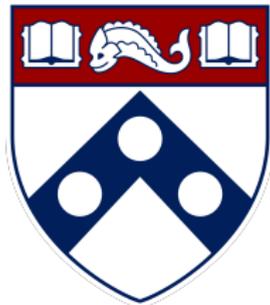
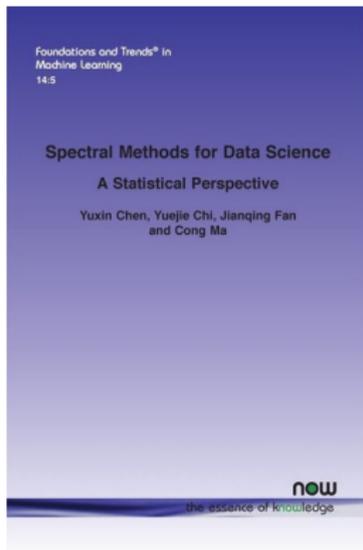


Spectral methods for data science: A statistical perspective



Yuxin Chen, Statistics & Data Science, UPenn

Summer school on theoretical stats., PKU BICMR & Math, 2023



“Spectral methods for data science: a statistical perspective,” Yuxin Chen, Yuejie Chi, Jianqing Fan, Cong Ma, Foundations and Trends in Machine Learning, 2021

Part 1: Introduction

- Motivating applications
 - community detection
 - matrix/tensor completion
 - ranking
- A general recipe for spectral methods

Motivating application: community detection

Graph clustering / community detection

Community structures are common in many social networks



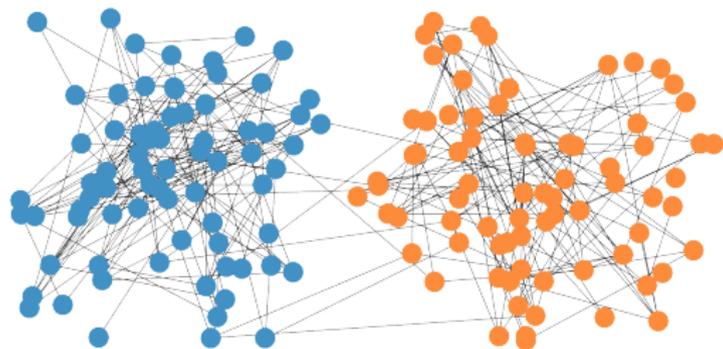
figure credit: The Future Buzz



figure credit: S. Papadopoulos

Goal: partition users into several clusters based on their friendships / similarities

An idealistic model: stochastic block model (SBM)

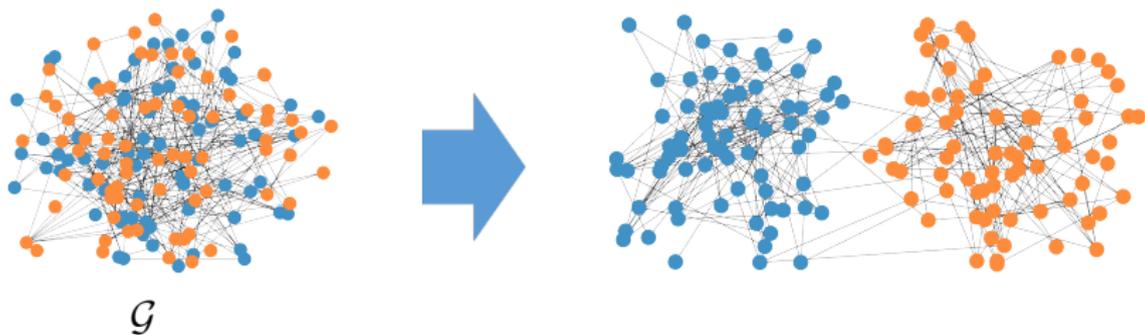


$x_i = 1$: 1st community

$x_i = -1$: 2nd community

- n nodes $\{1, \dots, n\}$
- 2 communities
- n unknown variables: $x_1, \dots, x_n \in \{1, -1\}$
 - encode community memberships

An idealistic model: stochastic block model (SBM)

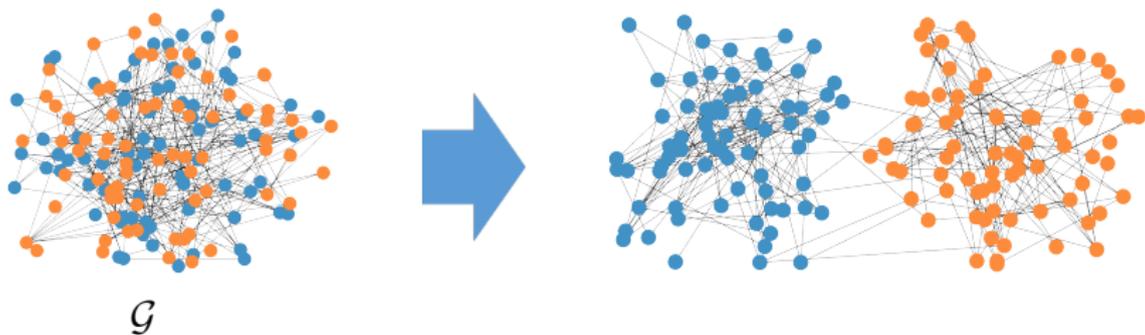


- **observation:** a (random) graph \mathcal{G}

$$(i, j) \in \mathcal{G} \text{ with prob. } \begin{cases} p, & \text{if } i \text{ and } j \text{ are from same community} \\ q, & \text{else} \end{cases}$$

- $p > q$ (i.e. more within-cluster edges than between-cluster edges)

An idealistic model: stochastic block model (SBM)



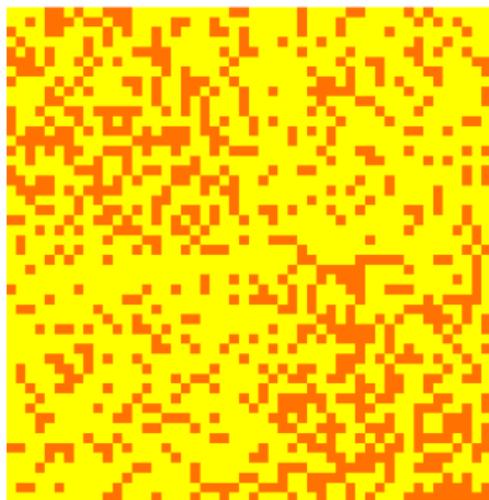
- **observation:** a (random) graph \mathcal{G}

$$(i, j) \in \mathcal{G} \text{ with prob. } \begin{cases} p, & \text{if } i \text{ and } j \text{ are from same community} \\ q, & \text{else} \end{cases}$$

◦ $p > q$ (i.e. more within-cluster edges than between-cluster edges)

- **goal:** recover community memberships of all nodes, i.e. $\{x_i^*\}$

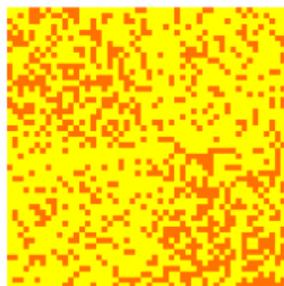
Key structure of adjacency matrix



The adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$ of \mathcal{G} :

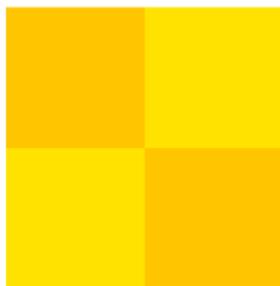
$$A_{i,j} = \begin{cases} 1, & \text{if } (i, j) \in \mathcal{G} \\ 0, & \text{else} \end{cases}$$

Key structure of adjacency matrix



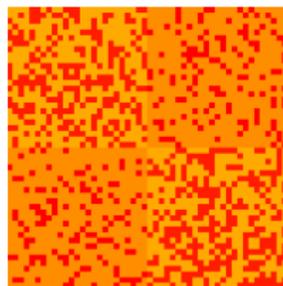
\mathbf{A}

=



$\underbrace{\mathbb{E}[\mathbf{A}]}_{\text{rank 2}}$

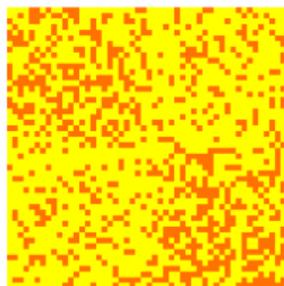
+



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

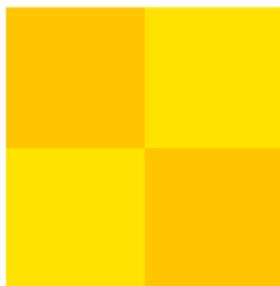
WLOG, suppose $x_1^* = \dots = x_{n/2}^* = 1$, $x_{n/2+1}^* = \dots = x_n^* = -1$:

Key structure of adjacency matrix



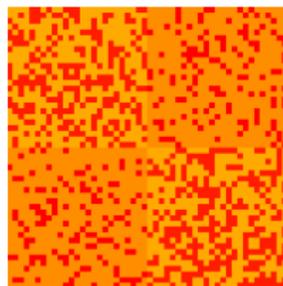
\mathbf{A}

=



$\underbrace{\mathbb{E}[\mathbf{A}]}_{\text{rank 2}}$

+



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

WLOG, suppose $x_1^* = \dots = x_{n/2}^* = 1$, $x_{n/2+1}^* = \dots = x_n^* = -1$:

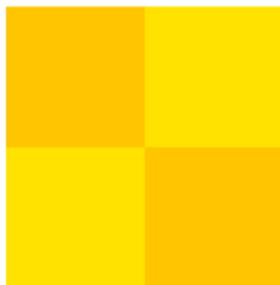
$$\mathbb{E}[\mathbf{A}] = \begin{bmatrix} p\mathbf{1}\mathbf{1}^\top & q\mathbf{1}\mathbf{1}^\top \\ q\mathbf{1}\mathbf{1}^\top & p\mathbf{1}\mathbf{1}^\top \end{bmatrix} = \underbrace{\frac{p+q}{2}\mathbf{1}\mathbf{1}^\top}_{\text{uninformative bias}} + \frac{p-q}{2} \underbrace{\begin{bmatrix} \mathbf{1} \\ -\mathbf{1} \end{bmatrix}}_{=: \mathbf{x}^* = [x_i^*]_{1 \leq i \leq n}} [\mathbf{1}^\top, -\mathbf{1}^\top]$$

Spectral clustering



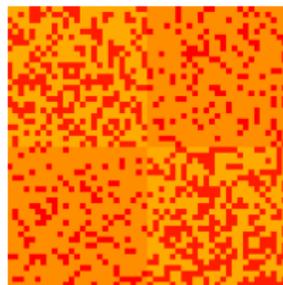
\mathbf{A}

=



$\underbrace{\mathbb{E}[\mathbf{A}]}_{\text{rank 2}}$

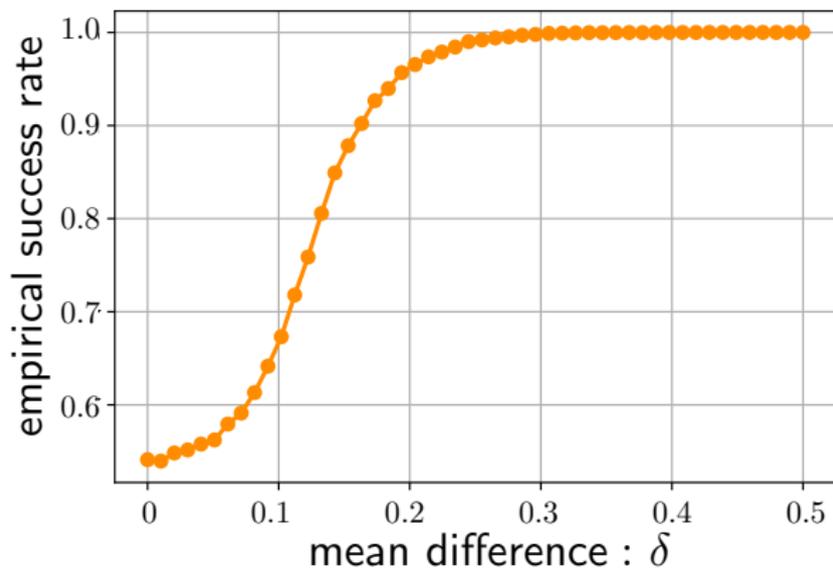
+



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

1. computing leading eigenvector $\mathbf{u} = [u_i]_{1 \leq i \leq n}$ of $\mathbf{A} - \frac{p+q}{2} \mathbf{1}\mathbf{1}^\top$
2. rounding: output $x_i = \begin{cases} 1, & \text{if } u_i > 0 \\ -1, & \text{if } u_i < 0 \end{cases}$

Empirical clustering accuracy



$$n = 100, p = \frac{1+\delta}{2}, q = \frac{1-\delta}{2}$$

Rationale: spectral clustering is reliable if $\underbrace{\mathbf{A} - \mathbb{E}[\mathbf{A}]}_{\text{perturbation}}$ is “small”

- if $\mathbf{A} - \mathbb{E}[\mathbf{A}] = \mathbf{0}$, then

$$\mathbf{u} \propto \pm \begin{bmatrix} \mathbf{1} \\ -\mathbf{1} \end{bmatrix} \implies \text{perfect clustering}$$

What we'll demonstrate: effect of perturbation $\mathbf{A} - \mathbb{E}[\mathbf{A}]$ on \mathbf{u}

Motivating application: matrix/tensor completion

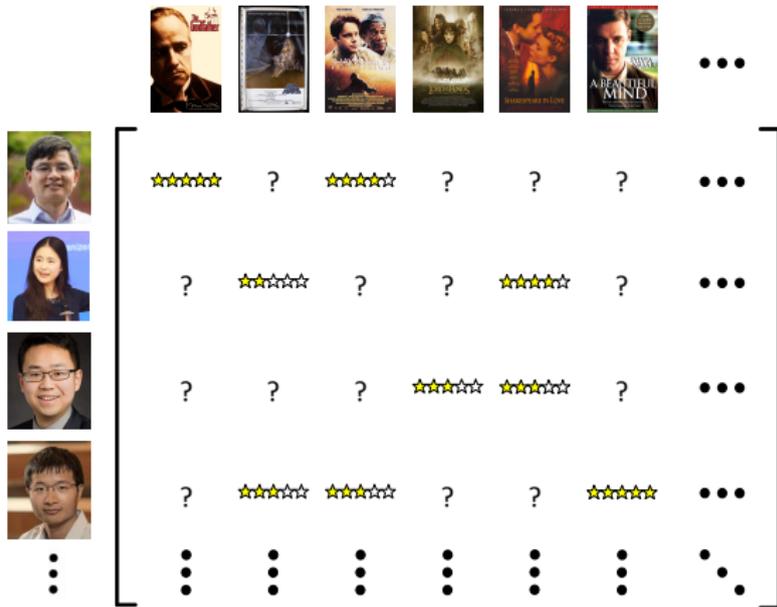


figure credit: Candes et al.

- Netflix challenge: Netflix provides highly incomplete ratings from 0.5 million users for & 17,770 movies
- How to predict unseen user ratings for movies?

Matrix completion

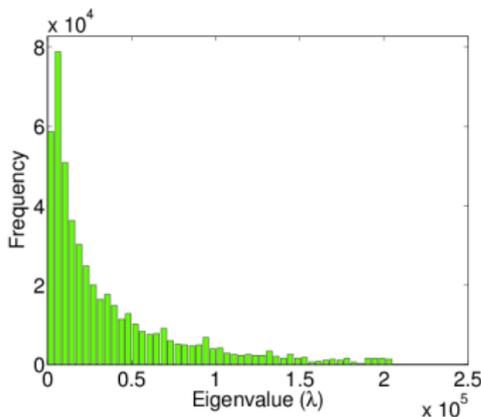
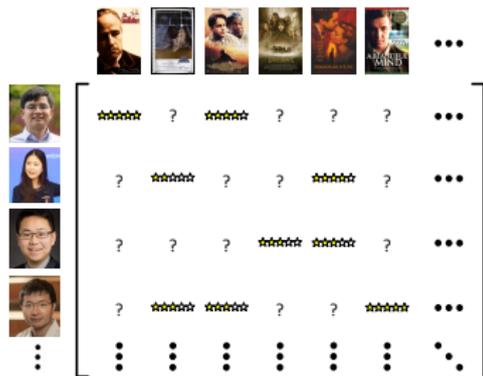
In general, we cannot infer missing ratings

$$\begin{bmatrix} \checkmark & ? & ? & ? & \checkmark & ? \\ ? & ? & \checkmark & \checkmark & ? & ? \\ \checkmark & ? & ? & \checkmark & ? & ? \\ ? & ? & \checkmark & ? & ? & \checkmark \\ \checkmark & ? & ? & ? & ? & ? \\ ? & \checkmark & ? & ? & \checkmark & ? \\ ? & ? & \checkmark & \checkmark & ? & ? \end{bmatrix}$$

— an underdetermined system (more unknowns than observations)

Matrix completion

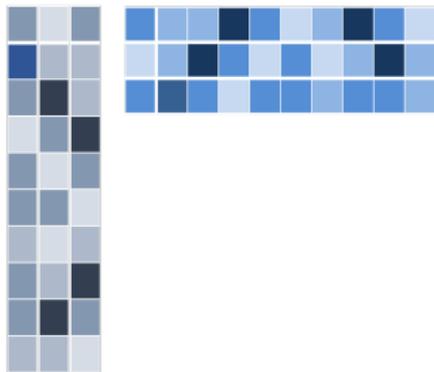
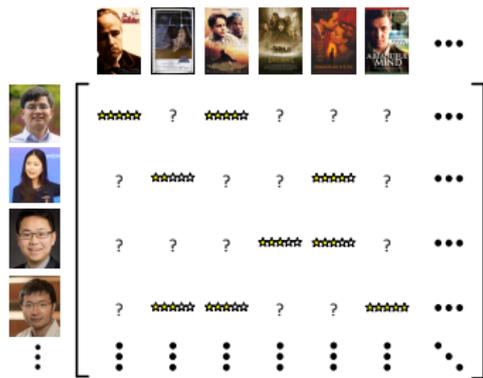
... unless rating matrix has other structure



A few factors explain most of the data

Matrix completion

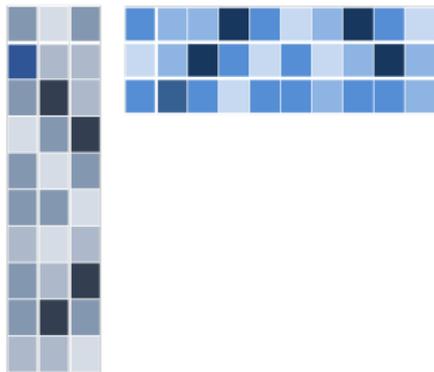
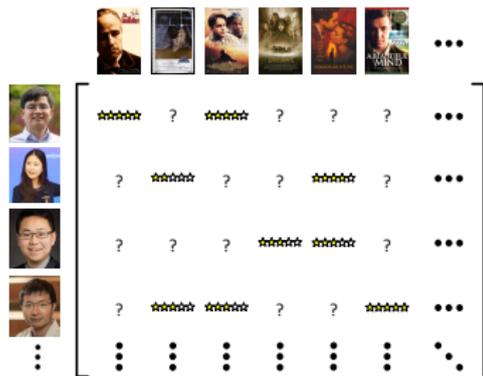
... unless rating matrix has other structure



A few factors explain most of the data \longrightarrow **low-rank** approximation

Matrix completion

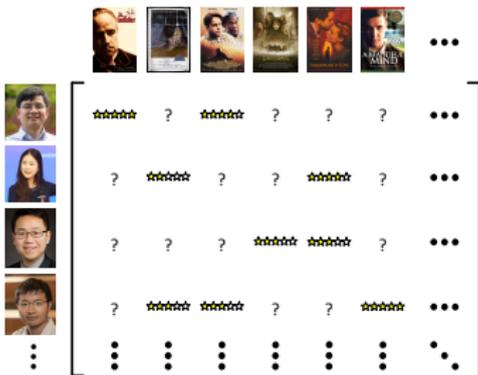
... unless rating matrix has other structure



A few factors explain most of the data \rightarrow **low-rank** approximation

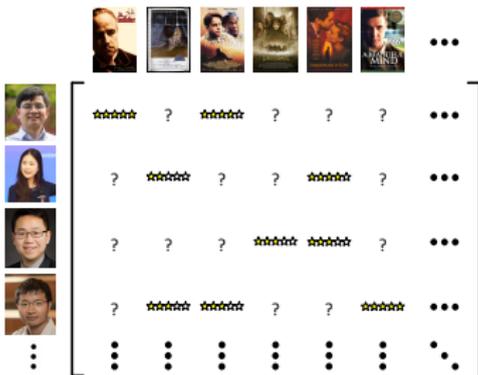
How to exploit (approx.) low-rank structure in prediction?

