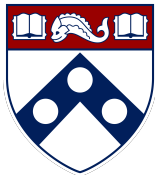


Estimation and Inference for Heteroskedastic PCA with Missing Data



Yuxin Chen

Wharton Statistics & Data Science

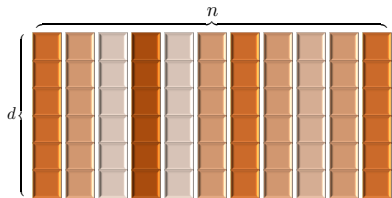


Yuling Yan
Princeton ORFE



Jianqing Fan
Princeton ORFE

Principal component analysis



$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$$

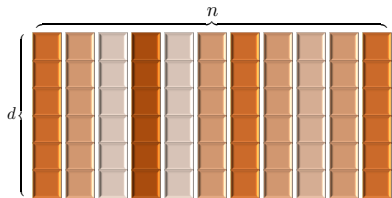
- Ground-truth data

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}, \quad \mathbf{x}_i \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{S}^*)$$

$$\text{where } \mathbf{S}^* = \mathbf{U}^* \mathbf{\Lambda}^* \mathbf{U}^{*\top} \in \mathbb{R}^{d \times d}$$

Principal component analysis

$\text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq U^*$ (r -dimensional)



$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$$

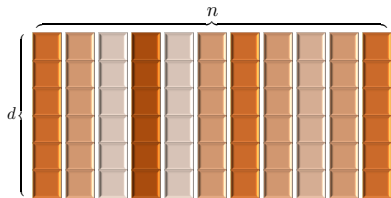
- Ground-truth data

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}, \quad \mathbf{x}_i \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{S}^*)$$

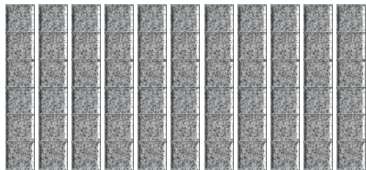
where $\mathbf{S}^* = \mathbf{U}^* \mathbf{\Lambda}^* \mathbf{U}^{*\top} \in \mathbb{R}^{d \times d}$ has **rank** $r \ll d$

Principal component analysis

$\text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq U^*$ (r -dimensional)



$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$$



noise matrix: \mathbf{E}

- Ground-truth data

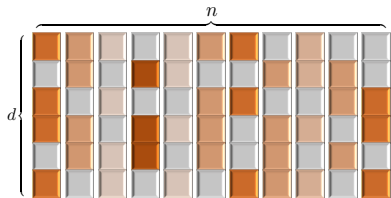
$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}, \quad \mathbf{x}_i \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{S}^*)$$

where $\mathbf{S}^* = \mathbf{U}^* \mathbf{\Lambda}^* \mathbf{U}^{*\top} \in \mathbb{R}^{d \times d}$ has **rank** $r \ll d$

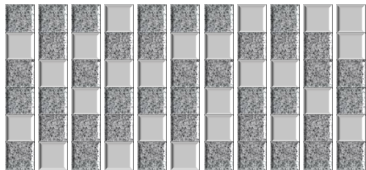
- Noisy observations: $\mathbf{X} + \mathbf{E}$ (a.k.a. spiked covariance model)

Principal component analysis

$\text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq U^*$ (r -dimensional)



$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$$



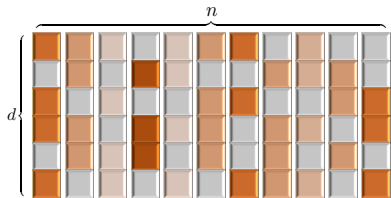
noise matrix: \mathbf{E}

- Incomplete observations \longrightarrow sampling set Ω :

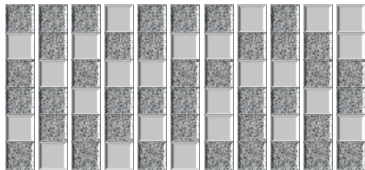
$$Y_{i,j} = \begin{cases} X_{i,j}^* + E_{i,j}, & (i,j) \in \Omega \\ 0, & \text{else} \end{cases} \quad \text{or} \quad \mathbf{Y} = \mathcal{P}_\Omega(\mathbf{X} + \mathbf{E})$$

Principal component analysis

$\text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq U^*$ (r -dimensional)



$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$$



noise matrix: \mathbf{E}

- **Goal:**

- Construct confidence regions for principal subspace U^*
- Construct entrywise confidence intervals for covariance matrix \mathbf{S}^*

What we consider here . . .

- **Heteroskedastic noise:** $\{E_{i,j}\}$ are ind. sub-Gaussian obeying

$$\mathbb{E}[E_{i,j}] = 0, \quad \mathbb{E}[E_{i,j}^2] = \omega_i^{*2} \in [\omega_{\min}^2, \omega_{\max}^2], \quad \underbrace{\|E_{i,j}\|_{\psi_2}}_{\text{sub-Gaussian norm}} = O(\omega_i^*)$$

- noise variance $\{\omega_i^{*2}\}$: **unknown**, location-varying

What we consider here . . .

