# CS-GY 9223 I: Lecture 10
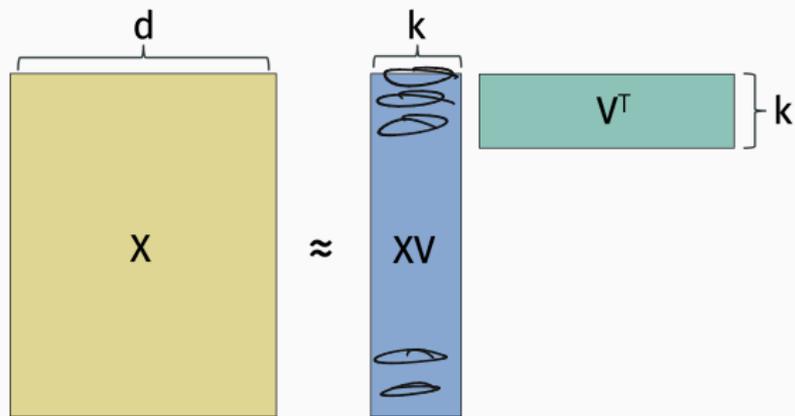# Spectral clustering, spectral graph theory.

NYU Tandon School of Engineering, Prof. Christopher Musco

- Project proposal feedback.
- Problem set.
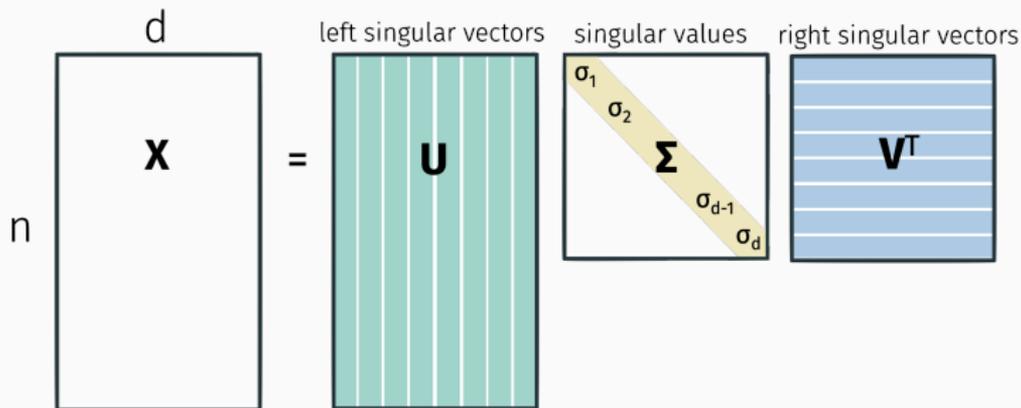
Write X as a rank *k* factorization by projecting onto the subspace spanned by an orthonormal matrix $V \in \mathbb{R}^{d \times k}$

One-stop shop for computing optimal low-rank approximations.
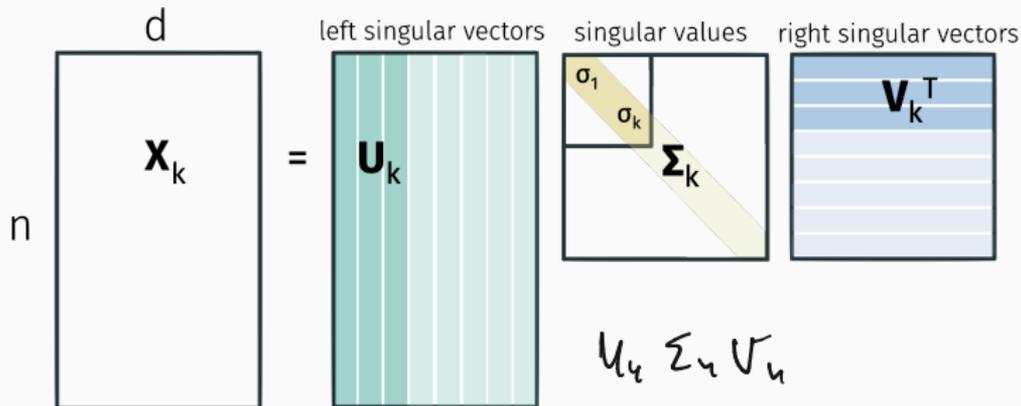
Any matrix **X** can be written:



Where $U^T U = I$, $V^T V = I$, and $\sigma_1 \geq \sigma_2 \geq \ldots \sigma_d \geq 0$.

