

Homework Set 3, CPSC 8420, Spring 2022

Last Name, First Name

Due 03/31/2022, Thursday, 11:59PM EST

Problem 1

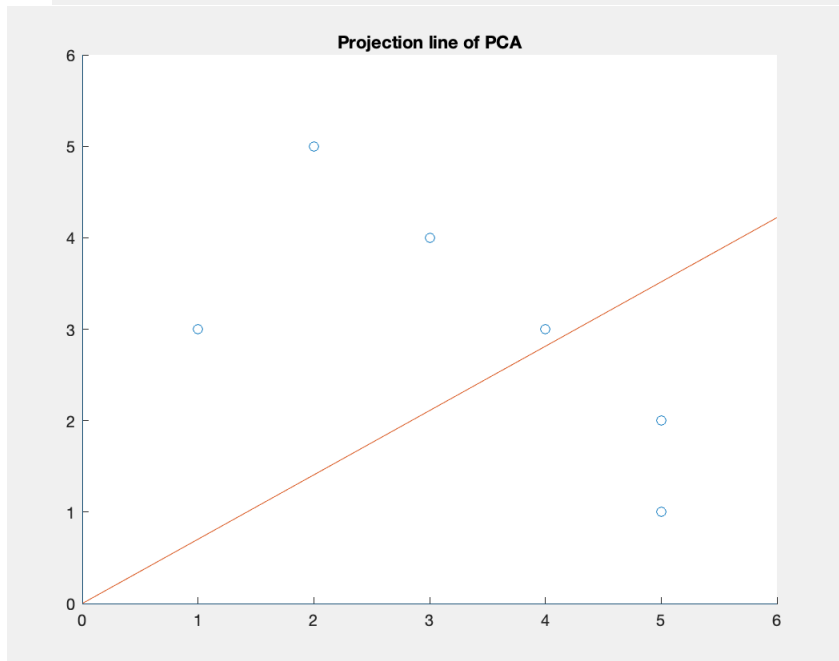
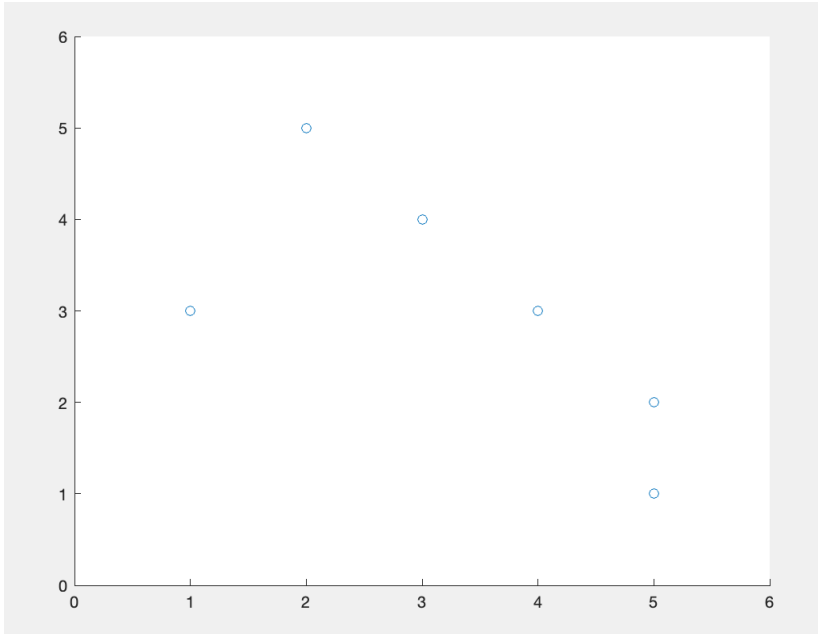
Given data-points $\{\{1, 3\}, \{2, 5\}, \{3, 4\}, \{4, 3\}, \{5, 2\}, \{5, 1\}\}$.

