

High-dimensional principal component analysis with heterogeneous missingness

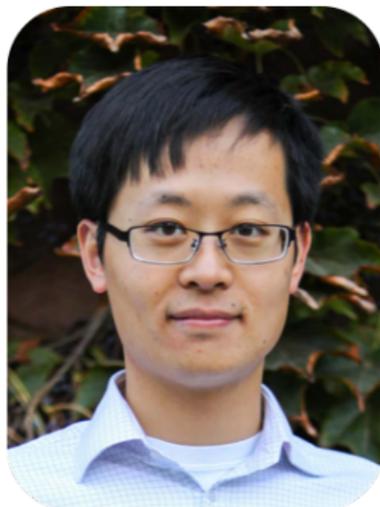
Tengyao Wang

University College London

Statistics Seminar, University of Kent

24 Oct 2019

Collaborators



Ziwei Zhu



Richard Samworth

Missing data

Missing data are ever more important in the Big Data era, because complete-case analysis is less feasible.

Consider a complete-case analysis with an $n \times d$ matrix, where each entry is observed independently with probability $p = 0.99$.

Missing data

Missing data are ever more important in the Big Data era, because complete-case analysis is less feasible.

Consider a complete-case analysis with an $n \times d$ matrix, where each entry is observed independently with probability $p = 0.99$.

- ▶ When $d = 5$, around 95% of observations are retained

Missing data

Missing data are ever more important in the Big Data era, because complete-case analysis is less feasible.

Consider a complete-case analysis with an $n \times d$ matrix, where each entry is observed independently with probability $p = 0.99$.

- ▶ When $d = 5$, around 95% of observations are retained
- ▶ When $d = 300$, only around 5% of observations are retained.

Missing data

Approaches to handle missing data:

- ▶ Imputation (Ford, 1983; Rubin, 2004)
- ▶ Factored likelihood (Anderson, 1957)
- ▶ Expectation-Maximisation (Dempster et al., 1977)

Missing data

Approaches to handle missing data:

- ▶ Imputation (Ford, 1983; Rubin, 2004)
- ▶ Factored likelihood (Anderson, 1957)
- ▶ Expectation-Maximisation (Dempster et al., 1977)

Recently, there has been increased emphasis on missing data in high-dimensional problems:

- ▶ Sparse regression (Loh and Wainwright, 2012; Belloni et al., 2017)
- ▶ Classification (Cai and Zhang, 2018b)
- ▶ Covariance and precision matrix estimation (Lounici, 2014; Loh and Tan, 2018)

High-dimensional PCA problem set-up

Suppose the (partially observed) matrix $\mathbf{Y} \in \mathbb{R}^{n \times d}$ is of the form

$$\mathbf{Y} = \mathbf{U}\mathbf{V}_K^\top + \mathbf{Z},$$

where $\mathbf{V}_K \in \mathbb{R}^{d \times K}$ has orthonormal columns and \mathbf{U} is a random $n \times K$ matrix (with $n > K$) having i.i.d. rows with mean zero.

High-dimensional PCA problem set-up

Suppose the (partially observed) matrix $\mathbf{Y} \in \mathbb{R}^{n \times d}$ is of the form

$$\mathbf{Y} = \mathbf{U}\mathbf{V}_K^\top + \mathbf{Z},$$

where $\mathbf{V}_K \in \mathbb{R}^{d \times K}$ has orthonormal columns and \mathbf{U} is a random $n \times K$ matrix (with $n > K$) having i.i.d. rows with mean zero.

Let $\Omega_{ij} := \{Y_{ij} \text{ is observed}\}$ and $\mathbf{Y}_\Omega := \mathbf{Y} \circ \Omega$. We observe the pair $(\mathbf{Y}_\Omega, \Omega)$ and wish to estimate $\text{Col}(\mathbf{V}_K)$.

High-dimensional PCA problem set-up

Suppose the (partially observed) matrix $\mathbf{Y} \in \mathbb{R}^{n \times d}$ is of the form

$$\mathbf{Y} = \mathbf{U}\mathbf{V}_K^\top + \mathbf{Z},$$

where $\mathbf{V}_K \in \mathbb{R}^{d \times K}$ has orthonormal columns and \mathbf{U} is a random $n \times K$ matrix (with $n > K$) having i.i.d. rows with mean zero.

Let $\Omega_{ij} := \{Y_{ij} \text{ is observed}\}$ and $\mathbf{Y}_\Omega := \mathbf{Y} \circ \Omega$. We observe the pair $(\mathbf{Y}_\Omega, \Omega)$ and wish to estimate $\text{Col}(\mathbf{V}_K)$.

Performance of an estimator $\widehat{\mathbf{V}}_K$ measured by the loss function

$$L(\widehat{\mathbf{V}}_K, \mathbf{V}_K) := \|\sin \Theta(\widehat{\mathbf{V}}_K, \mathbf{V}_K)\|_F,$$

where $\Theta(\mathbf{U}, \mathbf{V})$ is the matrix of principal angles between $\text{Col}(\mathbf{U})$ and $\text{Col}(\mathbf{V})$.

Inverse-Probability Weighted estimator

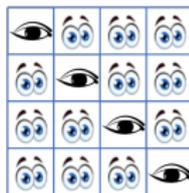
Consider the p -homogeneous setting, where $\Omega_{ij} \stackrel{\text{iid}}{\sim} \text{Bern}(p)$. Then

$$\mathbf{P} := \mathbb{E}(\boldsymbol{\omega}_1 \boldsymbol{\omega}_1^\top) = p^2 \{ \mathbf{1}_d \mathbf{1}_d^\top - (1 - p^{-1}) \mathbf{I}_d \}.$$

Its elementwise inverse is $\mathbf{W} := p^{-2} \{ \mathbf{1}_d \mathbf{1}_d^\top - (1 - p) \mathbf{I}_d \}$, and we can define the weighted sample covariance matrix

$$\mathbf{G} := \left(\frac{1}{n} \mathbf{Y}_\Omega^\top \mathbf{Y}_\Omega \right) \circ \mathbf{W}.$$

This ensures that $\mathbb{E}(\mathbf{G} \mid \mathbf{Y}) = n^{-1} \mathbf{Y}^\top \mathbf{Y}$.