Model: low-rank matrix completion



- ground truth: rank- r matrix $M^* = \underbrace{U^* \Sigma^* V^{*\top}}_{\text{rank-}r \text{ SVD}} \in \mathbb{R}^{n_1 \times n_2}$
- each entry $M_{i,j}^*$ is observed independently with prob. p
- **goal:** fill in unseen entries of M^*

Model: low-rank matrix completion



- ground truth: rank- r matrix $M^* = \underbrace{U^* \Sigma^* V^{*\top}}_{\text{rank-}r \text{ SVD}} \in \mathbb{R}^{n_1 \times n_2}$
- each entry $M_{i,j}^*$ is observed independently with prob. p
- **goal:** fill in unseen entries of M^*
- **intermediate step:** estimate U^*, V^*, Σ^*

Spectral method for matrix completion

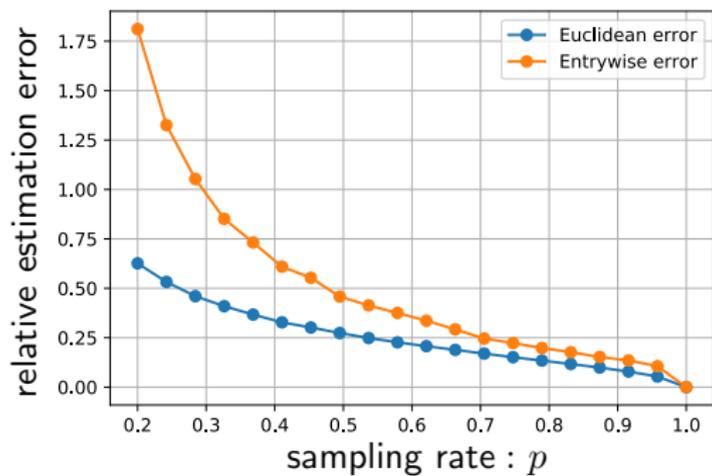
1. construct a rescaled zero-filled matrix $\mathbf{M} = [M_{i,j}] \in \mathbb{R}^{n_1 \times n_2}$ as

$$\forall(i, j) : M_{i,j} = \begin{cases} \frac{1}{p} M_{i,j}^*, & \text{if } M_{i,j}^* \text{ is observed} \\ 0, & \text{else} \end{cases}$$

- o **rationale:** ensures $\mathbb{E}[\mathbf{M}] = \mathbf{M}^*$

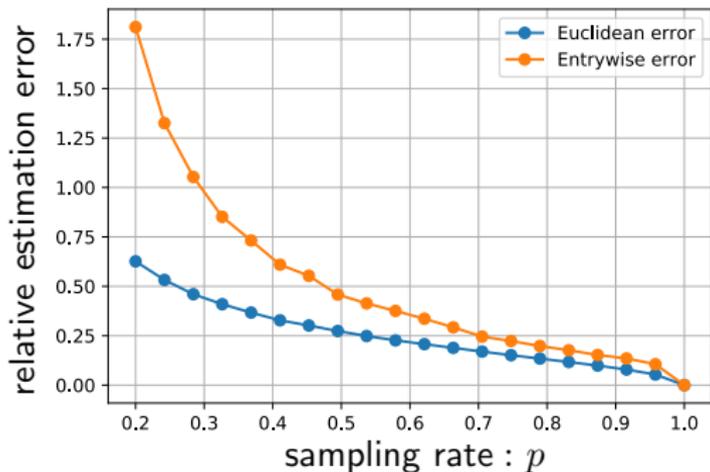
2. compute rank- r SVD $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ of \mathbf{M} , and return $\widehat{\mathbf{M}} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$

Empirical matrix estimation accuracy



$$n = 200, r = 5$$

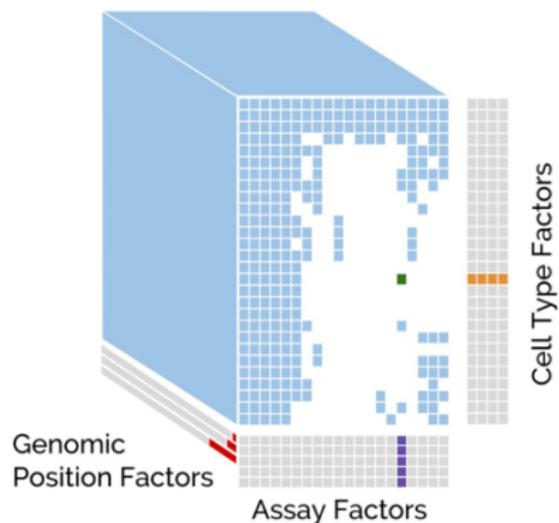
Empirical matrix estimation accuracy



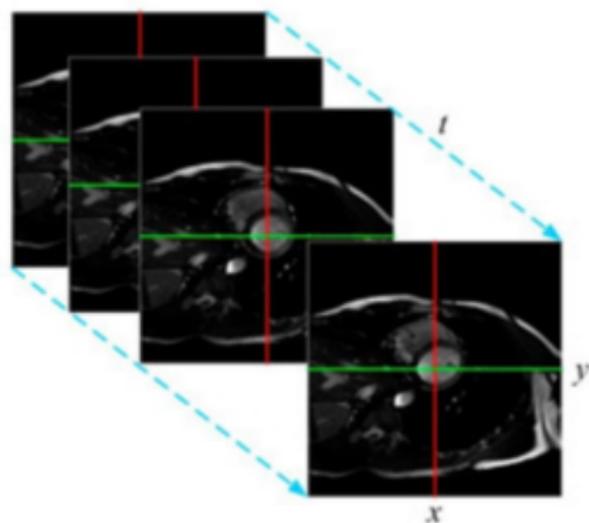
$$n = 200, r = 5$$

What we'll see: effect of sampling rate p upon estimation accuracy

Extension: tensor data

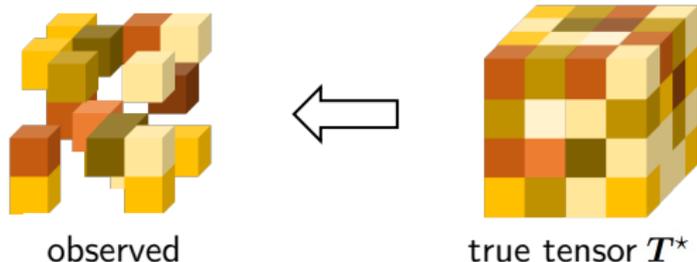


computational genomics
— *fig. credit: Schreiber et al. 19*



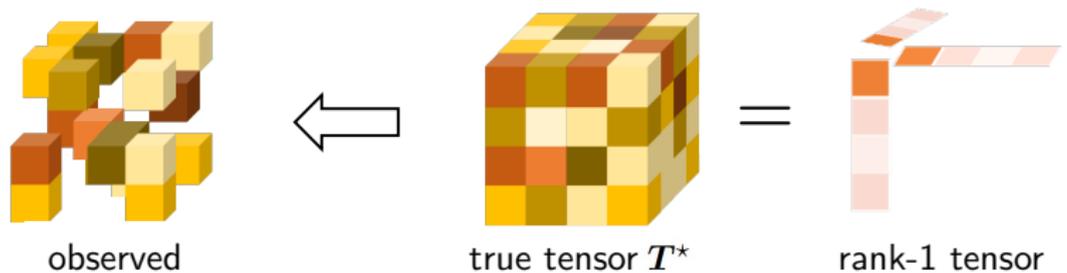
dynamic MRI
— *fig. credit: Liu et al. 17*

Extension: low-rank tensor completion



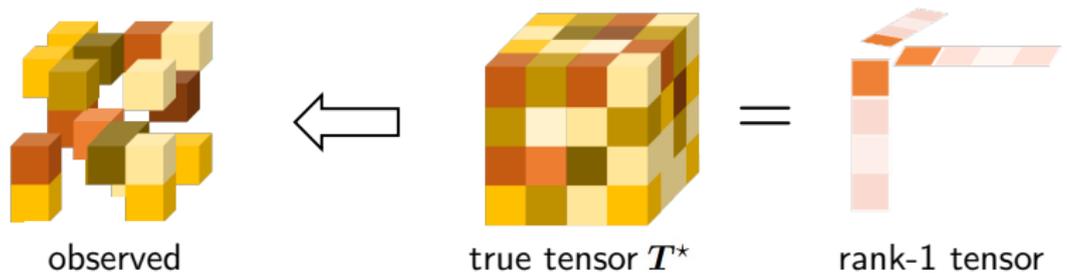
- ground truth: rank-1 tensor $T^* \in \mathbb{R}^{n \times n \times n}$
- each entry $T_{i,j,k}^*$ is observed independently with prob. p
- **goal:** fill in unseen entries of T^*

Extension: low-rank tensor completion



- ground truth: rank-1 tensor $T^* = \mathbf{u}^* \otimes \mathbf{u}^* \otimes \mathbf{u}^* \in \mathbb{R}^{n \times n \times n}$
- each entry $T_{i,j,k}^*$ is observed independently with prob. p
- **goal**: fill in unseen entries of T^*

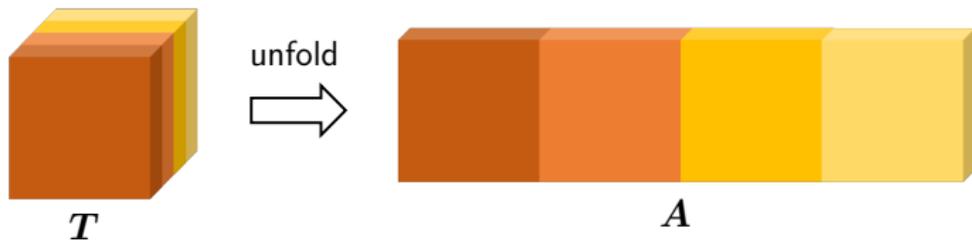
Extension: low-rank tensor completion



- ground truth: rank-1 tensor $T^* = \mathbf{u}^* \otimes \mathbf{u}^* \otimes \mathbf{u}^* \in \mathbb{R}^{n \times n \times n}$
- each entry $T_{i,j,k}^*$ is observed independently with prob. p
- **goal:** fill in unseen entries of T^*

Can we exploit low-rank tensor structure in prediction?

Spectral method for low-rank tensor completion



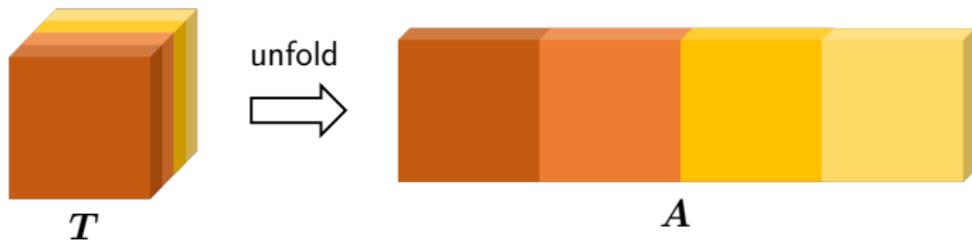
1. construct a rescaled zero-filled tensor $\mathbf{T} = [T_{i,j,k}] \in \mathbb{R}^{n \times n \times n}$ as

$$T_{i,j,k} = \begin{cases} \frac{1}{p} T_{i,j,k}^* & \text{if } T_{i,j,k}^* \text{ is observed} \\ 0 & \text{else} \end{cases}$$

- rescaling ensures $\mathbb{E}[\mathbf{T}] = \mathbf{T}^*$

2. matricization: $\mathbf{A} = \text{unfold}(\mathbf{T})$

Spectral method for low-rank tensor completion



3. compute spectral estimates (after diagonal deletion):

\mathbf{u} \leftarrow leading eigenvector
 λ \leftarrow leading eigenvalue of $\mathcal{P}_{\text{off-diag}}(\mathbf{A}\mathbf{A}^\top)$ (remove diagonal)

4. return $\hat{\mathbf{T}} = \lambda \mathbf{u} \otimes \mathbf{u} \otimes \mathbf{u}$

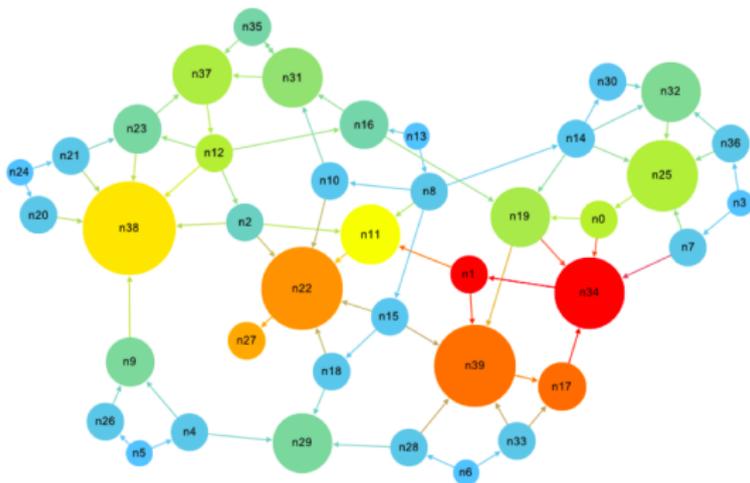
We will explain special treatments for diagonals

Motivating application: ranking

Ranking

A fundamental problem in a wide range of contexts

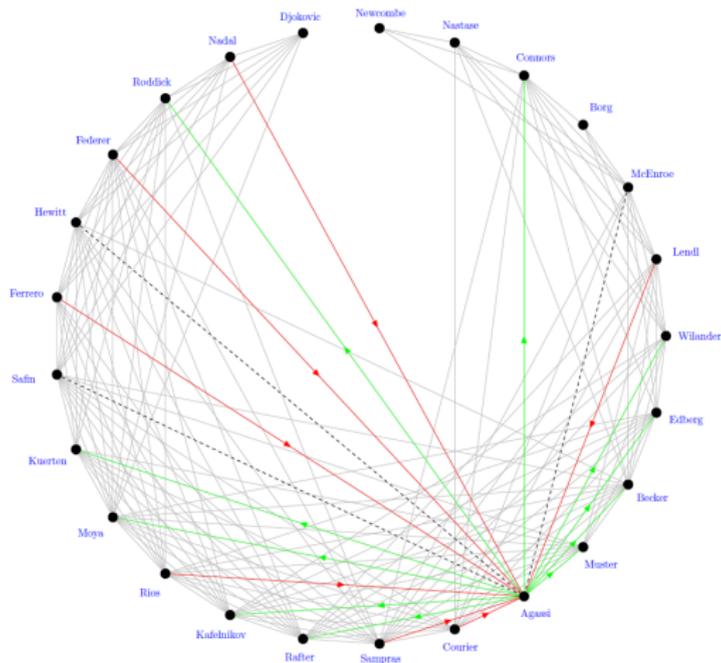
- web search, recommendation systems, admissions, sports competitions, voting, ...



PageRank

figure credit: Dzenan Hamzic

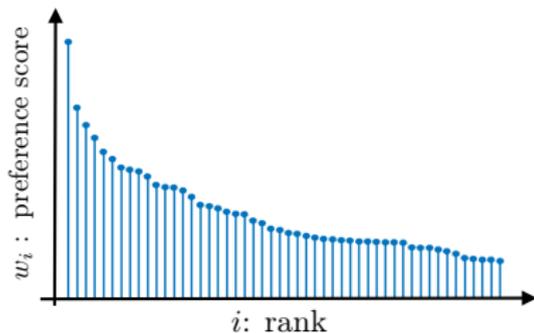
Ranking from pairwise comparisons



pairwise comparisons for ranking tennis players

figure credit: Bozóki, Csató, Temesi

Parametric models

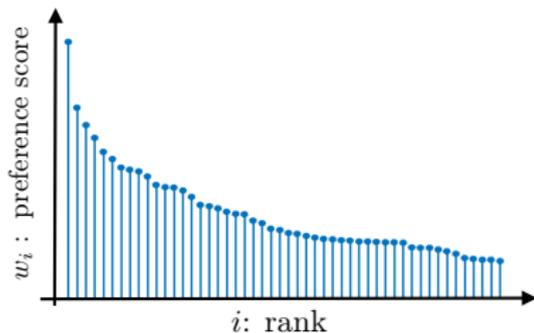


- n items to be ranked
- assign a latent score $\{w_i^*\}_{1 \leq i \leq n}$ to each item, so that

$$\text{item } i \succ \text{item } j \quad \text{if} \quad w_i^* > w_j^*$$

- rank items in accordance with score estimates

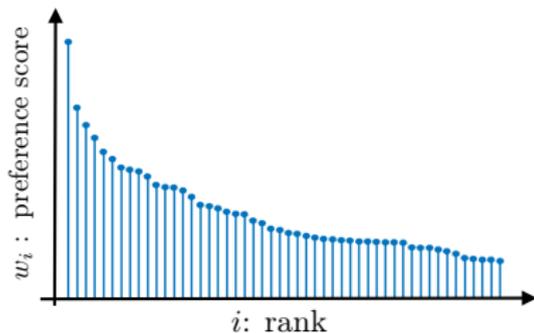
Bradley-Terry-Luce (logistic) model



- each pair of items (i, j) is compared independently

$$\mathbb{P} \{ \text{item } j \text{ beats item } i \} = \frac{w_j^*}{w_i^* + w_j^*}$$

Bradley-Terry-Luce (logistic) model



- each pair of items (i, j) is compared independently

$$\mathbb{P}\{\text{item } j \text{ beats item } i\} = \frac{w_j^*}{w_i^* + w_j^*}$$

$$\iff y_{i,j} \stackrel{\text{ind.}}{=} \begin{cases} 1, & \text{with prob. } \frac{w_j^*}{w_i^* + w_j^*} \\ 0, & \text{else} \end{cases}$$

Spectral ranking

Key idea: consider a probability transition matrix $P^* \in \mathbb{R}^{n \times n}$:

$$P_{i,j}^* = \begin{cases} \frac{1}{n} \cdot \frac{w_j^*}{w_i^* + w_j^*}, & \text{if } i \neq j \\ 1 - \sum_{l:l \neq i} P_{i,l}^*, & \text{if } i = j \end{cases}$$

Spectral ranking

Key idea: consider a probability transition matrix $P^* \in \mathbb{R}^{n \times n}$:

$$P_{i,j}^* = \begin{cases} \frac{1}{n} \cdot \frac{w_j^*}{w_i^* + w_j^*}, & \text{if } i \neq j \\ 1 - \sum_{l:l \neq i} P_{i,l}^*, & \text{if } i = j \end{cases}$$

- stationary distribution π^* of P^* : $\pi^* = \frac{1}{\sum_l w_l^*} \mathbf{w}^*$
leading left eigenvector of P^*
 - can be seen from **detailed balance** property: $w_i^* P_{i,j}^* = w_j^* P_{j,i}^*$

Spectral ranking

Key idea: consider a probability transition matrix $P^* \in \mathbb{R}^{n \times n}$:

$$P_{i,j}^* = \begin{cases} \frac{1}{n} \cdot \frac{w_j^*}{w_i^* + w_j^*}, & \text{if } i \neq j \\ 1 - \sum_{l:l \neq i} P_{i,l}^*, & \text{if } i = j \end{cases}$$

- stationary distribution π^* of P^* : $\pi^* = \frac{1}{\sum_l w_l^*} \mathbf{w}^*$
leading left eigenvector of P^*
 - can be seen from **detailed balance** property: $w_i^* P_{i,j}^* = w_j^* P_{j,i}^*$

