

## CS-GY 9223 I: Lecture 10

Krylov methods, spectral clustering, spectral graph theory.

---

NYU Tandon School of Engineering, Prof. Christopher Musco

# COMPUTATION IN LINEAR ALGEBRA

$$\min_{\underline{x}} \|Ax - b\|_2^2$$

Three classes of methods.

$\rightarrow n \times n$  matrix

- Direct Methods:

Exact computation of  $A^{-1}$ : Gaussian elimination

QR Algorithm:  $O(n^3)$  for SVD  
eigendecomposition  $O(n^3)$

- Iterative Methods:

Power Method: Top singular vectors.

$\|Ax - b\|_2^2$ : Gradient Descent:

Matrix-Vector product w/  $A$ .  $\rightarrow O(nd)$

- Randomized Methods:

Methods based on JL.

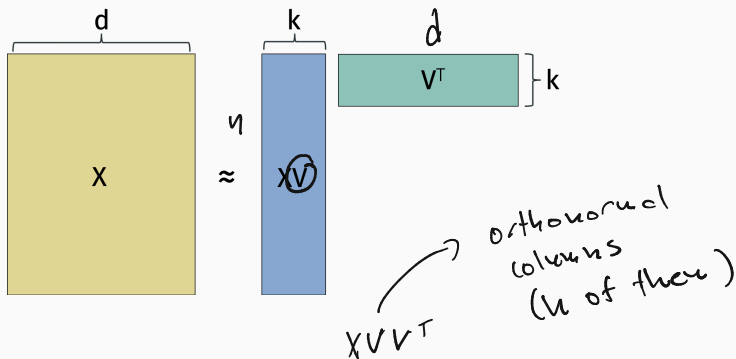
$\text{nnz}(A)$  = "number of nonzeros in  $A$ "

SGD, SCD

In general computing  $Ax$  takes  $O(\text{nnz}(A))$  time  $\leq O(nd)$

## LOW-RANK APPROXIMATION

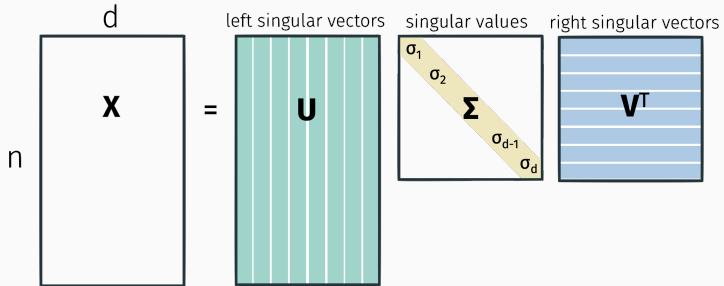
Write  $X$  as a rank  $k$  factorization by projecting onto the subspace spanned by an orthonormal matrix  $V \in \mathbb{R}^{d \times k}$



# SINGULAR VALUE DECOMPOSITION

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  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d \geq 0$ .

## COMPUTATIONAL QUESTION

Given a subspace  $\mathcal{V}$  spanned by the  $k$  columns in  $\mathbf{V}$ ,

$$\|\mathbf{X} - \mathbf{X}\mathbf{V}\mathbf{V}^T\|_F^2 = \min_{\mathbf{C}} \|\mathbf{X} - \mathbf{C}\mathbf{V}^T\|_F^2$$

We want to find the best  $\mathbf{V} \in \mathbb{R}^{d \times k}$ :

$$\min_{\text{orthonormal } \mathbf{V} \in \mathbb{R}^{d \times k}} \|\mathbf{X} - \mathbf{X}\mathbf{V}\mathbf{V}^T\|_F^2 \quad (1)$$

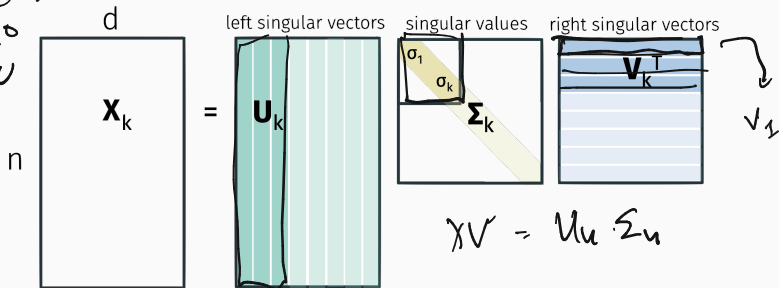
Note that  $\|\mathbf{X} - \mathbf{X}\mathbf{V}\mathbf{V}^T\|_F^2 = \|\mathbf{X}\|_F^2 - \|\mathbf{X}\mathbf{V}\mathbf{V}^T\|_F^2$  for all orthonormal  $\mathbf{V}$  (since  $\mathbf{V}\mathbf{V}^T$  is a projection). Equivalent form:

$$\max_{\text{orthonormal } \mathbf{V} \in \mathbb{R}^{d \times k}} \|\mathbf{X}\mathbf{V}\mathbf{V}^T\|_F^2 = \|\mathbf{X}\mathbf{V}\|_F^2 \quad (2)$$

# SINGULAR VALUE DECOMPOSITION

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

*$O(nd^2)$   
time to  
compute*



$$X_k = U_k U_k^T X = X V_k V_k^T$$

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

# POWER METHOD

Goal: Find some  $\mathbf{z} \approx \mathbf{v}_1$ .