Can read off optimal low-rank approximations from the SVD:



$$X_k = U_k U_k^T X = \boxed{X V_k V_k^T.}$$

$$V_k = \underset{\text{orthonormal } V \in \mathbb{R}^{d \times k}}{\arg\min} \|X - X V V^T\|_F^2 = \underset{\text{orthonormal } V \in \mathbb{R}^{d \times k}}{\arg\max} \|X V V^T\|_F^2$$

### Theorem (Power Method Convergence)

*Let $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$ be parameter capturing the "gap" between the first and second largest singular values. If Power Method is initialized with a random Gaussian vector then, with high probability, after $T = O\left(\frac{\log d/\epsilon}{\gamma}\right)$ steps, we have:*

$$\|\mathbf{v}_1 - \mathbf{z}^{(T)}\|_2 \leq \epsilon.$$

Total runtime: $O(T \cdot \text{nnz}(\mathbf{X})) \leq O(T \cdot nd)$

number of non-zeros

Block power method:

$$Z^T Z = I$$

- Choose $G \in \mathbb{R}^{d \times k}$ be a random Gaussian matrix.
- $Z_0 = \text{orth}(G)$.

$$g r (Z)$$

- For $i = 1, \ldots, T$
    - $Z^{(i)} = X^T \cdot (XZ^{(i-1)})$
    - $Z^{(i)} = \text{orth}(Z^{(i)})$

    Return $Z^{(T)}$

$$Z^{(T)} \approx V_k$$

$$\sim O\left(\frac{1}{\epsilon}\right)$$

**Convergence Guarantee**: $T = O\left(\frac{\log d/\epsilon}{\epsilon}\right)$ iterations to obtain a nearly optimal low-rank approximation:

$$\|A - AZZ^T\|_F^2 \leq (1 + \epsilon)\|A - AV_kV_k^T\|_F^2.$$

**Runtime**: $O(\text{nnz}(X) \cdot k \cdot T) \leq O(ndk \cdot T)$.

number of nonzeros in $X$

7

Possible to "accelerate" these methods. $\frac{1}{\sqrt{\epsilon}}$

**Convergence Guarantee**: $T = O\left(\frac{\log d/\epsilon}{\sqrt{\epsilon}}\right)$ iterations to obtain a nearly optimal low-rank approximation:

$$\|A - AZZ^T\|_F^2 \leq (1 + \epsilon)\|A - AV_kV_k^T\|_F^2.$$

**Runtime**: $O(\text{nnz}(X) \cdot k \cdot T) \leq O(ndk \cdot T)$.

Corpus of Documents

Term Document Matrix X

Low-Rank Approximation via SVD

X ≈ Y Z

**Corpus of Documents**

**Term Document Matrix X**

**Low-Rank Approximation via SVD**

$X \approx Y$

$z_a$ is *(column)* of $Z$

$y_i$ row of $Y$

- $\langle \vec{y_i}, \vec{z_a} \rangle \approx 1$ when $doc_i$ contains $word_a$.
- If $doc_i$ and $doc_j$ both contain $word_a$, $\langle \vec{y_i}, \vec{z_a} \rangle \approx \langle \vec{y_j}, \vec{z_a} \rangle = 1$.

**Term Document Matrix X**

**Low-Rank Approximation via SVD**

$X \approx Y \quad Z$

- The columns $\vec{z}_1, \vec{z}_2, \ldots$ give representations of words, with $\vec{z}_i$ and $\vec{z}_j$ tending to have high dot product if $word_i$ and $word_j$ appear in many of the same documents.

- Z corresponds to the top $k$ right singular vectors: the eigenvectors of $X^T X$. Intuitively, what is $XX^T$?

- $(X^T X)_{i,j} =$ $= \langle x_i, x_j \rangle$ for $i, j$ column of $X$

Not obvious how to convert a word into a feature vector that captures the meaning of that word. Approach suggested by LSA: build a $d \times d$ symmetric "similarity matrix" $M$ between words, and factorize: $M \approx FF^T$ for rank $k$ $F$.

- **Similarity measures:** How often do $word_i, word_j$ appear in the same sentence, in the same window of $w$ words, in similar positions of documents in different languages?
- Replacing $XX^T$ with these different metrics (sometimes appropriately transformed) leads to popular word embedding algorithms: `word2vec`, `GloVe`, etc.

`word2vec` was originally described as a neural-network method, but Levy and Goldberg show that it is simply low-rank approximation of a specific similarity matrix. *(Neural word embedding as implicit matrix factorization.)*

Main idea: Understand graph data by constructing natural matrix representations, and studying that matrix's spectrum (eigenvalues/eigenvectors).



$1, \cdots 8 \quad \to \quad (1,2), (2,5) \cdots (7,8)$

For now assume $G = (V, E)$ is an undirected, unweighted graph with *n* nodes.

Two most common representations: $n \times n$ <u>adjacency matrix</u> $\mathbf{A}$ and <u>graph Laplacian</u> $\mathbf{L} = \mathbf{D} - \mathbf{A}$ where $\mathbf{D}$ is the diagonal degree matrix.



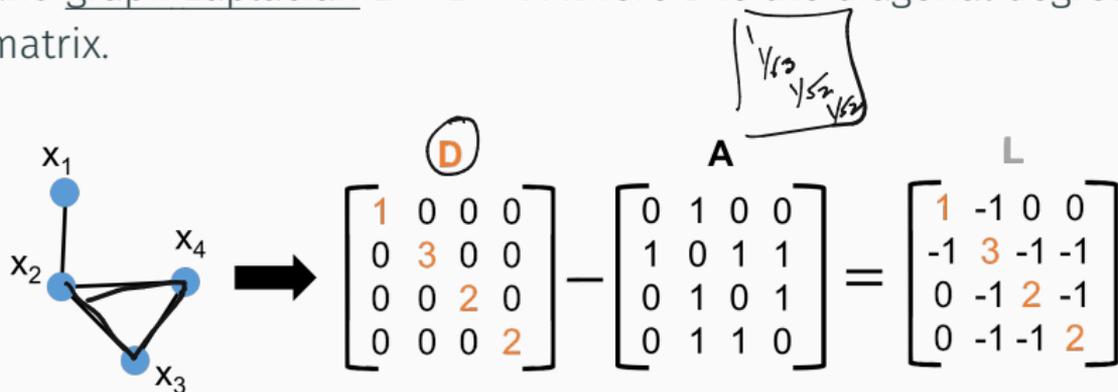$$\mathbf{D} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 2 \end{bmatrix} - \mathbf{A} \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{bmatrix} = \mathbf{L} \begin{bmatrix} 1 & -1 & 0 & 0 \\ -1 & 3 & -1 & -1 \\ 0 & -1 & 2 & -1 \\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

