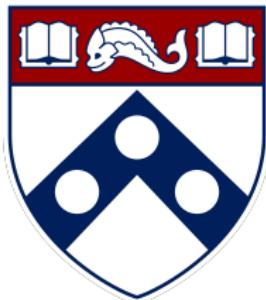


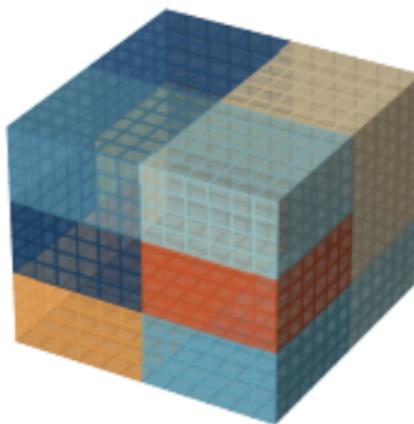
Towards more reliable tensor learning

- heteroskedastic tensor clustering and uncertainty quantification for low-rank tensors



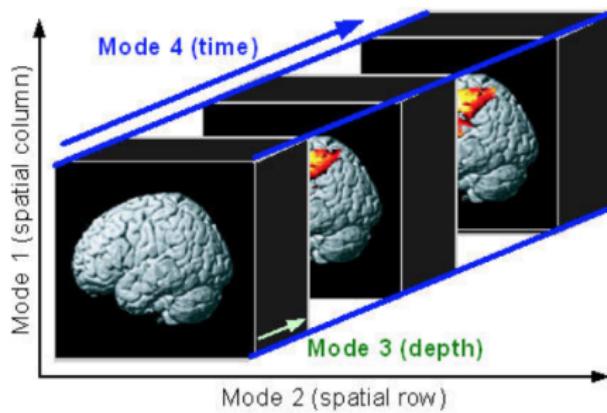
Yuchen Zhou, Wharton Statistics & Data Science

Tensors: high-order arrays



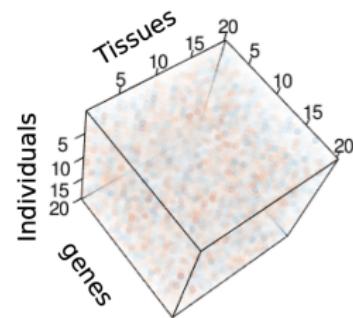
- $\mathcal{X} \in \mathbb{R}^{n_1 \times \dots \times n_d}$: order- d tensor
- vectors and matrices are order-1 and order-2 tensors

Tensors are everywhere



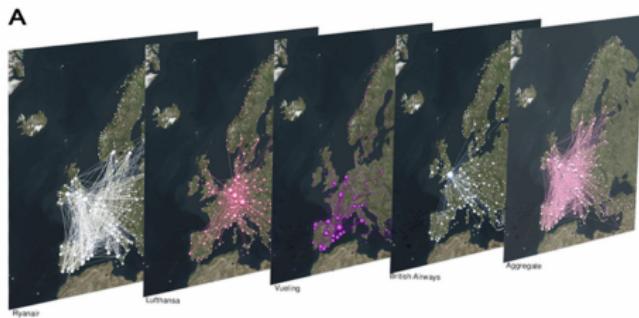
fMRI

fig. credit: Lu et al. '13



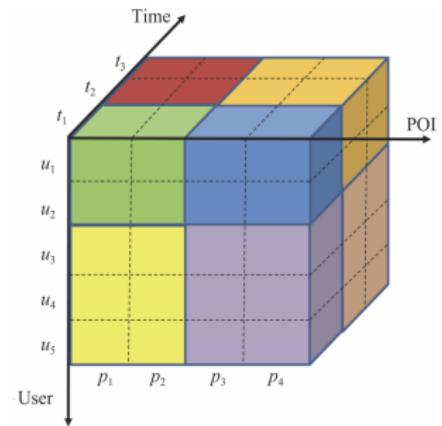
Multi-tissue gene expression
fig. credit: Wang et al. '19

Tensors are everywhere



Multilayer network analysis

fig. credit: Domenico et al. '15



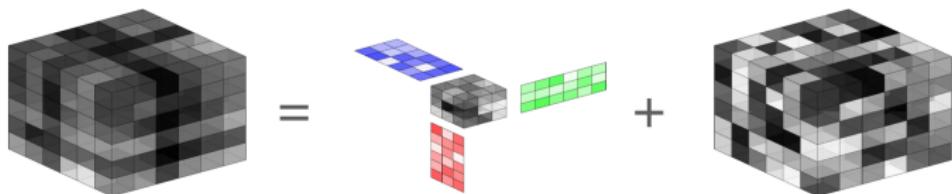
Recommender system

fig. credit: Luan et al. '17

Various tensor problems

— *Montanari and Richard '14, Zhang and Xia '18, Chen '19*

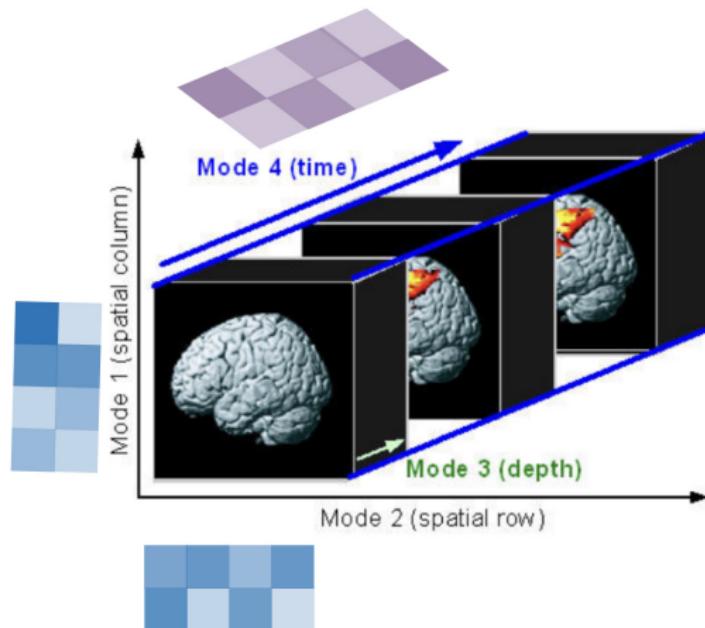
Tensor SVD/PCA



Various tensor problems

— Montanari and Richard '14, Zhang and Xia '18, Chen '19

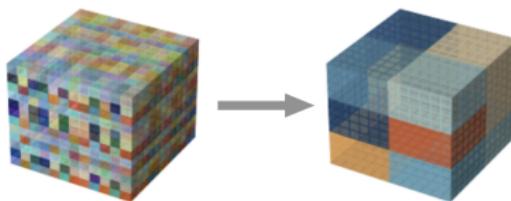
Tensor SVD/PCA



Various tensor problems

— Wang and Zeng '19, Chi et al. '20, Han et al. '22

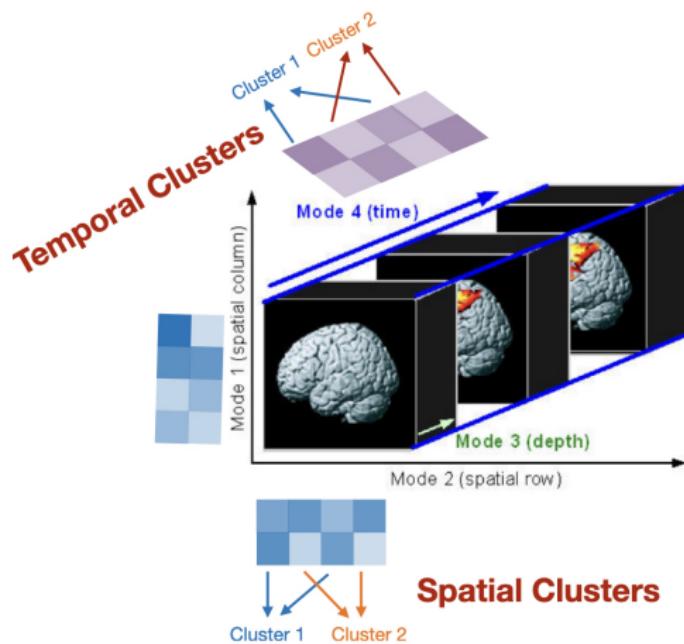
Tensor clustering



Various tensor problems

— Wang and Zeng '19, Chi et al. '20, Han et al. '22

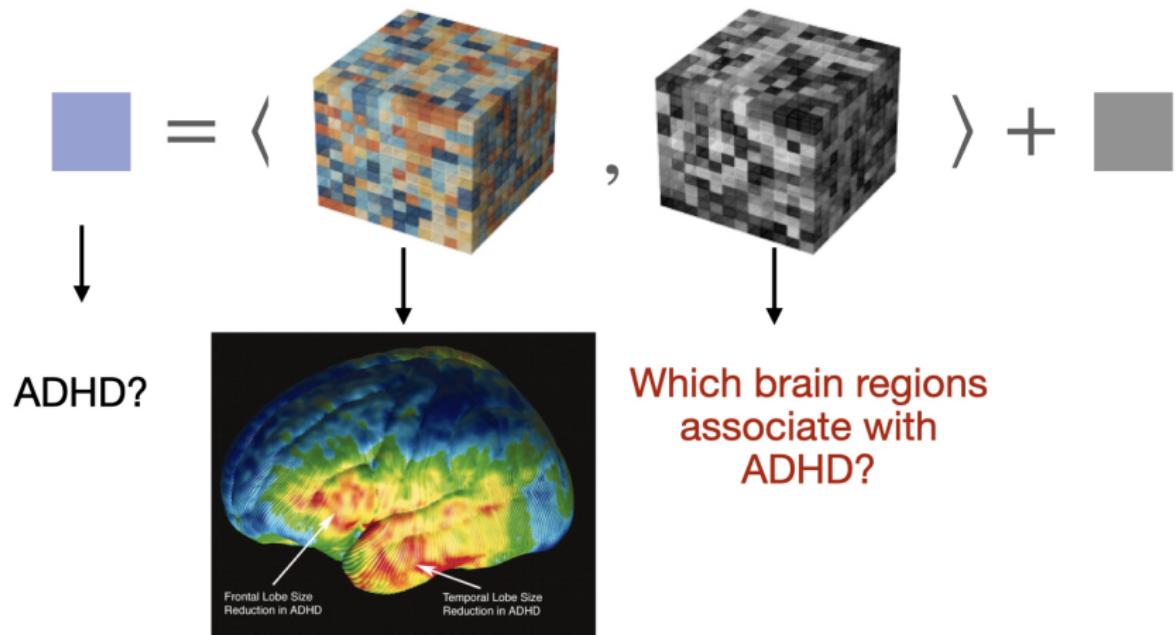
Tensor clustering



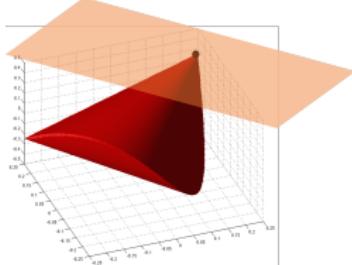
Various tensor problems

— Zhou et al. '13, Raskutti et al. '19, Zhang et al. '20

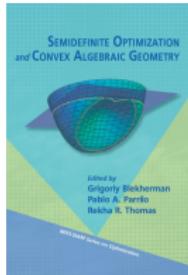
Tensor regression



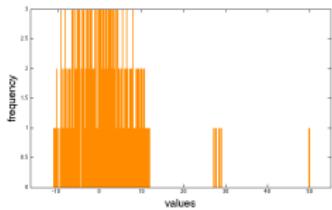
Prior arts



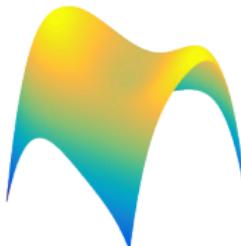
convex relaxation



sum-of-squares



spectral methods



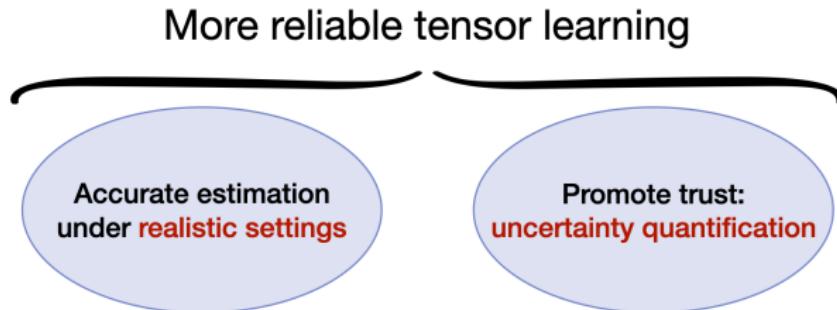
nonconvex optimization

- De Lathauwer et al. '00
- Tomioka and Suzuki '13
- Zhou et al. '13
- Anandkumaret al. '14
- Montanari and Richard '14
- Yuan and Zhang '16
- Rauhut et al. '17
- Sun and Li '17
- Li et al. '18
- Montanari and Sun '18
- Arous et al. '19
- Raskutti et al. '19
- Chi et al. '20
- Zhang et al. '20
- Cai et al. '22
- Han et al. '22
- ...

Prior arts

Challenges remain:

- ideal assumptions (e.g., i.i.d. noise) are often violated
- uncertainty quantification



Can we make tensor learning more reliable and realistic?

In this talk...

