# CS-GY 6763: Lecture 12
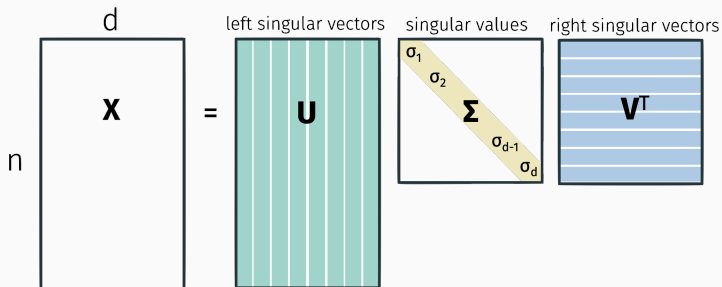# Power Method, Krylov Subspace Methods, Spectral Graph Partitioning

NYU Tandon School of Engineering, Prof. Christopher Musco

One of the most fundamental results in linear algebra.
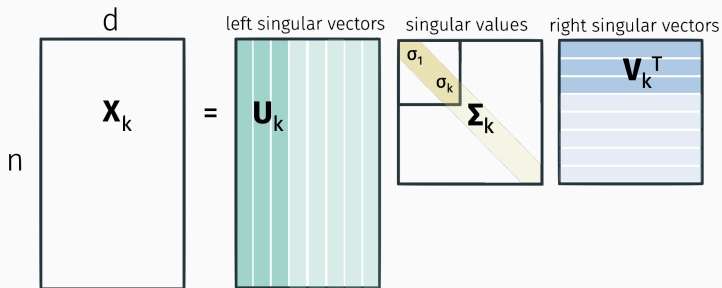
<u>Any</u> matrix $\mathbf{X}$ can be written:



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

Singular values are unique. Factors are not. E.g. would still get a valid SVD by multiplying both $i^{\text{th}}$ column of $\mathbf{V}$ and $\mathbf{U}$ by $-1$.

**Key result:** Can find the best low-rank approximation from the singular value decomposition.



$$\mathbf{U}_k = \underset{\text{orthogonal } \mathbf{Z} \in \mathbb{R}^{d \times k}}{\arg \min} \|\mathbf{X} - \mathbf{Z}\mathbf{Z}^T\mathbf{X}\|_F^2$$

$$\mathbf{V}_k = \underset{\text{orthogonal } \mathbf{W} \in \mathbb{R}^{d \times k}}{\arg \min} \|\mathbf{X} - \mathbf{X}\mathbf{W}\mathbf{W}^T\|_F^2$$

3

Suffices to compute right singular vectors $V$:

- Compute $X^T X$.
- Find eigendecomposition $V \Lambda V^T = X^T X$ using e.g. QR algorithm.
- Compute $L = XV$. Set $\sigma_i = \|L_i\|_2$ and $U_i = L_i / \|L_i\|_2$.

Total runtime $\approx$

How to go faster?

- Compute <u>approximate</u> solution.
- Only compute <u>top $k$ singular vectors/values</u>.
- <u>Iterative algorithms</u> achieve runtime $\approx O(ndk)$ vs. $O(nd^2)$ time.
    - **Krylov subspace methods** like the Lanczos method are most commonly used in practice.
    - **Power method** is the simplest Krylov subspace method, and still works very well.

**Today:** What about when $k = 1$?

**Goal:** Find some $z \approx v_1$.

**Input:** $X \in \mathbb{R}^{n \times d}$ with SVD $U\Sigma V^T$.

Power method:

- Choose $z^{(0)}$ randomly. $z_0 \sim \mathcal{N}(0, 1)$.
- $z^{(0)} = z^{(0)}/\|z^{(0)}\|_2$
- For $i = 1, \ldots, T$
    - $z^{(i)} = X^T \cdot (Xz^{(i-1)})$
    - $n_i = \|z^{(i)}\|_2$
    - $z^{(i)} = z^{(i)}/n_i$

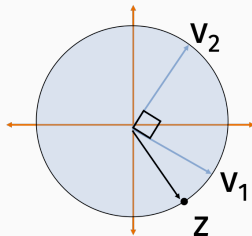    Return $z^{(T)}$

0 iterations     1 iterations     2 iterations

## Theorem (Basic 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 either:*

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

Total runtime: $O\left(nd \cdot \frac{\log d/\epsilon}{\gamma}\right)$

Write $z^{(i)}$ in the right singular vector basis:

$$z^{(0)} = c_1^{(0)}v_1 + c_2^{(0)}v_2 + \ldots + c_d^{(0)}v_d$$
$$z^{(1)} = c_1^{(1)}v_1 + c_2^{(1)}v_2 + \ldots + c_d^{(1)}v_d$$
$$\vdots$$
$$z^{(i)} = c_1^{(i)}v_1 + c_2^{(i)}v_2 + \ldots + c_d^{(i)}v_d$$

**Note:** $[c_1^{(i)}, \ldots, c_d^{(i)}] = c^{(i)} = V^T z^{(i)}$.

**Also:** Since $V$ is orthogonal and $\|z^{(i)}\|_2 = 1$, $\|c^{(i)}\|_2^2 = 1$.

**Claim:** After update $z^{(i)} = \frac{1}{n_i}X^T X z^{(i-1)}$,

$$c_j^{(i)} = \frac{1}{n_i}\sigma_j^2 c_j^{(i-1)}$$

$$z^{(i)} = \frac{1}{n_i}\left[c_1^{(i-1)}\sigma_1^2 \cdot v_1 + c_2^{(i-1)}\sigma_2^2 \cdot v_2 + \ldots + c_d^{(i-1)}\sigma_d^2 \cdot v_d\right]$$

Equivalently: $c^{(i)} = \frac{1}{n_i}\Sigma^2 c^{(i-1)}$.