True ranks are revealed by leading left eigenvector of P^*

Spectral ranking

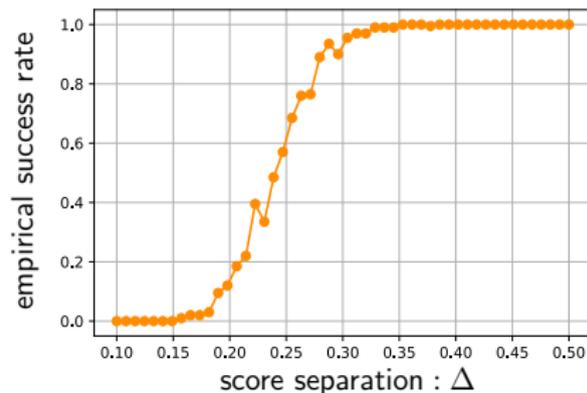
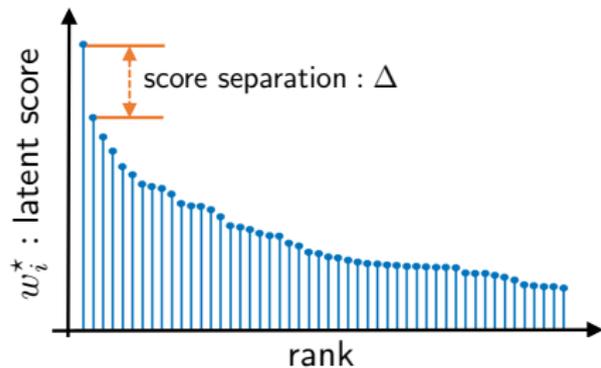
1. construct a surrogate matrix P obeying

$$P_{i,j} = \begin{cases} \frac{1}{n}y_{i,j}, & \text{if } i \neq j \\ 1 - \sum_{l:l \neq i} P_{i,l}, & \text{if } i = j \end{cases}$$

2. compute leading left eigenvector π of P as score estimate
3. rank in accordance with π

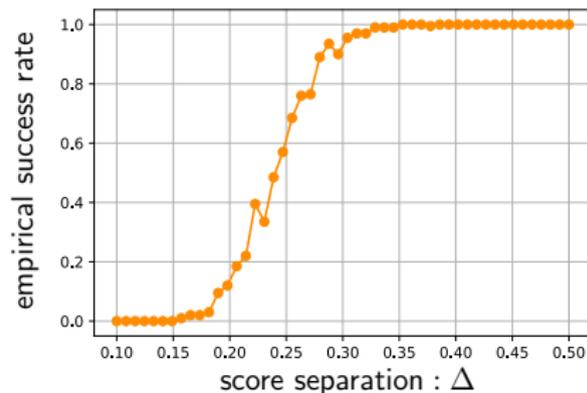
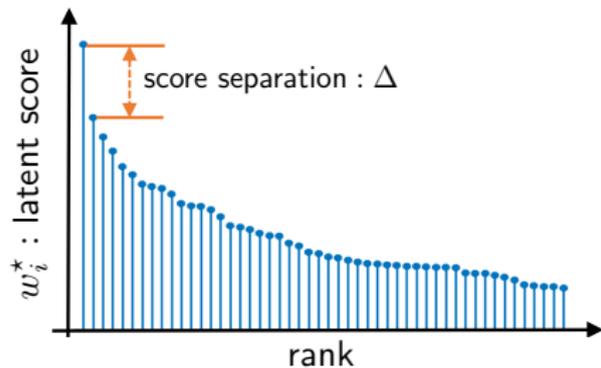
— closely related to PageRank

Empirical accuracy in finding top-ranked item



$n = 200$

Empirical accuracy in finding top-ranked item



$$n = 200$$

What we will demonstrate: efficacy of spectral ranking

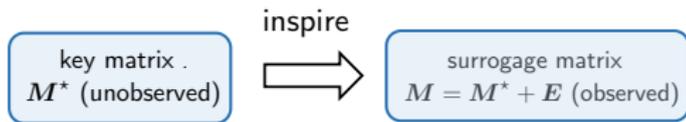
A unified recipe for spectral methods

A unified recipe

key matrix .
 M^* (unobserved)

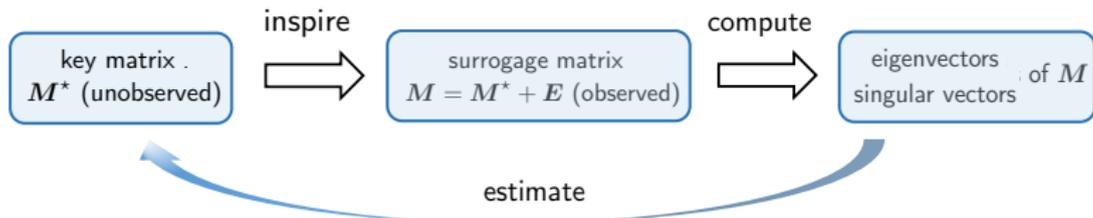
1. identify a key matrix M^* — typically unobserved — whose eigenvectors or singular vectors encode crucial information

A unified recipe



1. identify a key matrix M^* — typically unobserved — whose eigenvectors or singular vectors encode crucial information
2. construct a surrogate matrix M of M^* using data samples

A unified recipe



1. identify a key matrix M^* — typically unobserved — whose eigenvectors or singular vectors encode crucial information
2. construct a surrogate matrix M of M^* using data samples
3. compute corresponding eigenvectors or singular vectors of M

Key factors

A few factors that dictate the performance of spectral methods:

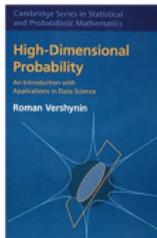
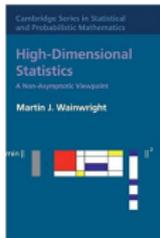
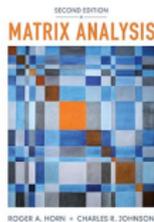
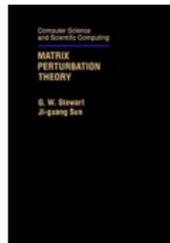
- proximity of M and M^* (e.g. $\|M - M^*\|$)
- spectrum (e.g. eigenvalues, singular values) of M^*
- ...

Key factors

A few factors that dictate the performance of spectral methods:

- proximity of M and M^* (e.g. $\|M - M^*\|$)
- spectrum (e.g. eigenvalues, singular values) of M^*
- ...

Aim of this tutorial: quantify influences of these factors



matrix perturbation theory

high-dimensional statistics
probability

- **algebraic tools:** matrix perturbation theory (Part 2)
- **statistical & probabilistic tools:**
 - matrix concentration bounds (Part 3: ℓ_2 analysis)
 - leave-one-out analysis (Part 4: fine-grained analysis)

Asymptotic notation used in this tutorial

- $f(n) \lesssim g(n)$ or $f(n) = O(g(n))$ means

$$\lim_{n \rightarrow \infty} \frac{|f(n)|}{|g(n)|} \leq \text{const}$$

- $f(n) \gtrsim g(n)$ means

$$\lim_{n \rightarrow \infty} \frac{|f(n)|}{|g(n)|} \geq \text{const}$$

- $f(n) \asymp g(n)$ means

$$\text{const}_1 \leq \lim_{n \rightarrow \infty} \frac{|f(n)|}{|g(n)|} \leq \text{const}_2$$

- $f(n) = o(g(n))$ means

$$\lim_{n \rightarrow \infty} \frac{|f(n)|}{|g(n)|} = 0$$

Part 2: Matrix perturbation theory

- Eigen-space perturbation theory
 - Distances and angles between two subspaces
 - The Davis-Kahan $\sin \Theta$ theorem
- Singular subspace perturbation theory (Wedin's theorem)
- Eigenvector perturbation for probability transition matrices

Eigen-space perturbation theory

Setup and notation

Consider 2 symmetric matrices M^* , $M = M^* + E \in \mathbb{R}^{n \times n}$ with eigen-decompositions

$$M^* = \sum_{i=1}^n \lambda_i^* \mathbf{u}_i^* \mathbf{u}_i^{*\top} = [\mathbf{U}^*, \mathbf{U}_\perp^*] \begin{bmatrix} \mathbf{\Lambda}^* & \\ & \mathbf{\Lambda}_\perp^* \end{bmatrix} \begin{bmatrix} \mathbf{U}^{*\top} \\ \mathbf{U}_\perp^{*\top} \end{bmatrix}$$

$$M = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^\top = [\mathbf{U}, \mathbf{U}_\perp] \begin{bmatrix} \mathbf{\Lambda} & \\ & \mathbf{\Lambda}_\perp \end{bmatrix} \begin{bmatrix} \mathbf{U}^\top \\ \mathbf{U}_\perp^\top \end{bmatrix}$$

- eigenvalues: $\lambda_1^* \geq \dots \geq \lambda_n^*$, $\lambda_1 \geq \dots \geq \lambda_n$
- $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_r] \in \mathbb{R}^{n \times r}$, $\mathbf{\Lambda} = \text{diag}([\lambda_1, \dots, \lambda_r]) \in \mathbb{R}^{r \times r}$, \dots

Setup and notation

$$M = \left[\underbrace{\mathbf{u}_1 \ \cdots \ \mathbf{u}_r}_{\mathbf{U}} \ \underbrace{\mathbf{u}_{r+1} \ \cdots \ \mathbf{u}_n}_{=:\mathbf{U}_\perp} \right]$$
$$\left[\begin{array}{ccc} \lambda_1 & & \\ & \ddots & \\ & & \lambda_r \\ \underbrace{\hspace{10em}}_{=:\mathbf{\Lambda}} & & \\ & & \lambda_{r+1} \\ & & & \ddots \\ & & & & \lambda_n \\ & & \underbrace{\hspace{10em}}_{=:\mathbf{\Lambda}_\perp} & & \end{array} \right] \left[\begin{array}{c} \mathbf{u}_1^\top \\ \vdots \\ \mathbf{u}_r^\top \\ \mathbf{u}_{r+1}^\top \\ \vdots \\ \mathbf{u}_n^\top \end{array} \right] \left. \vphantom{\begin{array}{c} \mathbf{u}_1^\top \\ \vdots \\ \mathbf{u}_r^\top \\ \mathbf{u}_{r+1}^\top \\ \vdots \\ \mathbf{u}_n^\top \end{array}} \right\} =:\mathbf{U}^\top$$
$$\left. \vphantom{\begin{array}{c} \mathbf{u}_1^\top \\ \vdots \\ \mathbf{u}_r^\top \\ \mathbf{u}_{r+1}^\top \\ \vdots \\ \mathbf{u}_n^\top \end{array}} \right\} =:\mathbf{U}_\perp^\top$$

Setup and notation

- $\mathcal{O}^{r \times r}$: set of all $r \times r$ orthonormal matrices
- $\|M\|$: spectral norm (largest singular value of M)
- $\|M\|_F$: Frobenius norm ($\|M\|_F = \sqrt{\text{tr}(M^\top M)} = \sqrt{\sum_{i,j} M_{i,j}^2}$)

Eigen-space perturbation theory

Main focus: how does perturbation matrix E affect “distance” between U^* and U ?

Question #0: how to define distance between two subspaces?

Eigen-space perturbation theory

Main focus: how does perturbation matrix E affect “distance” between U^* and U ?

Question #0: how to define distance between two subspaces?

- $\|U - U^*\|_F$ and $\|U - U^*\|$ are not appropriate, since they fall short of accounting for global orthonormal transformation

\forall orthonormal $R \in \mathbb{R}^{r \times r}$, U and UR represent same subspace

Distances and angles between two subspaces

Two valid choices of distance metrics

Key: taking care of global orthonormal transformation

- **Distance modulo optimal rotation:** adjust for rotation before computing distance:

$$\text{dist}(U, U^*) := \min_{R \in \mathcal{O}^{r \times r}} \|UR - U^*\| \quad (2.1)$$

Two valid choices of distance metrics

Key: taking care of global orthonormal transformation

- **Distance modulo optimal rotation:** adjust for rotation before computing distance:

$$\text{dist}(U, U^*) := \min_{R \in O^{r \times r}} \|UR - U^*\| \quad (2.1)$$

- **Distance using projection matrices:** replace U (resp. U^*) with its associated projection matrix before computing distance:

$$\text{dist}_p(U, U^*) := \left\| \underbrace{UU^T}_{\text{projection onto subspace } U} - U^*U^{*\top} \right\| \quad (2.2)$$

(Near)-equivalence of two distance metrics

$$\text{dist}(U, U^*) := \min_{R \in \mathcal{O}^{r \times r}} \|UR - U^*\|$$

$$\text{dist}_p(U, U^*) := \|UU^\top - U^*U^{*\top}\|$$

Lemma 2.1

Suppose $[U, U_\perp]$, $[U^*, U_\perp^*]$ are square orthonormal matrices. Then

$$\text{dist}_p(U, U^*) \leq \text{dist}(U, U^*) \leq \sqrt{2} \text{dist}_p(U, U^*)$$

- $\text{dist}(\cdot, \cdot)$ and $\text{dist}_p(\cdot, \cdot)$ are orderwise equivalent
- proof: see Section 2.6.3 of Chen et al. 21

An alternative expression for $\text{dist}_p(\cdot, \cdot)$

As it turns out, $\text{dist}_p(\cdot, \cdot)$ has several equivalent expressions:

Lemma 2.2

Recall that $[U, U_\perp]$, $[U^, U_\perp^*]$ are square orthonormal matrices. Then*

$$\text{dist}_p(U, U^*) = \|U^\top U_\perp^*\| = \|U^{*\top} U_\perp\|$$

An alternative expression for $\text{dist}_p(\cdot, \cdot)$

As it turns out, $\text{dist}_p(\cdot, \cdot)$ has several equivalent expressions:

Lemma 2.2

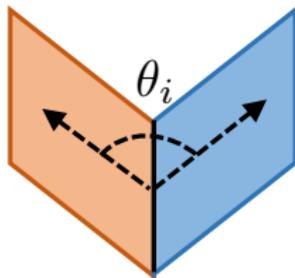
Recall that $[U, U_\perp], [U^, U_\perp^*]$ are square orthonormal matrices. Then*

$$\text{dist}_p(U, U^*) = \|U^\top U_\perp^*\| = \|U^{*\top} U_\perp\|$$

- sanity check: if $U = U^*$, then $\text{dist}(U, U^*) = \|U^\top U^*\| = 0$
- proof: see Slide 2-15

Principal angles between two subspaces

In addition to “distance”, one might also be interested in “angles”



We can quantify the similarity between two lines (represented resp. by unit vectors \mathbf{u} and \mathbf{u}^*) by an angle between them

$$\theta = \arccos\langle \mathbf{u}, \mathbf{u}^* \rangle$$

Principal angles between two eigen-spaces

For r -dimensional subspaces, one needs r angles

Specifically, given $\|U^\top U^*\| \leq 1$, write the SVD of $U^\top U^* \in \mathbb{R}^{r \times r}$ as

$$U^\top U^* = X \underbrace{\begin{bmatrix} \cos \theta_1 & & \\ & \ddots & \\ & & \cos \theta_r \end{bmatrix}}_{=: \cos \Theta} Y^\top =: X \cos \Theta Y^\top$$

- $X, Y \in \mathbb{R}^{r \times r}$: square orthonormal matrices
- $\{\theta_1, \dots, \theta_r\}$ are called the **principal angles** between U and U^*

Relations between principal angles and distance

As expected, principal angles and distances are closely related

Lemma 2.3

Suppose $[U, U_{\perp}]$, $[U^*, U_{\perp}^*]$ are square orthonormal matrices. Then

$$\|U^{\top} U_{\perp}^*\| = \|\sin \Theta\| = \max\{|\sin \theta_1|, \dots, |\sin \theta_r|\}$$

Lemmas 2.2 and 2.3 taken collectively give

$$\text{dist}_p(U, U^*) = \max\{|\sin \theta_1|, \dots, |\sin \theta_r|\} \quad (2.3)$$

Proof of Lemma 2.3

$$\begin{aligned}\|U^\top U_\perp^\star\| &= \|U^\top \underbrace{U_\perp^\star U_\perp^{\star\top}}_{=I-U^\star U^{\star\top}} U\|^{\frac{1}{2}} \\ &= \|U^\top U - U^\top U^\star U^{\star\top} U\|^{\frac{1}{2}} \\ &= \|I - X \cos^2 \Theta X^\top\|^{\frac{1}{2}} \quad (\text{since } U^\top U^\star = X \cos \Theta Y^\top) \\ &= \|I - \cos^2 \Theta\|^{\frac{1}{2}} \\ &= \|\sin \Theta^2\|^{\frac{1}{2}} \\ &= \|\sin \Theta\|\end{aligned}$$

Proof of Lemma 2.2

We first claim that the SVD of $U_{\perp}^{\top} U^{\star}$ can be written as

$$U_{\perp}^{\top} U^{\star} = \widetilde{X} \sin \Theta Y^{\top} \quad (2.4)$$

for some orthonormal \widetilde{X} (to be proved later). Armed w/ this claim, one has

$$U^{\star} = [U, U_{\perp}] \begin{bmatrix} U^{\top} \\ U_{\perp}^{\top} \end{bmatrix} U^{\star} = [U, U_{\perp}] \begin{bmatrix} X \cos \Theta Y^{\top} \\ \widetilde{X} \sin \Theta Y^{\top} \end{bmatrix}$$

$$\implies U^{\star} U^{\star \top} = [U, U_{\perp}] \begin{bmatrix} X \cos^2 \Theta X^{\top} & X \cos \Theta \sin \Theta \widetilde{X}^{\top} \\ \widetilde{X} \cos \Theta \sin \Theta X^{\top} & \widetilde{X} \sin^2 \Theta \widetilde{X}^{\top} \end{bmatrix} \begin{bmatrix} U^{\top} \\ U_{\perp}^{\top} \end{bmatrix}$$

As a consequence,

$$\begin{aligned} & UU^{\top} - U^{\star} U^{\star \top} \\ &= [U, U_{\perp}] \begin{bmatrix} I - X \cos^2 \Theta X^{\top} & -X \cos \Theta \sin \Theta \widetilde{X}^{\top} \\ -\widetilde{X} \cos \Theta \sin \Theta X^{\top} & -\widetilde{X} \sin^2 \Theta \widetilde{X}^{\top} \end{bmatrix} \begin{bmatrix} U^{\top} \\ U_{\perp}^{\top} \end{bmatrix} \end{aligned}$$

Proof of Lemma 2.2 (cont.)

This further gives

$$\begin{aligned} & \|UU^\top - U^*U^{*\top}\| \\ &= \left\| \begin{bmatrix} \mathbf{X} & \\ & \widetilde{\mathbf{X}} \end{bmatrix} \begin{bmatrix} \sin^2 \Theta & -\cos \Theta \sin \Theta \\ -\cos \Theta \sin \Theta & -\sin^2 \Theta \end{bmatrix} \begin{bmatrix} \mathbf{X}^\top & \\ & \widetilde{\mathbf{X}}^\top \end{bmatrix} \right\| \\ &= \left\| \underbrace{\begin{bmatrix} \sin^2 \Theta & -\cos \Theta \sin \Theta \\ -\cos \Theta \sin \Theta & -\sin^2 \Theta \end{bmatrix}}_{\text{each block is a diagonal matrix}} \right\| \quad (\|\cdot\| \text{ is rotationally invariant}) \\ &= \max_{1 \leq i \leq r} \left\| \begin{bmatrix} \sin^2 \theta_i & -\cos \theta_i \sin \theta_i \\ -\cos \theta_i \sin \theta_i & -\sin^2 \theta_i \end{bmatrix} \right\| \\ &= \max_{1 \leq i \leq r} \left\| \sin \theta_i \begin{bmatrix} \sin \theta_i & -\cos \theta_i \\ -\cos \theta_i & -\sin \theta_i \end{bmatrix} \right\| \\ &= \max_{1 \leq i \leq r} |\sin \theta_i| = \|\sin \Theta\| \end{aligned}$$

Proof of Lemma 2.2 (cont.)