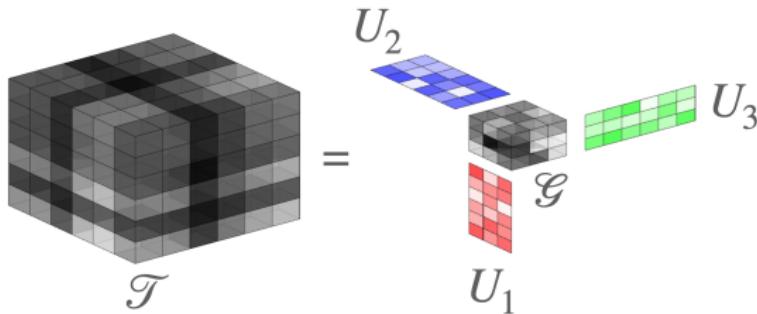
- **Heteroskedastic** tensor clustering
 - also discuss heteroskedastic PCA
- **Uncertainty quantification** for tensor learning

Tucker Low-rank Tensors

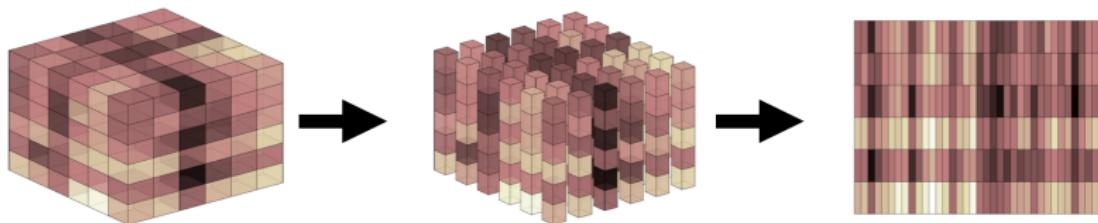
- Low-rank structures are commonly assumed!
- $\mathcal{T} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ has Tucker rank (r_1, r_2, r_3) if

$$\mathcal{T} = \mathcal{G} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \times_3 \mathbf{U}_3 = (\mathbf{U}_1, \mathbf{U}_2, \mathbf{U}_3) \cdot \mathcal{G},$$

where $\mathcal{G} \in \mathbb{R}^{r_1 \times r_2 \times r_3}$ and $\mathbf{U}_i \in \mathbb{O}_{n_i, r_i}$ for $i \in [3]$.

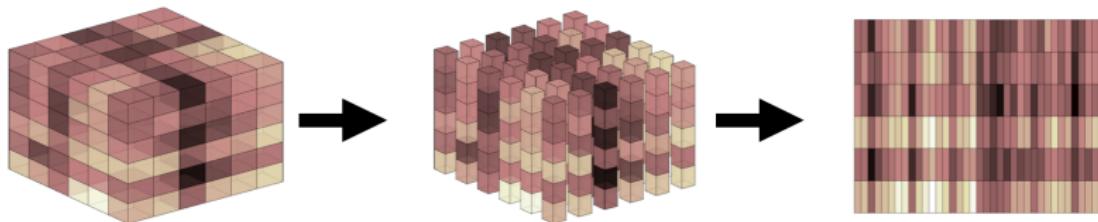


Tucker Low-rank Tensors



$\mathcal{M}_j(\mathcal{T}) \in \mathbb{R}^{n_j \times (n_1 n_2 n_3 / n_j)}$: *j*th matricization of \mathcal{T}

Tucker Low-rank Tensors



$\mathcal{M}_j(\mathcal{T}) \in \mathbb{R}^{n_j \times (n_1 n_2 n_3 / n_j)}$: *j*th matricization of \mathcal{T}

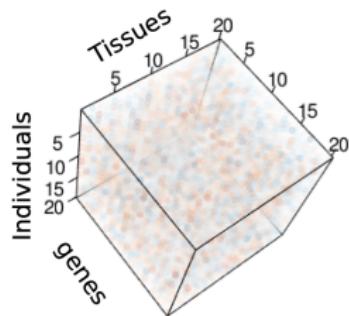
$$\text{rank}(\mathcal{T}) = (r_1, r_2, r_3) \iff \text{rank}(\mathcal{M}_j(\mathcal{T})) = r_j, \quad \forall 1 \leq j \leq 3.$$

Part 1: Heteroskedastic tensor clustering

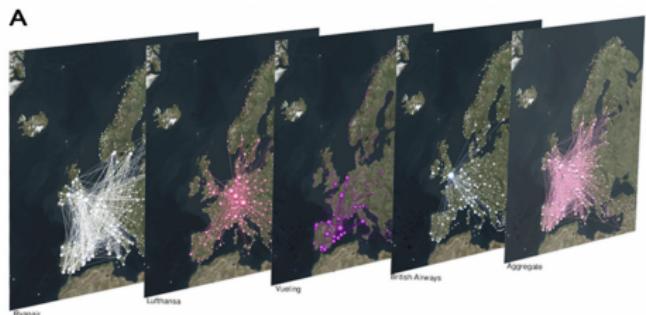


Yuxin Chen
Wharton Statistics & Data Science

Multiway clustering structure

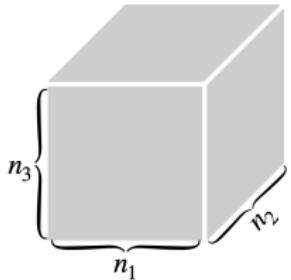


Multi-tissue gene expression
fig. credit: Wang et al. '19

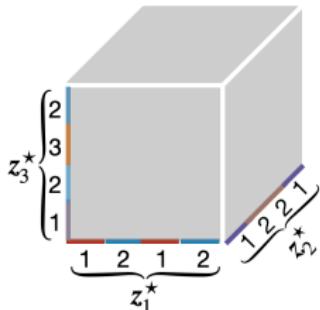


Multilayer network analysis
fig. credit: Domenico et al. '15

Tensor block model

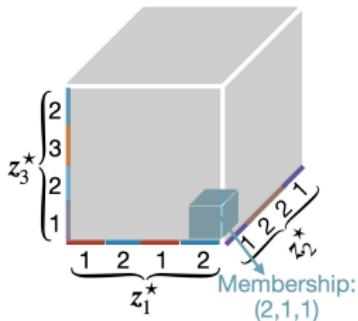


Tensor block model



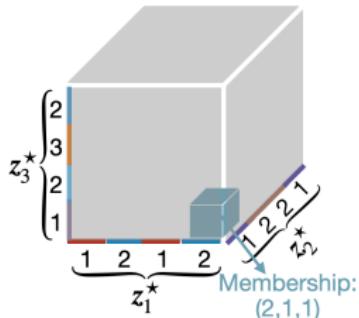
- $z_i^* \in [k_i]^{n_i}$: cluster assignment vector along the i -th mode
 - $z_{i,j}^* = \ell$ if the j th index falls within cluster ℓ

Tensor block model



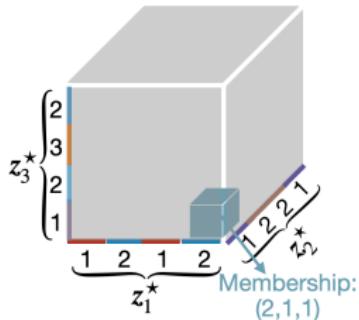
- $z_i^* \in [k_i]^{n_i}$: cluster assignment vector along the i -th mode
 - $z_{i,j}^* = \ell$ if the j th index falls within cluster ℓ
 - The membership of (i_1, i_2, i_3) is $(z_{1,i_1}^*, z_{2,i_2}^*, z_{3,i_3}^*)$

Tensor block model



- $z_i^* \in [k_i]^{n_i}$: cluster assignment vector along the i -th mode
- $\mathcal{S}^* \in \mathbb{R}^{k_1 \times k_2 \times k_3}$: block/clustering mean

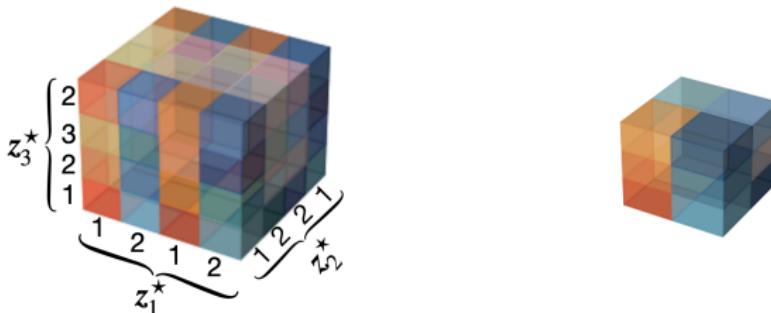
Tensor block model



- $z_i^* \in [k_i]^{n_i}$: cluster assignment vector along the i -th mode
- $\mathcal{S}^* \in \mathbb{R}^{k_1 \times k_2 \times k_3}$: block/clustering mean
- **Noisy observations:** for all $(i_1, i_2, i_3) \in [n_1] \times [n_2] \times [n_3]$,

$$Y_{i_1, i_2, i_3} = S_{z_{1,i_1}^*, z_{2,i_2}^*, z_{3,i_3}^*}^* + \underbrace{E_{i_1, i_2, i_3}}_{\text{zero-mean ind. noise}} .$$

Tensor block model

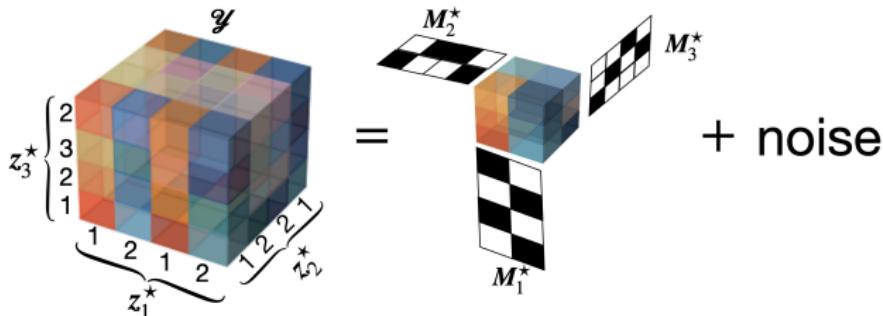


- $z_i^* \in [k_i]^{n_i}$: cluster assignment vector along the i -th mode
- $S^* \in \mathbb{R}^{k_1 \times k_2 \times k_3}$: block/clustering mean
- **Noisy observations:** for all $(i_1, i_2, i_3) \in [n_1] \times [n_2] \times [n_3]$,

$$Y_{i_1, i_2, i_3} = S^*_{z_{1, i_1}^*, z_{2, i_2}^*, z_{3, i_3}^*} + \underbrace{E_{i_1, i_2, i_3}}_{\text{zero-mean ind. noise}}.$$

- **Goal:** recover z_i^* , $i \in [3]$ from \mathcal{Y}

Tensor block model



Equivalently, we observe

$$\mathcal{Y} = \mathcal{X}^* + \underbrace{\mathcal{E}}_{\text{zero-mean ind. noise}} \in \mathbb{R}^{n_1 \times n_2 \times n_3},$$

- $\mathcal{X}^* = \mathcal{S}^* \times_1 M_1^* \times_2 M_2^* \times_3 M_3^*$, where $M_i^* \in \{0, 1\}^{n_i \times k_i}$ s.t.

$$(M_i^*)_{j,\ell} = \begin{cases} 1, & \text{if } z_{i,j}^* = \ell, \\ 0, & \text{else.} \end{cases}$$

- **Goal:** recover M_i^* , $i \in [3]$ from \mathcal{Y}

Two special cases

- **Sub-Gaussian tensor block models**
 - entries of \mathcal{E} are indep. zero-mean sub-Gaussian
 - generalization of (sub-)Gaussian mixture model

Two special cases

- **Sub-Gaussian tensor block models**

- entries of \mathcal{E} are indep. zero-mean sub-Gaussian
- generalization of (sub-)Gaussian mixture model

- **Stochastic tensor block models**

- generalization of bipartite stochastic block model
- each entry of \mathcal{S}^* is connection probability

- $$Y_{i,j,\ell} = \begin{cases} 1, & \text{with prob. } S^*_{z^*_{1,i}, z^*_{2,j}, z^*_{3,\ell}}, \text{ represents if there is a} \\ 0, & \text{o.w.} \end{cases}$$
 hyper-edge connecting (i, j, ℓ)

Two methods in the literature

Least-square estimator (Wang and Zeng '19):

$$(\hat{\mathcal{S}}, \hat{z}_1, \hat{z}_2, \hat{z}_3) := \arg \min_{\substack{\mathcal{S} \in [k_1] \times [k_2] \times [k_3] \\ z_i \in [k_i]^{n_i}}} \sum_{i,j,\ell} (Y_{i,j,\ell} - S_{z_1,i, z_2,j, z_3,\ell})^2 \quad (1)$$

- statistically accurate, **computationally intractable!**

Two methods in the literature

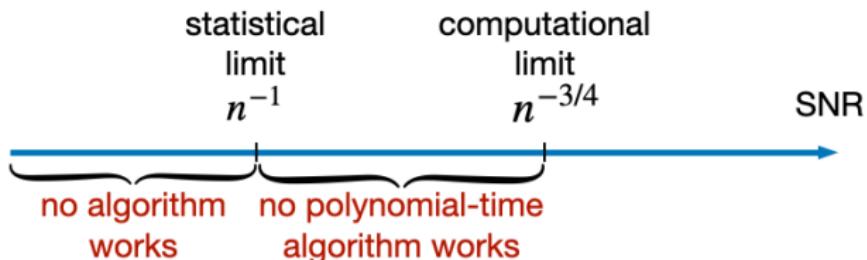
High-order Spectral Clustering + High-order Lloyd (Han et al. '22)

- estimate top singular subspaces along each mode (based on vanilla SVD)



- project tensor data onto estimated subspaces & approximate k-means procedure to cluster nodes
- High-order Lloyd (HLloyd): iterative further refinements

Statistical-computational gap



$\text{SNR} := \Delta_{\min}/\omega_{\max}$, where

- signal strength:

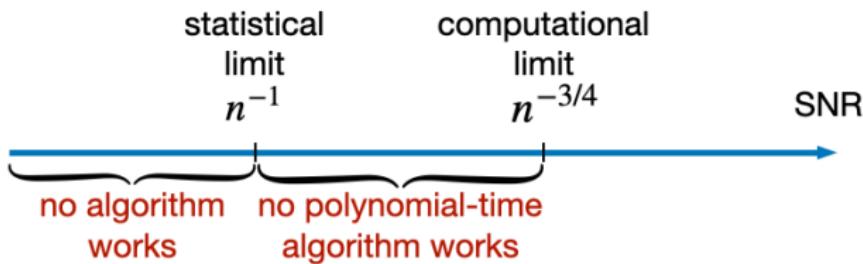
minimum separation distance along mode-1

$$\Delta_1 := \overbrace{\min_{1 \leq j_1 \neq j_2 \leq k_1} \|S_{j_1,:,:}^* - S_{j_2,:,:}^*\|_2}^{\text{minimum separation distance along mode-1}}$$

$$\Delta_{\min} := \min_{1 \leq i \leq 3} \Delta_i$$

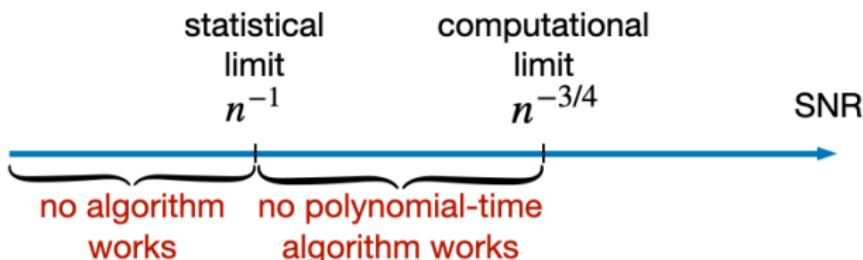
- noise level: ω_{\max}

Statistical-computational gap



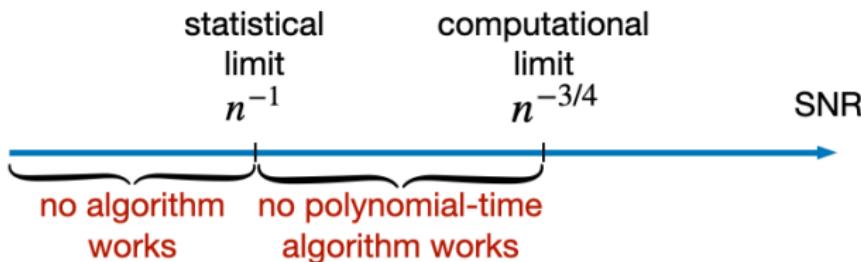
- Under i.i.d. sub-Gaussian noise, HSC + HLloyd achieves exact clustering if SNR exceeds $n^{-3/4}$ (up to log factors)

Statistical-computational gap



- Under i.i.d. sub-Gaussian noise, HSC + HLloyd achieves exact clustering if SNR exceeds $n^{-3/4}$ (up to log factors)
- A common scenario in practice: **heteroskedastic data**
 - noise variances $\{\mathbb{E}[E_{i,j,\ell}^2]\}$ are location-varying
 underlined and labeled *unknown a priori*
 - discrete-valued observations: multi-tissue gene expression data, multilayer network data

Statistical-computational gap



- Under i.i.d. sub-Gaussian noise, HSC + HLloyd achieves exact clustering if SNR exceeds $n^{-3/4}$ (up to log factors)
- A common scenario in practice: **heteroskedastic data**
 - noise variances $\underbrace{\{E[E_{i,j,\ell}^2]\}}_{\text{unknown } a \text{ priori}}$ are location-varying
 - discrete-valued observations: multi-tissue gene expression data, multilayer network data
- fail dramatically in the face of **heteroskedastic noise!**

Can we achieve exact clustering with heteroskedastic noise under essential conditions?

Spectral clustering

- Step 1: estimate “important” subspaces of $\mathbf{X}_i^* = \mathcal{M}_i(\mathcal{X}^*)$



- Step 2: apply approximate k -means to obtain cluster nodes

Spectral clustering

- Step 1: estimate “important” subspaces of $\mathbf{X}_i^* = \mathcal{M}_i(\mathcal{X}^*)$



- key challenge: how to handle heteroskedastic noise?
- Step 2: apply approximate k -means to obtain cluster nodes

Road map

Subspace estimation problem

.Key challenges: heteroskedastic noise, ...

.Propose Deflated HeteroPCA! 😎

More reliable solution to PCA problem!



Spectral tensor clustering

.Still need assumptions on least singular value

.Thresholded Deflated HeteroPCA + k -means!

A detour: a subspace estimation / model

$$Y = U^\star \Sigma^\star V^{\star\top} + E$$

The diagram shows a large green matrix \$Y\$ on the left, which is decomposed into four smaller matrices: \$U^\star\$ (blue), \$\Sigma^\star\$ (green), \$V^{\star\top}\$ (purple), and \$E\$ (grey). The dimensions of \$Y\$ are indicated by curly braces: \$n_1\$ vertically and \$n_2\$ horizontally.

- **Ground truth:** rank- r matrix X^* with SVD ($r \ll \min\{n_1, n_2\}$)

$$\mathbf{X}^{\star} = \mathbf{U}^{\star}\boldsymbol{\Sigma}^{\star}\mathbf{V}^{\star\top} = \sum_{i=1}^r \sigma_i^{\star} \mathbf{u}_i^{\star} \mathbf{v}_i^{\star\top} \in \mathbb{R}^{n_1 \times n_2}$$

where $U^* \in \mathbb{R}^{n_1 \times r}$, $\Sigma^* = \text{diag}\{\sigma_1^*, \dots, \sigma_r^*\}$, $V^* \in \mathbb{R}^{n_2 \times r}$

A detour: a subspace estimation / model

$$n_1 \underbrace{\begin{matrix} & n_2 \\ \overbrace{\hspace{1cm}}^{} & \end{matrix}}_{Y} = U^* \Sigma^* V^{*\top} + E$$

- **Ground truth:** rank- r matrix X^* with SVD ($r \ll \min\{n_1, n_2\}$)

$$X^* = U^* \Sigma^* V^{*\top} = \sum_{i=1}^r \sigma_i^* u_i^* v_i^{*\top} \in \mathbb{R}^{n_1 \times n_2}$$

- **Noisy observations:** $Y = X^* + \underbrace{E}_{\text{zero-mean ind. noise}}$

A detour: a subspace estimation / model

$$Y = U^* \Sigma^* V^{*\top} + E$$

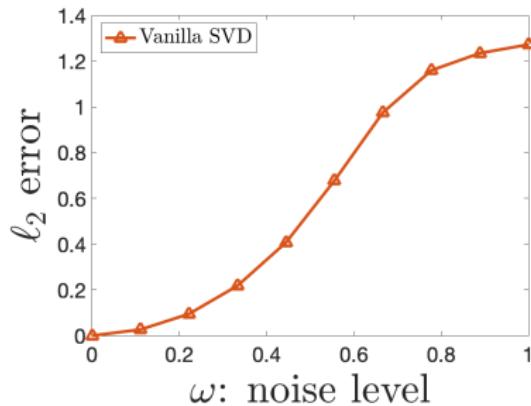
- **Ground truth:** rank- r matrix X^* with SVD ($r \ll \min\{n_1, n_2\}$)

$$\mathbf{X}^{\star} = \mathbf{U}^{\star}\boldsymbol{\Sigma}^{\star}\mathbf{V}^{\star\top} = \sum_{i=1}^r \sigma_i^{\star} \mathbf{u}_i^{\star} \mathbf{v}_i^{\star\top} \in \mathbb{R}^{n_1 \times n_2}$$

- **Noisy observations:** $Y = X^* + \underbrace{E}_{\text{zero-mean ind. noise}}$
 - **Goal:** estimate column subspace $U^* \in \mathbb{R}^{n_1 \times r}$ based on Y

Review of popular methods

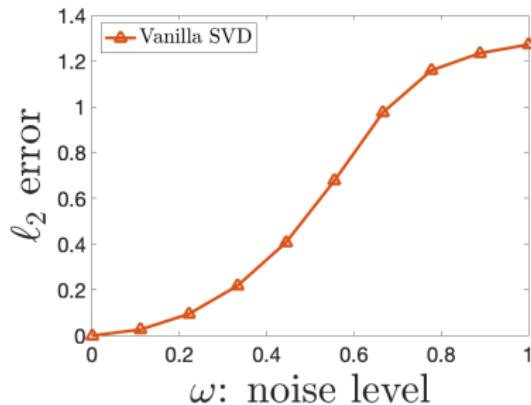
$n_1 = 100, n_2 = 10,000$
 $r = 2, \kappa := \sigma_1^*/\sigma_r^* = 2$
 $\omega_i \stackrel{\text{ind.}}{\sim} \text{Unif}([0, \omega])$,
 $E_{i,j} \stackrel{\text{ind.}}{\sim} \mathcal{N}(0, \omega_i^2)$



vanilla SVD: $U \leftarrow$ rank- r left singular subspace of $\mathbf{Y} = \mathbf{X}^* + \mathbf{E}$

Review of popular methods

$$\begin{aligned}n_1 &= 100, n_2 = 10,000 \\r &= 2, \kappa := \sigma_1^*/\sigma_r^* = 2 \\ \omega_i &\stackrel{\text{ind.}}{\sim} \text{Unif}([0, \omega]), \\ E_{i,j} &\stackrel{\text{ind.}}{\sim} \mathcal{N}(0, \omega_i^2)\end{aligned}$$



vanilla SVD: $\mathbf{U} \leftarrow$ rank- r left singular subspace of $\mathbf{Y} = \mathbf{X}^* + \mathbf{E}$

- often sub-optimal due to large bias in diagonal entries:

$$\mathbb{E}[\mathbf{Y}\mathbf{Y}^\top] = \underbrace{\mathbf{X}^*\mathbf{X}^{*\top}}_{\checkmark} + \underbrace{\text{diag}\left\{\left[\sum_j \mathbb{E}[E_{i,j}^{*2}]\right]_{1 \leq i \leq n_1}\right\}}_{\text{potentially large diagonal matrix!}}$$

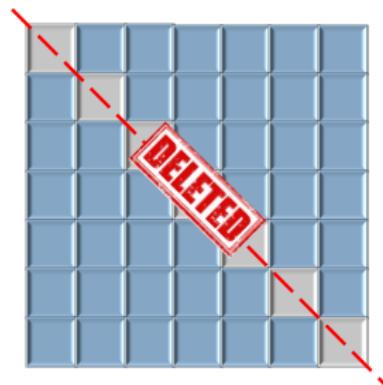
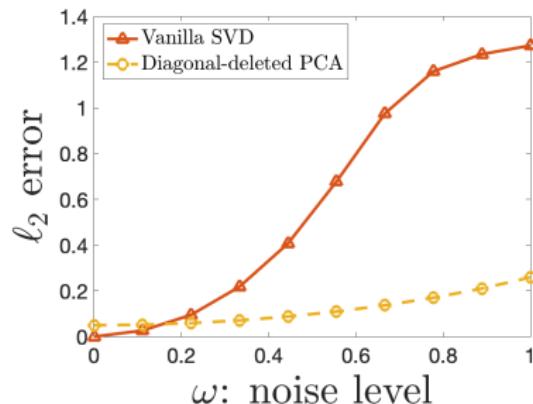
Review of popular methods

$$n_1 = 100, n_2 = 10,000$$

$$r = 2, \kappa := \sigma_1^*/\sigma_r^* = 2$$

$$\omega_i \stackrel{\text{ind.}}{\sim} \text{Unif}([0, \omega]),$$

$$E_{i,j} \stackrel{\text{ind.}}{\sim} \mathcal{N}(0, \omega_i^2)$$

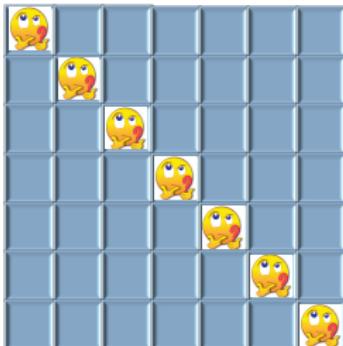
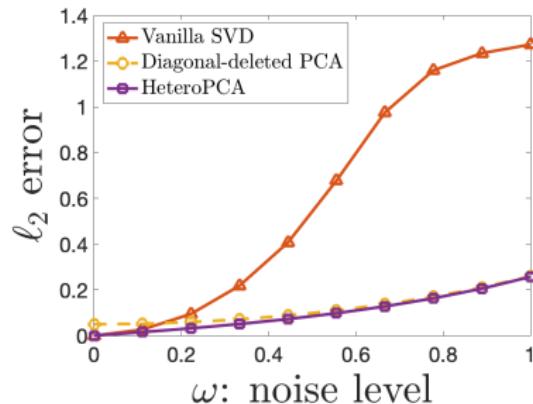


diagonal-deleted PCA:

- remove $\text{diag}(YY^\top)$
- compute top- r eigen-space

Review of popular methods

$$\begin{aligned}n_1 &= 100, n_2 = 10,000 \\r &= 2, \kappa := \sigma_1^*/\sigma_r^* = 2 \\ \omega_i &\stackrel{\text{ind.}}{\sim} \text{Unif}([0, \omega]), \\ E_{i,j} &\stackrel{\text{ind.}}{\sim} \mathcal{N}(0, \omega_i^2)\end{aligned}$$

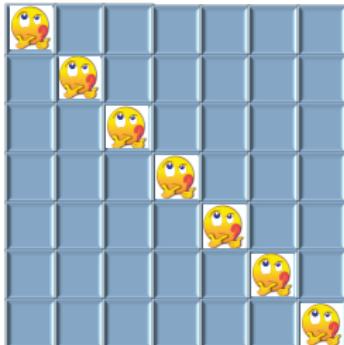
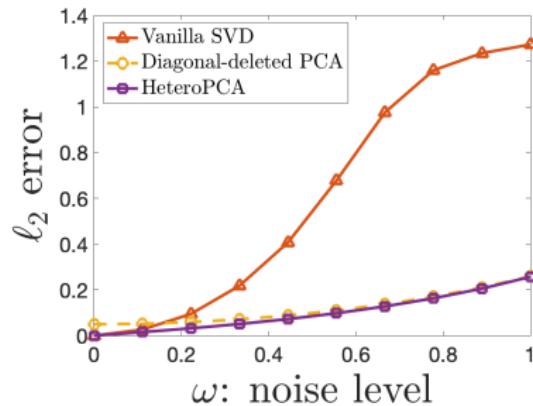


HeteroPCA (Zhang, Cai, Wu '22)

- iteratively estimate $\text{diag}(\mathbf{X}\mathbf{X}^\top)$
- compute top- r eigen-space

Review of popular methods

$$\begin{aligned}n_1 &= 100, n_2 = 10,000 \\r &= 2, \kappa := \sigma_1^*/\sigma_r^* = 2 \\ \omega_i &\stackrel{\text{ind.}}{\sim} \text{Unif}([0, \omega]), \\ E_{i,j} &\stackrel{\text{ind.}}{\sim} \mathcal{N}(0, \omega_i^2)\end{aligned}$$