**Claim:** After $T$ updates:

$$\mathbf{z}^{(T)} = \frac{1}{\prod_{i=1}^{T} n_i} \left[ c_1^{(0)} \sigma_1^{2T} \cdot \mathbf{v}_1 + c_2^{(0)} \sigma_2^{2T} \cdot \mathbf{v}_2 + \ldots + c_d^{(0)} \sigma_d^{2T} \cdot \mathbf{v}_d \right]$$

Let $\alpha_j = \frac{1}{\prod_{i=1}^{T} n_i} c_j^{(0)} \sigma_j^{2T}$. **Goal:** Show that $\alpha_j \ll \alpha_1$ for all $j \neq 1$.

Since $z^{(T)}$ is a unit vector, $\sum_{i=1}^{d} \alpha_i^2 = 1$.

If we can prove that $\left|\frac{\alpha_j}{\alpha_1}\right| \leq \sqrt{\frac{\epsilon}{2d}}$ then we will have that $\|v_1 - z^{(T)}\|_2^2 \leq \epsilon$.

$$\alpha_j^2 \leq \alpha_1^2 \cdot \frac{\epsilon}{2d}$$

$$1 = \alpha_1^2 + \sum_{j=2}^{d} \alpha_d^2 \leq \alpha_1^2 + \frac{\epsilon}{2}$$

$$\alpha_1^2 \geq 1 - \frac{\epsilon}{2}$$

$$|\alpha_1| \geq 1 - \frac{\epsilon}{2}$$

$$\|v_1 - z^{(T)}\|_2^2 = 2 - 2\langle v_1, z^{(T)} \rangle \leq \epsilon$$

Want to prove that after $T = O(\log(d/\epsilon)/\gamma)$ steps, $\left|\frac{\alpha_j}{\alpha_1}\right| \leq \sqrt{\frac{\epsilon}{2d}}$

where $\alpha_j = \frac{1}{\prod_{i=1}^{T} n_i} c_j^{(0)} \sigma_j^{2T}$. Recall that $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_2}$.

**Assumption:** Starting coeff. on first eigenvector is not too small:

$$\left| c_1^{(0)} \right| \geq O\left(\frac{1}{\sqrt{d}}\right).$$

We will prove shortly that this holds with probability 99/100.

$$\frac{|\alpha_j|}{|\alpha_1|} = \frac{\sigma_j^{2T}}{\sigma_1^{2T}} \cdot \frac{|c_j^{(0)}|}{|c_1^{(0)}|} \leq$$

Need $T =$

**Need to prove:** Starting coefficient on first eigenvector is not too small. I.e., with probability 99/100,

$$\left| c_1^{(0)} \right| \geq O\left( \frac{1}{\sqrt{d}} \right).$$

**Prove using Gaussian <u>anti</u>-concentration.** First use rotational invariance of Gaussian:

$$c^{(0)} = \frac{V^T z^{(0)}}{\|z^{(0)}\|_2} = \frac{V^T z^{(0)}}{\|V^T z^{(0)}\|_2} \sim \frac{g}{\|g\|_2},$$

where $g \sim \mathcal{N}(0,1)^d$.

Need to show that with high probability, first entry of
$\frac{\mathbf{g}}{\|\mathbf{g}\|_2} \geq c \cdot \frac{1}{\sqrt{d}}$.
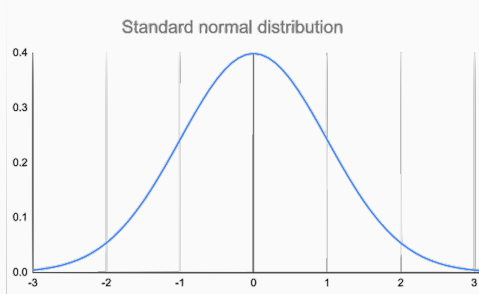
Part 1: With super high probability (e.g. 99/100),

$$\|\mathbf{g}\|_2^2 \leq$$

Need to show that with high probability, the magnitude of the first entry of $\mathbf{g} \geq c$ for a constant $c$. Think e.g. $c = 1/100$.