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

Power method:

• Choose  $\mathbf{z}^{(0)}$  randomly. E.g.  $\mathbf{z}_0 \sim \mathcal{N}(0, 1)$ .

• For  $i = 1, \dots, T$

•  $\mathbf{z}^{(i)} = \mathbf{X}^T \cdot (\mathbf{X}\mathbf{z}^{(i-1)})$

•  $n_i = \|\mathbf{z}^{(i)}\|_2$

•  $\mathbf{z}^{(i)} = \mathbf{z}^{(i)} / n_i$

Return  $\mathbf{z}_T$

"number of non-zeros in  $\mathbf{X}^T$ "

$\text{nnz}(\mathbf{X})$  time

$\text{nnz}(\mathbf{X}) + d$

## POWER METHOD CONVERGENCE

### 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:

$$\|v_1 - z^{(T)}\|_2 \leq \epsilon.$$

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

↓  
cost per  
iterations



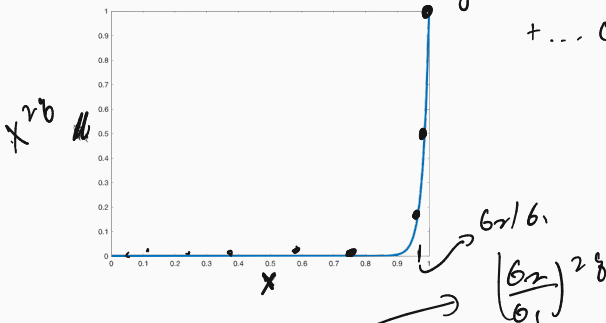
# KRYLOV SUBSPACE METHODS

svd's  
eigs

$q$  = # of iterations  
random gaussian

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

$$g = c_1 v_1 + c_2 v_2 + \dots + c_d v_d$$



$$\underline{z^{(q)}} = c \cdot \left[ c_1 \cdot \frac{\sigma_1^{2q}}{\sigma_1^{2q}} v_1 + c_2 \cdot \frac{\sigma_2^{2q}}{\sigma_1^{2q}} v_2 + \dots + c_d \cdot \frac{\sigma_d^{2q}}{\sigma_1^{2q}} v_d \right]$$

$$X^T \left( X \left( X^T (X g) \right) \right)$$

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

Along the way we computed:

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

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

**Idea behind Krylov methods:** Don't throw away everything before  $(X^T X)^q \cdot g$ . What you're using when you run `svds` or `eigs` in MATLAB or Python.

# KRYLOV SUBSPACE METHODS

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

Lanczos method:  $\mathcal{K} = \begin{array}{|c|} \hline \text{ } \\ \hline \end{array} \begin{array}{c} d \\ q \end{array} \rightarrow O(dq^2)$

- Let  $\mathbf{Q} \in \mathbb{R}^{d \times q}$  be an orthonormal span for the vectors in  $\mathcal{K}_q$
- Solve  $\min_{\mathbf{v}=\mathbf{Q}\mathbf{w}} \|\mathbf{X} - \mathbf{X}\mathbf{v}\mathbf{v}^T\|_F^2 \rightarrow$  need to compute SVD
  - Find best vector in the Krylov subspace, instead of just using last vector.  $\mathbf{Q}^T \mathbf{X}^T \mathbf{X} \mathbf{Q}$
  - Can be done in  $O(\text{nnz}(\mathbf{X}) \cdot q + dq^2)$  time.

$\min_{\mathbf{v}} \|\mathbf{X} - \mathbf{X}\mathbf{v}\mathbf{v}^T\|_F^2 \rightarrow$  optimal rank 1 approximation  
↓  
Requires SVD.

# LANCZOS METHOD ANALYSIS

**Claim 1:** For any degree  $q$  polynomial  $p$ , we can write  $p(X^T X) \cdot \underline{g}$  as Qw for some  $w$ .

$$\left( c_1 \underbrace{(X^T X)} + c_2 \underbrace{(X^T X)^2} + \dots + c_q \underbrace{(X^T X)^q} \right) \underline{g} = p(X^T X) \underline{g}$$

**Claim 2:**

$$\min_{v=Qw} \|X - Xvv^T\|_F^2 = \min_{\text{degree } q \text{ polynomial } p} \|X - Xv_p v_p^T\|_F^2$$

where  $v_p$  =  $p(X^T X) \cdot \underline{g}$ .