It remains to justify (2.4). To this end, observe that

$$\begin{aligned}U^{\star\top}U_{\perp}U_{\perp}^{\top}U^{\star} &= U^{\star\top}U^{\star} - U^{\star\top}UUU^{\top}U^{\star} \\ &= I - Y \cos^2 \Theta Y^{\top} \\ &= Y \sin^2 \Theta Y^{\top}\end{aligned}$$

and hence the right singular space (resp. singular values) of $U^{\star\top}U_{\perp}$ is given by Y (resp. $\sin \Theta$). This immediately implies (2.4).

Summary: four (almost) equivalent distance metrics

- 1) $\|UU^\top - U^*U^{*\top}\|$
- 2) $\|\sin \Theta\|$
- 3) $\|U_\perp^\top U^*\| = \|U^\top U_\perp^*\|$
- 4) $\min_{R \in \mathcal{O}^{r \times r}} \|UR - U^*\|$

Summary: four (almost) equivalent distance metrics

- 1) $\|UU^T - U^*U^{*\top}\|$
- 2) $\|\sin \Theta\|$
- 3) $\|U_{\perp}^T U^*\| = \|U^T U_{\perp}^*\|$
- 4) $\min_{R \in \mathcal{O}^{r \times r}} \|UR - U^*\|$

Near-equivalence of these metrics continue to hold if $\|\cdot\|$ is replaced by $\|\cdot\|_F$

The Davis-Kahan $\sin \Theta$ theorem

Warm-up example ($0 < \epsilon < 1$):

$$\mathbf{M}^* = \begin{bmatrix} 1 + \epsilon & 0 \\ 0 & 1 - \epsilon \end{bmatrix}, \quad \mathbf{E} = \begin{bmatrix} -\epsilon & \epsilon \\ \epsilon & \epsilon \end{bmatrix}, \quad \mathbf{M} = \begin{bmatrix} 1 & \epsilon \\ \epsilon & 1 \end{bmatrix}$$

Warm-up example ($0 < \epsilon < 1$):

$$\mathbf{M}^* = \begin{bmatrix} 1 + \epsilon & 0 \\ 0 & 1 - \epsilon \end{bmatrix}, \quad \mathbf{E} = \begin{bmatrix} -\epsilon & \epsilon \\ \epsilon & \epsilon \end{bmatrix}, \quad \mathbf{M} = \begin{bmatrix} 1 & \epsilon \\ \epsilon & 1 \end{bmatrix}$$

- leading eigenvectors of \mathbf{M}^* and \mathbf{M} :

$$\mathbf{u}_1^* = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \mathbf{u}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \Longrightarrow \quad \|\mathbf{u}_1 \mathbf{u}_1^\top - \mathbf{u}_1^* \mathbf{u}_1^{*\top}\|_2 = \frac{1}{\sqrt{2}}$$

— *eigenvector distance is large regardless of size of ϵ (or size of $\|\mathbf{E}\|$)*

Warm-up example ($0 < \epsilon < 1$):

$$\mathbf{M}^* = \begin{bmatrix} 1 + \epsilon & 0 \\ 0 & 1 - \epsilon \end{bmatrix}, \quad \mathbf{E} = \begin{bmatrix} -\epsilon & \epsilon \\ \epsilon & \epsilon \end{bmatrix}, \quad \mathbf{M} = \begin{bmatrix} 1 & \epsilon \\ \epsilon & 1 \end{bmatrix}$$

- leading eigenvectors of \mathbf{M}^* and \mathbf{M} :

$$\mathbf{u}_1^* = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \mathbf{u}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \Longrightarrow \quad \|\mathbf{u}_1 \mathbf{u}_1^\top - \mathbf{u}_1^* \mathbf{u}_1^{*\top}\|_2 = \frac{1}{\sqrt{2}}$$

— *eigenvector distance is large regardless of size of ϵ (or size of $\|\mathbf{E}\|$)*

Diagnosis: eigen-gap $\lambda_1^* - \lambda_2^* = 2\epsilon$ also small (proportional to $\|\mathbf{E}\|$)

— *both perturbation size and eigen-gap might play important roles*

Davis-Kahan $\sin \Theta$ Theorem: a simple case

— recall the setup in Page 2-3



Chandler Davis



William Kahan

Theorem 2.4

Suppose $M^* \succeq \mathbf{0}$ and has rank r . If $\|E\| < (1 - 1/\sqrt{2})\lambda_r^*$, then

$$\text{dist}_p(U, U^*) = \|U_{\perp}^{\top} U^*\| = \|\sin \Theta\| \leq \frac{\sqrt{2}\|EU^*\|}{\lambda_r^*} \leq \frac{\sqrt{2}\|E\|}{\lambda_r^*}$$

Interpretations

Suppose $M^* \succeq \mathbf{0}$ and has rank r . If $\|E\| < (1 - 1/\sqrt{2})\lambda_r^*$, then

$$\text{dist}_p(U, U^*) \leq \frac{\sqrt{2}\|EU^*\|}{\lambda_r^*} \leq \frac{\sqrt{2}\|E\|}{\lambda_r^*}$$

Key factors:

1. eigen-gap: $\lambda_r^* = \lambda_r^* - \underbrace{\lambda_{r+1}^*}_{=0}$

2. perturbation size: $\|E\|$

3. signal-to-noise ratio (SNR): $\frac{\lambda_r^*}{\|E\|}$

- the bound w/ $\|EU^*\|$ is sometimes useful (e.g. for ℓ_∞ analysis)

Proof of Theorem 2.4

We intend to control $U_{\perp}^{\top} U^*$ by studying their interactions through E :

$$\begin{aligned}\|U_{\perp}^{\top} E U^*\| &= \left\| U_{\perp}^{\top} \left(\underbrace{U \Lambda U^{\top} + U_{\perp} \Lambda_{\perp} U_{\perp}^{\top}}_{M^* + E} - \underbrace{U^* \Lambda^* U^{*\top}}_{M^*} \right) U^* \right\| \\ &= \|\Lambda_{\perp} U_{\perp}^{\top} U^* - U_{\perp}^{\top} U^* \Lambda^*\| \quad (\text{since } U_{\perp}^*{}^{\top} U^* = U_{\perp}^{\top} U = \mathbf{0}) \\ &\geq \|U_{\perp}^{\top} U^* \Lambda^*\| - \|\Lambda_{\perp} U_{\perp}^{\top} U^*\| \quad (\text{triangle inequality}) \\ &\geq \|U_{\perp}^{\top} U^*\| \lambda_r^* - \|U_{\perp}^{\top} U^*\| \|\Lambda_{\perp}\| \quad (2.5)\end{aligned}$$

Weyl's Theorem gives $\|\Lambda_{\perp}\| \leq \|E\|$, which combined with (2.5) yields

$$\|U_{\perp}^{\top} U^*\| \leq \frac{\|U_{\perp}^{\top} E U^*\|}{\lambda_r^* - \|E\|} \leq \frac{\|U_{\perp}\| \cdot \|E U^*\|}{\lambda_r^* - \|E\|} = \frac{\|E U^*\|}{\lambda_r^* - \|E\|}$$

This together with assumption $\|E\| \leq (1 - \sqrt{2}/2)\lambda_r^*$ and Lemmas 2.2-2.3 completes the proof

Davis-Kahan's $\sin \Theta$ theorem: general case

— $\text{eigenvalues}(A)$: set of eigenvalues of A

Theorem 2.5 (Davis-Kahan's $\sin \Theta$ theorem: general version)

Assume that

$$\text{eigenvalues}(\Lambda^*) \subseteq (-\infty, \alpha - \Delta] \cup [\beta + \Delta, \infty); \quad (2.6a)$$

$$\text{eigenvalues}(\Lambda_{\perp}^*) \subseteq [\alpha, \beta]. \quad (2.6b)$$

for some eigengap $\Delta > 0$. Suppose $\|E\| \leq (1 - \sqrt{2}/2)\Delta$. Then

$$\text{dist}(U, U^*) \leq \sqrt{2} \text{dist}_p(U, U^*) = \sqrt{2} \|\sin \Theta\| \leq \frac{2\|EU^*\|}{\Delta} \leq \frac{2\|E\|}{\Delta}$$

Davis-Kahan's $\sin \Theta$ theorem: general case

— $\text{eigenvalues}(\mathbf{A})$: set of eigenvalues of \mathbf{A}

Theorem 2.5 (Davis-Kahan's $\sin \Theta$ theorem: general version)

Assume that

$$\text{eigenvalues}(\mathbf{\Lambda}^*) \subseteq (-\infty, \alpha - \Delta] \cup [\beta + \Delta, \infty); \quad (2.6a)$$

$$\text{eigenvalues}(\mathbf{\Lambda}_{\perp}^*) \subseteq [\alpha, \beta]. \quad (2.6b)$$

for some eigengap $\Delta > 0$. Suppose $\|\mathbf{E}\| \leq (1 - \sqrt{2}/2)\Delta$. Then

$$\text{dist}(\mathbf{U}, \mathbf{U}^*) \leq \sqrt{2} \text{dist}_p(\mathbf{U}, \mathbf{U}^*) = \sqrt{2} \|\sin \Theta\| \leq \frac{2\|\mathbf{E}\mathbf{U}^*\|}{\Delta} \leq \frac{2\|\mathbf{E}\|}{\Delta}$$

- conclusion remains valid if Assumption (2.6) is reversed
- proof: see Section 2.3.4 of Chen et al. '21

Singular subspace perturbation theory

Setup and notation

Consider 2 matrices M^* , $M = M^* + E \in \mathbb{R}^{n_1 \times n_2}$ ($n_1 \leq n_2$) w/ SVD

$$M^* = \sum_{i=1}^{n_1} \sigma_i^* \mathbf{u}_i^* \mathbf{v}_i^{*\top} = \begin{bmatrix} U^* & U_{\perp}^* \end{bmatrix} \begin{bmatrix} \Sigma^* & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \Sigma_{\perp}^* & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}^{*\top} \\ \mathbf{V}_{\perp}^{*\top} \end{bmatrix}$$
$$M = \sum_{i=1}^{n_1} \sigma_i \mathbf{u}_i \mathbf{v}_i^{\top} = \begin{bmatrix} U & U_{\perp} \end{bmatrix} \begin{bmatrix} \Sigma & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \Sigma_{\perp} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}^{\top} \\ \mathbf{V}_{\perp}^{\top} \end{bmatrix}$$

- $\sigma_1 \geq \dots \geq \sigma_{n_1}$: singular values of M
- $\sigma_1^* \geq \dots \geq \sigma_{n_1}^*$: singular values of M^*
- $U = [\mathbf{u}_1, \dots, \mathbf{u}_r] \in \mathbb{R}^{n_1 \times r}$, $\Sigma = \text{diag}(\lambda_1, \dots, \lambda_r) \in \mathbb{R}^{r \times r}$, ...

Wedin's $\sin \Theta$ theorem

Davis-Kahan's theorem generalizes to singular subspace perturbation:

Theorem 2.6 (Wedin's $\sin \Theta$ theorem)

If $\|E\| < (1 - 1/\sqrt{2}) \underbrace{(\sigma_r^* - \sigma_{r+1}^*)}_{\text{spectral gap}}$, then one has

$$\begin{aligned} \max \{ \text{dist}(U, U^*), \text{dist}(V, V^*) \} &\leq \sqrt{2} \max \{ \text{dist}_p(U, U^*), \text{dist}_p(V, V^*) \} \\ &\leq \frac{2 \max \{ \|E^\top U^*\|, \|EV^*\| \}}{\sigma_r^* - \sigma_{r+1}^*} \leq \frac{2\|E\|}{\sigma_r^* - \sigma_{r+1}^*} \end{aligned}$$

- both EV^* and $E^\top U^*$ matter

Proof of Theorem 2.6

Similar to proof of Davis-Kahan theorem, we concentrate on $U_{\perp}^{\top} U^*$:

$$\begin{aligned}
 U_{\perp}^{\top} U^* &= U_{\perp}^{\top} (U^* \Sigma^* V^{*\top}) V^* \Sigma^{*-1} \\
 &= U_{\perp}^{\top} \left(\underbrace{M - E}_{=M^*} - U_{\perp}^* \Sigma_{\perp}^* V_{\perp}^{*\top} \right) V^* \Sigma^{*-1} \\
 &= U_{\perp}^{\top} \left(\underbrace{U \Sigma V^{\top} + U_{\perp} \Sigma_{\perp} V_{\perp}^{\top}}_{=M} - E - \cancel{U_{\perp}^* \Sigma_{\perp}^* V_{\perp}^{*\top}} \right) V^* \Sigma^{*-1} \\
 &= \Sigma_{\perp} V_{\perp}^{\top} V^* \Sigma^{*-1} - U_{\perp}^{\top} E V^* \Sigma^{*-1}
 \end{aligned}$$

Applying triangle inequality and Weyls' inequality yields

$$\begin{aligned}
 \|U_{\perp}^{\top} U^*\| &\leq \|\Sigma_{\perp}\| \cdot \|V_{\perp}^{\top} V^*\| \cdot \|\Sigma^{*-1}\| + \|U_{\perp}^{\top}\| \cdot \|E V^*\| \cdot \|\Sigma^{*-1}\| \\
 &= \sigma_{r+1} \cdot \|V_{\perp}^{\top} V^*\| \cdot \frac{1}{\sigma_r^*} + \|E V^*\| \cdot \frac{1}{\sigma_r^*} \\
 &\leq \frac{\sigma_{r+1}^* + \|E\|}{\sigma_r^*} \|V_{\perp}^{\top} V^*\| + \frac{\|E V^*\|}{\sigma_r^*}
 \end{aligned} \tag{2.7}$$

Proof of Theorem 2.6 (cont.)

Repeating the same argument yields

$$\|\mathbf{V}_\perp^\top \mathbf{V}^*\| \leq \frac{\sigma_{r+1}^* + \|\mathbf{E}\|}{\sigma_r^*} \|\mathbf{U}_\perp^\top \mathbf{U}^*\| + \frac{\|\mathbf{E}^\top \mathbf{U}^*\|}{\sigma_r^*} \quad (2.8)$$

Combine inequalities (2.7) and (2.8) to obtain

$$\begin{aligned} \max \{ \|\mathbf{U}_\perp^\top \mathbf{U}^*\|, \|\mathbf{V}_\perp^\top \mathbf{V}^*\| \} &\leq \frac{\max \{ \|\mathbf{E}^\top \mathbf{U}^*\|, \|\mathbf{E} \mathbf{V}^*\| \}}{\sigma_r^*} \\ &+ \frac{\sigma_{r+1}^* + \|\mathbf{E}\|}{\sigma_r^*} \max \{ \|\mathbf{U}_\perp^\top \mathbf{U}^*\|, \|\mathbf{V}_\perp^\top \mathbf{V}^*\| \} \end{aligned}$$

Rearrange terms and utilize $\|\mathbf{E}\| < (1 - \sqrt{2}/2)(\sigma_r^* - \sigma_{r+1}^*)$ and Lemma 2.1 to arrive at desired result

*Eigenvector perturbation for
probability transition matrices*

Eigen-decomposition for asymmetric matrices

Eigen-decomposition for asymmetric matrices is more tricky:

1. both eigenvalues & eigenvectors might be complex-valued
2. eigenvectors might not be orthogonal to each other

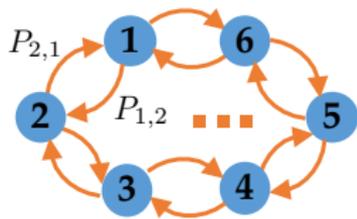
Eigen-decomposition for asymmetric matrices

Eigen-decomposition for asymmetric matrices is more tricky:

1. both eigenvalues & eigenvectors might be complex-valued
2. eigenvectors might not be orthogonal to each other

Let us look at a special case: **probability transition matrices**

Probability transition matrices



Consider a Markov chain $\{X_t\}_{t \geq 0}$

- n states
- transition probability $\mathbb{P}\{X_{t+1} = j \mid X_t = i\} = P_{i,j}$
- transition matrix $\mathbf{P} = [P_{i,j}]_{1 \leq i,j \leq n}$
- stationary distribution $\boldsymbol{\pi} = [\pi_i]_{1 \leq i \leq n}$ obeys

$$\boldsymbol{\pi} \geq \mathbf{0}, \quad \mathbf{1}^\top \boldsymbol{\pi} = 1, \quad \text{and} \quad \boldsymbol{\pi}^\top \mathbf{P} = \boldsymbol{\pi}^\top$$

- leading left eigenvector of \mathbf{P} with eigenvalue 1

Reversible Markov chains

Markov chain $\{X_t\}_{t \geq 0}$ with transition matrix P and stationary distribution π is said to be **reversible** if

$$\pi_i P_{i,j} = \pi_j P_{j,i} \quad \text{for all } i, j$$

— *detailed balance condition*

- If P represents reversible chain, then all eigenvalues of P are real

Setup and notation

- P^* : probability transition matrix of a **reversible** Markov chain
- $P = P^* + E$: (perturbed) probability transition matrix
- π^* (resp. π): leading left eigenvectors of P^* (resp. P)

Question: how does perturbation E affect leading left eigenvector?

Additional notation: for any probability vector $\pi = [\pi_i]_{1 \leq i \leq n} > \mathbf{0}$:

- vector norm: $\|\mathbf{x}\|_\pi := \sqrt{\sum_i \pi_i x_i^2}$ with $\mathbf{x} = [x_i]_{1 \leq i \leq n}$
- matrix norm: $\|\mathbf{A}\|_\pi := \sup_{\|\mathbf{x}\|_\pi=1} \|\mathbf{A}\mathbf{x}\|_\pi$ with $\mathbf{A} = [A_{i,j}]_{1 \leq i,j \leq n}$

Eigenvector perturbation for transition matrices

Theorem 2.7 (Chen, Fan, Ma, Wang '19)

Suppose P^* represents a *reversible* Markov chain, whose stationary distribution vector π^* is strictly positive. Assume

$$\|E\|_{\pi^*} < (1 - 1/\sqrt{2}) \left(1 - \max \{ \lambda_2(P^*), -\lambda_n(P^*) \} \right)$$

Then one has

$$\|\pi - \pi^*\|_{\pi^*} \leq \frac{\sqrt{2} \|\pi^{*\top} E\|_{\pi^*}}{1 - \max \{ \lambda_2(P^*), -\lambda_n(P^*) \}}$$

- similar to Davis-Kahan theorem
- eigengap: $1 - \max \{ \lambda_2(P^*), -\lambda_n(P^*) \}$ since $1 = \lambda_1(P^*)$
- perturbation size: $\|\pi^{*\top} E\|_{\pi^*}$

Part 3: Application of ℓ_2 perturbation theory

- Matrix tail bounds
- Community detection
- Matrix completion
- Ranking from pairwise comparisons

Matrix tail bounds

A hammer: matrix Bernstein inequality

Consider a sequence of independent random matrices $\{\mathbf{X}_l \in \mathbb{R}^{d_1 \times d_2}\}$

- $\mathbb{E}[\mathbf{X}_l] = \mathbf{0}$
- $\|\mathbf{X}_l\| \leq B$ for each l
- variance statistic:

$$v := \max \left\{ \left\| \mathbb{E} \left[\sum_l \mathbf{X}_l \mathbf{X}_l^\top \right] \right\|, \left\| \mathbb{E} \left[\sum_l \mathbf{X}_l^\top \mathbf{X}_l \right] \right\| \right\}$$

Theorem 3.8 (Matrix Bernstein inequality)

For all $\tau \geq 0$,

$$\mathbb{P} \left\{ \left\| \sum_l \mathbf{X}_l \right\| \geq \tau \right\} \leq (d_1 + d_2) \exp \left(\frac{-\tau^2/2}{v + B\tau/3} \right)$$

A hammer: matrix Bernstein inequality

$$\mathbb{P} \left\{ \left\| \sum_l \mathbf{X}_l \right\| \geq \tau \right\} \leq (d_1 + d_2) \exp \left(\frac{-\tau^2/2}{v + B\tau/3} \right)$$