Inverse-Probability Weighted estimator

Consider the p -homogeneous setting, where $\Omega_{ij} \stackrel{\text{iid}}{\sim} \text{Bern}(p)$. Then

$$\mathbf{P} := \mathbb{E}(\boldsymbol{\omega}_1 \boldsymbol{\omega}_1^\top) = p^2 \{ \mathbf{1}_d \mathbf{1}_d^\top - (1 - p) \mathbf{I}_d \}.$$

Its elementwise inverse is $\mathbf{W} := p^{-2} \{ \mathbf{1}_d \mathbf{1}_d^\top - (1 - p) \mathbf{I}_d \}$, and we can define the weighted sample covariance matrix

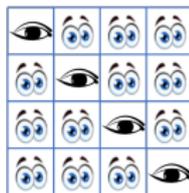
$$\mathbf{G} := \left(\frac{1}{n} \mathbf{Y}_\Omega^\top \mathbf{Y}_\Omega \right) \circ \mathbf{W}.$$

This ensures that $\mathbb{E}(\mathbf{G} \mid \mathbf{Y}) = n^{-1} \mathbf{Y}^\top \mathbf{Y}$. Define $\widehat{\mathbf{W}}$ by replacing p in \mathbf{W} with $\widehat{p} := (nd)^{-1} \|\boldsymbol{\Omega}\|_1$, and set

$$\widehat{\mathbf{G}} := \left(\frac{1}{n} \mathbf{Y}_\Omega^\top \mathbf{Y}_\Omega \right) \circ \widehat{\mathbf{W}}.$$

The IPW estimator of \mathbf{V}_K is given by the top K eigenvectors of $\widehat{\mathbf{G}}$, denoted $\widehat{\mathbf{V}}_K$

(Cai and Zhang, 2018a; Cho, Kim and Rohe, 2017).



Assumptions

For $r \in \mathbb{N}$ and a d -dimensional random vector \mathbf{x} , define its Orlicz norm

$$\|\mathbf{x}\|_{\psi_r} := \sup_{\mathbf{u} \in \mathcal{S}^{d-1}} \sup_{q \in \mathbb{N}} \frac{(\mathbb{E}|\mathbf{u}^\top \mathbf{x}|^q)^{1/q}}{q^{1/r}}$$

and a version that is invariant to invertible affine transformations:

$$\|\mathbf{x}\|_{\psi_r^*} := \sup_{\mathbf{u} \in \mathcal{S}^{d-1}} \frac{\|\mathbf{u}^\top (\mathbf{x} - \mathbb{E}\mathbf{x})\|_{\psi_r}}{\text{Var}^{1/2}(\mathbf{u}^\top \mathbf{x})}.$$

Assumptions

For $r \in \mathbb{N}$ and a d -dimensional random vector \mathbf{x} , define its Orlicz norm

$$\|\mathbf{x}\|_{\psi_r} := \sup_{\mathbf{u} \in \mathcal{S}^{d-1}} \sup_{q \in \mathbb{N}} \frac{(\mathbb{E}|\mathbf{u}^\top \mathbf{x}|^q)^{1/q}}{q^{1/r}}$$

and a version that is invariant to invertible affine transformations:

$$\|\mathbf{x}\|_{\psi_r^*} := \sup_{\mathbf{u} \in \mathcal{S}^{d-1}} \frac{\|\mathbf{u}^\top (\mathbf{x} - \mathbb{E}\mathbf{x})\|_{\psi_r}}{\text{Var}^{1/2}(\mathbf{u}^\top \mathbf{x})}.$$

- (A1) \mathbf{U} , \mathbf{Z} and $\mathbf{\Omega}$ are independent;
- (A2) $\|\mathbf{u}_1\|_{\psi_2^*} \leq \tau$;
- (A3) $\mathbf{Z} = (z_{ij})_{i \in [n], j \in [d]}$ has i.i.d. entries with $\mathbb{E}z_{11} = 0$, $\text{Var} z_{11} = 1$ and $\|z_{11}\|_{\psi_2^*} \leq \tau$;
- (A4) $\|y_{1j}^2\|_{\psi_1} \leq M$ for all $j \in [d]$.

Upper bound

Theorem. Assume (A1)–(A4) and that $n, d \geq 2, dp \geq 1$. Let λ_j denote the j th largest eigenvalue of Σ_u . If $n \geq d \log^2 d \log^2 n / (\lambda_1 p + \log d)$, then

$$\mathbb{E}L(\widehat{\mathbf{V}}_K, \mathbf{V}_K) \lesssim_{M, \tau} \frac{1}{\lambda_K} \left(\frac{Kd(\lambda_1 p + \log d) \log^2 d}{np^2} \right)^{1/2}.$$

Upper bound

Theorem. Assume (A1)–(A4) and that $n, d \geq 2, dp \geq 1$. Let λ_j denote the j th largest eigenvalue of Σ_u . If $n \geq d \log^2 d \log^2 n / (\lambda_1 p + \log d)$, then

$$\mathbb{E}L(\widehat{\mathbf{V}}_K, \mathbf{V}_K) \lesssim_{M, \tau} \frac{1}{\lambda_K} \left(\frac{Kd(\lambda_1 p + \log d) \log^2 d}{np^2} \right)^{1/2}.$$

The sample size requirement is reasonable: with no missing data and when $\lambda_1 \gg 1$, the top eigenvector of the sample covariance matrix estimator is consistent if and only if $d/(n\lambda_1) \rightarrow 0$ (Shen et al., 2016).

Upper bound

Theorem. Assume (A1)–(A4) and that $n, d \geq 2, dp \geq 1$. Let λ_j denote the j th largest eigenvalue of Σ_u . If $n \geq d \log^2 d \log^2 n / (\lambda_1 p + \log d)$, then

$$\mathbb{E}L(\widehat{\mathbf{V}}_K, \mathbf{V}_K) \lesssim_{M, \tau} \frac{1}{\lambda_K} \left(\frac{Kd(\lambda_1 p + \log d) \log^2 d}{np^2} \right)^{1/2}.$$

The sample size requirement is reasonable: with no missing data and when $\lambda_1 \gg 1$, the top eigenvector of the sample covariance matrix estimator is consistent if and only if $d/(n\lambda_1) \rightarrow 0$ (Shen et al., 2016).

The theorem reveals a phase transition depending on the relative magnitudes of $\lambda_1 p$ and $\log d$. In particular,

$$\mathbb{E}L(\widehat{\mathbf{V}}_K, \mathbf{V}_K) \lesssim_{M, \tau} \begin{cases} \frac{1}{\lambda_K} \left(\frac{Kd \log^3 d}{np^2} \right)^{1/2} & \text{if } \lambda_1 p \lesssim \log d, \\ \frac{\lambda_1^{1/2}}{\lambda_K} \left(\frac{Kd \log^2 d}{np} \right)^{1/2} & \text{if } \lambda_1 p \gtrsim \log d. \end{cases}$$

Minimax lower bound

Let $\mathcal{P}_{n,d}(\lambda_1, p)$ denote the class of distributions of pairs $(\mathbf{Y}_\Omega, \Omega)$ satisfying (A1), (A2), (A3) with $K = 1$. Since we are now working with vectors instead of matrices, we write \mathbf{v} in place of \mathbf{V}_1 .

