

Machine Learning for Data Science (CS4786)

Lecture 14

Spectral Clustering

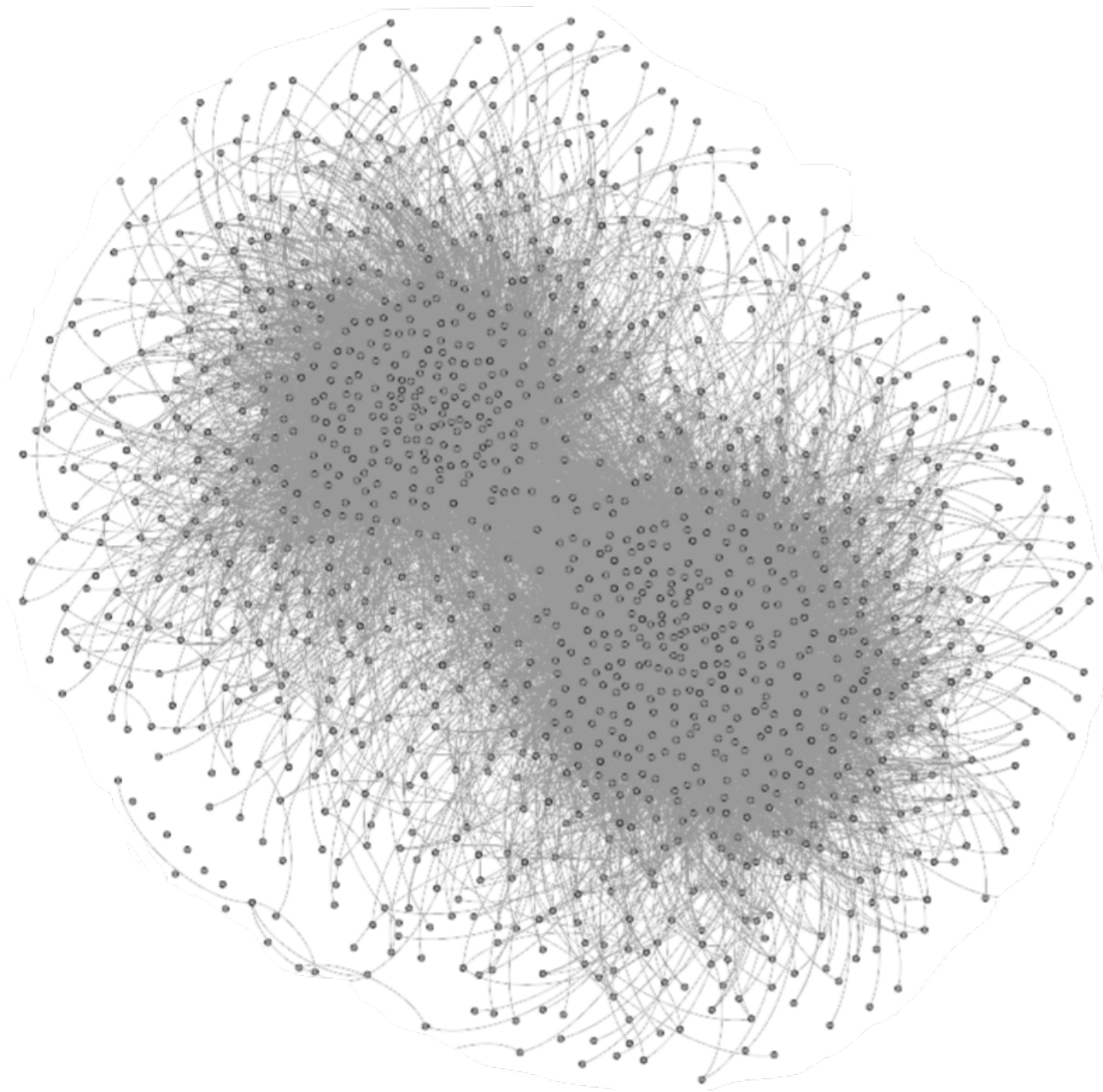
Course Webpage :

<http://www.cs.cornell.edu/Courses/cs4786/2017fa/>

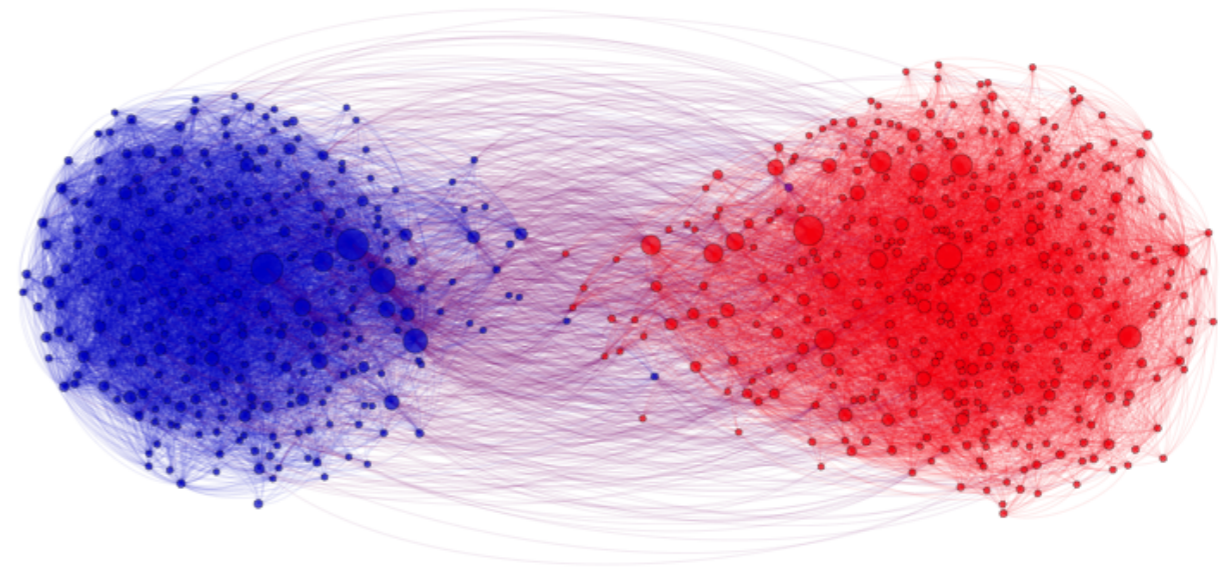
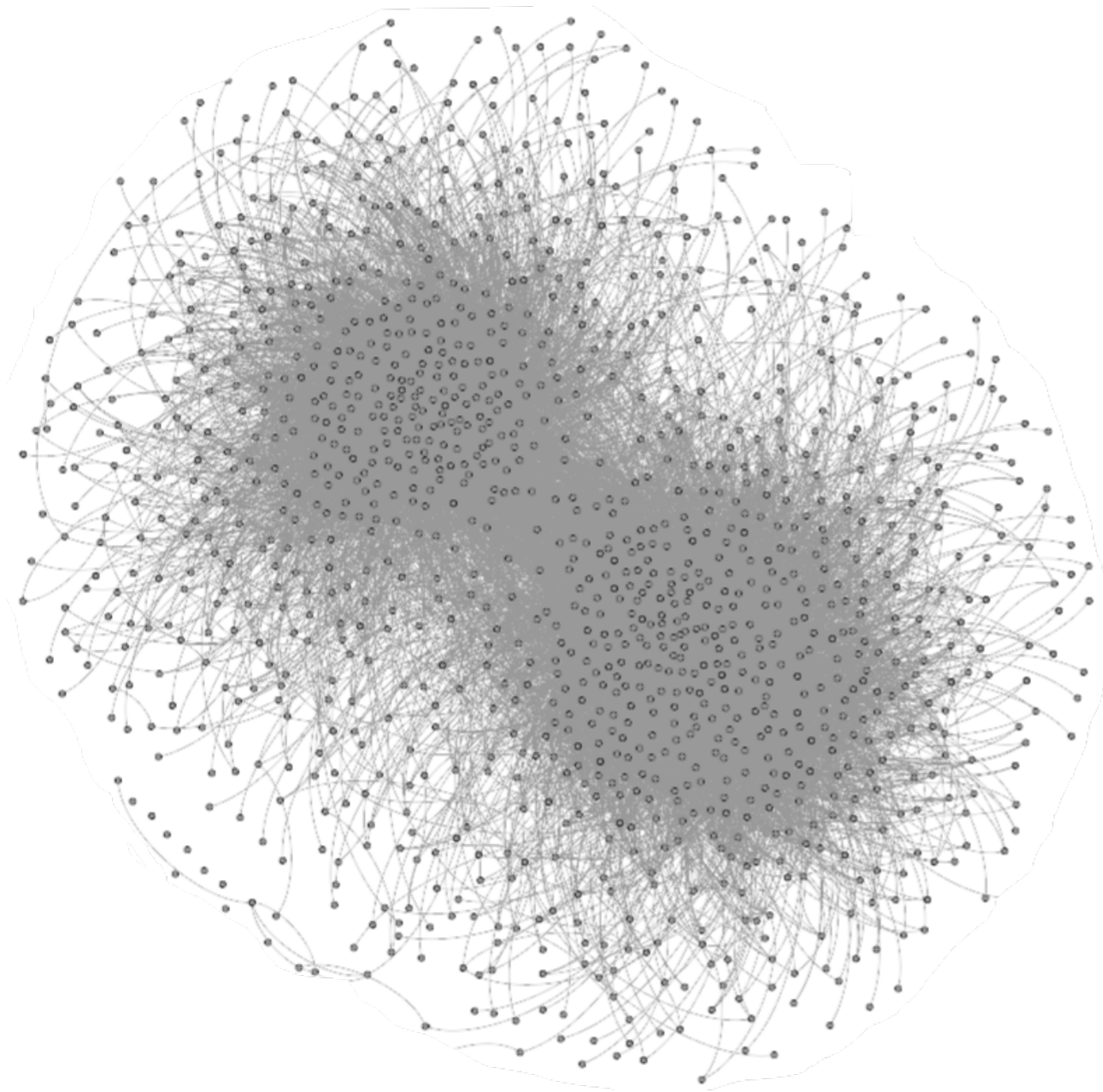
Announcement