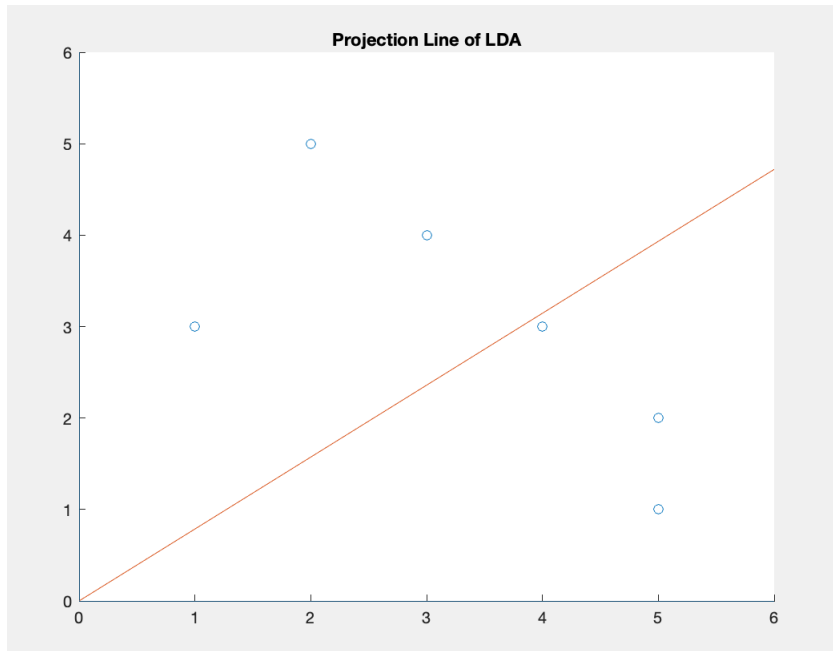
1. Please scatter-plot each data point within one figure (you can use Matlab, Python or any other programming language).
2. Now if we want to reduce the dimension from 2 to 1 by PCA, please determine the projection line which crosses the origin (please plot the line based on the scatter-plot figure above).

For the PCA projection line. I calculated the Covariance matrix and then performed an SVD on that matrix. I took the slope between the two points of the first column of U and got a slope of 0.7034.

3. Assume the first 4 data points belong to one class, while the rest 2 belong to the other. Now if we want to reduce the dimension from 2 to 1 by LDA, please determine the projection line which crosses the origin (you are expected to plot the line based on the scatter-plot figure).

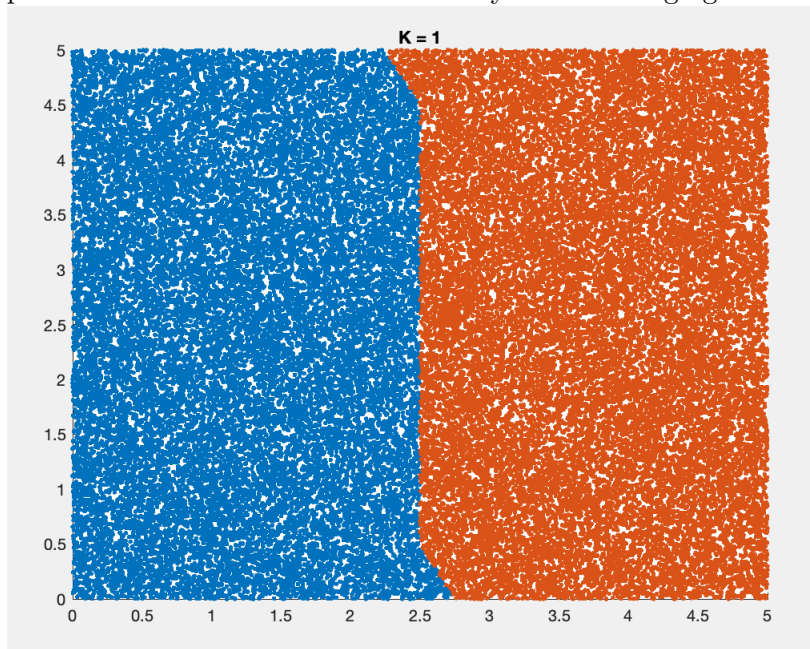
For the LDA projection line I calculated the weights (W) of the data, multiplied that with the data and then calculated the slope using the first and last data point. Which I got to be 0.78