**Part 2:** With probablility $1 - O(\alpha)$,

$$|g_1| \geq \alpha.$$


Standard normal distribution

> **Theorem ( Power Method Convergence, $k = 1$)**
>
> *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 either:*
>
> $$\|\mathbf{v}_1 - \mathbf{z}^{(T)}\|_2 \leq \epsilon \qquad \text{or} \qquad \|\mathbf{v}_1 - (-\mathbf{z}^{(T)})\|_2 \leq \epsilon.$$

The method truly won't converge if $\gamma$ is very small. Consider extreme case when $\gamma = 0$.

$$\mathbf{z}^{(T)} = \frac{1}{\prod_{i=1}^{T} n_i} \left[ c_1^{(0)} \sigma_1^{2T} \cdot \mathbf{v}_1 + c_2^{(0)} \sigma_2^{2T} \cdot \mathbf{v}_2 + \ldots + c_d^{(0)} \sigma_d^{2T} \cdot \mathbf{v}_d \right]$$

17

### Theorem (Gapless Power Method Convergence)

*If Power Method is initialized with a random Gaussian vector then, with high probability, after $T = O\left(\frac{\log d/\epsilon}{\epsilon}\right)$ steps, we obtain a z satisfying:*

$$\|X - Xzz^T\|_F^2 \leq (1 + \epsilon)\|X - Xv_1v_1^T\|_F^2$$

**Intuition:** For a good low-rank approximation, we don't actually need to converge to $v_1$ if $\sigma_1$ and $\sigma_2$ are the same or very close. Would suffice to return either $v_1$ or $v_2$, or some linear combination of the two.

- Block Power Method aka Simultaneous Iteration aka Subspace Iteration aka Orthogonal Iteration

**Power method:**

- Choose $G \in \mathbb{R}^{d \times k}$ be a random Gaussian matrix.
- $Z_0 = \text{orth}(G)$.
- For $i = 1, \ldots, T$
    - $Z^{(i)} = X^T \cdot (XZ^{(i-1)})$
    - $Z^{(i)} = \text{orth}(Z^{(i)})$

  Return $Z^{(T)}$

  **Guarantee:** After $O\left(\frac{\log d/\epsilon}{\epsilon}\right)$ iterations:

  $$\|X - XZZ^T\|_F^2 \leq (1 + \epsilon)\|X - XV_kV_k^T\|_F^2.$$

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

Possible to "accelerate" these methods.

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

$$\|X - XZZ^T\|_F^2 \leq (1 + \epsilon)\|X - XV_kV_k^T\|_F^2.$$

For a normalizing constant $c$, power method returns:

$$z^{(q)} = c \cdot \left(X^T X\right)^q \cdot g$$

Along the way we computed:

$$\mathcal{K}_q = \left[g, \left(X^T X\right) \cdot g, \left(X^T X\right)^2 \cdot g, \ldots, \left(X^T X\right)^q \cdot g\right]$$

$\mathcal{K}$ is called the Krylov subspace of degree $q$.

Idea behind Krlyov methods: Don't throw away everything before $\left(X^T X\right)^q \cdot g$.

Want to find $v$, which minimizes $\|X - Xvv^T\|_F^2$.

### Lanczos method:

- Let $Q \in \mathbb{R}^{d \times k}$ be an orthonormal span for the vectors in $\mathcal{K}$.
- Solve $\min_{v=Qw} \|X - Xvv^T\|_F^2$.
    - Find <u>best</u> vector in the Krylov subspace, instead of just using last vector.
    - Can be done in $O\left(ndk + dk^2\right)$ time.
    - What you're using when you run `svds` or `eigs` in MATLAB or Python.

For a degree $T$ polynomial $p$, let $\mathbf{v}_p = \frac{p(\mathbf{X}^T\mathbf{X})\mathbf{g}}{\|p(\mathbf{X}^T\mathbf{X})\mathbf{g}\|_2}$. We always have that $\mathbf{v}_p \in \mathcal{K}_T$, the Krylov subspace constructed with $T$ iterations.
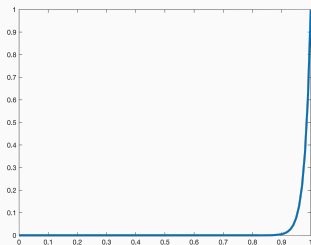
Power method returns:

$$\mathbf{v}_{x^T}.$$

Lanczos method returns $\mathbf{v}_{p^*}$ where:

$$p^* = \underset{\text{degree } T\ p}{\arg\min}\ \|\mathbf{X} - \mathbf{X}\mathbf{v}_p\mathbf{v}_p^T\|_F^2.$$

**Claim:** There is a $q = O\left(\sqrt{T \log \frac{1}{\Delta}}\right)$ degree polynomial $\hat{p}$ approximating $x^T$ up to error $\Delta$ on $[0, 1]$.



$$\|X - Xv_{p*}v_{p*}^T\|_F^2 \leq \|X - Xv_{\hat{p}}v_{\hat{p}}^T\|_F^2 \approx \|X - Xv_{x^T}v_{x^T}^T\|_F^2 \approx \|X - Xv_1v_1^T\|_F^2$$