- Competition I data is out:
- <https://confluence.cornell.edu/x/f3zHF>

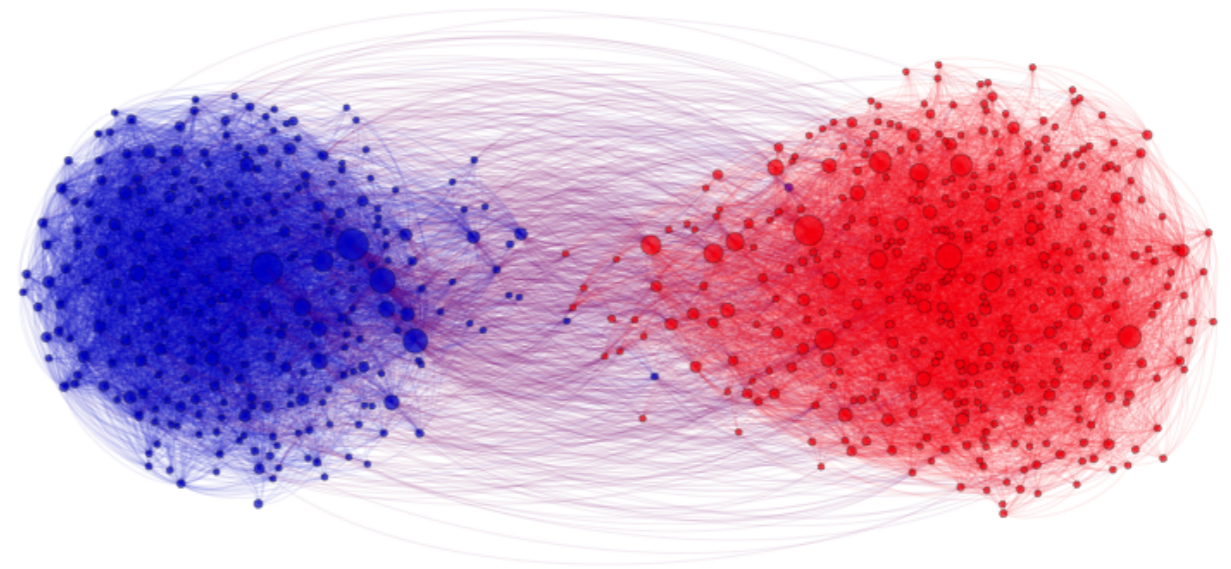
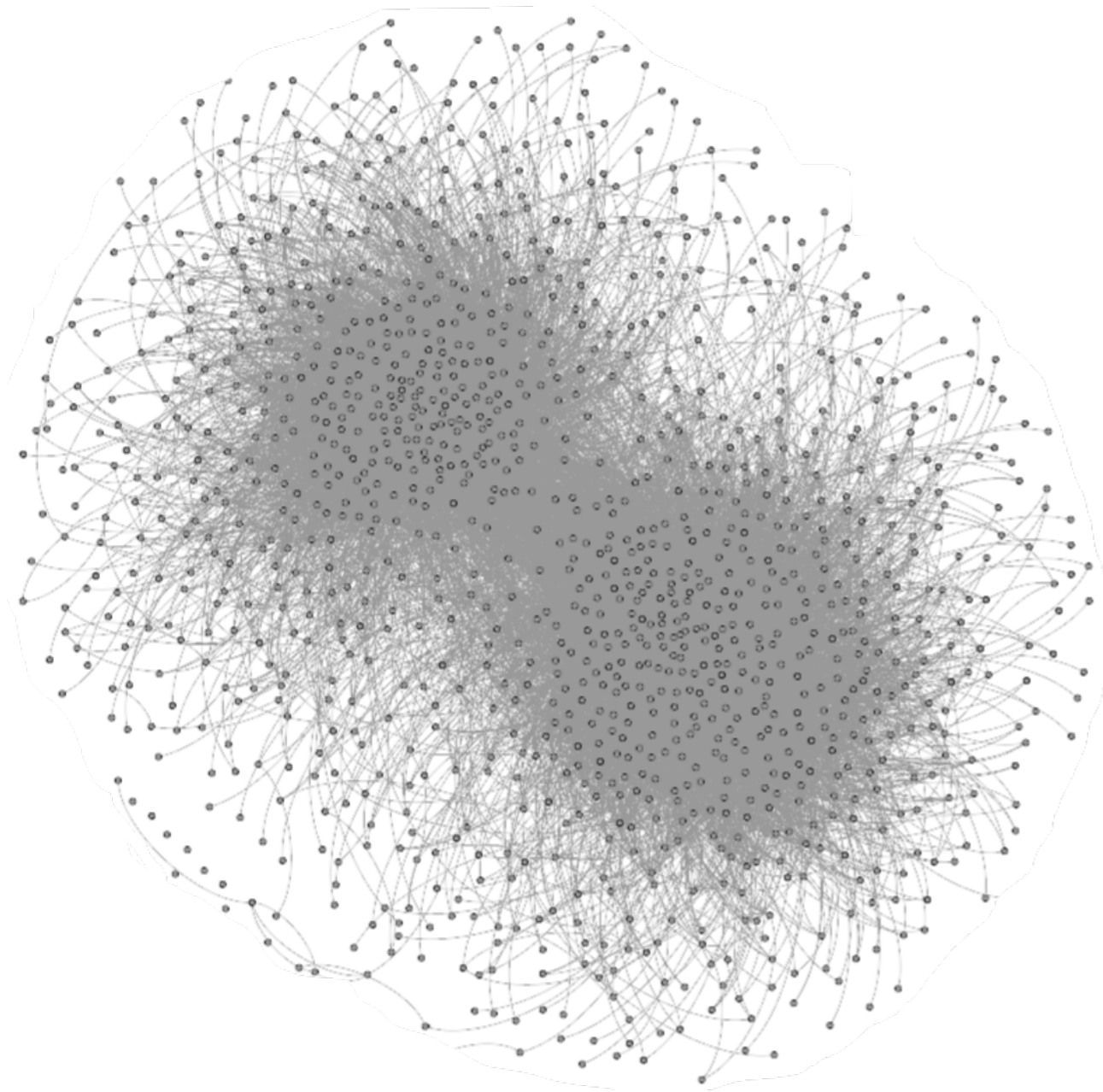
SPECTRAL CLUSTERING



SPECTRAL CLUSTERING



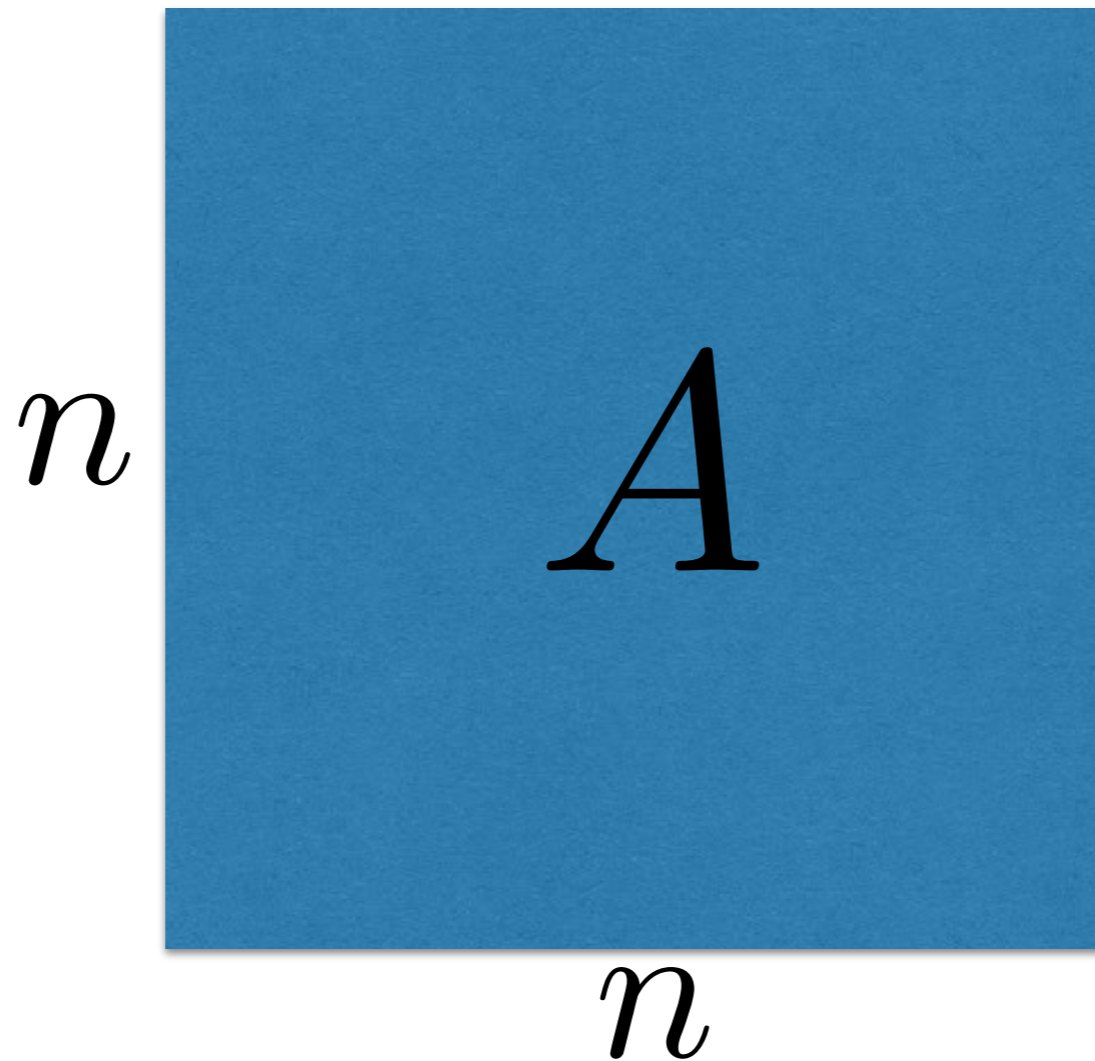
SPECTRAL CLUSTERING



- Cluster nodes in a graph.
- Analysis of social network data.

SPECTRAL CLUSTERING

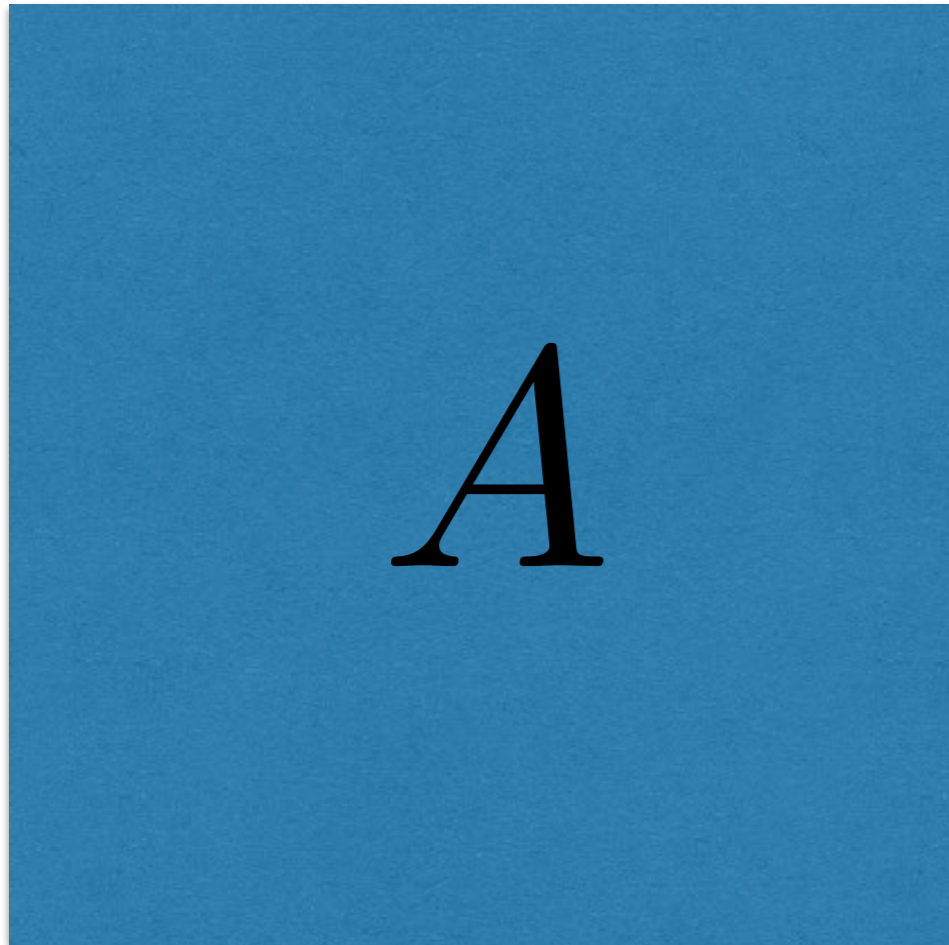
$$A_{i,j} = \begin{cases} 1 & \text{if } (i,j) \in E \\ 0 & \text{otherwise} \end{cases}$$