- **moderate-deviation regime** (τ is small):
 - sub-Gaussian tail behavior $\exp(-\tau^2/2v)$
- **large-deviation regime** (τ is large):
 - sub-exponential tail behavior $\exp(-3\tau/2B)$ (slower decay)
- **user-friendly form** (exercise): with prob. $1 - O((d_1 + d_2)^{-10})$

$$\left\| \sum_l \mathbf{X}_l \right\| \lesssim \sqrt{v \log(d_1 + d_2)} + B \log(d_1 + d_2) \quad (3.9)$$

Another hammer: spectral norm of random matrices w/ independent entries

Consider a symmetric random matrix $\mathbf{X} = [X_{i,j}]_{1 \leq i,j \leq n}$ with independent entries s.t. $\forall (i, j)$:

- $\mathbb{E}[X_{i,j}] = 0$
- $\text{Var}(X_{i,j}) \leq \sigma^2$
- $|X_{i,j}| \leq B$

Theorem 3.9 (Bandeira, van Handel '16)

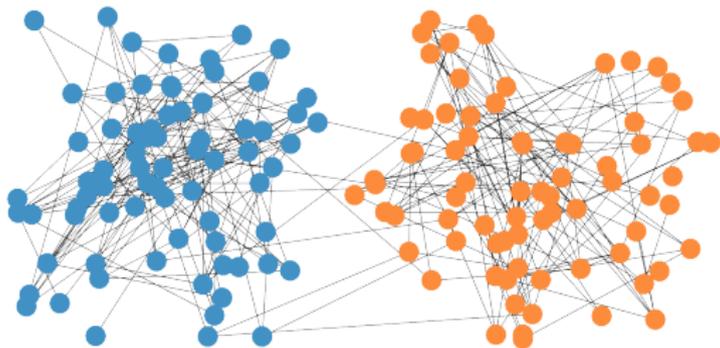
With probability exceeding $1 - O(n^{-10})$,

$$\|\mathbf{X}\| \leq 4\sigma\sqrt{n} + O(B\sqrt{\log n})$$

- often tighter than matrix Bernstein by some log factor

Community detection

Recap: spectral clustering for SBMs



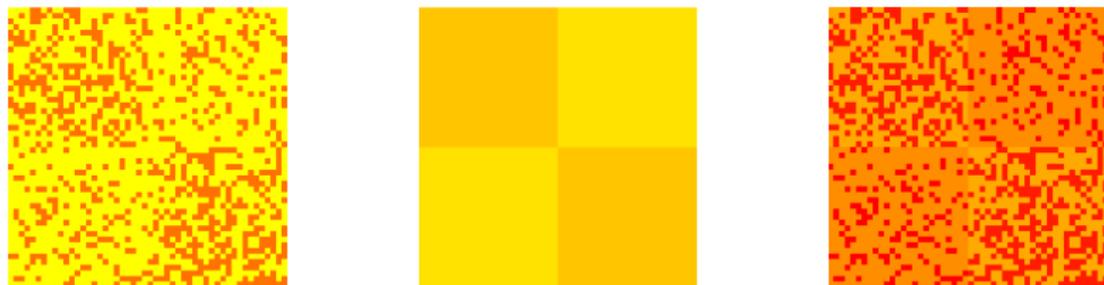
$x_i = 1$: 1st community

$x_i = -1$: 2nd community

- n nodes $\{1, \dots, n\}$
- 2 communities
- community memberships to recover: $\mathbf{x} = [x_i]_{1 \leq i \leq n} \in \{1, -1\}^n$
- observed: an adjacency matrix \mathbf{A} of a random graph s.t.

$$\mathbb{P}(A_{i,j} = 1) = \begin{cases} p, & \text{if } x_i = x_j \\ q, & \text{else} \end{cases}$$

Recap: spectral clustering for SBMs


$$\mathbf{A} = \underbrace{\mathbb{E}[\mathbf{A}]}_{= \frac{p+q}{2} \mathbf{1}\mathbf{1}^\top + \frac{p-q}{2} \mathbf{x}\mathbf{x}^\top} + \mathbf{A} - \mathbb{E}[\mathbf{A}]$$

1. computing leading eigenvector $\mathbf{u} = [u_i]_{1 \leq i \leq n}$ of $\mathbf{A} - \frac{p+q}{2} \mathbf{1}\mathbf{1}^\top$
2. rounding: output $x_i = \begin{cases} 1, & \text{if } u_i > 0 \\ -1, & \text{if } u_i < 0 \end{cases}$

Analysis via Davis-Kahan's theorem

Let us apply Davis-Kahan to analyze accuracy of spectral clustering:

- take $M^* = \underbrace{\mathbb{E}[\mathbf{A}] - \frac{p+q}{2}\mathbf{1}\mathbf{1}^\top}_{=\frac{p-q}{2}\mathbf{x}\mathbf{x}^\top}$
 - leading eigenvector (resp. value) $\mathbf{u}^* = \frac{1}{\sqrt{n}}\mathbf{x}^*$ (resp. $\lambda^* = \frac{(p-q)n}{2}$)
- take $M = \mathbf{A} - \frac{p+q}{2}\mathbf{1}\mathbf{1}^\top$ w/ leading eigenvector \mathbf{u}
- Theorem 2.4 yields

$$\text{dist}(\mathbf{u}, \mathbf{u}^*) \leq \frac{2\|\mathbf{M} - \mathbf{M}^*\|}{\lambda_1^*} = \frac{2\|\mathbf{A} - \mathbb{E}[\mathbf{A}]\|}{\frac{(p-q)n}{2}} \quad (3.10)$$

Analysis via Davis-Kahan's theorem

Let us apply Davis-Kahan to analyze accuracy of spectral clustering:

- take $M^* = \underbrace{\mathbb{E}[\mathbf{A}] - \frac{p+q}{2}\mathbf{1}\mathbf{1}^\top}_{=\frac{p-q}{2}\mathbf{x}\mathbf{x}^\top}$
 - leading eigenvector (resp. value) $\mathbf{u}^* = \frac{1}{\sqrt{n}}\mathbf{x}^*$ (resp. $\lambda^* = \frac{(p-q)n}{2}$)
- take $M = \mathbf{A} - \frac{p+q}{2}\mathbf{1}\mathbf{1}^\top$ w/ leading eigenvector \mathbf{u}
- Theorem 2.4 yields

$$\text{dist}(\mathbf{u}, \mathbf{u}^*) \leq \frac{2\|\mathbf{M} - \mathbf{M}^*\|}{\lambda_1^*} = \frac{2\|\mathbf{A} - \mathbb{E}[\mathbf{A}]\|}{\frac{(p-q)n}{2}} \quad (3.10)$$

Question: how to bound $\|\mathbf{A} - \mathbb{E}[\mathbf{A}]\|$?

Bounding $\|\mathbf{A} - \mathbb{E}[\mathbf{A}]\|$

Lemma 3.10

Consider SBM with $p > q \gtrsim \frac{\log n}{n}$. Then with prob. $1 - O(n^{-10})$,

$$\|\mathbf{A} - \mathbb{E}[\mathbf{A}]\| \lesssim \sqrt{np} \quad (3.11)$$

proof: note that

- $\text{Var}(A_{i,j}) = \begin{cases} p(1-p) & \text{if } x_i = x_j \\ q(1-q) & \text{else} \end{cases} \leq p =: \sigma^2$
- $A_{i,j} \leq 1 =: B$

Applying Theorem 3.9 and using $p \gtrsim \frac{\log n}{n}$ conclude the proof

Statistical accuracy of spectral clustering

Substitute (3.11) into (3.10) to reach

$$\text{dist}(\mathbf{u}, \mathbf{u}^*) \leq \frac{2\|\mathbf{A} - \mathbb{E}[\mathbf{A}]\|}{\frac{(p-q)n}{2}} \lesssim \frac{\sqrt{np}}{(p-q)n}$$

provided that $(p-q)n \gg \sqrt{np}$

Thus, under condition $\frac{p-q}{\sqrt{p}} \gg \frac{1}{\sqrt{n}}$, with high prob. one has

$$\text{dist}(\mathbf{u}, \mathbf{u}^*) \ll 1 \quad \implies \quad \text{almost exact clustering}$$

Statistical accuracy of spectral clustering

$$\frac{p - q}{\sqrt{p}} \gg \frac{1}{\sqrt{n}} \implies \text{almost exact clustering} \quad (3.12)$$

- **dense regime:** if $p \asymp q \asymp 1$, then this condition reads

$$p - q \gg \frac{1}{\sqrt{n}}$$

- **“sparse” regime:** if $p = \frac{\alpha \log n}{n}$ and $q = \frac{\beta \log n}{n}$ for $\alpha, \beta \asymp 1$, then

$$\alpha - \beta \gg \frac{1}{\sqrt{\log n}}$$

This condition is information-theoretically optimal (up to log factor)
— Mossel, Neeman, Sly '15, Abbe '18

Matrix completion

Recap: spectral method for matrix completion

$$\begin{bmatrix} \checkmark & ? & ? & ? & \checkmark & ? \\ ? & ? & \checkmark & \checkmark & ? & ? \\ \checkmark & ? & ? & \checkmark & ? & ? \\ ? & ? & \checkmark & ? & ? & \checkmark \\ \checkmark & ? & ? & ? & ? & ? \\ ? & \checkmark & ? & ? & \checkmark & ? \\ ? & ? & \checkmark & \checkmark & ? & ? \end{bmatrix}$$

							...
	★★★★★	?	★★★★★	?	?	?	...
	?	★★★★★	?	?	★★★★★	?	...
	?	?	?	★★★★★	★★★★★	?	...
	?	★★★★★	★★★★★	?	?	★★★★★	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

- ground truth: $M^* = \mathbf{u}^* \mathbf{v}^{*\top} \in \mathbb{R}^{n \times n}$ (a simple case)

$$\mathbf{u}^* = \frac{1}{\|\tilde{\mathbf{u}}\|_2} \tilde{\mathbf{u}}, \quad \mathbf{v}^* = \frac{1}{\|\tilde{\mathbf{v}}\|_2} \tilde{\mathbf{v}}, \quad \tilde{\mathbf{u}}, \tilde{\mathbf{v}} \stackrel{\text{indep.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$$

- each entry $M_{i,j}^*$ is observed independently with prob. p
- goal:** fill in unseen entries of M^*

Recap: spectral method for matrix completion

1. construct a rescaled zero-filled matrix $\mathbf{M} = [M_{i,j}] \in \mathbb{R}^{n_1 \times n_2}$ as

$$\forall(i, j) : M_{i,j} = \begin{cases} \frac{1}{p} M_{i,j}^*, & \text{if } M_{i,j}^* \text{ is observed} \\ 0, & \text{else} \end{cases}$$

- o **rationale:** ensures $\mathbb{E}[\mathbf{M}] = \mathbf{M}^*$

2. compute rank-1 SVD $\sigma \mathbf{u} \mathbf{v}^\top$ of \mathbf{M} , and return $\widehat{\mathbf{M}} = \sigma \mathbf{u} \mathbf{v}^\top$

How does sampling rate p affect estimation accuracy?

Statistical accuracy of spectral estimate

From Wedin's Theorem: if $\|M - M^*\| \leq (1 - \frac{\sqrt{2}}{2})\sigma_1^* = 1 - \frac{\sqrt{2}}{2}$, then

$$\begin{aligned} \max \{ \text{dist}(\mathbf{u}, \mathbf{u}^*), \text{dist}(\mathbf{v}, \mathbf{v}^*) \} &\leq \frac{2\|M - M^*\|}{\sigma_1^*} \asymp \|M - M^*\| \\ &\lesssim \sqrt{\frac{\log^2 n}{np}} \end{aligned} \quad (3.13)$$

where last inequality is a consequence of:

Lemma 3.11

Suppose $p \gg \frac{\log n}{n}$. Then with high prob.,

$$\|M - M^*\| \lesssim \sqrt{\frac{\log^2 n}{np}} = o(1) \quad (3.14)$$

Sample complexity

For rank-1 matrix completion, (3.13) implies

$$p \gg \frac{\log^2 n}{n} \implies \text{nearly accurate estimates of } \mathbf{u}^* \text{ \& } \mathbf{v}^*$$
$$\implies \text{nearly accurate estimates of } \mathbf{M}^*$$

To yield reliable spectral estimates, it suffices to have sample size

$$\underbrace{n^2 p \asymp n \log^2 n}_{\text{optimal up to log factor}}$$

Proof of inequality (3.14)

- First, based on Gaussianity, we have

$$\frac{1}{p} \max_{i,j} |M_{i,j}^*| \lesssim \frac{\log n}{pn} =: B \quad (\text{check})$$

- Next,

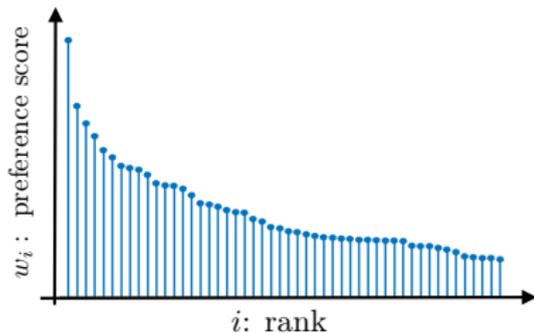
$$\max_{i,j} \text{Var}(M_{i,j}) = \frac{1-p}{p} \max_{i,j} (M_{i,j}^*)^2 \lesssim \frac{\log^2 n}{n^2 p} =: \sigma^2$$

Applying Theorem 3.9 w/ dilation trick $\|\mathbf{A}\| = \left\| \begin{bmatrix} & \mathbf{A} \\ \mathbf{A}^\top & \end{bmatrix} \right\|$ gives

$$\|\mathbf{M} - \mathbf{M}^*\| \lesssim \sigma\sqrt{n} + B\sqrt{\log n} \asymp \frac{\log n}{\sqrt{np}} + \frac{\log^{3/2} n}{np} \asymp \frac{\log n}{\sqrt{np}}$$

Ranking from pairwise comparisons

Recap: spectral ranking for BTL model



- n items with latent scores w_1^*, \dots, w_n^*
- each pair of items (i, j) is compared independently

$$\iff y_{i,j} \stackrel{\text{ind.}}{=} \begin{cases} 1, & \text{with prob. } \frac{w_j^*}{w_i^* + w_j^*} \\ 0, & \text{else} \end{cases}$$

- estimate $\mathbf{w}^* = [w_i^*]_{1 \leq i \leq n}$ (and rank items)

Recap: spectral ranking for BTL model

A key probability transition matrix $\mathbf{P}^* \in \mathbb{R}^{n \times n}$:

$$P_{i,j}^* = \begin{cases} \frac{1}{n} \cdot \frac{w_j^*}{w_i^* + w_j^*}, & \text{if } i \neq j \\ 1 - \sum_{l:l \neq i} P_{i,l}^*, & \text{if } i = j \end{cases}$$

1. construct a surrogate matrix \mathbf{P} obeying

$$P_{i,j} = \begin{cases} \frac{1}{n} y_{i,j}, & \text{if } i \neq j \\ 1 - \sum_{l:l \neq i} P_{i,l}, & \text{if } i = j \end{cases}$$

2. compute leading left eigenvector $\boldsymbol{\pi}$ of \mathbf{P} as score estimate
3. rank in accordance with $\boldsymbol{\pi}$

Can we characterize the accuracy of spectral estimates?

Analysis of spectral ranking

Apply Theorem 2.7 to yield

$$\|\pi - \pi^*\|_{\pi^*} \lesssim \frac{\|\pi^{*\top} \mathbf{E}\|_{\pi^*}}{1 - \max\{\lambda_2(\mathbf{P}^*), -\lambda_n(\mathbf{P}^*)\}}$$

with $\mathbf{E} = \mathbf{P} - \mathbf{P}^*$, provided that

$$\|\mathbf{E}\|_{\pi^*} \leq (1 - 1/\sqrt{2}) \left(1 - \max\{\lambda_2(\mathbf{P}^*), -\lambda_n(\mathbf{P}^*)\}\right)$$

Analysis of spectral ranking

Apply Theorem 2.7 to yield

$$\|\pi - \pi^*\|_{\pi^*} \lesssim \frac{\|\pi^{*\top} \mathbf{E}\|_{\pi^*}}{1 - \max\{\lambda_2(\mathbf{P}^*), -\lambda_n(\mathbf{P}^*)\}}$$

with $\mathbf{E} = \mathbf{P} - \mathbf{P}^*$, provided that

$$\|\mathbf{E}\|_{\pi^*} \leq (1 - 1/\sqrt{2}) \left(1 - \max\{\lambda_2(\mathbf{P}^*), -\lambda_n(\mathbf{P}^*)\}\right)$$

— *need to understand spectral gap and noise size*

Analysis of spectral ranking (cont.)

$$\text{condition number: } \kappa := \frac{\max_{1 \leq i \leq n} w_i^*}{\min_{1 \leq i \leq n} w_i^*}$$

Lemma 3.12 (spectral gap)

$$1 - \max \{ \lambda_2(\mathbf{P}^*), -\lambda_n(\mathbf{P}^*) \} \geq \frac{1}{2\kappa^2}$$

- proof is based on comparison between two reversible Markov chains; see Section 3.6.4 of Chen et al. '21

Lemma 3.13 (noise size)

With probability at least $1 - O(n^{-8})$,

$$\|\mathbf{E}\|_{\pi^*} \leq \sqrt{\kappa} \|\mathbf{E}\| \lesssim \sqrt{\frac{\kappa \log n}{n}}$$

Analysis of spectral ranking (cont.)

Recall perturbation bound

$$\begin{aligned}\|\pi - \pi^*\|_{\pi^*} &\leq \frac{\|\pi^{*\top} \mathbf{E}\|_{\pi^*}}{1 - \max\{\lambda_2(\mathbf{P}^*), -\lambda_n(\mathbf{P}^*)\} - \|\mathbf{E}\|_{\pi^*}} \\ &\leq 4\kappa^2 \|\pi^{*\top} \mathbf{E}\|_{\pi^*} \quad (\text{provided that } n \gg \kappa^5 \log n)\end{aligned}$$

Note that for any \mathbf{v} , one has

$$\|\mathbf{v}\|_{\pi^*} \leq \sqrt{\pi_{\max}^*} \|\mathbf{v}\|_2, \quad \text{and} \quad \|\mathbf{v}\|_2 \leq \frac{1}{\sqrt{\pi_{\min}^*}} \|\mathbf{v}\|_{\pi^*}$$

As a result, one has

$$\begin{aligned}\|\pi - \pi^*\|_2 &\leq \frac{1}{\sqrt{\pi_{\min}^*}} \|\pi - \pi^*\|_{\pi^*} \leq \frac{4\kappa^2}{\sqrt{\pi_{\min}^*}} \|\pi^{*\top} \mathbf{E}\|_{\pi^*} \\ &\leq 4\kappa^{2.5} \|\pi^{*\top} \mathbf{E}\|_2 \leq 4\kappa^{2.5} \|\mathbf{E}\| \|\pi^*\|_2\end{aligned}$$

Analysis of spectral ranking (cont.)

