

MATH 285 HW2

Xiaoyan Chong

November 23, 2015

# 1 Problem 1

(a) Figure 1 shows the map of Chinese Cities using build-in function *cmdscale.m*.

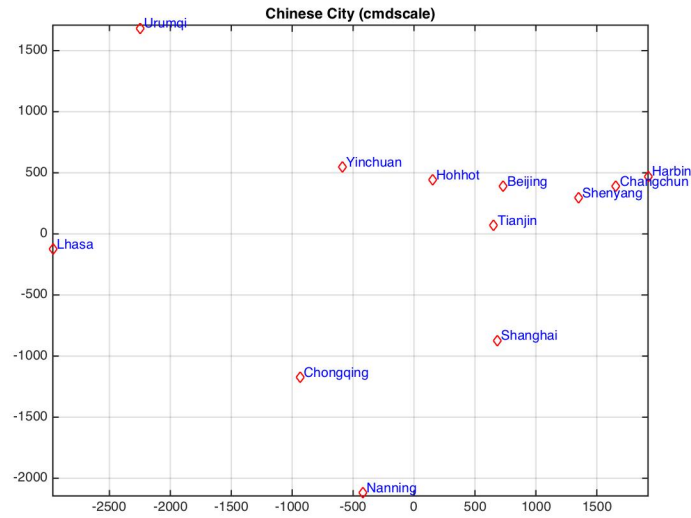


Figure 1: Chinese Cities (cmdscale)

Figure 2 shows the map of Chinese Cities using build-in function *xcmds.m*. By this function, I also get  $Stress = 0.0805$ . Since the stress is less than 0.1, we conclude the result is pretty good.

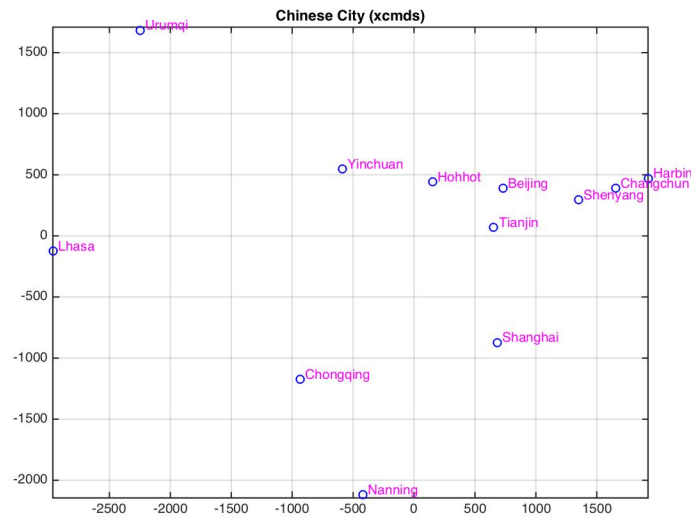


Figure 2: Chinese Cities (xcmds)

The following is the Matlab code:

```

6      %% load data
7      load ChineseCityData.mat
8
9      %% plot with cmdscale
10     [Y, e] = cmdscale(dists)
11
12     f1_1 = figure;
13     labels = Cities;
14     plot(-Y(:,1),-Y(:,2),'rd');
15     axis equal;
16     text(-Y(:,1)+30,-Y(:,2)+30,labels,'Color','b','HorizontalAlignment','left');
17     title('Chinese City (cmdscale)');
18     grid on;
19
20     saveas(f1_1,'/Users/XC/Dropbox/SJSU/Courses/Math 285 Data Modeling/HW/hw2/1_cmdscale.jpg','jpg')
21
22     %% plot with my own function
23
24     [Y, stress] = xcmds(dists,2)
25     f1_2 = figure;
26     labels = Cities;
27     plot(Y(:,1),Y(:,2),'bo')
28     axis equal;
29     text(Y(:,1)+40,Y(:,2)+30,labels,'Color','m','HorizontalAlignment','left');
30     title('Chinese City (xcmds)');
31     grid on;
32
33     saveas(f1_2,'/Users/XC/Dropbox/SJSU/Courses/Math 285 Data Modeling/HW/hw2/1_xcmds.jpg','jpg')
34