- **Heteroskedastic noise:** $\{E_{i,j}\}$ are ind. sub-Gaussian obeying

$$\mathbb{E}[E_{i,j}] = 0, \quad \mathbb{E}[E_{i,j}^2] = \omega_i^{*2} \in [\omega_{\min}^2, \omega_{\max}^2], \quad \underbrace{\|E_{i,j}\|_{\psi_2}}_{\text{sub-Gaussian norm}} = O(\omega_i^*)$$

- noise variance $\{\omega_i^{*2}\}$: **unknown**, location-varying

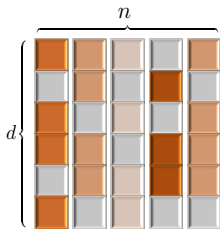
- **Random sampling:** $(i, j) \in \Omega$ independently with prob. p

What we consider here . . .

Our focus: estimating/infering column subspace when $\underbrace{n \gg d}$
more challenging regime

What we consider here ...

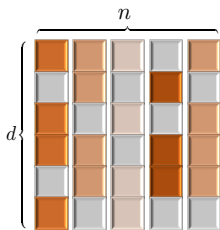
Our focus: estimating/infering column subspace when $\underbrace{n \gg d}$
more challenging regime



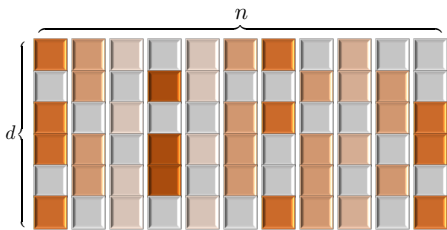
$n \lesssim d$: solvable via *matrix completion* methods
(e.g., Chen, Fan, Ma, Yan '19)

What we consider here ...

Our focus: estimating/infering column subspace when $n \gg d$
more challenging regime



$n \lesssim d$: solvable via *matrix completion* methods
(e.g., Chen, Fan, Ma, Yan '19)



$n \gg d$: sometimes it's only feasible to estimate col-space instead of whole matrix

Applications beyond PCA

- Tensor completion

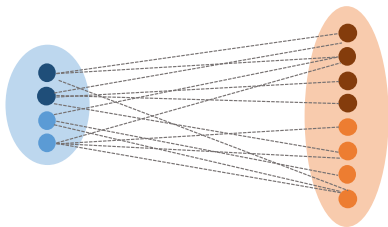


Applications beyond PCA

- Tensor completion



- One-sided community recovery in bipartite random graphs



A natural SVD-based algorithm

- **Compute:** rank- r SVD $U\Sigma V^T$ of $Y = \mathcal{P}_\Omega(X + E)$
- **Output:** $U \rightarrow$ estimate of U^*

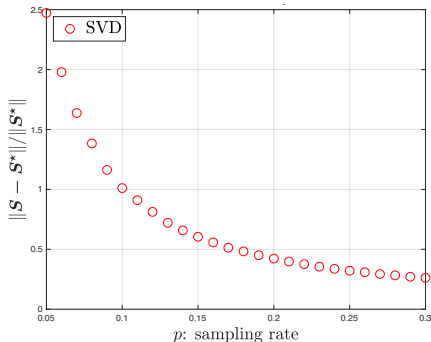
A natural SVD-based algorithm

- **Compute:** rank- r SVD $U\Sigma V^\top$ of $Y = \mathcal{P}_\Omega(\mathbf{X} + \mathbf{E})$
- **Output:** $U \rightarrow$ estimate of U^*

Rationale: under zero-mean noise and random sampling, we have

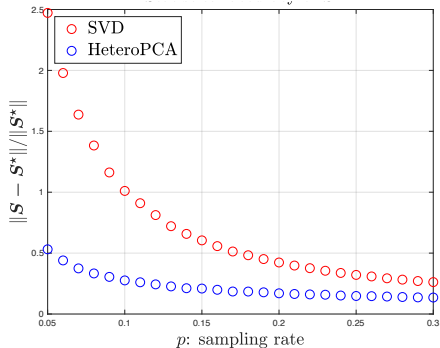
$$\text{col-space}(\mathbb{E}[\mathbf{Y}]) = \text{col-space}(\mathbf{X}) = U^*$$

Numerical suboptimality of SVD-based approach



$n = 2000$, $d = 100$, $r = 3$, $\omega_1^*, \dots, \omega_d^* \stackrel{\text{i.i.d.}}{\sim} \text{Unif}[0.025, 0.1]$

Numerical suboptimality of SVD-based approach



$$n = 2000, \quad d = 100, \quad r = 3, \quad \omega_1^*, \dots, \omega_d^* \stackrel{\text{i.i.d.}}{\sim} \text{Unif}[0.025, 0.1]$$

Plain SVD is suboptimal in the presence of missing data if $n \gg d$

Diagnosis: dagonal entries need special treatment

$$\text{col-space}(\mathbf{Y}) = \text{eig-space}(\mathbf{Y}\mathbf{Y}^{\top})$$

Diagnosis: diagonal entries need special treatment

col-space(\mathbf{Y}) = eig-space($\mathbf{Y}\mathbf{Y}^\top$)

Large bias in diagonal entries:

$$\frac{1}{p^2} \mathbb{E}[\mathbf{Y}\mathbf{Y}^\top] = \underbrace{\mathbf{X}\mathbf{X}^\top}_{\checkmark} + \underbrace{\left(\frac{1}{p} - 1\right) \mathcal{P}_{\text{diag}}(\mathbf{X}\mathbf{X}^\top) + \frac{n}{p} \text{diag}\{\omega_i^{*2}\}}_{\text{potentially large diagonal matrix!}}$$

Diagnosis: diagonal entries need special treatment

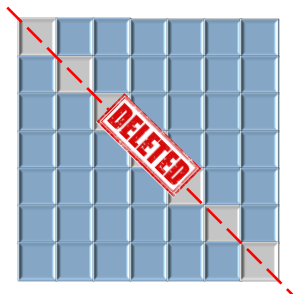
col-space(\mathbf{Y}) = eig-space($\mathbf{Y}\mathbf{Y}^\top$)