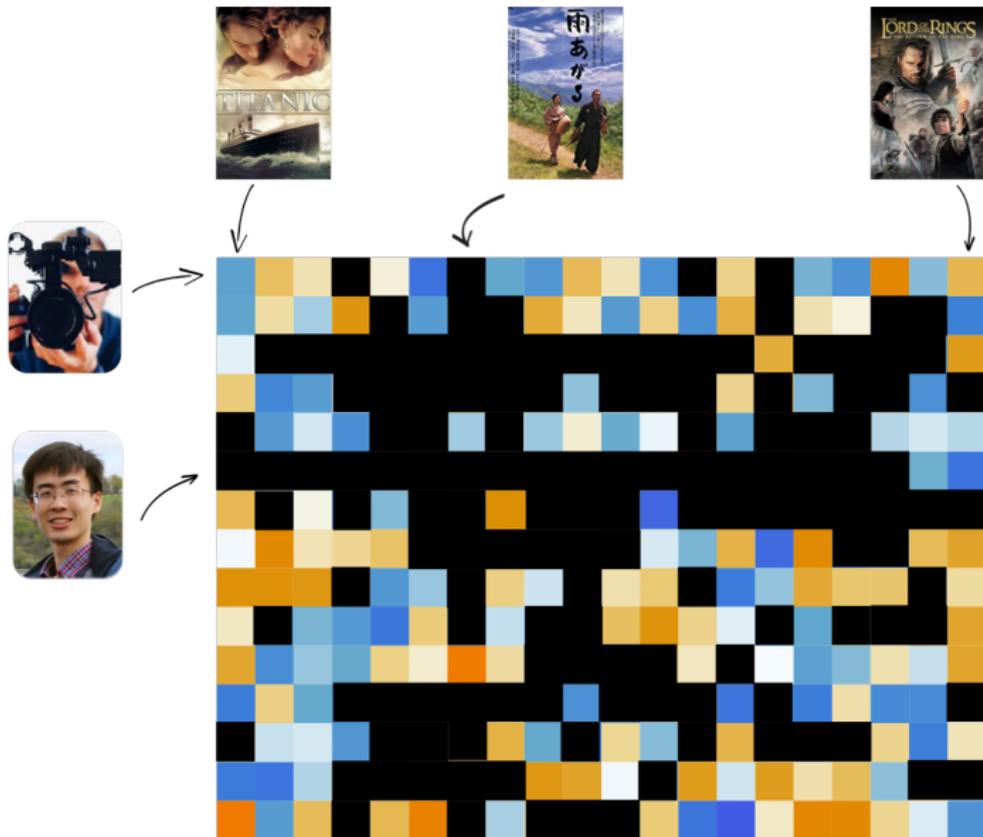
Theorem. There exists a universal constant $c > 0$ such that

$$\inf_{\hat{\mathbf{v}}} \sup_{P \in \mathcal{P}_{n,d}(\lambda_1, p)} \mathbb{E}_P L(\hat{\mathbf{v}}, \mathbf{v}) \geq c \min \left\{ \frac{1}{\lambda_1} \left(\frac{d(\lambda_1 p + 1)}{np^2} \right)^{1/2}, 1 \right\},$$

where the infimum is taken over all estimators $\hat{\mathbf{v}} = \hat{\mathbf{v}}(\mathbf{Y}_\Omega, \Omega)$ of \mathbf{v} .

Thus $\hat{\mathbf{V}}_1$ achieves the minimax optimal rate of estimation up to a poly-logarithmic factor when M and τ are regarded as constants and $K = 1$.

General observation mechanisms



General observation mechanisms: example

Suppose that

$$\mathbb{P}\{\boldsymbol{\omega}_1 = (1, 0, 1, \dots, 1)^\top\} = \mathbb{P}\{\boldsymbol{\omega}_1 = (0, 1, 1, \dots, 1)^\top\} = 1/2.$$

Consider $\boldsymbol{\Sigma} = \mathbf{I}_d + \boldsymbol{\alpha}\boldsymbol{\alpha}^\top$, where $\boldsymbol{\alpha} = (2^{-1/2}, 2^{-1/2}, 0, \dots, 0)^\top \in \mathbb{R}^d$, and $\boldsymbol{\Sigma}' = \mathbf{I}_d + \boldsymbol{\alpha}'(\boldsymbol{\alpha}')^\top$, where $\boldsymbol{\alpha}' = (2^{-1/2}, -2^{-1/2}, 0, \dots, 0)^\top \in \mathbb{R}^d$.

Suppose that $\mathbf{y} \sim N_d(\mathbf{0}, \boldsymbol{\Sigma})$ and $\mathbf{y}' \sim N_d(\mathbf{0}, \boldsymbol{\Sigma}')$. Then $(\mathbf{y} \circ \boldsymbol{\omega}, \boldsymbol{\omega})$ and $(\mathbf{y}' \circ \boldsymbol{\omega}, \boldsymbol{\omega})$ are identically distributed.

General observation mechanisms: example

Suppose that

$$\mathbb{P}\{\boldsymbol{\omega}_1 = (1, 0, 1, \dots, 1)^\top\} = \mathbb{P}\{\boldsymbol{\omega}_1 = (0, 1, 1, \dots, 1)^\top\} = 1/2.$$

Consider $\boldsymbol{\Sigma} = \mathbf{I}_d + \boldsymbol{\alpha}\boldsymbol{\alpha}^\top$, where $\boldsymbol{\alpha} = (2^{-1/2}, 2^{-1/2}, 0, \dots, 0)^\top \in \mathbb{R}^d$, and $\boldsymbol{\Sigma}' = \mathbf{I}_d + \boldsymbol{\alpha}'(\boldsymbol{\alpha}')^\top$, where $\boldsymbol{\alpha}' = (2^{-1/2}, -2^{-1/2}, 0, \dots, 0)^\top \in \mathbb{R}^d$.

Suppose that $\mathbf{y} \sim N_d(\mathbf{0}, \boldsymbol{\Sigma})$ and $\mathbf{y}' \sim N_d(\mathbf{0}, \boldsymbol{\Sigma}')$. Then $(\mathbf{y} \circ \boldsymbol{\omega}, \boldsymbol{\omega})$ and $(\mathbf{y}' \circ \boldsymbol{\omega}, \boldsymbol{\omega})$ are identically distributed.

But the respective leading eigenvectors of $\boldsymbol{\Sigma}$ and $\boldsymbol{\Sigma}'$ are $\boldsymbol{\alpha}$ and $\boldsymbol{\alpha}'$, which are orthogonal!

primePCA: a single iteration of refinement

primePCA (projected refinement for imputation of missing entries in PCA) iteratively refines a warm initialiser. We write $\tilde{\mathbf{y}}_i := \mathbf{y}_i \circ \boldsymbol{\omega}_i$.