```

Figure 3: Problem 1 Code

```

1      function [Y, stress] = xcmds(X,k)
2      % this is a function for mds
3      % we can compare this with cmdscale in Matlab
4      [D,N] = size(X);
5      Xsquare = X.* X;
6      unit1 = ones(D,1);
7      v = ones(D,1);
8      I = diag(v);
9      J = I- 1/D * unit1 * unit1';
10     B = (-1/2)*J * Xsquare * J;
11
12     % extract engivalues and eigenvectors of B
13     [V,D] = eig(B);
14     D2 = diag(sort(diag(D),'descend'));
15     e=diag(D2);
16     [c,ind] = sort(diag(D),'descend');
17     V2 = V(:,ind)
18     eigvaluek = D2(1:k,1:k);eigvectork = V2(:, 1:k);
19     eigvaluek
20
21     % Build Y
22     Y = eigvectork * sqrt(eigvaluek);
23
24     % stress
25     d = L2_distance(Y',Y')
26     stress = sqrt(sum(sum((X-d).^2))/sum(sum(X.^2)))
27     end

```

Figure 4: xcmds.m

- (b) No, we could not use these distances to construct a world map. In this case, when we project the distance from spherical onto 2 dimensional space, we could not preserve the distance between two cities. If on spherical, the distance between two cities are pretty close, when projecting to 2-D plane, we cannot preserve the closeness, they may be far away from each other.

## 2 Problem 2

- (1) Data set:Glassdata
- (2) These data belong to a Glass Identification Database and were downloaded from the UCI Machine Learning Repository [Newman, et al., 1998].
- (3) This data has 214 observations, with 9 variables: refractive index, sodium, magnesium, aluminum, silicon, potassium, calcium, barium, and iron.
- (4) The data set also includes a class label: building windows float processed(type 1), building windows not float processed(type 2), vehicle windows float processed(type 3), building windows not float processed (none in this dataset)(type 4), containers(type 5), tableware(type 6), and headlamps(type 7).
- (5) One could use this data set to develop a classifier that could be used in forensic science applications.
- (6) By the description on (4), we may see type 1,2,3 are very similar, and type 6, 7 are very similar. According to Figure 7 (Isomap result in 2-D), we could see that 1,2 , 3 are grouped together, and 6, 7 are in a different side of (type 1,2, and 3). Since different types are classified based on the concentrations of variables in introduced in (3), we could say there is a pattern, because similar types of glass contains similar concentrations of Na, Mg, Al, Si, K, Ca, Ba, and Fe. And in our case, these similar types are very close to each other in Isomap result.
- (7) Although this data set is not as good as figure data set, we could still see there is a pattern. But if we use a figure dataset, it is much easier to detect a pattern.

