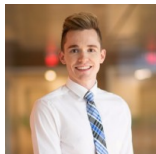


Probabilistic PCA for Heteroscedastic Data



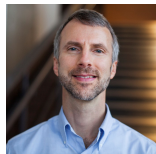
David Hong



Kyle Gilman



Laura Balzano



Jeffrey A. Fessler

EECS, University of Michigan

UM BIGDATA 2020-01-15

Introduction

Weighted PCA

Homoscedastic PPCA (review)

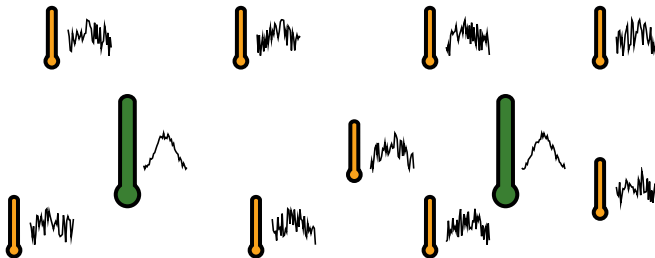
Heteroscedastic PPCA: known variances (2019)

Heteroscedastic PPCA: unknown variances (2020)

Extensions (2021)

Modern datasets increasingly contain data of varying quality.

For example: mixture of a few high quality (costly) sensors with many low quality (inexpensive) sensors.



Question: How should we account for heterogeneous quality?

Many settings: medical imaging, astronomy, sensor networks, ...



varying radiation levels



varying atmosphere



varying sensor quality

<http://www.medicalnewstoday.com/articles/153201.php>

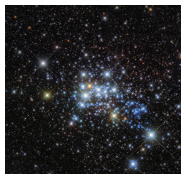
<https://www.nasa.gov/multimedia/imagegallery/iotd.html>

<http://www.livescience.com/27992-portable-pollution-sensors-improve-data-nsf-ria.html>

Many settings: medical imaging, astronomy, sensor networks, ...



varying radiation levels
(millions of voxels)



varying atmosphere
(thousands of pixels)



varying sensor quality
(thousands of locations)

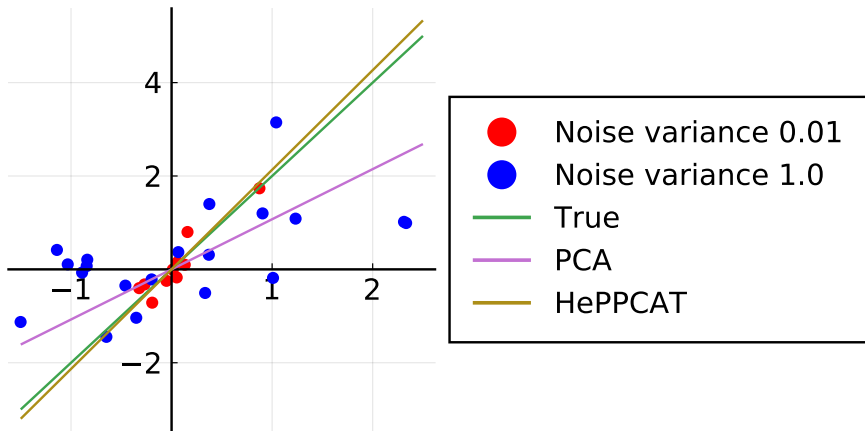
Much modern data is high-dimensional...with low-dimensional structure.
A standard (well-understood?) tool: Principal Component Analysis

How should we account for heterogeneous quality in PCA?

<http://www.medicalnewstoday.com/articles/153201.php>

<https://www.nasa.gov/multimedia/imagegallery/iotd.html>

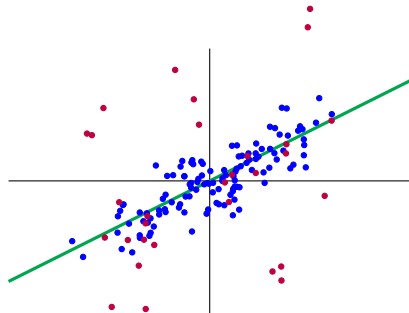
<http://www.livescience.com/27992-portable-pollution-sensors-improve-data-nsf-ria.html>



HePPCAT: Heteroscedastic probabilistic PCA technique

Model samples $y_1, \dots, y_n \in \mathbb{R}^d$ as

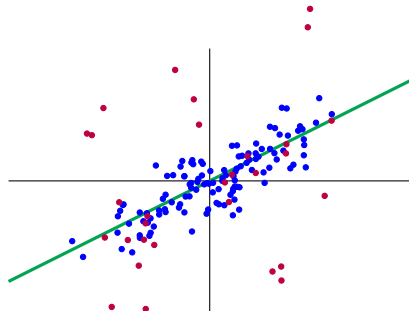
$$y_i = \mathbf{F}z_i + \sigma_i \varepsilon_i.$$



Model samples $y_1, \dots, y_n \in \mathbb{R}^d$ as

$$y_i = \mathbf{F}z_i + \sigma_i \varepsilon_i.$$

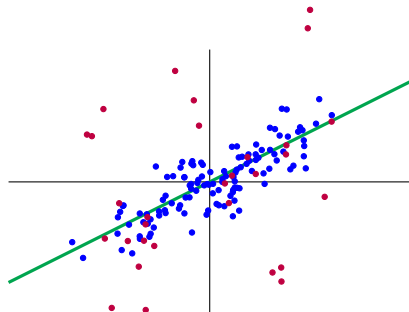
\mathbf{F} \swarrow \nwarrow σ_i
 k latent factors coefficients
 $d \times k$ matrix



