

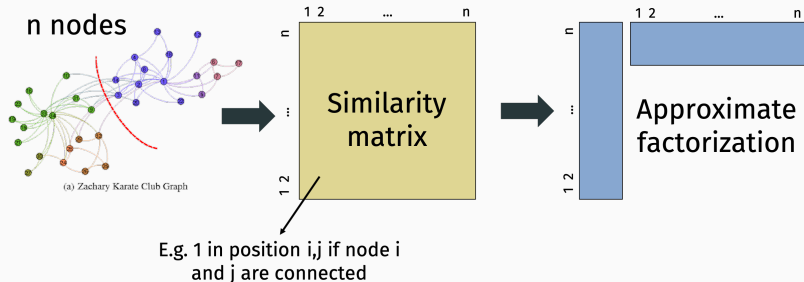
CS-GY 9223 I: Lecture 11

Spectral graph theory + randomized numerical linear algebra.

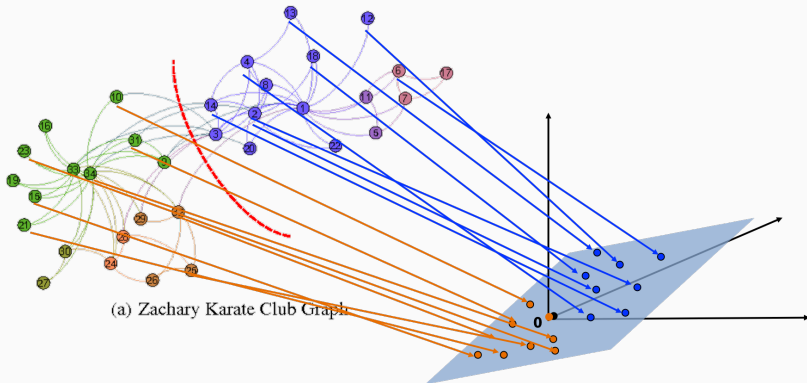
NYU Tandon School of Engineering, Prof. Christopher Musco

ENCODING GRAPH SIMILARITY

Often data is represented as a graph and similarities can be obtained from that graph:



ENCODING GRAPH SIMILARITY

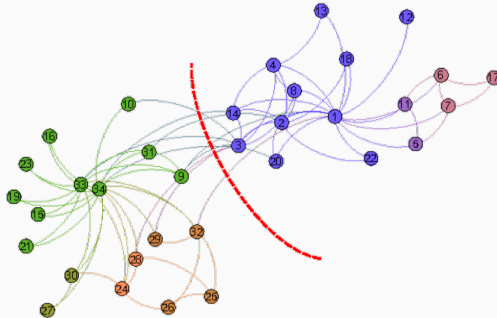


Spectral graph theory lets us formalize this heuristic idea.

CUT MINIMIZATION

Goal: 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.

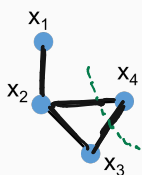


(a) Zachary Karate Club Graph

THE LAPLACIAN VIEW

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

$$\|Bc\|_2^2 = c^T L c = c^T B^T B c$$



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

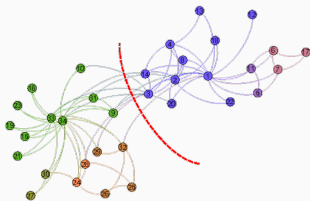
(cut indicator vector $[1 \ 1 \ -1 \ -1]$)

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

$$\begin{matrix} & \text{node 1} & \text{node 2} & 3 & 4 \dots \\ \text{edges} & \begin{bmatrix} 1 & -1 & 0 & 0 \\ 0 & 1 & -1 & 0 \\ 0 & 1 & 0 & -1 \\ 0 & 0 & 1 & -1 \end{bmatrix} & \begin{bmatrix} 0 \\ 0 \\ 2 \\ 2 \end{bmatrix} & = & \begin{bmatrix} 0 \\ 0 \\ 2 \\ 2 \end{bmatrix}
 \end{matrix}$$

$BC = \begin{cases} 0 & \text{for edges not across the cut} \\ 2 & \text{for edges across the cut} \end{cases}$

$$\|Bc\|_2^2 = 4 \cdot (\# \text{ edges across cut})^2 = 4 \cdot 5 = 20$$



(a) Zachary Karate Club Graph

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

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

- $\mathbf{c}^T \mathbf{1} = |T| - |S|$. \approx "imbalance in cut"

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

Courant–Fischer min-max principle

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

↓
for symmetric
matrix

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

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

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

⋮

$$v_n = \arg \max_{\|v\|=1, v \perp v_1, \dots, v_{n-1}} v^T L v$$

Courant–Fischer min-max principle

Let $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_n]$ be the eigenvectors of \mathbf{L} .