Runtime: $O\left(\frac{\log(d/\epsilon)}{\sqrt{\epsilon}} \cdot \text{nnz}(X)\right)$ vs. $O\left(\frac{\log(d/\epsilon)}{\epsilon} \cdot \text{nnz}(X)\right)$
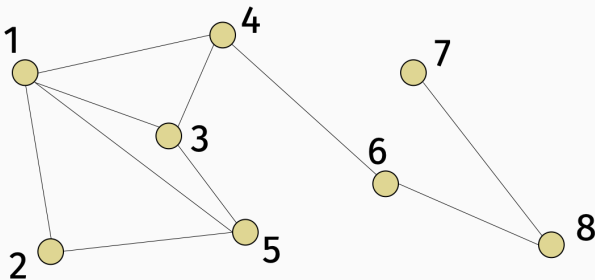
- Block Krylov methods

- Let $G \in \mathbb{R}^{d \times k}$ be a random Gaussian matrix.
- $\mathcal{K}_q = \left[ G, (X^TX) \cdot G, (X^TX)^2 \cdot G, \ldots, (X^TX)^q \cdot G \right]$

Runtime: $O\left( \text{nnz}(X) \cdot k \cdot \frac{\log d/\epsilon}{\sqrt{\epsilon}} \right)$ to obtain a nearly optimal low-rank approximation.
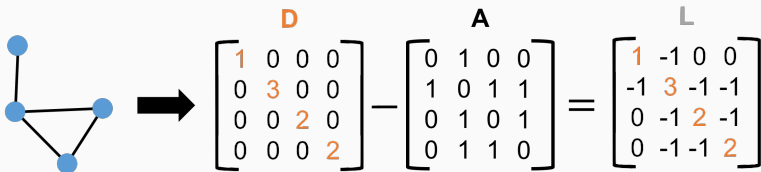
BREAK

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



For now assume $G = (V, E)$ is an undirected, unweighted graph with $n$ nodes.
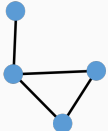
26

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.



$$\mathbf{D} \qquad \mathbf{A} \qquad \mathbf{L}$$

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

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 $\bar{\mathbf{L}} = \mathbf{I} - \mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$.

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

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

$B$ has a row for every edge in $G$. The row for edge $(i, j)$ has a $+1$ at position $i$, a $-1$ at position $j$, and zeros elsewhere.
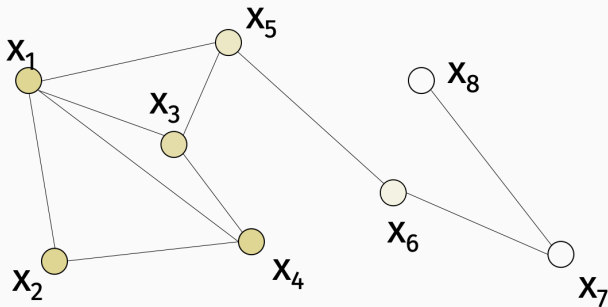
## Conclusions from $L = B^T B$

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

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

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

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

$x^T L x = \sum_{(i,j) \in E} (x(i) - x(j))^2$. So $x^T L x$ is small if $x$ is a "smooth" function with respect to 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 $L$.

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

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

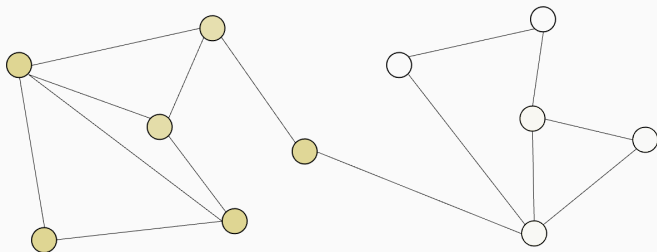$$v_3 = \underset{\|v\|=1, v \perp v_1, v_2}{\arg\max}\, v^T L v$$

$$\vdots$$

$$v_n = \underset{\|v\|=1, v \perp v_1, \ldots, v_{n-1}}{\arg\max}\, v^T L v$$

Another conclusion from $L = B^T B$:

For a <u>cut indicator vector</u> $\mathbf{c} \in \{-1, 1\}^n$ with $\mathbf{c}(i) = -1$ for $i \in S$ and $\mathbf{c}(i) = 1$ for $i \in T = V \setminus S$:
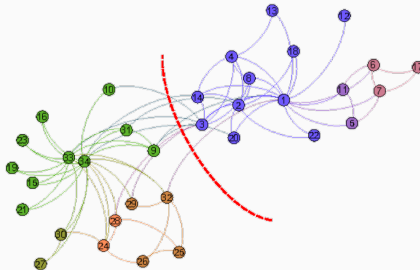
$$\mathbf{c}^T L \mathbf{c} = \sum_{(i,j) \in E} (\mathbf{c}(i) - \mathbf{c}(j))^2 = 4 \cdot \text{cut}(S, T). \tag{1}$$



33

- Introduce NP-hard graph partitioning prob. important in:
    - Understanding social networks.
    - Unsupervised machine learning (spectral clustering).
    - Graph visualization.
    - Mesh partitioning.
- See how this problem can be solved heuristically using Laplacian eigenvectors.
- Give an "average case" analysis of the method for a common random graph model.
- Use two tools: matrix concentration and eigenvector perturbation bounds.