Large bias in diagonal entries:

$$\frac{1}{p^2}\mathbb{E}[\mathbf{Y}\mathbf{Y}^\top] = \underbrace{\mathbf{X}\mathbf{X}^\top}_{\checkmark} + \underbrace{\left(\frac{1}{p} - 1\right) \mathcal{P}_{\text{diag}}(\mathbf{X}\mathbf{X}^\top) + \frac{n}{p} \text{diag}\{\omega_i^{*2}\}}_{\text{potentially large diagonal matrix!}}$$

- a common issue under missing data or heteroskedastic noise

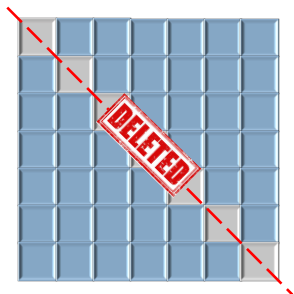
Two spectral algorithms that take care of diagonals



diagonal-deleted/reweighted PCA

- remove/reweight $\text{diag}(\mathbf{Y}\mathbf{Y}^\top)$

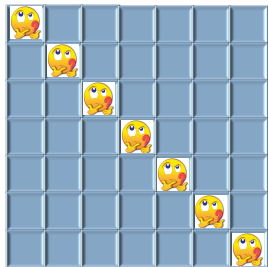
Two spectral algorithms that take care of diagonals



diagonal-deleted/reweighted PCA

- remove/reweight $\text{diag}(\mathbf{Y}\mathbf{Y}^\top)$

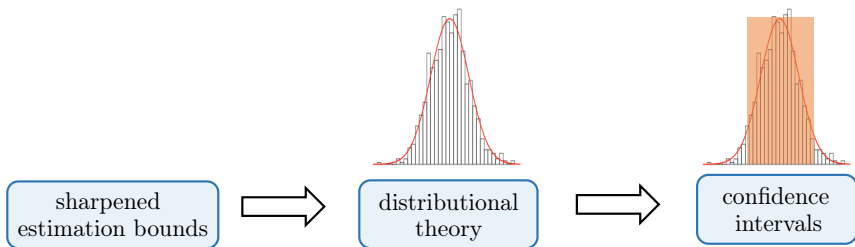
- Loh, Wainwright '12
- Lounici '13 '14
- Florescu and Perkins '16
- Montanari and Sun '18
- Zhu, Wang, Samworth '19
- Cai, Li, Chi, Poor, Chen '19
- ...



HeteroPCA (Zhang et al '18)

- iteratively estimate $\text{diag}(\mathbf{Y}\mathbf{Y}^\top)$

Our contributions: estimation and inference based on HeteroPCA

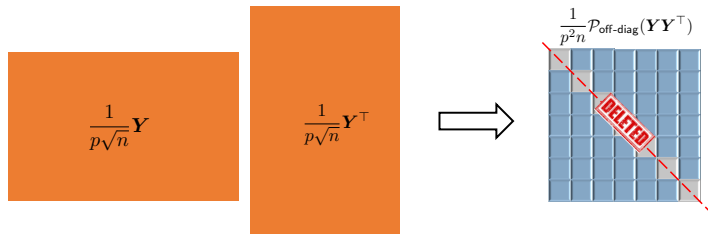


HeteroPCA (Zhang, Cai, Wu '18)

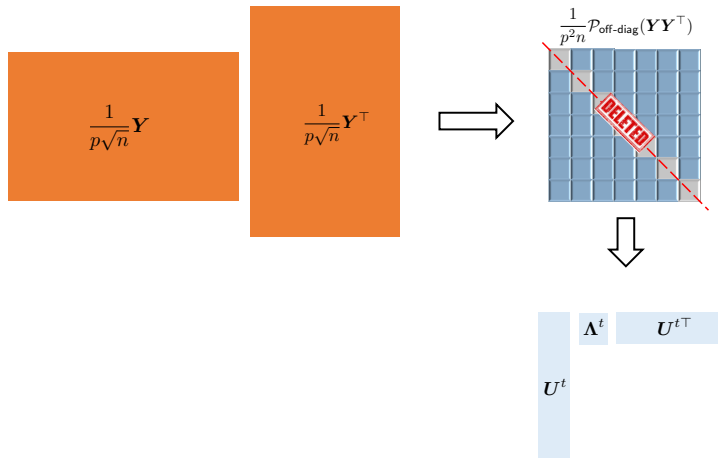
$$\frac{1}{p\sqrt{n}}\mathbf{Y}$$

$$\frac{1}{p\sqrt{n}}\mathbf{Y}^\top$$

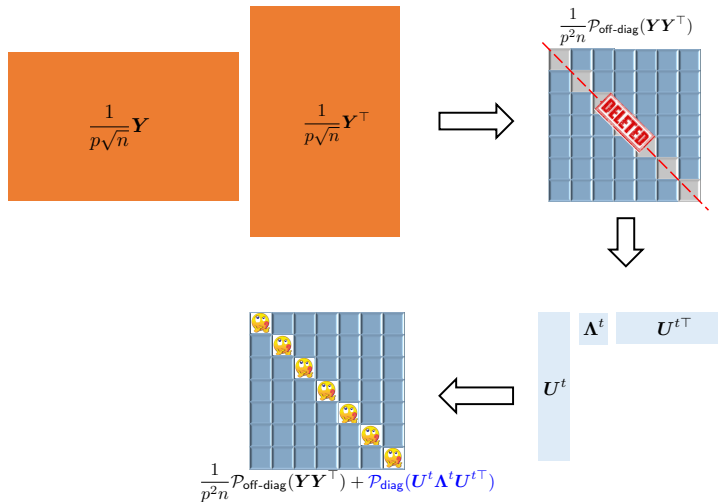
HeteroPCA (Zhang, Cai, Wu '18)