$$\mathbf{v}_n = \arg \min_{\|\mathbf{v}\|=1} \mathbf{v}^T \mathbf{L} \mathbf{v}$$

$$\mathbf{v}_{n-1} = \arg \min_{\|\mathbf{v}\|=1, \mathbf{v} \perp \mathbf{v}_n} \mathbf{v}^T \mathbf{L} \mathbf{v}$$

$$\mathbf{v}_{n-2} = \arg \min_{\|\mathbf{v}\|=1, \mathbf{v} \perp \mathbf{v}_n, \mathbf{v}_{n-1}} \mathbf{v}^T \mathbf{L} \mathbf{v}$$

$$\vdots$$

$$\mathbf{v}_1 = \arg \min_{\|\mathbf{v}\|=1, \mathbf{v} \perp \mathbf{v}_n, \dots, \mathbf{v}_2} \mathbf{v}^T \mathbf{L} \mathbf{v}$$

SMALLEST LAPLACIAN EIGENVECTOR

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

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

~~arg~~ \rightarrow for $\frac{1}{\sqrt{n}} \cdot \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$

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

1) All of L 's eigenvalues are ≥ 0 .

L is positive semidefinite: $B^T B = L$.

$B^T B$

$$2) \mathbf{1} \cdot \mathbf{1} = B^T B \cdot \mathbf{1} = B^T \mathbf{0} = \mathbf{0}$$

$$\begin{bmatrix} 1 & -1 \\ -1 & 1 \\ \vdots & \vdots \\ -1 & 1 \end{bmatrix} \cdot \mathbf{1} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$L \cdot \left(\frac{1}{\sqrt{n}} \mathbf{1} \right) = 0 \left(\frac{1}{\sqrt{n}} \right) \mathbf{1}$$

SECOND SMALLEST LAPLACIAN EIGENVECTOR

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

$$\mathbf{v}_{n-1} = \operatorname{argmin}_{\|\mathbf{v}\|=1, \mathbf{v}_n^T \mathbf{v}=0} \mathbf{v}^T \mathbf{L} \mathbf{v}$$

$$v_n = \frac{1}{\sqrt{n}} \cdot \mathbf{1}$$

cut value for
cut indicators
add constraint
 $v_{n-1} \in \{ \frac{1}{\sqrt{n}}, -\frac{1}{\sqrt{n}} \}$

If \mathbf{v}_{n-1} were binary, i.e. $\in \{-1, 1\}^n$, scaled by $\frac{1}{\sqrt{n}}$, it would have:

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

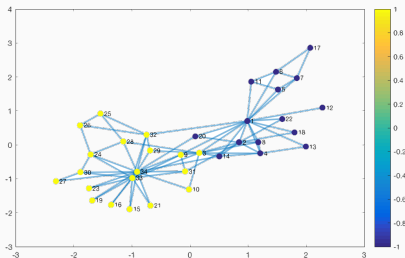
$\mathbf{v}_{n-1} \in \mathbb{R}^n$ is not generally binary, but still satisfies a 'relaxed' version of this property.

CUTTING WITH THE SECOND LAPLACIAN EIGENVECTOR

Find a good partition of the graph by using an eigendecomposition to compute

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

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



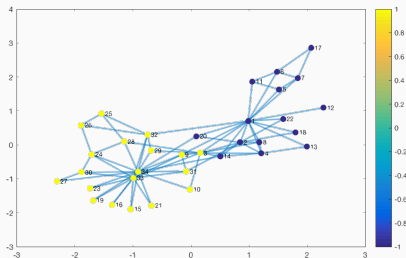
$\mathbf{v}_{n-1} \approx$

CUTTING WITH THE SECOND LAPLACIAN EIGENVECTOR

Find a good partition of the graph by using an eigendecomposition to compute

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

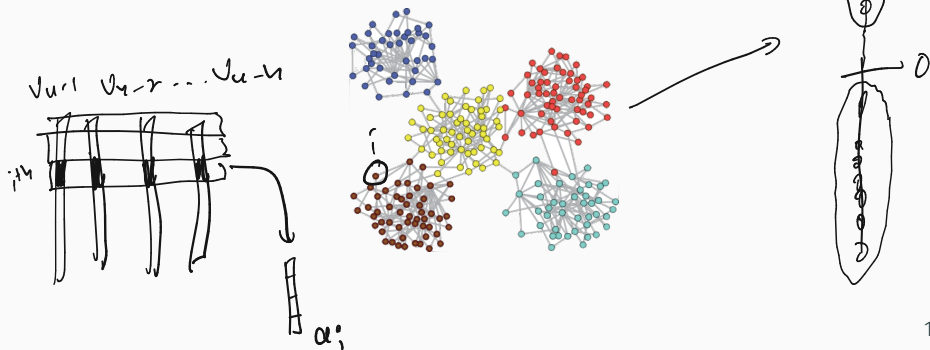
Set S to be all nodes with $\mathbf{v}_{n-1}(i) < 0$, and T to be all with $\mathbf{v}_{n-1}(i) \geq 0$.