Algorithm 1 $\text{refine}(K, \widehat{\mathbf{V}}_K^{(\text{in})}, \boldsymbol{\Omega}, \mathbf{Y}_\Omega)$, a single step of refinement of current iterate $\widehat{\mathbf{V}}_K^{(\text{in})}$

Input: $K \in [d]$, $\widehat{\mathbf{V}}_K^{(\text{in})} \in \mathbb{O}^{d \times K}$, $\boldsymbol{\Omega} \in \{0, 1\}^{n \times d}$ with $\min_i \|\boldsymbol{\omega}_i\|_1 \geq 1$, $\mathbf{Y}_\Omega \in \mathbb{R}^{n \times d}$

Output: $\widehat{\mathbf{V}}_K^{(\text{out})} \in \mathbb{O}^{d \times K}$

- 1: **for** i in $[n]$ **do**
 - 2: $\mathcal{J}_i \leftarrow \{j \in [d] : \omega_{ij} = 1\}$
 - 3: $\hat{\mathbf{u}}_i \leftarrow (\widehat{\mathbf{V}}_K^{(\text{in})})_{\mathcal{J}_i}^\dagger \tilde{\mathbf{y}}_{i, \mathcal{J}_i}$
 - 4: $\hat{\mathbf{y}}_{i, \mathcal{J}_i^c} \leftarrow \widehat{\mathbf{V}}_K^{(\text{in})} \hat{\mathbf{u}}_{i, \mathcal{J}_i^c}$
 - 5: $\hat{\mathbf{y}}_{i, \mathcal{J}_i} \leftarrow \mathbf{y}_{i, \mathcal{J}_i}$
 - 6: **end for**
 - 7: $\hat{\mathbf{Y}} \leftarrow (\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_n)^\top$
 - 8: $\widehat{\mathbf{V}}_K^{(\text{out})} \leftarrow$ top K right singular vectors of $\hat{\mathbf{Y}}$
-

Two-to-infinity subspace distance

For $U, V \in \mathbb{O}^{d \times K}$, let $W_1 D_{U,V} W_2^\top$ be an SVD of $V^\top U$ and let $W_{U,V} := W_1 W_2^\top$. Then $W_{U,V}$ solves the Procrustes problem in the sense that

$$W_{U,V} \in \arg \min_{W \in \mathbb{O}^{K \times K}} \|U - VW\|_F.$$

The two-to-infinity distance between $\text{Col}(U)$ and $\text{Col}(V)$ is then defined to be

$$\mathcal{T}(U, V) := \|U - VW_{U,V}\|_{2 \rightarrow \infty},$$

where $\|A\|_{2 \rightarrow \infty} := \sup_{\mathbf{x}: \|\mathbf{x}\|_2=1} \|A\mathbf{x}\|_\infty$.

Contraction of each iteration

Proposition. Let $\widehat{\mathbf{V}}_K^{(\text{out})} := \text{refine}(K, \widehat{\mathbf{V}}_K^{(\text{in})}, \mathbf{\Omega}, \mathbf{Y}_\Omega)$. We assume that $\min_{i \in [n]} \|\boldsymbol{\omega}_i\|_1 > K$, that $\min_{i \in [n]} \frac{d^{1/2} \sigma_K((\widehat{\mathbf{V}}_K^{(\text{in})})_{\mathcal{J}_i})}{|\mathcal{J}_i|^{1/2}} \geq 1/\sigma_* > 0$, and write the SVD of \mathbf{Y} as $\mathbf{L}\mathbf{\Gamma}\mathbf{R}^\top$.

Suppose that $\mathbf{Z} = \mathbf{0}$, and that both $\|\mathbf{L}\|_{2 \rightarrow \infty} \leq \mu_1(K/n)^{1/2}$ and $\|\mathbf{R}\|_{2 \rightarrow \infty} \leq \mu_2(K/d)^{1/2}$ hold for some $\mu_1, \mu_2 \geq 1$. Then there exist $c_1, C > 0$, depending only on μ_1, μ_2 and σ_* , such that whenever

$$(i) \quad \mathcal{T}(\widehat{\mathbf{V}}_K^{(\text{in})}, \mathbf{V}_K) \leq \frac{c_1 \sigma_K(\mathbf{\Gamma})}{K^2 \sigma_1(\mathbf{\Gamma}) \sqrt{d}},$$

$$(ii) \quad \rho := \frac{CK^2 \sigma_1(\mathbf{\Gamma}) \|\mathbf{\Omega}^c\|_{1 \rightarrow 1}}{\sigma_K(\mathbf{\Gamma}) n} < 1,$$

we have that

$$\mathcal{T}(\widehat{\mathbf{V}}_K^{(\text{out})}, \mathbf{V}_K) \leq \rho \mathcal{T}(\widehat{\mathbf{V}}_K^{(\text{in})}, \mathbf{V}_K).$$

Algorithm 2 primePCA, an iterative algorithm for estimating \mathbf{V}_K given initialiser $\widehat{\mathbf{V}}_K^{(0)}$

Input: $K \in [d]$, $\widehat{\mathbf{V}}_K^{(0)} \in \mathbb{O}^{d \times K}$, $\mathbf{\Omega} \in \{0, 1\}^{n \times d}$, $\mathbf{Y}_{\mathbf{\Omega}} \in \mathbb{R}^{n \times d}$, $n_{\text{iter}} \in \mathbb{N}$, $\sigma_* \in (0, \infty)$, $\kappa^* \in [0, \infty)$

Output: $\widehat{\mathbf{V}}_K \in \mathbb{R}^{d \times K}$

```

1: for  $i$  in  $[n]$  do
2:    $\mathcal{J}_i \leftarrow \{j \in [d] : \omega_{ij} = 1\}$ 
3: end for
4: for  $t$  in  $[n_{\text{iter}}]$  do
5:    $\mathcal{I}^{(t-1)} \leftarrow \{i : \|\boldsymbol{\omega}_i\|_1 > K, \sigma_K((\widehat{\mathbf{V}}_K^{(t-1)})_{\mathcal{J}_i}) \geq \frac{|\mathcal{J}_i|^{1/2}}{d^{1/2}\sigma_*}\}$ 
6:    $\widehat{\mathbf{V}}_K^{(t)} \leftarrow \text{refine}(K, \widehat{\mathbf{V}}_K^{(t-1)}, \mathbf{\Omega}_{\mathcal{I}^{(t-1)}}, (\mathbf{Y}_{\mathbf{\Omega}})_{\mathcal{I}^{(t-1)}})$  # refine is defined in Algorithm 1.
7:   if  $L(\widehat{\mathbf{V}}_K^{(t)}, \widehat{\mathbf{V}}_K^{(t-1)}) < \kappa^*$  then break
8:   end if
9: end for
10: return  $\widehat{\mathbf{V}}_K = \widehat{\mathbf{V}}_K^{(t)}$ 