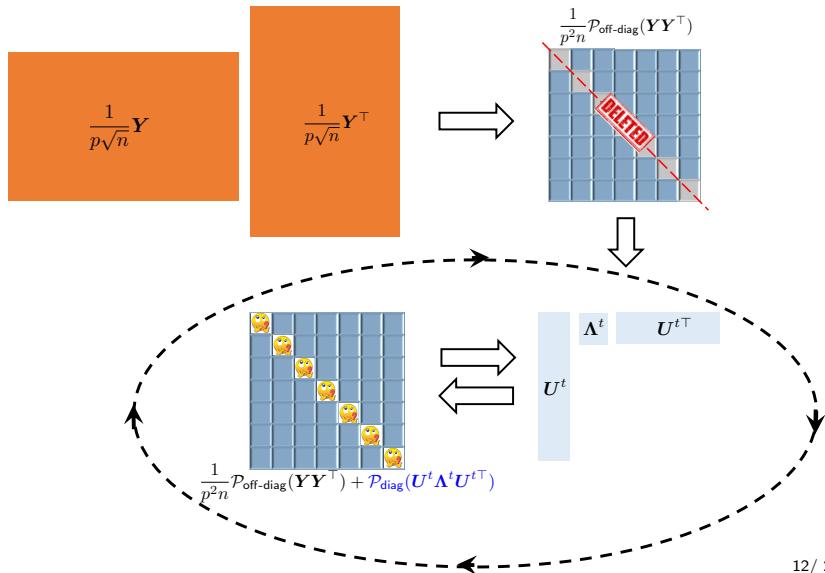
HeteroPCA (Zhang, Cai, Wu '18)



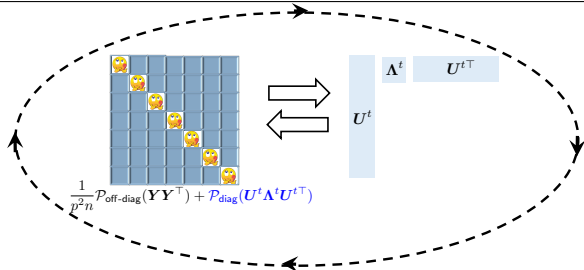
HeteroPCA (Zhang, Cai, Wu '18)



HeteroPCA (Zhang, Cai, Wu '18)

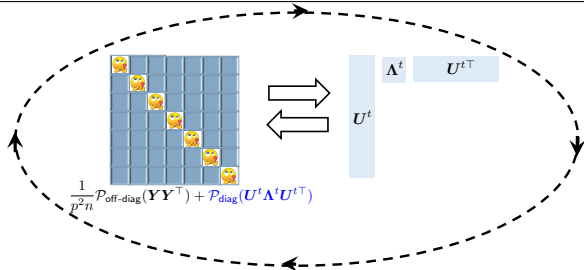


HeteroPCA (Zhang, Cai, Wu '18)



- **Initialize:** $\mathbf{G}^0 = \frac{1}{np^2} \mathcal{P}_{\text{off-diag}}(\mathbf{Y}\mathbf{Y}^\top)$
- **Iterative update:** for $t = 0, 1, \dots, t_0$
 $(\mathbf{U}^t, \mathbf{\Lambda}^t) = \text{eigs}(\mathbf{G}^t, r)$
 $\mathbf{G}^{t+1} = \mathbf{G}^0 + \mathcal{P}_{\text{diag}}(\mathbf{U}^t \mathbf{\Lambda}^t \mathbf{U}^{t\top})$

HeteroPCA (Zhang, Cai, Wu '18)



- **Initialize:** $\mathbf{G}^0 = \frac{1}{np^2} \mathcal{P}_{\text{off-diag}}(\mathbf{Y}\mathbf{Y}^\top)$
- **Iterative update:** for $t = 0, 1, \dots, t_0$

$$(\mathbf{U}^t, \mathbf{\Lambda}^t) = \text{eigs}(\mathbf{G}^t, r)$$

$$\mathbf{G}^{t+1} = \mathbf{G}^t + \mathcal{P}_{\text{diag}}(\mathbf{U}^t \mathbf{\Lambda}^t \mathbf{U}^{t\top})$$
- **Output:** $\mathbf{U} := \mathbf{U}^{t_0} \rightarrow$ estimate of \mathbf{U}^*
 $\mathbf{S} := \mathbf{U}^{t_0} \mathbf{\Lambda}^{t_0} \mathbf{U}^{t_0\top} \rightarrow$ estimate of $\mathbf{S}^* = \mathbf{U}^* \mathbf{\Lambda}^* \mathbf{U}^{*\top}$

Sharpened estimation guarantees for HeteroPCA

Assumptions (omitting log factors)