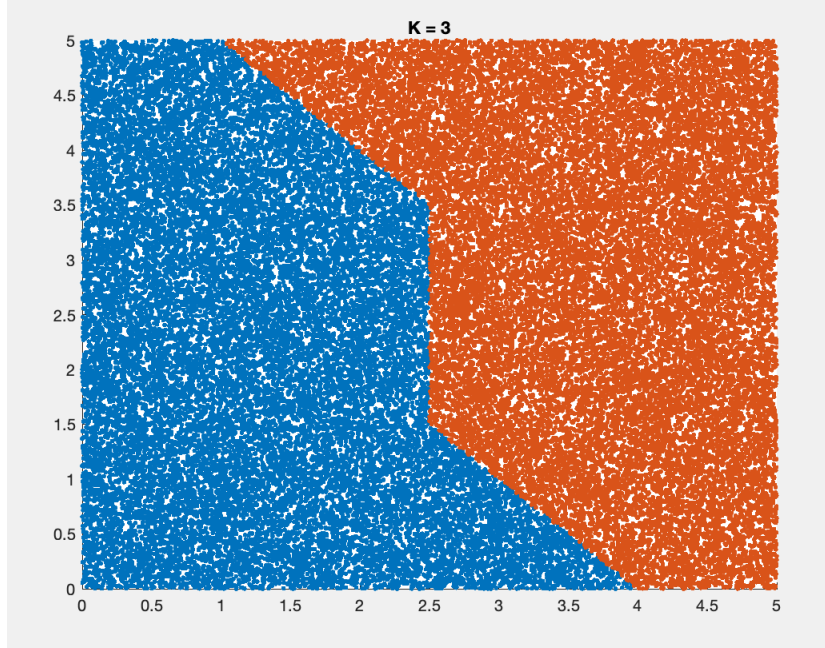




Problem 2

Given positive data-set $\{\{1, 1\}, \{2, 2\}, \{2, 3\}\}$, as well as negative data-set $\{\{3, 2\}, \{3, 3\}, \{4, 4\}\}$, please determine the decision boundary when leveraging k -NN where $k = 1$ and $k = 3$ respectively.





Problem 3

Given X, Y, Z , now please follow the idea/method used in LDA/PCA to find the best solution to:

$$\begin{aligned} & \underbrace{\arg \max}_{a,b} a^T Z b \\ & s.t. \quad a^T X a = 1, \quad b^T Y b = 1 \end{aligned} \quad (1)$$

Since X, Y, Z are SPD we can expand our constraints $s.t. \quad a^T X^{1/2} X^{1/2} a = 1, \quad b^T Y^{1/2} Y^{1/2} b = 1$

Let $u = X^{1/2} a$ and $v = Y^{1/2} b$, therefore we have that $u^T u = 1$ and $v^T v = 1$

We can arrange to say that $a^T = u^T X^{1/2}$ and $b = Y^{1/2} v$ and therefore we $\max u^T X^{-1/2} Z Y^{-1/2} v$

If we use SVD we find that $[U, \Sigma, V] = \text{svd}(X^{-1/2} Z Y^{-1/2})$. We can then let u be the first column of U , and v be the first column of V , and then find a and b .

Code

NOTE: For plotting data points on graph I just removed some aspects of the PCA projection code.

PCA Projection Line Code

```

1  clc; clear;
2  a = [[1, 3]; [2, 5]; [3, 4]; [4, 3]; [5, 2]; [5, 1]];
3

```

```

4 M = mean(a);
5 lower_bound = 0;
6 upper_bound = 6;
7
8 a_new = a - M;
9 C = cov(a)
10 [U,S,V] = svd(C);
11 points = a * U(:, 1)
12 slope_of_projection_line = (U(2,1) - U(1,1))/2;
13
14 x = [lower_bound, upper_bound];
15 y = slope_of_projection_line * x;
16
17
18 h = figure;
19 scatter(a(:, 1), a(:, 2));
20 hold on
21 plot(x, y)
22 xlim([lower_bound, upper_bound])
23 ylim([lower_bound, upper_bound])
24 line([0 0], ylim);
25 line(xlim, [0 0]);
26 title("Projection line of PCA")
27 hold off
28 waitfor(h)