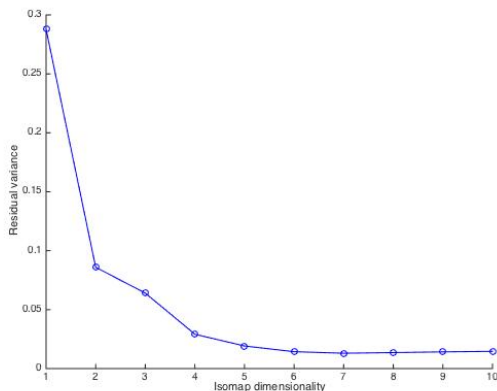


Figure 5: Scree Plot

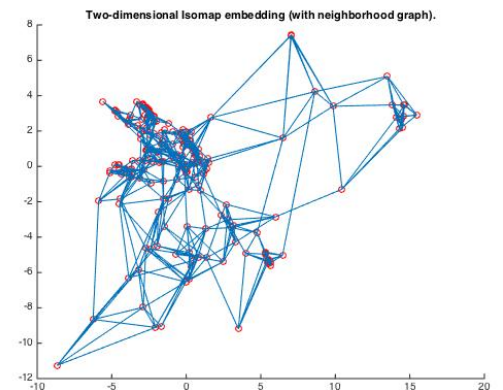


Figure 6: 2-D Isomap embedding

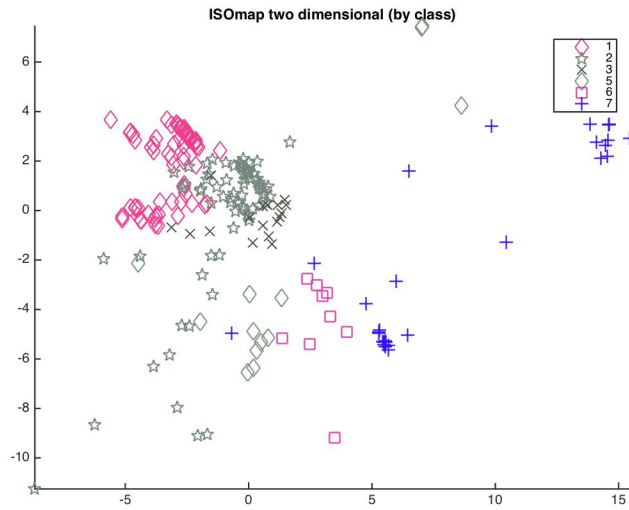


Figure 7: ISomap

The code is as follows.

```

1 %% Problem 2
2 clear
3 close all
4 clc
5 %% load glass dataset
6 glassdata = importdata('glass.txt',' ', 0);
7
8 %%
9 % Then get the interpoint dissimilarity matrix.
10 % We will use standardized Euclidean distance.
11 tmp = pdist(glassdata,'seuclidean');
12
13 % Now put it into a square matrix.
14 D = squareform(tmp);
15
16 % Now we do ISOMAP.
17 % We will define the neighborhood using the number of nearest neighbors, k = 5.
18 [Y,R,E] = Isomap(D,'k',5);
19
20 lables = zeros(size(glassdata,1),1);
21 lables(glassdata(:,11)==1) = '1';
22 lables(glassdata(:,11)==2) = '2';
23 lables(glassdata(:,11)==3) = '3';
24 lables(glassdata(:,11)==5) = '5';
25 lables(glassdata(:,11)==6) = '6';
26 lables(glassdata(:,11)==7) = '7';
27
28 f2_1 = figure; gcpplot(Y.coords{2}',lables); axis equal
29 legend('1','2','3','5','6','7')
30 title('ISomap two dimensional (by class)')
31 saveas(f2_1,'/Users/XC/Dropbox/SJSU/Courses/Math 285 Data Modeling/HW/hw2/2_isomap.jpg','jpg')

```

Figure 8: Code part 1

```
33 %% compare with mds
34
35 - Y_mds = mds(D, 3);
36 - lables = zeros(size(glassdata,1),1);
37 - lables(glassdata(:,11)==1) = '1';
38 - lables(glassdata(:,11)==2) = '2';
39 - lables(glassdata(:,11)==3) = '3';
40 - lables(glassdata(:,11)==5) = '5';
41 - lables(glassdata(:,11)==6) = '6';
42 - lables(glassdata(:,11)==7) = '7';
43 - f2_2 = figure; gcplot(Y_mds,lables);
44 - legend('1','2','3','5','6','7')
45 - title 'mds two dimensional (by class)'
46
47 - saveas(f2_2, '/Users/XC/Dropbox/SJSU/Courses/Math 285 Data Modeling/HW/hw2/2_mds.jpg', 'jpg')
48
```

Figure 9: Code part 2

### 3 Problem 2: Digit Dataset

- (1) Data: I use the dataset from homework one, the digit dataset (mnist-digit1.mat). I picked the first 1000 data points, and did the following analysis based on these 1000 points. Figure 10 is the scatter plot of original data.
- (2) According to Figure 11, we conclude there is a pattern in this figure. For digits on the upper left, they are thinner, and pretty much vertical (only a few have a direction: positive slope). For the digits on the upper right, they are also thinner, but they have a direction (negative slope). For digits on the bottom left, they are thicker, and they are vertical (only several have a direction: positive slope). For digits on the bottom right, they are thicker, and they also have a direction (negative slope).