```

Geometric convergence in noiseless case

Theorem. For $t \in [n_{\text{iter}}]$, let $\widehat{\mathbf{V}}_K^{(t)}$ be the t^{th} iterate of Algorithm 2 with input K , $\widehat{\mathbf{V}}_K^{(0)}$, $\boldsymbol{\Omega} \in \{0, 1\}^{n \times d}$, $\mathbf{Y}_{\boldsymbol{\Omega}} \in \mathbb{R}^{n \times d}$, $n_{\text{iter}} \in \mathbb{N}$, $\sigma_* \in (0, \infty)$ and $\kappa^* = 0$. Let

$$\mathcal{I} := \left\{ i : \|\boldsymbol{\omega}_i\|_1 > K, \sigma_K((\mathbf{V}_K)_{\mathcal{J}_i}) \geq |\mathcal{J}_i|^{1/2} / (d^{1/2} \sigma_*) \right\},$$

where $\mathcal{J}_i := \{j : \omega_{ij} = 1\}$. Let $\mathbf{Y}_{\mathcal{I}} = \mathbf{L}\boldsymbol{\Gamma}\mathbf{R}^\top$ be an SVD of $\mathbf{Y}_{\mathcal{I}}$. Suppose that both $\|\mathbf{L}\|_{2 \rightarrow \infty} \leq \mu_1(K/|\mathcal{I}|)^{1/2}$ and $\|\mathbf{R}\|_{2 \rightarrow \infty} \leq \mu_2(K/d)^{1/2}$. Let

$$\mathcal{Z} := \left\{ \sigma_K((\mathbf{V}_K)_{\mathcal{J}_i}) d^{1/2} / |\mathcal{J}_i|^{1/2} : i \in [n], \|\boldsymbol{\omega}_i\|_1 > K \right\},$$

and assume that $\epsilon := \min_{z \in \mathcal{Z}} |z - \sigma_*^{-1}| > 0$. Then there exist $c_1, C > 0$, depending only on μ_1, μ_2, σ_* and ϵ , such that whenever

$$\mathcal{T}(\widehat{\mathbf{V}}_K^{(0)}, \mathbf{V}_K) \leq \frac{c_1 \sigma_K(\mathbf{Y}_{\mathcal{I}})}{K^2 \sigma_1(\mathbf{Y}_{\mathcal{I}}) \sqrt{d}} \quad \text{and} \quad \rho := \frac{CK^2 \sigma_1(\mathbf{Y}_{\mathcal{I}}) \|\boldsymbol{\Omega}_{\mathcal{I}}^c\|_{1 \rightarrow 1}}{\sigma_K(\mathbf{Y}_{\mathcal{I}}) |\mathcal{I}|} < 1,$$

we have $\mathcal{T}(\widehat{\mathbf{V}}_K^{(t)}, \mathbf{V}_K) \leq \rho^t \mathcal{T}(\widehat{\mathbf{V}}_K^{(0)}, \mathbf{V}_K)$ for every $t \in [n_{\text{iter}}]$.

Initialisation

Consider the following modified weighted sample covariance matrix

$$\tilde{\mathbf{G}} := \frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^\top \circ \tilde{\mathbf{W}},$$

where for any $j, k \in [d]$,

$$\tilde{\mathbf{W}}_{jk} := \begin{cases} \frac{n}{\sum_{i=1}^n \omega_{ij} \omega_{ik}} & \text{if } \sum_{i=1}^n \omega_{ij} \omega_{ik} > 0, \\ 0, & \text{otherwise.} \end{cases}$$

We take as our initial estimator of \mathbf{V}_K the matrix of top K eigenvectors of $\tilde{\mathbf{G}}$, denoted $\tilde{\mathbf{V}}_K$.

Performance of initialiser

Proposition. Assume the same conditions as in the previous theorem. Then there exists a universal constant $C > 0$ such that for any $\xi > 1$, if

$$\lambda_K > C \left\{ \left(\frac{M\tau^2 R \|\widetilde{\mathbf{W}}\|_{1 \rightarrow 1} \xi \log d}{n} \right)^{1/2} + \frac{M \|\widetilde{\mathbf{W}}\|_{\text{op}} \xi \log^2 d}{n} \right\},$$

then with \mathbb{P}^Ω -probability at least $1 - (2K + 4)d^{-(\xi-1)}$, we have

$$L(\widetilde{\mathbf{V}}_K, \mathbf{V}_K) \leq \frac{2^{9/2} e \tau \mu}{\lambda_K} \left(\frac{KMR}{d} \right)^{1/2} \left(\frac{\xi^{1/2} \|\widetilde{\mathbf{W}}\|_1^{1/2} \log^{1/2} d}{n^{1/2}} + \frac{\xi \|\widetilde{\mathbf{W}}\|_{\text{F}} \log d}{n} \right).$$

Performance of initialiser

Proposition. Assume the same conditions as in the previous theorem. Then there exists a universal constant $C > 0$ such that for any $\xi > 1$, if

$$\lambda_K > C \left\{ \left(\frac{M\tau^2 R \|\widetilde{\mathbf{W}}\|_{1 \rightarrow 1} \xi \log d}{n} \right)^{1/2} + \frac{M \|\widetilde{\mathbf{W}}\|_{\text{op}} \xi \log^2 d}{n} \right\},$$

then with \mathbb{P}^Ω -probability at least $1 - (2K + 4)d^{-(\xi-1)}$, we have

$$L(\widetilde{\mathbf{V}}_K, \mathbf{V}_K) \leq \frac{2^{9/2} e \tau \mu}{\lambda_K} \left(\frac{KMR}{d} \right)^{1/2} \left(\frac{\xi^{1/2} \|\widetilde{\mathbf{W}}\|_1^{1/2} \log^{1/2} d}{n^{1/2}} + \frac{\xi \|\widetilde{\mathbf{W}}\|_{\text{F}} \log d}{n} \right).$$

N.B. The bound depends on $\widetilde{\mathbf{W}}$ only through the *entrywise* ℓ_1 and ℓ_2 norms of the whole matrix.