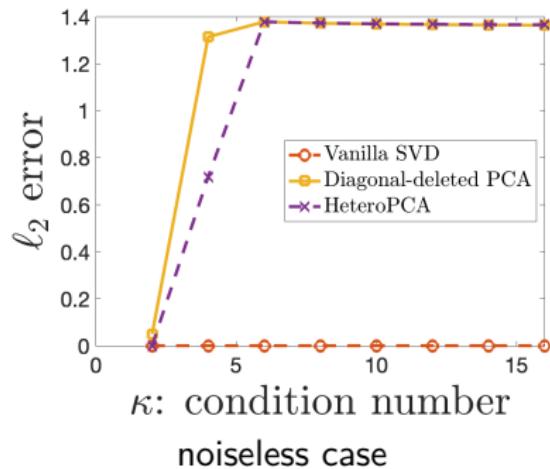


HeteroPCA (Zhang, Cai, Wu '22)

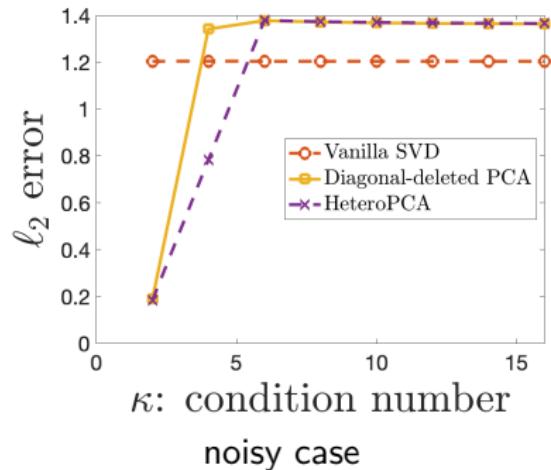
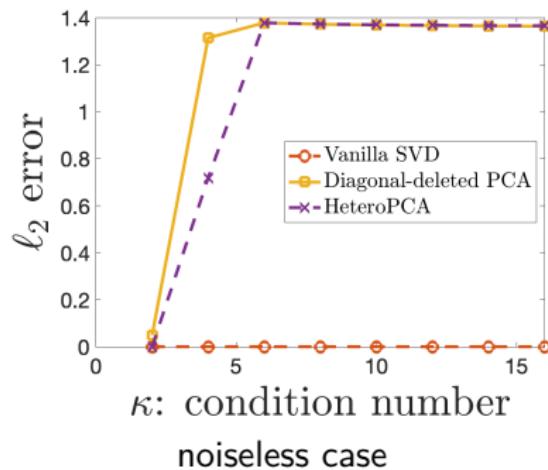
- **initialize:** $\mathbf{G}^0 = \mathcal{P}_{\text{off-diag}}(\mathbf{Y}\mathbf{Y}^\top)$
- for $t = 0, 1, \dots$
 $(\mathbf{U}^t, \Lambda^t) = \text{eigs}(\mathbf{G}^t, r)$
 $\mathbf{G}^{t+1} = \mathbf{G}^0 + \mathcal{P}_{\text{diag}}(\mathbf{U}^t \Lambda^t \mathbf{U}^{t\top})$

A curious phenomenon: curse of ill-conditioning

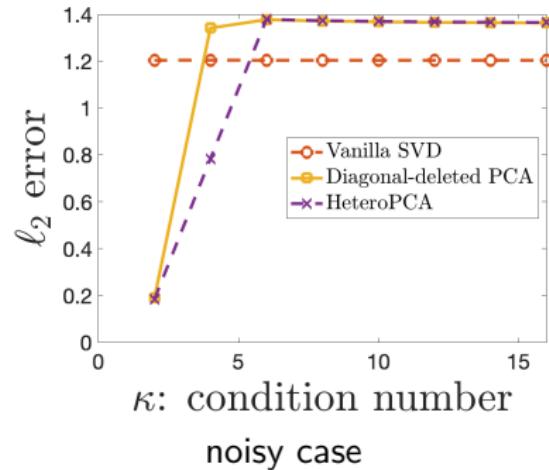
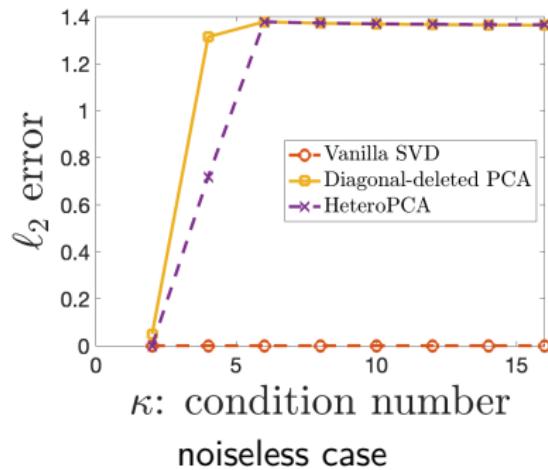
Somewhat surprising numerical example: $r = 2, n_1 = 200, n_2 = 40,000$



Somewhat surprising numerical example: $r = 2, n_1 = 200, n_2 = 40,000$



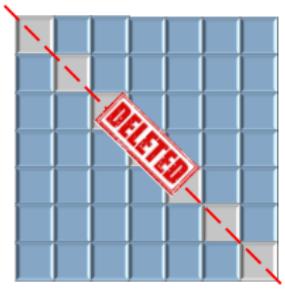
Somewhat surprising numerical example: $r = 2, n_1 = 200, n_2 = 40,000$



Previous methods degrade as condition number of X^* increases!

but this actually makes problem info-theoretically easier ...

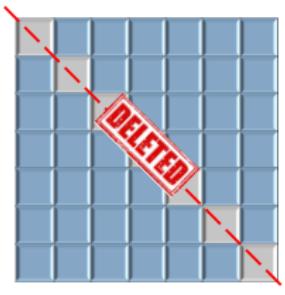
Diagonals: influences of diagonal deletion



$$\mathbb{E} \left[\mathcal{P}_{\text{off-diag}} (\mathbf{Y} \mathbf{Y}^\top) \right] = \mathbf{X}^* \mathbf{X}^{*\top} - \underbrace{\mathcal{P}_{\text{diag}} (\mathbf{X}^* \mathbf{X}^{*\top})}_{\text{ideally negligible compared to } \sigma_r^{*2}}$$

- ideally, we hope diagonal deletion has negligible influences

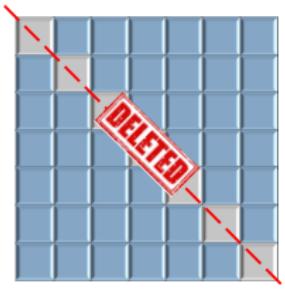
Diagonals: influences of diagonal deletion



$$\mathbb{E} \left[\mathcal{P}_{\text{off-diag}}(\mathbf{Y}\mathbf{Y}^\top) \right] = \mathbf{X}^* \mathbf{X}^{*\top} - \underbrace{\mathcal{P}_{\text{diag}}(\mathbf{X}^* \mathbf{X}^{*\top})}_{\text{but larger than } \sigma_r^{*2} \text{ if } \sigma_1^*/\sigma_r^* \text{ is too large}}$$

- ideally, we hope diagonal deletion has negligible influences
- non-negligible for ill-conditioned case though ...

Diagonals: influences of diagonal deletion



$$\mathbb{E}[\mathcal{P}_{\text{off-diag}}(\mathbf{Y}\mathbf{Y}^\top)] = \mathbf{X}^* \mathbf{X}^{*\top} - \underbrace{\mathcal{P}_{\text{diag}}(\mathbf{X}^* \mathbf{X}^{*\top})}_{\text{but larger than } \sigma_r^{*2} \text{ if } \sigma_1^*/\sigma_r^* \text{ is too large}}$$

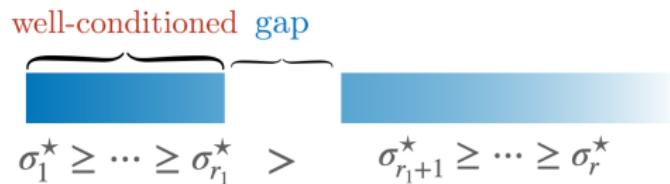
- ideally, we hope diagonal deletion has negligible influences
- non-negligible for ill-conditioned case though ...

Both diagonal-deleted PCA & HeteroPCA become ineffective
initialized by diagonal-deleted PCA
in the presence of ill-conditioning!

*Can we break the curse of ill-conditioning while
accommodating widest SNR range?*

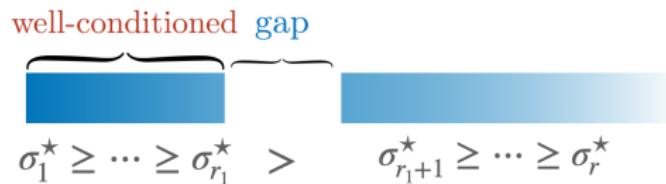
revisit HeteroPCA theory: works well if

- X^* is well-conditioned
- least singular value σ_r^* (or spectral gap) is not buried by noise



revisit HeteroPCA theory: works well if

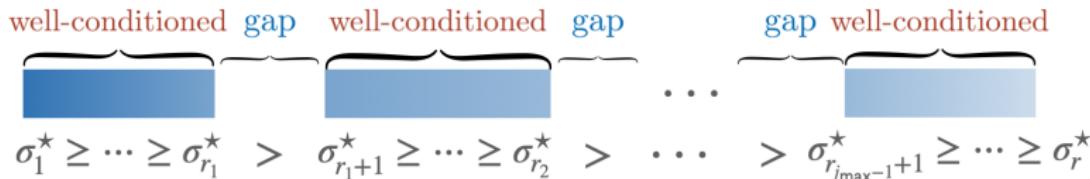
- X^* is well-conditioned
- least singular value σ_r^* (or spectral gap) is not buried by noise



solution:

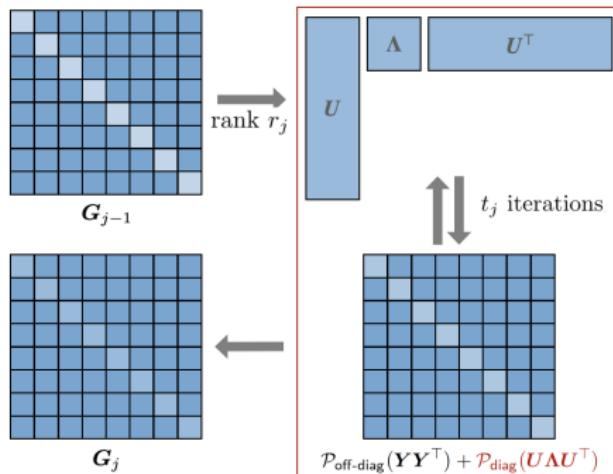
- divide eigenvalues into well-conditioned & well-separated subblocks
- estimate subblocks sequentially

Proposed algorithm: deflated-HeteroPCA



- sequentially choose ranks $r_0 = 0 < r_1 < \dots < r_{j_{\max}} = r$ s.t.
 - $\sigma_{r_{j-1}+1}^* / \sigma_{r_j}^*$ is small
 - sufficient gap between $\sigma_{r_j}^*$ and $\sigma_{r_j+1}^*$

Proposed algorithm: deflated-HeteroPCA



- sequentially choose ranks $r_0 = 0 < r_1 < \dots < r_{j_{\max}} = r$ s.t.
 - $\sigma_{r_{j-1}+1}^*/\sigma_{r_j}^*$ is small
 - sufficient gap between $\sigma_{r_j}^*$ and $\sigma_{r_j+1}^*$
- invoke HeteroPCA($\underbrace{G_{j-1}}_{\text{input}}, \underbrace{r_j}_{\text{rank}}$) to impute diagonals & obtain \mathbf{G}_k

Proposed algorithm: deflated-HeteroPCA

- **Initialize:** $\mathbf{G}^0 = \mathcal{P}_{\text{off-diag}}(\mathbf{Y}\mathbf{Y}^\top)$, $k = 0$, $r_0 = 0$
- **Sequential updates:** while $r_k < r$

$$j = j + 1$$

select r_j in a data-driven manner

$$(\mathbf{G}_j, \mathbf{U}_j) = \text{HeteroPCA}(\underbrace{\mathbf{G}_{j-1}}_{\text{input}}, \underbrace{r_j}_{\text{rank}})$$

- **Output:** $\mathbf{U} := \mathbf{U}_j \longrightarrow$ estimate of \mathbf{U}^*

Proposed algorithm: deflated-HeteroPCA

- **Initialize:** $\mathbf{G}^0 = \mathcal{P}_{\text{off-diag}}(\mathbf{Y}\mathbf{Y}^\top)$, $k = 0$, $r_0 = 0$
- **Sequential updates:** while $r_k < r$

$$j = j + 1$$

select r_j in a data-driven manner

$$(\mathbf{G}_j, \mathbf{U}_j) = \text{HeteroPCA}(\underbrace{\mathbf{G}_{j-1}}_{\text{input}}, \underbrace{r_j}_{\text{rank}})$$

- **Output:** $\mathbf{U} := \mathbf{U}_j \longrightarrow$ estimate of \mathbf{U}^*

$$\text{Select } r_j = \begin{cases} \max \mathcal{R}_j, & \text{if } \mathcal{R}_j \neq \emptyset, \\ r, & \text{otherwise.} \end{cases} \text{ Here,}$$

$$\mathcal{R}_j := \{r' : r_{j-1} < r' \leq r, \underbrace{\sigma_{r_{j-1}+1}(\mathbf{G}_{j-1}) / \sigma_{r'}(\mathbf{G}_{j-1})}_{\text{well-conditioned}} \leq 4 \& \underbrace{\sigma_{r'}(\mathbf{G}_{j-1}) - \sigma_{r'+1}(\mathbf{G}_{j-1})}_{\text{gap}} \geq \sigma_{r'}(\mathbf{G}_{j-1}) / r\}.$$

Assumptions (ignoring log factors)

- **heteroskedasticity:** $E'_{i,j} s$ are indep. obeying
 - $\mathbb{E}[E_{i,j}] = 0, \quad \text{Var}[E_{i,j}] \leq \omega^2$
 - $|E_{i,j}| \lesssim \omega \min\left\{\left(n_1 n_2\right)^{1/4}, \sqrt{n_2}\right\}$ with high prob.
 - examples: sub-Gaussian, centered Poisson, centered Bernoulli