**Claim 3:**

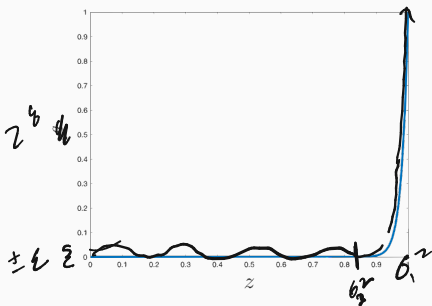
$$\underline{p(X^T X)} \underline{g} = c \cdot \left[ c_1 \cdot \frac{p(\sigma_1^2) v_1}{p(b_1^2)} + c_2 \cdot \frac{p(\sigma_2^2) v_2}{p(b_2^2)} + \dots + c_n \cdot p(\sigma_n^2) v_n \right]$$

$$(X^T X)^q \underline{g} = c \cdot \left[ c_1 b_1^{2q} v_1 + \dots + c_n b_n^{2q} v_n \right]$$

# LANCZOS METHOD ANALYSIS

**Claim:** There is an  $O\left(\sqrt{q \log \frac{1}{\epsilon}}\right)$  degree polynomial  $\hat{p}$  approximating  $x^q$  up to error  $\epsilon$  on  $[0, \sigma_1^2]$ .

$$p(z) = z^q$$



$\rightarrow \approx z$   
 $\frac{p(b_1^2)}{p(b_1)}$   
 $\frac{p(b_2)}{p(b_2)}$   
 $\dots$   
 $\frac{p(b_d^2)}{p(b_1)^r}$   
 $\rightarrow$  close to zero

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

**Runtime:**  $O\left(\frac{\log(d/\epsilon)}{\sqrt{\gamma}} \cdot \text{nnz}(X)\right)$  vs.  $O\left(\frac{\log(d/\epsilon)}{\gamma} \cdot \text{nnz}(X)\right)$

## POWER METHOD – NO GAP DEPENDENCE

Convergence is slow when  $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$  is small.  $z^{(q)}$  has large components of both  $v_1$  and  $v_2$ . But in this case:



$$\|X - Xv_1v_1^T\|_F^2 = \sum_{i \neq 1} \sigma_i^2 \approx \sum_{i \neq 2} \sigma_i^2$$

$\nearrow \leq 6_1^2 - 6_2^2$   
 $6_1^2 = \|X - Xv_2v_2^T\|_F^2$   
 $\rightarrow \leq 6_1^2 - 6_2^2$

So we don't care! Either  $v_1$  or  $v_2$  give good rank-1 approximations.

**Claim:** To achieve

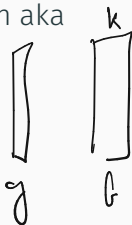
$$\|X - Xzz^T\|_F^2 \leq (1 + \epsilon)\|X - Xv_1v_1^T\|_F^2$$

we need  $O\left(\frac{\log(d/\epsilon)}{\epsilon}\right)$  power method iterations or  $O\left(\frac{\log(d/\epsilon)}{\sqrt{\epsilon}}\right)$  Lanczos iterations.

## GENERALIZATIONS TO LARGER $k$

$$(X - XV_1V_1^T) \xrightarrow{\text{Power Method}} \approx V_2$$

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



- Block Krylov methods

- Let  $G \in \mathbb{R}^{d \times k}$  be a random Gaussian matrix.

$$\mathcal{K}_q = [G, (X^T X) \cdot G, (X^T X)^2 \cdot G, \dots, (X^T X)^q \cdot G]$$

Orthogonalize  $\rightarrow \approx V_k$

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

$$[(X^T X)g \quad (X^T X)^2 g \quad \dots \quad ]$$

$$[G, \underline{(X^T X) \cdot G}]$$

$$[G, \text{orth}((X^T X) \cdot G)] \quad 15$$

# RANDOMIZED METHODS

What do you think a stochastic version of Krylov subspace method would look like?