- rank $r = O(1)$, incoherence $\mu = O(1)$, cond. number $\kappa = O(1)$
- sampling rate exceeds certain threshold

$$p \gtrsim \max \left\{ \frac{1}{\sqrt{nd}}, \frac{1}{n} \right\}$$

- per-entry signal-to-noise ratio (SNR) cannot be too low:

$$\frac{\omega_{\max}^2}{\lambda_r(\mathbf{S}^*)/d} \lesssim \min \{ pn, p\sqrt{nd} \}$$

Sharpened estimation guarantees for HeteroPCA

Theorem 1 (Yan, Chen, Fan '21)

With high prob., we have

$$\begin{aligned}\|U \operatorname{sgn}(U^\top U^*) - U^*\| &\lesssim \zeta_{\text{op}}, & \|U \operatorname{sgn}(U^\top U^*) - U^*\|_{2,\infty} &\lesssim \frac{1}{\sqrt{d}} \zeta_{\text{op}} \\ \|\mathbf{S} - \mathbf{S}^*\| &\lesssim \zeta_{\text{op}} \lambda_1^*, & \|\mathbf{S} - \mathbf{S}^*\|_\infty &\lesssim \frac{1}{d} \zeta_{\text{op}} \lambda_1^*\end{aligned}$$

where $\zeta_{\text{op}} := \frac{1}{\sqrt{ndp}} + \frac{\omega_{\max}^2}{p \lambda_r^*} \sqrt{\frac{d}{n}} + \sqrt{\frac{1}{np}} + \frac{\omega_{\max}}{\sqrt{\lambda_r^*}} \sqrt{\frac{d}{np}}$

Sharpened estimation guarantees for HeteroPCA

Theorem 1 (Yan, Chen, Fan '21)

With high prob., we have

$$\begin{aligned}\|U \operatorname{sgn}(U^\top U^*) - U^*\| &\lesssim \zeta_{\text{op}}, & \|U \operatorname{sgn}(U^\top U^*) - U^*\|_{2,\infty} &\lesssim \frac{1}{\sqrt{d}} \zeta_{\text{op}} \\ \|\mathbf{S} - \mathbf{S}^*\| &\lesssim \zeta_{\text{op}} \lambda_1^*, & \|\mathbf{S} - \mathbf{S}^*\|_\infty &\lesssim \frac{1}{d} \zeta_{\text{op}} \lambda_1^*\end{aligned}$$

where $\zeta_{\text{op}} := \frac{1}{\sqrt{ndp}} + \frac{\omega_{\max}^2}{p \lambda_r^*} \sqrt{\frac{d}{n}} + \sqrt{\frac{1}{np}} + \frac{\omega_{\max}}{\sqrt{\lambda_r^*}} \sqrt{\frac{d}{np}}$

- fine-grained estimation guarantees ($\ell_{2,\infty}$ and ℓ_∞ bounds)

Sharpened estimation guarantees for HeteroPCA

Theorem 1 (Yan, Chen, Fan '21)

With high prob., we have

$$\begin{aligned}\|U \operatorname{sgn}(U^\top U^*) - U^*\| &\lesssim \zeta_{\text{op}}, & \|U \operatorname{sgn}(U^\top U^*) - U^*\|_{2,\infty} &\lesssim \frac{1}{\sqrt{d}} \zeta_{\text{op}} \\ \|\mathbf{S} - \mathbf{S}^*\| &\lesssim \zeta_{\text{op}} \lambda_1^*, & \|\mathbf{S} - \mathbf{S}^*\|_\infty &\lesssim \frac{1}{d} \zeta_{\text{op}} \lambda_1^*\end{aligned}$$

where $\zeta_{\text{op}} := \frac{1}{\sqrt{ndp}} + \frac{\omega_{\max}^2}{p \lambda_r^*} \sqrt{\frac{d}{n}} + \sqrt{\frac{1}{np}} + \frac{\omega_{\max}}{\sqrt{\lambda_r^*}} \sqrt{\frac{d}{np}}$

- fine-grained estimation guarantees ($\ell_{2,\infty}$ and ℓ_∞ bounds)
- estimation errors are spread out across entries

Sharpened estimation guarantees for HeteroPCA

Theorem 1 (Yan, Chen, Fan '21)

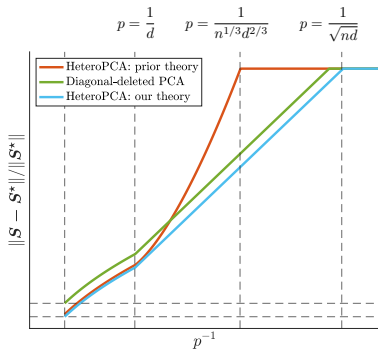
With high prob., we have

$$\begin{aligned}\|U \operatorname{sgn}(U^\top U^*) - U^*\| &\lesssim \zeta_{\text{op}}, & \|U \operatorname{sgn}(U^\top U^*) - U^*\|_{2,\infty} &\lesssim \frac{1}{\sqrt{d}} \zeta_{\text{op}} \\ \|\mathbf{S} - \mathbf{S}^*\| &\lesssim \zeta_{\text{op}} \lambda_1^*, & \|\mathbf{S} - \mathbf{S}^*\|_\infty &\lesssim \frac{1}{d} \zeta_{\text{op}} \lambda_1^*\end{aligned}$$