Assumptions (ignoring log factors)

- **heteroskedasticity:** $E'_{i,j} s$ are indep. obeying
 - $\mathbb{E}[E_{i,j}] = 0, \quad \text{Var}[E_{i,j}] \leq \omega^2$
 - $|E_{i,j}| \lesssim \omega \min\left\{(n_1 n_2)^{1/4}, \sqrt{n_2}\right\}$ with high prob.
 - examples: sub-Gaussian, centered Poisson, centered Bernoulli
- **signal-to-noise ratio (SNR):**

$$\frac{\sigma_r^*}{\omega} \gtrsim (n_1 n_2)^{1/4} + n_1^{1/2}$$

- necessary for existence of consistent estimators (Cai et al. '21)

Assumptions (ignoring log factors)

- **heteroskedasticity:** $E'_{i,j} s$ are indep. obeying
 - $\mathbb{E}[E_{i,j}] = 0, \quad \text{Var}[E_{i,j}] \leq \omega^2$
 - $|E_{i,j}| \lesssim \omega \min\left\{(n_1 n_2)^{1/4}, \sqrt{n_2}\right\}$ with high prob.
 - examples: sub-Gaussian, centered Poisson, centered Bernoulli
- **signal-to-noise ratio (SNR):**

$$\frac{\sigma_r^*}{\omega} \gtrsim (n_1 n_2)^{1/4} + n_1^{1/2}$$

- necessary for existence of consistent estimators (Cai et al. '21)
- rank $r = O(1)$
- incoherence $\mu := \max\left\{\frac{n_1}{r} \|\boldsymbol{U}^*\|_{2,\infty}^2, \frac{n_2}{r} \|\boldsymbol{V}^*\|_{2,\infty}^2\right\} = O(1)$

Theoretical guarantees

Theorem 1 (Zhou, Chen '23)

With high prob., Deflated-HeteroPCA yields

$$\|\mathbf{U} \mathbf{R}_{\mathbf{U}} - \mathbf{U}^*\| \lesssim \zeta_{\text{op}}$$

for some rotation matrix $\mathbf{R}_{\mathbf{U}}$, where $\zeta_{\text{op}} = \frac{\sqrt{n_1 n_2} \omega^2}{\sigma_r^{*2}} + \frac{\sqrt{n_1} \omega}{\sigma_r^*}$

- match minimax lower bounds in Zhang et al. '22 & Cai et al. '21
- condition-number-free

Theoretical guarantees

Theorem 1 (Zhou, Chen '23)

With high prob., Deflated-HeteroPCA yields

$$\|\mathbf{U}\mathbf{R}_U - \mathbf{U}^*\| \lesssim \zeta_{\text{op}}$$

$$\|\mathbf{U}\mathbf{R}_U - \mathbf{U}^*\|_{2,\infty} \lesssim \frac{1}{\sqrt{n_1}} \zeta_{\text{op}} \quad (\text{fine-grained } \ell_{2,\infty})$$

for some rotation matrix \mathbf{R}_U , where $\zeta_{\text{op}} = \frac{\sqrt{n_1 n_2} \omega^2}{\sigma_r^{*2}} + \frac{\sqrt{n_1} \omega}{\sigma_r^*}$

- match minimax lower bounds in Zhang et al. '22 & Cai et al. '21
- condition-number-free

Come back to tensor clustering...

Subspace estimation problem

.Key challenges: heteroskedastic noise,
not well-conditioning

.Propose Deflated HeteroPCA! 😎

More reliable solution to PCA problem!



Spectral tensor clustering

- .Still need assumptions on least singular value
- .Thresholded Deflated HeteroPCA + k -means!

Come back to tensor clustering...

Subspace estimation problem

.Key challenges: heteroskedastic noise,
not well-conditioning

.Propose Deflated HeteroPCA! 😎

More reliable solution to PCA problem!

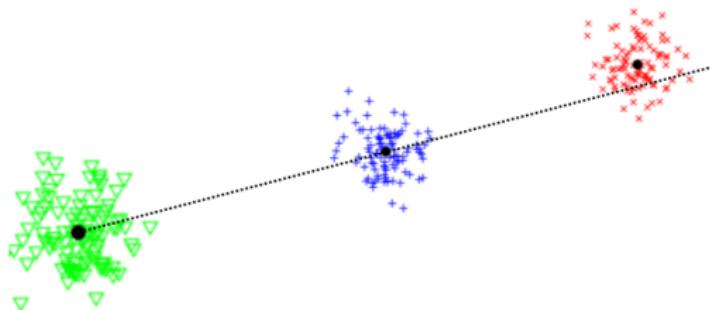


Spectral tensor clustering

- .Still need assumptions on least singular value
- .Thresholded Deflated HeteroPCA + k -means!

Come back to tensor clustering...

- Least singular value assumption is **unnecessary** for clustering!



- Add a thresholding procedure!

Thresholded Deflated-HeteroPCA(\mathbf{Y}, r, τ)

- **Initialize:** $\mathbf{G}^0 = \mathcal{P}_{\text{off-diag}}(\mathbf{Y}\mathbf{Y}^\top), j = 0, r_0 = 0$
- **Sequential updates:** while $r_j < r$ and $\sigma_{r_j+1}(\mathbf{G}_j) > \underbrace{\tau}_{\text{threshold}}$
 - $j = j + 1$
 - select r_j in a data-driven manner
 - $(\mathbf{G}_j, \mathbf{U}_j) = \text{HeteroPCA}(\underbrace{\mathbf{G}_{j-1}}_{\text{input}}, \underbrace{r_j}_{\text{rank}})$
- **Output:** $\mathbf{U} := \mathbf{U}_j \longrightarrow \text{estimate of } \mathbf{U}^*$

Thresholded Deflated-HeteroPCA(\mathbf{Y}, r, τ)

- **Initialize:** $\mathbf{G}^0 = \mathcal{P}_{\text{off-diag}}(\mathbf{Y}\mathbf{Y}^\top), j = 0, r_0 = 0$
- **Sequential updates:** while $r_j < r$ and $\sigma_{r_j+1}(\mathbf{G}_j) > \underbrace{\tau}_{\text{threshold}}$
 - $j = j + 1$
 - select r_j in a data-driven manner
 - $(\mathbf{G}_j, \mathbf{U}_j) = \text{HeteroPCA}(\underbrace{\mathbf{G}_{j-1}}_{\text{input}}, \underbrace{r_j}_{\text{rank}})$
- **Output:** $\mathbf{U} := \mathbf{U}_j \longrightarrow \text{estimate of } \mathbf{U}^*$

Select $r_j = \begin{cases} \max \mathcal{R}_j, & \text{if } \mathcal{R}_j \neq \emptyset, \\ r, & \text{otherwise.} \end{cases}$ Here,

$$\mathcal{R}_j := \{r' : r_{j-1} < r' \leq r, \underbrace{\sigma_{r_{j-1}+1}(\mathbf{G}_{j-1}) / \sigma_{r'}(\mathbf{G}_{j-1}) \leq 4}_{\text{well-conditioned}} \& \underbrace{\sigma_{r'}(\mathbf{G}_{j-1}) - \sigma_{r'+1}(\mathbf{G}_{j-1}) \geq \sigma_{r'}(\mathbf{G}_{j-1}) / r}_{\text{gap}}\}.$$

Proposed algorithm: High-order HeteroClustering

- Step 1: estimate “important” subspaces of $\mathbf{X}_i^* = \mathcal{M}_i(\mathcal{X}^*)$



- Step 1: estimate “important” subspaces of $\mathbf{X}_i^* = \mathcal{M}_i(\mathcal{X}^*)$
 - Apply Thresholded Deflated-HeteroPCA($\mathbf{X}_i^*, k_i, \tau$) to obtain subspace estimate $\mathbf{U}_i \in \mathbb{R}^{n_i \times r_i}$
- Step 2: apply approximate k -means on the rows of
$$\widehat{\mathbf{B}}_i = \mathcal{M}_i(\underbrace{\mathcal{Y} \times_i \mathbf{U}_i \mathbf{U}_i^\top}_{\text{denoising}} \quad \underbrace{\times_{i+1} \mathbf{U}_{i+1} \times_{i+2} \mathbf{U}_{i+2}}_{\text{dimension reduction & denoising}} \quad) \in \mathbb{R}^{n_i \times r_{i+1} r_{i+2}}$$
 - can be efficiently done by using k -means++!

Assumptions (ignoring log factors)

- dimension $n_1 \asymp n_2 \asymp n_3 \asymp n$
- **heteroskedasticity:** \mathcal{E} has indep. zero-mean entries obeying
 - $\text{Var}[E_{i,j,\ell}] \leq \omega_{\max}^2$
 - $|E_{i,j,\ell}| \lesssim \omega_{\max} n^{3/4}$ with high prob.
 - **examples:** sub-Gaussian, centered Poisson, stochastic tensor block model

Assumptions (ignoring log factors)

- dimension $n_1 \asymp n_2 \asymp n_3 \asymp n$
- **heteroskedasticity:** \mathcal{E} has indep. zero-mean entries obeying
 - $\text{Var}[E_{i,j,\ell}] \leq \omega_{\max}^2$
 - $|E_{i,j,\ell}| \lesssim \omega_{\max} n^{3/4}$ with high prob.
 - **examples:** sub-Gaussian, centered Poisson, stochastic tensor block model
- **signal-to-noise ratio (SNR):**
$$\frac{\Delta_{\min}}{\omega_{\max}} \gtrsim n^{-3/4}$$
 - match the computational limit (Han et al. '22)

Assumptions (ignoring log factors)

- dimension $n_1 \asymp n_2 \asymp n_3 \asymp n$
- **heteroskedasticity:** \mathcal{E} has indep. zero-mean entries obeying
 - $\text{Var}[E_{i,j,\ell}] \leq \omega_{\max}^2$
 - $|E_{i,j,\ell}| \lesssim \omega_{\max} n^{3/4}$ with high prob.
 - **examples:** sub-Gaussian, centered Poisson, stochastic tensor block model
- **signal-to-noise ratio (SNR):**
$$\frac{\Delta_{\min}}{\omega_{\max}} \gtrsim n^{-3/4}$$
 - match the computational limit (Han et al. '22)
- balanced cluster sizes: $|\{j \in [n_i] : (\mathbf{z}_i^\star)_j = \ell\}| \asymp n_i/k_i$
- number of clusters $k_i = O(1)$

Theoretical guarantees

Theorem 2 (Zhou, Chen '23)

Suppose that $\tau \asymp n^{3/4} \omega_{\max}^2$. Then with high prob., High-order HeteroClustering (HHC) yields

$$\text{MCR}(\hat{z}_i, z_i^*) = 0, \quad \forall 1 \leq i \leq 3.$$

Here, $\text{MCR}(z, \hat{z}) := \inf_{\text{permutation } \phi: [k] \rightarrow [k]} \frac{1}{n} \sum_{j=1}^n \mathbb{1}\{\hat{z}_j \neq \phi(z_j)\}.$

Theoretical guarantees

Theorem 2 (Zhou, Chen '23)

Suppose that $\tau \asymp n^{3/4} \omega_{\max}^2$. Then with high prob., High-order HeteroClustering (HHC) yields

$$\text{MCR}(\hat{z}_i, z_i^*) = 0, \quad \forall 1 \leq i \leq 3.$$

Here, $\text{MCR}(z, \hat{z}) := \inf_{\text{permutation } \phi: [k] \rightarrow [k]} \frac{1}{n} \sum_{j=1}^n \mathbb{1}\{\hat{z}_j \neq \phi(z_j)\}.$

- exact clustering under (almost) necessary SNR condition among poly-time algorithms
- handle heteroskedastic noise, no superfluous assumptions on \mathcal{S}^*

Theoretical guarantees

Theorem 2 (Zhou, Chen '23)

Assume that either (1) or (2) is satisfied:

- (1). noise is not too spiky (e.g., $\text{Var}[E_{i,j,\ell}] \asymp \omega_{\max}^2$);
- (2). The observation model is the stochastic tensor block model.

Then we can choose τ in a data-driven manner, w.h.p., HHC yields

$$\text{MCR}(\hat{z}_i, z_i^\star) = 0, \quad \forall 1 \leq i \leq 3.$$

Theoretical guarantees

Theorem 2 (Zhou, Chen '23)

Assume that either (1) or (2) is satisfied:

- (1). noise is not too spiky (e.g., $\text{Var}[E_{i,j,\ell}] \asymp \omega_{\max}^2$);
- (2). The observation model is the stochastic tensor block model.