Assuming $\kappa = O(1)$, we arrive at

$$\|\boldsymbol{\pi} - \boldsymbol{\pi}^*\|_2 \lesssim \sqrt{\frac{\log n}{n}} \|\boldsymbol{\pi}^*\|_2$$

- vanishingly small error as $n \rightarrow \infty$
- optimal error up to log factor

— Negahban, Oh, Shah '16, Chen, Fan, Ma, Wang '19

Proof of Lemma 3.13

By construction of P and P^* , we see that

$$E_{i,j} = P_{i,j} - P_{i,j}^* = \frac{1}{n} (y_{i,j} - \mathbb{E}[y_{i,j}])$$

for any $i \neq j$. In addition, for all $1 \leq i \leq n$, it follows that

$$E_{i,i} = P_{i,i} - P_{i,i}^* = - \sum_{j:j \neq i} E_{i,j} = -\frac{1}{n} \sum_{j:j \neq i} (y_{i,j} - \mathbb{E}[y_{i,j}])$$

We shall decompose E into three parts: upper triangular, diagonal, and lower triangular parts:

$$\|E\| \leq \|E_{\text{upper}}\| + \|E_{\text{diag}}\| + \|E_{\text{lower}}\|$$

— we will upper bound $\|E_{\text{upper}}\|$

Proof of Lemma 3.13 (controlling $\|E_{\text{diag}}\|$)

Observe that

$$\|E_{\text{diag}}\| = \max_{1 \leq i \leq n} |E_{i,i}| = \max_{1 \leq i \leq n} \frac{1}{n} \left| \underbrace{\sum_{j:j \neq i} (y_{i,j} - \mathbb{E}[y_{i,j}])}_{=: X_j} \right|$$

To invoke Bernstein's inequality, note that

- $|X_j| \leq 1 =: B$
- $\sum_{j:j \neq i} \mathbb{E}[X_j^2] = \sum_{j:j \neq i} \text{Var}(y_{i,j}) \leq n =: v$

Bernstein's inequality + union bound reveal that: with high prob.

$$\max_i |E_{i,i}| \lesssim \frac{1}{n} (\sqrt{v \log n} + B \log n) \asymp \sqrt{\frac{\log n}{n}}$$

Proof of Lemma 3.13 (controlling $\|\mathbf{E}_{\text{upper}}\|$)

- $\frac{1}{n}|y_{i,j}| \leq \frac{1}{n} =: B$
- $\text{Var}(\frac{1}{n}y_{i,j}) \leq \frac{1}{n^2} =: \sigma^2$

Applying Theorem 3.9 w/ dilation trick $\|\mathbf{A}\| = \left\| \begin{bmatrix} & \mathbf{A}^\top \\ \mathbf{A} & \end{bmatrix} \right\|$ gives

$$\|\mathbf{E}_{\text{upper}}\| \lesssim \sigma\sqrt{n} + B\sqrt{\log n} \asymp \frac{1}{\sqrt{n}} + \frac{\sqrt{\log n}}{n} \asymp \frac{1}{\sqrt{n}}$$

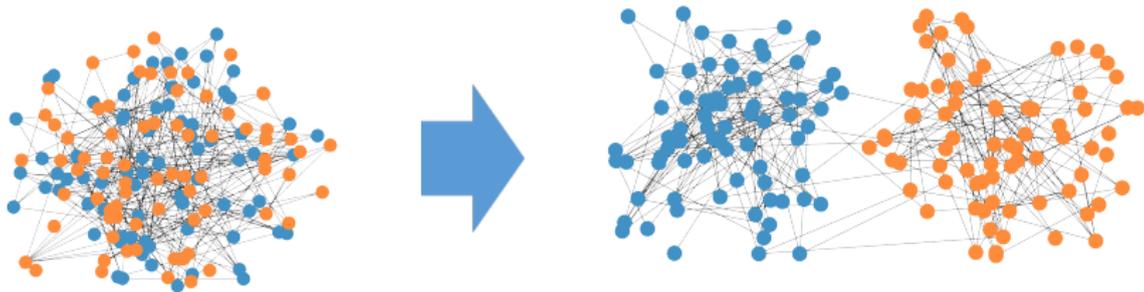
— same bound holds for $\|\mathbf{E}_{\text{lower}}\|$

Part 4: ℓ_∞ and $\ell_{2,\infty}$ perturbation theory

- Motivation
 - exact community recovery
 - top- K ranking
- Leave-one-out analysis: an illustrative example
- ℓ_∞ eigenvector perturbation theory (rank-1)
- Application: exact recovery in community detection
- $\ell_{2,\infty}$ eigen-space perturbation theory (rank- r)

Motivation: exact community recovery

Revisiting spectral clustering for SBMs



1. computing the leading eigenvector $\mathbf{u} = [u_i]_{1 \leq i \leq n}$ of $\mathbf{A} - \frac{p+q}{2} \mathbf{1}\mathbf{1}^\top$
2. rounding: output $x_i = \begin{cases} 1, & \text{if } u_i \geq 0 \\ -1, & \text{if } u_i < 0 \end{cases}$

Revisiting spectral clustering for SBMs

It has been shown in (3.12) that: if $p \asymp q \asymp \frac{\log n}{n}$, then

$$\delta := p - q \gg \frac{\sqrt{\log n}}{n} \implies \text{almost exact recovery}$$

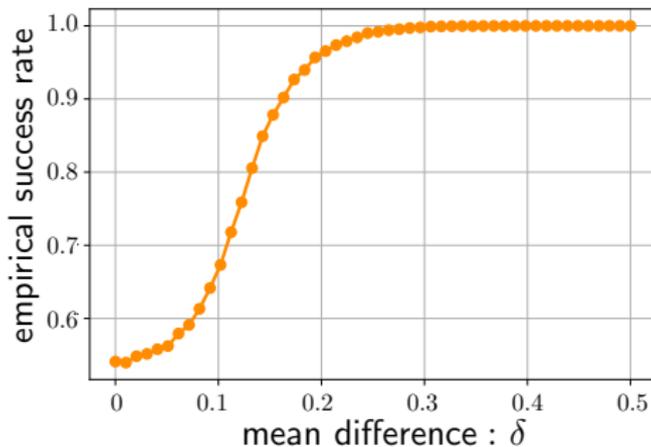
- Almost exact recovery means

$$\min \left\{ \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i \neq x_i^*\}, \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i \neq -x_i^*\} \right\} = o(1)$$

Exact recovery of all community memberships?

When $\delta := p - q$ increases, exact recovery becomes possible:

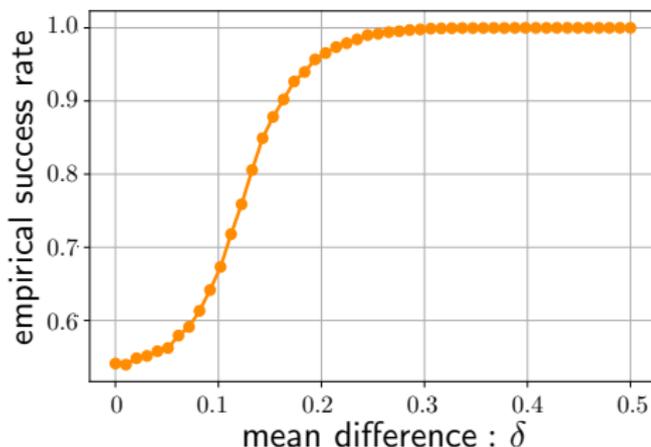
$$\min \left\{ \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i \neq x_i^*\}, \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i \neq -x_i^*\} \right\} = 0$$



Exact recovery of all community memberships?

When $\delta := p - q$ increases, exact recovery becomes possible:

$$\min \left\{ \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i \neq x_i^*\}, \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i \neq -x_i^*\} \right\} = 0$$



ℓ_2 perturbation theory alone falls short of explaining exact recovery

— calls for more **fine-grained** analysis

Exact recovery $\leftarrow \ell_\infty$ theory

exact recovery means $u_i u_i^* > 0, \forall i$ (or $u_i u_i^* < 0, \forall i$)

Exact recovery $\leftarrow \ell_\infty$ theory

exact recovery means $u_i u_i^* > 0, \forall i$ (or $u_i u_i^* < 0, \forall i$)



$$\|\mathbf{u} - \mathbf{u}^*\|_\infty < 1/\sqrt{n} \quad \text{or} \quad \|\mathbf{u} + \mathbf{u}^*\|_\infty < 1/\sqrt{n}$$

Exact recovery $\leftarrow \ell_\infty$ theory

exact recovery means $u_i u_i^* > 0, \forall i$ (or $u_i u_i^* < 0, \forall i$)



$$\|\mathbf{u} - \mathbf{u}^*\|_\infty < 1/\sqrt{n} \quad \text{or} \quad \|\mathbf{u} + \mathbf{u}^*\|_\infty < 1/\sqrt{n}$$



ℓ_∞ eigenvector perturbation theory

Motivation: top- K ranking

Top- K ranking

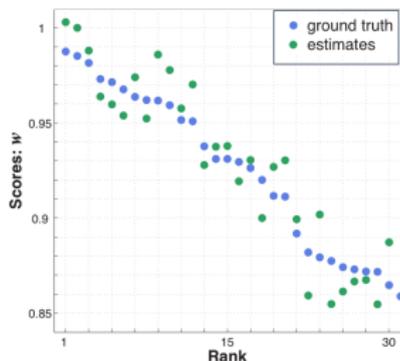
Goal: identify the set of top- K ranked items



Typical ranking procedure:

- estimate latent scores
- return top- K items in accordance with score estimates

Top- K ranking for BTL model



Top 3: {1, 2, 3}

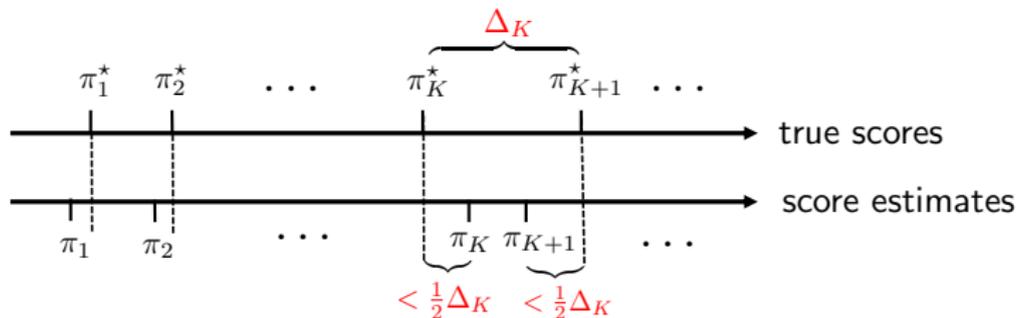


1. construct a surrogate matrix \mathbf{P} obeying

$$P_{i,j} = \begin{cases} \frac{1}{n} y_{i,j}, & \text{if } i \neq j \\ 1 - \sum_{l:l \neq i} P_{i,l}, & \text{if } i = j \end{cases}$$

2. compute leading left eigenvector $\boldsymbol{\pi}$ of \mathbf{P} as score estimate
3. return K items associated with largest score estimate π_i

Controlling entrywise estimation error



exact top- K ranking



$$2\|\pi - \pi^*\|_\infty < \Delta_K := \pi_{(K)}^* - \pi_{(K+1)}^*$$



ℓ_∞ eigenvector perturbation theory

Leave-one-out analysis: an illustrative example

from random matrix theory, stat. physics, etc

Setup and algorithm

- Ground truth: $M^* = \lambda^* \mathbf{u}^* \mathbf{u}^{*\top} \in \mathbb{R}^{n \times n}$, with $\lambda^* > 0$
 - $\|\mathbf{u}^*\|_2 = 1$, $\|\mathbf{u}^*\|_\infty = \sqrt{\mu/n}$ (μ : incoherence parameter)
- Observation: $M = M^* + E$
 - E : symmetric, entries in upper triangular part are i.i.d. $\mathcal{N}(0, \sigma^2)$
- Estimate \mathbf{u}^* using leading eigenvector \mathbf{u} of M

Question: can we characterize entrywise estimation error of \mathbf{u} , i.e.

$$\text{dist}_\infty(\mathbf{u}, \mathbf{u}^*) := \min \{ \|\mathbf{u} - \mathbf{u}^*\|_\infty, \|\mathbf{u} + \mathbf{u}^*\|_\infty \}$$

ℓ_2 guarantees

Davis-Kahan's $\sin \Theta$ theorem together with Theorem 3.9 gives

$$\text{dist}(\mathbf{u}, \mathbf{u}^*) \leq \frac{2\|\mathbf{E}\|}{\lambda^*} \leq \frac{10\sigma\sqrt{n}}{\lambda^*}$$

with high prob., as long as $\sigma\sqrt{n} \leq \frac{1-1/\sqrt{2}}{5}\lambda^*$

- as an immediate (but very crude) consequence

$$\text{dist}_\infty(\mathbf{u}, \mathbf{u}^*) \leq \text{dist}(\mathbf{u}, \mathbf{u}^*) \lesssim \frac{\sigma\sqrt{n}}{\lambda^*} \quad (4.15)$$

l_∞ guarantees for matrix denoising

Theorem 4.14

Suppose that $\sigma\sqrt{n} \leq c_0\lambda^*$ for some sufficiently small constant $c_0 > 0$. Then with high prob.,

$$\text{dist}_\infty(\mathbf{u}, \mathbf{u}^*) \lesssim \frac{\sigma(\sqrt{\log n} + \sqrt{\mu})}{\lambda^*}$$

- When $\mu \lesssim \log n$ (i.e. energy of \mathbf{u}^* is spread out):

$$\text{dist}_\infty(\mathbf{u}, \mathbf{u}^*) \lesssim \frac{\sigma\sqrt{\log n}}{\lambda^*}$$

- Much sharper (i.e. $\sqrt{n/\log n}$ times better) than (4.15)

Technical hurdle: statistical dependency

Let's take close inspection of l -th entry u_l of \mathbf{u} :

- Given that \mathbf{u} is an eigenvector of \mathbf{M} , we have

$$\mathbf{M}\mathbf{u} = \lambda\mathbf{u},$$

$$\implies u_l = \frac{1}{\lambda}[\mathbf{M}]_{l,:}\mathbf{u} = \frac{1}{\lambda}[\mathbf{M}^* + \mathbf{E}]_{l,:}\mathbf{u}$$

- challenge:** \mathbf{u} is statistically dependent on \mathbf{E} (in a complicated way)!

How to decouple complicated dependency between \mathbf{u} and \mathbf{E} ?

Decomposition w/ the aid of an independent proxy

Suppose we have access to a proxy $\mathbf{u}^{(l)}$ independent of $\mathbf{E}_{l,:}$, then

$$\underbrace{\mathbf{E}_{l,:}\mathbf{u}}_{\text{a term of interest}} = \underbrace{\mathbf{E}_{l,:}\mathbf{u}^{(l)}}_{=: \mathcal{J}_1} + \underbrace{\mathbf{E}_{l,:}(\mathbf{u} - \mathbf{u}^{(l)})}_{=: \mathcal{J}_2}$$

- \mathcal{J}_1 : controllable using independence btw $\mathbf{u}^{(l)}$ & $\mathbf{E}_{l,:}$
- \mathcal{J}_2 : small if $\mathbf{u}^{(l)} \approx \mathbf{u}$

How to construct a useful proxy?

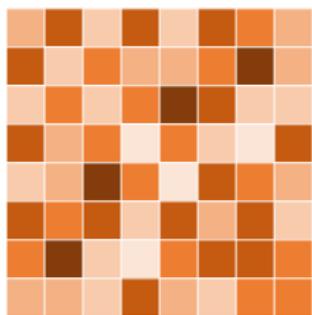
Leave-one-out auxiliary estimates

For each $1 \leq l \leq n$, construct an auxiliary matrix $M^{(l)}$

$$M^{(l)} := \lambda^* \mathbf{u}^* \mathbf{u}^{*\top} + \mathbf{E}^{(l)},$$

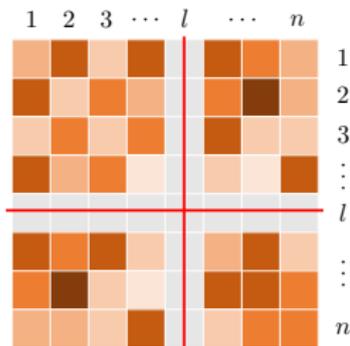
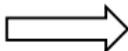
where the noise matrix $\mathbf{E}^{(l)}$ is generated according to

$$E_{i,j}^{(l)} := \begin{cases} E_{i,j}, & \text{if } i \neq l \text{ and } j \neq l \\ 0, & \text{else} \end{cases} \quad (\text{removing } l\text{-th row/col})$$



M

leave one
row/column out



$M^{(l)}$

Leave-one-out auxiliary estimates

For each $1 \leq l \leq n$, construct an auxiliary matrix $\mathbf{M}^{(l)}$

$$\mathbf{M}^{(l)} := \lambda^* \mathbf{u}^* \mathbf{u}^{*\top} + \mathbf{E}^{(l)},$$

where the noise matrix $\mathbf{E}^{(l)}$ is generated according to

$$E_{i,j}^{(l)} := \begin{cases} E_{i,j}, & \text{if } i \neq l \text{ and } j \neq l \\ 0, & \text{else} \end{cases} \quad (\text{removing } l\text{-th row/col})$$

leave-one-out estimates: $\mathbf{u}^{(l)}$ (resp. $\lambda^{(l)}$) is leading eigenvalue (resp. eigenvector) of $\mathbf{M}^{(l)}$

- **key property:** $\mathbf{u}^{(l)}$ is independent of $\mathbf{E}_{l,:}$

Intuition

WLOG, suppose $\mathbf{u}^\top \mathbf{u}^* > 0$ and $\mathbf{u}^{(l)\top} \mathbf{u}^* > 0 \dots$

- **proximity of $\mathbf{u}^{(l)}$ and \mathbf{u} :** since $\mathbf{u}^{(l)}$ is obtained by dropping only a tiny fraction of data, we expect $\mathbf{u} \approx \mathbf{u}^{(l)}$
- **proximity of $u_l^{(l)}$ and u_l^* :** by construction,

$$\begin{aligned} u_l^{(l)} &= \frac{1}{\lambda^{(l)}} \mathbf{M}_{l,\cdot}^{(l)} \mathbf{u}^{(l)} = \frac{1}{\lambda^{(l)}} \mathbf{M}_{l,\cdot}^* \mathbf{u}^{(l)} = \frac{\lambda^*}{\lambda^{(l)}} u_l^* \mathbf{u}^{*\top} \mathbf{u}^{(l)} \\ &\approx u_l^* \end{aligned}$$

Proof of Theorem 4.14

What we have learned from ℓ_2 analysis

$$\|\mathbf{E}\| \leq 5\sigma\sqrt{n}$$

$$\text{dist}(\mathbf{u}, \mathbf{u}^*) \leq \frac{10\sigma\sqrt{n}}{\lambda^*}$$

$$|\lambda - \lambda^*| \leq 5\sigma\sqrt{n}$$

$$\max_{j:j \geq 2} |\lambda_j(\mathbf{M})| \leq 5\sigma\sqrt{n}$$

$$\|\mathbf{E}^{(l)}\| \leq \|\mathbf{E}\| \leq 5\sigma\sqrt{n}$$

$$\text{dist}(\mathbf{u}^{(l)}, \mathbf{u}^*) \leq \frac{10\sigma\sqrt{n}}{\lambda^*}$$

$$|\lambda^{(l)} - \lambda^*| \leq 5\sigma\sqrt{n}$$

$$\max_{j:j \geq 2} |\lambda_j(\mathbf{M}^{(l)})| \leq 5\sigma\sqrt{n}$$

Addressing ambiguity

WLOG, assume

$$\begin{aligned}\|\mathbf{u} - \mathbf{u}^*\|_2 &= \text{dist}(\mathbf{u}, \mathbf{u}^*), \\ \|\mathbf{u}^{(l)} - \mathbf{u}^*\|_2 &= \text{dist}(\mathbf{u}^{(l)}, \mathbf{u}^*), \quad 1 \leq l \leq n\end{aligned}$$

A useful byproduct: if $c_0\sigma\sqrt{n} < \lambda^*$ for some small constant $c_0 > 0$, then one necessarily has (exercise)

$$\|\mathbf{u} - \mathbf{u}^{(l)}\|_2 = \text{dist}(\mathbf{u}, \mathbf{u}^{(l)}), \quad 1 \leq l \leq n$$

Key steps

To bound $u_l - u_l^*$ ($1 \leq l \leq n$), we see from triangle inequality that

$$|u_l - u_l^*| \leq |u_l^{(l)} - u_l^*| + \|\mathbf{u} - \mathbf{u}^{(l)}\|_\infty \leq |u_l^{(l)} - u_l^*| + \|\mathbf{u} - \mathbf{u}^{(l)}\|_2$$

Key steps

To bound $u_l - u_l^*$ ($1 \leq l \leq n$), we see from triangle inequality that

$$|u_l - u_l^*| \leq |u_l^{(l)} - u_l^*| + \|\mathbf{u} - \mathbf{u}^{(l)}\|_\infty \leq |u_l^{(l)} - u_l^*| + \|\mathbf{u} - \mathbf{u}^{(l)}\|_2$$

- control $\|\mathbf{u} - \mathbf{u}^{(l)}\|_2$ (Davis-Kahan)
- control $|u_l^{(l)} - u_l^*|$ ($u_l^{(l)}$ is independent from $\mathbf{E}_{l,:}$)

Bounding proximity $\|\mathbf{u} - \mathbf{u}^{(l)}\|_2$

Key: view \mathbf{M} as perturbation of $\mathbf{M}^{(l)}$, apply Davis-Kahan

$$\|\mathbf{u} - \mathbf{u}^{(l)}\|_2 \leq \frac{2\|(\mathbf{M} - \mathbf{M}^{(l)})\mathbf{u}^{(l)}\|_2}{\lambda^{(l)} - \max_{j \geq 2} |\lambda_j(\mathbf{M}^{(l)})|} \leq \frac{4\|(\mathbf{M} - \mathbf{M}^{(l)})\mathbf{u}^{(l)}\|_2}{\lambda^*}$$

as long as

$$\|\mathbf{M} - \mathbf{M}^{(l)}\| \leq (1 - 1/\sqrt{2}) \left(\lambda^{(l)} - \max_{j \geq 2} |\lambda_j(\mathbf{M}^{(l)})| \right),$$

$$\lambda^{(l)} - \max_{j \geq 2} |\lambda_j(\mathbf{M}^{(l)})| \geq \lambda^*/2$$

Bounding $\|(M - M^{(l)})\mathbf{u}^{(l)}\|_2$

By design,

$$(M - M^{(l)})\mathbf{u}^{(l)} = \mathbf{e}_l \mathbf{E}_{l,\cdot} \mathbf{u}^{(l)} + u_l^{(l)} (\mathbf{E}_{\cdot,l} - E_{l,l} \mathbf{e}_l),$$

which together with triangle inequality yields

$$\begin{aligned} \|(M - M^{(l)})\mathbf{u}^{(l)}\|_2 &\leq \underbrace{\|\mathbf{E}_{l,\cdot} \mathbf{u}^{(l)}\|_2}_{\mathbf{E}_{l,\cdot} \text{ and } \mathbf{u}^{(l)} \text{ are independent}} + \|\mathbf{E}_{\cdot,l}\|_2 \cdot |u_l^{(l)}| \\ &\leq 5\sigma \sqrt{\log n} + \|\mathbf{E}_{\cdot,l}\|_2 (|u_l| + \|\mathbf{u} - \mathbf{u}^{(l)}\|_\infty) \\ &\leq 5\sigma \sqrt{\log n} + 5\sigma \sqrt{n} \|\mathbf{u}\|_\infty + 5\sigma \sqrt{n} \|\mathbf{u} - \mathbf{u}^{(l)}\|_2 \end{aligned}$$

Bounding $\|\mathbf{u} - \mathbf{u}^{(l)}\|_2$ (cont.)