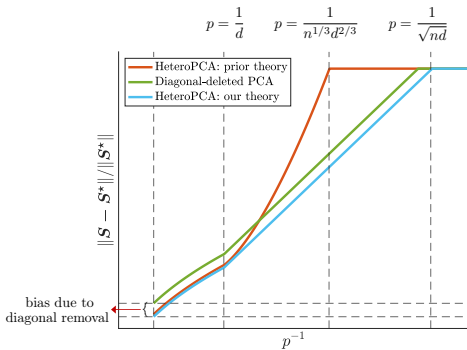
where $\zeta_{\text{op}} := \frac{1}{\sqrt{ndp}} + \frac{\omega_{\max}^2}{p \lambda_r^*} \sqrt{\frac{d}{n}} + \sqrt{\frac{1}{np}} + \frac{\omega_{\max}}{\sqrt{\lambda_r^*}} \sqrt{\frac{d}{np}}$

- fine-grained estimation guarantees ($\ell_{2,\infty}$ and ℓ_∞ bounds)
- estimation errors are spread out across entries
- our sample size and SNR conditions are **minimax-optimal** (in terms of achieving **consistent estimation**)

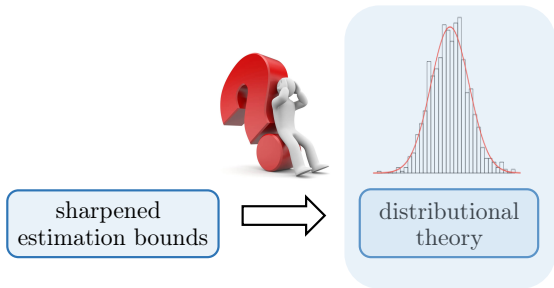
Sharpened estimation guarantees for HeteroPCA



Sharpened estimation guarantees for HeteroPCA



- diagonal-deleted PCA incurs some bias due to diagonal deletion
- HeteroPCA achieves bias correction via iterative refinement
method of choice
- first $\ell_{2,\infty}$ and ℓ_∞ theory for HeteroPCA



Given HeteroPCA is an appealing estimator, can we take one step further to obtain distributional characterizations?

Distributional theory for U

Theorem 2 (Yan, Chen, Fan '21)

Consider any $1 \leq l \leq d$ s.t. $\|U_{l,\cdot}^*\|_2$ is not too small. Under previous assumptions, we have

$$\sup_{\text{cvx set } \mathcal{C}} \left| \mathbb{P} \left(\left[U \underbrace{\text{sgn}(U^\top U^*)}_{\text{global rotation}} - U^* \right]_{l,\cdot} \in \mathcal{C} \right) - \mathcal{N}(\mathbf{0}, \Sigma_{U,l}^*) \{ \mathcal{C} \} \right| = o(1)$$

Distributional theory for U

Theorem 2 (Yan, Chen, Fan '21)

Consider any $1 \leq l \leq d$ s.t. $\|U_{l,\cdot}^*\|_2$ is not too small. Under previous assumptions, we have

$$\sup_{\text{cvx set } \mathcal{C}} \left| \mathbb{P} \left(\left[U \underbrace{\text{sgn}(U^\top U^*)}_{\text{global rotation}} - U^* \right]_{l,\cdot} \in \mathcal{C} \right) - \mathcal{N}(\mathbf{0}, \Sigma_{U,l}^*) \{ \mathcal{C} \} \right| = o(1)$$

- Each row of U is approximately Gaussian
 - nearly unbiased + tractable covariance

Distributional theory for U

Theorem 2 (Yan, Chen, Fan '21)

Consider any $1 \leq l \leq d$ s.t. $\|U_{l,\cdot}^*\|_2$ is not too small. Under previous assumptions, we have

$$\sup_{\text{cvx set } \mathcal{C}} \left| \mathbb{P} \left(\left[\underbrace{U \operatorname{sgn}(U^\top U^*)}_{\text{global rotation}} - U^* \right]_{l,\cdot} \in \mathcal{C} \right) - \mathcal{N}(\mathbf{0}, \Sigma_{U,l}^*) \{ \mathcal{C} \} \right| = o(1)$$

- Each row of U is approximately Gaussian
 - nearly unbiased + tractable covariance

$$\begin{aligned} \Sigma_{U,l}^* &:= \left(\frac{1-p}{np} S_{l,l}^* + \frac{\omega_l^{*2}}{np} \right) (\Lambda^*)^{-1} + \frac{2(1-p)}{np} U_{l,\cdot}^{*\top} U_{l,\cdot}^* \\ &\quad + (\Lambda^*)^{-1} U^{*\top} \operatorname{diag} \{ [d_{l,i}^*]_{1 \leq i \leq d} \} U^* (\Lambda^*)^{-1} \\ d_{l,i}^* &:= \frac{1}{np^2} \left[\omega_l^{*2} + (1-p) S_{l,l}^{*2} \right] \left[\omega_i^{*2} + (1-p) S_{i,i}^{*2} \right] + \frac{2(1-p)^2}{np^2} S_{l,i}^{*2} \end{aligned}$$

Distributional theory for U

Theorem 2 (Yan, Chen, Fan '21)

Consider any $1 \leq l \leq d$ s.t. $\|U_{l,\cdot}^*\|_2$ is not too small. Under previous assumptions, we have

$$\sup_{\text{cvx set } \mathcal{C}} \left| \mathbb{P} \left(\left[U \underbrace{\text{sgn}(U^\top U^*)}_{\text{global rotation}} - U^* \right]_{l,\cdot} \in \mathcal{C} \right) - \mathcal{N}(\mathbf{0}, \Sigma_{U,l}^*) \{ \mathcal{C} \} \right| = o(1)$$

- Key observations:

$$U \text{sgn}(U^\top U^*) - U^* \approx \left[\underbrace{EX^\top}_{\text{linear term}} + \underbrace{\mathcal{P}_{\text{off-diag}}(EE^\top)}_{\text{quadratic term}} \right] U^* (\Lambda^*)^{-1}$$