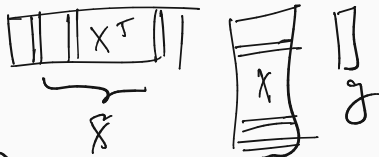
**Gauss Method**

$$\mathcal{K}_q = \left[ \mathbf{g}, (X^T X) \cdot \mathbf{g}, (X^T X)^2 \cdot \mathbf{g}, \dots, (X^T X)^{q-1} \cdot \mathbf{g} \right]$$

$$\underbrace{2A^T A x - A^T b}_{\text{Krylov subspace method}}$$

$$X^T X \mathbf{g}$$

$$O(nd) \quad O(\text{nnz}(X))$$



$$\tilde{X}^T \tilde{X} \mathbf{g}$$

$$X^T X \mathbf{g} \approx X_i X_i^T \mathbf{g}$$

↓  
one row from X







$$A^T A$$

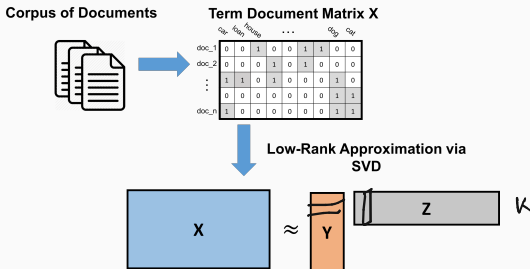
$$A^{-1} x$$

$$A x$$

Applications of (partial) singular value decomposition:

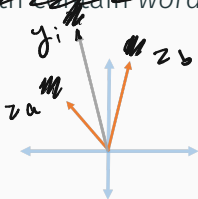
- Low-rank approximation (data compression)
- Denoising, in-painting, matrix completion
- Semantic embeddings

# EXAMPLE: LATENT SEMANTIC ANALYSIS



- $\langle \vec{y}_i, \vec{z}_a \rangle \approx 1$  when  $doc_i$  contains word  $a$ .
- If  $doc_i$  and  $doc_j$  both contain word  $a$ ,  $\langle \vec{y}_i, \vec{z}_a \rangle \approx \langle \vec{y}_j, \vec{z}_a \rangle = 1$ .

If  $doc_i$  contains word  $a$  and word  $b$ :

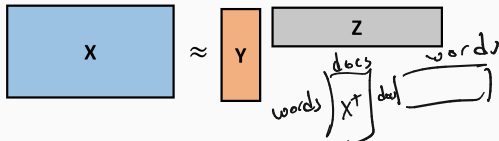


# EXAMPLE: LATENT SEMANTIC ANALYSIS

Term Document Matrix X

	car	loan	house	...	dog	cat
doc_1	0	0	1	0	0	1
doc_2	0	0	0	1	0	1
...	1	1	0	1	0	0
doc_n	1	0	0	0	0	1

Low-Rank Approximation via SVD

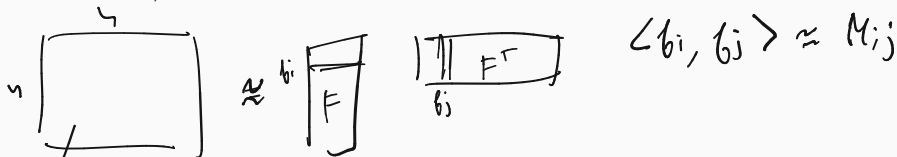


- The columns  $\vec{z}_1, \vec{z}_2, \dots$  give representations of words, with  $\vec{z}_i$  and  $\vec{z}_j$  tending to have high dot product if  $word_i$  and  $word_j$  appear in many of the same documents.
- Z corresponds to the top  $k$  right singular vectors: the eigenvectors of ~~the matrix~~. Intuitively, what is ~~the matrix~~?

$(X^T X)_{i,j} =$ 
 $(X^T X)_{i,j} =$  # of documents both word  $i$  and  $j$  appear in

## EXAMPLE: WORD EMBEDDING

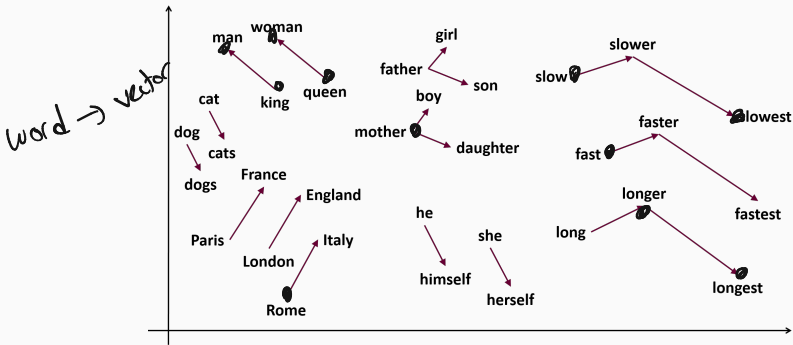
Not obvious how to convert a word into a feature vector that captures the meaning of that word. Approach suggested by LSA: build a  $d \times d$  symmetric “similarity matrix”  $\underline{M}$  between words, and factorize:  $\mathbf{M} \approx \mathbf{F}\mathbf{F}^T$  for rank  $k$   $\mathbf{F}$ .



similarity  
between  
words

- **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  $\mathbf{X}\mathbf{X}^T$  with these different metrics (sometimes appropriately transformed) leads to popular word embedding algorithms: **word2vec**, **GloVe**, etc.