Also common to look at normalized versions of both of these: $\bar{\mathbf{A}} = \mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$ and $\bar{\mathbf{L}} = \mathbf{I} - \mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$.

- If **L** have $k$ eigenvalues equal to 0, then $G$ has $k$ connected components.
- Sum of cubes of **A**'s eigenvalues is equal to number of triangles in the graph.
- Sum of eigenvalues to the power $q$ is the number of $q$ cycles.

$$\left\{ B^T B \right\}_{ii} = \left[ B^{(i)\,T} B^{(i)} \right] = \| B^{(i)} \|_2^2$$

**D**        **A**        **L**

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

$$= \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$L = B^T B$ where $B$ is the signed "edge-vertex incidence" matrix.

$B =$

$1 \ldots n$

$B \in \mathbb{B}^{m \times n}$

$x_1 - x_2$

$= \rightarrow x_2 - x_3$

$x$

$m$ edges

$B^T B \quad (n \times m) \times (m \times n)$

$= (n \times n)$

17

$$PSD \qquad x^T B^T B x = x^T L x$$
$$= \| Bx \|_2^2$$

$$B^T B = b_1 b_1^T + b_2 b_2^T + \ldots + b_m b_m^T,$$

where $b_i$ is the $i^{th}$ row of $B$ (each row corresponds to a single edge).



$$= -1 \cdot \overline{1} \cdot b_1 b_1^T$$
$$= b_1 b_1^T$$

Conclusions from $L = B^T B$

- $L$ is positive semidefinite: $x^T L x \geq 0$ <u>for all</u> x.

- $L = \underline{V \Sigma^2 V^T}$ where $U \Sigma V^T$ is B's SVD. Columns of V are <u>eigenvectors</u> of L.

$$V \Sigma U^T U \Sigma V^T$$

- For any vector $x \in \mathbb{R}^n$,

"smoothness of $x$"

$$x^T L x = \sum_{(i,j) \in E} (x(i) - x(j))^2.$$

$x \in \mathbb{R}^n$

$= \|Bx\|_2^2$

$x^T L x$

19

$\widehat{x^T L x}$ is small if x is a "smooth" function with respect to the graph.

$$x^T L x = 0 \quad \text{if} \quad x = \vec{1}$$



Eigenvectors of the Laplacian with small eigenvalues correspond to smooth functions over the graph.

## Courant–Fischer min-max principle

Let $V = [v_1, \ldots, v_n]$ be the eigenvectors of $L$.

$$v_n = \underset{\|v\|=1}{\arg\min} \, v^T L v$$

$$v_{n-1} = \underset{\|v\|=1, v \perp v_n}{\arg\min} \, v^T L v$$

$$v_{n-2} = \underset{\|v\|=1, v \perp v_n, v_{n-1}}{\arg\min} \, v^T L v$$

$$\vdots$$

$$v_1 = \underset{\|v\|=1, v \perp v_n, \ldots, v_2}{\arg\min} \, v^T L v$$

## Courant–Fischer min-max principle

Let $V = [v_1, \ldots, v_n]$ be the eigenvectors of $L$.

$$v_1 = \underset{\|v\|=1}{\arg\max}\ v^T L v$$

$$v_2 = \underset{\|v\|=1, v \perp v_1}{\arg\max}\ v^T L v$$

$$v_3 = \underset{\|v\|=1, v \perp v_1, v_2}{\arg\max}\ v^T L v$$

$$\vdots$$

$$v_n = \underset{\|v\|=1, v \perp v_1, \ldots, v_{n-1}}{\arg\max}\ v^T L v$$
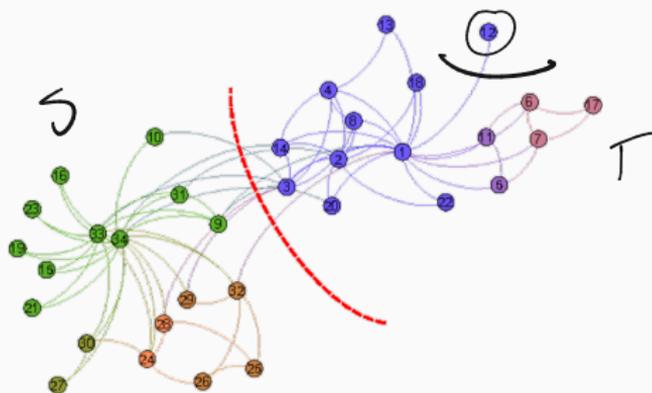
## EXAMPLE APPLICATION OF SPECTRAL GRAPH THEORY

- Study graph partitioning problem important in 1) understanding social networks 2) nonlinear clustering in unsupervised machine learning (spectral clustering).
- See how this problem can be solved approximately using Laplacian eigenvectors.
- Give a full analysis of the method for a common random graph model.
- Use two tools: matrix concentration and eigenvector perturbation bounds.