**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

**Example application:** Understanding <u>community structure</u> in social networks.
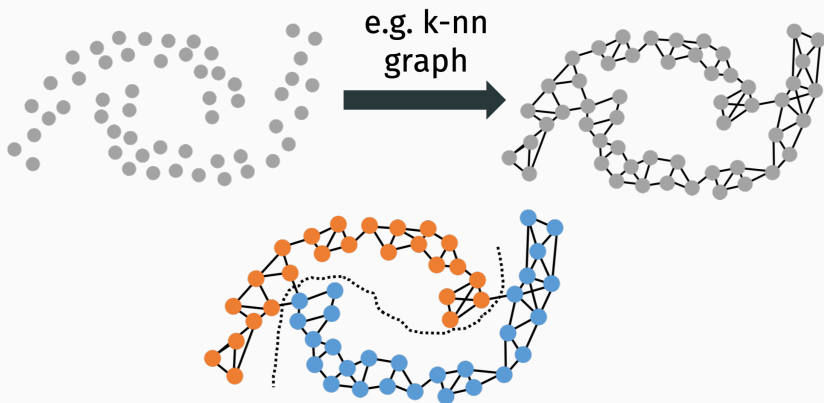
Wayne W. Zachary (1977). An Information Flow Model for Conflict and Fission in Small Groups.

"At the beginning of the study there was an incipient conflict between the club president, John A., and Mr. Hi over the price of karate lessons. Mr. Hi, who wished to raise prices, claimed the authority to set his own lesson fees, since he was the instructor. John A., who wished to stabilize prices, claimed the authority to set the lesson fees since he was the club's chief administrator. As time passed the entire club became divided over this issue, and the conflict became translated into ideological terms by most club members."

Zachary constructed a social network by hand and used a minimum cut algorithm to correctly predict who sided with who in the conflict. Beautiful paper – definitely worth checking out!
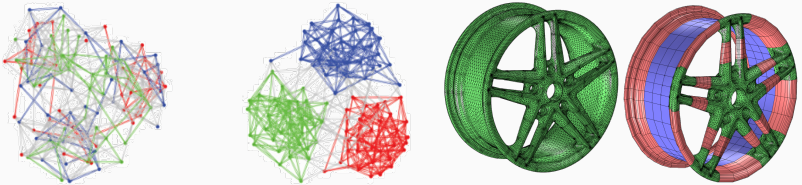
**Idea:** Construct synthetic graph for data that is hard to cluster.



e.g. k-nn graph

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

Balanced cut algorithms are also use in distributing data in graph databases, for partitioning finite element meshes in scientific computing (e.g., that arise when solving differential equations), and more.



Lots of good software packages (e.g. METIS).

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

### $\beta$-Balanced Cut:

$$\min_S \text{cut}(S, V \setminus S) \quad \text{such that} \quad \min(|S|, |V \setminus S|) \geq \beta \cdot n \text{ for } \beta \leq .5$$

### Sparsest Cut:

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

All natural formalizations lead to NP-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.
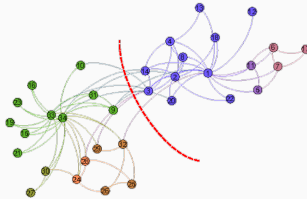
Spectral methods run no more than $O(n^3)$ time (must faster if you use iterative methods for computing eigenvectors).

39

Basic spectral clustering method:

- Compute <u>second</u> smallest eigenvector of graph, $\mathbf{v}_{n-1}$.
- $\mathbf{v}_{n-1}$ has an entry for every node $i$ in the graph.
- If the $i^{\text{th}}$ entry is positive, put node $i$ in $T$.
- Otherwise if the $i^{\text{th}}$ entry is negative, put $i$ in $S$.

This shouldn't make much sense yet! We will see that is a "relax and round" algorithm in disguise.
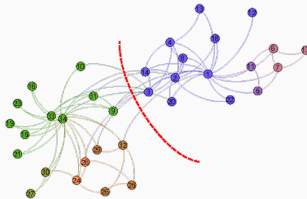
(a) Zachary Karate Club Graph

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

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

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

(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:

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

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

Perfectly balanced balanced cut problem:

$$\min_{\mathbf{c} \in \{-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\}^n} \mathbf{c}^T \mathbf{L} \mathbf{c} \text{ such that } \mathbf{c}^T \mathbf{1} = 0.$$

Relaxed perfectly balanced balanced cut problem:

$$\min_{\|\mathbf{c}\|_2 = 1} \mathbf{c}^T \mathbf{L} \mathbf{c} \text{ such that } \mathbf{c}^T \mathbf{1} = 0.$$

Claim: The relaxed problem is exactly minimized by the second smallest eigenvector $\mathbf{v}_{n-1}$ of $\mathbf{L}$.
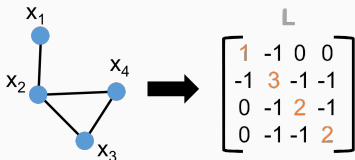
Approach: Relax, find $\mathbf{v}_{n-1}$, then round back to a vector with $-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}$ entries.