A is adjacency matrix of a graph

Steps

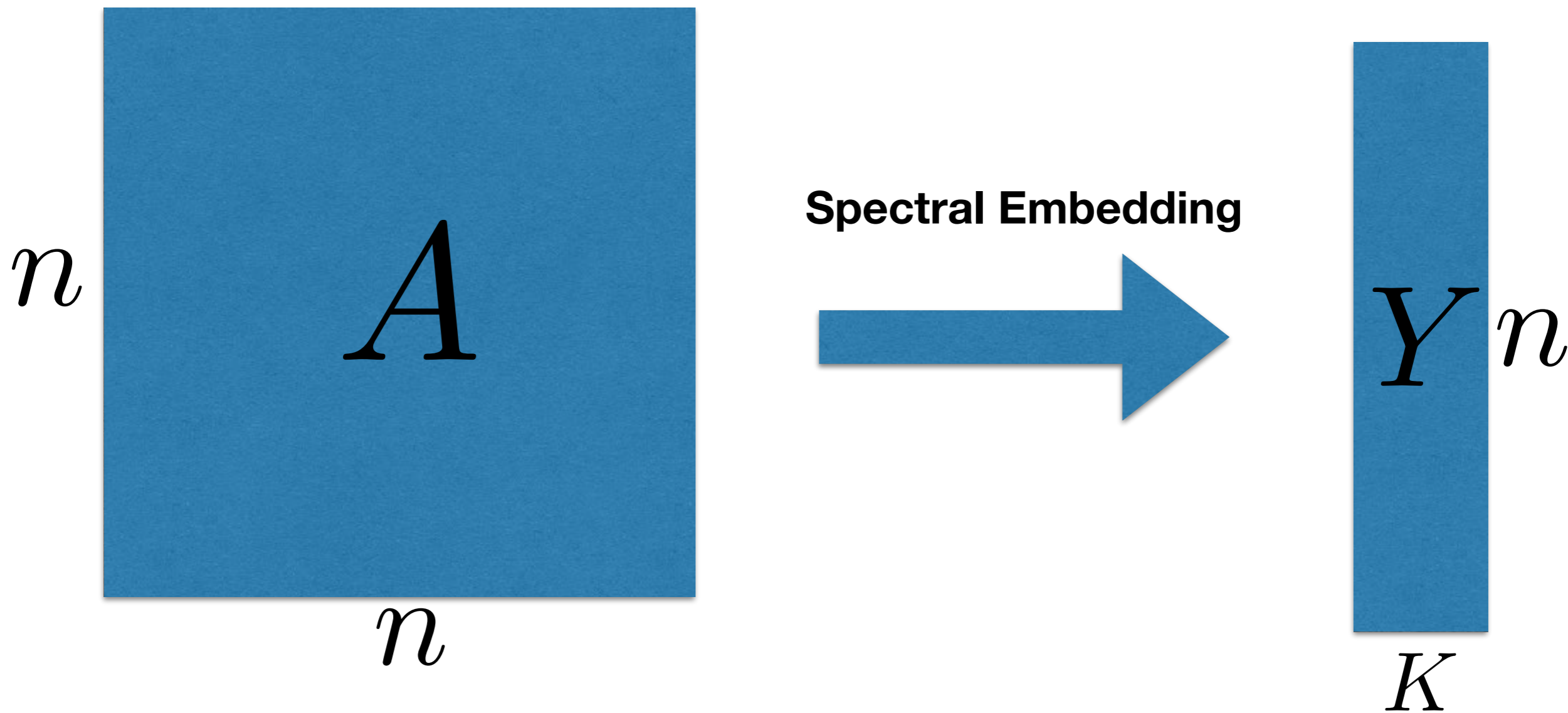
n



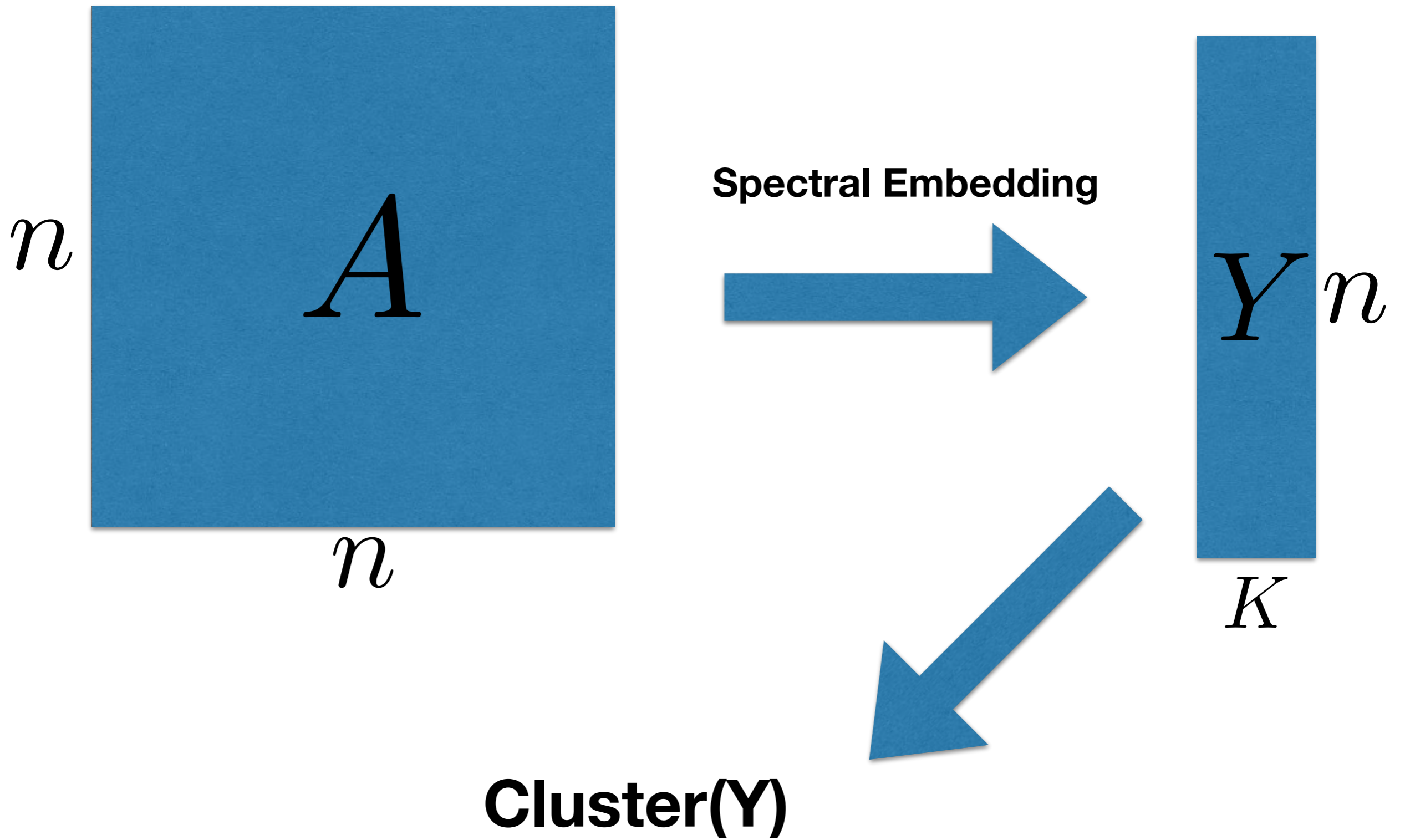
A

n

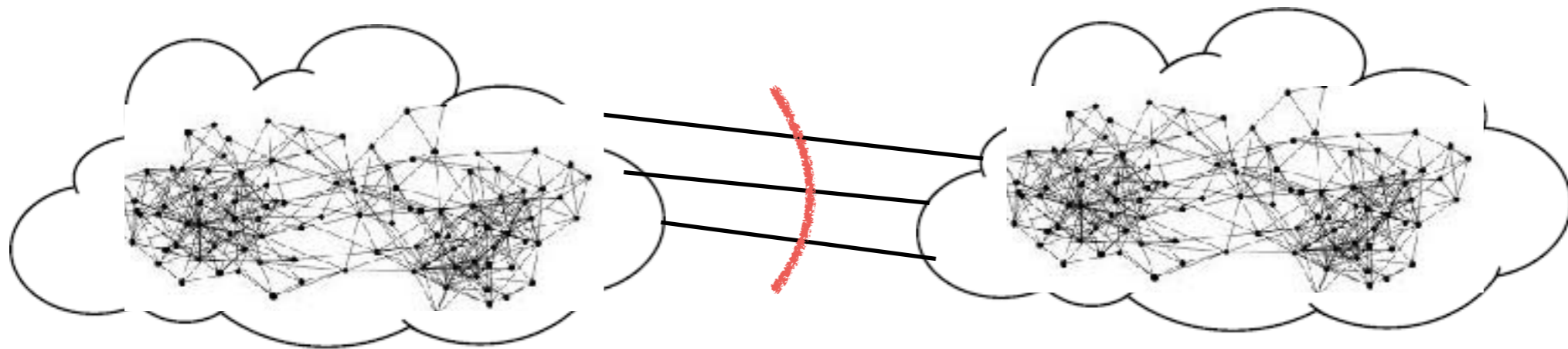
Steps



Steps

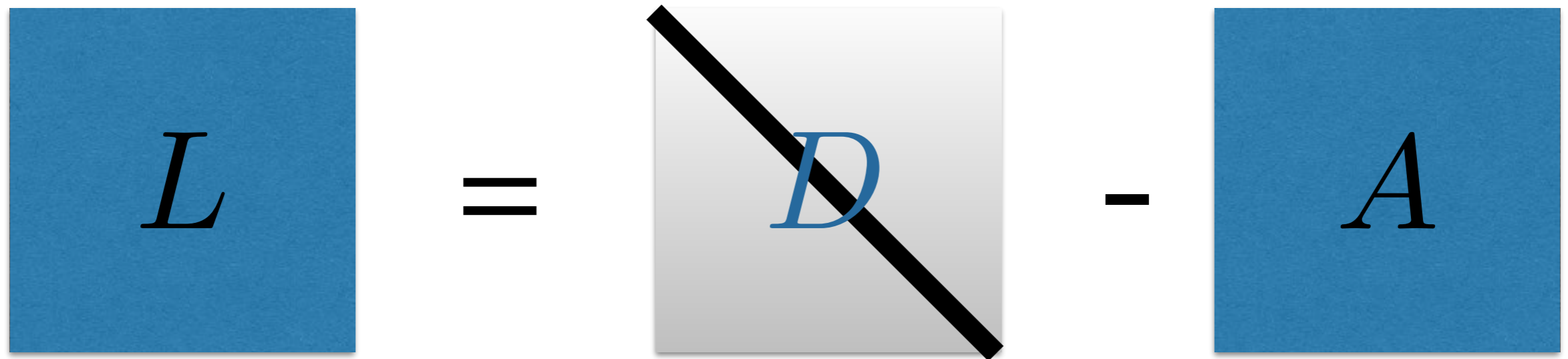


What is the Embedding?



- Map each node in V to \mathbb{R}^k
- Nodes linked to each other are close
- Disconnected groups of nodes are far from each other

SPECTRAL CLUSTERING

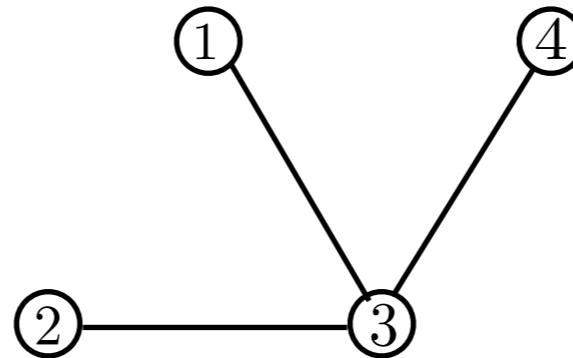
$$L = D - A$$
The diagram illustrates the equation $L = D - A$ using three square boxes. The first box on the left is blue and contains the letter L . To its right is an equals sign. The second box is gray and contains the letter D , but it is crossed out with a thick black diagonal line from the top-left to the bottom-right. To the right of this box is a minus sign. The final box on the right is blue and contains the letter A .

SPECTRAL CLUSTERING

A diagram illustrating the relationship between the Laplacian matrix L , the degree matrix D , and the adjacency matrix A . On the left is a blue square containing the letter L . To its right is an equals sign. Next is a light gray square containing the letter D , which is crossed out with a thick black diagonal line. To the right of this is a minus sign, followed by a blue square containing the letter A .