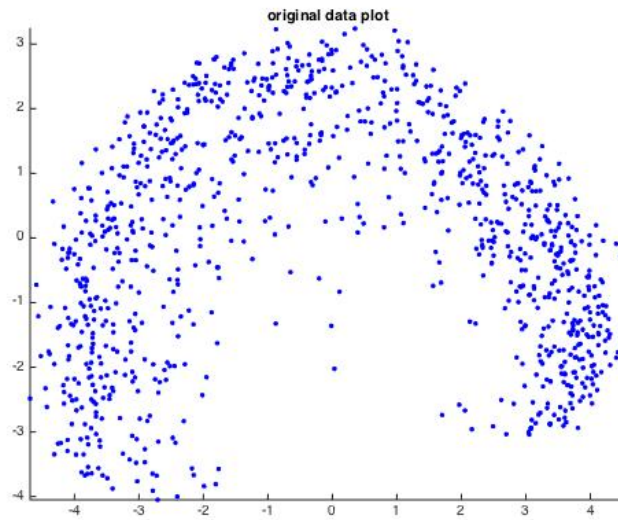


Figure 10: Original Data Scatter

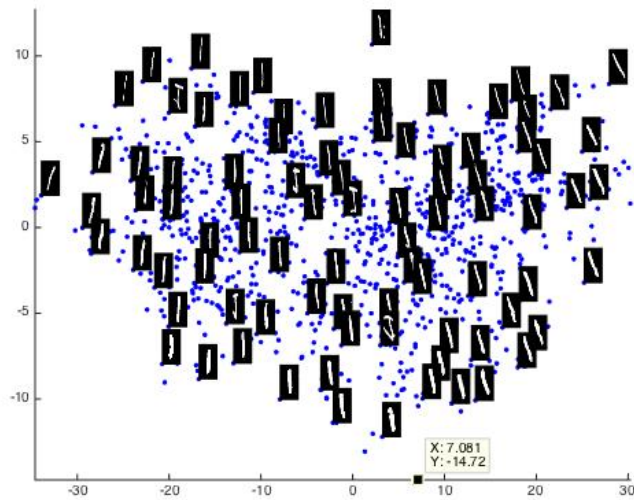


Figure 11: ISomap

(3) Code, as shown in Figure 12

```

50
51 %% Problem 2 Digit Dataset
52 clear
53 close all
54 clc
55
56 %% Load data
57 load mnist_digit1.mat
58
59 %% pick 1000 observations
60 X_test = X(1:1000,:);
61
62 % original data plot
63 figure; gcpplot(X_test)
64 title 'original data plot'
65
66 %% isomap
67 D = L2_distance(X_test',X_test',1);
68 [Y,R,E] = Isomap(D,'k',5);
69 figure; gcpplot(Y.coords{2}');
70
71 %% add image to the figure
72 cursor_info; a = cursor_info.Position; hold on; ...
73 imagesc([a(1) a(1)+2], [a(2) a(2)+2], reshape(X_test(cursor_info.DataIndex,:), 28,28)); ...
74 figure(gcf); colormap gray
75

```

Figure 12: Code part



## 4 Problem 3

(1)  $O \rightarrow A: 2$

(2)  $O \rightarrow B: 4$

(3)  $O \rightarrow C: 4$

(4)  $O \rightarrow D: 8$

(5)  $O \rightarrow E: 7$

(6)  $O \rightarrow F: 14$

(7)  $O \rightarrow T: 13$

## 5 Problem 4

Figure 13 is the original data scatter plot.

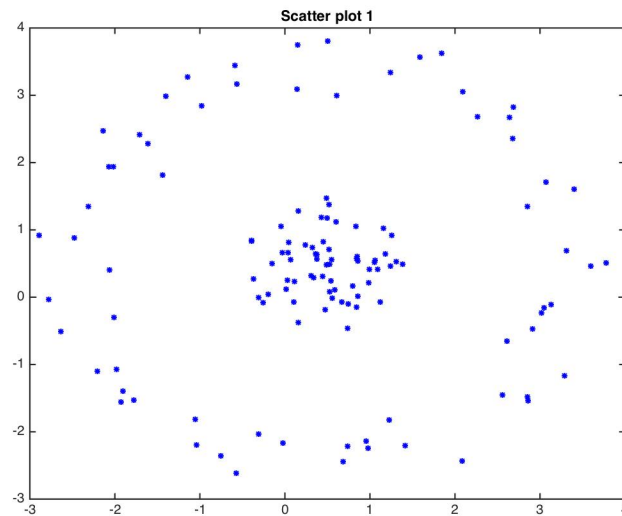


Figure 13: Scatter plot of original

Figure 14 is the plot for Gaussian Kernel PCA 2-D.

