## COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Cameron Musco

University of Massachusetts Amherst. Fall 2019.
Lecture 15

Last Class:

Last Class:

- Entity embeddings (e.g., word embeddings).
- Dimensionality reduction for data not lying close to a low-dimensional subspace (non-linear dimensionality reduction).
- Approach via low-rank approximation of a graph based similarity matrix (adjacency matrix).
- Spectral graph theory, spectral clustering, graph Laplacian.

### Last Class:

- Entity embeddings (e.g., word embeddings).
- Dimensionality reduction for data not lying close to a low-dimensional subspace (non-linear dimensionality reduction).
- Approach via low-rank approximation of a graph based similarity matrix (adjacency matrix).
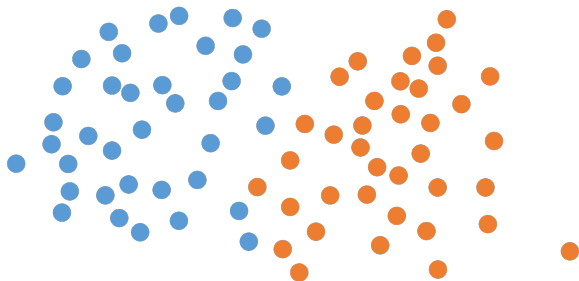- Spectral graph theory, spectral clustering, graph Laplacian.

### This Class: Finish up spectral clustering.

- Clustering non-linearly separable data via graph eigenvectors.
- Application to the *stochastic block model* and community detection.

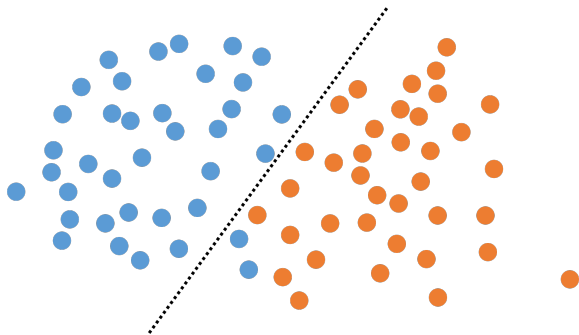Goal: Partition or cluster vertices in a graph based on 'similarity'.

**Goal:** Partition or cluster vertices in a graph based on 'similarity'.
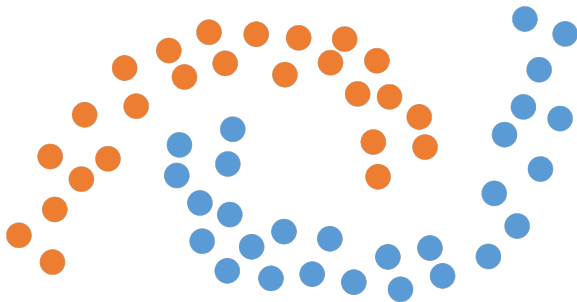
**Linearly separable data.**

**Goal:** Partition or cluster vertices in a graph based on 'similarity'.
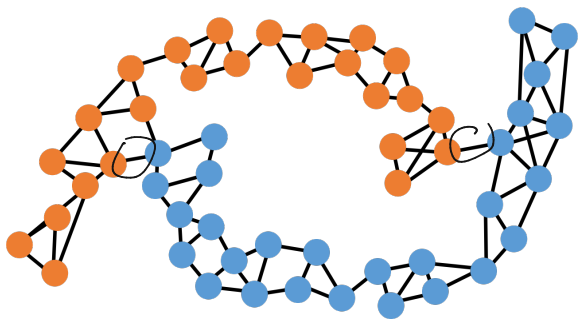
Linearly separable data.

**Goal:** Partition or cluster vertices in a graph based on 'similarity'.

**Non-linearly separable data $k$-nearest neighbor graph.**

Goal: Partition or cluster vertices in a graph based on 'similarity'.

Non-linearly separable data $k$-nearest neighbor graph.
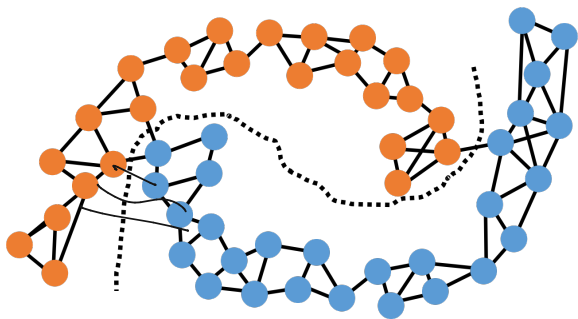
**Goal:** Partition or cluster vertices in a graph based on 'similarity'.