SPECTRAL PARTITIONING IN PRACTICE

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

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



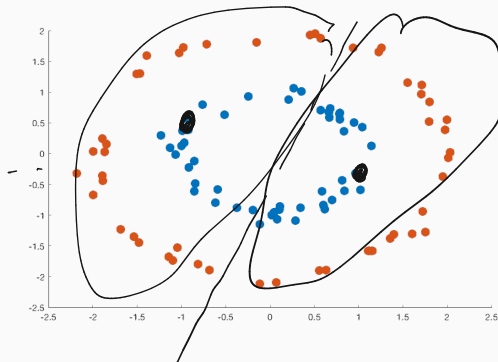
Spectral Clustering:

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

LAPLACIAN EMBEDDING

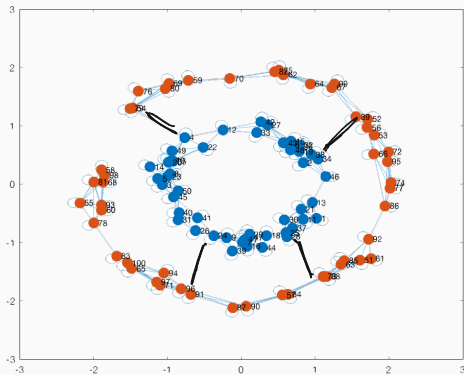
2-means

Original Data: (not linearly separable)



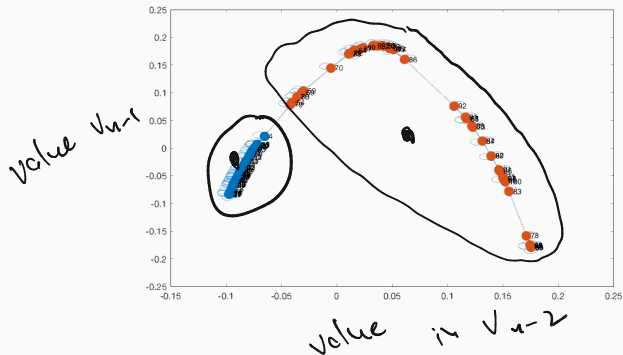
LAPLACIAN EMBEDDING

→ 5 or 10
 k -Nearest Neighbors Graph:



LAPLACIAN EMBEDDING

Embedding with eigenvectors v_{n-1}, v_{n-2} : (linearly separable)

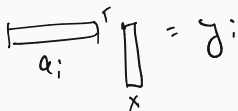


So far: Spectral clustering partitions a graph along a small cut between large pieces.

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

Common approach: Give a natural **generative model** for which produces random but realistic 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 ℓ_2 linear regression, k -means clustering, PCA, etc.)


$$a_i \quad x = \gamma$$

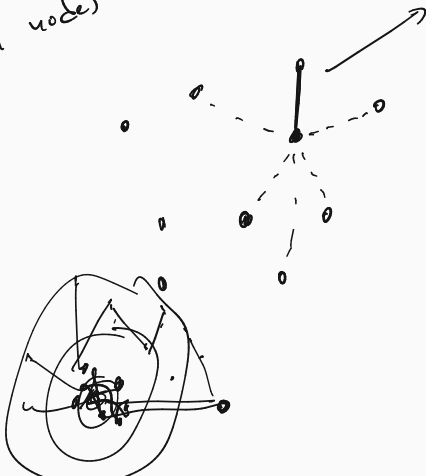

$$A \quad x$$

$$x = \min \|Ax - b\|_{\ell_1, \ell_0}$$

STOCHASTIC BLOCK MODEL

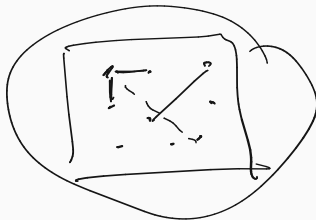
Ideas for a generative model for graphs that would allow us to understand partitioning?

n nodes



connect with prob
 $p = 0.1$

add with prob
 $p = 0.15$



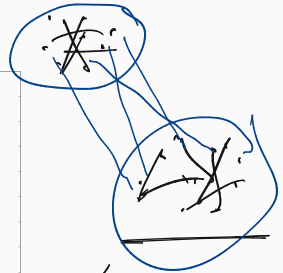
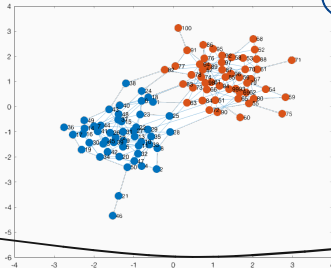
STOCHASTIC BLOCK MODEL

Stochastic Block Model (Planted Partition Model):

Let $G_n(p, q)$ be a distribution over graphs on n nodes, split equally into two groups \underline{B} and \underline{C} , each with $\underline{n/2}$ nodes.