Distributional theory for S

Theorem 3 (Yan, Chen, Fan '21)

Consider any (i, j) s.t. $\|\mathbf{U}_{i,\cdot}^*\|_2$ and $\|\mathbf{U}_{j,\cdot}^*\|_2$ are not too small. Under previous assumptions, we have

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P} \left(\frac{S_{i,j} - S_{i,j}^*}{\sqrt{v_{i,j}^*}} \leq t \right) - \Phi(t) \right| = o(1)$$

where $\Phi(\cdot)$ is the CDF of $\mathcal{N}(0, 1)$

Distributional theory for S

Theorem 3 (Yan, Chen, Fan '21)

Consider any (i, j) s.t. $\|U_{i,\cdot}^*\|_2$ and $\|U_{j,\cdot}^*\|_2$ are not too small. Under previous assumptions, we have

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P} \left(\frac{S_{i,j} - S_{i,j}^*}{\sqrt{v_{i,j}^*}} \leq t \right) - \Phi(t) \right| = o(1)$$

where $\Phi(\cdot)$ is the CDF of $\mathcal{N}(0, 1)$

- Each entry of S is approximately Gaussian
— nearly unbiased + tractable variance

Distributional theory for \mathcal{S}

Theorem 3 (Yan, Chen, Fan '21)

Consider any (i, j) s.t. $\|\mathbf{U}_{i,\cdot}^*\|_2$ and $\|\mathbf{U}_{j,\cdot}^*\|_2$ are not too small. Under previous assumptions, we have

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P} \left(\frac{S_{i,j} - S_{i,j}^*}{\sqrt{v_{i,j}^*}} \leq t \right) - \Phi(t) \right| = o(1)$$

where $\Phi(\cdot)$ is the CDF of $\mathcal{N}(0, 1)$

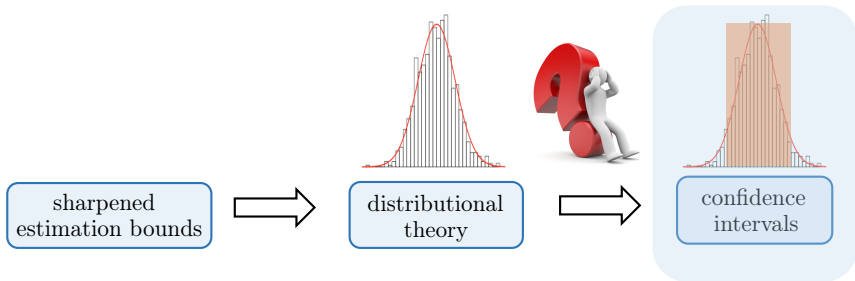
- Each entry of \mathcal{S} is approximately Gaussian
— nearly unbiased + tractable variance

For any $i \neq j$:

$$\begin{aligned} v_{i,j}^* := & \frac{2-p}{np} S_{i,\cdot}^* S_{j,\cdot}^* + \frac{4-3p}{np} S_{i,j}^{*2} + \frac{1}{np} (\omega_i^2 S_{i,j}^* + \omega_j^2 S_{i,j}^*) \\ & + \frac{2(1-p)^2}{np^2} \left[\sum_{k=1}^d S_{i,k}^* (\mathbf{U}_k, \mathbf{U}_i^*)^2 + S_{j,k}^* (\mathbf{U}_k, \mathbf{U}_j^*)^2 \right] \\ & + \frac{1}{np^2} \sum_{k=1}^d [\omega_i^2 + (1-p) S_{i,k}^*] [\omega_j^2 + (1-p) S_{j,k}^*] (\mathbf{U}_k, \mathbf{U}_i^*)^2 \\ & + \frac{1}{np^2} \sum_{k=1}^d [\omega_j^2 + (1-p) S_{j,k}^*] [\omega_i^2 + (1-p) S_{i,k}^*] (\mathbf{U}_k, \mathbf{U}_j^*)^2 \end{aligned}$$

For any $1 \leq i \leq d$:

$$\begin{aligned} v_{i,i}^* := & \frac{12-9p}{np} S_{i,\cdot}^{*2} + \frac{4}{np} \omega_i^2 S_{i,\cdot}^* + \frac{8(1-p)^2}{np^2} \sum_{k=1}^d S_{i,k}^* (\mathbf{U}_k, \mathbf{U}_i^*)^2 \\ & + \frac{4}{np^2} \sum_{k=1}^d [\omega_i^2 + (1-p) S_{i,k}^*] [\omega_i^2 + (1-p) S_{i,k}^*] (\mathbf{U}_k, \mathbf{U}_i^*)^2 \end{aligned}$$



How to compute confidence intervals in a data-driven manner (e.g., without prior knowledge of noise levels)?