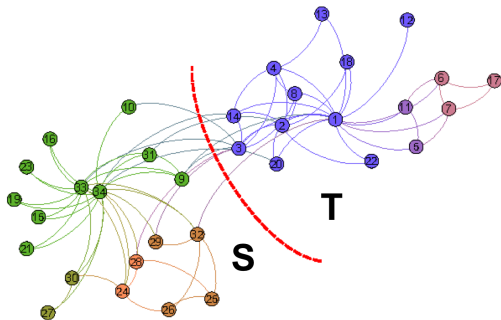**Non-linearly separable data $k$-nearest neighbor graph.**

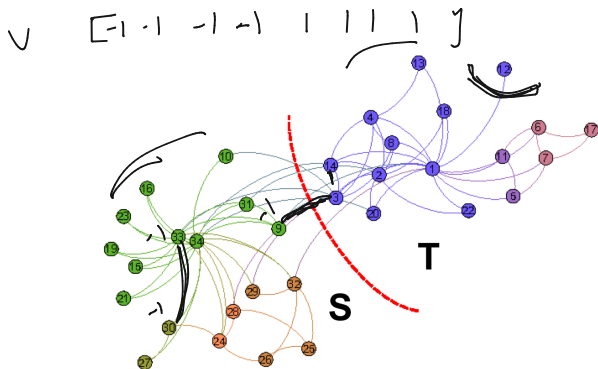**Goal:** Partition or cluster vertices in a graph based on 'similarity'.

**Community detection in naturally occurring networks.**



(a) Zachary Karate Club Graph

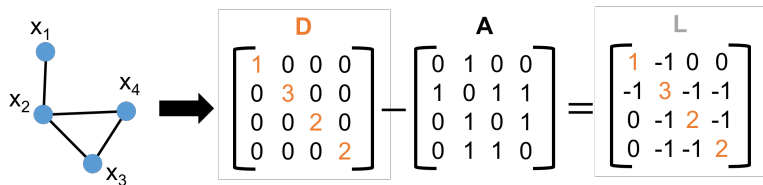**Main Idea:** Partition clusters along a cut that:

1. Has few edges crossing it: $|\{(u, v) \in E : u \in S, v \in T\}|$ is small.
2. Separates large sections of the graph: $|S|, |T|$ are not too small.



(a) Zachary Karate Club Graph

For a graph with adjacency matrix **A** and degree matrix **D**, $L = D - A$ is the graph Laplacian.



$$D = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 2 \end{bmatrix}$$

$$A = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{bmatrix}$$

$$L = \begin{bmatrix} 1 & -1 & 0 & 0 \\ -1 & 3 & -1 & -1 \\ 0 & -1 & 2 & -1 \\ 0 & -1 & -1 & 2 \end{bmatrix}$$

For a cut indicator vector $\vec{v} \in \{-1, 1\}^n$ with $\vec{v}(i) = -1$ for $i \in S$ and $\vec{v}(i) = 1$ for $i \in T$:

$(1 \cdot -1) := 2$   $2^2 : 4$

$(-1 \cdot 1) := 2$

Hedges between $S + T$

1. $\vec{v}^T L \vec{v} = \sum_{(i,j) \in E} (\vec{v}(i) - \vec{v}(j))^2 = 4 \cdot cut(S, T)$.

2. $\vec{v}^T \vec{1} = |T| - |S|$.

$$\sum_{i=1}^{n} v(i) \cdot 1 \; : \; \sum v(i)$$

For a graph with adjacency matrix **A** and degree matrix **D**, $L = D - A$ is the graph Laplacian.



For a cut indicator vector $\vec{v} \in \{-1, 1\}^n$ with $\vec{v}(i) = -1$ for $i \in S$ and $\vec{v}(i) = 1$ for $i \in T$:

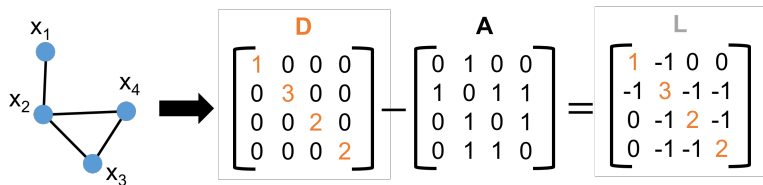1. $\vec{v}^T L \vec{v} = \sum_{(i,j) \in E} (\vec{v}(i) - \vec{v}(j))^2 = 4 \cdot cut(S, T)$.
2. $\vec{v}^T \vec{1} = |V| - |S|$.

Want to minimize both $\vec{v}^T L \vec{v}$ (cut size) and $\vec{v}^T \vec{1}$ (imbalance).

4

The smallest eigenvector of the Laplacian is:

$$\vec{v}_n = \frac{1}{\sqrt{n}} \cdot \vec{1} = \underset{v \in \mathbb{R}^n \text{ with } \|\vec{v}\|=1}{\arg\min} \vec{v}^T L \vec{V}$$

$$\vec{v}_n = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

with $\vec{v}_n^T L \vec{v}_n = 0$.

---

$n$: number of nodes in graph, $A \in \mathbb{R}^{n \times n}$: adjacency matrix, $D \in \mathbb{R}^{n \times n}$: diagonal degree matrix, $L \in \mathbb{R}^{n \times n}$: Laplacian matrix $L = A - D$.