**Claim:** The smallest eigenvector/singular vector of any graph Laplacian $L$ always equals:

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

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

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

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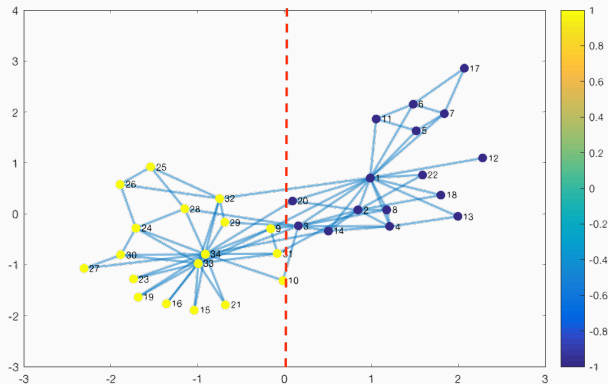
which is equivalent to

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

**Final relax and round algorithm:** Compute

$$\mathbf{v}_{n-1} = \underset{v \in \mathbb{R}^n \text{ with } \|\mathbf{v}\|=1, \; \mathbf{v}^T \mathbf{1}=0}{\arg\min} \mathbf{v}^T L \mathbf{v}$$
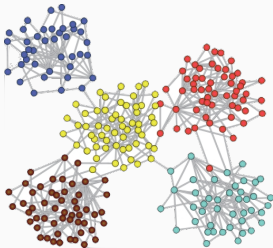
Set $S$ to be all nodes with $\mathbf{v}_{n-1}(i) < 0$, and $T$ to be all with $\mathbf{v}_{n-1}(i) \geq 0$. I.e. set $\mathbf{c} = \text{sign}(\mathbf{v}_{n-1})$



46

Lots of different variants used in practice:

- Often do some sort of normalization of edge weights by degree. E.g. the Shi-Malik normalized cuts algorithm use the normalized Laplacian $\overline{L} = D^{-1/2}LD^{-1/2}$.
- Different methods for how to choose the threshold to partition the second smallest eigenvector.
- Lots of variants to split the graph into more than two parts.
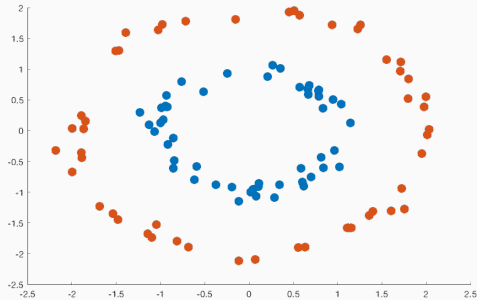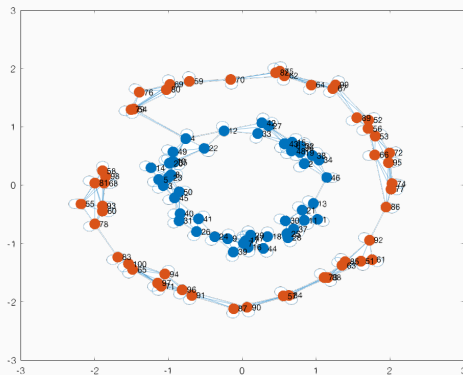
Multiway spectral partitioning:

- Compute smallest $\ell$ eigenvectors $\mathbf{v}_{n-1}, \ldots, \mathbf{v}_{n-\ell}$ of L.
- Represent each node by its corresponding row in $\mathsf{V} \in \mathbb{R}^{n \times \ell}$ whose rows are $\mathbf{v}_{n-1}, \ldots \mathbf{v}_{n-\ell}$.
- Cluster these rows using $k$-means clustering (or really any clustering method).
- Often we choose $\ell = k$, but not necessarily.

Let $\tilde{\mathbf{x}}_i \in \mathbb{R}^\ell$ denote this embedding for a particullar node $i$.
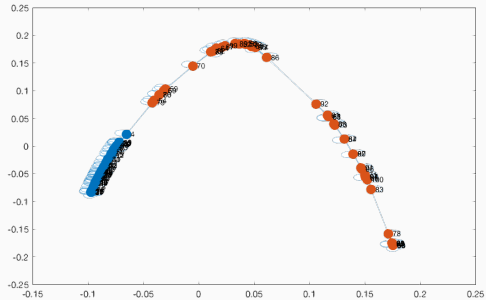
**Original Data:** (not linearly separable)

$k$-Nearest Neighbors Graph:

Embedding with eigenvectors $v_{n-1}, v_{n-2}$: (linearly separable)
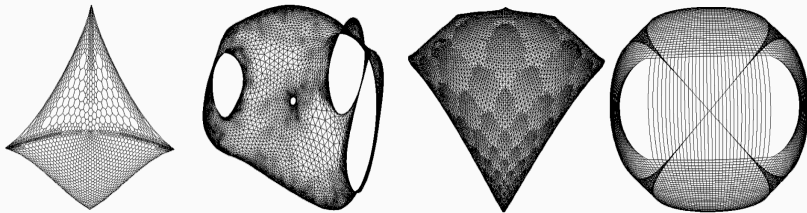
Formally, this choice of embedding minimizes:

$$\sum_{i,j \in E} \|\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j\|^2$$

I.e., we 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.
- Can fail for worst case input graphs.