Then we can choose τ in a data-driven manner, w.h.p., HHC yields

$$\text{MCR}(\hat{z}_i, z_i^*) = 0, \quad \forall 1 \leq i \leq 3.$$

Theorem 2 + Han et al. '22 Theorem 2: HHC + HLlloyd can achieve same theoretical guarantees

Comparisons with prior theory

	heteroskedastic noise	no superfluous spectrum assumptions	SNR for exact recovery
HSC Han et al. '22	✗	✓	$n^{-3/4}$
HSC + HLloyd Han et al. '22	✗	✓	$n^{-3/4}$
Successive Projection Agterberg and Zhang '22	✓	$\kappa \lesssim n^{1/8}$	$\kappa n^{-3/4}$
HHC	✓	✓	$n^{-3/4}$
HHC + HLloyd	✓	✓	$n^{-3/4}$

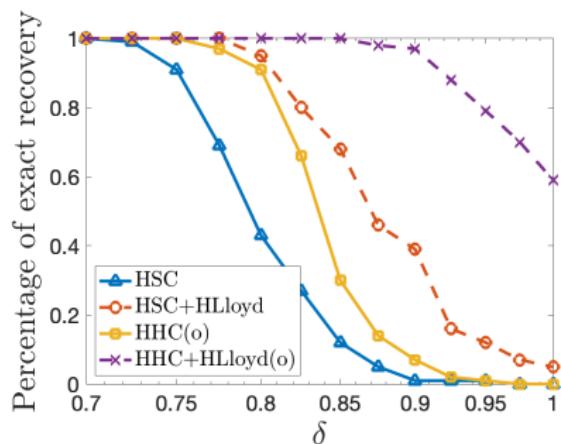
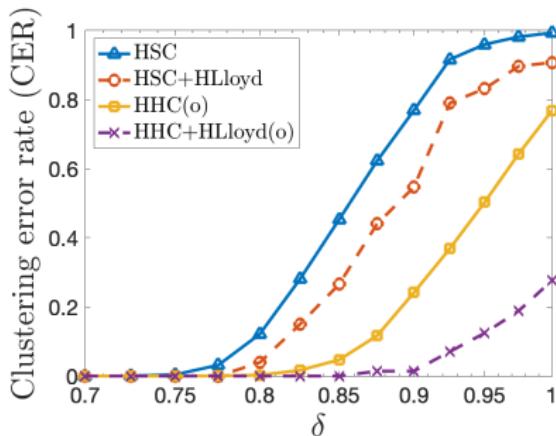
Comparisons with prior theory

	heteroskedastic noise	no superfluous spectrum assumptions	SNR for exact recovery
HSC Han et al. '22	✗	✓	$n^{-3/4}$
HSC + HLloyd Han et al. '22	✗	✓	$n^{-3/4}$
Successive Projection Agterberg and Zhang '22	✓	$\kappa \lesssim n^{1/8}$	$\kappa n^{-3/4}$
HHC	✓	✓	$n^{-3/4}$
HHC + HLloyd	✓	✓	$n^{-3/4}$

$\kappa :=$ condition number of \mathcal{S}^*

Numerical experiments

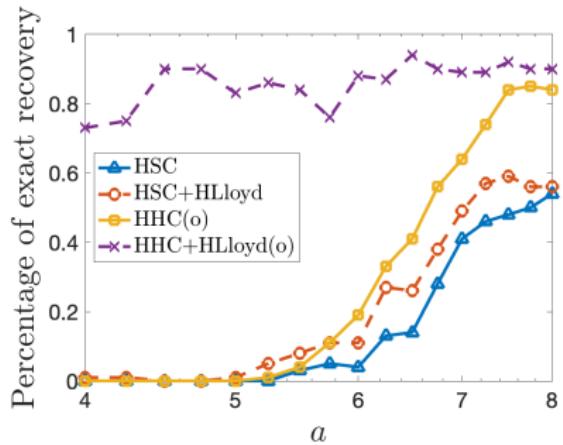
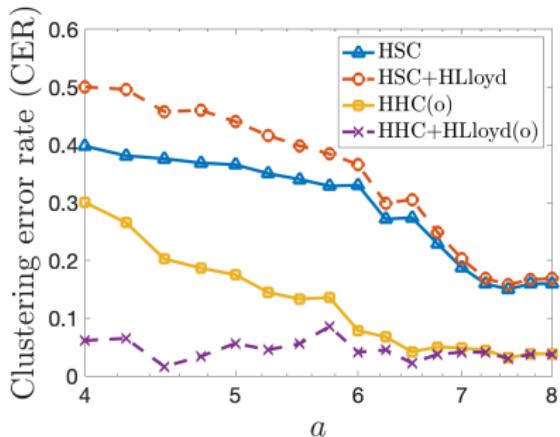
Sub-Gaussian tensor block models



$$n = 100, k = 3, \Delta_{\min} = 40n^{-\delta}, \omega_{\max} = 1$$

Numerical experiments

Stochastic tensor block models



- $n = 100, k = 3$
- $S_{i_1, i_2, i_3}^* = \begin{cases} 10a \cdot n^{-3/2} \left(1 - \frac{i_1-1}{2(k-1)}\right), & i_1 = i_2 = i_3, \\ 0.1a \cdot n^{-3/2}, & \text{otherwise.} \end{cases}$

Real data example: the flight route network

- OpenFlights Airports Database: global flight information*
- consider top 50 airports based on the number of flights → 39 airlines†
- data $\mathcal{Y} \in \{0, 1\}^{39 \times 50 \times 50}$:

$$Y_{i,j,\ell} = \begin{cases} 1, & \text{if airline } i \text{ operates a flight route b/w airports } j, \ell, \\ 0, & \text{otherwise.} \end{cases}$$

*original data: <https://openflights.org/data>

†processed data: https://github.com/RungangHLloyd/blob/master/experiment/flight_route.RData

Real data example: the flight route network

- Select the clustering sizes based on BIC: $(k_1, k_2, k_3) = (5, 5, 5)$

Real data example: the flight route network

- Select the clustering sizes based on BIC: $(k_1, k_2, k_3) = (5, 5, 5)$
- Airport clustering results using HHC + HLlloyd (our method)
 1. Mixture
 2. USA
 3. Europe: LHR (London), MAD (Madrid), CDG (Paris), FCO (Rome), USA: ATL (Atlanta)
 4. China
 5. North America + FRA (Germany)
- Airport clustering results using HSC + HLlloyd (Han et al. '22)
 1. Mixture
 2. USA + LHR (UK)
 3. CDG (France), ATL (USA)
 4. China
 5. North America + Europe

Summary

a novel method called High-order HeteroClustering that

- can handle heteroskedastic noise
- achieves exact clustering if SNR exceeds computational limit

Summary

a novel method called High-order HeteroClustering that

- can handle heteroskedastic noise
- achieves exact clustering if SNR exceeds computational limit

subspace estimation: new algorithm called **Deflated-HeteroPCA** that

- breaks curse of ill-conditioning & achieves near-optimal statistical guarantees (ℓ_2 and $\ell_{2,\infty}$)
- can also lead to remarkable improvement for tensor PCA

Summary

a novel method called High-order HeteroClustering that

- can handle heteroskedastic noise
- achieves exact clustering if SNR exceeds computational limit

subspace estimation: new algorithm called **Deflated-HeteroPCA** that

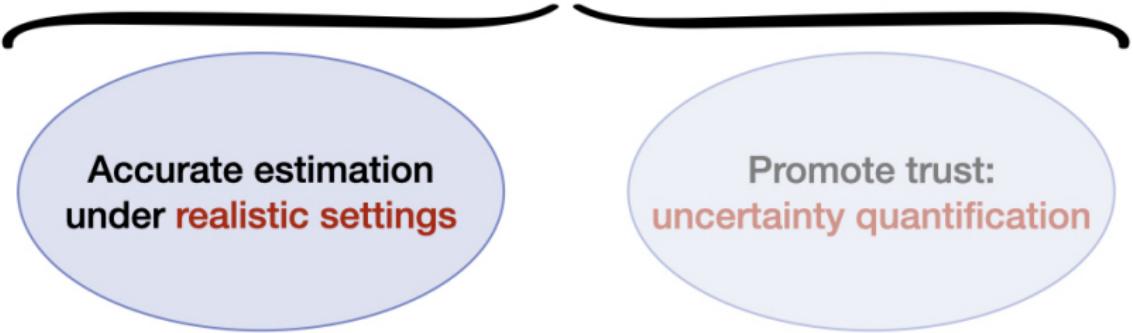
- breaks curse of ill-conditioning & achieves near-optimal statistical guarantees (ℓ_2 and $\ell_{2,\infty}$)
- can also lead to remarkable improvement for tensor PCA

papers:

Y. Zhou, Y. Chen, "Deflated HeteroPCA: Overcoming the Curse of Ill-conditioning in Heteroskedastic PCA," under review at *Annals of Statistics*

Y. Zhou, Y. Chen, "Heteroskedastic Tensor Clustering," under review at *JRSSB*

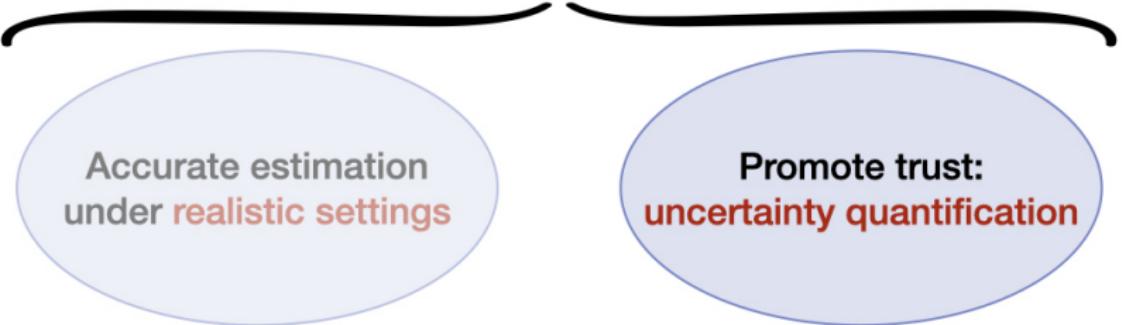
More reliable tensor learning



Accurate estimation
under **realistic settings**

Promote trust:
uncertainty quantification

More reliable tensor learning



```
graph TD; A([Accurate estimation under realistic settings]); B([Promote trust: uncertainty quantification]); A --- C; B --- C;
```

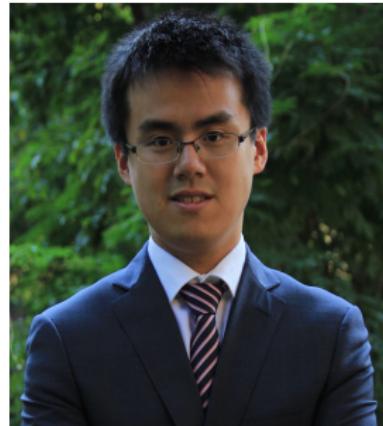
Accurate estimation
under **realistic settings**

Promote trust:
uncertainty quantification

Part 2:
Uncertainty quantification for tensor learning



Dong Xia
HKUST



Anru Zhang
Duke University

Uncertainty Quantification for Tensor Learning?

- Many methods on low-rank tensor estimation
 - tensor PCA
 - tensor clustering
 - tensor regression
 - ...

Uncertainty Quantification for Tensor Learning?

- Many methods on low-rank tensor estimation
 - tensor PCA
 - tensor clustering
 - tensor regression
 - ...
- **The statistical inference or uncertainty quantification** for low-rank tensor models remains largely unexplored!

Uncertainty Quantification for Tensor Learning?

We focus on

- inference of principle components of Tucker low-rank tensors under tensor PCA model

Uncertainty Quantification for Tensor Learning?

We focus on

- **inference of principle components** of Tucker low-rank tensors under tensor PCA model

Covered in the paper but not the talk:

- **entry-wise inference** for rank-1 tensors
- **inference of principle components** of Tucker low-rank tensors under tensor regression model

Uncertainty Quantification for Tensor Learning?

We focus on

- **inference of principle components** of Tucker low-rank tensors under tensor PCA model

Covered in the paper but not the talk:

- **entry-wise inference** for rank-1 tensors
- **inference of principle components** of Tucker low-rank tensors under tensor regression model

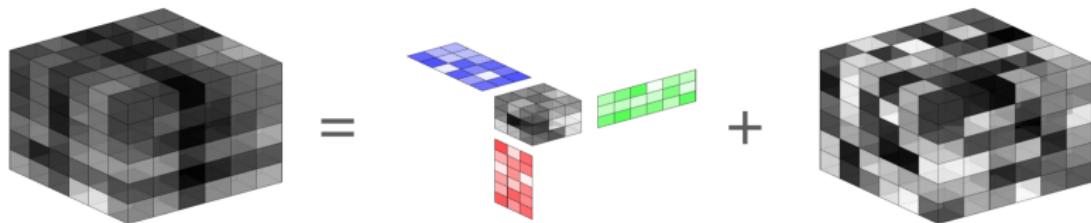
First inference results for Tucker low-rank tensors!

Inference for Tucker low-rank tensor PCA

Tensor PCA model: Observe

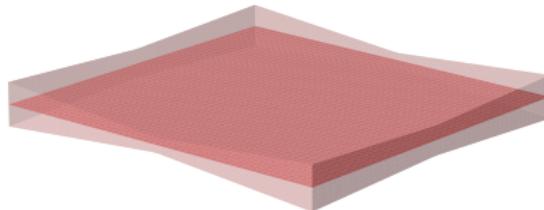
$$\mathcal{A} = \mathcal{T}^* + \mathcal{Z} \in \mathbb{R}^{n_1 \times n_2 \times n_3}, \quad (1)$$

where the signal $\mathcal{T}^* = \mathcal{G}^* \times_1 \mathbf{U}_1^* \times_2 \mathbf{U}_2^* \times_3 \mathbf{U}_3^*$ and the noise \mathcal{Z} contains i.i.d. entries $N(0, \omega^2)$.

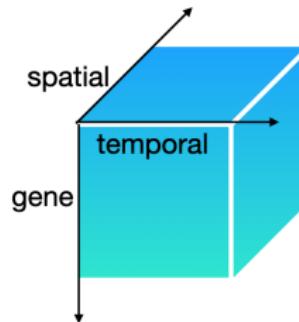


Inference for Tucker low-rank tensor PCA

Goal: Inference for $\{U_i\}_{i=1}^3$



- Spatial and temporal patterns of gene regulation during brain development (Liu et al. '17)



- Hidden components in Gaussian mixture models (Anandkumar et al. '14)

Review: estimation

Estimation

- Least squares estimator:

$$\begin{aligned}\hat{\mathcal{T}} &= \arg \min_{\text{rank}(\mathcal{T}) \leq (r_1, r_2, r_3)} \|\mathcal{T} - \mathcal{A}\|_{\text{F}}^2 \\ \iff (\widehat{\mathbf{U}}_1, \widehat{\mathbf{U}}_2, \widehat{\mathbf{U}}_3) &= \arg \max_{\mathbf{U}_j \in \mathbb{O}_{n_j, r_j}} \|\mathcal{A} \times_1 \mathbf{U}_1^\top \times_2 \mathbf{U}_2^\top \times_3 \mathbf{U}_3^\top\|_{\text{F}}^2.\end{aligned}$$

Highly non-convex!