The smallest eigenvector of the Laplacian is:

$$\vec{v}_n = \frac{1}{\sqrt{n}} \cdot \vec{1} = \operatorname*{arg\,min}_{v \in \mathbb{R}^n \text{ with } \|\vec{v}\|=1} \vec{v}^T L \vec{V}$$

with $\vec{v}_n^T L \vec{v}_n = 0$. Why? Use that $L = D - A$.
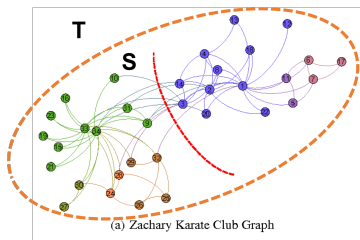
$$V_n^T D V_n - V_n^T A V_n = 0$$

---

$n$: number of nodes in graph, $A \in \mathbb{R}^{n \times n}$: adjacency matrix, $D \in \mathbb{R}^{n \times n}$: diagonal degree matrix, $L \in \mathbb{R}^{n \times n}$: Laplacian matrix $L = A - D$.

The smallest eigenvector of the Laplacian is:

$$\vec{v}_n = \frac{1}{\sqrt{n}} \cdot \vec{1} = \underset{v \in \mathbb{R}^n \text{ with } \|\vec{v}\|=1}{\arg\min} \vec{v}^T L \vec{V}$$

with $\vec{v}_n^T L \vec{v}_n = 0$. Why? Use that $L = D - A$.



(a) Zachary Karate Club Graph

---

$n$: number of nodes in graph, $\mathbf{A} \in \mathbb{R}^{n \times n}$: adjacency matrix, $\mathbf{D} \in \mathbb{R}^{n \times n}$: diagonal degree matrix, $\mathbf{L} \in \mathbb{R}^{n \times n}$: Laplacian matrix $\mathbf{L} = \mathbf{A} - \mathbf{D}$.

$$\left(\widehat{V_1 \quad V_2 \quad V_3}\right) \cdots \left(\widehat{V_n}\right)$$

By Courant-Fischer, the second smallest eigenvector is given by:

$$\vec{v}_{n-1} = \underset{v \in \mathbb{R}^n \text{ with } \|\vec{v}\|=1, \; \underbrace{\vec{v}_n^T \vec{v}=0}}{\arg \min} \vec{v}^T L \vec{v} \; = \; \sum \vec{v}(i) \; = \; 0$$

n: number of nodes in graph, $A \in \mathbb{R}^{n \times n}$: adjacency matrix, $D \in \mathbb{R}^{n \times n}$: diagonal degree matrix, $L \in \mathbb{R}^{n \times n}$: Laplacian matrix $L = A - D$. $S, T$: vertex sets on different sides of cut.

By Courant-Fischer, the second smallest eigenvector is given by:

$$\vec{v}_{n-1} = \underset{v \in \mathbb{R}^n \text{ with } \|\vec{v}\|=1, \vec{v}_n^T \vec{v}=0}{\arg\min} \vec{v}^T L \vec{v}$$

If $\vec{v}_{n-1}$ were in $\{-1, 1\}^n$ it would have:

- $\vec{v}_{n-1}^T L \vec{v}_{n-1} = 4 cut(S, T)$ as small as possible given that $\vec{v}_{n-1}^T \vec{1} = |T| - |S| = 0$.

---

*n*: number of nodes in graph, $\mathbf{A} \in \mathbb{R}^{n \times n}$: adjacency matrix, $\mathbf{D} \in \mathbb{R}^{n \times n}$: diagonal degree matrix, $\mathbf{L} \in \mathbb{R}^{n \times n}$: Laplacian matrix $\mathbf{L} = \mathbf{A} - \mathbf{D}$. $S, T$: vertex sets on different sides of cut.

By Courant-Fischer, the second smallest eigenvector is given by:

$$\vec{v}_{n-1} = \underset{v \in \mathbb{R}^n \text{ with } \|\vec{v}\|=1, \; \vec{v}_n^T \vec{v}=0}{\arg\min} \vec{v}^T L \vec{V}$$

If $\vec{v}_{n-1}$ were in $\{-1, 1\}^n$ it would have:

- $\vec{v}_{n-1}^T L \vec{v}_{n-1} = cut(S, T)$ as small as possible given that $\vec{v}_{n-1}^T \vec{1} = |T| - |S| = 0$.
- I.e., $\vec{v}_{n-1}$ would indicate the smallest perfectly balanced cut.

---

$n$: number of nodes in graph, $A \in \mathbb{R}^{n \times n}$: adjacency matrix, $D \in \mathbb{R}^{n \times n}$: diagonal degree matrix, $L \in \mathbb{R}^{n \times n}$: Laplacian matrix $L = A - D$. $S, T$: vertex sets on different sides of cut.

By Courant-Fischer, the second smallest eigenvector is given by:

$$\vec{v}_{n-1} = \underset{v \in \mathbb{R}^n \text{ with } \|\vec{v}\|=1, \ \vec{v}_n^T \vec{v}=0}{\arg\min} \vec{v}^T L \vec{v}$$

If $\vec{v}_{n-1}$ were in $\{-1, 1\}^n$ it would have:

· $\vec{v}_{n-1}^T L \vec{v}_{n-1} = cut(S, T)$ as small as possible given that $\vec{v}_{n-1}^T \vec{1} = |T| - |S| = 0$.

· I.e., $\vec{v}_{n-1}$ would indicate the smallest perfectly balanced cut.

· The eigenvector $\vec{v}_{n-1} \in \mathbb{R}^n$ is not generally binary, but still satisfies a 'relaxed' version of this property.