Model samples $y_1, \dots, y_n \in \mathbb{R}^d$ as

$$y_i = \mathbf{F}z_i + \sigma_i \varepsilon_i.$$

noise std. dev. noise



80% $\sigma_i^2 = 0.1$, 20% $\sigma_i^2 = 1.9$

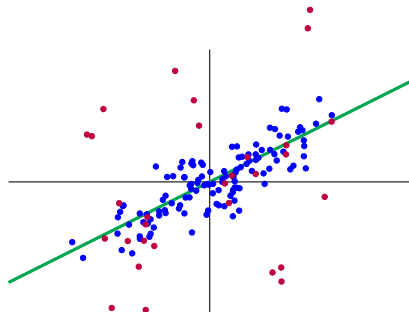
Model samples $y_1, \dots, y_n \in \mathbb{R}^d$ as

$$y_i = \mathbf{F}z_i + \sigma_i \varepsilon_i.$$

 ↑ ↑
noise std. dev. noise

In matrix form

$$\mathbf{Y} = [y_1, \dots, y_n]$$
$$= \mathbf{F} [z_1, \dots, z_n] + [\sigma_1 \varepsilon_1, \dots, \sigma_n \varepsilon_n].$$



80% $\sigma_i^2 = 0.1$, 20% $\sigma_i^2 = 1.9$

Recall:

$$y_i = \mathbf{F}z_i + \varepsilon_i \sim \mathcal{N}(\mathbf{F}z_i, \sigma_i^2 \mathbf{I})$$

Negative log-likelihood:

$$L(\mathbf{F}, \mathbf{Z}) = \sum_i \frac{1}{2\sigma_i^2} \|y_i - \mathbf{F}z_i\|_2^2.$$

Change of variables: $\tilde{y}_i = y_i/\sigma_i$, $\tilde{z}_i = z_i/\sigma_i$

$$\min_{\mathbf{F}, \mathbf{X}} \|\tilde{\mathbf{Y}} - \mathbf{F}\tilde{\mathbf{Z}}\|_{\mathbb{F}}^2 \implies \hat{\mathbf{F}} = \text{first } k \text{ left singular vectors of } \tilde{\mathbf{Y}}$$

This is “sample-weighted” low-rank approximation / factorization.
aka ML factor analysis (Young, 1941)

See Hong et al. arXiv 1810.12862 for optimally weighted PCA

Introduction

Weighted PCA

Homoscedastic PPCA (review)

Heteroscedastic PPCA: known variances (2019)

Heteroscedastic PPCA: unknown variances (2020)

Extensions (2021)

Model: $y_i = \mathbf{F}z_i + \sigma\varepsilon_i \in \mathbb{R}^d$ with $z_i \sim \mathcal{N}(0, \mathbf{I}_k)$ and $\varepsilon_i \sim \mathcal{N}(0, \mathbf{I}_d)$

Model: $y_i = \mathbf{F}z_i + \sigma\varepsilon_i \in \mathbb{R}^d$ with $z_i \sim \mathcal{N}(0, \mathbf{I}_k)$ and $\varepsilon_i \sim \mathcal{N}(0, \mathbf{I}_d)$

Equivalently, $y_i \sim \mathcal{N}(0, \mathbf{F}\mathbf{F}' + \sigma^2\mathbf{I}_d)$.

Model: $y_i = \mathbf{F}z_i + \sigma\varepsilon_i \in \mathbb{R}^d$ with $z_i \sim \mathcal{N}(0, \mathbf{I}_k)$ and $\varepsilon_i \sim \mathcal{N}(0, \mathbf{I}_d)$

Equivalently, $y_i \sim \mathcal{N}(0, \mathbf{F}\mathbf{F}' + \sigma^2\mathbf{I}_d)$.

Approach: Maximize the log-likelihood of \mathbf{F} (dropping constants)

$$\begin{aligned}\mathcal{L}(\mathbf{F}) &:= \frac{1}{2} \sum_{i=1}^n \left\{ \ln \det(\mathbf{F}\mathbf{F}' + \sigma^2\mathbf{I}_d)^{-1} - y_i'(\mathbf{F}\mathbf{F}' + \sigma^2\mathbf{I}_d)^{-1}y_i \right\} \\ &= \frac{n}{2} \left\{ \ln \det(\mathbf{F}\mathbf{F}' + \sigma^2\mathbf{I}_d)^{-1} - \frac{1}{n} \sum_{i=1}^n y_i'(\mathbf{F}\mathbf{F}' + \sigma^2\mathbf{I}_d)^{-1}y_i \right\} \\ &= \frac{n}{2} \left[\ln \det(\mathbf{F}\mathbf{F}' + \sigma^2\mathbf{I}_d)^{-1} - \text{tr} \left\{ \frac{1}{n} \sum_{i=1}^n y_i y_i' (\mathbf{F}\mathbf{F}' + \sigma^2\mathbf{I}_d)^{-1} \right\} \right].\end{aligned}$$

Data appears only through the sample correlation matrix!

$$\mathcal{L}(\mathbf{F}) = \frac{n}{2} \left[\ln \det(\mathbf{F}\mathbf{F}' + \sigma^2 \mathbf{I}_d)^{-1} - \text{tr} \left\{ \frac{1}{n} \sum_{i=1}^n y_i y_i' (\mathbf{F}\mathbf{F}' + \sigma^2 \mathbf{I}_d)^{-1} \right\} \right].$$

$$\mathcal{L}(\mathbf{F}) = \frac{n}{2} \left[\ln \det(\mathbf{F}\mathbf{F}' + \sigma^2 \mathbf{I}_d)^{-1} - \text{tr} \left\{ \frac{1}{n} \sum_{i=1}^n y_i y_i' (\mathbf{F}\mathbf{F}' + \sigma^2 \mathbf{I}_d)^{-1} \right\} \right].$$

Fact (Tipping & Bishop 1999): The likelihood is maximized by

$$\hat{\mathbf{F}} = \mathbf{V} \text{diag} \left(\sqrt{\lambda_1 - \bar{\lambda}}, \dots, \sqrt{\lambda_k - \bar{\lambda}} \right),$$

where in terms of the **sample correlation matrix** $\frac{1}{n} \mathbf{Y}\mathbf{Y}'$:

- ▶ $\mathbf{V} \in \mathbb{R}^{d \times k}$ contains the k principal eigenvectors,
- ▶ $\lambda_1, \dots, \lambda_k$ are the k principal eigenvalues, and
- ▶ $\bar{\lambda}$ is the average of the remaining eigenvalues.

