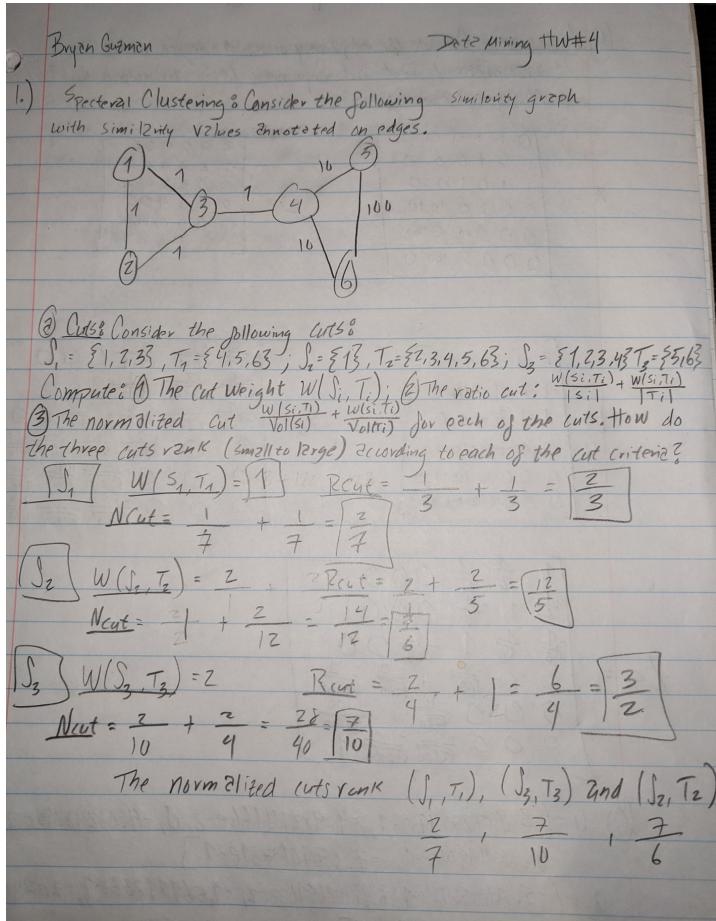


Bryan Guzman
 ICSI-435: Data Mining
 Homework #4
 001265918

1.)



(b) Lepelians: Compute the adjacency matrix A , degree matrix D , Laplacian matrix $L = D - A$ and Symmetric Lepelian matrix $L_S = D^{-1/2} L D^{-1/2}$ of the graph above.

$$A = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 10 \\ 0 & 0 & 0 & 10 & 0 \end{bmatrix} \quad D = \begin{bmatrix} 2 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 4 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

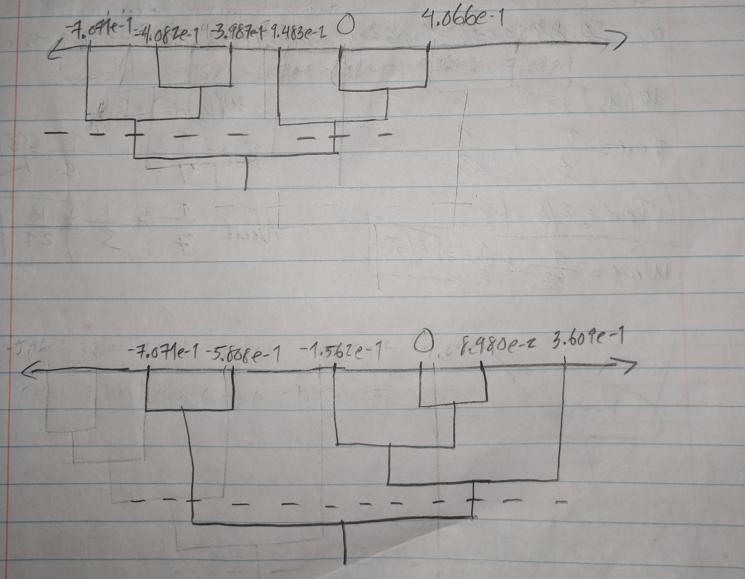
$$L = \begin{bmatrix} 2 & -1 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 \\ -1 & -1 & 3 & -1 & 0 \\ 0 & 0 & -1 & 10 & -10 \\ 0 & 0 & -10 & -10 & 0 \end{bmatrix} \quad L_S = \begin{bmatrix} \frac{2}{\sqrt{2}} & \frac{-1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} & 0 & 0 \\ \frac{-1}{\sqrt{2}} & \frac{2}{\sqrt{2}} & \frac{-1}{\sqrt{2}} & 0 & 0 \\ \frac{-1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} & \frac{3}{\sqrt{3}} & \frac{-1}{\sqrt{3}} & 0 \\ 0 & 0 & \frac{-1}{\sqrt{10}} & \frac{2}{\sqrt{10}} & \frac{-10}{\sqrt{10}} \\ 0 & 0 & 0 & \frac{1}{\sqrt{10}} & \frac{-10}{\sqrt{10}} \\ 0 & 0 & 0 & \frac{-1}{\sqrt{10}} & \frac{-10}{\sqrt{10}} \end{bmatrix}$$