---

*$n$: number of nodes in graph, $\mathbf{A} \in \mathbb{R}^{n \times n}$: adjacency matrix, $\mathbf{D} \in \mathbb{R}^{n \times n}$: diagonal degree matrix, $\mathbf{L} \in \mathbb{R}^{n \times n}$: Laplacian matrix $\mathbf{L} = \mathbf{A} - \mathbf{D}$. $S, T$: vertex sets on different sides of cut.*

Find a good partition of the graph by computing

$$\vec{v}_{n-1} = \underset{v \in \mathbb{R}^n \text{ with } \|\vec{v}\|=1, \ \vec{v}^T \vec{1}=0}{\arg \min} \vec{v}^T L \vec{V}$$

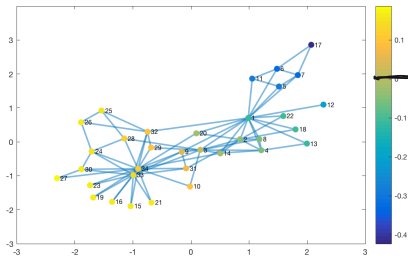Set $S$ to be all nodes with $\vec{v}_{n-1}(i) < 0$, $T$ to be all with $\vec{v}_{n-1}(i) \geq 0$.

---

$n$: number of nodes in graph, $\mathbf{A} \in \mathbb{R}^{n \times n}$: adjacency matrix, $\mathbf{D} \in \mathbb{R}^{n \times n}$: diagonal degree matrix, $\mathbf{L} \in \mathbb{R}^{n \times n}$: Laplacian matrix $\mathbf{L} = \mathbf{A} - \mathbf{D}$. $S, T$: vertex sets on different sides of cut.

Find a good partition of the graph by computing

$$\vec{v}_{n-1} = \underset{v \in \mathbb{R}^n \text{ with } \|\vec{v}\|=1,\ \vec{v}^T\vec{1}=0}{\arg\min} \vec{v}^T L \vec{V}$$

Set $S$ to be all nodes with $\vec{v}_{n-1}(i) < 0$, $T$ to be all with $\vec{v}_{n-1}(i) \geq 0$.



---

$n$: number of nodes in graph, $A \in \mathbb{R}^{n \times n}$: adjacency matrix, $D \in \mathbb{R}^{n \times n}$: diagonal degree matrix, $L \in \mathbb{R}^{n \times n}$: Laplacian matrix $L = A - D$. $S, T$: vertex sets on different sides of cut.

7

Find a good partition of the graph by computing

$$\vec{v}_{n-1} = \underset{v \in \mathbb{R}^n \text{ with } \|\vec{v}\|=1, \ \vec{v}^T\vec{1}=0}{\arg\min} \vec{v}^T L \vec{v}$$

Set $S$ to be all nodes with $\vec{v}_{n-1}(i) < 0$, $T$ to be all with $\vec{v}_{n-1}(i) \geq 0$.



*n*: number of nodes in graph, $A \in \mathbb{R}^{n \times n}$: adjacency matrix, $D \in \mathbb{R}^{n \times n}$: diagonal degree matrix, $L \in \mathbb{R}^{n \times n}$: Laplacian matrix $L = A - D$. $S, T$: vertex sets on different sides of cut.

7

The Shi-Malik normalized cuts algorithm is one of the most commonly used variants of this approach, using the normalized Laplacian $\overline{L} = D^{-1/2}LD^{-1/2}$.
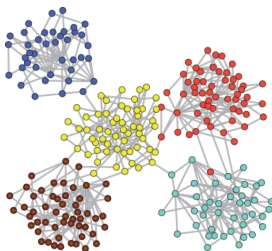
$n$: number of nodes in graph, $A \in \mathbb{R}^{n \times n}$: adjacency matrix, $D \in \mathbb{R}^{n \times n}$: diagonal degree matrix, $L \in \mathbb{R}^{n \times n}$: Laplacian matrix $L = A - D$.

The Shi-Malik normalized cuts algorithm is one of the most commonly used variants of this approach, using the normalized Laplacian $\overline{\mathsf{L}} = \mathsf{D}^{-1/2}\mathsf{L}\mathsf{D}^{-1/2}$.

**Important Consideration:** What to do when we want to split the graph into more than two parts?



$n$: number of nodes in graph, $\mathsf{A} \in \mathbb{R}^{n \times n}$: adjacency matrix, $\mathsf{D} \in \mathbb{R}^{n \times n}$: diagonal degree matrix, $\mathsf{L} \in \mathbb{R}^{n \times n}$: Laplacian matrix $\mathsf{L} = \mathsf{A} - \mathsf{D}$.

The Shi-Malik normalized cuts algorithm is one of the most commonly used variants of this approach, using the normalized Laplacian $\overline{\mathsf{L}} = \mathsf{D}^{-1/2}\mathsf{L}\mathsf{D}^{-1/2}$.

**Important Consideration:** What to do when we want to split the graph into more than two parts?

**Spectral Clustering:**

---

*n*: number of nodes in graph, $\mathsf{A} \in \mathbb{R}^{n \times n}$: adjacency matrix, $\mathsf{D} \in \mathbb{R}^{n \times n}$: diagonal degree matrix, $\mathsf{L} \in \mathbb{R}^{n \times n}$: Laplacian matrix $\mathsf{L} = \mathsf{A} - \mathsf{D}$.