$$\mathcal{L}(\mathbf{F}) = \frac{n}{2} \left[\ln \det(\mathbf{F}\mathbf{F}' + \sigma^2 \mathbf{I}_d)^{-1} - \text{tr} \left\{ \frac{1}{n} \sum_{i=1}^n y_i y_i' (\mathbf{F}\mathbf{F}' + \sigma^2 \mathbf{I}_d)^{-1} \right\} \right].$$

Fact (Tipping & Bishop 1999): The likelihood is maximized by

$$\hat{\mathbf{F}} = \mathbf{V} \text{diag} \left(\sqrt{\lambda_1 - \bar{\lambda}}, \dots, \sqrt{\lambda_k - \bar{\lambda}} \right),$$

where in terms of the **sample correlation matrix** $\frac{1}{n} \mathbf{Y}\mathbf{Y}'$:

- ▶ $\mathbf{V} \in \mathbb{R}^{d \times k}$ contains the k principal eigenvectors,
- ▶ $\lambda_1, \dots, \lambda_k$ are the k principal eigenvalues, and
- ▶ $\bar{\lambda}$ is the average of the remaining eigenvalues.

Clean solution: Eigendecomposition + shrinkage

$$\mathcal{L}(\mathbf{F}) = \frac{n}{2} \left[\ln \det(\mathbf{F}\mathbf{F}' + \sigma^2 \mathbf{I}_d)^{-1} - \text{tr} \left\{ \frac{1}{n} \sum_{i=1}^n y_i y_i' (\mathbf{F}\mathbf{F}' + \sigma^2 \mathbf{I}_d)^{-1} \right\} \right].$$

Fact (Tipping & Bishop 1999): The likelihood is maximized by

$$\hat{\mathbf{F}} = \mathbf{V} \text{diag} \left(\sqrt{\lambda_1 - \bar{\lambda}}, \dots, \sqrt{\lambda_k - \bar{\lambda}} \right),$$

where in terms of the **sample correlation matrix** $\frac{1}{n} \mathbf{Y}\mathbf{Y}'$:

- ▶ $\mathbf{V} \in \mathbb{R}^{d \times k}$ contains the k principal eigenvectors,
- ▶ $\lambda_1, \dots, \lambda_k$ are the k principal eigenvalues, and
- ▶ $\bar{\lambda}$ is the average of the remaining eigenvalues.

*Clean solution: Eigendecomposition + shrinkage
Components \mathbf{V} are same as in ordinary PCA*

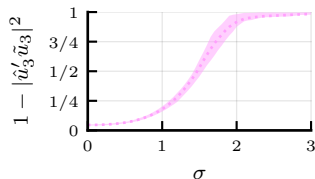
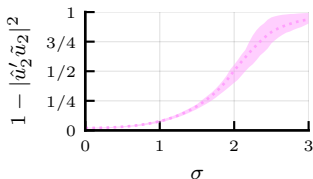
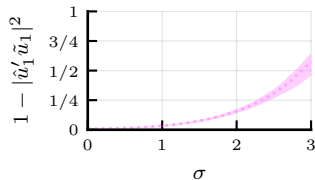
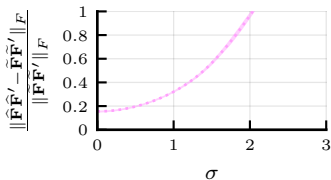
(Homoscedastic) PPCA degrades with heteroscedasticity



Setup: $n_1 = 200$ samples with $\sigma_i = 1$, $n_2 = 800$ samples with $\sigma_i = \sigma$.

$$y_i \sim \mathcal{N}(0, \mathbf{F}_* \mathbf{F}_*^{\prime} + \sigma_i^2 \mathbf{I}_{100}), \quad \mathbf{F}_* = [\tilde{u}_1, \dots, \tilde{u}_3] \text{diag}(4, 2, 1),$$

PPCA: Full data



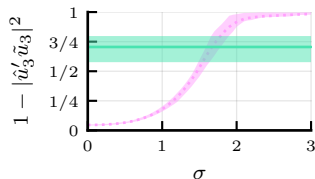
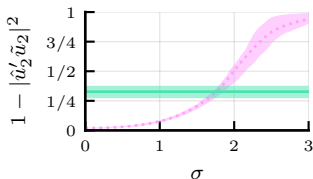
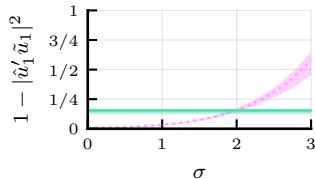
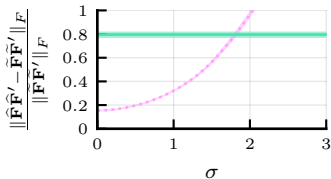
(Homoscedastic) PPCA degrades with heteroscedasticity

Setup: $n_1 = 200$ samples with $\sigma_i = 1$, $n_2 = 800$ samples with $\sigma_i = \sigma$.

$$y_i \sim \mathcal{N}(0, \mathbf{F}_* \mathbf{F}_*' + \sigma_i^2 \mathbf{I}_{100}), \quad \mathbf{F}_* = [\tilde{u}_1, \dots, \tilde{u}_3] \text{diag}(4, 2, 1),$$

• PPCA: Full data

■ PPCA: Group 1



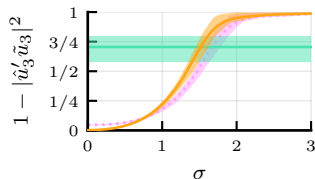
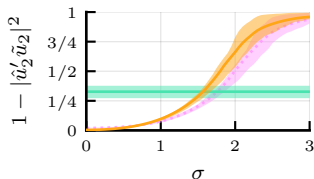
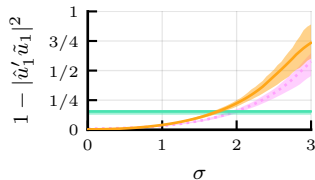
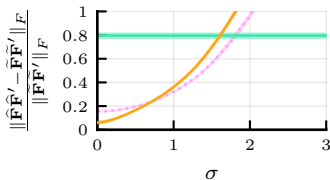
(Homoscedastic) PPCA degrades with heteroscedasticity