- Any two nodes in the **same group** are connected with probability \underline{p} (including self-loops). $p = 0.1$
- Any two nodes in **different groups** are connected with prob. $q < p$.

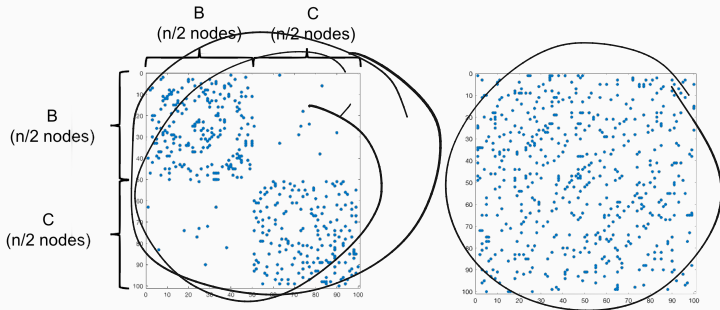
Erdos-Renyi
 $G_{n/2}(p)$



LINEAR ALGEBRAIC VIEW

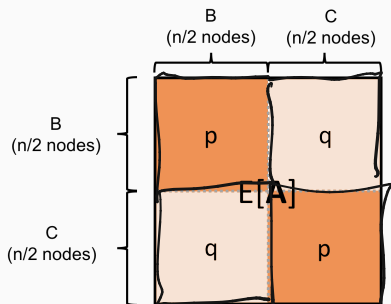
Let G be a stochastic block model graph drawn from $G_n(p, q)$.

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



EXPECTED ADJACENCY SPECTRUM

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 $\mathbb{E}[\mathbf{A}]$?

EXPECTED ADJACENCY SPECTRUM

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 $\mathbb{E}[A]$?

$$\mathbb{E}[A] = \begin{array}{|c|c|} \hline \dots & \dots \\ \hline p & q \\ \hline q & p \\ \hline \end{array}$$

$$\mathbb{E}[A] \cdot \mathbf{1} =$$

$$\begin{bmatrix} \frac{(p+q)n}{2} \\ \vdots \\ \frac{(p+q)n}{2} \\ \vdots \\ \frac{(p+q)n}{2} \end{bmatrix} = \frac{(p+q)n}{2} \mathbf{1}$$

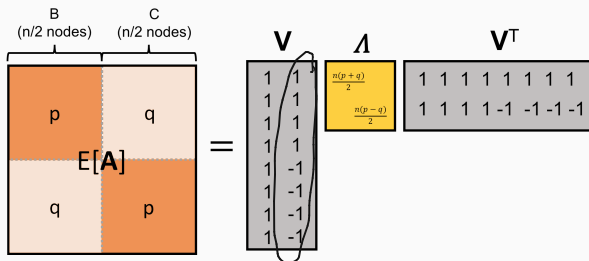
$$\text{Eigenvector 1: } \frac{1}{\sqrt{n}} \mathbf{1}$$

$$\mathbb{E}[A] \cdot u = \begin{bmatrix} (p-q)n/2 \\ \vdots \\ (p-q)n/2 \\ \hline (q-p)n/2 \end{bmatrix} = \frac{(p-q)n}{2} u$$

$$\text{Eigenvalue 1: } \frac{(p+q)n}{2}$$

$$\text{Eigenvector 2: } \frac{1}{\sqrt{n}} \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} \rightarrow u$$

EXPECTED ADJACENCY SPECTRUM



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

EXPECTED LAPLACIAN SPECTRUM

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

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

- If the random graph G (equivalently \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.

Matrix Concentration Inequality: If $p \geq O\left(\frac{\log^4 n}{n}\right)$, then with high probability

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

→ random

where $\|\cdot\|_2$ is the matrix **spectral** norm (operator norm).

For $X \in \mathbb{R}^{n \times d}$, $\|X\|_2 = \max_{z \in \mathbb{R}^d: \|z\|_2=1} \|Xz\|_2$. *→ vector* $\max (z^T A z)$

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

1974

Davis-Kahan Eigenvector Perturbation Theorem: Suppose $\mathbf{A}, \bar{\mathbf{A}} \in \mathbb{R}^{d \times d}$ are symmetric with $\|\mathbf{A} - \bar{\mathbf{A}}\|_2 \leq \epsilon$ and eigenvectors $\underline{v_1, v_2, \dots, v_d}$ and $\underline{\bar{v}_1, \bar{v}_2, \dots, \bar{v}_d}$. Letting $\theta(v_i, \bar{v}_i)$ denote the angle between v_i and \bar{v}_i , for all i :

measure of closeness between v_i, \bar{v}_i

$$\sin[\theta(v_i, \bar{v}_i)] \leq \frac{\epsilon}{\min_{j \neq i} |\lambda_i - \lambda_j|}$$

where $\lambda_1, \dots, \lambda_d$ are the eigenvalues of $\bar{\mathbf{A}}$.