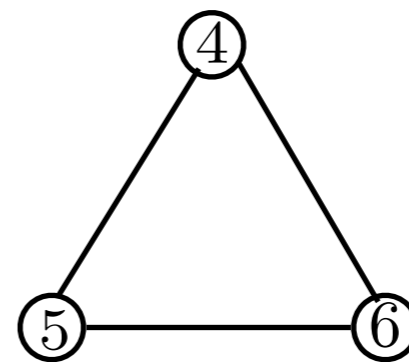
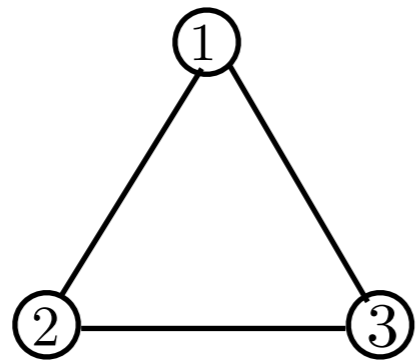
$$D_{i,i} = \sum_{j=1}^n A_{i,j}$$

GRAPH CLUSTERING



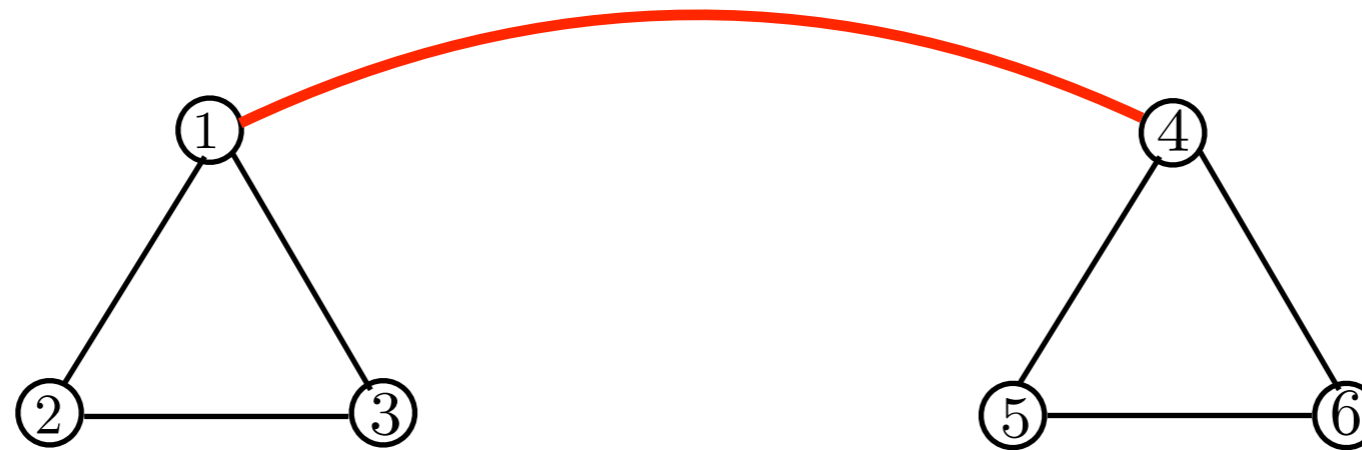
- Fact: For a connected graph, exactly one, the smallest of eigenvalues is 0 , corresponding eigenvector is $\mathbf{1} = (1, \dots, 1)^\top$
Proof: Sum of each row of L is 0 because $D_{i,i} = \sum_{j=1}^n A_{i,j}$ and $L = D - A$

GRAPH CLUSTERING



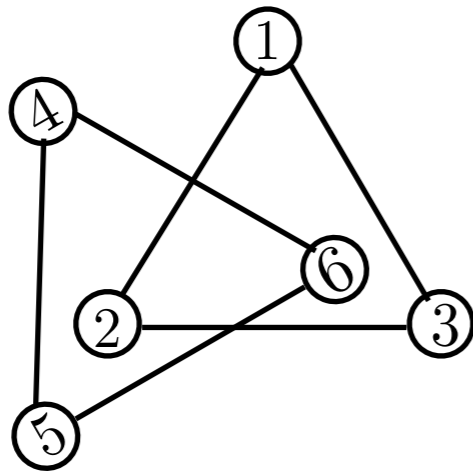
- Fact: For general graph, number of 0 eigenvalues correspond to number of connected components. The corresponding eigenvectors are all 1's on the nodes of connected components
Proof: L is block diagonal. Use connected graph result on each component.

GRAPH CLUSTERING

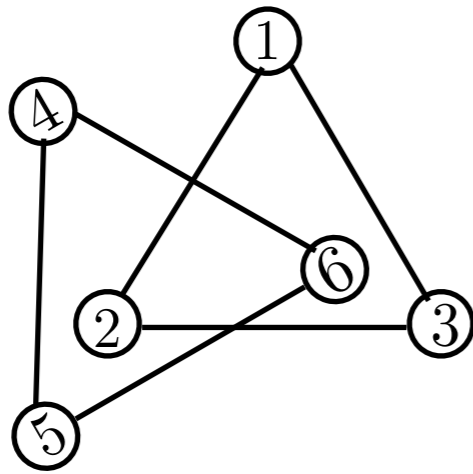


- Fact: For general graph, number of 0 eigenvalues correspond to number of connected components. The corresponding eigenvectors are all 1's on the nodes of connected components
Proof: L is block diagonal. Use connected graph result on each component.

Examples



Examples

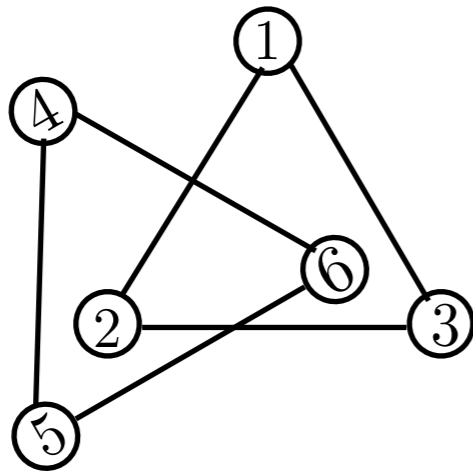


1,2,3

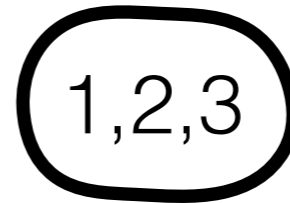
1D

4,5,6

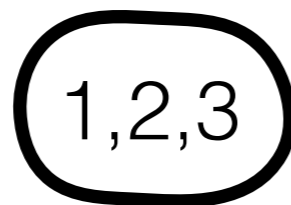
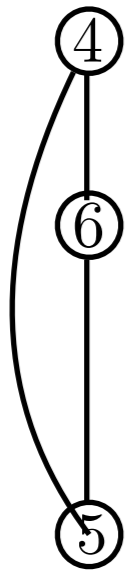
Examples



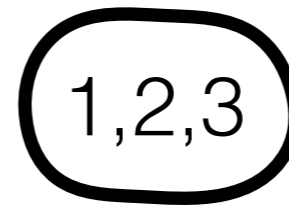
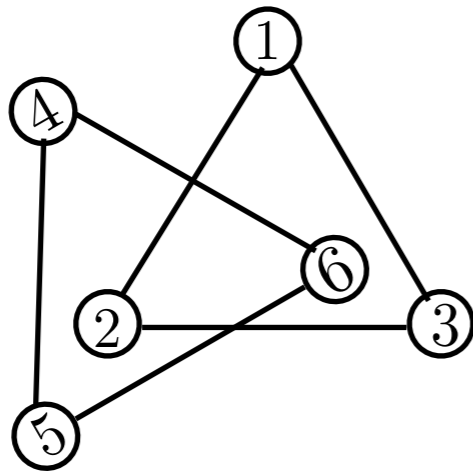
1D



2D



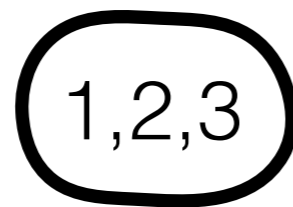
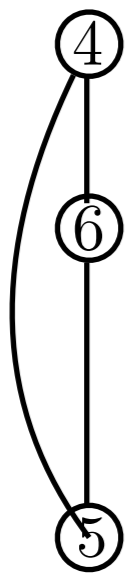
Examples



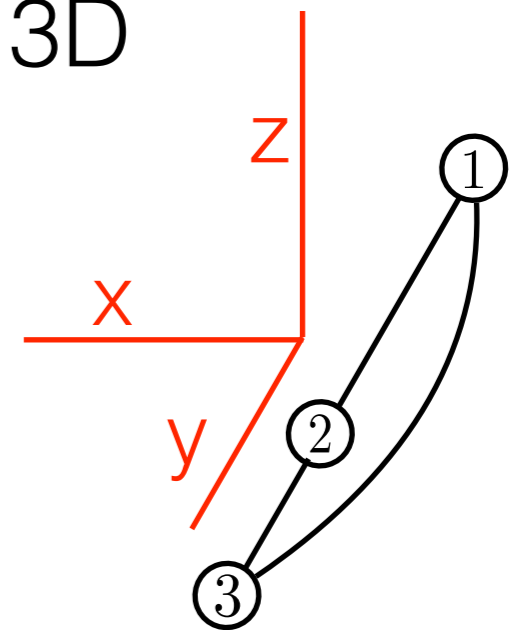
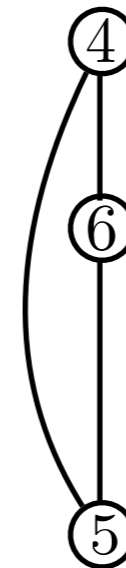
1D



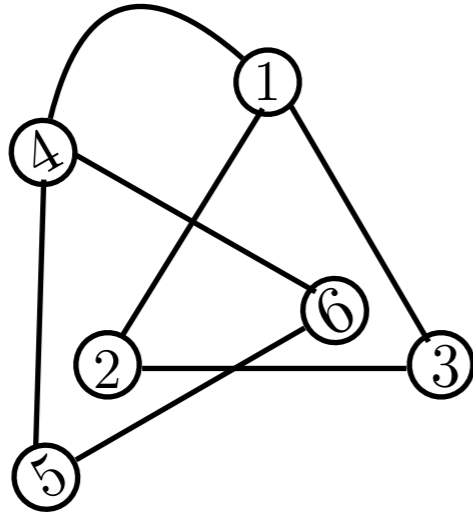
2D



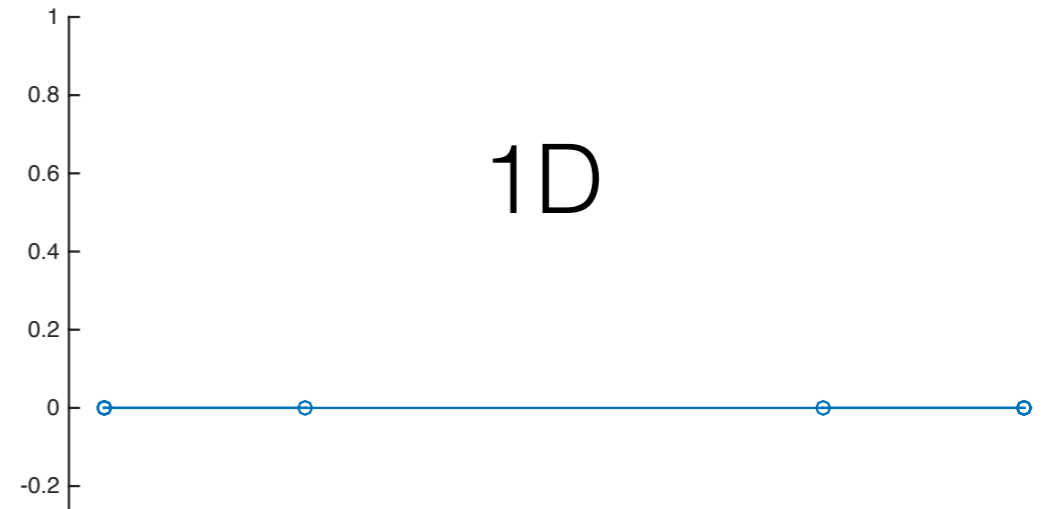
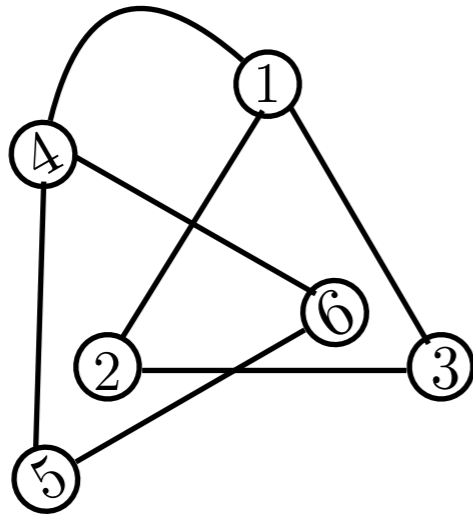
3D