Setup: $n_1 = 200$ samples with $\sigma_i = 1$, $n_2 = 800$ samples with $\sigma_i = \sigma$.

$$y_i \sim \mathcal{N}(0, \mathbf{F}_* \mathbf{F}_*' + \sigma_i^2 \mathbf{I}_{100}), \quad \mathbf{F}_* = [\tilde{u}_1, \dots, \tilde{u}_3] \text{diag}(4, 2, 1),$$

■ PPCA: Full data ■ PPCA: Group 1 ■ PPCA: Group 2

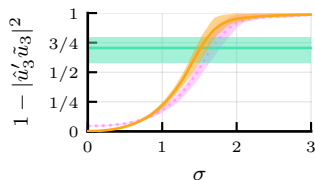
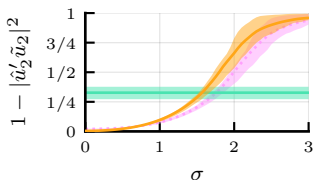
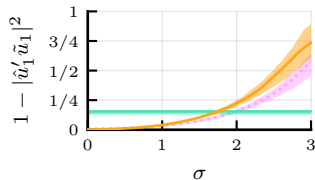
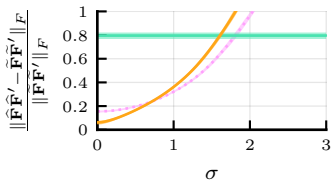


(Homoscedastic) PPCA degrades with heteroscedasticity

Setup: $n_1 = 200$ samples with $\sigma_i = 1$, $n_2 = 800$ samples with $\sigma_i = \sigma$.

$$y_i \sim \mathcal{N}(0, \mathbf{F}_* \mathbf{F}_*' + \sigma_i^2 \mathbf{I}_{100}), \quad \mathbf{F}_* = [\tilde{u}_1, \dots, \tilde{u}_3] \text{diag}(4, 2, 1),$$

■ PPCA: Full data
 ■ PPCA: Group 1
 ■ PPCA: Group 2



*PPCA degrades as data becomes heteroscedastic.
 PPCA may work better using only part of the data!?*

Introduction

Weighted PCA

Homoscedastic PPCA (review)

Heteroscedastic PPCA: known variances (2019)

Heteroscedastic PPCA: unknown variances (2020)

Extensions (2021)

This talk: Develop a heteroscedastic PPCA algorithm.

Version 1: Assume noise variances $\{v_i = \sigma_i^2\}$ are known.

This talk: Develop a heteroscedastic PPCA algorithm.

Version 1: Assume noise variances $\{v_i = \sigma_i^2\}$ are known.

Challenge: Log-likelihood no longer separates so nicely:

$$\mathcal{L}(\mathbf{F}) = \frac{1}{2} \sum_{i=1}^n \left\{ \ln \det(\mathbf{F}\mathbf{F}' + v_i \mathbf{I}_d)^{-1} - y_i'(\mathbf{F}\mathbf{F}' + v_i \mathbf{I}_d)^{-1} y_i \right\} \quad (1)$$

Apparently no closed-form for ML estimate of \mathbf{F} .

Deriving a heteroscedastic PPCA

This talk: Develop a heteroscedastic PPCA algorithm.

Version 1: Assume noise variances $\{v_i = \sigma_i^2\}$ are known.

Challenge: Log-likelihood no longer separates so nicely:

$$\mathcal{L}(\mathbf{F}) = \frac{1}{2} \sum_{i=1}^n \left\{ \ln \det(\mathbf{F}\mathbf{F}' + v_i \mathbf{I}_d)^{-1} - y_i'(\mathbf{F}\mathbf{F}' + v_i \mathbf{I}_d)^{-1} y_i \right\} \quad (1)$$

Apparently no closed-form for ML estimate of \mathbf{F} .

Approach: Instead, derive an Expectation Maximization (EM) algorithm using the complete data likelihood with complete data $\{(y_i, z_i)\}$

$$\mathcal{L}_c(\mathbf{F}) := - \sum_{i=1}^n \left(\frac{\|y_i - \mathbf{F}z_i\|_2^2}{2v_i} + \frac{\|z_i\|_2^2}{2} \right).$$

(Tipping & Bishop 1999) derived it for homoscedastic PPCA.

$$\mathcal{L}_c(\mathbf{F}) := - \sum_{i=1}^n \left(\frac{\|y_i - \mathbf{F}z_i\|_2^2}{2v_i} + \frac{\|z_i\|_2^2}{2} \right).$$

E-Step: Expectation with respect to $z_1, \dots, z_n | y_1, \dots, y_n, \mathbf{F}_t$ is

$$\bar{\mathcal{L}}(\mathbf{F}; \mathbf{F}_t) = \sum_{i=1}^n \left[\frac{1}{v_i} y_i' \mathbf{F} \bar{z}_{t,i} - \frac{1}{2v_i} \text{tr} \{ \mathbf{F}' \mathbf{F} (\bar{z}_{t,i} \bar{z}_{t,i}' + v_i \mathbf{M}_{t,i}) \} \right],$$

up to constants where

$$\mathbf{M}_{t,i} := (\mathbf{F}_t' \mathbf{F}_t + v_i \mathbf{I}_k)^{-1}, \quad \bar{z}_{t,i} := \mathbf{M}_{t,i} \mathbf{F}_t' y_i.$$

M-step: $\bar{\mathcal{L}}(\mathbf{F}; \mathbf{F}_t)$ is maximized with respect to \mathbf{F} by

$$\mathbf{F}_{t+1} = \left(\sum_{i=1}^n \frac{1}{v_i} y_i \bar{z}_{t,i}' \right) \left(\sum_{i=1}^n \frac{1}{v_i} \bar{z}_{t,i} \bar{z}_{t,i}' + \mathbf{M}_{t,i} \right)^{-1}.$$

$$\mathcal{L}_c(\mathbf{F}) := - \sum_{i=1}^n \left(\frac{\|y_i - \mathbf{F}z_i\|_2^2}{2v_i} + \frac{\|z_i\|_2^2}{2} \right).$$