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

EIGENVECTOR PERTURBATION

$$\begin{array}{ccc} \mathbf{A} & \mathbf{\bar{A}} & \mathbf{A-\bar{A}} \\ \begin{bmatrix} 1+\varepsilon & 0 \\ 0 & 1 \end{bmatrix} & \begin{bmatrix} 1 & 0 \\ 0 & 1+\varepsilon \end{bmatrix} & = \begin{bmatrix} \varepsilon & 0 \\ 0 & \varepsilon \end{bmatrix} \\ \downarrow & \downarrow & \\ \begin{array}{cc} v_1 \begin{bmatrix} 1 \\ 0 \end{bmatrix} & v_2 \begin{bmatrix} 0 \\ 1 \end{bmatrix} \\ \lambda_1 = 1+\varepsilon & \lambda_2 = 1 \end{array} & \begin{array}{cc} \begin{bmatrix} 0 \\ 1 \end{bmatrix} & \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \lambda_1 = 1+\varepsilon & \lambda_2 = 1 \end{array} & \|\mathbf{A-\bar{A}}\|_L = \varepsilon \end{array}$$

APPLICATION TO STOCHASTIC BLOCK MODEL

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

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

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

$$\sin \theta(v_2, \underline{\underline{\bar{v}}_2}) \leq \frac{O(\sqrt{pn})}{\min_{j \neq 2} |\lambda_2 - \lambda_j|} \leq \frac{O(\sqrt{pn})}{(p-q)n/2} = O\left(\frac{\sqrt{p}}{(p-q)\sqrt{n}}\right)$$

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

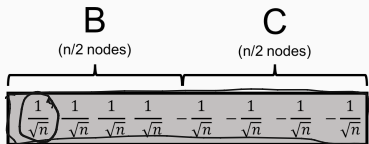
$$\frac{(p+q)n}{2} \quad \frac{(p-q)n}{2} \quad 0 \quad 0 \quad 0 \quad 0 \quad 0$$

$$\frac{(p-q)n}{2}$$

APPLICATION TO STOCHASTIC BLOCK MODEL

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.



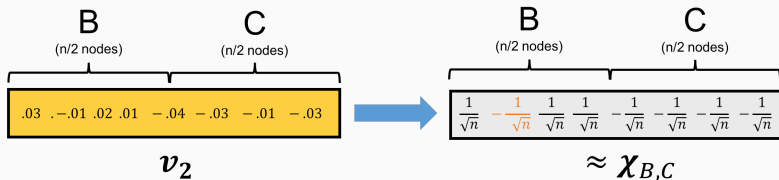
v_2

~~$v_2 = [0.1, 0.2, -0.001, -0.1, -0.5]$~~

- 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.

APPLICATION TO STOCHASTIC BLOCK MODEL

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.

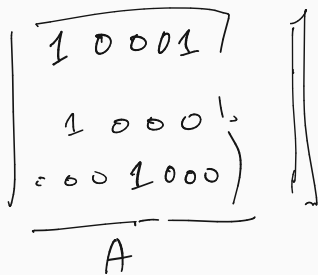
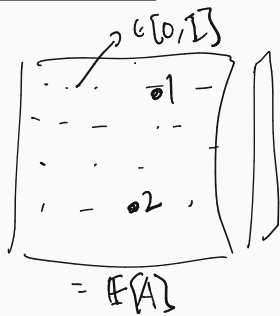
$$p - q = O(1/\sqrt{n}) \Rightarrow \frac{p}{(p-q)^2} = 0.1 \checkmark$$

RANDOMIZED NUMERICAL LINEAR ALGEBRA

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 approximate top eigenvectors, what could you do?