```

LDA Projection Line

```

1 clc; clear;
2 x_lower_lim = 0;
3 x_upper_lim = 6;
4 a = [[1,3];[2,5];[3,4];[4,3];[5,2];[5,1]];
5 %classes for data
6 M = mean(a);
7 %center data on axis
8 a_new = a - M;
9 a_val = [4;2];
10
11 [W, Y] = LDA(a_new, a_val);
12
13 slope = (Y(6) - Y(1))/5;
14 x = [x_lower_lim, x_upper_lim];
15 y = slope * x;
16
17
18 h = figure;

```

```

19 scatter(a(:, 1), a(:, 2));
20 hold on
21 %plot([-3, -2, -1, 1, 2, 3], Y)
22 plot(x, y);
23 xlim([x_lower_lim, x_upper_lim])
24 ylim([x_lower_lim, x_upper_lim])
25 line([0 0], ylim);
26 line(xlim, [0 0]);
27 title("Projection Line of LDA")
28 hold off
29 waitfor(h)

```

LDA Algorithm

```

1 function [W,y] = LDA(Data, a)
2 % This code is written by Alaa Tharwat (Frankfurt University of Applied
   Science – Germany)
3 % For more details about the code of the numerical example(s) see our
   paper "Tharwat, A., Gaber, T., Ibrahim, A.,
4 % & Hassanien, A. E. Linear discriminant analysis: A detailed tutorial.
   AI Communications,
5 % (Preprint), 1–22.?
6 % engalaatharwat@hotmail.com
7 %% Examples
8 % c1=[1 2;2 3;3 3;4 5;5 5] % the first class (5 observations)
9 % c2=[4 2;5 0;5 2;3 2;5 3;6 3] % the second class (6 observations)
10 % data=[c1;c2] % te whole data
11 % a=[5;6] % the number of samples in each class
12 % [W,y] = LDA_ClassIndependent(data, a)
13 % Calculate the total number of all classes
14 c=unique(a);
15 pos=1;
16 for i=1:length(c)
17     for j=1:a(i)
18         Labels(pos,1)=i;
19         pos=pos+1;
20     end
21 end
22 % Calculate Mean of each class
23 for i=1:length(c)
24     mu(i,1:size(Data,2))=mean(Data(Labels==i,:));
25 end
26 % Calculate the total mean of all classes
27 muTotal=zeros(size(mu(1,:)));
28 for i=1:length(c)
29     muTotal=muTotal+a(i)*mu(i,:);

```

```

30 end
31 muTotal=muTotal/(sum(a));
32 % Subtract the originla data from the mean
33 D=zeros(size(Data,1),size(Data,2));
34 for i=1:length(c)
35     D(Labels==i,:)=Data(Labels==i,:)-repmat(mu(i,:),a(i),1);
36 end
37 % Calculate the within class variance (SW)
38 SW=zeros(size(Data,2),size(Data,2));
39 for i=1:c
40     SW=SW+D(Labels==i,:)'*D(Labels==i,:);
41 end
42 % Calculate the Between-class variance (SB)
43 SB=zeros(size(Data,2),size(Data,2));
44 for i=1:length(c)
45     SB= SB+a(i)*(mu(i,:)-muTotal)'*(mu(i,:)-muTotal);
46 end
47 % Calculate J(W)
48 J=inv(SW)*SB;
49 % Calculate the eignevalues and eigenvectors of (J)
50 [evec,eval]=eig(J);
51 % Sort the eigenvectors according to their corresponding eigenvalues (
    descending order)
52 eval = diag(eval);
53 [junk, index] = sort(-eval);
54 eval = eval(index);
55 evec = evec(:, index);
56 % Slect the most largest c eigenvectors as a lower dimensional space
57 W=evec(:,1:length(c)-1);
58 % project the original data on the lower dimensional space (W)
59 y=Data*W;

```

KNN Boundary Code

```

1  clc;
2  rng(1)
3  positive = [[1,1];[2,2];[2,3]];
4  negative = [[3,2];[3,3];[4,4]];
5  data = [positive; negative];
6  value = [1;1;1;2;2;2];
7  x = rand(50000, 2) * 5;
8
9
10 k_list = [1,3];
11
12 for i = 1:2