E-Step: Expectation with respect to $z_1, \dots, z_n | y_1, \dots, y_n, \mathbf{F}_t$ is

$$\bar{\mathcal{L}}(\mathbf{F}; \mathbf{F}_t) = \sum_{i=1}^n \left[\frac{1}{v_i} y_i' \mathbf{F} \bar{z}_{t,i} - \frac{1}{2v_i} \text{tr} \{ \mathbf{F}' \mathbf{F} (\bar{z}_{t,i} \bar{z}_{t,i}' + v_i \mathbf{M}_{t,i}) \} \right],$$

up to constants where

$$\mathbf{M}_{t,i} := (\mathbf{F}_t' \mathbf{F}_t + v_i \mathbf{I}_k)^{-1}, \quad \bar{z}_{t,i} := \mathbf{M}_{t,i} \mathbf{F}_t' y_i.$$

M-step: $\bar{\mathcal{L}}(\mathbf{F}; \mathbf{F}_t)$ is maximized with respect to \mathbf{F} by

$$\mathbf{F}_{t+1} = \left(\sum_{i=1}^n \frac{1}{v_i} y_i \bar{z}_{t,i}' \right) \left(\sum_{i=1}^n \frac{1}{v_i} \bar{z}_{t,i} \bar{z}_{t,i}' + \mathbf{M}_{t,i} \right)^{-1}.$$

Now the factor estimates $\hat{\mathbf{F}}$ differ from PCA eigenvectors.

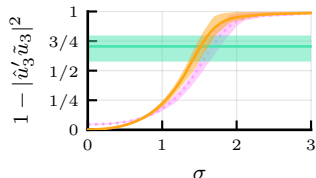
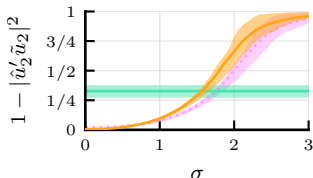
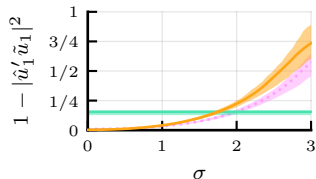
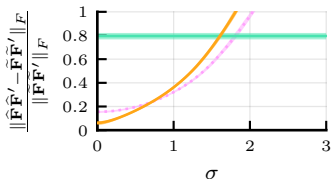
Heteroscedastic PPCA effectively uses all data

Setup: $n_1 = 200$ samples with $\sigma_i = 1$, $n_2 = 800$ samples with $\sigma_i = \sigma$.

$$y_i \sim \mathcal{N}(0, \mathbf{F}_* \mathbf{F}_*^T + v_i \mathbf{I}_{100}),$$

$$\mathbf{F}_* = [\tilde{u}_1, \dots, \tilde{u}_3] \text{diag}(4, 2, 1),$$

PPCA: Full data PPCA: Group 1 PPCA: Group 2



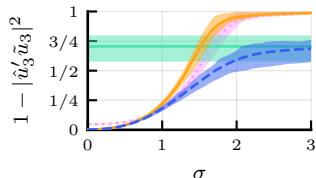
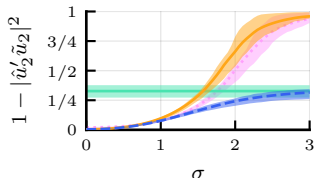
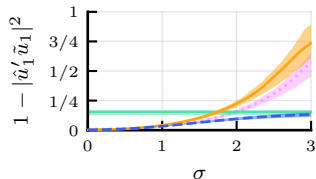
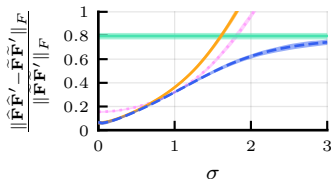
Heteroscedastic PPCA effectively uses all data

Setup: $n_1 = 200$ samples with $\sigma_i = 1$, $n_2 = 800$ samples with $\sigma_i = \sigma$.

$$y_i \sim \mathcal{N}(0, \mathbf{F}_* \mathbf{F}_*^T + v_i \mathbf{I}_{100}),$$

$$\mathbf{F}_* = [\tilde{u}_1, \dots, \tilde{u}_3] \text{diag}(4, 2, 1),$$

PPCA: Full data PPCA: Group 1 PPCA: Group 2 HeteroPPCA

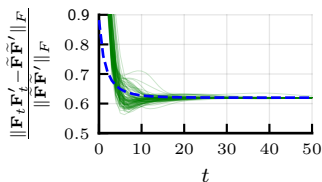
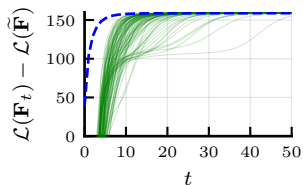


But how to initialize for this nonconcave problem?

Setup: $n_1 = 200$ samples with $v_i = 1$, $n_2 = 800$ samples with $v_i = 4$.

$$y_i \sim \mathcal{N}(0, \mathbf{F}_* \mathbf{F}_*^' + v_i \mathbf{I}_{100}), \quad \mathbf{F}_* = [\tilde{u}_1, \dots, \tilde{u}_3] \text{diag}(4, 2, 1)^{1/2},$$

Rand init (true eigvals, rand eigvecs) vs. Homoscedastic PPCA init



*Interesting phenomenon:
iterates do not seem to hit bad local maxima! Why?*

What about doing a weighted PCA?

Idea: Give noisier samples less weight in PCA.