- High-order orthogonal iteration (HOOI, De Lathauwer et al. '00):
alternating minimization

High-order orthogonal iteration (HOOI)

- **Spectral initialization:** $\widehat{\mathbf{U}}_i^{(0)}$, $1 \leq i \leq 3$
- **Sequential updates:** while $t \leq t_{\max}$

$$\widehat{\mathbf{U}}_1^{(t)} = \text{SVD}_{r_1} \left(\mathcal{M}_1 \left(\mathcal{A} \times_2 \widehat{\mathbf{U}}_2^{(t-1)\top} \times_3 \widehat{\mathbf{U}}_3^{(t-1)\top} \right) \right)$$

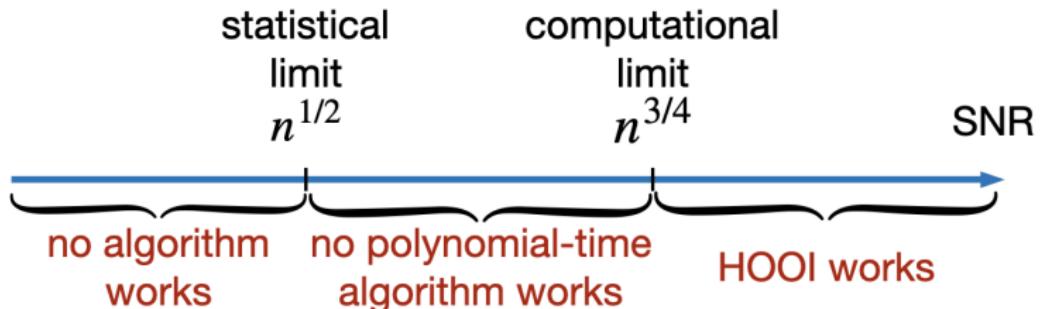
$$\widehat{\mathbf{U}}_2^{(t)} = \text{SVD}_{r_2} \left(\mathcal{M}_2 \left(\mathcal{A} \times_3 \widehat{\mathbf{U}}_3^{(t-1)\top} \times_1 \widehat{\mathbf{U}}_1^{(t-1)\top} \right) \right)$$

$$\widehat{\mathbf{U}}_3^{(t)} = \text{SVD}_{r_3} \left(\mathcal{M}_3 \left(\mathcal{A} \times_1 \widehat{\mathbf{U}}_1^{(t-1)\top} \times_2 \widehat{\mathbf{U}}_2^{(t-1)\top} \right) \right)$$

$$t \leftarrow t + 1$$

- **Output:** $\widehat{\mathbf{U}}_i := \widehat{\mathbf{U}}_i^{(t_{\max})} \longrightarrow \text{estimate of } \mathbf{U}_i^*$
 $\widehat{\mathcal{G}} := \mathcal{A} \times_1 \widehat{\mathbf{U}}_1^\top \times_2 \widehat{\mathbf{U}}_2^\top \times_3 \widehat{\mathbf{U}}_3^\top \longrightarrow \text{estimate of } \mathcal{G}^*$

Statistical-computational gap



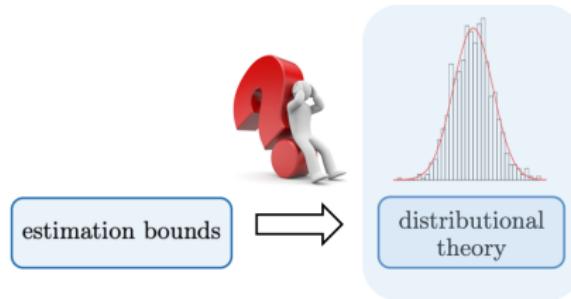
signal strength $\lambda_{\min}^* := \lambda_{\min}(\mathcal{T}^*) = \min_j \sigma_{r_j}(\mathcal{M}_j(\mathcal{T}^*))$

noise level $\omega :=$ standard deviation of entries of \mathcal{Z}

SNR := λ_{\min}^*/ω

Inference for Tucker Low-rank Tensor PCA

- How to make inference?
- Challenge: distribution theory for subspace estimates from non-convex iterations



- The same (or slightly stronger) SNR condition enables us to make inference and no need to debias!

Two-iteration algorithm for inference

- **Input:** Estimators $\widehat{\boldsymbol{U}}_1^{(0)}, \widehat{\boldsymbol{U}}_2^{(0)}, \widehat{\boldsymbol{U}}_3^{(0)}$
- **Sequential updates:** for $t = 0, 1$

$$\widehat{\boldsymbol{U}}_1^{(t+1)} = \text{SVD}_{r_1} \left(\mathcal{M}_1 \left(\mathcal{A} \times_2 \widehat{\boldsymbol{U}}_2^{(t)\top} \times_3 \widehat{\boldsymbol{U}}_3^{(t)\top} \right) \right)$$

$$\widehat{\boldsymbol{U}}_2^{(t+1)} = \text{SVD}_{r_2} \left(\mathcal{M}_2 \left(\mathcal{A} \times_3 \widehat{\boldsymbol{U}}_3^{(t)\top} \times_1 \widehat{\boldsymbol{U}}_1^{(t)\top} \right) \right)$$

$$\widehat{\boldsymbol{U}}_3^{(t+1)} = \text{SVD}_{r_3} \left(\mathcal{M}_3 \left(\mathcal{A} \times_1 \widehat{\boldsymbol{U}}_1^{(t)\top} \times_2 \widehat{\boldsymbol{U}}_2^{(t)\top} \right) \right)$$

- **Output:** Test statistic $\widehat{\boldsymbol{U}}_1 := \widehat{\boldsymbol{U}}_1^{(2)}, \widehat{\boldsymbol{U}}_2 := \widehat{\boldsymbol{U}}_2^{(2)}, \widehat{\boldsymbol{U}}_3 := \widehat{\boldsymbol{U}}_3^{(2)}, \widehat{\mathcal{G}} := \mathcal{A} \times_1 \widehat{\boldsymbol{U}}_1^\top \times_2 \widehat{\boldsymbol{U}}_2^\top \times_3 \widehat{\boldsymbol{U}}_3^\top$

Assumptions

- $n_j \asymp n$ for $1 \leq j \leq 3$.
- **Initialization Error:** $\|\sin \Theta(\widehat{\mathbf{U}}_j^{(0)}, \mathbf{U}_j^*)\| \lesssim \sqrt{n}\omega/\lambda_{\min}^*$ w.h.p.
 - Attainable by HOOI under the essential SNR condition $\lambda_{\min}^*/\omega \gtrsim n^{3/4}$ (Zhang and Xia, 2018, Theorem 1)!
- **signal-to-noise ratio (SNR):** $\frac{\lambda_{\min}^*}{\omega} \gg n^{3/4}$ (ignoring log factors)
- ranks $r_j = O(1)$, condition number $\kappa = O(1)$

Here, $\|\sin \Theta(\widehat{\mathbf{U}}_j^{(0)}, \mathbf{U}_j^*)\| = \sqrt{1 - \sigma_{r_j}^2(\mathbf{U}_j^{*\top} \widehat{\mathbf{U}}_j^{(0)})}$.

Theoretical guarantees

Theorem 3 (Xia, Zhang, Zhou '22)

$$\frac{\|\sin \Theta(\widehat{\mathbf{U}}_j, \mathbf{U}_j^*)\|_{\text{F}}^2 - n_j \omega^2 \|\Lambda_j^{\star-1}\|_{\text{F}}^2}{\sqrt{2n_j} \omega^2 \|\Lambda_j^{\star-2}\|_{\text{F}}} \sim N(0, 1) + \text{negligible terms},$$

where $\Lambda_j^{\star} = \text{diag}(\lambda_1^{\star(j)}, \dots, \lambda_{r_j}^{\star(j)})$ is the diagonal matrix containing the singular values of $\mathcal{M}_j(\mathcal{T}^*)$.

$$\|\sin \Theta(\widehat{\mathbf{U}}_j^{(0)}, \mathbf{U}_j^*)\|_{\text{F}} = \sqrt{r_j - \sum_{i=1}^{r_j} \sigma_i^2(\mathbf{U}_j^{\star\top} \widehat{\mathbf{U}}_j^{(0)})}.$$

Theoretical guarantees

Theorem 3 (Xia, Zhang, Zhou '22)

$$\frac{\|\sin \Theta(\widehat{U}_j, U_j^*)\|_F^2 - n_j \omega^2 \|\Lambda_j^{*-1}\|_F^2}{\sqrt{2n_j} \omega^2 \|\Lambda_j^{*-2}\|_F} \sim N(0, 1) + \text{negligible terms},$$

where $\Lambda_j^* = \text{diag}(\lambda_1^{*(j)}, \dots, \lambda_{r_j}^{*(j)})$ is the diagonal matrix containing the singular values of $\mathcal{M}_j(\mathcal{T}^*)$.

- SNR condition $\lambda_{\min}^*/\omega \gg n^{3/4}$ is slightly stronger than the one for estimation ($\lambda_{\min}^*/\omega \gtrsim n^{3/4}$)

Theoretical guarantees

Theorem 3 (Xia, Zhang, Zhou '22)

$$\frac{\|\sin \Theta(\widehat{U}_j, U_j^*)\|_F^2 - n_j \omega^2 \|\Lambda_j^{*-1}\|_F^2}{\sqrt{2n_j} \omega^2 \|\Lambda_j^{*-2}\|_F} \sim N(0, 1) + \text{negligible terms},$$

where $\Lambda_j^* = \text{diag}(\lambda_1^{*(j)}, \dots, \lambda_{r_j}^{*(j)})$ is the diagonal matrix containing the singular values of $\mathcal{M}_j(\mathcal{T}^*)$.

- SNR condition $\lambda_{\min}^*/\omega \gg n^{3/4}$ is slightly stronger than the one for estimation ($\lambda_{\min}^*/\omega \gtrsim n^{3/4}$)
- To make inference, we still need to estimate Λ_j^* and ω^2

Theoretical guarantees

Estimates for Λ_j^* and ω^2 :

$\widehat{\Lambda}_j$: diagonal matrix containing the top r_j singular values of

$$\mathcal{M}_j(\mathcal{A} \times_{j+1} \widehat{\mathbf{U}}_{j+1}^\top \times_{j+2} \widehat{\mathbf{U}}_{j+2}^\top),$$

$$\widehat{\omega} = \|\mathcal{A} - \underbrace{\mathcal{A} \times_1 \widehat{\mathbf{U}}_1 \widehat{\mathbf{U}}_1^\top \times_2 \widehat{\mathbf{U}}_2 \widehat{\mathbf{U}}_2^\top \times_3 \widehat{\mathbf{U}}_3 \widehat{\mathbf{U}}_3^\top}_{\approx \mathcal{T}}\|_{\text{F}} / \sqrt{n_1 n_2 n_3}.$$

Theoretical guarantees

Estimates for Λ_j^* and ω^2 :

$\widehat{\Lambda}_j$: diagonal matrix containing the top r_j singular values of

$$\mathcal{M}_j(\mathcal{A} \times_{j+1} \widehat{\mathbf{U}}_{j+1}^\top \times_{j+2} \widehat{\mathbf{U}}_{j+2}^\top),$$

$$\widehat{\omega} = \left\| \mathcal{A} - \underbrace{\mathcal{A} \times_1 \widehat{\mathbf{U}}_1 \widehat{\mathbf{U}}_1^\top \times_2 \widehat{\mathbf{U}}_2 \widehat{\mathbf{U}}_2^\top \times_3 \widehat{\mathbf{U}}_3 \widehat{\mathbf{U}}_3^\top}_{\approx \mathcal{T}} \right\|_{\text{F}} / \sqrt{n_1 n_2 n_3}.$$

Theorem 3 (Xia, Zhang, Zhou '22)

$$\frac{\|\sin \Theta(\widehat{\mathbf{U}}_j, \mathbf{U}_j^*)\|_{\text{F}}^2 - n_j \widehat{\omega}^2 \|\widehat{\Lambda}_j^{-1}\|_{\text{F}}^2}{\sqrt{2n_j} \widehat{\omega}^2 \|\widehat{\Lambda}_j^{-2}\|_{\text{F}}} \sim N(0, 1) + \text{negligible terms.}$$

Theoretical guarantees

Estimates for Λ_j^* and ω^2 :

$\widehat{\Lambda}_j$: diagonal matrix containing the top r_j singular values of

$$\mathcal{M}_j(\mathcal{A} \times_{j+1} \widehat{\mathbf{U}}_{j+1}^\top \times_{j+2} \widehat{\mathbf{U}}_{j+2}^\top),$$

$$\widehat{\omega} = \|\mathcal{A} - \underbrace{\mathcal{A} \times_1 \widehat{\mathbf{U}}_1 \widehat{\mathbf{U}}_1^\top \times_2 \widehat{\mathbf{U}}_2 \widehat{\mathbf{U}}_2^\top \times_3 \widehat{\mathbf{U}}_3 \widehat{\mathbf{U}}_3^\top}_{\approx \mathcal{T}}\|_{\text{F}} / \sqrt{n_1 n_2 n_3}.$$

Theoretical guarantees

Estimates for Λ_j^* and ω^2 :

$\widehat{\Lambda}_j$: diagonal matrix containing the top r_j singular values of

$$\mathcal{M}_j(\mathcal{A} \times_{j+1} \widehat{\mathbf{U}}_{j+1}^\top \times_{j+2} \widehat{\mathbf{U}}_{j+2}^\top),$$

$$\widehat{\omega} = \left\| \mathcal{A} - \underbrace{\mathcal{A} \times_1 \widehat{\mathbf{U}}_1 \widehat{\mathbf{U}}_1^\top \times_2 \widehat{\mathbf{U}}_2 \widehat{\mathbf{U}}_2^\top \times_3 \widehat{\mathbf{U}}_3 \widehat{\mathbf{U}}_3^\top}_{\approx \mathcal{T}} \right\|_{\text{F}} / \sqrt{n_1 n_2 n_3}.$$

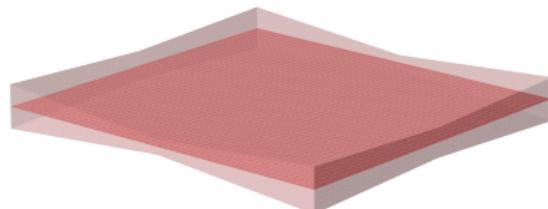
Theorem 4 (Xia, Zhang, Zhou '22)

$$\frac{\|\sin \Theta(\widehat{\mathbf{U}}_j, \mathbf{U}_j^*)\|_{\text{F}}^2 - n_j \widehat{\omega}^2 \|\widehat{\Lambda}_j^{-1}\|_{\text{F}}^2}{\sqrt{2n_j} \widehat{\omega}^2 \|\widehat{\Lambda}_j^{-2}\|_{\text{F}}} \sim N(0, 1) + \text{negligible terms.}$$