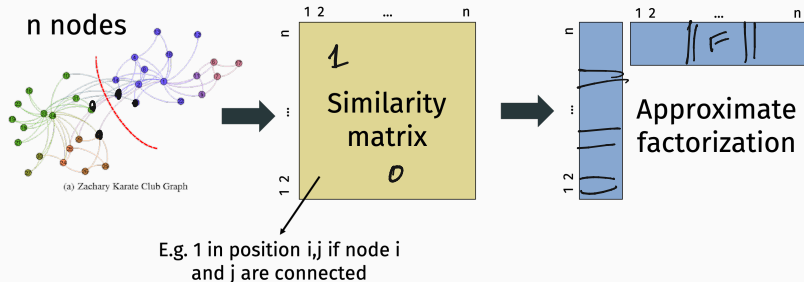
## EXAMPLE: ORD EMBEDDING



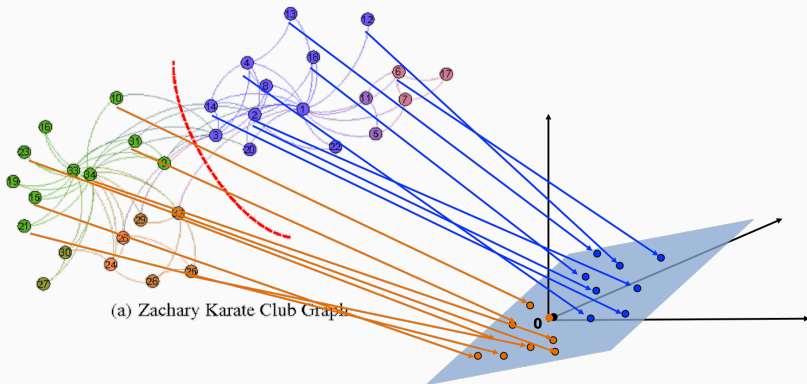
word2vec was originally described as a neural-network method, but Levy and Goldberg show that it is simply low-rank approximation of a specific similarity matrix. *Neural word embedding as implicit matrix factorization.*

# ENCODING GRAPH SIMILARITY

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



# ENCODING GRAPH SIMILARITY

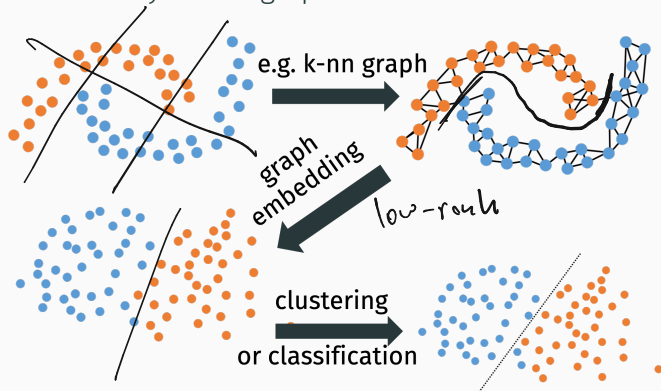


**Social networks in 1970:** “The network captures 34 members of a karate club, documenting links between pairs of members who interacted outside the club. During the study a conflict arose between the administrator “John A” and instructor “Mr. Hi” (pseudonyms), which led to the split of the club into 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



## SPECTRAL CLUSTERING

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



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

Spectral graph theory lets us formalize this heuristic idea.