**Common approach:** Design a natural <span style="color:orange">generative model</span> that produces <u>random but realistic</u> 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. Often our best way to understand why some algorithms "just work" in practice.
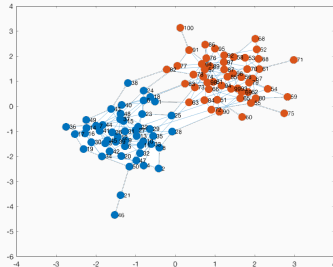- Similar approach to Bayesian modeling in machine learning.

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

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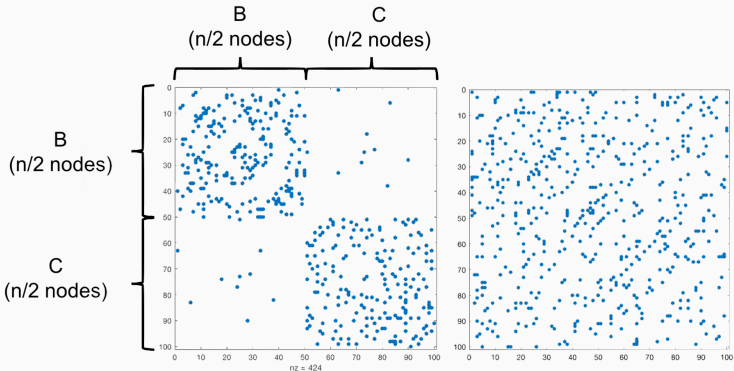
- Any two nodes in the **same group** are connected with probability $p$ (including self-loops).
- Any two nodes in **different groups** are connected with prob. $q < p$.

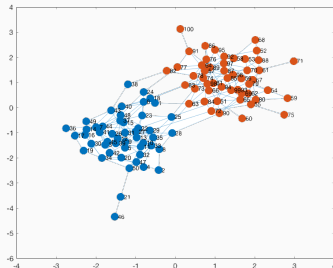Let $G$ be a stochastic block model graph drawn from $G_n(p, q)$.

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



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

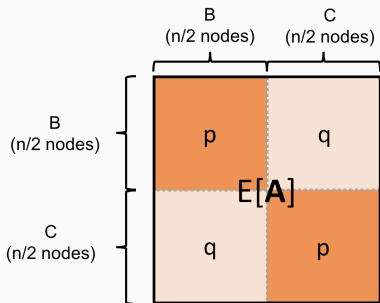Goal is to find the "ground truth" balanced partition $B, C$ using our standard spectal method.



To do so, we need to understand the second smallest eigenvector of $L = D - A$. We will start by considering the expected value of these matrices:

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

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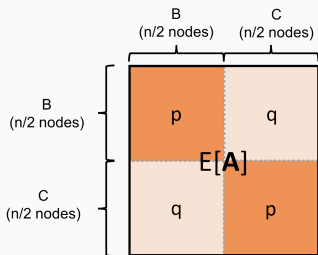
What is the expected Laplacian of $G_n(p, q)$?

$\mathbb{E}[\mathsf{A}]$ and $\mathbb{E}[\mathsf{L}]$ have the same eigenvectors and eigenvalues are equal up to a shift/inversion. So second largest eigenvector of $\mathbb{E}[\mathsf{A}]$ is the same as the second smallest of $\mathbb{E}[\mathsf{L}]$

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

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

If we compute $\bar{\mathbf{v}}_2$ then we <u>exactly recover</u> the communities $B$ and $C$!

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

- If the random graph $G$ (equivilantly $\mathbf{A}$ and $\mathbf{L}$) 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., $\mathbf{A}$) is close to its expectation? **Matrix concentration inequalities.**

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

Alon, Krivelevich, Vu, 2002:

> **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).

Recall that $\|X\|_2 = \max_{z \in \mathbb{R}^d : \|z\|_2 = 1} \|Xz\|_2 = \sigma_1(X)$.

$\|A\|_2$ is on the order of $O(p\sqrt{n})$ so another way of thinking about the right hand side is $\frac{\|A\|_2}{\sqrt{p}}$. I.e. get's better with $p$.

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

> **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 $v_1, v_2, \ldots, v_n$ and $\bar{v}_1, \bar{v}_2, \ldots, \bar{v}_n$. Letting $\theta(v_i, \bar{v}_i)$ denote the angle between $v_i$ and $\bar{v}_i$, for all $i$:
>
> $$\sin[\theta(v_i, \bar{v}_i)] \leq \frac{\epsilon}{\min_{j \neq i} |\lambda_i - \lambda_j|}$$
>
> where $\lambda_1, \ldots, \lambda_n$ are the eigenvalues of $\overline{A}$.