Inference for Tucker Low-rank Tensor PCA

$(1 - \alpha)$ -level confidence region for \mathbf{U}_j^* : Let $z_\alpha = \Phi^{-1}(1 - \alpha)$,

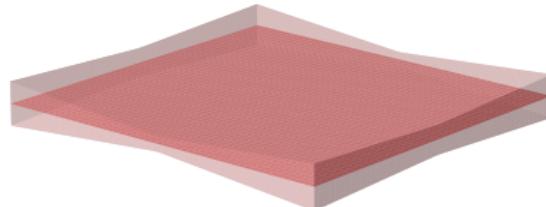
$$\text{CR}_\alpha(\widehat{\mathbf{U}}_j) := \left\{ \mathbf{V} \in \mathbb{O}_{n_j, r_j} : \|\sin \Theta(\widehat{\mathbf{U}}_j, \mathbf{V})\|_{\text{F}}^2 \leq n_j \widehat{\omega}^2 \|\widehat{\boldsymbol{\Lambda}}_j^{-1}\|_{\text{F}}^2 + z_\alpha \sqrt{2n_j} \widehat{\omega}^2 \|\widehat{\boldsymbol{\Lambda}}_j^{-2}\|_{\text{F}} \right\}.$$



Inference for Tucker Low-rank Tensor PCA

$(1 - \alpha)$ -level confidence region for \mathbf{U}_j^* : Let $z_\alpha = \Phi^{-1}(1 - \alpha)$,

$$\text{CR}_\alpha(\widehat{\mathbf{U}}_j) := \left\{ \mathbf{V} \in \mathbb{O}_{n_j, r_j} : \|\sin \Theta(\widehat{\mathbf{U}}_j, \mathbf{V})\|_{\text{F}}^2 \leq n_j \widehat{\omega}^2 \|\widehat{\boldsymbol{\Lambda}}_j^{-1}\|_{\text{F}}^2 + z_\alpha \sqrt{2n_j} \widehat{\omega}^2 \|\widehat{\boldsymbol{\Lambda}}_j^{-2}\|_{\text{F}} \right\}.$$



Corollary 1 (Confidence region for tensor PCA)

If $\lambda_{\min}^*/\omega \gg n^{3/4}$, then $\lim_{p \rightarrow \infty} \mathbb{P}(\mathbf{U}_j^* \in \text{CR}_\alpha(\widehat{\mathbf{U}}_j)) = 1 - \alpha$.

Proof Sketch

Assume $\omega = 1$.

$$\begin{aligned} 2\|\sin \Theta(\widehat{\mathbf{U}}_1, \mathbf{U}_1^*)\|_{\text{F}}^2 &= \|\widehat{\mathbf{U}}_1 \widehat{\mathbf{U}}_1^\top - \mathbf{U}_1^* \mathbf{U}_1^{*\top}\|_{\text{F}}^2 = 2r_1 - 2\langle \widehat{\mathbf{U}}_1 \widehat{\mathbf{U}}_1^\top, \mathbf{U}_1^* \mathbf{U}_1^{*\top} \rangle \\ &= -2\langle \mathbf{U}_1^* \mathbf{U}_1^{*\top}, \widehat{\mathbf{U}}_1 \widehat{\mathbf{U}}_1^\top - \mathbf{U}_1^* \mathbf{U}_1^{*\top} \rangle. \end{aligned}$$

Step 1: By the [spectral representation formula](#) (Xia, 2019, Theorem 1):

$$\widehat{\mathbf{U}}_1 \widehat{\mathbf{U}}_1^\top - \mathbf{U}_1^* \mathbf{U}_1^{*\top} = \mathbf{S}_1 + \mathbf{S}_2 + \mathbf{S}_3 + \sum_{k \geq 4} \mathbf{S}_k.$$

Step 2: $\langle \mathbf{U}_1^* \mathbf{U}_1^{*\top}, \mathbf{S}_1 \rangle = 0$.

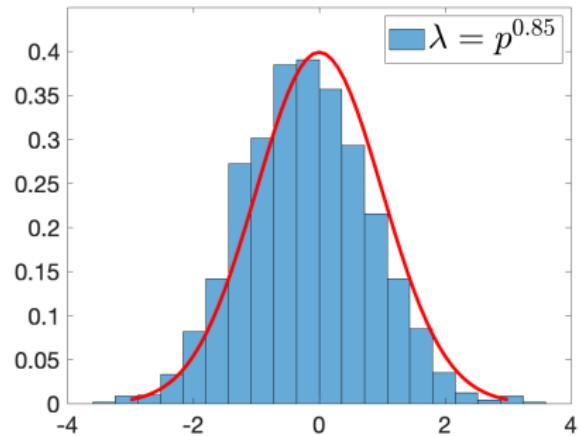
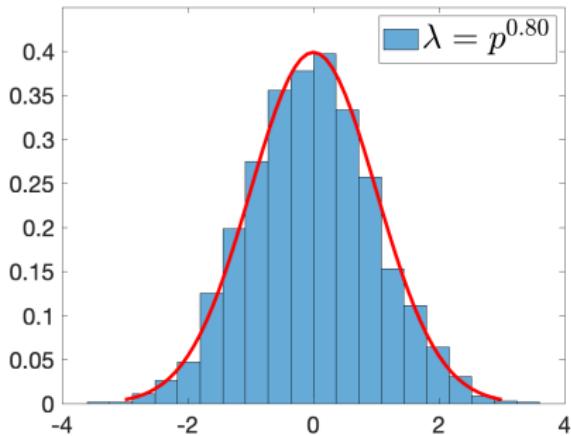
Step 3: $\langle \mathbf{U}_1^* \mathbf{U}_1^{*\top}, \mathbf{S}_2 \rangle = \text{sum of i.i.d. random variables} + \text{high-order terms}$.

Step 4: $|\langle \mathbf{U}_1^* \mathbf{U}_1^{*\top}, \mathbf{S}_3 \rangle| = \text{high-order terms w.h.p.}$

Step 5: $\|\mathbf{S}_k\| \leq (C_1 \sqrt{n}/\lambda_{\min})^k$ w.h.p.

$\Rightarrow \langle \mathbf{U}_1^* \mathbf{U}_1^{*\top}, \sum_{k \geq 4} \mathbf{S}_k \rangle = \text{high-order terms.}$

Numerical simulations



Normal approximation of $\frac{\|\sin \Theta(\widehat{\mathbf{U}}_j, \mathbf{U}_j^\star)\|_F^2 - n_j \widehat{\omega}^2 \|\widehat{\boldsymbol{\Lambda}}_j^{-1}\|_F^2}{\sqrt{2n_j \widehat{\omega}^2 \|\widehat{\boldsymbol{\Lambda}}_j^{-2}\|_F}}$

$n_1 = n_2 = n_3 = n = 200$, $r = 3$, and $\omega = 1$

Summary

- Two-iteration algorithm for inference
- The inference of principle components can be efficiently done when a computationally feasible optimal estimate is achievable!

Summary

- Two-iteration algorithm for inference
- The inference of principle components can be efficiently done when a computationally feasible optimal estimate is achievable!

paper:

D. Xia, A. Zhang, Y. Zhou, "Inference for Low-rank Tensors – No Need to Debias," *Annals of Statistics*, 2022

Highlight of other works

My research focuses on

- designing computationally fast and statistically accurate algorithms for modern statistical problems
- establishing theoretical foundations for new and existing methods

High-order tensor SVD

$$\mathcal{X}_{i_1 \dots i_d} = G_1 \times \cdots \times G_{d-1} \times G_d$$

Diagram illustrating the decomposition of a high-order tensor \mathcal{X} into a product of matrices G_1, G_2, \dots, G_d . The tensor \mathcal{X} has dimensions i_1, i_2, \dots, i_d . The decomposition is shown as:

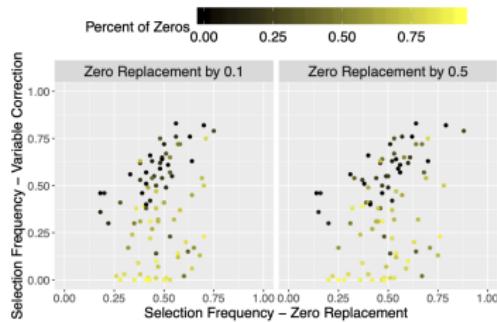
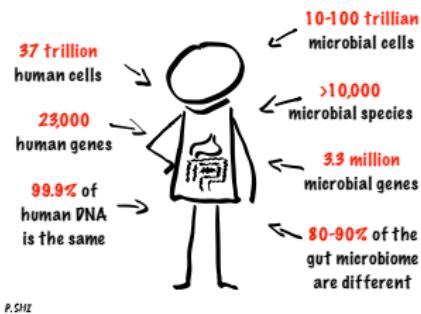
$$\mathcal{X}_{i_1 \dots i_d} = \begin{matrix} G_1 \\ \vdots \\ \text{Color-coded blocks} \end{matrix} \times \begin{matrix} G_2 \\ \vdots \\ \text{Color-coded blocks} \end{matrix} \times \cdots \times \begin{matrix} G_{d-1} \\ \vdots \\ \text{Color-coded blocks} \end{matrix} \times \begin{matrix} G_d \\ \vdots \\ \text{Color-coded blocks} \end{matrix}$$

where $i_1 = 1, i_2 = 2, \dots, i_{d-1} = 1, i_d = 3$.

- Novel method: tensor-train orthogonal iteration (TTOI)
- Theory: *minimax optimal* estimation error

Y. Zhou, A. Zhang, L. Zheng, Y. Wang, "Optimal High-Order Tensor SVD via Tensor-Train Orthogonal Iteration," *IEEE Transactions on Information Theory*, 2022

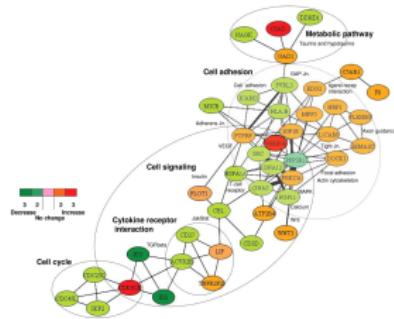
Microbial compositional data



- Goal: Identifying microbial taxa or genes associated with clinical phenotypes
- Propose **log-error-in-variable model** & new method
- Theory: **optimal guarantee** for estimation error

P. Shi, Y. Zhou, A. Zhang, "High-dimensional Log-Error-in-Variable Regression with Applications to Microbial Compositional Data Analysis," *Biometrika*, 2022

Double sparse regression



- Theory: sparse group Lasso achieves **optimal** estimation guarantees
- **Multi-objective regularizers can help!**

T. Cai, A. Zhang, Y. Zhou, "Sparse Group Lasso: Optimal Sample Complexity, Convergence Rate, and Statistical Inference ,," *IEEE Transactions on Information Theory*, 2022

Future directions



- Tensor problems under more realistic assumptions
- Federated statistical learning
- Theoretical foundations of deep learning
- ...

Entry-wise Inference for Rank-1 Tensors

Model: Observe

$$\mathcal{A} = \mathcal{T}^* + \mathcal{Z}, \quad \mathcal{T}^* = \lambda^* \cdot \mathbf{u}^* \otimes \mathbf{v}^* \otimes \mathbf{w}^*.$$

Here, $\mathbf{u}^* \in \mathbb{S}^{n_1-1}$, $\mathbf{v}^* \in \mathbb{S}^{n_2-1}$, $\mathbf{w}^* \in \mathbb{S}^{n_3-1}$, the singular value $\lambda^* > 0$, and $\mathcal{Z} \stackrel{i.i.d.}{\sim} N(0, \omega^2)$.

Goal: Inference for $T_{i,j,k}^*$

Entry-wise Inference for Rank-1 Tensors

Theorem 5

Suppose $\lambda^*/\omega \gg n^{3/4}$. Apply HOOI with $t_{\max} \geq C_1 \log n$. If $|u_i^*|, |v_j^*|, |w_k^*| \ll \min\{\lambda^*/(n\omega), 1\}$, then

$$\left(\frac{\lambda}{\omega}(\hat{u}_i - u_i^*), \frac{\lambda}{\omega}(\hat{v}_j - v_j^*), \frac{\lambda}{\omega}(\hat{w}_k - w_k^*) \right)^\top \xrightarrow{\text{d}} N(0, I_3) \quad \text{as } p \rightarrow \infty.$$

If $\omega/\lambda^* \ll |u_i^*|, |v_j^*|, |w_k^*| \ll \min\{\lambda^*/(n\omega), 1/\sqrt{\log n}\}$, then

$$\frac{\hat{T}_{ijk} - \mathcal{T}_{ijk}}{\hat{\omega} \sqrt{\hat{u}_i^2 \hat{v}_j^2 + \hat{v}_j^2 \hat{w}_k^2 + \hat{w}_k^2 \hat{u}_i^2}} \xrightarrow{\text{d}} N(0, 1) \quad \text{as } n \rightarrow \infty.$$

Entry-wise Inference for Rank-1 Tensors

- upper bound conditions on $|u_i^*|, |v_j^*|, |w_k^*|$ are weaker than the incoherence condition
- lower bound condition is essential for asymptotic normality: if $u_i^* = v_j^* = w_k^* = 0$, then
$$\frac{\lambda^2 \widehat{T}_{ijk}}{\omega^3} \xrightarrow{d} G_1 G_2 G_3 \text{ as } n \rightarrow \infty, \quad (G_1, G_2, G_3)^\top \sim N(0, I_3).$$
- Asymptotic $(1 - \alpha)$ confidence interval for T_{ijk} :
$$[\widehat{T}_{ijk} - z_{\alpha/2} \sigma \sqrt{\widehat{u}_i^2 \widehat{v}_j^2 + \widehat{v}_j^2 \widehat{w}_k^2 + \widehat{w}_k^2 \widehat{u}_i^2}, \widehat{T}_{ijk} + z_{\alpha/2} \sigma \sqrt{\widehat{u}_i^2 \widehat{v}_j^2 + \widehat{v}_j^2 \widehat{w}_k^2 + \widehat{w}_k^2 \widehat{u}_i^2}].$$