**Common goal:** Given a graph $G = (V, E)$, partition nodes along a cut that:

- Has few crossing edges: $|\{(u, v) \in E : u \in S, v \in T\}|$ is small.
- Separates large partitions: $|S|, |T|$ are not too small.



(a) Zachary Karate Club Graph

Important in understanding community structure in social networks.

24

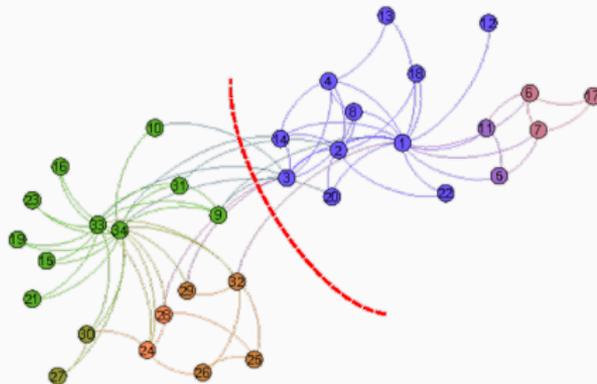Wayne W. Zachary (1977). An Information Flow Model for Conflict and Fission in Small Groups.

"The network captures 34 members of a karate club, documenting links between pairs of members who interacted outside the club. During the study a conflict arose between the administrator "John A" and instructor "Mr. Hi" (pseudonyms), which led to the split of the club into two. Half of the members formed a new club around Mr. Hi; members from the other part found a new instructor or gave up karate. Based on collected data Zachary correctly assigned all but one member of the club to the groups they actually joined after the split." – Wikipedia

Beautiful paper – definitely worth checking out!

**Common goal:** Given a graph $G = (V, E)$, partition nodes along a cut that:

- Has few crossing edges: $|\{(u, v) \in E : u \in S, v \in T\}|$ is small.
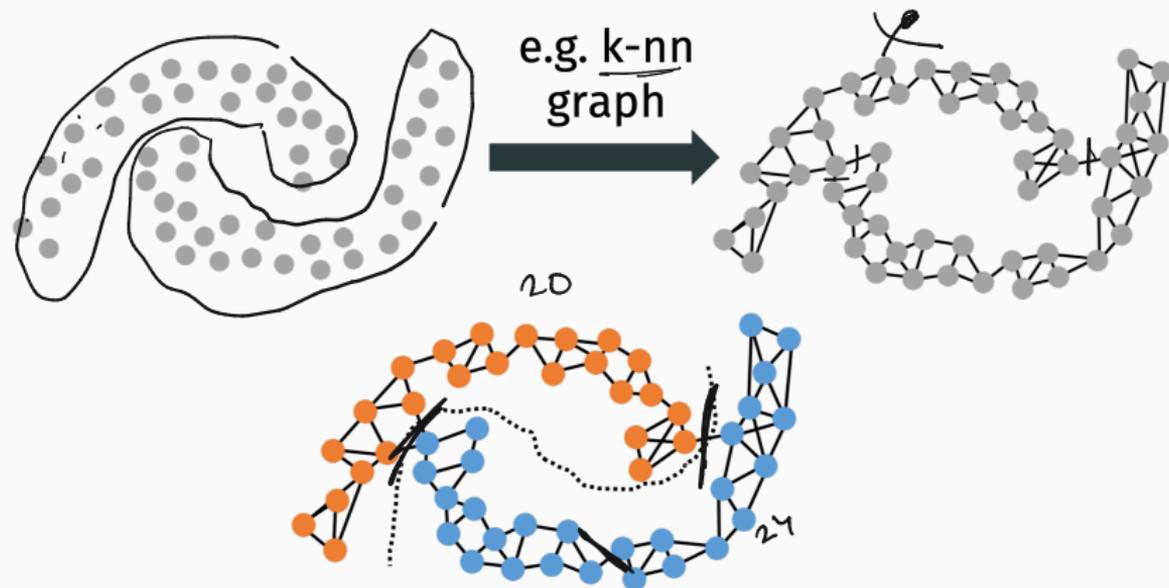- Separates large partitions: $|S|, |T|$ are not too small.



(a) Zachary Karate Club Graph

Important in understanding <u>community structure</u> in social networks.

**Idea:** Construct synthetic graph for data that is hard to cluster.



e.g. k-nn graph

Spectral Clustering, Laplacian Eigenmaps, Locally linear embedding, Isomap, etc.