$$L_S = \begin{bmatrix} 1 & -\frac{1}{2} & -\frac{1}{16} & 0 & 0 \\ \frac{1}{2} & 1 & -\frac{1}{6} & 0 & 0 \\ -\frac{1}{8} & -\frac{1}{8} & 1 & \frac{1}{\sqrt{63}} & 0 \\ 0 & 0 & \frac{1}{\sqrt{63}} & 1 & -\frac{10}{\sqrt{2310}} \\ 0 & 0 & 0 & \frac{-10}{\sqrt{2310}} & 1 & -\frac{10}{\sqrt{2310}} \\ 0 & 0 & 0 & \frac{-10}{\sqrt{2310}} & \frac{-10}{\sqrt{2310}} & 1 \end{bmatrix}$$

$U = [-3.98727149e-1, -9.48314666e-2, 0, -4.8829889e-1, 9.06641104e-1, -7.0710f781e-1]$

$U_S = [-5.80805048e-1, 8.98076570e-2, -1.56151223e-1, -7.07106789e-1, 3.60857002e-1, 0]$

(a) Hierarchical Cluster from eigenvectors: Here we will do agglomerative single link clustering of the nodes represented by their U values and by their U_S values computed in the previous part. Note that we will get two dendograms, one for U_S and one for U , and in each the nodes will be one-dimensional points. Draw the two dendograms.



(b) Hierarchical Clusters from eigenvectors: Cut the dendograms above to obtain two clusters. Show the resulting partitions in your solution file. What are their cut weights, normalized cuts and ratio cuts? Discuss why we obtain these partitions by relating the cut measures to the specific cut measures.

$$U = \{0, 4.066e-1, 9.483e-2\} \quad S = \{-4.682e-1, -3.987e-1, -7.071e-1\}$$

$$U_S = \{0, 8.980e-2, -1.562e-1, 3.609e-1\} \quad T = \{-7.071e-1, -5.608e-1\}$$

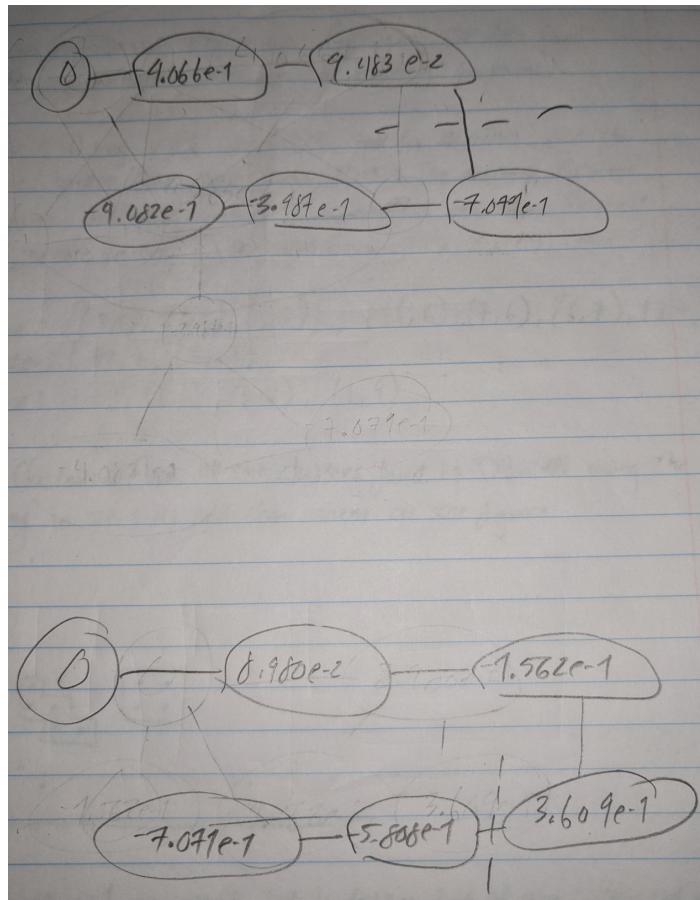
$$w(U, S) = 1 \quad w(U_S, T) = 1$$

$$R_{cut} = \frac{7}{3} + \frac{1}{23} = \frac{7}{3} + \frac{1}{23} = \frac{2}{2} = \frac{2}{2} = \boxed{\frac{2}{2}}$$