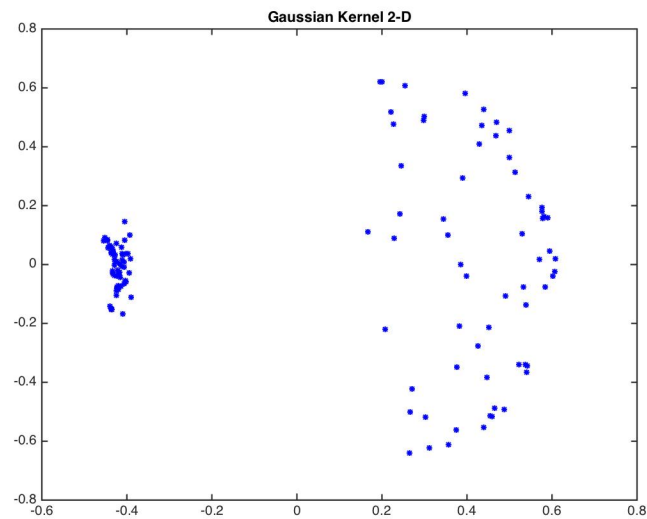


Figure 14: Gaussian Kernel

Comparing the above two figures, we could find that by using Gaussian Kernel PCA, the two groups could be separated very clearly, and we can find a linear boundary between these two groups.

The code is as follows.

```
Problem_4.m x +
1 - clear
2 - close all
3 - clc
4
5 - %% Load data
6 - load kernelpca_data.mat
7
8 - %% Gaussian Kernel
9
10 - [D,N] = size(X);
11
12 - % calculating distance
13 - D1 = L2_distance(X',X')
14 - D2 = D1.^2;
15
16 - % calculating sigma
17 - [idx,di] = knnsearch(X,X,'k',9);
18 - eighth_idx = idx(:,9);
19 - eighth_nb = di(:,9);
20 - sigma = mean(eighth_nb);
21
22 - % compute kernel matrix
23 - K = exp(-D2./(2*sigma^2));
24 - J = ones(D,D)/D;
25
26 - % Centering kernel matrix (non-linearly mapped data).
27 - Kc = K - J*K - K*J + J*K*J;
28
```

Figure 15: Code part 1

```
28
29 - %% eigendecomposition
30 - [V,S] = eig(Kc);
31 - [dsort, idum]=sort(diag(S),'descend');
32 - l=abs(dsort);
33 - V=V(:,idum);
34 - l2 = l(1:2);
35
36 - % Two dimensional representation of the data obtained by Kernel PCA
37 - X_proj = V(:,1:2) * diag(sqrt(l2));
38
39 - f4_1 = figure;plot(X_proj(:,1),X_proj(:,2),'b*','MarkerSize',4);
40 - title 'Gaussian Kernel 2-D'
41 - saveas(f4_1,'/Users/XC/Dropbox/SJSU/Courses/Math 285 Data Modeling/HW/hw2/4_GKernel.jpg','jpg')
42
43 - %% compare original plot and kernel PCA plot
44
45 - % original figure
46 - f4_2=figure; plot(X(:,1), X(:,2),'b*','MarkerSize',4)
47 - title 'Scatter plot 1'
48 - saveas(f4_2,'/Users/XC/Dropbox/SJSU/Courses/Math 285 Data Modeling/HW/hw2/4_ORG.jpg','jpg')
49
```

Figure 16: Code part 2

## 6 Problem 5

- (a) (i) Figure 17 shows three true clusters, and Figure 18 shows the results using k-means.
- (ii) The error percentage of my clustering is  $1 - 0.8933 = 0.1067 = 10.67\%$ . I think the error percentage is a little high. So the result is not very good.
- (iii) From the following two figures, we can also tell that the boundary of the left two groups are not clear (figure 17), which leads to many misclassification in that small area (figure 18)
- (iv) By wikipedia, we find that "A key limitation of k-means is its cluster model. The concept is based on spherical clusters that are separable in a way so that the mean value converges towards the cluster center. " In this case, the shape of each group in our data is ellipse, so this method is not a good one for Iris data.