Simulations: Noiseless case

Fix $n = 2000$, $d = 500$, $K = 2$ and $\mathbf{u}_i \sim N_d(0, \Sigma_u)$ where $\Sigma_u = 100\mathbf{I}_2$. Set

$$\mathbf{V}_K = \sqrt{\frac{1}{500}} \begin{pmatrix} \mathbf{1}_{250} & \mathbf{1}_{250} \\ \mathbf{1}_{250} & -\mathbf{1}_{250} \end{pmatrix} \in \mathbb{R}^{500 \times 2}.$$

Simulations: Noiseless case

Fix $n = 2000$, $d = 500$, $K = 2$ and $\mathbf{u}_i \sim N_d(0, \Sigma_u)$ where $\Sigma_u = 100\mathbf{I}_2$. Set

$$\mathbf{V}_K = \sqrt{\frac{1}{500}} \begin{pmatrix} \mathbf{1}_{250} & \mathbf{1}_{250} \\ \mathbf{1}_{250} & -\mathbf{1}_{250} \end{pmatrix} \in \mathbb{R}^{500 \times 2}.$$

- (H1) Homogeneous: $\mathbb{P}(\omega_{ij} = 1) = 0.05$ for all $i \in [n], j \in [d]$;
- (H2) Mildly heterogeneous: $\mathbb{P}(\omega_{ij} = 1) = P_i Q_j$ for $i \in [n], j \in [d]$, where $P_1, \dots, P_n \stackrel{\text{iid}}{\sim} U[0, 0.2]$ and $Q_1, \dots, Q_d \stackrel{\text{iid}}{\sim} U[0.05, 0.95]$ independently;
- (H3) Highly heterogeneous columns: $\mathbb{P}(\omega_{ij} = 1) = 0.19$ for $i \in [n]$ and all odd $j \in [d]$ and $\mathbb{P}(\omega_{ij} = 1) = 0.01$ for $i \in [n]$ and all even $j \in [d]$.
- (H4) Highly heterogeneous rows: $\mathbb{P}(\omega_{ij} = 1) = 0.18$ for $j \in [d]$ and all odd $i \in [n]$ and $\mathbb{P}(\omega_{ij} = 1) = 0.02$ for $j \in [d]$ and all even $i \in [n]$.

Simulations: Noiseless case

Fix $n = 2000$, $d = 500$, $K = 2$ and $\mathbf{u}_i \sim N_d(0, \Sigma_u)$ where $\Sigma_u = 100\mathbf{I}_2$. Set

$$\mathbf{V}_K = \sqrt{\frac{1}{500}} \begin{pmatrix} \mathbf{1}_{250} & \mathbf{1}_{250} \\ \mathbf{1}_{250} & -\mathbf{1}_{250} \end{pmatrix} \in \mathbb{R}^{500 \times 2}.$$

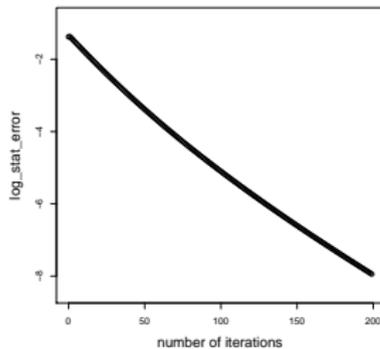
- (H1) Homogeneous: $\mathbb{P}(\omega_{ij} = 1) = 0.05$ for all $i \in [n], j \in [d]$;
- (H2) Mildly heterogeneous: $\mathbb{P}(\omega_{ij} = 1) = P_i Q_j$ for $i \in [n], j \in [d]$, where $P_1, \dots, P_n \stackrel{\text{iid}}{\sim} U[0, 0.2]$ and $Q_1, \dots, Q_d \stackrel{\text{iid}}{\sim} U[0.05, 0.95]$ independently;
- (H3) Highly heterogeneous columns: $\mathbb{P}(\omega_{ij} = 1) = 0.19$ for $i \in [n]$ and all odd $j \in [d]$ and $\mathbb{P}(\omega_{ij} = 1) = 0.01$ for $i \in [n]$ and all even $j \in [d]$.
- (H4) Highly heterogeneous rows: $\mathbb{P}(\omega_{ij} = 1) = 0.18$ for $j \in [d]$ and all odd $i \in [n]$ and $\mathbb{P}(\omega_{ij} = 1) = 0.02$ for $j \in [d]$ and all even $i \in [n]$.

Compare with `softImpute`: fix $\lambda > 0$ and take the top K eigenvectors of

$$\hat{\mathbf{Y}}^{\text{soft}} := \arg \min_{\mathbf{X} \in \mathbb{R}^{n \times d}} \left\{ \frac{1}{2} \|\mathbf{Y}_\Omega - \mathbf{X}_\Omega\|_F^2 + \lambda \|\mathbf{X}\|_* \right\}$$

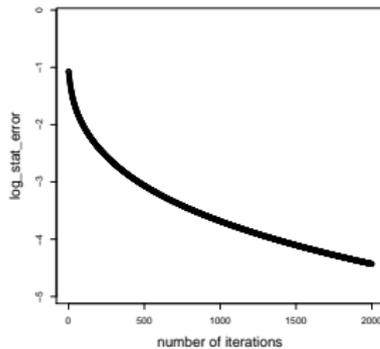
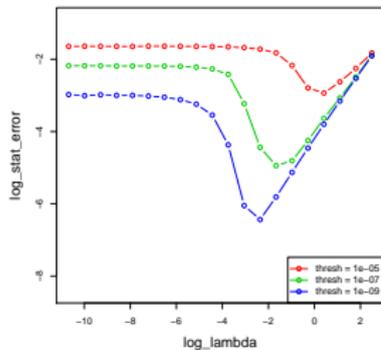
Noiseless case

primePCA

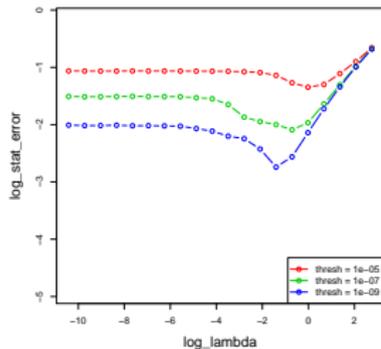


(H1)

softImpute

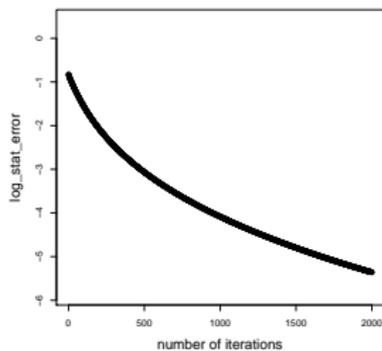


(H2)

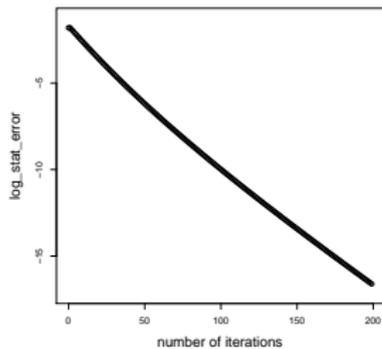


Noiseless case

primePCA

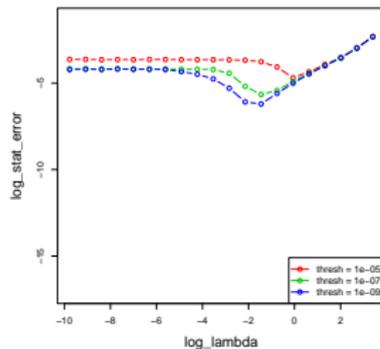
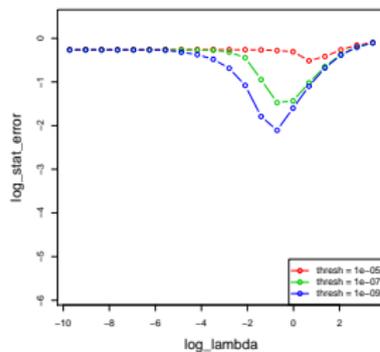