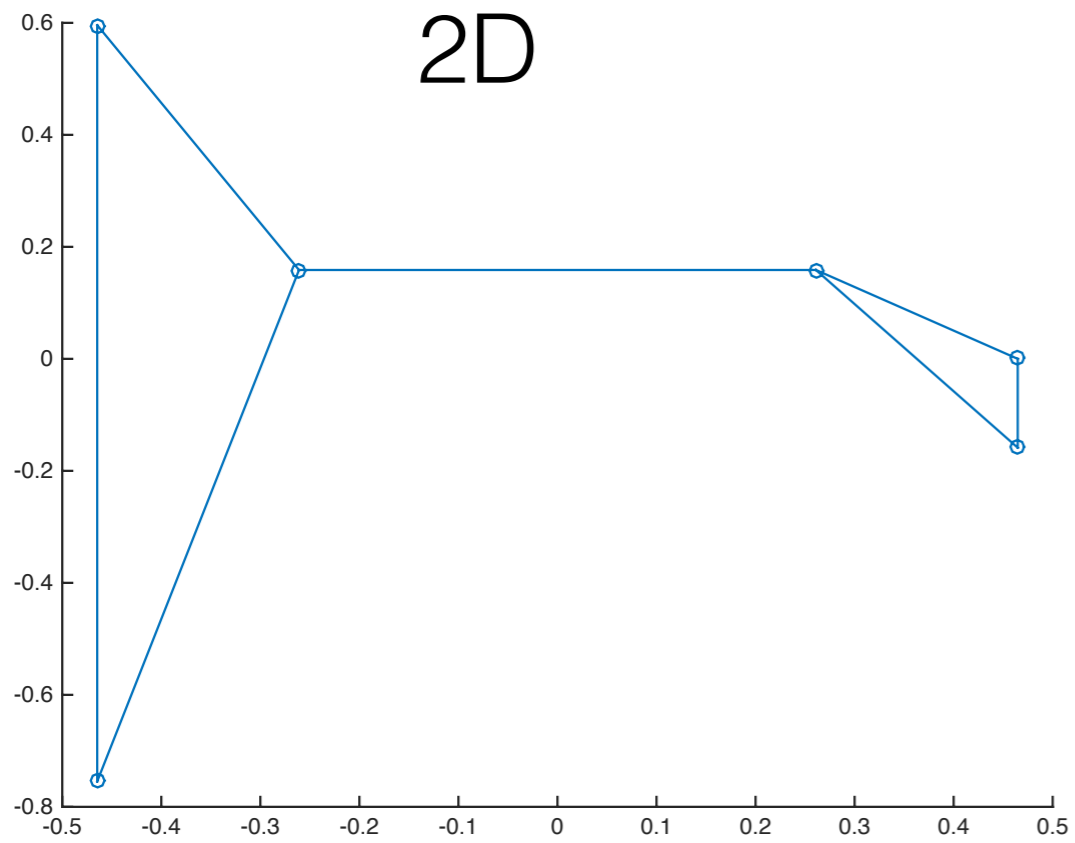
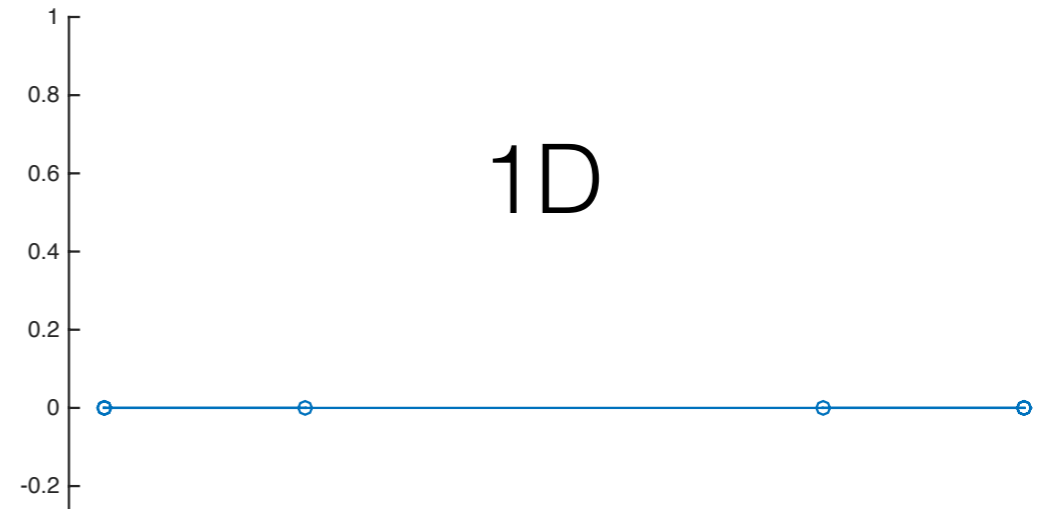
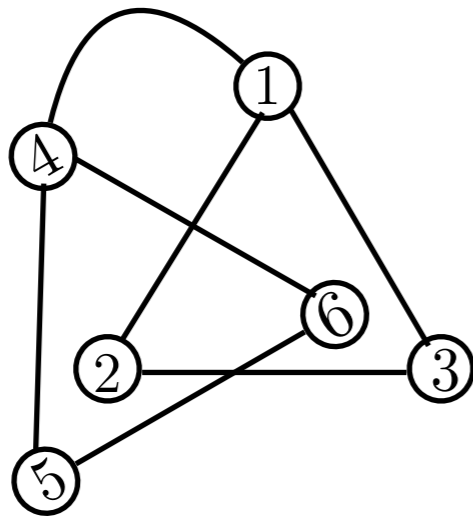
Examples



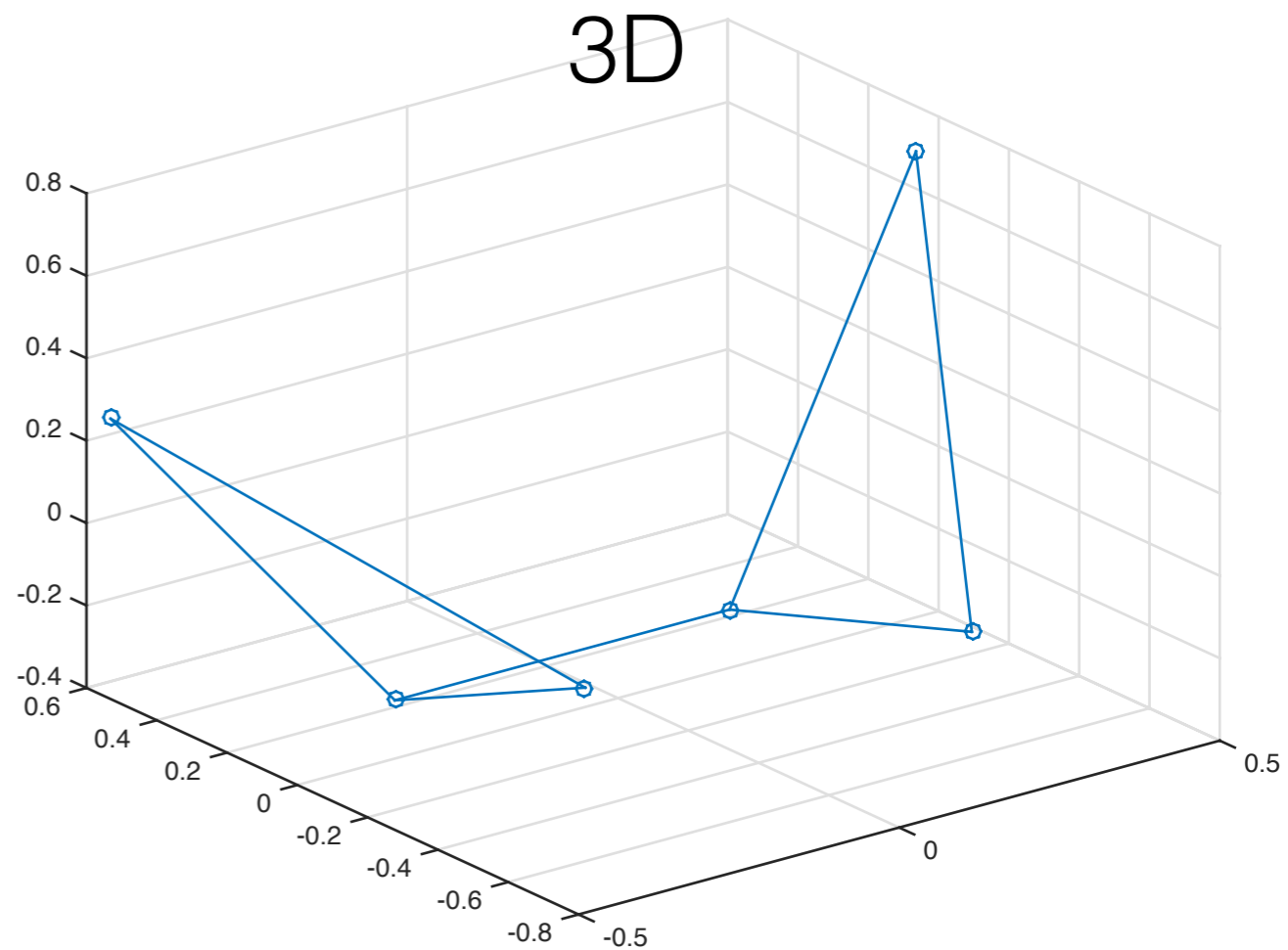
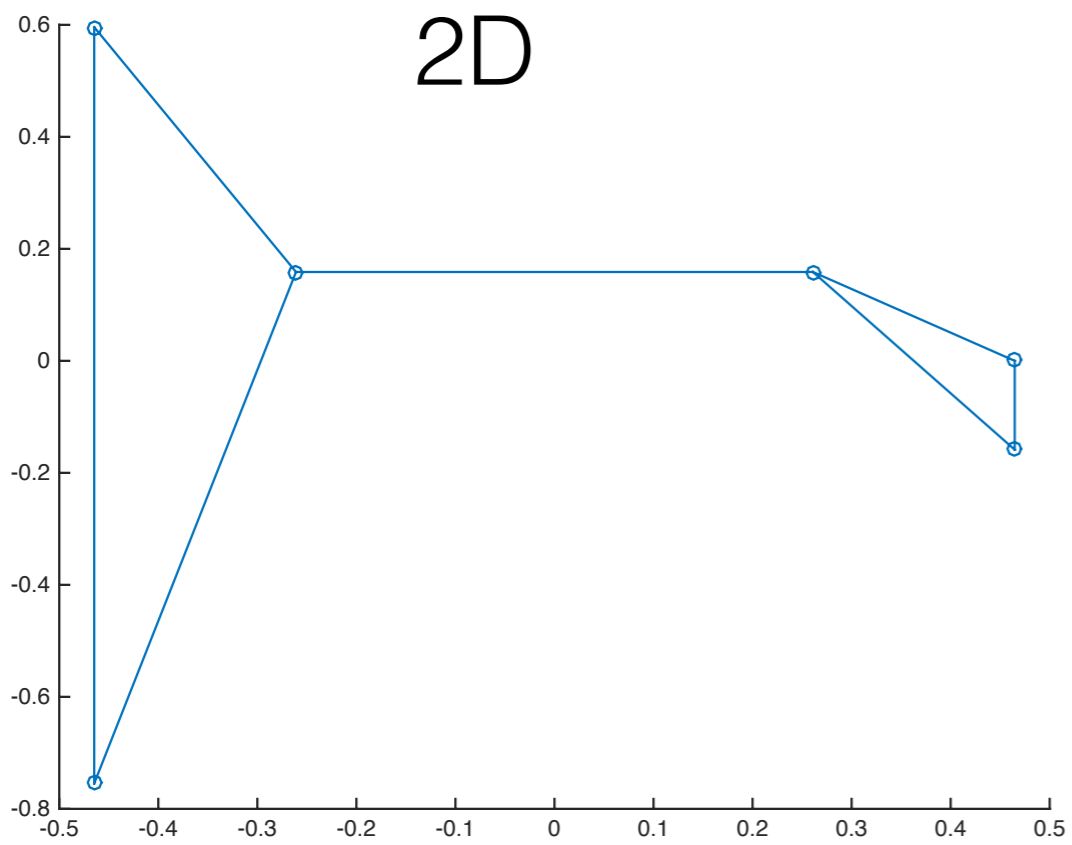
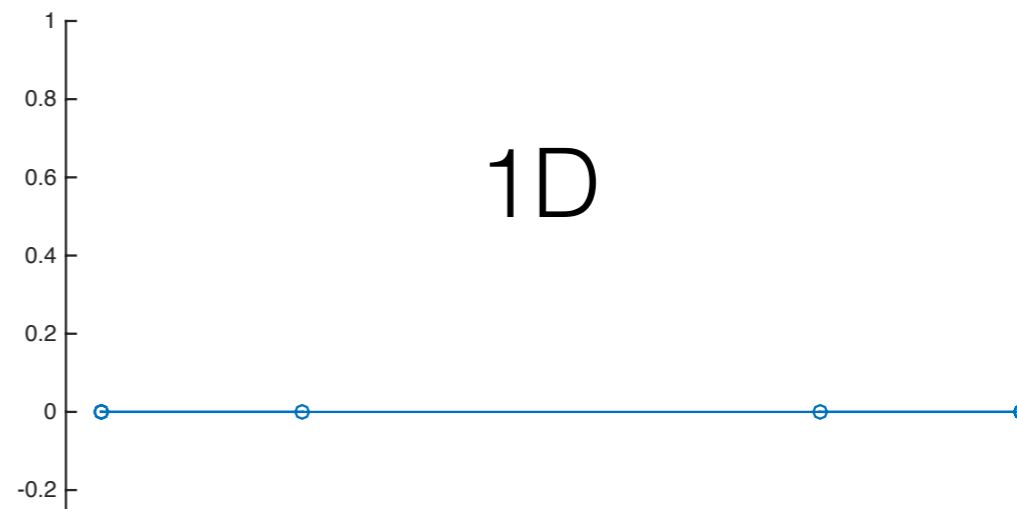
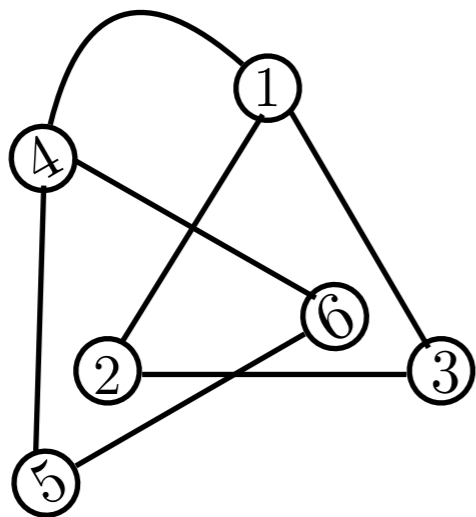
Examples