There are many way's to formalize Zachary's problem:

**Sparsest Cut:**

$$\min_{S} \frac{\text{cut}(S, V \setminus S)}{\min(|S|, |V \setminus S|)}$$

$\rightarrow 3$

$\rightarrow 20$

**$\beta$-Balanced Cut:**

$|S|$

$\beta \leq 1/2$

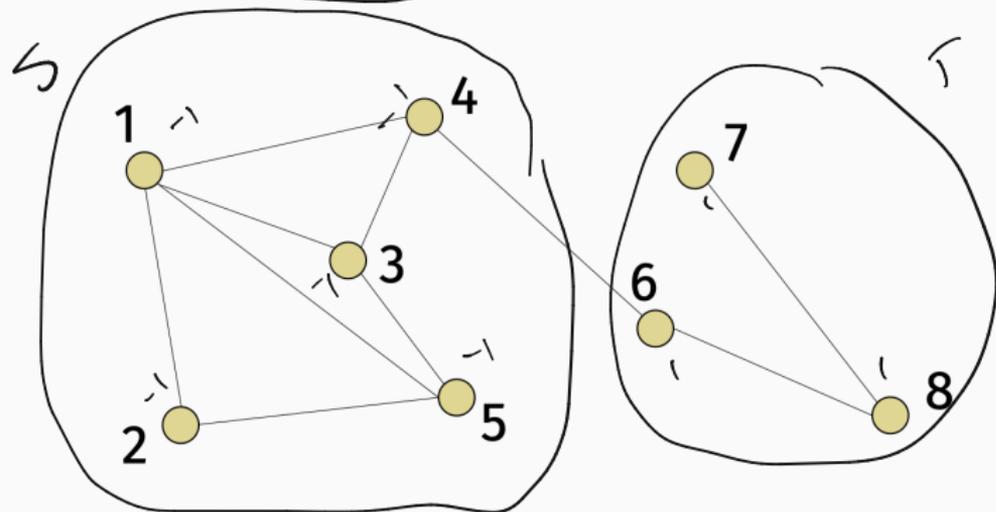$$\min_{S} \text{cut}(S, V \setminus S) \quad \text{such that} \quad \min(|S|, |V \setminus S|) \geq \beta n$$

Most formalizations lead to computationally hard problems. Lots of interest in designing polynomial time approximation algorithms, but tend to be slow. In practice, much simpler methods based on the graph spectrum are used.
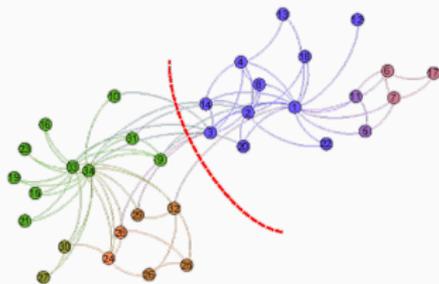
Another conclusion from $L = B^T B$:

For a underline{cut indicator vector} $\mathbf{c} \in \{-1, 1\}^n$ with $\mathbf{c}(i) = -1$ for $i \in S$ and $\mathbf{c}(i) = 1$ for $i \in T = V \setminus S$:

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

(a) Zachary Karate Club Graph

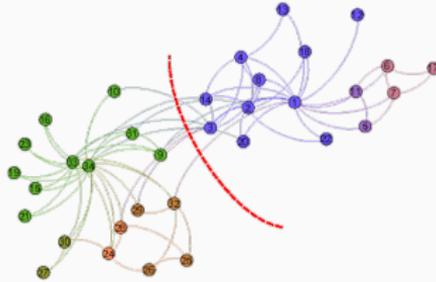For a underline{cut indicator vector} $\mathbf{c} \in \{-1, 1\}^n$ with $\mathbf{c}(i) = -1$ for $i \in S$ and $\mathbf{c}(i) = 1$ for $i \in T$:

- $\mathbf{c}^T L \mathbf{c} = 4 \cdot cut(S, T)$.
- $|\mathbf{c}^T \mathbf{1}| = ||T| - |S||.$

$$c^T 1 = \sum_{i=1}^{m} c_i$$

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

30

(a) Zachary Karate Club Graph

Equivalent formulation if we divide everything by $\sqrt{n}$ so that $\mathbf{c}$ has norm 1. Then $\mathbf{c} \in \{-\frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}}\}^n$ and:

- $\mathbf{c}^T L \mathbf{c} = \frac{4}{n} \cdot cut(S, T)$.

$$\| c \|_2 = 1$$

- $\mathbf{c}^T \mathbf{1} = \frac{1}{\sqrt{n}}(|T| - |S|)$.

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

31

The smallest eigenvector/singular vector $\mathbf{v}_n$ satisfies:

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

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

By Courant-Fischer, $\mathbf{v}_{n-1}$ is given by:

$$\mathbf{v}_{n-1} = \underset{\|\mathbf{v}\|=1,\ \mathbf{v}_n^T\mathbf{v}=0}{\arg\min} \ \mathbf{v}^T L \mathbf{v}$$

If $\mathbf{v}_{n-1}$ were <u>binary</u> $\{-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\}^n$ it would have:

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

- $\mathbf{v}_{n-1}$ would indicate the smallest <u>perfectly balanced</u> cut.

$\mathbf{v}_{n-1} \in \mathbb{R}^n$ is not generally binary, but a natural approach is to 'round' the vector to obtain a cut.