Combining previous bounds, we arrive at

$$\begin{aligned}\|\mathbf{u} - \mathbf{u}^{(l)}\|_2 &\leq \frac{20\sigma\sqrt{\log n} + 20\sigma\sqrt{n}\|\mathbf{u}\|_\infty + 20\sigma\sqrt{n}\|\mathbf{u} - \mathbf{u}^{(l)}\|_2}{\lambda^*} \\ &\leq \frac{20\sigma\sqrt{\log n} + 20\sigma\sqrt{n}\|\mathbf{u}\|_\infty}{\lambda^*} + \frac{1}{2}\|\mathbf{u} - \mathbf{u}^{(l)}\|_2\end{aligned}$$

provided that $40\sigma\sqrt{n} \leq \lambda^*$

Rearranging terms and taking union bound give: with high prob.

$$\|\mathbf{u} - \mathbf{u}^{(l)}\|_2 \leq \frac{40\sigma\sqrt{\log n} + 40\sigma\sqrt{n}\|\mathbf{u}\|_\infty}{\lambda^*} \quad 1 \leq l \leq n$$

Analyzing leave-one-out iterates

Recall that

$$u_i^{(l)} = \frac{1}{\lambda^{(l)}} \mathbf{M}_{i,\cdot}^{(l)} \mathbf{u}^{(l)} = \frac{1}{\lambda^{(l)}} \mathbf{M}_{i,\cdot}^* \mathbf{u}^{(l)} = \frac{\lambda^*}{\lambda^{(l)}} u_i^* \mathbf{u}^{*\top} \mathbf{u}^{(l)}$$

This implies

$$\begin{aligned} u_i^{(l)} - u_i^* &= u_i^* \left(\frac{\lambda^*}{\lambda^{(l)}} \mathbf{u}^{*\top} \mathbf{u}^{(l)} - \mathbf{u}^{*\top} \mathbf{u}^* \right) \\ &= u_i^* \left(\frac{\lambda^* - \lambda^{(l)}}{\lambda^{(l)}} \mathbf{u}^{*\top} \mathbf{u}^{(l)} \right) + u_i^* \mathbf{u}^{*\top} (\mathbf{u}^{(l)} - \mathbf{u}^*) \end{aligned}$$

Analyzing leave-one-out iterates (cont.)

Triangle inequality gives

$$\begin{aligned} |u_l^{(l)} - u_l^*| &\leq |u_l^*| \cdot \frac{|\lambda^* - \lambda^{(l)}|}{|\lambda^{(l)}|} \cdot \|\mathbf{u}^*\|_2 \cdot \|\mathbf{u}^{(l)}\|_2 \\ &\quad + |u_l^*| \cdot \|\mathbf{u}^*\|_2 \cdot \|\mathbf{u}^{(l)} - \mathbf{u}^*\|_2 \\ &\leq |u_l^*| \cdot \frac{10\sigma\sqrt{n}}{\lambda^*} + |u_l^*| \cdot \frac{10\sigma\sqrt{n}}{\lambda^*} \\ &\leq \frac{20\sigma\sqrt{n}}{\lambda^*} \|\mathbf{u}^*\|_\infty \end{aligned}$$

Putting all pieces together

Now we come to conclude that

$$\begin{aligned}\|\mathbf{u} - \mathbf{u}^*\|_\infty &= \max_l |u_l - u_l^*| \leq \max_l \left\{ |u_l^{(l)} - u_l^*| + \|\mathbf{u} - \mathbf{u}^{(l)}\|_2 \right\} \\ &\leq \frac{20\sigma\sqrt{n}}{\lambda^*} \|\mathbf{u}^*\|_\infty + \frac{40\sigma\sqrt{\log n} + 40\sigma\sqrt{n}\|\mathbf{u}\|_\infty}{\lambda^*}\end{aligned}$$

One more triangle inequality gives

$$\|\mathbf{u} - \mathbf{u}^*\|_\infty \leq \frac{40\sigma\sqrt{\log n} + 60\sigma\sqrt{n}\|\mathbf{u}^*\|_\infty}{\lambda^*} + \frac{1}{2}\|\mathbf{u} - \mathbf{u}^*\|_\infty$$

provided that $80\sigma\sqrt{n} \leq \lambda^*$. Rearranging terms yields

$$\|\mathbf{u} - \mathbf{u}^*\|_\infty \leq \frac{80\sigma\sqrt{\log n} + 120\sigma\sqrt{n}\|\mathbf{u}^*\|_\infty}{\lambda^*} = \underbrace{\frac{80\sigma\sqrt{\log n} + 120\sigma\sqrt{\mu}}{\lambda^*}}_{\text{from definition of } \mu}$$

l_∞ eigenvector perturbation theory (rank-1)

Setup and algorithm

- Ground truth: $\mathbf{M}^* = \lambda^* \mathbf{u}^* \mathbf{u}^{*\top} \in \mathbb{R}^{n \times n}$, with $\lambda^* > 0$
 - $\|\mathbf{u}^*\|_2 = 1$, $\|\mathbf{u}^*\|_\infty = \sqrt{\mu/n}$ (μ : incoherence parameter)
- Observation: $\mathbf{M} = \mathbf{M}^* + \mathbf{E}$ with symmetric \mathbf{E}
- Estimate \mathbf{u}^* using leading eigenvector \mathbf{u} of \mathbf{M}

Question: can we accommodate more general noise distributions beyond Gaussian?

Noise assumptions

Entries in lower triangular part of $\mathbf{E} = [E_{i,j}]_{1 \leq i, j \leq n}$ are independently generated obeying

$$\mathbb{E}[E_{i,j}] = 0, \quad \mathbb{E}[E_{i,j}^2] \leq \sigma^2, \quad |E_{i,j}| \leq B, \quad \text{for all } i \geq j$$

Further, assume that

$$c_b := \frac{B}{\sigma \sqrt{n/(\mu \log n)}} = O(1)$$

- in general, B is allowed to be significantly larger than σ

Theorem 4.15

With high prob, there exists $z \in \{1, -1\}$ such that

$$\|z\mathbf{u} - \mathbf{u}^*\|_\infty \lesssim \frac{\sigma\sqrt{\mu} + \sigma\sqrt{\log n}}{\lambda^*} \quad (4.18a)$$

$$\left\| z\mathbf{u} - \frac{1}{\lambda^*} \mathbf{M}\mathbf{u}^* \right\|_\infty \lesssim \frac{\sigma\sqrt{\mu}}{\lambda^*} + \frac{\sigma^2\sqrt{n\log n} + \sigma B\sqrt{\mu\log^3 n}}{(\lambda^*)^2} \quad (4.18b)$$

provided $\sigma\sqrt{n\log n} \leq c_\sigma\lambda^*$ for some small enough constant $c_\sigma > 0$

- estimation error is delocalized (recall that $\text{dist}(\mathbf{u}, \mathbf{u}^*) \lesssim \sigma\sqrt{n}$)

First-order expansion

Theorem 4.15 reveals tightness of first-order approximation

$$\mathbf{u} = \frac{M\mathbf{u}}{\lambda} \approx \frac{M\mathbf{u}^*}{\lambda^*} \approx \frac{M^*\mathbf{u}^*}{\lambda^*} = \mathbf{u}^*$$

- (4.18b) often leads to tighter approximation than (4.18a)
- important in certain applications such as SBM

Application: exact recovery in community detection

Exact recovery using spectral methods

Consider the case where

$$p = \frac{\alpha \log n}{n}, \quad q = \frac{\beta \log n}{n}$$

Theorem 4.16

Fix any constant $\varepsilon > 0$. Suppose $\alpha > \beta > 0$ are large enough^a and

$$(\sqrt{\alpha} - \sqrt{\beta})^2 \geq 2(1 + \varepsilon)$$

With probability $1 - o(1)$, spectral method achieves exact recovery

^athis assumption can be removed

Optimality of spectral method

Lower bound: if

$$(\sqrt{\alpha} - \sqrt{\beta})^2 \leq 2(1 - \varepsilon) \quad (4.19)$$

for any constant $\varepsilon > 0$, then no method can achieve exact recovery

- taking this w/ Theorem 4.16 reveals information-theoretic optimality of spectral method

What is the operational meaning of $(\sqrt{\alpha} - \sqrt{\beta})^2$ or $(\sqrt{p} - \sqrt{q})^2$?

Squared Hellinger distance

Definition 4.17

Consider two distributions P and Q over a finite alphabet \mathcal{Y} . The squared Hellinger distance $H^2(P, Q)$ between P and Q is

$$H^2(P, Q) := \frac{1}{2} \sum_{y \in \mathcal{Y}} \left(\sqrt{P(y)} - \sqrt{Q(y)} \right)^2$$

- squared Hellinger distance between $\text{Bern}(p)$ and $\text{Bern}(q)$:

$$\begin{aligned} H^2(\text{Bern}(p), \text{Bern}(q)) &:= \frac{1}{2}(\sqrt{p} - \sqrt{q})^2 + \frac{1}{2}(\sqrt{1-p} - \sqrt{1-q})^2 \\ &= (1 + o(1)) \frac{1}{2}(\sqrt{p} - \sqrt{q})^2 \end{aligned}$$

when $p = o(1)$ and $q = o(1)$

Optimality of spectral method (cont.)

Theorem 4.16 and lower bound (4.19) reveal sharp phase transition:

spectral method works if $H^2(\text{Bern}(p), \text{Bern}(q)) \geq (1 + \varepsilon) \frac{\log n}{n}$

no algorithm works if $H^2(\text{Bern}(p), \text{Bern}(q)) \leq (1 - \varepsilon) \frac{\log n}{n}$

for arbitrarily small constant $\varepsilon > 0$

Fine-grained analysis of spectral clustering

WLOG, assume $x_1^* = \dots = x_{n/2}^* = 1$ and $x_{n/2+1}^* = \dots = x_n^* = -1$, and recall that

$$\mathbf{M}^* := \mathbb{E}[\mathbf{A}] - \frac{p+q}{2} \mathbf{1}\mathbf{1}^\top = \frac{p-q}{2} \mathbf{x}^* \mathbf{x}^{*\top}$$

These imply

$$\begin{aligned} \lambda^* &= \frac{n(p-q)}{2}, & \mu &= 1, \\ B &= 1, & \sigma^2 &\leq \max\{p, q\} = p \end{aligned}$$

Applying ℓ_∞ perturbation theory

ℓ_∞ perturbation bound (4.18b) yields: for some constant $C > 0$,

$$\begin{aligned}\|z\lambda^* \mathbf{u} - \mathbf{M}\mathbf{u}^*\|_\infty &\lesssim \sigma + \frac{\sigma^2 \sqrt{n \log n}}{\lambda^*} + \frac{\sigma B \log^{3/2} n}{\lambda^*} \\ &\leq C \left(\sqrt{p} + \frac{p \sqrt{\log n}}{\sqrt{n}(p-q)} + \frac{\sqrt{p} \log^{3/2} n}{n(p-q)} \right) =: \Delta\end{aligned}\tag{4.20}$$

It boils down to characterizing entrywise behavior of $\mathbf{M}\mathbf{u}^*$
— what happens at phase transition point $(\sqrt{\alpha} - \sqrt{\beta})^2 = 2$?

Bounding entries in $M\mathbf{u}^*$

Lemma 4.18

Suppose that $(\sqrt{\alpha} - \sqrt{\beta})^2 \geq 2(1 + \varepsilon)$ for some constant $\varepsilon > 0$.
Then with prob. $1 - o(1)$,

$$M_{l,\cdot}\mathbf{u}^* \geq \frac{\eta \log n}{\sqrt{n}} \text{ for all } l \leq \frac{n}{2}, \quad M_{l,\cdot}\mathbf{u}^* \leq -\frac{\eta \log n}{\sqrt{n}} \text{ for all } l > \frac{n}{2}$$

where $\eta > 0$ obeys $(\sqrt{\alpha} - \sqrt{\beta})^2 - \eta \log(\alpha/\beta) > 2$

key message: entries of $M\mathbf{u}^*$ are bounded away from 0 with correct signs if $(\sqrt{\alpha} - \sqrt{\beta})^2 > 2$

Completing the picture

Combine Lemma 4.18 with (4.20) leads to a claim: if

$$\frac{\eta \log n}{\sqrt{n}} > \Delta \tag{4.21}$$

then it follows that

$$z u_l u_l^* > 0 \quad \text{for all } 1 \leq l \leq n \quad \implies \quad \text{exact recovery}$$

Proof of relation (4.21)

Lemma 4.18 and (4.20) tell us that: it suffices to show

$$\frac{\eta \log n}{\sqrt{n}} \geq C \left(\sqrt{p} + \frac{p\sqrt{\log n}}{\sqrt{n}(p-q)} + \frac{\sqrt{p} \log^{3/2} n}{n(p-q)} \right)$$

- 1st term: $\sqrt{p} \asymp \sqrt{\frac{\log n}{n}} \ll \frac{\eta \log n}{\sqrt{n}}$
- 2nd term: $\frac{p\sqrt{\log n}}{\sqrt{n}(p-q)} \asymp \sqrt{\frac{\log n}{n}} \ll \frac{\eta \log n}{\sqrt{n}}$
- 3rd term: divide discussion into two cases $\alpha/\beta \leq 2$, and $\alpha/\beta \geq 2$

Comparing two sets of Bernoulli r.v.s

Lemma 4.19

Suppose $\alpha > \beta$, $\{W_i\}_{1 \leq i \leq n/2}$ are i.i.d. $\text{Bern}(\frac{\alpha \log n}{n})$, and $\{Z_i\}_{1 \leq i \leq n/2}$ are i.i.d. $\text{Bern}(\frac{\beta \log n}{n})$, which are independent of W_i . For any $t > 0$,

$$\mathbb{P} \left(\sum_{i=1}^{n/2} W_i - \sum_{i=1}^{n/2} Z_i \leq t \log n \right) \leq n^{-(\sqrt{\alpha} - \sqrt{\beta})^2/2 + t \log(\alpha/\beta)/2}$$

Comparing two sets of Bernoulli r.v.s

Lemma 4.19

Suppose $\alpha > \beta$, $\{W_i\}_{1 \leq i \leq n/2}$ are i.i.d. $\text{Bern}(\frac{\alpha \log n}{n})$, and $\{Z_i\}_{1 \leq i \leq n/2}$ are i.i.d. $\text{Bern}(\frac{\beta \log n}{n})$, which are independent of W_i . For any $t > 0$,

$$\mathbb{P} \left(\sum_{i=1}^{n/2} W_i - \sum_{i=1}^{n/2} Z_i \leq t \log n \right) \leq n^{-(\sqrt{\alpha} - \sqrt{\beta})^2/2 + t \log(\alpha/\beta)/2}$$

- roughly speaking,

$$\mathbb{P} \left(\sum_{i=1}^{n/2} W_i - \sum_{i=1}^{n/2} Z_i \leq 0 \right) \leq n^{-(\sqrt{\alpha} - \sqrt{\beta})^2/2}$$

Comparing two sets of Bernoulli r.v.s

Lemma 4.19

Suppose $\alpha > \beta$, $\{W_i\}_{1 \leq i \leq n/2}$ are i.i.d. $\text{Bern}(\frac{\alpha \log n}{n})$, and $\{Z_i\}_{1 \leq i \leq n/2}$ are i.i.d. $\text{Bern}(\frac{\beta \log n}{n})$, which are independent of W_i . For any $t > 0$,

$$\mathbb{P} \left(\sum_{i=1}^{n/2} W_i - \sum_{i=1}^{n/2} Z_i \leq t \log n \right) \leq n^{-(\sqrt{\alpha} - \sqrt{\beta})^2/2 + t \log(\alpha/\beta)/2}$$

- $(\sqrt{\alpha} - \sqrt{\beta})^2 > 2$ guarantees $\mathbb{P}(\sum_{i=1}^{n/2} W_i - \sum_{i=1}^{n/2} Z_i \leq 0) < n^{-1}$
 - probability of error $o(n^{-1})$ is crucial, since in $\mathbf{M}\mathbf{u}^*$ we have n independent groups of $\{W_i\}$ and $\{Z_i\}$ (need union bound)

Proof of Lemma 4.18

Note that $M\mathbf{u}^* = (\mathbf{A} - \frac{p+q}{2}\mathbf{1}\mathbf{1}^\top)\mathbf{u}^* = \mathbf{A}\mathbf{u}^*$. Hence

$$M_{1,:}\mathbf{u}^* = \mathbf{A}_{1,:}\mathbf{u}^* = \frac{1}{\sqrt{n}} \left(\sum_{j=1}^{n/2} A_{1,j} - \sum_{j=n/2+1}^n A_{1,j} \right)$$

Apply Lemma 4.19 to obtain with probability at least $1 - n^{-(\sqrt{a}-\sqrt{b})^2/2+\eta \log(a/b)/2} = 1 - o(n^{-1})$

$$M_{1,:}\mathbf{u}^* \geq \frac{\eta \log n}{\sqrt{n}}$$

Invoke union bound to complete proof

Proof of Lemma 4.19

We apply the Laplace transform method: for any $\lambda < 0$

$$\begin{aligned} & \mathbb{P} \left(\sum_{i=1}^{n/2} W_i - \sum_{i=1}^{n/2} Z_i \leq t \log n \right) \\ &= \mathbb{P} \left(\exp \left(\lambda \left(\sum_{i=1}^{n/2} W_i - \sum_{i=1}^{n/2} Z_i \right) \right) \geq \exp(\lambda t \log n) \right) \\ &\leq \frac{\mathbb{E} \left[\exp \left(\lambda \left(\sum_{i=1}^{n/2} W_i - \sum_{i=1}^{n/2} Z_i \right) \right) \right]}{\exp(\lambda t \log n)} \end{aligned}$$

By independence, one has

$$\mathbb{E} \left[\exp \left(\lambda \left(\sum_{i=1}^{n/2} W_i - \sum_{i=1}^{n/2} Z_i \right) \right) \right] = \prod_{i=1}^{n/2} \mathbb{E} [\exp(\lambda W_i)] \mathbb{E} [\exp(-\lambda Z_i)]$$

Proof of Lemma 4.19 (cont.)