The Shi-Malik normalized cuts algorithm is one of the most commonly used variants of this approach, using the normalized Laplacian $\overline{L} = D^{-1/2}LD^{-1/2}$.

**Important Consideration:** What to do when we want to split the graph into more than two parts?

### Spectral Clustering:

· Compute smallest $k$ nonzero eigenvectors $\underbrace{\vec{v}_{n-1}}, \ldots, \underbrace{\vec{v}_{n-k}}$ of $\overline{L}$.

---

*n*: number of nodes in graph, $A \in \mathbb{R}^{n \times n}$: adjacency matrix, $D \in \mathbb{R}^{n \times n}$: diagonal degree matrix, $L \in \mathbb{R}^{n \times n}$: Laplacian matrix $L = A - D$.

8

The Shi-Malik normalized cuts algorithm is one of the most commonly used variants of this approach, using the normalized Laplacian $\overline{L} = D^{-1/2}LD^{-1/2}$.

**Important Consideration:** What to do when we want to split the graph into more than two parts?

**Spectral Clustering:**

"spectral embedding"

- Compute smallest $k$ nonzero eigenvectors $\vec{v}_{n-1}, \ldots, \vec{v}_{n-k}$ of $\overline{L}$.
- Represent each node by its corresponding row in $V \in \mathbb{R}^{n \times k}$ whose rows are $\vec{v}_{n-1}, \ldots \vec{v}_{n-k}$.

columns

---

$n$: number of nodes in graph, $A \in \mathbb{R}^{n \times n}$: adjacency matrix, $D \in \mathbb{R}^{n \times n}$: diagonal degree matrix, $L \in \mathbb{R}^{n \times n}$: Laplacian matrix $L = A - D$.

8

The Shi-Malik normalized cuts algorithm is one of the most commonly used variants of this approach, using the normalized Laplacian $\overline{\mathbf{L}} = \mathbf{D}^{-1/2}\mathbf{L}\mathbf{D}^{-1/2}$.

$L = D - A$

$L = I \cdot d - A$

Important Consideration: What to do when we want to split the graph into more than two parts?

Spectral Clustering:

· Compute smallest $k$ nonzero eigenvectors $\vec{v}_{n-1}, \ldots, \vec{v}_{n-k}$ of $\overline{\mathbf{L}}$.
· Represent each node by its corresponding row in $\mathbf{V} \in \mathbb{R}^{n \times k}$ whose rows are $\vec{v}_{n-1}, \ldots \vec{v}_{n-k}$.    $V_1, V_2 \ldots V_k$
· Cluster these rows using $k$-means clustering (or really any clustering method).

---

*n*: number of nodes in graph, $\mathbf{A} \in \mathbb{R}^{n \times n}$: adjacency matrix, $\mathbf{D} \in \mathbb{R}^{n \times n}$: diagonal degree matrix, $\mathbf{L} \in \mathbb{R}^{n \times n}$: Laplacian matrix $\mathbf{L} = \mathbf{A} - \mathbf{D}$.

The smallest eigenvectors of $L = D - A$ give the orthogonal 'functions' that are smoothest over the graph. I.e., minimize

$$\vec{v}^T L \vec{v} = \sum_{(i,j) \in E} [\vec{v}(i) - \vec{v}(j)]^2.$$

## LAPLACIAN EMBEDDING

The smallest eigenvectors of $L = D - A$ give the orthogonal 'functions' that are smoothest over the graph. I.e., minimize
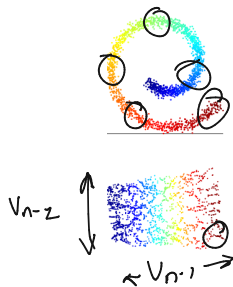
$$\vec{v}^T L \vec{v} = \sum_{(i,j) \in E} [\vec{v}(i) - \vec{v}(j)]^2.$$

Embedding points with coordinates given by $[\vec{v}_{n-1}(j), \vec{v}_{n-2}(j), \ldots, \vec{v}_{n-k}(j)]$ ensures that coordinates connected by edges have minimum total squared Euclidean distance.

The smallest eigenvectors of $L = D - A$ give the orthogonal 'functions' that are smoothest over the graph. I.e., minimize

$$\vec{v}^T L \vec{v} = \sum_{(i,j) \in E} [\vec{v}(i) - \vec{v}(j)]^2.$$

Embedding points with coordinates given by $[\vec{v}_{n-1}(j), \vec{v}_{n-2}(j), \ldots, \vec{v}_{n-k}(j)]$ ensures that coordinates connected by edges have minimum total squared Euclidean distance.