Find a good partition of the graph by computing

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

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

*Handwritten annotations:*

$V \Lambda V^T$

$P \Lambda P^{-1}$

$V_{n-1}^T V_n = 0$

$v_n$

Assume $v$ is an eigenvector:

$Lv = \lambda v$

$v^T \lambda v$
$= \lambda v^T v \rightarrow \|v\|_2 = 1$
$= \lambda$
$\underbrace{\quad}$
$\boxed{\lambda}$

$\|v\|_2 = 1$

$\boxed{v^T L v = 0}$

$\rightarrow v^T B^T B v$
$= \|Bv\|_2^2$
$\geq 0$
$L$ is PSD



$x = 1$

$x^T L x = x^T B^T B x = \sum_{(i,j) \in E} (x(i) - x(j))^2$

Find a good partition of the graph by computing
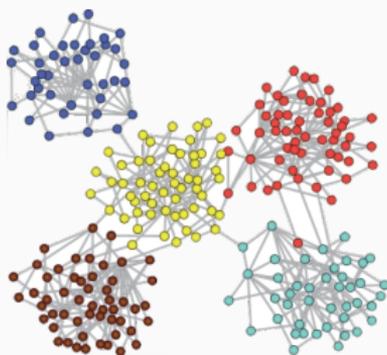
$$\mathbf{v}_{n-1} = \underset{v \in \mathbb{R}^n \text{ with } \|\mathbf{v}\|=1, \ \mathbf{v}^T\mathbf{1}=0}{\arg\min} \mathbf{v}^T L \mathbf{v}$$

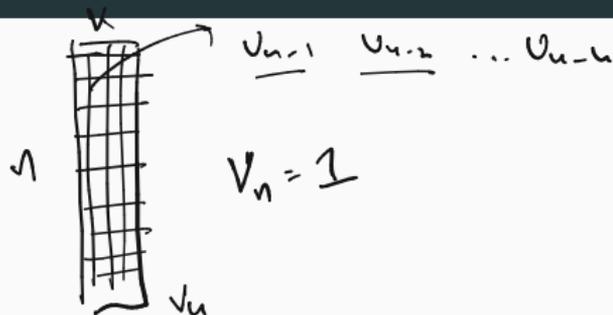Set $S$ to be all nodes with $\mathbf{v}_{n-1}(i) < 0$, and $T$ to be all with $\mathbf{v}_{n-1}(i) \geq 0$.

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?
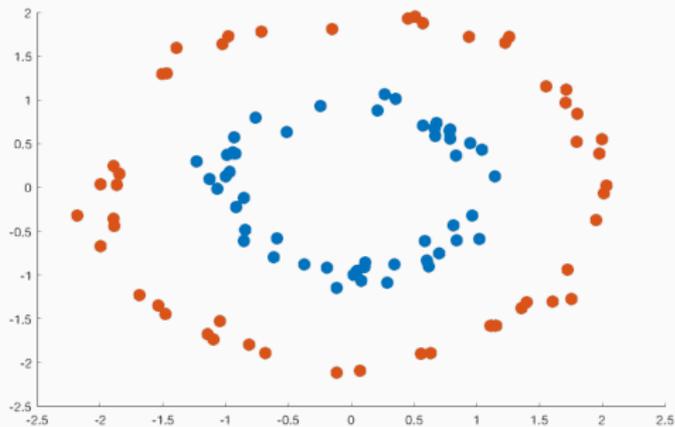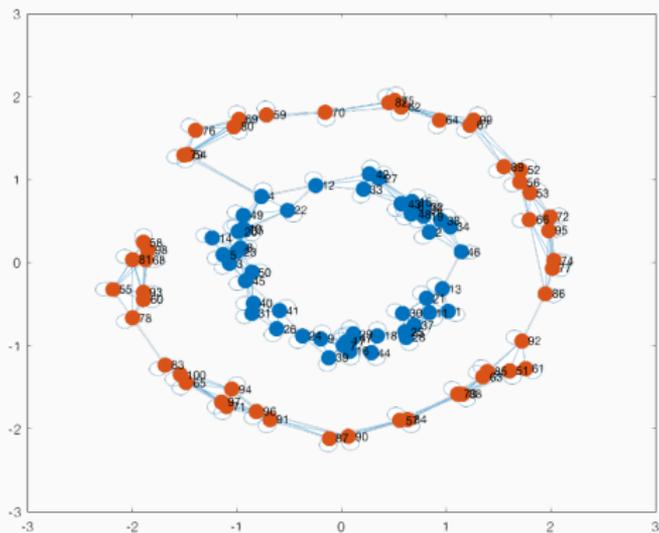
Spectral Clustering:

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

$$\min_{v} \; v^\top L v \qquad \vec{1}^\top L \vec{1} = \sum_{(i,j) \in E} (\vec{1}_i - \vec{1}_j)^2 = 0$$

$$v^\top L v = 0$$

36

**Original Data:** (not linearly separable)

$k$-Nearest Neighbors Graph:

Embedding with eigenvectors $v_{n-1}$, $v_{n-2}$: (linearly separable)
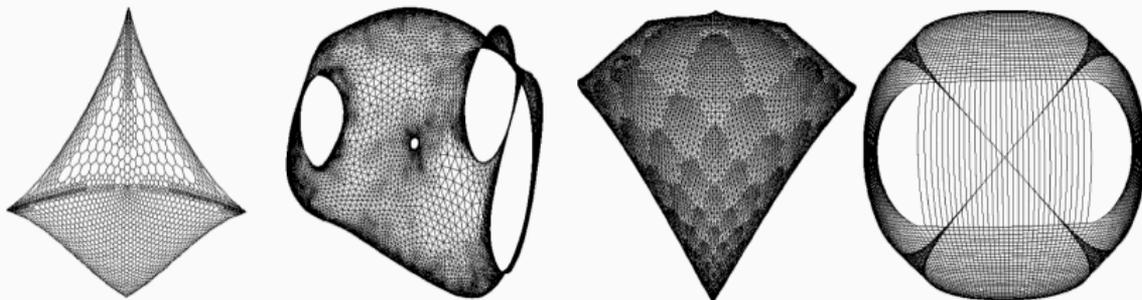
Intuitively, since $\mathbf{v} \in \mathbf{v}_1, \ldots \mathbf{v}_k$ are smooth over the graph,

$$\sum_{i,j \in E} (\mathbf{v}[i] - \mathbf{v}[j])^2$$

is small for each coordinate. I.e. this embedding explicitly encourages nodes connected by an edge to be placed in nearby locations in the embedding.

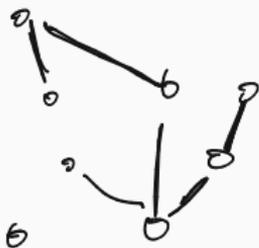

Also useful e.g., in graph drawing.

So far: Showed that spectral clustering partitions a graph along a small cut between large pieces.

- No formal guarantee on the 'quality' of the partitioning.
- Would be difficult to analyze for general input graphs.

Common approach: Design a natural generative model that produces random but realistic inputs and analyze how the algorithm performs on inputs drawn from this model.

- Very common in algorithm design and analysis. Great way to start approaching a problem.
- This is also the whole idea behind Bayesian Machine Learning (can be used to justify $\ell_2$ linear regression, $k$-means clustering, PCA, etc.)

Ideas for a generative model for **social network graphs** that would allow us to understand partitioning?



Erdos - Rengi

$p \in \{0, 1\}$

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.

$\to \in [0, 1]$

- 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 $q < p$.

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

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



B (n/2 nodes)    C (n/2 nodes)

B (n/2 nodes)

C (n/2 nodes)

Note that we are <u>arbitrarily</u> ordering the nodes in A by group.
In reality A would look "scrambled" as on the right.

$$L = D - A \qquad A = D - L \qquad \mathbb{E}[A] = \mathbb{E}[D] - \mathbb{E}[L]$$

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.



$$\begin{pmatrix} \left(\frac{p+b}{2}\right) & & \\ & \ddots & \\ & & \ddots \end{pmatrix}$$

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

45

What is the expected Laplacian $G_n(p, q)$?

A and L have the same eigenvectors and eigenvalue are equal up to a shift.

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

- $\mathbf{v}_1 \sim \mathbf{1}$ with eigenvalue $\lambda_1 = \frac{(p+q)n}{2}$.
- $\mathbf{v}_2 \sim \chi_{B,C}$ with eigenvalue $\lambda_2 = \frac{(p-q)n}{2}$.
- $\chi_{B,C}(i) = 1$ if $i \in B$ and $\chi_{B,C}(i) = -1$ for $i \in C$.

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

**Upshot:** The second smallest eigenvector of $\mathbb{E}[\mathsf{L}]$, equivalently the second largest of $\mathbb{E}[\mathsf{A}]$, is $\chi_{B,C}$ – the indicator vector for the cut between the communities.

- If the random graph $G$ (equivilantly $\mathsf{A}$ and $\mathsf{L}$) were exactly equal to its expectation, partitioning using this eigenvector would exactly recover communities $B$ and $C$.

How do we show that a matrix (e.g., $\mathsf{A}$) is close to its expectation? **Matrix concentration inequalities.**

- Analogous to scalar concentration inequalities like Markovs, Chebyshevs, Bernsteins.

> **Matrix Concentration Inequality:** If $p \geq O\left(\frac{\log^4 n}{n}\right)$, then with high probability
>
> $$\|A - \mathbb{E}[A]\|_2 \leq O(\sqrt{pn}).$$
>
> where $\|\cdot\|_2$ is the matrix spectral norm (operator norm).

For $X \in \mathbb{R}^{n \times d}$, $\|X\|_2 = \max_{z \in \mathbb{R}^d : \|z\|_2 = 1} \|Xz\|_2$.

**Exercise:** Show that $\|X\|_2$ is equal to the largest singular value of $X$. For symmetric $X$ (like $A - \mathbb{E}[A]$) show that it is equal to the magnitude of the largest magnitude eigenvalue.

For the stochastic block model application, we want to show that the second eigenvectors of $A$ and $\mathbb{E}[A]$ are close. How does this relate to their difference in spectral norm?