We will apply with $\overline{A} = \mathbb{E}[A]$.

| **A** | | **Ā** | | **A-Ā** | |
|---|---|---|---|---|---|
| 1+ε | 0 | 1 | 0 | ε | 0 |
| 0 | 1 | 0 | 1+ε | 0 | ε |

$$\mathbf{A} - \mathbf{\bar{A}} = \mathbf{A} - \mathbf{\bar{A}}$$

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

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

**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 $\frac{(p-q)n}{2}$ will be the minimum of these two gaps.

**Claim 2 (Davis-Kahan):** For $p \geq O\left(\frac{\log^4 n}{n}\right)$,
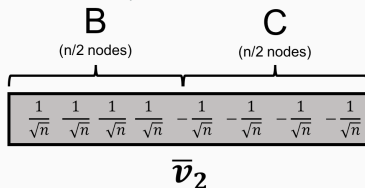
$$\sin\theta(v_2, \bar{v}_2) \leq \frac{O(\sqrt{pn})}{\min_{j \neq i} |\lambda_i - \lambda_j|} \leq \frac{O(\sqrt{pn})}{(p-q)n/2} = O\left(\frac{\sqrt{p}}{(p-q)\sqrt{n}}\right)$$

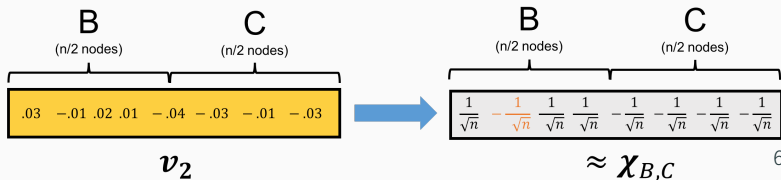(A slightly trickier analysis can remove the $qn$ term entirely.)

**So far:** $\sin\theta(\mathbf{v}_2, \bar{\mathbf{v}}_2) \leq O\left(\frac{\sqrt{p}}{(p-q)\sqrt{n}}\right)$. What does this give us?

- Can show that this implies $\|\mathbf{v}_2 - \bar{\mathbf{v}}_2\|_2^2 \leq O\left(\frac{p}{(p-q)^2 n}\right)$ (exercise).

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



B
(n/2 nodes)

C
(n/2 nodes)

| $\frac{1}{\sqrt{n}}$ | $\frac{1}{\sqrt{n}}$ | $\frac{1}{\sqrt{n}}$ | $\frac{1}{\sqrt{n}}$ | $-\frac{1}{\sqrt{n}}$ | $-\frac{1}{\sqrt{n}}$ | $-\frac{1}{\sqrt{n}}$ | $-\frac{1}{\sqrt{n}}$ |

$\overline{\boldsymbol{v}}_2$

- We want to show that $\text{sign}(\mathbf{v}_2)$ and $\bar{\mathbf{v}}_2$ are close. They only differ at locations where $\mathbf{v}_2$ and $\bar{\mathbf{v}}_2$ differ in sign.



B
(n/2 nodes)

C
(n/2 nodes)

| .03 | $-.01$ | .02 | .01 | $-.04$ | $-.03$ | $-.01$ | $-.03$ |

$\boldsymbol{v}_2$

B
(n/2 nodes)

C
(n/2 nodes)

| $\frac{1}{\sqrt{n}}$ | $-\frac{1}{\sqrt{n}}$ | $\frac{1}{\sqrt{n}}$ | $\frac{1}{\sqrt{n}}$ | $-\frac{1}{\sqrt{n}}$ | $-\frac{1}{\sqrt{n}}$ | $-\frac{1}{\sqrt{n}}$ | $-\frac{1}{\sqrt{n}}$ |

$\approx \boldsymbol{\chi}_{B,C}$

67

Main argument:

- Every $i$ where $v_2(i)$, $\bar{v}_2(i)$ differ in sign contributes $\geq \frac{1}{n}$ to $\|\mathbf{v}_2 - \bar{\mathbf{v}}_2\|_2^2$.
- We know that $\|\mathbf{v}_2 - \bar{\mathbf{v}}_2\|_2^2 \leq O\left(\frac{p}{(p-q)^2 n}\right)$.
- So $\mathbf{v}_2$ and $\bar{\mathbf{v}}_2$ 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 largest 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.

- Hard case: Suppose $q = .8p$ so $\frac{p}{(p-q)^2} = 25/p$.

If $p = c/n$, the number of mistakes is $\frac{25}{c} \cdot n$. I.e., $< 10\%$ error if average node has roughly 250 connections.

Forget about the previous problem, but still consider the matrix $\mathbf{M} = \mathbb{E}[\mathbf{A}]$.

- Dense $n \times n$ matrix.
- Computing top eigenvectors takes $\approx O(n^2/\sqrt{\epsilon})$ time.

If someone asked you to speed this up and return <u>approximate</u> top eigenvectors, what could you do?

**Main idea:** If you want to compute singular vectors, multiply two matrices, solve a regression problem, etc.:

1. Compress your matrices using a randomized method (e.g. subsampling).
2. Solve the problem on the smaller or sparser matrix.
   - $\tilde{A}$ called a "sketch" or "coreset" for $A$.