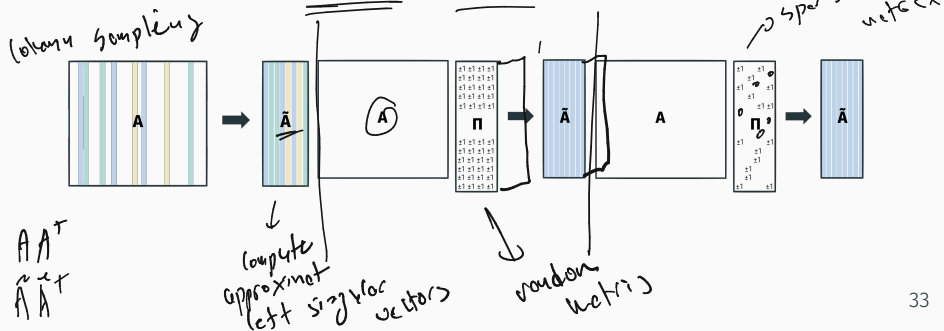


RANDOMIZED NUMERICAL LINEAR ALGEBRA

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

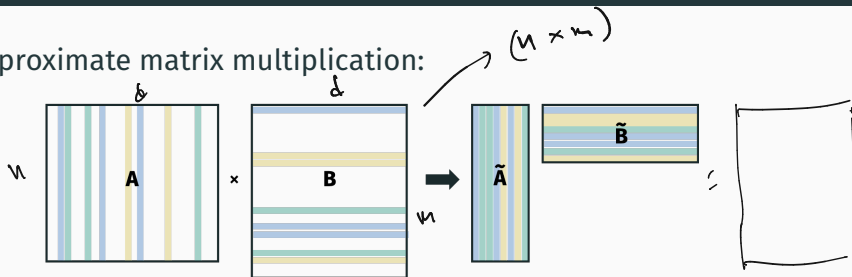
1. Compress your matrices using a randomized method.
2. Solve the problem on the smaller or sparser matrix.

• \tilde{A} called a “sketch” or “coreset” for A .

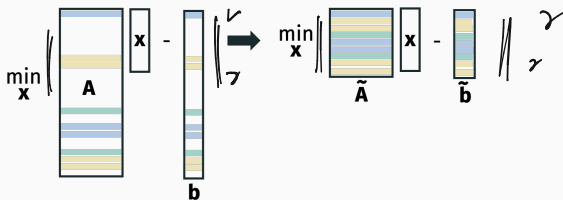


RANDOMIZED NUMERICAL LINEAR ALGEBRA

Approximate matrix multiplication:



Approximate regression:

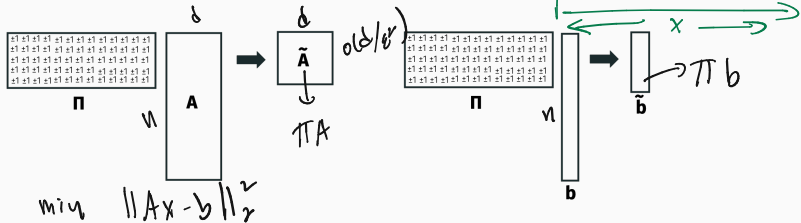


COMPARISON

	Direct Gaussian Elimination QR Algorithms	Iterative Gradient Descent + Power Method + Krylov Methods	Randomized sketch ↓ sampling, random projection
Method:			
Speed:	$O(n^3)$ for inverse, eigendecomp Slow	$O(n^2) \cdot (\# \text{ of iterations})$ Fast	$\approx O(n^2)$ + accuracy matrix pars etc
Accuracy:	Exact <u>$\log \log(1/\epsilon)$</u>	$\log(1/\epsilon)$ <u>$1/\epsilon$ $1/\sqrt{\epsilon}$</u>	Faster (sometimes) $1/\epsilon^2$ $1/\epsilon$

SKETCHED REGRESSION

Randomized approximate regression using a Johnson-Lindenstrauss Matrix:



Input: $A \in \mathbb{R}^{n \times d}$, $b \in \mathbb{R}^n$. $\tilde{x}^* = \arg \min_x \|\tilde{A}x - \tilde{b}\|_2^2$

Goal: Want $\|Ax - b\|_2^2 \leq (1 + \epsilon) \|\tilde{A}x - \tilde{b}\|_2^2$

to prove: $\|A\tilde{x}^* - b\|_2^2 \leq (1 + \epsilon) \|A\tilde{x}^* - b\|_2^2$

Claim: Suffices to prove that for all $x \in \mathbb{R}^d$, where $x^* =$

$\arg \min \|Ax - b\|_2^2$

$$(6-6) \left(\|Ax - b\|_2^2 \leq \|\Pi Ax - \Pi b\|_2^2 \leq (1 + \epsilon) \|Ax - b\|_2^2 \right)$$

Lemma (Distributional JL)

If Π is chosen to a properly scaled random Gaussian matrix, sign matrix, sparse random matrix, etc., with $O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ rows then for any fixed y ,

$$(1 - \epsilon) \|y\|_2^2 \leq \|\Pi y\|_2^2 \leq (1 + \epsilon) \|y\|_2^2$$

with probability $(1 - \delta)$.

$$y = Ax - b$$

Corollary: For any fixed x , with probability $(1 - \delta)$,

$$(1 - \epsilon) \|Ax - b\|_2^2 \leq \|\Pi Ax - \Pi b\|_2^2 \leq (1 + \epsilon) \|Ax - b\|_2^2.$$

FOR ANY TO FOR ALL

How do we go from "for any fixed \mathbf{x} " to "for all $\mathbf{x} \in \mathbb{R}^d$ ".