Simple tweak: Replace the sample covariance with a *weighted* version

$$\frac{1}{n} \sum_{i=1}^n y_i y_i' \quad \longrightarrow \quad \frac{1}{n} \sum_{i=1}^n \omega_i^2 y_i y_i',$$

Choices for weights:

- ▶ unweighted: $\omega_i^2 = 1$
- ▶ inverse noise variance $\omega_i^2 = 1/\sigma_i^2$
 - ▶ rescales data to make noise homoscedastic
 - ▶ corresponds to MLE for low-rank signal $\mathbf{F} [z_1, \dots, z_n]$
- ▶ square inverse noise variance $\omega_i^2 = 1/\sigma_i^4$
 - ▶ more aggressive down-weighting, can help for low SNRs

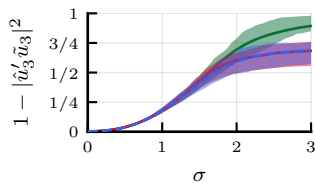
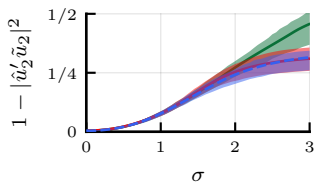
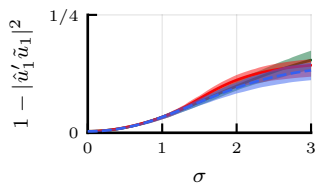
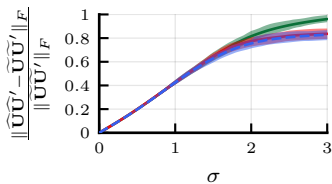
High-dim asymptotics studied in (Hong, Fessler, Balzano 2019).

What about doing a weighted PCA? (Jolliffe 2002)

Setup: $n_1 = 200$ samples with $\sigma_i = 1$, $n_2 = 800$ samples with $\sigma_i = \sigma$.

$$y_i \sim \mathcal{N}(0, \mathbf{F}_* \mathbf{F}_*' + \sigma_i^2 \mathbf{I}_{100}), \quad \mathbf{F}_* = [\tilde{u}_1, \dots, \tilde{u}_3] \text{diag}(4, 2, 1),$$

■ Inv. noise var. ■ Sq. inv. noise var. ■ HeteroPPCA



- ▶ PPCA for data with heteroscedastic noise that uses all the data effectively, no matter how noisy
- ▶ practical EM algorithm
- ▶ interesting phenomenon: do not seem to get stuck in bad local maxima!

- ▶ PPCA for data with heteroscedastic noise that uses all the data effectively, no matter how noisy
- ▶ practical EM algorithm
- ▶ interesting phenomenon: do not seem to get stuck in bad local maxima!

Ongoing/future: (from Dec. 2019 CAMSAP talk)

- ▶ joint estimation of noise variance
- ▶ optimization landscape
- ▶ other optimization approaches (manifold optimization?, minorize maximize?)
- ▶ analysis of asymptotic performance

Introduction

Weighted PCA

Homoscedastic PPCA (review)

Heteroscedastic PPCA: known variances (2019)

Heteroscedastic PPCA: unknown variances (2020)

Extensions (2021)

Model $n_1 + \dots + n_L = n$ data samples in \mathbb{R}^d from L noise level groups:

$$\mathbf{y}_{\ell,i} = \mathbf{F}\mathbf{z}_{\ell,i} + \boldsymbol{\varepsilon}_{\ell,i}, \quad i \in \{1, \dots, n_\ell\}, \ell \in \{1, \dots, L\}, \quad (2)$$

$\mathbf{F} \in \mathbb{R}^{d \times k}$: deterministic factor matrix to estimate

$\mathbf{z}_{\ell,i} \sim \mathcal{N}(\mathbf{0}_k, \mathbf{I}_k)$: random coefficients

$\boldsymbol{\varepsilon}_{\ell,i} \sim \mathcal{N}(\mathbf{0}_d, v_\ell \mathbf{I}_d)$: noise vectors

v_1, \dots, v_L : deterministic noise variances to estimate.

Equivalently:

$$\mathbf{y}_{\ell,i} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{F}\mathbf{F}' + v_\ell \mathbf{I}_d)$$

Joint log-likelihood

Joint log-likelihood for factors \mathbf{F} and variances \mathbf{v} (dropping a constant):

$$\mathcal{L}(\mathbf{F}, \mathbf{v}) \triangleq \frac{1}{2} \sum_{\ell=1}^L \left[n_{\ell} \ln \det(\mathbf{F}\mathbf{F}' + v_{\ell}\mathbf{I}_d)^{-1} - \text{tr} \{ \mathbf{Y}'_{\ell}(\mathbf{F}\mathbf{F}' + v_{\ell}\mathbf{I}_d)^{-1} \mathbf{Y}_{\ell} \} \right],$$

$\mathbf{Y}_{\ell} \triangleq [\mathbf{y}_{\ell,1}, \dots, \mathbf{y}_{\ell,n_{\ell}}] \in \mathbb{R}^{d \times n_{\ell}}$ for $\ell \in \{1, \dots, L\}$:
 data matrices for each of the L groups.

Similar to earlier log-likelihood (1).

Apparently no closed form ML estimates for \mathbf{F} and \mathbf{v}

No closed-form ML estimate of \mathbf{F} given \mathbf{v}

No closed-form ML estimate of \mathbf{v} given \mathbf{F}

\implies Alternating ascent algorithms

\mathbf{F} update : EM update similar to Version 1 HePPCAT

\mathbf{v} update : separates into L univariate maximizations (over $v_\ell \geq 0$)

$$\mathcal{L}_\ell(\mathbf{v}_\ell) \triangleq - \sum_{j=0}^k \left\{ \alpha_j \ln(\gamma_j + v_\ell) + \frac{\beta_j}{\gamma_j + v_\ell} \right\}, \quad (3)$$

where $\alpha_0 \triangleq d - k$, $\beta_0 \triangleq \|(\mathbf{I}_d - \mathbf{U}_t \mathbf{U}_t') \mathbf{Y}_\ell\|_F^2 / n_\ell$, $\gamma_0 \triangleq 0$,

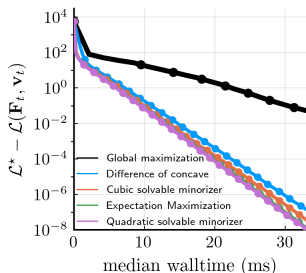
$$j \geq 1 : \quad \alpha_j \triangleq 1, \quad \beta_j \triangleq \|\mathbf{Y}'_\ell \boldsymbol{\mu}_{t,j}\|_2^2 / n_\ell, \quad \gamma_j \triangleq \lambda_{t,j},$$

and $\mathbf{U}_t = [\boldsymbol{\mu}_{t,1}, \dots, \boldsymbol{\mu}_{t,k}]$ and $\boldsymbol{\lambda}_t = (\lambda_{t,1}, \dots, \lambda_{t,k})$ are the eigenvectors and eigenvalues of $\mathbf{F}_t \mathbf{F}_t'$ at iteration t .

