# Empirical Bayes PCA in high dimensions

Xinyi Zhong Department of Statistics and Data Science Yale University

> Berkeley CDAR seminar Sept. 27<sup>th</sup>, 2022

## Joint work with



Adviser: Dr. Zhou Fan Department of Statistics and Data Science Yale University Chang Su Department of Biostatistics Yale University

## Motivation

- Model:  $\mathbf{X} = \mathbf{U} \mathbf{V}^{\top} + \mathbf{W}$
- Goal: estimate  $\mathbf{U}\mathbf{V}^{\mathsf{T}}$

- Naïve PCA
  - High dimensional noise
  - Doesn't utilize PC structure
- Bayes estimate
  - Intensive computation
  - Unavailable prior information

Empirical Bayes PCA (EBPCA) uses less computational cost to achieve more accurate estimate.

- polynomial time
- Bayes optimal





East-Asian

an 🔸

Hispanic

South Asian

# Outline

- Method and results
- Relation to other methods
- Extensions of the noise model

## The EB-PCA method

There are three statistical ideas underlying EB-PCA:

- Random matrix theory for the sample PCs
- Empirical Bayes estimation for normal means
- Bayesian inference using Approximate Message Passing (AMP)

#### Rank-one signal-plus-noise model

$$\mathbf{X} = \frac{s}{n} * \mathbf{u}\boldsymbol{v}^{\mathsf{T}} + \mathbf{W} \approx \frac{\lambda}{n} * \boldsymbol{f}\boldsymbol{g}^{\mathsf{T}}$$





#### Empirical Bayes in compound decision

$$g_j \sim N(\mu_* * \boldsymbol{v_j}, \sigma_*^2)$$
 for all  $j \in [n] \neq \mathbf{g} \sim N(\mu_* * \boldsymbol{v}, \sigma_*^2)$ 

#### Nonparametric empirical

We estimate  $\pi$  through nonparametric MLE:

-1

Model:  $\boldsymbol{x} \sim N(\boldsymbol{\theta}, \sigma^2), \ \boldsymbol{\theta} \sim \pi$ 

Inference:  $\hat{\pi} = \operatorname{argmax}_{\pi} \sum_{i=1}^{n} \log f_{\pi * N(0, \sigma^2)}(x_i)$ Denoise:  $\hat{\theta} = \operatorname{E}[\theta | x] \coloneqq \theta(\theta | \mu, \sigma^2, \pi)$ 

#### Empirical Bayes in compound decision

Given  $\mathbf{g} \sim N(\mu_* * \boldsymbol{\nu}, \sigma_*^2)$ 

1. Infer  $v_i \sim \pi$ 

2. Denoise *g* by taking  $\hat{\boldsymbol{v}} = E[\boldsymbol{v}|\boldsymbol{g}] \coloneqq \theta(\boldsymbol{g}|\mu_*, \sigma_*^2, \pi)$ 

#### Iterative refinement via AMP

As  $\widehat{v} \prec g$  for v, we expect  $Y\widehat{v} \prec Yg \propto f$  for u.

$$\boldsymbol{v}_{t} = \theta(\boldsymbol{g}_{t} | \boldsymbol{\mu}_{*}, \sigma_{*}^{2}, \pi)$$
$$\boldsymbol{f}_{t} = \boldsymbol{Y} \boldsymbol{v}_{t}$$
$$\boldsymbol{u}_{t} = \theta(\boldsymbol{f}_{t} | \bar{\boldsymbol{\mu}}_{*}, \bar{\sigma}_{*}^{2}, \bar{\pi})$$
$$\boldsymbol{g}_{t+1} = \boldsymbol{Y}^{\mathsf{T}} \boldsymbol{u}_{t}$$

#### Iterative refinement via AMP

As  $\widehat{v} \prec g$  for v, we expect  $Y\widehat{v} \prec Yg \propto f$  for u.

$$\boldsymbol{v}_{t} = \theta(\boldsymbol{g}_{t} | \boldsymbol{\mu}_{*}, \sigma_{*}^{2}, \pi)$$
$$\boldsymbol{f}_{t} = \boldsymbol{Y} \boldsymbol{v}_{t} - \boldsymbol{b}_{t} \boldsymbol{u}_{t-1}$$
$$\boldsymbol{u}_{t} = \theta(\boldsymbol{f}_{t} | \boldsymbol{\mu}_{*}, \boldsymbol{\sigma}_{*}^{2}, \boldsymbol{\pi})$$
$$\boldsymbol{g}_{t+1} = \boldsymbol{Y}^{\mathsf{T}} \boldsymbol{u}_{t} - \boldsymbol{\overline{b}}_{t} \boldsymbol{v}_{t}$$

# Theoretical guarantee of EB-PCA (informal)

Suppose  $u_1, \ldots, u_m \sim \overline{\pi}_*$  and  $v_1, \ldots, v_m \sim \pi_*$  for some priors  $\pi_*, \overline{\pi}_*$ , and the noise W has i.i.d. entries.

As  $n, d \to \infty$ , EB-PCA consistently estimates  $\pi_*, \overline{\pi}_*$ . Furthermore, for every iterate  $f_t$  and  $g_t$  of EB-PCA,

$$\mathbf{f}_{t} \sim N(\bar{\mu}_{t} * \boldsymbol{u}, \bar{\sigma}_{t}^{2}), \quad \mathbf{g}_{t} \sim N(\mu_{t} * \boldsymbol{v}, \sigma_{t}^{2}),$$

where  $\bar{\mu}_t$ ,  $\bar{\sigma}_t^2$ ,  $\mu_t$  and  $\sigma_t^2$  are deterministic sequence which can be prescribed by AMP state evolution.

#### EB-PCA method

- 1. Get samples PCs and estimate their alignments to the true PCs.
- 2. Use (NP)MLE to estimate the prior distribution of the true PCs.
- 3. Follow the AMP algorithm to get improved estimates iteratively.

#### Comparison

- Standard PCA
- Sparse PCA
- Mean-field Variational Bayes [Wang, Stephens 2021]
- James Stein denoisers



#### **Simulation Studies**



#### Rotationally invariant ensemble

#### **Stock returns**



#### HapMap3



# Universality results

- Heavy-tailed entries
- Heterogeneous entries