```

```

13     k = k_list(i);
14     [pred, idx, acc] = KNN(k, data, value, x);
15
16     class_1 = [];
17     class_2 = [];
18     for i=1:length(pred)
19         if pred(i) == 1
20             class_1 = [class_1; x(i, 1:2)];
21         else
22             class_2 = [class_2; x(i, 1:2)];
23         end
24     end
25
26     h= figure;
27     %scatter(data(:, 1), data(:, 2),500, '.');
28     hold on
29     scatter(class_1(:, 1), class_1(:, 2), 50, '.');
30     scatter(class_2(:, 1), class_2(:, 2), 50, '.');
31     xlim([0, 5])
32     ylim([0, 5])
33     t = "K = " + string(k);
34     title(t)
35     hold off
36     waitfor(h)
37 end

```

KNN Algorithm

```

1 function [predicted_labels, nn_index, accuracy] = KNN(k, data, labels,
    t_data, t_labels)
2 %KNN: classifying using k-nearest neighbors algorithm. The nearest
    neighbors
3 %search method is euclidean distance
4 %Usage:
5 %     [predicted_labels, nn_index, accuracy] = KNN(3, training,
    training_labels, testing, testing_labels)
6 %     predicted_labels = KNN(3, training, training_labels, testing)
7 %Input:
8 %     - k: number of nearest neighbors
9 %     - data: (NxD) training data; N is the number of samples and D
    is the
10 %     dimensionality of each data point
11 %     - labels: training labels
12 %     - t_data: (MxD) testing data; M is the number of data points
    and D
13 %     is the dimensionality of each data point

```



```

14 %      – t_labels: testing labels (default = [])
15 %Output:
16 %      – predicted_labels: the predicted labels based on the k-NN
17 %      algorithm
18 %      – nn_index: the index of the nearest training data point for
      each training sample (Mx1).
19 %      – accuracy: if the testing labels are supported, the accuracy
      of
20 %      the classification is returned, otherwise it will be zero.
21 %Author: Mahmoud Afifi – York University
22 %checks
23 if nargin < 4
24     error('Too few input arguments.')
```

```

25 elseif nargin < 5
26     t_labels=[];
27     accuracy=0;
28 end
29 if size(data,2)~=size(t_data,2)
30     error('data should have the same dimensionality');
31 end
32 if mod(k,2)==0
33     error('to reduce the chance of ties, please choose odd k');
34 end
35 %initialization
36 predicted_labels=zeros(size(t_data,1),1);
37 ed=zeros(size(t_data,1),size(data,1)); %ed: (MxN) euclidean distances
38 ind=zeros(size(t_data,1),size(data,1)); %corresponding indices (MxN)
39 k_nn=zeros(size(t_data,1),k); %k-nearest neighbors for testing sample (
      Mxk)
40 %calc euclidean distances between each testing data point and the
      training
41 %data samples
42 for test_point=1:size(t_data,1)
43     for train_point=1:size(data,1)
44         %calc and store sorted euclidean distances with corresponding
            indices
45         ed(test_point,train_point)=sqrt(...
46             sum((t_data(test_point,:)-data(train_point,:)).^2));
47     end
48     [ed(test_point,:),ind(test_point,:)]=sort(ed(test_point,:));
49 end
50 %find the nearest k for each data point of the testing data
51 k_nn=ind(:,1:k);
52 nn_index=k_nn(:,1);
53 %get the majority vote

```

```

54 for i=1:size(k_nn,1)
55     options=unique(labels(k_nn(i,:)'));
56     max_count=0;
57     max_label=0;
58     for j=1:length(options)
59         L=length(find(labels(k_nn(i,:)')==options(j)));
60         if L>max_count
61             max_label=options(j);
62             max_count=L;
63         end
64     end
65     predicted_labels(i)=max_label;
66 end
67 %calculate the classification accuracy
68 if isempty(t_labels)==0
69     accuracy=length(find(predicted_labels==t_labels))/size(t_data,1);
70 end

```