This statement requires establishing a Johnson-Lindenstrauss type bound for an infinity of possible vectors $(\mathbf{Ax} - \mathbf{b})$, which obviously can't be tackled with a union bound argument.

$x_1, x_2, \dots, x_{100}, \dots$

$$O\left(\frac{\log 1/\delta}{\epsilon^2}\right)$$

Then for each x_i

$$\|\mathbf{Ax}_i - \mathbf{b}\|_2 = (1 \pm \epsilon) \|\mathbf{Tx}_i - \mathbf{b}\|_2$$

with prob. $1 - \delta$,

so for x_1, \dots, x_{100} with prob. $1 - 100\delta$ this holds

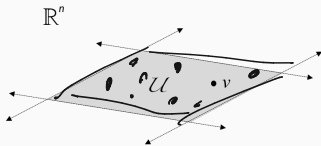
SUBSPACE EMBEDDINGS

Theorem (Subspace Embedding from JL)

Let $\mathcal{U} \subset \mathbb{R}^n$ be a d -dimensional linear subspace in \mathbb{R}^n . If $\Pi \in \mathbb{R}^{m \times d}$ is chosen from any distribution \mathcal{D} satisfying the Dist. JL Lemma, then with probability $1 - \delta$,

$$(1 - \epsilon) \|\mathbf{v}\|_2 \leq \|\Pi \mathbf{v}\|_2 \leq (1 + \epsilon) \|\mathbf{v}\|_2$$

for all $\mathbf{v} \in \mathcal{U}$, as long as $m = O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon^2}\right)^1$.



¹It's possible to obtain a slightly tighter bound of $O\left(\frac{d + \log(1/\delta)}{\epsilon^2}\right)$. It's a nice challenge to try proving this.

SUBSPACE EMBEDDING TO APPROXIMATE REGRESSION

Corollary: If we choose Π and properly scale then, with $O(d/\epsilon^2)$ rows, then with high probability,

$$(1 - \epsilon) \|Ax - b\|_2^2 \leq \|\Pi Ax - \Pi b\|_2^2 \leq (1 + \epsilon) \|Ax - b\|_2^2$$

for all x and thus

for any x can be written as $Ax - b$, want $(1 - \epsilon) \|y\|_2^2 \leq \|\Pi y\|_2^2 \leq (1 + \epsilon) \|y\|_2^2$

$$\min_x \|\Pi Ax - \Pi b\|_2^2 \leq (1 + \epsilon) \min_x \|Ax - b\|_2^2.$$

i.e. we can solve linear regression approximately using the $O(d/\epsilon^2) \times d$ matrix ΠA in place of A .

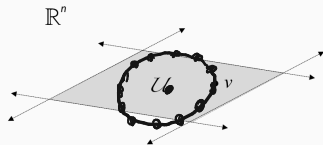
Proof: Apply Subspace Embedding Thm. to the $(d + 1)$ dimensional subspace spanned by A 's d columns and b . Every vector $Ax - b$ lies in this subspace.

Theorem (Subspace Embedding from JL)

Let $\mathcal{U} \subset \mathbb{R}^n$ be a d -dimensional linear subspace in \mathbb{R}^n . If $\Pi \in \mathbb{R}^{m \times d}$ is chosen from any distribution \mathcal{D} satisfying the Dist. JL Lemma, then with probability $1 - \delta$,

$$(1 - \epsilon) \|\mathbf{v}\|_2 \leq \|\Pi \mathbf{v}\|_2 \leq (1 + \epsilon) \|\mathbf{v}\|_2 \quad (1)$$

for all $\mathbf{v} \in \mathcal{U}$, as long as $m = O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon^2}\right)^2$.



²It's possible to obtain a slightly tighter bound of $O\left(\frac{d + \log(1/\delta)}{\epsilon^2}\right)$. It's a nice challenge to try proving this.

SUBSPACE EMBEDDING PROOF

Observation: The theorem holds as long as (1) holds for all \mathbf{w} on the unit sphere in \mathcal{U} . Denote the sphere $S_{\mathcal{U}}$:

$$S_{\mathcal{U}} = \{\mathbf{w} \mid \underline{\mathbf{w}} \in \mathcal{U} \text{ and } \|\underline{\mathbf{w}}\|_2 = 1\}.$$

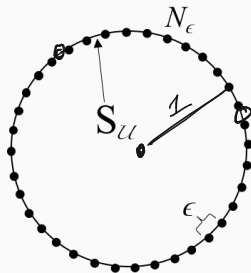
Follows from linearity: Any point $\mathbf{v} \in \mathcal{U}$ can be written as $c\mathbf{w}$ for some scalar c and some point $\mathbf{w} \in S_{\mathcal{U}}$.