Laplacians  $\rightarrow$  matrix representation of  
a graph

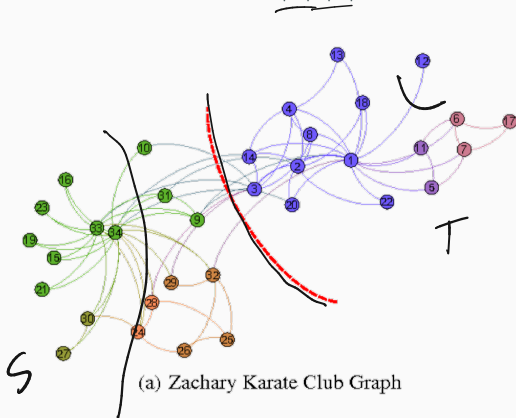
Adjacency Matrix

Powers of Adjacency matrix

# CUT MINIMIZATION

**Goal:** Partition nodes along a cut that:

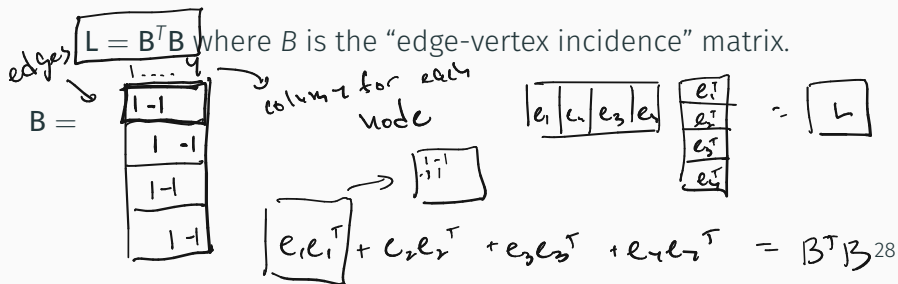
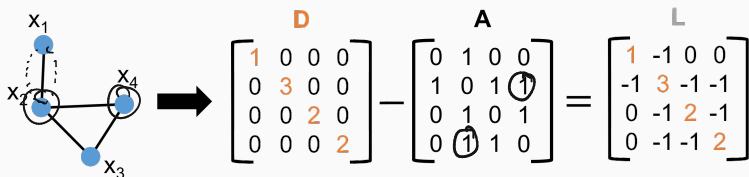
- Has few crossing edges:  $|\{(u, v) \in E : u \in S, v \in T\}|$  is small.
- Separates large partitions:  $|S|, |T|$  are not too small.



(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**.



# THE LAPLACIAN VIEW

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

Conclusions from  $L = B^T B$

•  $L$  is positive semidefinite:  $x^T L x \geq 0$  for all  $x$ .

•  $L = V \Sigma^2 V^T$  where  $U \Sigma V^T$  is  $B$ 's SVD. Columns of  $V$  are eigenvectors of  $L$ .

$$L = \underbrace{V U^T U}_{I} V^T = V \Sigma^2 V^T$$

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

$$T: V \setminus S$$

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

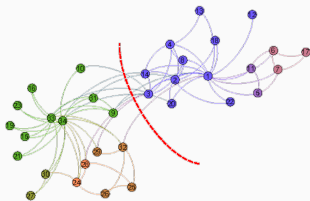
$$c^T B^T B c = \|B c\|_2^2$$

# of edges between  $S$  and  $T$



0 if edge is "inside"  $S$  or  $T$

2 otherwise



(a) Zachary Karate Club Graph

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

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

- $\underline{c}^T \mathbf{1} = |T| - |S|$ .

want small  $\rightarrow$  small cut

$$|\underline{c}^T \mathbf{1}| = ||T| - |S||$$

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

## SMALLEST LAPLACIAN EIGENVECTOR

### Courant-Fischer min-max principle

Let  $V = [v_1, \dots, v_n]$  be the <sup>singular vector</sup> eigenvectors of  $L$ .

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

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

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

$\vdots$

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

$$\begin{aligned} & \|L V V^T\|_F^2 \\ &= \|L\|_F^2 - \|L - L V V^T\|_F^2 \end{aligned}$$

## Courant–Fischer min-max principle

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

$$\begin{aligned}
 \mathbf{v}_n &= \arg \min_{\|\mathbf{v}\|=1} \mathbf{v}^T \mathbf{L} \mathbf{v} \\
 \mathbf{v}_{n-1} &= \arg \min_{\|\mathbf{v}\|=1, \mathbf{v} \perp \mathbf{v}_n} \mathbf{v}^T \mathbf{L} \mathbf{v} \\
 \mathbf{v}_{n-2} &= \arg \min_{\|\mathbf{v}\|=1, \mathbf{v} \perp \mathbf{v}_n, \mathbf{v}_{n-1}} \mathbf{v}^T \mathbf{L} \mathbf{v} \\
 &\vdots \\
 \mathbf{v}_1 &= \arg \min_{\|\mathbf{v}\|=1, \mathbf{v} \perp \mathbf{v}_n, \dots, \mathbf{v}_2} \mathbf{v}^T \mathbf{L} \mathbf{v}
 \end{aligned}$$



## SMALLEST LAPLACIAN EIGENVECTOR

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

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

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

$$L \mathbf{v}_n = \vec{0}$$

$$\mathbf{v}_n = \begin{bmatrix} \frac{1}{\sqrt{n}} \\ \vdots \\ \frac{1}{\sqrt{n}} \end{bmatrix}$$

## SECOND SMALLEST LAPLACIAN EIGENVECTOR

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

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

*scaled all ones*

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

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

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

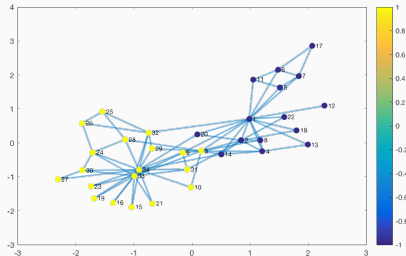
$\mathbf{v}_{n-1} \in \mathbb{R}^n$  is not generally binary, but still satisfies a ‘relaxed’ version of this property.