By definition and using $1 + x \leq e^x$, one has

$$\begin{aligned}\mathbb{E}[\exp(\lambda W_i)] &= \frac{\alpha \log n}{n} \exp(\lambda) + \left(1 - \frac{\alpha \log n}{n}\right) \\ &\leq \exp\left(\frac{\alpha \log n}{n} \exp(\lambda) - \frac{\alpha \log n}{n}\right)\end{aligned}$$

Similarly, for Z_i one has

$$\mathbb{E}[\exp(-\lambda W_i)] \leq \exp\left(\frac{\beta \log n}{n} \exp(-\lambda) - \frac{\beta \log n}{n}\right)$$

Combine these two to see that

$$\begin{aligned}\mathbb{E}[\exp(\lambda W_i)] \mathbb{E}[\exp(-\lambda Z_i)] \\ \leq \exp\left(\frac{\log n}{n} (\alpha \exp(\lambda) + \beta \exp(-\lambda) - \alpha - \beta)\right)\end{aligned}$$

Proof of Lemma 4.19 (cont.)

Combine previous two pages to see

$$\begin{aligned} \log \mathbb{P} \left(\sum_{i=1}^{n/2} W_i - \sum_{i=1}^{n/2} Z_i \leq t \log n \right) \\ \leq -\lambda t \log n + \frac{n \log n}{2} (\alpha \exp(\lambda) + \beta \exp(-\lambda) - \alpha - \beta) \end{aligned}$$

Set $\lambda = -\log(\alpha/\beta)/2$ to obtain

$$\alpha \exp(\lambda) + \beta \exp(-\lambda) - \alpha - \beta = \alpha \sqrt{\frac{\beta}{\alpha}} + \beta \sqrt{\frac{\alpha}{\beta}} - \alpha - \beta = -(\sqrt{\alpha} - \sqrt{\beta})^2$$

thus concluding the proof

$\ell_{2,\infty}$ eigen-space perturbation theory (rank- r)

Setup and algorithm

- Ground truth: $M^* = U^* \Sigma^* V^{*\top} \in \mathbb{R}^{n_1 \times n_2}$, with singular values $\sigma_1^* \geq \sigma_2^* \geq \dots \geq \sigma_r^* > 0$ (assume $n_1 \leq n_2$)
- Observation: $M = M^* + E$
- Convenient notation:

$$\kappa := \frac{\sigma_1^*}{\sigma_r^*}, \quad n := n_1 + n_2$$

- Estimate U^* (resp. V^*) using rank- r leading left (resp. right) singular subspace U (resp. V) of M

Question: can we characterize entrywise estimation error of U , i.e.

$$\text{dist}_{2,\infty}(U, U^*) := \min_{R \in \mathcal{O}^{r \times r}} \|UR - U^*\|_{2,\infty}$$

Assumptions

Noise assumptions: entries of $E = [E_{i,j}]_{1 \leq i \leq n_1, 1 \leq j \leq n_2}$ are independent obeying

$$\mathbb{E}[E_{i,j}] = 0, \quad \mathbb{E}[E_{i,j}^2] \leq \sigma^2, \quad |E_{i,j}| \leq B, \quad \text{for all } i, j$$

and it is assumed that

$$c_b := \frac{B}{\sigma \sqrt{n_1} / (\mu \log n)} = O(1)$$

Assumptions

Noise assumptions: entries of $\mathbf{E} = [E_{i,j}]_{1 \leq i \leq n_1, 1 \leq j \leq n_2}$ are independent obeying

$$\mathbb{E}[E_{i,j}] = 0, \quad \mathbb{E}[E_{i,j}^2] \leq \sigma^2, \quad |E_{i,j}| \leq B, \quad \text{for all } i, j$$

and it is assumed that

$$c_b := \frac{B}{\sigma \sqrt{n_1}/(\mu \log n)} = O(1)$$

Incoherence parameter of orthonormal matrix $\mathbf{U}^* \in \mathbb{R}^{n \times r}$ is

$$\mu(\mathbf{U}^*) := \frac{n \|\mathbf{U}^*\|_{2,\infty}^2}{r}$$

and for $\mathbf{M}^* = \mathbf{U}^* \mathbf{\Sigma}^* \mathbf{V}^{*\top}$ we define $\mu := \max\{\mu(\mathbf{U}^*), \mu(\mathbf{V}^*)\}$

$l_{2,\infty}$ distance between U and U^*

Need to take into account rotational ambiguity

— *which rotation matrix to use?*

$\ell_{2,\infty}$ distance between U and U^*

Need to take into account rotational ambiguity

— *which rotation matrix to use?*

Definition 4.20

For any square matrix Z with SVD $Z = U_Z \Sigma_Z V_Z^\top$, define

$$\text{sgn}(Z) := U_Z V_Z^\top \quad (4.22)$$

to be matrix sign function of Z (solution to Procrustes problem)

Let us employ $\text{sgn}(U^\top U^*)$ and look at

$$\|U \text{sgn}(U^\top U^*) - U^*\|_{2,\infty}$$

Theorem 4.21

With probability at least $1 - O(n^{-5})$, one has

$$\begin{aligned} \max \left\{ \|U \operatorname{sgn}(U^\top U^*) - U^*\|_{2,\infty}, \|V \operatorname{sgn}(V^\top V^*) - V^*\|_{2,\infty} \right\} \\ \lesssim \frac{\sigma \sqrt{r} (\kappa \sqrt{\frac{n_2}{n_1}} \mu + \sqrt{\log n})}{\sigma_r^*} \end{aligned}$$

provided that $\sigma \sqrt{n \log n} \leq c_1 \sigma_r^*$ for some small constant $c_1 > 0$

Entrywise matrix reconstruction error

Recall $M = \underbrace{U\Sigma V^T}_{\text{rank } r \text{ approx.}} + U_{\perp}\Sigma_{\perp}V_{\perp}^T$

Corollary 4.22

In addition, if $\sigma\kappa\sqrt{n\log n} \leq c_2\sigma_r^$ for some small enough constant $c_2 > 0$, then the following holds with probability at least $1 - O(n^{-5})$:*

$$\|U\Sigma V^T - M^*\|_{\infty} \lesssim \sigma\kappa^2\mu r \sqrt{\frac{(n_2/n_1)\log n}{n_1}}$$

De-localization of estimation error

For simplicity, let us consider the case where $\mu, \kappa, n_2/n_1 = O(1)$. Davis-Kahan theorem results in the following ℓ_2 estimation guarantees

$$\text{dist}_{\mathbb{F}}(\mathbf{U}, \mathbf{U}^*) \leq \sqrt{r} \text{dist}(\mathbf{U}, \mathbf{U}^*) \lesssim \frac{\sigma \sqrt{nr}}{\sigma_r^*}$$

In comparison, the $\ell_{2,\infty}$ bound derived in Theorem 4.21 simplifies to

$$\min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{U}\mathbf{R} - \mathbf{U}^*\|_{2,\infty} \leq \|\mathbf{U} \text{sgn}(\mathbf{U}^\top \mathbf{U}^*) - \mathbf{U}^*\|_{2,\infty} \lesssim \frac{\sigma \sqrt{r \log n}}{\sigma_r^*}$$

De-localization of estimation error (cont.)

For the matrix reconstruction error, one has

$$\|U\Sigma V^T - M^*\| \leq 2\|M - M^*\| \lesssim \sigma\sqrt{n},$$

which implies $\|U\Sigma V^T - M^*\|_F \lesssim \sigma\sqrt{nr}$

In comparison, one has

$$\|U\Sigma V^T - M^*\|_\infty \lesssim \sigma r \sqrt{\frac{\log n}{n}}$$

Concluding remarks

Concluding remarks: spectral methods

- A powerful family that permeates data science applications
 - community detection
 - matrix/tensor completion
 - ranking
 - phase retrieval
 - (joint) graph matching
 - robust PCA
 - clustering in mixture models
 - ...
- Simple yet efficient; sometimes optimal (in some weak sense)
- Commonly used to initialize nonconvex optimization algorithms
 - *see overview article Chi, Lu, Chen '19*

Concluding remarks: leave-one-out analysis

A power fine-grained analysis framework that proves effective for problems far beyond spectral methods

- robust M-estimation (e.g. El Karoui et al. '13, El Karoui '18)
- generalized power methods (e.g. Zhong, Boumal '18)
- likelihood ratio test in logistic regression (e.g. Sur et al. '19)
- nonconvex optimization (e.g. Ma et al. '20, Chen et al. '19, Cai et al. '22)
- convex relaxation (e.g. Chen et al. '19, '20, '21)
- reinforcement learning (e.g. Agarwal et al. '19, Pananjady et al. '20, Li et al. '23)
- ...

References: general

- “*Spectral methods for data science: A statistical perspective*,” Y. Chen, Y. Chi, J. Fan, C. Ma, *Foundations and Trends in Machine Learning*, 2021.
- “*Inference, estimation, and information processing, EE 378B lecture notes*,” A. Montanari, Stanford University.
- “*COMS 4772 lecture notes*,” D. Hsu, Columbia University.
- “*Spectral algorithms*,” R. Kannan, S. Vempala, *Foundations and Trends in Theoretical Computer Science*, 2009.
- “*Principal component analysis for big data*,” J. Fan, Q. Sun, W. Zhou, Z. Zhu, *arXiv:1801.01602*, 2018.
- “*Nonconvex optimization meets low-rank matrix factorization: An overview*,” Y. Chi, Y. M. Lu, Y. Chen, *IEEE Transactions on Signal Processing*, 2019.

References: matrix perturbation theory

- “*The rotation of eigenvectors by a perturbation*,” C. Davis, W. Kahan, *SIAM Journal on Numerical Analysis*, 1970.
- “*Perturbation bounds in connection with singular value decomposition*,” P. Wedin, *BIT Numerical Mathematics*, 1972.
- “*Matrix perturbation theory*,” G. W. Stewart, J. Sun, 1990.
- “*A useful variant of the Davis-Kahan theorem for statisticians*,” Y. Yu, T. Wang, R. J. Samworth, *Biometrika*, vol. 102, no. 2, pp. 315-323, 2015.
- “*Spectral method and regularized MLE are both optimal for top- K ranking*,” Y. Chen, J. Fan, C. Ma, K. Wang, *Annals of Statistics*, 2019.

References: matrix tail bounds

- “*An introduction to matrix concentration inequalities*,” J. Tropp, *Foundations and Trends in Machine Learning*, 2015.
- “*Sharp nonasymptotic bounds on the norm of random matrices with independent entries*,” A. S. Bandeira, R. van Handel, *Annals of Probability*, 2016.
- “*High-dimensional probability: An introduction with applications in data science*,” R. Vershynin, 2018.
- “*High-dimensional statistics: a non-asymptotic viewpoint*,” M. Wainwright, *Cambridge University Press*, 2019.

References: applications

- “*Community detection and stochastic block models*,” E. Abbe, *Foundations and Trends in Communications and Information Theory*, 2018.
- “*Consistency thresholds for the planted bisection model*,” E. Mossel, J. Neeman, A. Sly, *ACM Symposium on Theory of Computing*, 2015.
- “*Exact recovery in the stochastic block model*,” E. Abbe, A. S. Bandeira, G. Hall, *IEEE Transactions on information theory*, vol. 62, no. 1, pp. 471-487, 2015.
- “*Exact matrix completion via convex optimization*,” E. Candes, B. Recht, *Foundations of Computational Mathematics*, 2019.
- “*Matrix completion from a few entries*,” R. Keshavan, A. Montanari, S. Oh, *IEEE Transactions on Information Theory*, 2010.
- “*Low-rank matrix completion using alternating minimization*,” P. Jain, P. Netrapalli, S. Sanghavi, *Symposium on Theory of computing*, 2013.

References: applications

- “*Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval and matrix completion*,” C. Ma, K. Wang, Y. Chi, Y. Chen, *Foundations of Computational Mathematics*, 2020.
- “*Inference and uncertainty quantification for noisy matrix completion*,” Y. Chen, J. Fan, C. Ma, Y. Yan, *Proceedings of the National Academy of Sciences*, vol. 116, no. 46, pp. 22931-22937, 2019.
- “*Spectral algorithms for tensor completion*,” A. Montanari, N. Sun, *Communications on Pure and Applied Mathematics*, vol. 71, no. 11, pp. 2381-2425, 2018.
- “*Nonconvex low-rank tensor completion from noisy data*,” C. Cai, G. Li, H. V. Poor, Y. Chen, *Operations Research*, vol. 70, no. 2, 2022.
- “*The PageRank citation ranking: bringing order to the web*,” L. Page, S. Brin, R. Motwani, T. Winograd, 1999.

References: applications

- “*Rank centrality: ranking from pairwise comparisons*,” S. Negahban, S. Oh, D. Shah, *Operations Research*, 2017.
- “*Heteroskedastic PCA: Algorithm, optimality, and applications*,” A. Zhang, T. Cai, Y. Wu, *Annals of Statistics*, 2022.
- “*Normal approximation and confidence region of singular subspaces*,” D. Xia, *Electronic Journal of Statistics*, vol. 15, no. 2, pp. 3798-3851.
- “*Phase retrieval via Wirtinger flow: Theory and algorithms*,” E. J. Candes, X. Li, M. Soltanolkotabi, *IEEE Transactions on Information Theory*, vol. 61, no. 4, pp. 1985-2007, 2015.
- “*Solving random quadratic systems of equations is nearly as easy as solving linear systems*,” Y. Chen, E. J. Candes, *Communications on pure and applied mathematics*, vol. 70, no. 5, pp. 822-883, 2017.

References: applications

- “Rate-optimal perturbation bounds for singular subspaces with applications to high-dimensional statistics,” T. Cai, A. Zhang, *Annals of Statistics*, vol. 46, no. 1, pp. 60-89, 2018.
- “Optimality of spectral clustering in the Gaussian mixture model,” M. Loffler, A. Y. Zhang, H. H. Zhou, *Annals of Statistics*, vol. 49, no. 5, pp. 2506-2530, 2021.
- “Angular synchronization by eigenvectors and semidefinite programming,” A. Singer, *Applied and computational harmonic analysis*, vol. 30, no. 1, pp. 20-36, 2011.

References: l_∞ and $l_{2,\infty}$ theory

- “Near-optimal bounds for phase synchronization,” Y. Zhong, N. Boumal, *SIAM Journal on Optimization*, 2018.
- “Entrywise eigenvector analysis of random matrices with low expected rank,” E. Abbe, J. Fan, K. Wang, Y. Zhong, *Annals of Statistics*, 2020.
- “Spectral method and regularized MLE are both optimal for top- K ranking,” Y. Chen, J. Fan, C. Ma, K. Wang, *Annals of Statistics*, 2019.
- “The two-to-infinity norm and singular subspace geometry with applications to high-dimensional statistics,” J. Cape, M. Tang, C. E. Priebe, *Annals of Statistics*, vol. 47, no. 5, pp. 2405-2439, 2019.
- “Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval and matrix completion,” C. Ma, K. Wang, Y. Chi, Y. Chen, *Foundations of Computational Mathematics*, 2020.

References: l_∞ and $l_{2,\infty}$ theory

- “Subspace estimation from unbalanced and incomplete data matrices: $l_{2,\infty}$ statistical guarantees,” C. Cai, G. Li, Y. Chi, H. V. Poor, Y. Chen, vol. 49, no. 2, 2021.
- “Unified $l_{2 \rightarrow \infty}$ eigenspace perturbation theory for symmetric random matrices,” L. Lei, arXiv:1909.04798, 2019.
- “SIMPLE: Statistical inference on membership profiles in large networks,” J. Fan, Y. Fan, X. Han, J. Lv, *Journal of the Royal Statistical Society Series B: Statistical Methodology*, vol. 84, no. 2, pp. 630-653, 2022.
- “Entrywise Estimation of Singular Vectors of Low-Rank Matrices with Heteroskedasticity and Dependence,” J. Agterberg, Z. Lubberts, C. Priebe, *IEEE Transactions on Information Theory*, vol. 68, no. 7, pp. 4618-4650, 2022.

References: l_∞ and $l_{2,\infty}$ theory

- “*Partial recovery for top- K ranking: optimality of MLE and suboptimality of the spectral method*,” P. Chen, C. Gao, and A. Y. Zhang, *Annals of Statistics*, vol. 50, no. 3, pp. 1618-1652, 2022.
- “*Uncertainty quantification in the Bradley-Terry-Luce model*,” C. Gao, Y. Shen, A. Y. Zhang, *Information and Inference: A Journal of the IMA*, vol. 12, no. 2, pp. 1073-1140, 2023.
- “*Inference for heteroskedastic PCA with missing data*,” Y. Yan, Y. Chen, J. Fan, arXiv:2107.12365, 2021.

References: leave-one-out analysis (other methods)

- “*On robust regression with high-dimensional predictors*,” N. El Karoui, D. Bean, P. J. Bickel, C. Lim, B. Yu, *Proceedings of the National Academy of Sciences*, vol. 110, no. 36, pp. 14557-14562, 2013.
- “*On the impact of predictor geometry on the performance on high-dimensional ridge-regularized generalized robust regression estimators*,” N. El Karoui, *Probability Theory and Related Fields*, vol. 170, pp. 95-175, 2018.
- “*Near-optimal bounds for phase synchronization*,” Y. Zhong, N. Boumal, *SIAM Journal on Optimization*, 2018.
- “*The likelihood ratio test in high-dimensional logistic regression is asymptotically a rescaled chi-square*,” P. Sur, Y. Chen, E. Candes, *Probability theory and related fields*, vol. 175, pp. 487-558, 2019.

References: leave-one-out analysis (other methods)

- “*Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval and matrix completion*,” C. Ma, K. Wang, Y. Chi, Y. Chen, *Foundations of Computational Mathematics*, 2020.
- “*Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval*,” Y. Chen, Y. Chi, J. Fan, C. Ma, *Mathematical Programming*, vol. 176, pp. 5-37, 2019.
- “*Nonconvex low-rank tensor completion from noisy data*,” C. Cai, G. Li, H. V. Poor, Y. Chen, *Operations Research*, vol. 70, no. 2, 2022.
- “*Nonconvex rectangular matrix completion via gradient descent without $\ell_{2,\infty}$ regularization*,” J. Chen, D. Liu, X. Li, *IEEE Transactions on Information Theory*, vol. 66, no. 9, pp. 5806-5841, 2020.

References: leave-one-out analysis (other methods)

- “*Inference and uncertainty quantification for noisy matrix completion*,” Y. Chen, J. Fan, C. Ma, Y. Yan, *Proceedings of the National Academy of Sciences*, vol. 116, no. 46, pp. 22931-22937, 2019.
- “*Noisy matrix completion: Understanding statistical guarantees for convex relaxation via nonconvex optimization*,” Y. Chen, Y. Chi, J. Fan, C. Ma, Y. Yan, *SIAM journal on optimization*, vol. 30, no. 4, pp. 3098-3121.
- “*Bridging convex and nonconvex optimization in robust PCA: Noise, outliers and missing data*,” Y. Chen, J. Fan, C. Ma, Y. Yan, *Annals of Statistics*, vol. 49, no. 5, pp. 2948-2971, 2021.
- “*Model-based reinforcement learning with a generative model is minimax optimal*,” A. Agarwal, S. Kakade, L. Yang, *Conference on Learning Theory*, 2020.

References: leave-one-out analysis (other methods)

- “*Breaking the sample size barrier in model-based reinforcement learning with a generative model*,” G. Li, Y. Wei, Y. Chi, Y. Chen, *Operations Research*, 2023+.
- “*Instance-dependent ℓ_∞ -bounds for policy evaluation in tabular reinforcement learning*,” A. Pananjady, M. Wainwright, *IEEE Transactions on Information Theory*, vol. 67, no. 1, pp. 566-585, 2020.