- If $\underline{(1 - \epsilon)\|\mathbf{w}\|_2} \leq \|\Pi\mathbf{w}\|_2 \leq \underline{(1 + \epsilon)\|\mathbf{w}\|_2}$.
- then $c(1 - \epsilon)\|\mathbf{w}\|_2 \leq c\|\Pi\mathbf{w}\|_2 \leq c(1 + \epsilon)\|\mathbf{w}\|_2$,
- and thus $\underline{(1 - \epsilon)\|c\mathbf{w}\|_2} \leq \|\Pi c\mathbf{w}\|_2 \leq \underline{(1 + \epsilon)\|c\mathbf{w}\|_2}$.

$c \cdot \mathbf{w}$
 $\mathbf{v} = \|\mathbf{v}\|_2 \cdot \frac{\mathbf{v}}{\|\mathbf{v}\|_2}$
 $\in S_{\mathcal{U}}$

SUBSPACE EMBEDDING PROOF

Intuition: There are not too many “different” points on a d -dimensional sphere:



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

$$O\left(\frac{1}{\epsilon^2}\right)$$

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

N_ϵ is called an “ ϵ ”-net.

If we can prove

$$\|w\| (1 - \epsilon) \leq \|\Pi w\|_2 \leq (1 + \epsilon) \|w\|_r$$

for all points $w \in \underline{N_\epsilon}$, we can hopefully extend to all of S_U .

Lemma (ϵ -net for the sphere)

For any $\epsilon \leq 1$, there exists a set $N_\epsilon \subset S_U$ with $|N_\epsilon| = \left(\frac{4}{\epsilon}\right)^d$ such that $\forall v \in S_U$,

$$\min_{w \in N_\epsilon} \|v - w\| \leq \epsilon.$$

Want (i) $\rightarrow (1+\epsilon)\|w\|_2 \leq \|\Pi w\|_2 \leq \|w\|_2 \cdot (1+\epsilon)$
to hold for all $w \in N_\epsilon$.

$$\text{Set } \delta = \frac{1}{|N_\epsilon|} \quad \log(1/\delta) \quad \log(|N_\epsilon|)$$
$$= \log\left(\frac{4}{\epsilon}\right)^d$$
$$= d \log\left(\frac{4}{\epsilon}\right)$$

1. Set $\delta = \left(\frac{\epsilon}{8}\right)^d$. By a union bound, with high probability, for all $\mathbf{w} \in N_\epsilon$,

$$(1 - \epsilon) \leq \|\Pi\mathbf{w}\|_2 \leq (1 + \epsilon).$$

as long as Π has $O\left(\frac{\log(1/\delta)}{\epsilon^2}\right) = O\left(\frac{d \log(1/\epsilon)}{\epsilon^2}\right)$ rows.

2. Consider any $\mathbf{v} \in S_{\mathcal{U}}$. You can check that, for some $\mathbf{w}_0, \mathbf{w}_1, \mathbf{w}_2 \dots \in N_{\epsilon}$, \mathbf{v} can be written:

$$\mathbf{v} = \mathbf{w}_0 + c_1 \mathbf{w}_1 + c_2 \mathbf{w}_2 + \dots$$

for constants c_1, c_2, \dots where $|c_i| \leq \epsilon^i$.

3. Applying triangle inequality, we have

$$\begin{aligned} \|\Pi \mathbf{v}\|_2 &= \|\Pi \mathbf{w}_0 + c_1 \Pi \mathbf{w}_1 + c_2 \Pi \mathbf{w}_2 + \dots\| \\ &\leq \|\Pi \mathbf{w}_0\| + \epsilon \|\Pi \mathbf{w}_1\| + \epsilon^2 \|\Pi \mathbf{w}_2\| + \dots \\ &\leq (1 + \epsilon) + \epsilon(1 + \epsilon) + \epsilon^2(1 + \epsilon) + \dots \\ &\leq 1 + O(\epsilon). \end{aligned}$$

4. Similarly,

$$\begin{aligned}
 \|\Pi v\|_2 &= \|\Pi w_0 + c_1 \Pi w_1 + c_2 \Pi w_2 + \dots\| \\
 &\geq \|\Pi w_0\| - \epsilon \|\Pi w_1\| - \epsilon^2 \|\Pi w_2\| - \dots \\
 &\geq (1 - \epsilon) - \epsilon(1 + \epsilon) - \epsilon^2(1 + \epsilon) - \dots \\
 &\geq 1 - O(\epsilon).
 \end{aligned}$$

So we have proven

$$1 - O(\epsilon) \leq \|\Pi v\|_2 \leq 1 + O(\epsilon)$$

for all v in $S_{\mathcal{U}}$.

Adjusting ϵ proves the Subspace Embedding theorem.