 subplot(1,2,1);
17 scatter(z(1,I==1),z(2,I==1),'b','k.','SizeData',200);hold on;
18 scatter(z(1,I==2),z(2,I==2),'r','k.','SizeData',200);
19 scatter(z(1,I==3),z(2,I==3),'k','k.','SizeData',200);
20 xlabel('LDA component 1');
21 ylabel('LDA component 2');
22 title('LDA');
23
24 [u,~,~] = svd(Xc);
25 z = u(:,1:2)'*Xc;
26 subplot(1,2,2);
27 scatter(z(1,I==1),z(2,I==1),'b','k.','SizeData',200);hold on;

```

```

28 scatter(z(1,I==2),z(2,I==2),'r','k.','SizeData',200);
29 scatter(z(1,I==3),z(2,I==3),'k','k.','SizeData',200);
30 xlabel('Principal component 1');
31 ylabel('Principal component 2');
32 title('PCA');

```

6.4 Problem 3

```

1 clear all;clc;
2
3 % Given data and partition
4 load WisconsinBreastCancerData
5 X = Data.WCD.Matrix;
6 I = I.Label;
7 [n,p]=size(X);
8 Xbar = mean(X,2);
9 Xc = X - Xbar;
10
11 k = 2;
12 [z,Q] = lda(k, Xc, I);
13
14 subplot(1,2,1);
15 histogram(z(I==1),50);hold on;
16 histogram(z(I==2),50);
17 title('LDA','FontSize',16);
18 subplot(1,2,2);
19 [u,~,~] = svd(Xc);
20 z0 = u(:,1)'*Xc;
21 histogram(z0(I==1),50);hold on;
22 histogram(z0(I==2),50);
23 title('PCA','FontSize',16);
24
25 % (c)
26 % select the largest components of Q:
27 % component 6,15,17,18,20,28,30
28 choose = [6,15,17,18,20,28,30];
29 Xc2 = Xc(choose,:);
30
31 % repeat
32 [z,Q0] = lda(k, Xc2, I);
33 figure();
34 subplot(1,2,1);
35 histogram(z(I==1),50);hold on;
36 histogram(z(I==2),50);
37 title('LDA','FontSize',16);
38 subplot(1,2,2);
39 [u,~,~] = svd(Xc2);
40 z0 = u(:,1)'*Xc2;
41 histogram(z0(I==1),50);hold on;
42 histogram(z0(I==2),50);
43 title('PCA','FontSize',16);

```

6.5 Problem 4

```

1 clear all;clc;
2
3 load WineData
4 Xbar = mean(X,2);
5 Xc = X - Xbar;
6
7 k = 3;

```



```

8  [z,Q] = lda(k, Xc, I);
9
10 scatter(z(1,I==1),-z(2,I==1),'b','k.','SizeData',200);hold on;
11 scatter(z(1,I==2),-z(2,I==2),'r','k.','SizeData',200);
12 scatter(z(1,I==3),-z(2,I==3),'k','k.','SizeData',200);
13 xlabel('LDA component 1');
14 ylabel('LDA component 2');
15 title('LDA');
16
17 % select the largest components:
18 % component 1,6,7,8,11,12
19 choose = [1,6,7,8,11,12];
20 Xc2 = Xc(choose,:);
21 [z,Q2] = lda(k, Xc2, I);
22
23 figure();
24 scatter(z(1,I==1),z(2,I==1),'b','k.','SizeData',200);hold on;
25 scatter(z(1,I==2),z(2,I==2),'r','k.','SizeData',200);
26 scatter(z(1,I==3),z(2,I==3),'k','k.','SizeData',200);
27 xlabel('LDA component 1');
28 ylabel('LDA component 2');
29 title('LDA');

```

6.6 Problem 5

```

1  clear all;clc;
2
3  load ForestSpectra
4  k = 4;
5  I=Itype;
6  [z,Q] = lda(k, X, I);
7
8  % LDA components 1,2
9  scatter(z(1,I==1),z(2,I==1),'b','k.','SizeData',200);hold on;
10 scatter(z(1,I==2),z(2,I==2),'r','k.','SizeData',200);
11 scatter(z(1,I==3),z(2,I==3),'k','k.','SizeData',200);
12 scatter(z(1,I==4),z(2,I==4),'m','k.','SizeData',200);
13 xlabel('LDA component 1');
14 ylabel('LDA component 2');
15 title('LDA component 1 and 2');
16 % LDA components 1,3 and 2,3
17 figure();
18 subplot(1,2,1);
19 scatter(z(1,I==1),z(3,I==1),'b','k.','SizeData',200);hold on;
20 scatter(z(1,I==2),z(3,I==2),'r','k.','SizeData',200);
21 scatter(z(1,I==3),z(3,I==3),'k','k.','SizeData',200);
22 scatter(z(1,I==4),z(3,I==4),'m','k.','SizeData',200);
23 xlabel('LDA component 1');
24 ylabel('LDA component 3');
25 title('LDA component 1 and 3');
26 subplot(1,2,2);
27 scatter(z(2,I==1),z(3,I==1),'b','k.','SizeData',200);hold on;
28 scatter(z(2,I==2),z(3,I==2),'r','k.','SizeData',200);
29 scatter(z(2,I==3),z(3,I==3),'k','k.','SizeData',200);
30 scatter(z(2,I==4),z(3,I==4),'m','k.','SizeData',200);
31 xlabel('LDA component 2');
32 ylabel('LDA component 3');
33 title('LDA component 2 and 3');

```