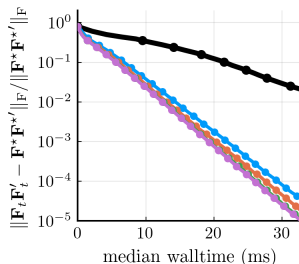
Univariate but non-convex.

Methods developed and investigated:

root finding, EM, minorize-maximize (MM), difference of concave...



(a) Convergence w.r.t. \mathcal{L} .



(b) Convergence w.r.t. \mathbf{F} .

Figure: Comparison of update methods for $n = 10^3$ samples in $d = 10^2$ dimensions with $k = 3$ underlying factors $\lambda_* = (4, 2, 1)$. The noise is heteroscedastic: the first $n_1 = 200$ samples have noise variance $\tilde{v}_1 = 1$ and the remaining $n_2 = 800$ have $\tilde{v}_2 = 4$. Walltimes are medians taken over 100 runs of the algorithm to reduce the effect of experimental noise. Markers are placed every five iterations.

Statistical performance example

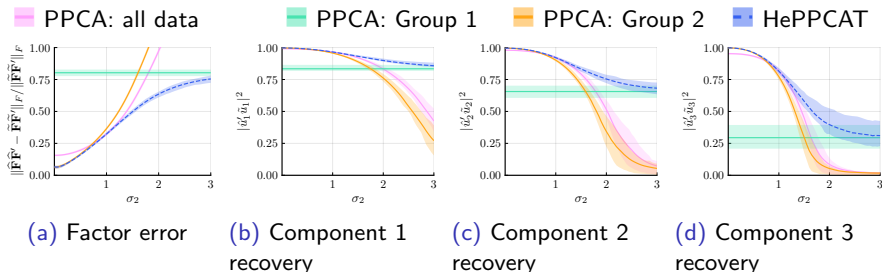


Figure: Comparison with homoscedastic PPCA applied on: i) full data, ii) only group 1, i.e., the $n_1 = 200$ samples with noise variance $\tilde{v}_1 = 1$, and iii) only group 2, i.e., the $n_2 = 800$ samples with noise variance $\tilde{v}_2 = \sigma_2^2$. Lower is better in (a), and higher is better in (b)-(d). HePPCAT outperforms the homoscedastic methods on all four metrics.

Here, problem size is large enough that HePPCAT with unknown variance work essentially as well as with known variance.

Effect of block size on variance estimates

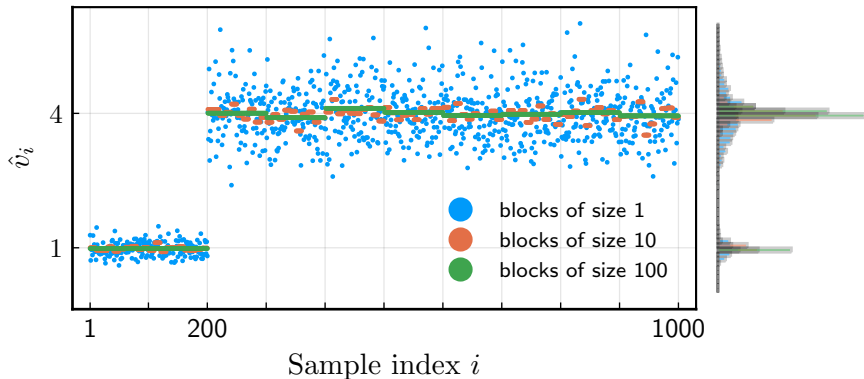


Figure: Estimated noise variances for various block sizes.

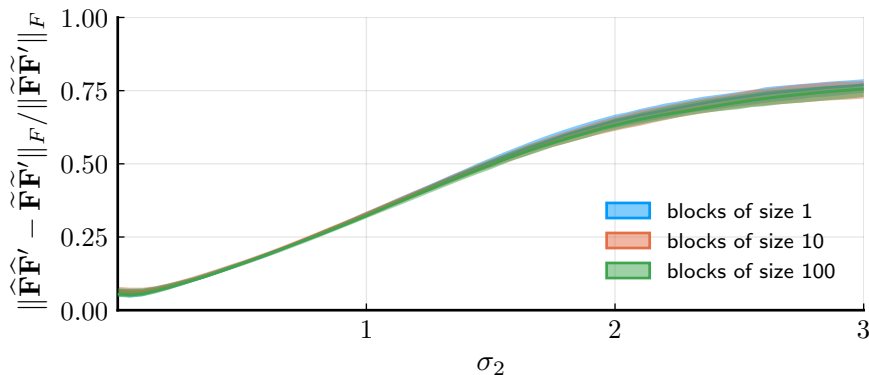
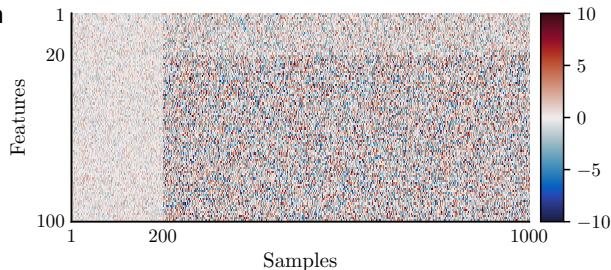


Figure: Normalized factor estimation error (median and interquartile intervals) for varying block sizes.

The three block sizes lead to practically identical performance here.

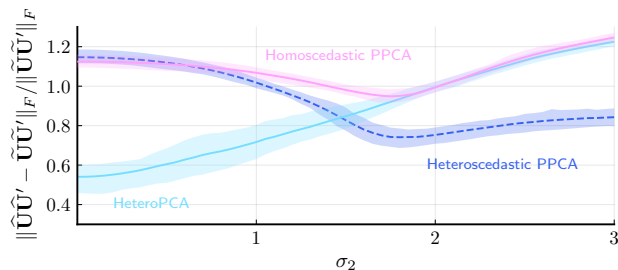
Heteroscedasticity across features *and* samples

Single data realization
(heatmap).