The smallest eigenvectors of $\mathbf{L} = \mathbf{D} - \mathbf{A}$ give the orthogonal 'functions' that are smoothest over the graph. I.e., minimize
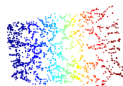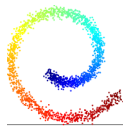
$$\vec{v}^T L \vec{v} = \sum_{(i,j) \in E} [\vec{v}(i) - \vec{v}(j)]^2.$$

$$v_{n-1} L v_{n-1} + v_{n-2} L v_{n-2}^t \cdots$$
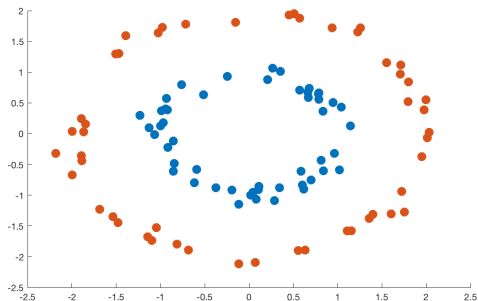$$= \sum_{(i,j) \in E} \| x(i) - x(j) \|_2^2$$

Embedding points with coordinates given by
$x(j) [\vec{v}_{n-1}(j), \vec{v}_{n-2}(j), \ldots, \vec{v}_{n-k}(j)]$ ensures that coordinates connected by edges have minimum total squared Euclidean distance.



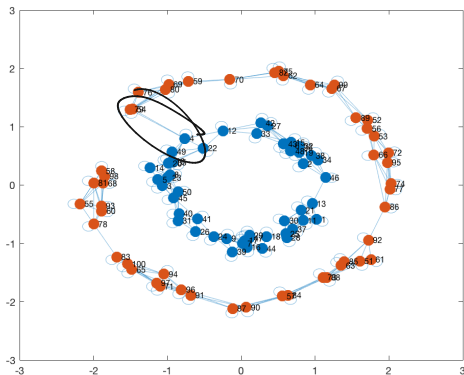- Spectral Clustering
- Laplacian Eigenmaps
- Locally linear embedding
- Isomap
- Etc…

Original Data: (not linearly separable)

$k$-Nearest Neighbors Graph:
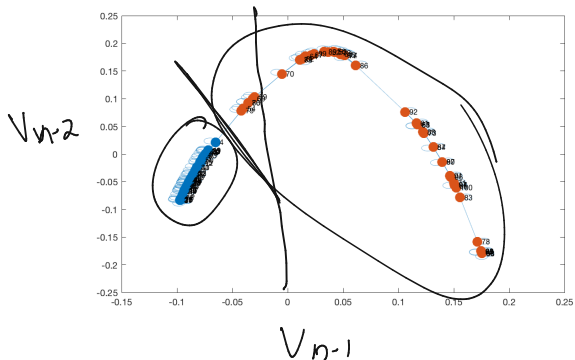
Embedding with eigenvectors $\vec{v}_{n-1}, \vec{v}_{n-2}$: (linearly separable)



$V_{n-2}$

$V_{n-1}$

So Far: Have argued that spectral clustering partitions a graph effectively, along a small cut that separates the graph into large pieces.

So Far: Have argued that spectral clustering partitions a graph effectively, along a small cut that separates the graph into large pieces.

· Haven't given any formal guarantee on the 'quality' of the partitioning.

So Far: Have argued that spectral clustering partitions a graph effectively, along a small cut that separates the graph into large pieces.

· Haven't given any formal guarantee on the 'quality' of the partitioning.
· This is difficult to do for general input graphs.

**So Far:** Have argued that spectral clustering partitions a graph effectively, along a small cut that separates the graph into large pieces.

- Haven't given any formal guarantee on the 'quality' of the partitioning.
- This is difficult to do for general input graphs.

**Common Approach:** Give a natural generative model for random inputs and analyze how the algorithm performs on inputs drawn from this model.

**So Far:** Have argued that spectral clustering partitions a graph effectively, along a small cut that separates the graph into large pieces.

- Haven't given any formal guarantee on the 'quality' of the partitioning.
- This is difficult to do for general input graphs.

**Common Approach:** Give a natural generative model for random inputs and analyze how the algorithm performs on inputs drawn from this model.

- Very common in algorithm design for data analysis/machine learning (can be used to justify $\ell_2$ linear regression, $k$-means clustering, PCA, etc.)

11

**Stochastic Block Model (Planted Partition Model):** Let $G_n(p, q)$ be a distribution over graphs on $n$ nodes, split equally into two groups $B$ and $C$, each with $n/2$ nodes.
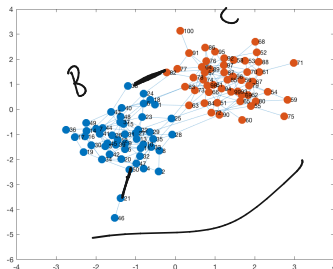
**Stochastic Block Model (Planted Partition Model):** Let $G_n(p, q)$ be a distribution over graphs on $n$ nodes, split equally into two groups $B$ and $C$, each with $n/2$ nodes.

- Any two nodes in the same group are connected with probability $p$ (including self-loops).
- Any two nodes in different groups are connected with prob. $q < p$.
- Connections are independent.

**Stochastic Block Model (Planted Partition Model):** Let $G_n(p, q)$ be a distribution over graphs on $n$ nodes, split equally into two groups $B$ and $C$, each with $n/2$ nodes.

- Any two nodes in the same group are connected with probability $p$ (including self-loops).
- Any two nodes in different groups are connected with prob. $q < p$.
- Connections are independent.

Let $G$ be a stochastic block model graph drawn from $G_n(p, q)$.