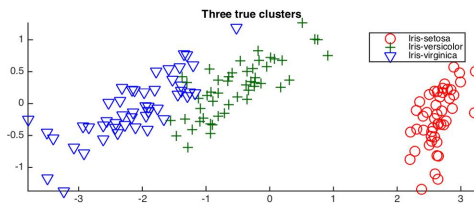


Figure 17: Three true clusters

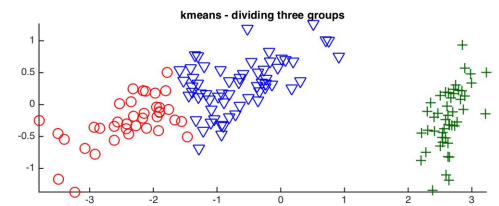


Figure 18: kmeans - dividing three groups

- (b) Figure 19 shows the scree plot. From this figure, we could find that  $k=2$  or  $k=3$  are both fine. If we choose  $k=2$ , we could find that the two classes ( now Virginica and Versicolor are grouped into one class, setosa is another class) are seperatd clearly. In this case, the misclassification error should be very small.

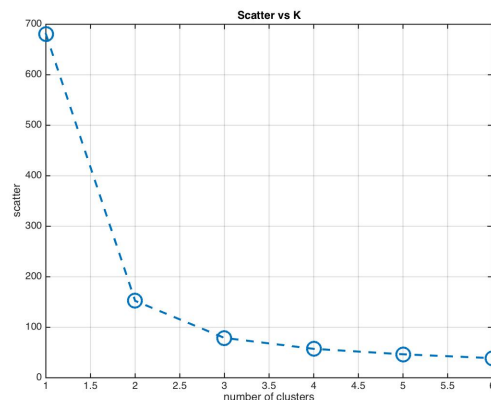


Figure 19: scatter vs k.

The code is as follows.

```

1 %% Problem 5
2 clear
3 close all
4 clc
5 %% read data
6 fileID = fopen('iris.txt');
7 C = textscan(fileID, '%3.1f, %3.1f, %3.1f, %3.1f, %s');
8 fclose(fileID);
9
10 %% form data matrix and plot data with true labels
11 X = [C{1:4}]; % concatenate the first four cells to form a matrix
12
13 % true labels
14 labels = zeros(size(X,1),1);
15 labels(strcmp(C{5}, 'Iris-setosa')) = 1;
16 labels(strcmp(C{5}, 'Iris-versicolor')) = 2;
17 labels(strcmp(C{5}, 'Iris-virginica')) = 3;
18
19 % display the three true clusters
20 f5_1 = figure; gcpplot(X, labels); axis equal
21 legend('Iris-setosa', 'Iris-versicolor', 'Iris-virginica')
22 title 'Three true clusters'
23 saveas(f5_1, '/Users/XC/Dropbox/SJSU/Courses/Math 285 Data Modeling/HW/hw2/5_true.jpg', 'jpg')
24
25 %% perform kmeans
26 labels_kmeans = kmeans(X, 3, 'Replicates', 10);
27 f5_2=figure; gcpplot(X, labels_kmeans); axis equal
28 title 'kmeans - dividing three groups'
29 saveas(f5_2, '/Users/XC/Dropbox/SJSU/Courses/Math 285 Data Modeling/HW/hw2/5_kmeans3.jpg', 'jpg')
30

```

Figure 20: Code part 1

```

30
31 %% part a: compute clustering accuracy rate
32 trueLabels = labels;
33 accuracy = 1 - computing_percentage_of_misclassified_points(labels_kmeans, trueLabels)
34
35 %% part b: select number of clusters when not given
36 k=3
37 scatter = zeros(1,2*k); % vector of total scatter
38 scatter(1) = sum(sum((X - repmat(mean(X,1), size(X,1), 1)).^2, 2)); % j = 1
39 for j = 2 : 2*k
40     [~, ~, intravars] = kmeans(X, j, 'Replicates', 10);
41     scatter(j) = sum(intravars);
42 end
43
44 f5_3=figure;
45 plot(scatter, 'o--', 'markersize', 14, 'linewidth', 2)
46 xlabel 'number of clusters'
47 ylabel 'scatter'
48 title 'Scatter vs K'
49 grid on
50 saveas(f5_3, '/Users/XC/Dropbox/SJSU/Courses/Math 285 Data Modeling/HW/hw2/scatter_k.jpg', 'jpg')
51

```

Figure 21: Code part 2