Normalized error in
subspace estimation
(lower is better).

Zhang, Cai, Wu. 2018
ArXiv 1810.08316



Apparently favorable convergence

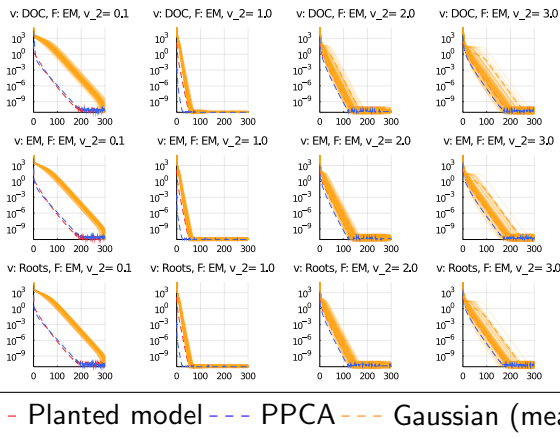


Figure: Convergence gaps of each algorithm to the maximum converged log-likelihood per heteroscedastic noise experiment. $n = [200, 800]$ and $v_1 = 1$. Three different types of initializations.

Introduction

Weighted PCA

Homoscedastic PPCA (review)

Heteroscedastic PPCA: known variances (2019)

Heteroscedastic PPCA: unknown variances (2020)

Extensions (2021)

- ▶ Incremental / online algorithms

- ▶ Incremental / online algorithms
- ▶ Acceleration via momentum

- ▶ Incremental / online algorithms
- ▶ Acceleration via momentum
- ▶ convergence theory / performance guarantees

- ▶ Incremental / online algorithms
- ▶ Acceleration via momentum
- ▶ convergence theory / performance guarantees
- ▶ Sketching to reduce computation (e.g., for \mathbf{v})

- ▶ Incremental / online algorithms
- ▶ Acceleration via momentum
- ▶ convergence theory / performance guarantees
- ▶ Sketching to reduce computation (e.g., for \mathbf{v})
- ▶ Nonnegative PCA, sparse PCA, robust PCA, etc.,
all generalized for heteroscedastic noise

- ▶ Incremental / online algorithms
- ▶ Acceleration via momentum
- ▶ convergence theory / performance guarantees
- ▶ Sketching to reduce computation (e.g., for \mathbf{v})
- ▶ Nonnegative PCA, sparse PCA, robust PCA, etc., all generalized for heteroscedastic noise
- ▶ Supervised HePPCAT?

- ▶ Incremental / online algorithms
- ▶ Acceleration via momentum
- ▶ convergence theory / performance guarantees
- ▶ Sketching to reduce computation (e.g., for \mathbf{v})
- ▶ Nonnegative PCA, sparse PCA, robust PCA, etc., all generalized for heteroscedastic noise
- ▶ Supervised HePPCAT?
- ▶ missing data, heterogeneity across features

- ▶ Incremental / online algorithms
- ▶ Acceleration via momentum
- ▶ convergence theory / performance guarantees
- ▶ Sketching to reduce computation (e.g., for \mathbf{v})
- ▶ Nonnegative PCA, sparse PCA, robust PCA, etc., all generalized for heteroscedastic noise
- ▶ Supervised HePPCAT?
- ▶ missing data, heterogeneity across features
- ▶ bias correction for variance estimates

- ▶ Incremental / online algorithms
- ▶ Acceleration via momentum
- ▶ convergence theory / performance guarantees
- ▶ Sketching to reduce computation (e.g., for \mathbf{v})
- ▶ Nonnegative PCA, sparse PCA, robust PCA, etc., all generalized for heteroscedastic noise
- ▶ Supervised HePPCAT?
- ▶ missing data, heterogeneity across features
- ▶ bias correction for variance estimates
- ▶ probabilistic dictionary learning (\mathbf{F} wide, \mathbf{z} sparse, so non-normal)

Data $\mathbf{Y} \in \mathbb{R}^{d \times n}$; Factors $\mathbf{F} \in \mathbb{R}^{d \times k}$; Clay's version:

$$\min_{\beta \in \mathbb{R}^k} \min_{\mathbf{F} \in \mathbb{R}^{d \times k}} L(\mathbf{Y}'\mathbf{F}\beta) + \lambda \|\mathbf{Y} - \mathbf{F}\mathbf{F}'\mathbf{Y}\|_{\mathbf{F}}^2. \quad (4)$$

(Any statistical model associated with norm? If so, then homoscedastic.)

Data $\mathbf{Y} \in \mathbb{R}^{d \times n}$; Factors $\mathbf{F} \in \mathbb{R}^{d \times k}$; Clay's version:

$$\min_{\beta \in \mathbb{R}^k} \min_{\mathbf{F} \in \mathbb{R}^{d \times k}} L(\mathbf{Y}'\mathbf{F}\beta) + \lambda \|\mathbf{Y} - \mathbf{F}\mathbf{F}'\mathbf{Y}\|_{\mathbf{F}}^2. \quad (4)$$

(Any statistical model associated with norm? If so, then homoscedastic.)

Supervised HePPCAT for known noise variances:

$$\min_{\beta \in \mathbb{R}^k} \min_{\mathbf{F} \in \mathbb{R}^{d \times k}} L(\mathbf{Y}'\mathbf{F}\beta) + \underbrace{\frac{\lambda}{2} \sum_i -\ln \det(\mathbf{F}\mathbf{F}' + \sigma_i^2) + \mathbf{y}'_i (\mathbf{F}\mathbf{F}' + \sigma_i^2)^{-1} \mathbf{y}_i}_{\text{neg. log likelihood}}.$$

EM majorizer for \mathbf{F} update is quadratic, not “quartic” like (4):

$$L(\mathbf{Y}'\mathbf{F}\beta) + \sum_i \frac{\lambda}{2\sigma_i^2} [-2\mathbf{y}'_i \mathbf{F} \bar{\mathbf{z}}_i - \bar{\mathbf{z}}'_i \mathbf{F}' \mathbf{F} \bar{\mathbf{z}}_i + \text{tr}(\mathbf{F} \bar{\mathbf{M}} \mathbf{F}')] .$$