# CUTTING WITH THE SECOND LAPLACIAN EIGENVECTOR

Find a good partition of the graph by computing

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

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

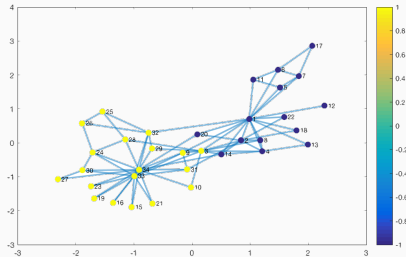


# CUTTING WITH THE SECOND LAPLACIAN EIGENVECTOR

Find a good partition of the graph by computing

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

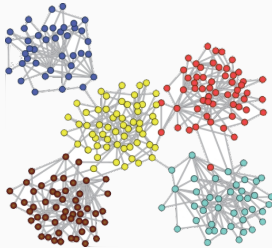
Set  $S$  to be all nodes with  $\mathbf{v}_{n-1}(i) < 0$ , and  $T$  to be all with  $\mathbf{v}_{n-1}(i) \geq 0$ .



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

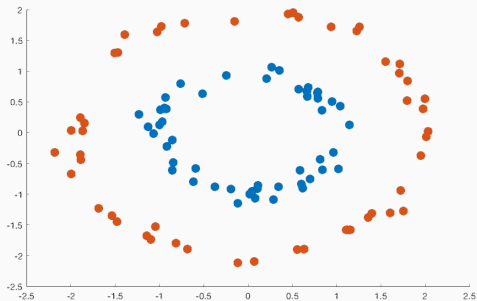
**Important consideration:** What to do when we want to split the graph into more than two parts?

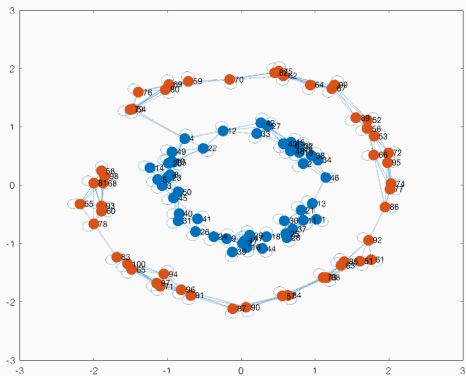


## Spectral Clustering:

- 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).

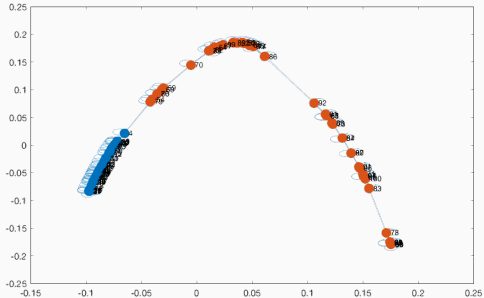
Original Data: (not linearly separable)



$k$ -Nearest Neighbors Graph:



Embedding with eigenvectors  $\mathbf{v}_{n-1}, \mathbf{v}_{n-2}$ : (linearly separable)



**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:** Give a natural **generative model** for which produces random but realistic inputs and analyze how the algorithm performs on inputs drawn from this model.

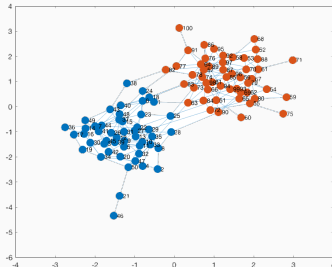
- Very common in algorithm design for data analysis/machine learning (can be used to justify  $\ell_2$  linear regression,  $k$ -means clustering, PCA, etc.)

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

## Stochastic Block Model (Planted Partition Model):

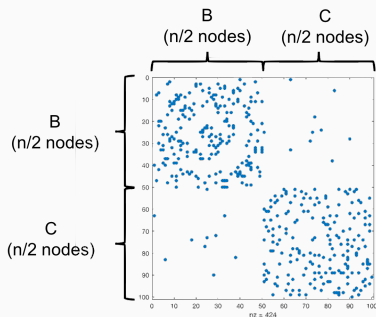
Let  $G_n(p, q)$  be a distribution over graphs on  $n$  nodes, split equally into two groups  $B$  and  $C$ , each with  $n/2$  nodes.

- 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.  $q < p$ .



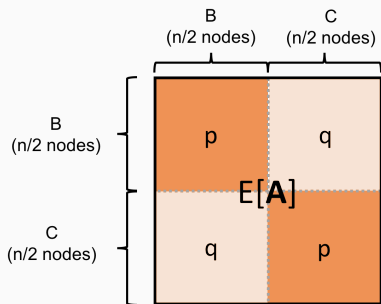
Let  $G$  be a stochastic block model graph drawn from  $G_n(p, q)$ .