(H3)



(H4)

softImpute



Simulations: Noisy case

Now generate $\mathbf{z}_i \sim N_d(\mathbf{0}, \mathbf{I}_d)$, independent of all other data, and set $\Sigma_{\mathbf{u}} = \nu^2 \mathbf{I}_2$ where $\nu \in \{20, 40, 60\}$, corresponding to

$$\text{SNR} := \text{tr Cov}(\mathbf{y}_1) / \text{tr Cov}(\mathbf{z}_1) \in \{1.6, 6.4, 14.4\}.$$

Also compare with `hardImpute` (Mazumder, Hastie and Tibshirani, 2010), which retains only a fixed number of top singular values in each iteration of matrix imputation; i.e. `softImpute` with $\lambda = 0$.

For `softImpute`, use oracle choice of λ for each repetition.

Simulations: Noisy case

		$\nu = 20$	$\nu = 40$	$\nu = 60$
(H1)	hardImpute	0.444 _(0.001)	0.251 _(0.001)	0.186 _(0.0005)
	softImpute(oracle)	0.186 _(0.0004)	0.095 _(0.0002)	0.064 _(0.0002)
	primePCA_init	0.306 _(0.001)	0.266 _(0.001)	0.259 _(0.001)
	primePCA	0.171 _(0.0004)	0.084 _(0.0002)	0.056 _(0.0001)
(H2)	hardImpute	0.473 _(0.001)	0.291 _(0.001)	0.236 _(0.001)
	softImpute(oracle)	0.308 _(0.001)	0.185 _(0.001)	0.141 _(0.001)
	primePCA_init	0.399 _(0.002)	0.357 _(0.001)	0.349 _(0.001)
	primePCA	0.232 _(0.001)	0.115 _(0.001)	0.077 _(0.0005)
(H3)	hardImpute	0.479 _(0.001)	0.385 _(0.001)	0.427 _(0.001)
	softImpute(oracle)	0.374 _(0.001)	0.222 _(0.001)	0.170 _(0.001)
	primePCA_init	0.486 _(0.001)	0.449 _(0.001)	0.442 _(0.001)
	primePCA	0.290 _(0.001)	0.145 _(0.001)	0.097 _(0.0004)
(H4)	hardImpute	0.174 _(0.0005)	0.089 _(0.0003)	0.062 _(0.0003)
	softImpute(oracle)	0.121 _(0.0002)	0.062 _(0.0001)	0.042 _(0.0001)
	primePCA_init	0.203 _(0.001)	0.175 _(0.0005)	0.169 _(0.0004)
	primePCA	0.116 _(0.0003)	0.058 _(0.0002)	0.038 _(0.0001)

Million Song Dataset

- ▶ Original data has 110,000 users (rows) and 163,206 songs (columns); entries represent number of times a song was played by a particular user.
- ▶ Proportion of non-missing entries in the matrix is 0.008%.

Million Song Dataset

- ▶ Original data has 110,000 users (rows) and 163,206 songs (columns); entries represent number of times a song was played by a particular user.
- ▶ Proportion of non-missing entries in the matrix is 0.008%.
- ▶ Restrict attention to songs that have at least 100 listeners (1,777 songs in total). This improves the proportion of non-missing entries to 0.23%.

Million Song Dataset

- ▶ Original data has 110,000 users (rows) and 163,206 songs (columns); entries represent number of times a song was played by a particular user.
- ▶ Proportion of non-missing entries in the matrix is 0.008%.
- ▶ Restrict attention to songs that have at least 100 listeners (1,777 songs in total). This improves the proportion of non-missing entries to 0.23%.
- ▶ Quantiles of the number of listeners for each song:

0%	50%	60%	70%	80%	90%	100%
100	154	178	214	272.8	455.6	5043

Quantiles of the total play counts of each user:

0%	50%	60%	70%	80%	90%	100%
0	6	9	14	21	38	1114

Million Song Dataset

- ▶ Original data has 110,000 users (rows) and 163,206 songs (columns); entries represent number of times a song was played by a particular user.
- ▶ Proportion of non-missing entries in the matrix is 0.008%.
- ▶ Restrict attention to songs that have at least 100 listeners (1,777 songs in total). This improves the proportion of non-missing entries to 0.23%.
- ▶ Quantiles of the number of listeners for each song:

0%	50%	60%	70%	80%	90%	100%
100	154	178	214	272.8	455.6	5043

Quantiles of the total play counts of each user:

0%	50%	60%	70%	80%	90%	100%
0	6	9	14	21	38	1114

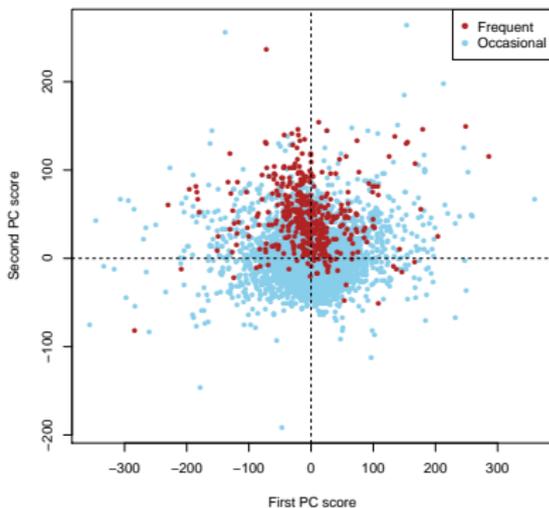
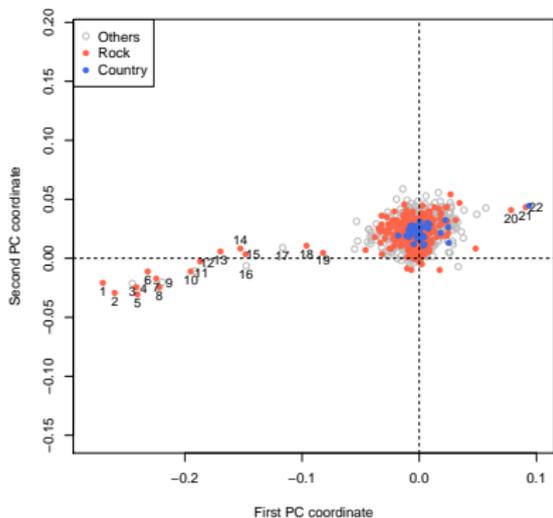
- ▶ Quantiles of non-missing matrix entry values:

0%	50%	60%	70%	80%	90%	100%
1	1	2	3	5	8	500

To guard against excessive influence from outliers, discretise play counts:

Play count	1	2 - 3	4 - 6	7 - 10	≥ 11
Level of interest	1	2	3	4	5

Million Song Dataset



Plots of the first two principal components $\widehat{V}_2^{\text{prime}}$ (left) and the associated scores $\{\widehat{u}_i\}_{i=1}^n$ (right).