Examples



Examples



More Examples

Spectral Embedding

- Nodes linked to each other are close in embedded space
- What has this got to do with Laplacian matrix?

CUTS AND LAPLACIAN

K = 1

$$\text{Obj}(c) = \frac{1}{2} \sum_{(i,j) \in E} (c_i - c_j)^2$$

CUTS AND LAPLACIAN

K = 1

$$\begin{aligned}\text{Obj}(c) &= \frac{1}{2} \sum_{(i,j) \in E} (c_i - c_j)^2 \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n A_{i,j} (c_i - c_j)^2\end{aligned}$$

CUTS AND LAPLACIAN

K = 1

$$\begin{aligned}\text{Obj}(c) &= \frac{1}{2} \sum_{(i,j) \in E} (c_i - c_j)^2 \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n A_{i,j} (c_i - c_j)^2 \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n A_{i,j} (c_i^2 + c_j^2 - 2c_i c_j)\end{aligned}$$

CUTS AND LAPLACIAN

K = 1

$$\begin{aligned}\text{Obj}(c) &= \frac{1}{2} \sum_{(i,j) \in E} (c_i - c_j)^2 \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n A_{i,j} (c_i - c_j)^2 \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n A_{i,j} (c_i^2 + c_j^2 - 2c_i c_j) \\ &= \frac{1}{2} \sum_{i=1}^n \left(\sum_{j=1}^n A_{i,j} \right) c_i^2 + \frac{1}{2} \sum_{j=1}^n \left(\sum_{i=1}^n A_{i,j} \right) c_j^2 - \sum_{i=1}^n \sum_{j=1}^n A_{i,j} c_i c_j\end{aligned}$$

CUTS AND LAPLACIAN

K = 1

$$\begin{aligned}\text{Obj}(c) &= \frac{1}{2} \sum_{(i,j) \in E} (c_i - c_j)^2 \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n A_{i,j} (c_i - c_j)^2 \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n A_{i,j} (c_i^2 + c_j^2 - 2c_i c_j) \\ &= \frac{1}{2} \sum_{i=1}^n \left(\sum_{j=1}^n A_{i,j} \right) c_i^2 + \frac{1}{2} \sum_{j=1}^n \left(\sum_{i=1}^n A_{i,j} \right) c_j^2 - \sum_{i=1}^n \sum_{j=1}^n A_{i,j} c_i c_j \\ &= \frac{1}{2} \sum_{i=1}^n D_{i,i} c_i^2 + \frac{1}{2} \sum_{j=1}^n D_{j,j} c_j^2 - \sum_{i=1}^n \sum_{j=1}^n A_{i,j} c_i c_j\end{aligned}$$

CUTS AND LAPLACIAN

K = 1

$$\begin{aligned}\text{Obj}(c) &= \frac{1}{2} \sum_{(i,j) \in E} (c_i - c_j)^2 \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n A_{i,j} (c_i - c_j)^2 \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n A_{i,j} (c_i^2 + c_j^2 - 2c_i c_j) \\ &= \frac{1}{2} \sum_{i=1}^n \left(\sum_{j=1}^n A_{i,j} \right) c_i^2 + \frac{1}{2} \sum_{j=1}^n \left(\sum_{i=1}^n A_{i,j} \right) c_j^2 - \sum_{i=1}^n \sum_{j=1}^n A_{i,j} c_i c_j \\ &= \frac{1}{2} \sum_{i=1}^n D_{i,i} c_i^2 + \frac{1}{2} \sum_{j=1}^n D_{j,j} c_j^2 - \sum_{i=1}^n \sum_{j=1}^n A_{i,j} c_i c_j \\ &= c^\top D c - c^\top A c = c^\top L c\end{aligned}$$

SPECTRAL CLUSTERING, $K = 1$

Hence to find the solution we need to solve for

$$\text{Minimize } c^T L c \quad \text{s.t.} \quad \|c\| = 1$$

SPECTRAL CLUSTERING, $K = 1$

Hence to find the solution we need to solve for

$$\text{Minimize } c^T L c \quad \text{s.t.} \quad \|c\| = 1$$

Hence solution c to above is an Eigen vector, first smallest one is the all 1's vector (for connected graph), second smallest one is our solution

To get clustering assignment we simply threshold at 0

SPECTRAL CLUSTERING, $K > 1$

- Solution obtained by considering the second smallest up to K^{th} smallest eigenvectors

$$\text{Obj}(c) = \sum_{k=1}^K c^k \top L c^k$$

c^k 's are orthogonal to each other and the all ones vector

SPECTRAL CLUSTERING ALGORITHM (UNNORMALIZED)

- 1 Given matrix A calculate diagonal matrix D s.t. $D_{i,i} = \sum_{j=1}^n A_{i,j}$
- 2 Calculate the Laplacian matrix $L = D - A$
- 3 Find eigen vectors $\mathbf{v}_1, \dots, \mathbf{v}_n$ of L (ascending order of eigenvalues)
- 4 Pick the K eigenvectors with smallest eigenvalues to get $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^K$
- 5 Use K-means clustering algorithm on $\mathbf{y}_1, \dots, \mathbf{y}_n$

SPECTRAL CLUSTERING ALGORITHM (UNNORMALIZED)

- 1 Given matrix A calculate diagonal matrix D s.t. $D_{i,i} = \sum_{j=1}^n A_{i,j}$
- 2 Calculate the Laplacian matrix $L = D - A$
- 3 Find eigen vectors $\mathbf{v}_1, \dots, \mathbf{v}_n$ of L (ascending order of eigenvalues)
- 4 Pick the K eigenvectors with smallest eigenvalues to get $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^K$
- 5 Use K-means clustering algorithm on $\mathbf{y}_1, \dots, \mathbf{y}_n$

$\mathbf{y}_1, \dots, \mathbf{y}_n$ are called spectral embedding

SPECTRAL CLUSTERING ALGORITHM (UNNORMALIZED)

- 1 Given matrix A calculate diagonal matrix D s.t. $D_{i,i} = \sum_{j=1}^n A_{i,j}$
- 2 Calculate the Laplacian matrix $L = D - A$
- 3 Find eigen vectors $\mathbf{v}_1, \dots, \mathbf{v}_n$ of L (ascending order of eigenvalues)
- 4 Pick the K eigenvectors with smallest eigenvalues to get $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^K$
- 5 Use K-means clustering algorithm on $\mathbf{y}_1, \dots, \mathbf{y}_n$

$\mathbf{y}_1, \dots, \mathbf{y}_n$ are called spectral embedding

Embeds the n nodes into $K-1$ dimensional vectors

- Unnormalized Spectral clustering aims to cluster based on minimizing cut

- Unnormalized Spectral clustering aims to cluster based on minimizing cut
- cut: Number of edges that need to be deleted to have no links between the cluster and other nodes outside

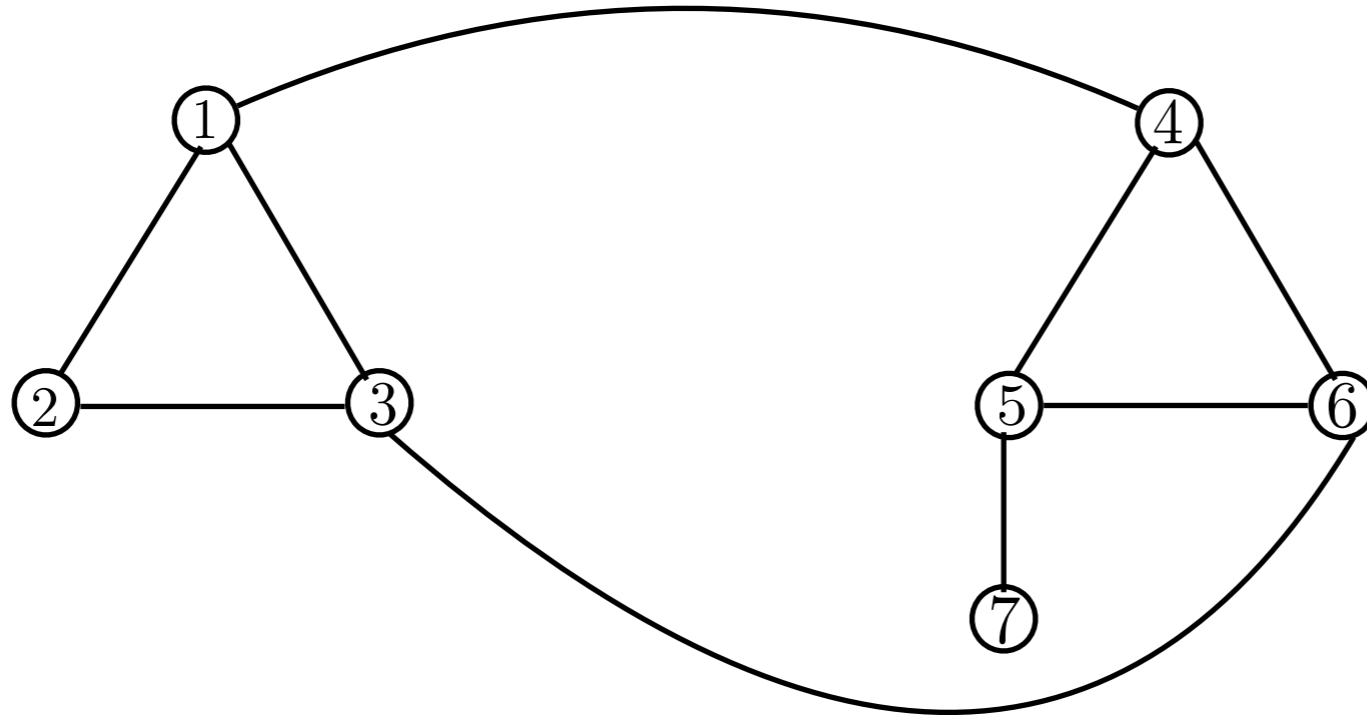
- Unnormalized Spectral clustering aims to cluster based on minimizing cut
- cut: Number of edges that need to be deleted to have no links between the cluster and other nodes outside
- But is cut the right metric?