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



## EXPECTED ADJACENCY SPECTRUM

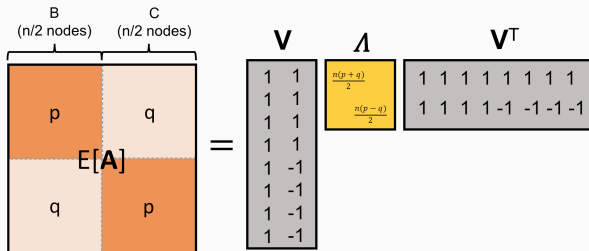
Letting  $G$  be a stochastic block model graph drawn from  $G_n(p, q)$  and  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be its adjacency matrix.  $(\mathbb{E}[\mathbf{A}])_{i,j} = p$  for  $i, j$  in same group,  $(\mathbb{E}[\mathbf{A}])_{i,j} = q$  otherwise.



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

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

## EXPECTED ADJACENCY SPECTRUM



- $\vec{v}_1 = \vec{1}$  with eigenvalue  $\lambda_1 = \frac{(p+q)n}{2}$ .
- $\vec{v}_2 = \chi_{B,C}$  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  $\vec{v}_2$  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  $G$  (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.
- Random matrix theory is a very recent and cutting edge subfield of mathematics that is being actively applied in computer science, statistics, and machine learning.

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

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

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

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

**Exercise:** Show that  $\|\mathbf{X}\|_2$  is equal to the largest singular value of  $\mathbf{X}$ . For symmetric  $\mathbf{X}$  (like  $\mathbf{A} - \mathbb{E}[\mathbf{A}]$ ) 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  $\mathbf{A}$  and  $\mathbb{E}[\mathbf{A}]$  are close. How does this relate to their difference in spectral norm?

**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  $v_1, v_2, \dots, v_d$  and  $\bar{v}_1, \bar{v}_2, \dots, \bar{v}_d$ . Letting  $\theta(v_i, \bar{v}_i)$  denote the angle between  $v_i$  and  $\bar{v}_i$ , for all  $i$ :

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

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

$$\begin{array}{c} \mathbf{A} \\ \boxed{\begin{array}{cc} 1+\varepsilon & 0 \\ 0 & 1 \end{array}} \end{array} - \begin{array}{c} \bar{\mathbf{A}} \\ \boxed{\begin{array}{cc} 1 & 0 \\ 0 & 1+\varepsilon \end{array}} \end{array} = \begin{array}{c} \mathbf{A}-\bar{\mathbf{A}} \\ \boxed{\begin{array}{cc} \varepsilon & 0 \\ 0 & \varepsilon \end{array}} \end{array}$$

## APPLICATION TO STOCHASTIC BLOCK MODEL

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

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

**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:**  $\mathbb{E}[\mathbf{A}]$ , has eigenvalues  $\lambda_1 = \frac{(p+q)n}{2}$ ,  $\lambda_2 = \frac{(p-q)n}{2}$ ,  $\lambda_i = 0$  for  $i \geq 3$ .

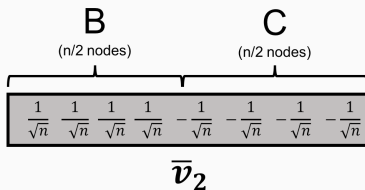
$$\min_{j \neq i} |\lambda_i - \lambda_j| = \min\left(qn, \frac{(p-q)n}{2}\right).$$

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

## APPLICATION TO STOCHASTIC BLOCK MODEL

So Far:  $\sin \theta(v_2, \bar{v}_2) \leq O\left(\frac{\sqrt{p}}{(\rho-q)\sqrt{n}}\right)$ . What does this give us?

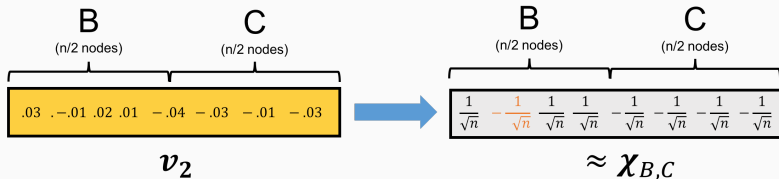
- Can show that this implies  $\|v_2 - \bar{v}_2\|_2^2 \leq O\left(\frac{p}{(\rho-q)^2 n}\right)$  (exercise).
- $\bar{v}_2$  is  $\frac{1}{\sqrt{n}}\chi_{B,C}$ : the community indicator vector.



- 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  $O\left(\frac{p}{(\rho-q)^2}\right)$  positions.

## APPLICATION TO STOCHASTIC BLOCK MODEL

**Upshot:** If  $G$  is a stochastic block model graph with adjacency matrix  $\mathbf{A}$ , 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.