Estimating unknown model parameters

- Compute estimate (U, Λ, S) for (U^*, Λ^*, S^*) via HeteroPCA

¹ $\{y_{i,j} : (i, j) \in \Omega\}$ are zero-mean r.v.s with common variance $S_{i,i}^* + \omega_i^{*2}$

Estimating unknown model parameters

- Compute estimate (U, Λ, S) for (U^*, Λ^*, S^*) via HeteroPCA
- Estimate noise variances $\{\omega_i^{*2}\}_{i=1}^d$ in a data-driven manner¹

$$\omega_i^2 := \frac{\sum_{j=1}^n y_{i,j}^2 \mathbf{1}_{(i,j) \in \Omega}}{\sum_{j=1}^n \mathbf{1}_{(i,j) \in \Omega}} - S_{i,i}$$

¹ $\{y_{i,j} : (i,j) \in \Omega\}$ are zero-mean r.v.s with common variance $S_{i,i}^* + \omega_i^{*2}$

Estimating unknown model parameters

- Compute estimate (U, Λ, S) for (U^*, Λ^*, S^*) via HeteroPCA
- Estimate noise variances $\{\omega_i^{*2}\}_{i=1}^d$ in a data-driven manner¹

$$\omega_i^2 := \frac{\sum_{j=1}^n y_{i,j}^2 \mathbf{1}_{(i,j) \in \Omega}}{\sum_{j=1}^n \mathbf{1}_{(i,j) \in \Omega}} - S_{i,i}$$

- Compute “plug-in” estimate $v_{i,j}$ for $v_{i,j}^*$

¹ $\{y_{i,j} : (i,j) \in \Omega\}$ are zero-mean r.v.s with common variance $S_{i,i}^* + \omega_i^{*2}$

Entrywise confidence intervals for S^*

For any target coverage level $1 - \alpha$ and each (i, j) , compute

$$CI_{i,j}^{1-\alpha} := \left[S_{i,j} \pm \Phi^{-1}(1 - \alpha/2) \sqrt{v_{i,j}} \right]$$

Entrywise confidence intervals for S^*

For any target coverage level $1 - \alpha$ and each (i, j) , compute

$$\text{CI}_{i,j}^{1-\alpha} := \left[S_{i,j} \pm \Phi^{-1}(1 - \alpha/2) \sqrt{v_{i,j}} \right]$$

Theorem 4 (Yan, Chen, Fan '21)

Suppose previous conditions hold and $\frac{\omega_{\max}}{\omega_{\min}} = O(1)$. Then we have

$$\mathbb{P} \left(S_{i,j}^* \in \text{CI}_{i,j}^{1-\alpha} \right) = 1 - \alpha + o(1)$$

Entrywise confidence intervals for S^*

For any target coverage level $1 - \alpha$ and each (i, j) , compute

$$\text{CI}_{i,j}^{1-\alpha} := \left[S_{i,j} \pm \Phi^{-1}(1 - \alpha/2) \sqrt{v_{i,j}} \right]$$

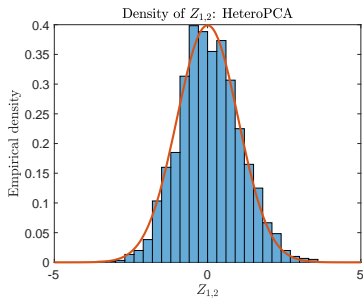
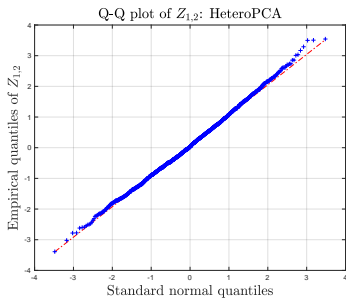
Theorem 4 (Yan, Chen, Fan '21)

Suppose previous conditions hold and $\frac{\omega_{\max}}{\omega_{\min}} = O(1)$. Then we have

$$\mathbb{P} \left(S_{i,j}^* \in \text{CI}_{i,j}^{1-\alpha} \right) = 1 - \alpha + o(1)$$

- adaptive to unknown noise levels
- adaptive to noise heteroskedasticity

Numerical verification



$n = 2000, d = 100, p = 0.6, r = 3, \omega_1^*, \dots, \omega_d^* \stackrel{\text{i.i.d.}}{\sim} \text{Unif}[0.025, 0.1],$

$$Z_{1,2} = \frac{S_{1,2} - S_{1,2}^*}{\sqrt{v_{1,2}}}$$

Concluding remarks

- Missing data and heterogeneous noise require special treatment
- HeteroPCA is provably effective for estimation & inference
minimax optimal in some sense

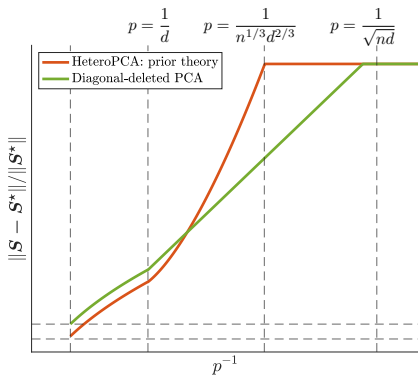
papers:

Y. Yan, Y. Chen, J. Fan, “Inference for Heteroskedastic PCA with Missing Data,” arxiv:2107.12365, 2021

C. Cai, G. Li, Y. Chi, H. V. Poor, Y. Chen, “Subspace Estimation from Unbalanced and Incomplete Data Matrices: $\ell_{2,\infty}$ Statistical Guarantees,” *Annals of Statistics*, 2021

Backup slides

prior theory
(noiseless, $n > d$)



	$\ \cdot\ $ estimation error bounds	min sample size requirement
HeteroPCA (Zhang et al. '18)	$\frac{1}{\sqrt{nd^2p^3}} + \frac{1}{\sqrt{np}}$	$n^{\frac{2}{3}}d^{\frac{1}{3}}$
diagonal-deleted PCA (Cai et al. 19)	$\frac{1}{\sqrt{ndp^2}} + \frac{1}{\sqrt{np}} + \frac{1}{d}$	\sqrt{nd}