NORMALIZED CUT

- Why cut is perhaps not a good measure?

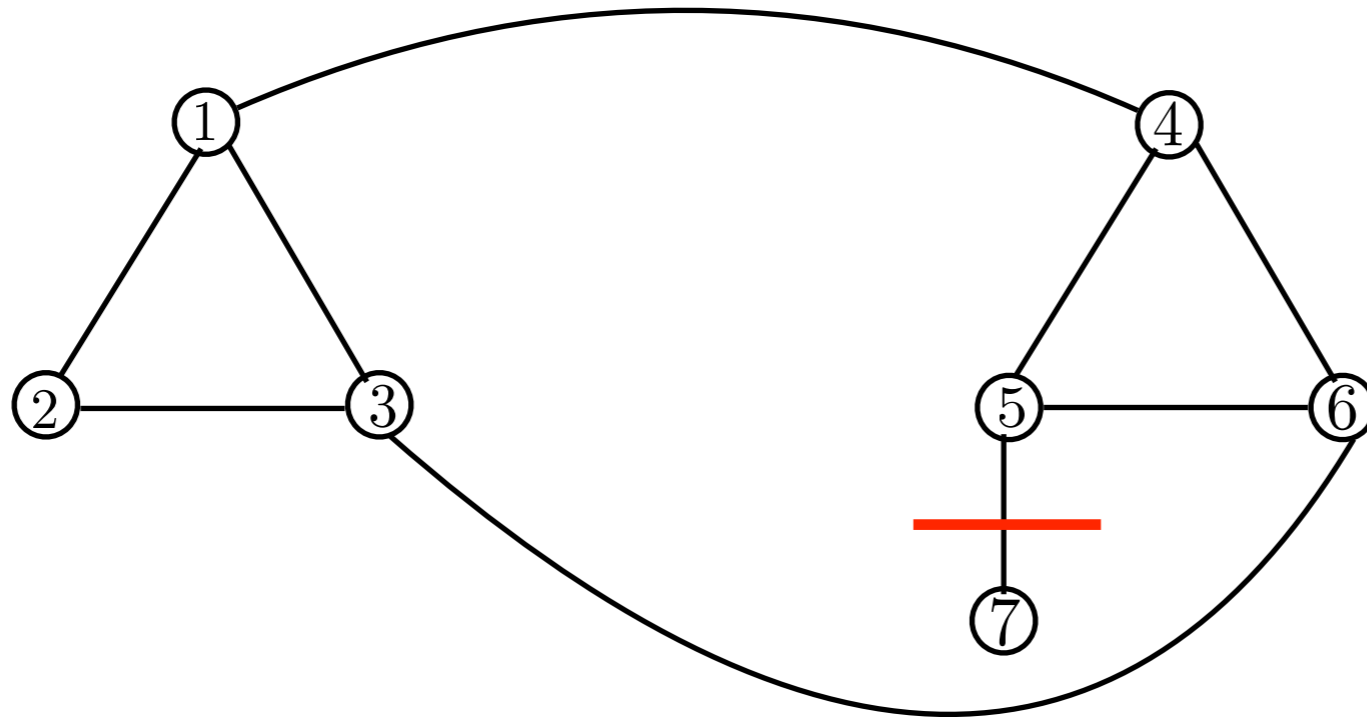
NORMALIZED CUT

- Why cut is perhaps not a good measure?



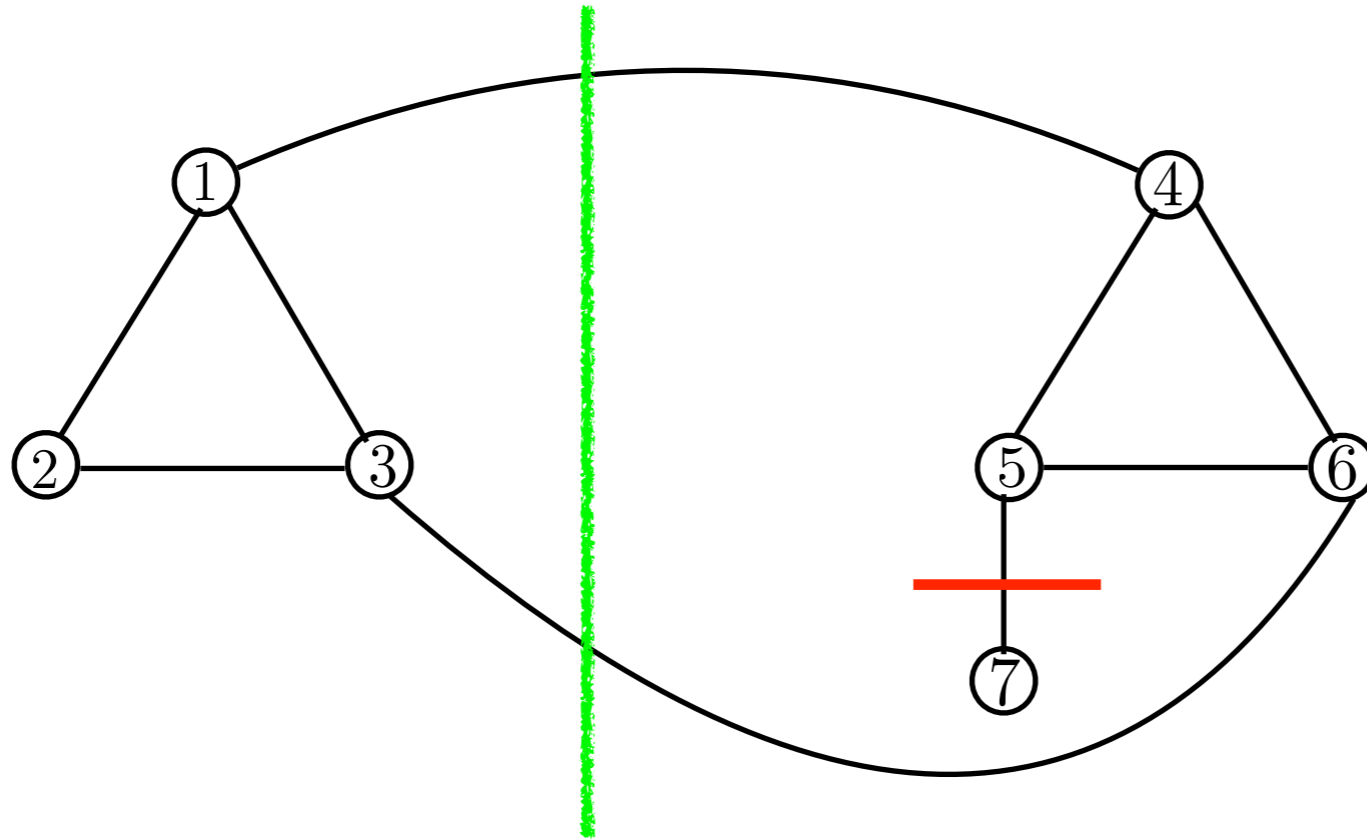
NORMALIZED CUT

- Why cut is perhaps not a good measure?



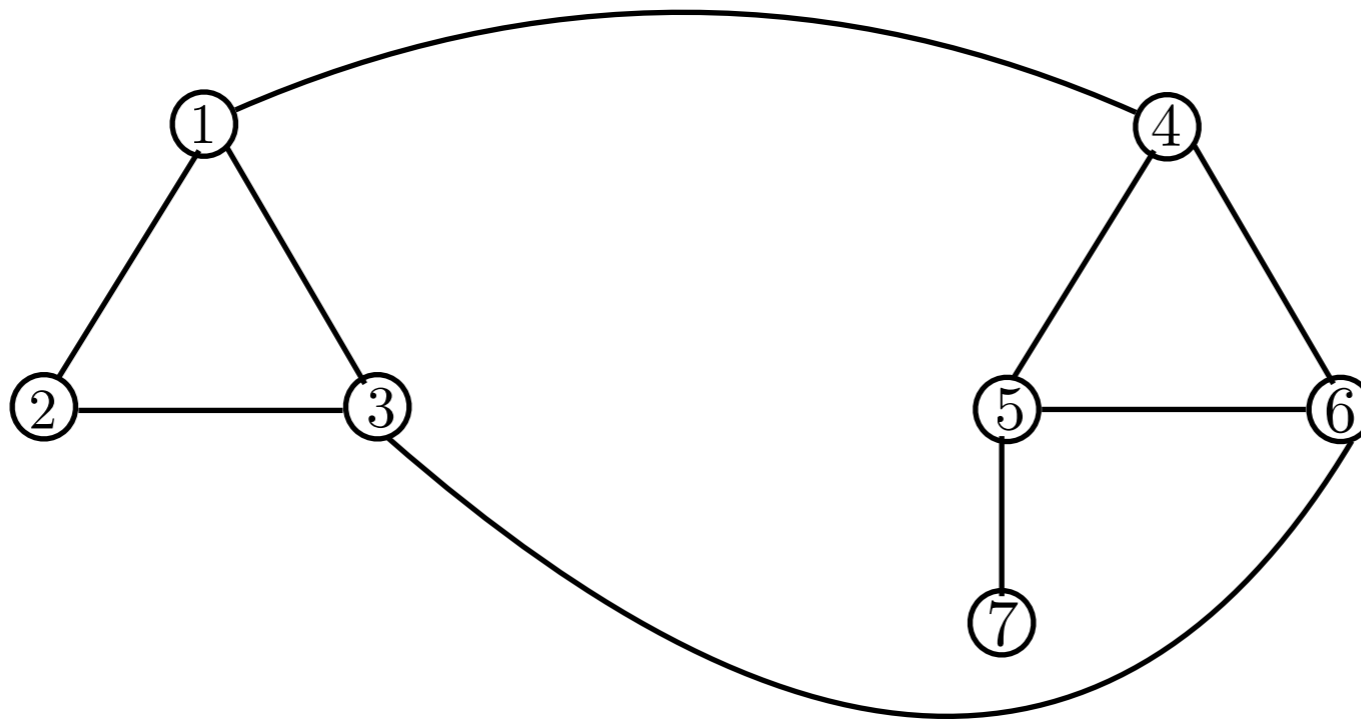
NORMALIZED CUT

- Why cut is perhaps not a good measure?



RATIO CUT

- Why cut is perhaps not a good measure?
- Fixes?