Outlier songs

ID	Title	Artist	Genre
1	Your Hand In Mine	Explosions In The Sky	Rock
2	All These Things That I've Done	The Killers	Rock
3	Lady Marmalade	Christina Aguilera / Lil' Kim/ Mya / Pink	Pop
4	Here It Goes Again	Ok Go	Rock
5	I Hate Pretending (Album Version)	Secret Machines	Rock
6	No Rain	Blind Melon	Rock
7	Comatose (Comes Alive Version)	Skillet	Rock
8	Life In Technicolor	Coldplay	Rock
9	New Soul	Yael Naïm	Pop
10	Blurry	Puddle Of Mudd	Rock
11	Give It Back	Polly Paulusma	Pop
12	Walking On The Moon	The Police	Rock
13	Face Down (Album Version)	The Red Jumpsuit Apparatus	Rock
14	Savior	Rise Against	Rock
15	Swing Swing	The All-American Rejects	Rock
16	Without Me	Eminem	Rap
17	Almaz	Randy Crawford	Pop
18	Hotel California	Eagles	Rock
19	Hey There Delilah	Plain White T's	Rock
20	Revelry	Kings Of Leon	Rock
21	Undo	Björk	Rock
22	You're The One	Dwight Yoakam	Country

Summary

- ▶ Heterogeneous missingness is ubiquitous.

Summary

- ▶ Heterogeneous missingness is ubiquitous.
- ▶ The way in which the heterogeneity interacts with the underlying structure of interest is crucial.

Summary

- ▶ Heterogeneous missingness is ubiquitous.
- ▶ The way in which the heterogeneity interacts with the underlying structure of interest is crucial.
- ▶ `primePCA` iteratively projects observed entries of data matrix onto column space of current estimate to impute missing entries, then updates estimate by computing leading right singular space of imputed matrix.

Summary

- ▶ Heterogeneous missingness is ubiquitous.
- ▶ The way in which the heterogeneity interacts with the underlying structure of interest is crucial.
- ▶ `primePCA` iteratively projects observed entries of data matrix onto column space of current estimate to impute missing entries, then updates estimate by computing leading right singular space of imputed matrix.
- ▶ With an incoherence condition, the error of `primePCA` converges to zero at geometric rate in the noiseless setting.

Summary

- ▶ Heterogeneous missingness is ubiquitous.
- ▶ The way in which the heterogeneity interacts with the underlying structure of interest is crucial.
- ▶ `primePCA` iteratively projects observed entries of data matrix onto column space of current estimate to impute missing entries, then updates estimate by computing leading right singular space of imputed matrix.
- ▶ With an incoherence condition, the error of `primePCA` converges to zero at geometric rate in the noiseless setting.
- ▶ Theoretical guarantees depend on average, as opposed to worst-case, properties of the missingness mechanism.

Summary

- ▶ Heterogeneous missingness is ubiquitous.
- ▶ The way in which the heterogeneity interacts with the underlying structure of interest is crucial.
- ▶ `primePCA` iteratively projects observed entries of data matrix onto column space of current estimate to impute missing entries, then updates estimate by computing leading right singular space of imputed matrix.
- ▶ With an incoherence condition, the error of `primePCA` converges to zero at geometric rate in the noiseless setting.
- ▶ Theoretical guarantees depend on average, as opposed to worst-case, properties of the missingness mechanism.

Main reference:

- ▶ Zhu, Z., Wang, T. and Samworth, R. J. (2019) High-dimensional principal component analysis with heterogeneous missingness.
<https://arxiv.org/abs/1906.10125>.

Other references

- ▶ Anderson, T. W. (1957) Maximum likelihood estimates for a multivariate normal distribution when some observations are missing. *J. Amer. Statist. Assoc.*, **52**, 200–203.
- ▶ Belloni, A., Rosenbaum, M. and Tsybakov, A. B. (2017) Linear and conic programming estimators in high dimensional errors-in-variables models. *J. Roy. Statist. Soc., Ser. B*, **79**, 939–956.
- ▶ Cai, T. T. and Zhang, A. (2018a) Rate-optimal perturbation bounds for singular subspaces with applications to high-dimensional statistics. *Ann. Statist.*, **46**, 60–89.
- ▶ Cai, T. T. and Zhang, L. (2018b) High-dimensional linear discriminant analysis: optimality, adaptive algorithm, and missing data. [arXiv:1804.03018](https://arxiv.org/abs/1804.03018).
- ▶ Cho, J., Kim, D. and Rohe, K. (2017) Asymptotic theory for estimating the singular vectors and values of a partially-observed low rank matrix with noise. *Statist. Sinica*, **27**, 1921–1948.
- ▶ Dempster, A. P., Laird, N. M. and Rubin, D. B. (1977) Maximum likelihood from incomplete data via the EM algorithm. *J. Roy. Statist. Soc., Ser. B*, **39**, 1–38.
- ▶ Elsener, A and van de Geer, S. (2018) Sparse spectral estimation with missing and corrupted measurements. [arXiv:1811.10443](https://arxiv.org/abs/1811.10443).

Other references

- ▶ Ford, B. L. (1983) An overview of hot-deck procedures. In W. G. Madow, I. Olkin and D. B. Rubin (Eds.) *Incomplete Data in Sample Surveys, Vol. 2: Theory and Bibliographies*, 185–207. Academic Press, New York.
- ▶ Loh, P.-L. and Tan, X. L. (2018) High-dimensional robust precision matrix estimation: Cellwise corruption under ϵ -contamination. *Electron. J. Statist.*, **12**, 1429–1467.
- ▶ Loh, P.-L. and Wainwright, M. J. (2012) High-dimensional regression with noisy and missing data: provable guarantees with nonconvexity. *Ann. Statist.*, **40**, 1637–1664.
- ▶ Lounici, K. (2014) High-dimensional covariance matrix estimation with missing observations. *Bernoulli*, **20**, 1029–1058.
- ▶ Mazumder, R., Hastie, T. and Tibshirani R. (2010) Spectral regularization algorithms for learning large incomplete matrices. *J. Mach. Learn. Res.*, **11**, 2287–2322.
- ▶ Rubin, D. B. (2004) *Multiple Imputation for Nonresponse in Surveys*. John Wiley & Sons, Hoboken.
- ▶ Shen, D., Shen, H., Zhu, H. and Marron, J. (2016) The statistics and mathematics of high dimension low sample size asymptotics. *Statist. Sinica*, **26**, 1747–1770.