$G_n(p, q)$: stochastic block model distribution. $B, C$: groups with $n/2$ nodes each. Connections are independent with probability $p$ between nodes in the same group, and probability $q$ between nodes not in the same group.

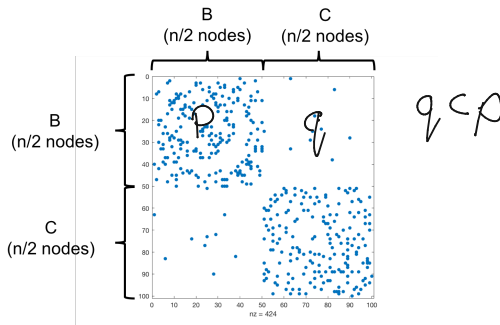Let $G$ be a stochastic block model graph drawn from $G_n(p, q)$.

· Let $A \in \mathbb{R}^{n \times n}$ be the adjacency matrix of $G$.

> $G_n(p, q)$: stochastic block model distribution. $B, C$: groups with $n/2$ nodes each. Connections are independent with probability $p$ between nodes in the same group, and probability $q$ between nodes not in the same group.

Let *G* be a stochastic block model graph drawn from $G_n(p, q)$.

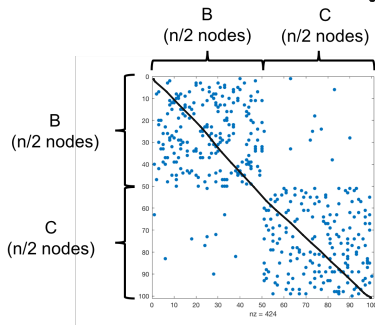- Let $A \in \mathbb{R}^{n \times n}$ be the adjacency matrix of *G*.



$G_n(p, q)$: stochastic block model distribution. *B*, *C*: groups with $n/2$ nodes each. Connections are independent with probability *p* between nodes in the same group, and probability *q* between nodes not in the same group.

Let *G* be a stochastic block model graph drawn from $G_n(p, q)$.

- Let $A \in \mathbb{R}^{n \times n}$ be the adjacency matrix of *G*. What is $\mathbb{E}[A]$?



$G_n(p, q)$: stochastic block model distribution. *B*, *C*: groups with $n/2$ nodes each. Connections are independent with probability *p* between nodes in the same group, and probability *q* between nodes not in the same group.

13

Letting $G$ be a stochastic block model graph drawn from $G_n(p, q)$ and $\mathbf{A} \in \mathbb{R}^{n \times n}$ be its adjacency matrix. What is $\mathbb{E}[\mathbf{A}]$? $= \overline{A}$



$$B \begin{bmatrix} P & q \\ q & P \end{bmatrix}$$

$\overline{A}_{ij} = \mathbb{E} A_{ij}$    $i, j$ in same group

$= P$

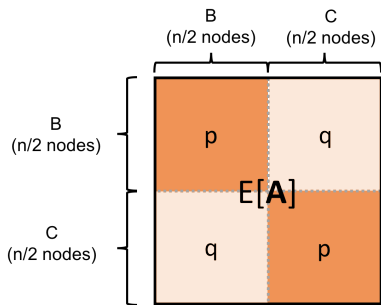$\overline{A}_{ij} = \mathbb{E} A_{ij} = q$    if $i, j$ in diff. groups

---

$G_n(p, q)$: stochastic block model distribution. $B, C$: groups with $n/2$ nodes each. Connections are independent with probability $p$ between nodes in the same group, and probability $q$ between nodes not in the same group.
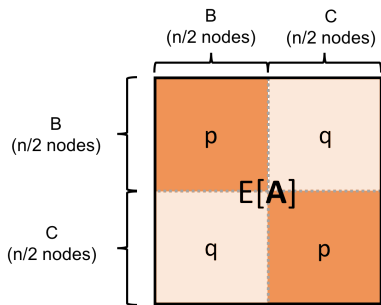
14

Letting $G$ be a stochastic block model graph drawn from $G_n(p, q)$ and $\mathbf{A} \in \mathbb{R}^{n \times n}$ be its adjacency matrix. $(\mathbb{E}[\mathbf{A}])_{i,j} = p$ for $i, j$ in same group, $(\mathbb{E}[\mathbf{A}])_{i,j} = q$ otherwise.



$G_n(p, q)$: stochastic block model distribution. $B, C$: groups with $n/2$ nodes each. Connections are independent with probability $p$ between nodes in the same group, and probability $q$ between nodes not in the same group.

Letting $G$ be a stochastic block model graph drawn from $G_n(p, q)$ and $\mathbf{A} \in \mathbb{R}^{n \times n}$ be its adjacency matrix. $(\mathbb{E}[\mathbf{A}])_{i,j} = p$ for $i, j$ in same group, $(\mathbb{E}[\mathbf{A}])_{i,j} = q$ otherwise.



What are the eigenvectors and eigenvalues of $\mathbb{E}[\mathbf{A}]$?