$$R_{cut} = \frac{1}{4} + \frac{1}{2} = \boxed{\frac{3}{4}}$$

$$N_{cut} = \frac{1}{7} + \frac{1}{3} = \boxed{\frac{10}{21}}$$

$$W_{cut} = \frac{1}{5} + \frac{1}{5} = \boxed{\frac{2}{5}}$$



2.)

(a) $\frac{0+1+2+2+10}{5} = 3 \text{ M} \quad \text{median} = 2$

$$\sum_i (x_i - \mu)^2 = (0-3)^2 + (1-3)^2 + (2-3)^2 + (2-3)^2 + (10-3)^2 = 64$$

$$\sum_i (x_i - m)^2 = (0-2)^2 + (1-2)^2 + (2-2)^2 + (2-2)^2 + (10-2)^2 = 69$$

The median is larger.

(b) Show that the sum of squared distances from the mean is smaller than that to the median in general, i.e.

Proof by cases: $\sum_i (x_i - \mu)^2 \leq \sum_i (x_i - m)^2$ Sum of Deviations

Case 1: Assume $\mu \geq m$.

$$\begin{aligned} \sum_i (x_i - \mu)^2 &\leq \sum_i (x_i - m)^2 \\ &= \sum_i x_i(x_i - \mu) \leq \sum_i (x_i - m)^2 \\ &= \sum_i x_i(x_i - \mu) - \mu \sum_i x_i \leq \sum_i (x_i - m)^2 \\ &= \sum_i x_i(x_i - \mu) - \mu(\sum_i x_i - \mu) \leq \sum_i (x_i - m)^2 \\ &= \sum_i x_i(x_i - \mu) \leq \sum_i (x_i - m)^2 \\ &= \sum_i x_i(x_i - \mu) + \mu \sum_i x_i \leq \sum_i (x_i - m)^2 \\ &= \sum_i (x_i - \mu)(x_i - \mu) \leq \sum_i (x_i - m)^2 \\ &= \sum_i (x_i^2 - \mu^2) \leq \sum_i (x_i - m)^2 \end{aligned}$$

3.)

3.) Density-based Clusterings: Considering the data points in the figure below, answer the following questions.

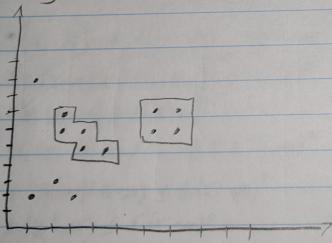
- (a) L_{∞} : Using $\epsilon=1$, $\text{minpts}=3$, find all core, border and noise points. Hint: The distance between points according to L_{∞} is defined as maximum absolute difference in any of the dimensions, i.e. $d_{\infty}(x_i, x_j) = \max_{i=1 \dots d} |x_{i,j} - x_{j,i}|$

Core: $\{(2,6), (3,6), (3,5)\}, \{(6,6), (7,6), (6,7), (7,7)\}$

Border: $\{(2,7), (4,5)\}$

Noise: $\{(1,2), (2,3), (3,2), (1,9)\}$

- (b) L_2 Clusters: Find all the clusters found by DBSCAN using the setting in part (a) and show them on the figure.



- (c) L_2 : Using $\epsilon=1$, $\text{minpts}=3$, find all core, border and noise points. Compare your results with part (a). What difference do you see? Why?

Core: $\{(2,6), (3,6), (3,5)\}, \{(6,6), (7,6), (6,7), (7,7)\}$

Border: $\{(2,7), (4,5)\}$

Noise: $\{(1,2), (2,3), (3,2), (1,9)\}$

I see no difference with any of the core, border and noise points. Nothing changed since when using Euclidean distance of 1, it introduce no new points into the clusters. The noise points and already present clusters had 2 distance greater than 1.

- (d) Find all the clusters found by DBSCAN using the setting in part (c) and show them on the figure. Compare your results with part (b). What difference do you see? Why?

* See question 3b for the figure.

There is no difference in figure 3d when compared to figure 3b. They are exactly the same due to what I discussed in question 3c.