## Eigenvector Perturbation

> **Davis-Kahan** Eigenvector Perturbation Theorem: Suppose $A, \overline{A} \in \mathbb{R}^{d \times d}$ are symmetric with $\|A - \overline{A}\|_2 \leq \epsilon$ and eigenvectors $v_1, v_2, \ldots, v_d$ and $\overline{v}_1, \overline{v}_2, \ldots, \overline{v}_d$. Letting $\theta(v_i, \overline{v}_i)$ denote the angle between $v_i$ and $\overline{v}_i$, for all $i$:
>
> $$\sin[\theta(v_i, \overline{v}_i)] \leq \frac{\epsilon}{\min_{j \neq i} |\lambda_i - \lambda_j|}$$
>
> eigenvectors of $A$.
>
> where $\lambda_1, \ldots, \lambda_d$ are the eigenvalues of $\overline{A}$.

The error gets larger if there are eigenvalues with similar magnitudes.

$$\|\mathbf{A-\bar{A}}\|_2 = \varepsilon$$

$$\mathbf{A} - \mathbf{\bar{A}} = \begin{bmatrix} \varepsilon & 0 \\ 0 & \varepsilon \end{bmatrix}$$

$$\mathbf{A} = \begin{bmatrix} 1+\varepsilon & 0 \\ 0 & 1 \end{bmatrix} \qquad \mathbf{\bar{A}} = \begin{bmatrix} 1 & 0 \\ 0 & 1+\varepsilon \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \qquad \begin{bmatrix} 0 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

**Claim 1 (Matrix Concentration):** For $p \geq O\left(\frac{\log^4 n}{n}\right)$,

$$\|A - \mathbb{E}[A]\|_2 \leq O(\sqrt{pn}).$$

**Claim 2 (Davis-Kahan):** For $p \geq O\left(\frac{\log^4 n}{n}\right)$,

$$\sin \theta(v_2, \bar{v}_2) \leq \frac{O(\sqrt{pn})}{\min_{j \neq i} |\lambda_i - \lambda_j|} \leq \frac{O(\sqrt{pn})}{(p-q)n/2} = O\left(\frac{\sqrt{p}}{(p-q)\sqrt{n}}\right)$$

**Recall:** $\mathbb{E}[A]$, has eigenvalues $\lambda_1 = \frac{(p+q)n}{2}$, $\lambda_2 = \frac{(p-q)n}{2}$, $\lambda_i = 0$ for $i \geq 3$.
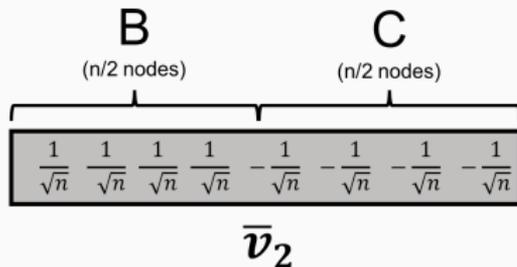
$$\min_{j \neq i} |\lambda_i - \lambda_j| = \min\left(qn, \frac{(p-q)n}{2}\right).$$

Assume $\left|\frac{(p-q)n}{2} - 0\right|$ will be the minimum of the two gaps. I.e. smaller than $\left|\frac{(p+q)n}{2} - \frac{(p-q)n}{2}\right| = qn$.

**So Far:** $\sin\theta(v_2, \bar{v}_2) \leq O\left(\frac{\sqrt{p}}{(p-q)\sqrt{n}}\right)$. What does this give us?
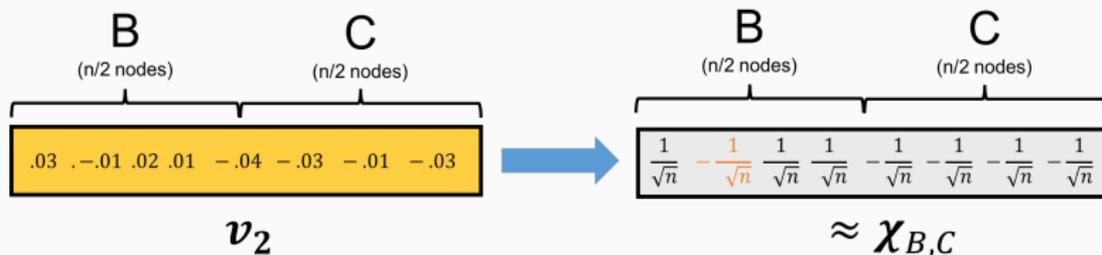
- Can show that this implies $\|v_2 - \bar{v}_2\|_2^2 \leq O\left(\frac{p}{(p-q)^2 n}\right)$ (exercise).

- $\bar{v}_2$ is $\frac{1}{\sqrt{n}}\chi_{B,C}$: the community indicator vector.



- Every $i$ where $v_2(i)$, $\bar{v}_2(i)$ differ in sign contributes $\geq \frac{1}{n}$ to $\|v_2 - \bar{v}_2\|_2^2$.

- So they differ in sign in at most $O\left(\frac{p}{(p-q)^2}\right)$ positions.

**Upshot:** If $G$ is a stochastic block model graph with adjacency matrix $A$, if we compute its second large eigenvector $v_2$ and assign nodes to communities according to the sign pattern of this vector, we will correctly assign all but $O\left(\frac{p}{(p-q)^2}\right)$ nodes.



- Why does the error increase as $q$ gets close to $p$?
- Even when $p - q = O(1/\sqrt{n})$, assign all but an $O(n)$ fraction of nodes correctly. E.g., assign 99% of nodes correctly.