$G_n(p, q)$: stochastic block model distribution. $B, C$: groups with $n/2$ nodes each. Connections are independent with probability $p$ between nodes in the same group, and probability $q$ between nodes not in the same group.
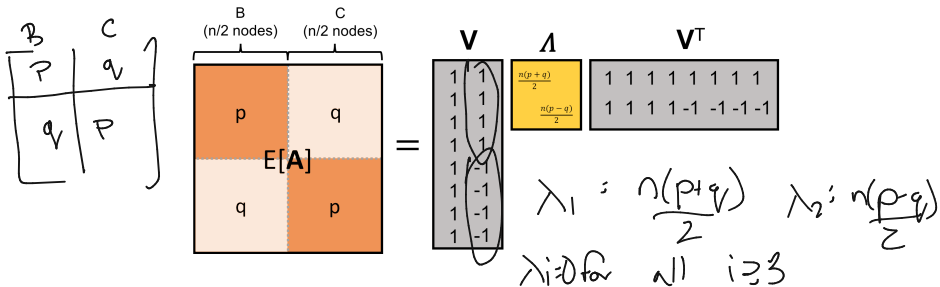
Letting $G$ be a stochastic block model graph drawn from $G_n(p, q)$ and $A \in \mathbb{R}^{n \times n}$ be its adjacency matrix, what are the eigenvectors and eigenvalues of $\mathbb{E}[A]$?



$$\overline{A} = \mathbb{E}A$$

$$\begin{bmatrix} P & q \\ \hline q & P \end{bmatrix} \begin{bmatrix} v_1 \\ \vdots \end{bmatrix} = \overline{A}v_1 = \begin{bmatrix} (p+q)\frac{n}{2} \\ (p+q)\frac{n}{2} \\ \vdots \\ (p+q)\frac{n}{2} \end{bmatrix} = \frac{(p+q)n}{2} \cdot v_1 \qquad \lambda_1$$

$$\begin{bmatrix} P & q \\ \hline q & P \end{bmatrix} \begin{bmatrix} v_2 \\ \vdots \\ - \end{bmatrix} = \overline{A}v_2 = \begin{bmatrix} (p-q)\frac{n}{2} \\ \vdots \\ (p-q)\frac{n}{2} \\ (q-p)\frac{n}{2} \\ \vdots \\ (q-p)\frac{n}{2} \end{bmatrix} = \frac{(p-q)n}{2} \cdot v_2 \qquad \lambda_2$$

16

Letting $G$ be a stochastic block model graph drawn from $G_n(p, q)$ and $A \in \mathbb{R}^{n \times n}$ be its adjacency matrix, what are the eigenvectors and eigenvalues of $\mathbb{E}[A]$?

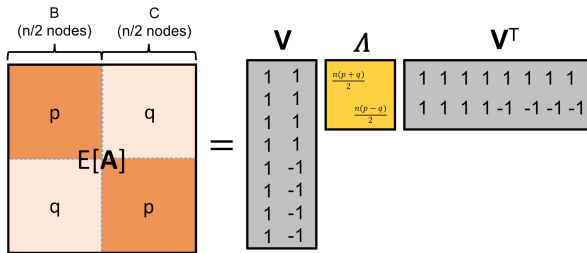If we compute $\vec{v}_2$ then we recover the communities $B$ and $C$!

$$\lambda_1 = \frac{n(p+q)}{2} \qquad \lambda_2 = \frac{n(p-q)}{2}$$

$$\lambda_i = 0 \text{ for all } i \geq 3$$

$$V_1 = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ \vdots \\ 1 \end{pmatrix} \qquad V_2 = \begin{pmatrix} 1 \\ -1 \\ -1 \\ \vdots \\ 1 \end{pmatrix}$$

If we compute $\vec{v}_2$ then we recover the communities *B* and *C*!

- Can show that for $G \sim G_n(p, q)$, **A** is close to $\mathbb{E}[\mathbf{A}]$ with high probability.
- Thus, the true second eigenvector of *A* is close to $[1, 1, 1, \ldots, -1, -1, -1]$ and gives a good estimate of the communities.

Letting $G$ be a stochastic block model graph drawn from $G_n(p, q)$, $A \in \mathbb{R}^{n \times n}$ be its adjacency matrix and $L$ be its Laplacian, what are the eigenvectors and eigenvalues of $\mathbb{E}[L]$?

$$\mathbb{E}[L] = \mathbb{E}[D - A] = \mathbb{E}D - \left[\begin{array}{c|c} p & q \\ \hline q & p \end{array}\right]$$

$$\left[ p \binom{n}{2} + q \binom{n}{2} \right] I - \left[\begin{array}{c|c} p & q \\ \hline q & p \end{array}\right]$$

$$\overline{L} = \mathbb{E}L, \quad \overline{A} = \mathbb{E}A$$

$$\overline{L} = \frac{(p+q)}{2}n \, I - \overline{A} \qquad v_i \text{ which is eigenvector of } \overline{A}$$

$$\overline{L} v_i = \frac{(p+q)}{2}n \, I v_i - \overline{A} v_i = \left[ \frac{(p+q)n}{2} - \lambda_i(A) \right] v_i$$

19

Letting $G$ be a stochastic block model graph drawn from $G_n(p, q)$, $A \in \mathbb{R}^{n \times n}$ be its adjacency matrix and $L$ be its Laplacian, what are the eigenvectors and eigenvalues of $\mathbb{E}[L]$?

Questions?