NORMALIZED CUT

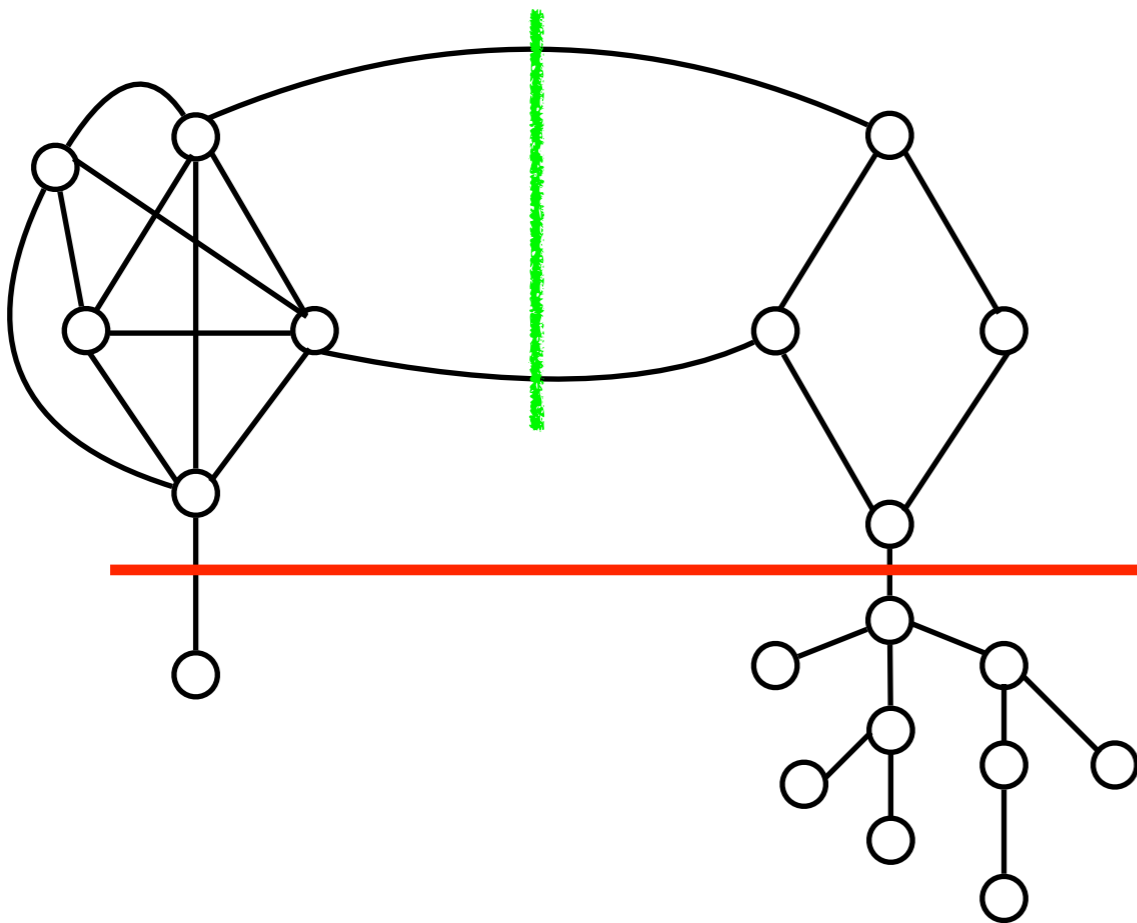
- Normalized cut: Minimize sum of ratio of number of edges cut per cluster and number of edges within cluster

$$\text{NCUT} = \sum_j \frac{\text{CUT}(C_j)}{\text{Edges}(C_j)}$$

NORMALIZED CUT

- Normalized cut: Minimize sum of ratio of number of edges cut per cluster and number of edges within cluster

$$\text{NCUT} = \sum_j \frac{\text{CUT}(C_j)}{\text{Edges}(C_j)}$$



NORMALIZED CUT

- Normalized cut: Minimize sum of ratio of number of edges cut per cluster and number of edges within cluster

$$\text{NCUT} = \sum_j \frac{\text{CUT}(C_j)}{\text{Edges}(C_j)}$$

- Example $K = 2$

$$\text{CUT}(C_1, C_2) \left(\frac{1}{\text{Edges}(C_1)} + \frac{1}{\text{Edges}(C_2)} \right)$$

NORMALIZED CUT

- Normalized cut: Minimize sum of ratio of number of edges cut per cluster and number of edges within cluster

$$\text{NCUT} = \sum_j \frac{\text{CUT}(C_j)}{\text{Edges}(C_j)}$$

- Example $K = 2$

$$\text{CUT}(C_1, C_2) \left(\frac{1}{\text{Edges}(C_1)} + \frac{1}{\text{Edges}(C_2)} \right)$$

- This is an NP hard problem! ...so relax

NORMALIZED SPECTRAL CLUSTERING

- As before, we want to minimize $\sum_{(i,j) \in E} (c_i - c_j)^2 = c^\top Lc$

NORMALIZED SPECTRAL CLUSTERING

- As before, we want to minimize $\sum_{(i,j) \in E} (c_i - c_j)^2 = c^T L c$
- But we also want to weight the values of c_i 's based on degree. We want high degree nodes to have larger c magnitude

NORMALIZED SPECTRAL CLUSTERING

- As before, we want to minimize $\sum_{(i,j) \in E} (c_i - c_j)^2 = c^\top Lc$
- But we also want to weight the values of c_i 's based on degree. We want high degree nodes to have larger c magnitude
- That is we want to simultaneously maximize $\sum_{i=1}^n c_i^2 D_{i,i} = c^\top Dc$

NORMALIZED SPECTRAL CLUSTERING

- As before, we want to minimize $\sum_{(i,j) \in E} (c_i - c_j)^2 = c^\top Lc$
- But we also want to weight the values of c_i 's based on degree. We want high degree nodes to have larger c magnitude
- That is we want to simultaneously maximize $\sum_{i=1}^n c_i^2 D_{i,i} = c^\top Dc$
- Find c so as to:

$$\begin{aligned} & \text{minimize } \frac{c^\top Lc}{c^\top Dc} \\ & \equiv \text{minimize } c^\top Lc \text{ subject to } c^\top Dc = 1 \end{aligned}$$

NORMALIZED SPECTRAL CLUSTERING

- As before, we want to minimize $\sum_{(i,j) \in E} (c_i - c_j)^2 = c^\top Lc$
- But we also want to weight the values of c_i 's based on degree. We want high degree nodes to have larger c magnitude
- That is we want to simultaneously maximize $\sum_{i=1}^n c_i^2 D_{i,i} = c^\top Dc$
- Find c so as to:

$$\text{minimize } \frac{c^\top Lc}{c^\top Dc}$$

$$\equiv \text{minimize } c^\top Lc \text{ subject to } c^\top Dc = 1$$

$$\equiv \text{minimize } u^\top D^{-1/2} L D^{-1/2} u \text{ subject to } \|u\| = 1$$

SPECTRAL CLUSTERING

Minimize $c^\top \tilde{L}c$ s.t. $c \perp \mathbf{1}$

SPECTRAL CLUSTERING

Minimize $c^\top \tilde{L}c$ s.t. $c \perp \mathbf{1}$

Approximately Minimize normalized cut!

SPECTRAL CLUSTERING

Minimize $c^\top \tilde{L}c$ s.t. $c \perp \mathbf{1}$

Approximately Minimize normalized cut!

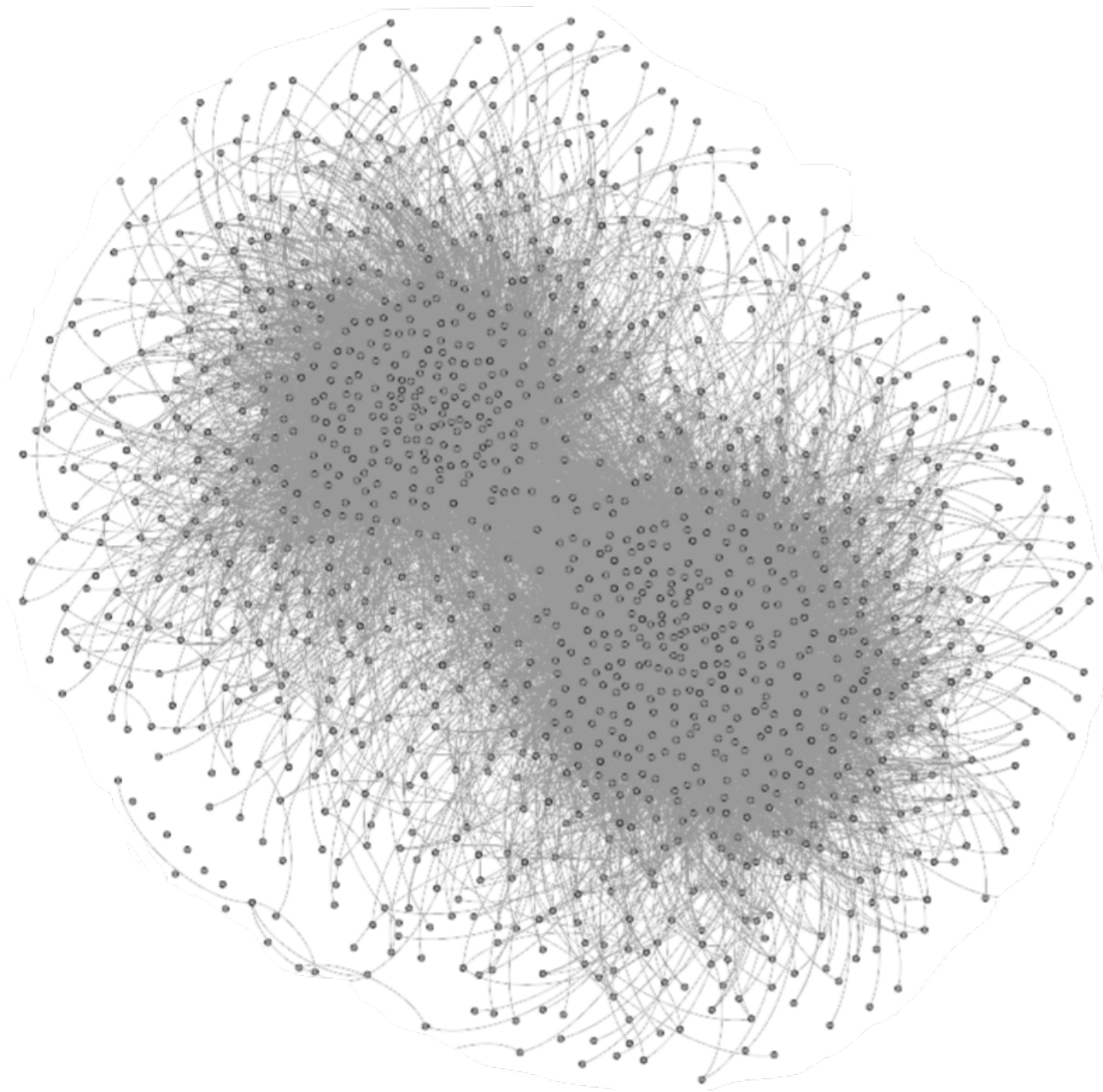
- Solution: Find second smallest eigenvectors of $\tilde{L} = I - D^{-1/2}AD^{-1/2}$

SPECTRAL CLUSTERING ALGORITHM (NORMALIZED)

- 1 Given matrix A calculate diagonal matrix D s.t. $D_{i,i} = \sum_{j=1}^n A_{i,j}$
- 2 Calculate the normalized Laplacian matrix $\tilde{L} = I - D^{-1/2}AD^{-1/2}$
- 3 Find eigen vectors $\mathbf{v}_1, \dots, \mathbf{v}_n$ of \tilde{L} (ascending order of eigenvalues)
- 4 Pick the K eigenvectors with smallest eigenvalues to get $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^K$
- 5 Use K-means clustering algorithm on $\mathbf{y}_1, \dots, \mathbf{y}_n$

Demo

SPECTRAL CLUSTERING



SPECTRAL CLUSTERING

