

# Matrix Denoising with **Doubly Heteroscedastic Noise**

## Fundamental Limits and Optimal Spectral Methods

**Yihan Zhang**\*



**Marco Mondelli**\*



\*Institute of Science and Technology Austria

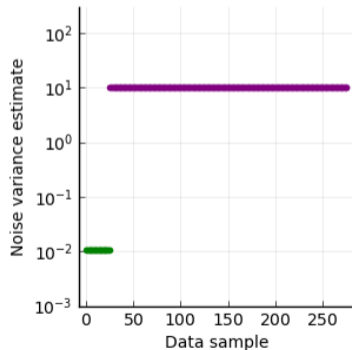
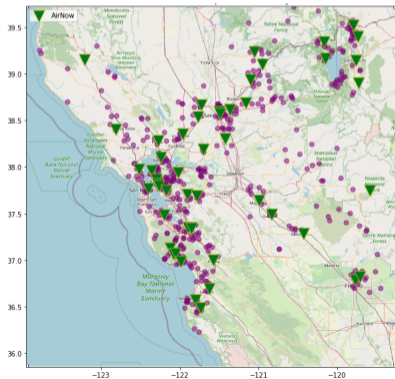
NeurIPS 2024



## Example 1: spatial heteroscedasticity

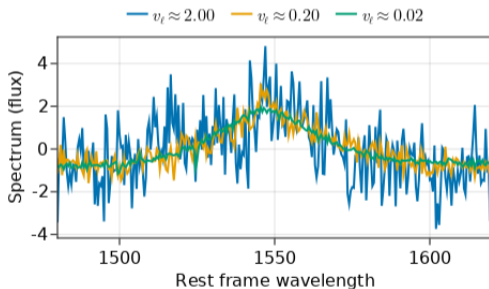
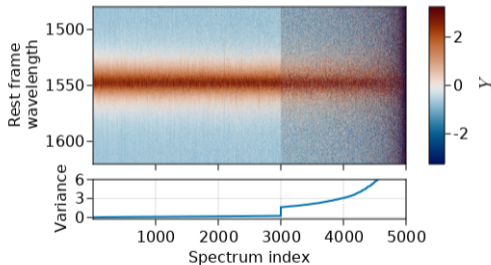
**Air quality** [HGBF21]: hourly readings of  $\text{PM}_{2.5}$  particulate density ( $\mu\text{g}/\text{m}^3$ ) in central California during Feb 9-13, 2021

- AirNow: 25 high-precision sensors
- PurpleAir: 250 low-precision sensors



## Example 2: temporal heteroscedasticity

**Star spectra** [APA<sup>+</sup>20, LHM<sup>+</sup>20, HYFB23]:  $n = 5000$  spectra for a fixed star, each a  $d = 281$ -dimensional vector recording the flux at  $d$  wavelengths in  $[1480, 1620]$  nm



# Matrix denoising with doubly heteroscedastic noise

**Yihan Zhang** and Marco Mondelli. "Matrix Denoising with Doubly Heteroscedastic Noise: Fundamental Limits and Optimal Spectral Methods." arXiv:2405.13912

$$A = \frac{\lambda}{n} u^* v^{*\top} + W \in \mathbb{R}^{n \times d}$$

where

$$W = \Xi^{1/2} \widetilde{W} \Sigma^{1/2}, \quad \widetilde{W}_{i,j} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1/n)$$

**proportional** regime:  $\frac{n}{d} \rightarrow \delta$

# Bayes-optimality

Theorem (“Information-theoretic limit” [ZM24])

$$\lim_{n \rightarrow \infty} \sup_{\tilde{u}} \mathbb{E} \left[ \frac{|\langle \Xi^{-1/2} u^*, \tilde{u} \rangle|^2}{\|\Xi^{-1/2} u^*\|_2^2 \|\tilde{u}\|_2^2} \right] = \frac{q_u^*}{\mathbb{E}[\bar{\Xi}^{-1}]},$$
$$\lim_{d \rightarrow \infty} \sup_{\tilde{v}} \mathbb{E} \left[ \frac{|\langle \Sigma^{-1/2} v^*, \tilde{v} \rangle|^2}{\|\Sigma^{-1/2} v^*\|_2^2 \|\tilde{v}\|_2^2} \right] = \frac{q_v^*}{\mathbb{E}[\bar{\Sigma}^{-1}]},$$

where  $q_u^*, q_v^*$  solve

$$q_u^* = \mathbb{E} \left[ \frac{\lambda^2 q_v^* \bar{\Xi}^{-2}}{\delta + \lambda^2 q_v^* \bar{\Xi}^{-1}} \right], \quad q_v^* = \mathbb{E} \left[ \frac{\lambda^2 q_u^* \bar{\Sigma}^{-2}}{1 + \lambda^2 q_u^* \bar{\Sigma}^{-1}} \right]$$

# Weak recovery threshold

## Corollary

Positive overlaps with  $u^*, v^*$  are achievable if and only if

$$\lambda > \left( \delta \mathbb{E}[\bar{\Xi}^{-2}]^{-1} \mathbb{E}[\bar{\Sigma}^{-2}]^{-1} \right)^{\frac{1}{4}}$$

## Remark

BBP threshold [BBAP05] for whitened  $\Xi^{-1/2} A \Sigma^{-1/2}$ :

$$\lambda > \left( \delta \mathbb{E}[\bar{\Xi}^{-1}]^{-2} \mathbb{E}[\bar{\Sigma}^{-1}]^{-2} \right)^{\frac{1}{4}}$$

# Weak recovery threshold

## Corollary

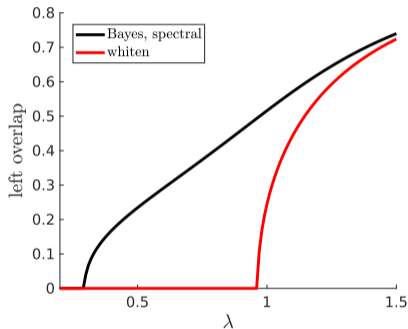
Positive overlaps with  $u^*, v^*$  are achievable if and only if

$$\lambda > \left( \delta \mathbb{E}[\bar{\Xi}^{-2}]^{-1} \mathbb{E}[\bar{\Sigma}^{-2}]^{-1} \right)^{\frac{1}{4}}$$

## Remark

BBP threshold [BBAP05] for whitened  $\Xi^{-1/2} A \Sigma^{-1/2}$ :

$$\lambda > \left( \delta \mathbb{E}[\bar{\Xi}^{-1}]^{-2} \mathbb{E}[\bar{\Sigma}^{-1}]^{-2} \right)^{\frac{1}{4}}$$



## A spectral estimator on steroids

### Theorem (Efficient algorithm [ZM24])

The estimator

$$\hat{u}^{\text{spec}} = G_2(\Xi)u_1\left(\underbrace{G_1(\Xi)AF_1(\Sigma)}_{A^*}\right)$$

- achieves the Bayes-optimal weak recovery threshold;



## A spectral estimator on steroids

### Theorem (Efficient algorithm [ZM24])

The estimator

$$\hat{u}^{\text{spec}} = G_2(\Xi)u_1\left(\underbrace{G_1(\Xi)AF_1(\Sigma)}_{A^*}\right)$$

- achieves the Bayes-optimal weak recovery threshold;
- achieves the Bayes-optimal overlap if  $\Xi = I_n$ .

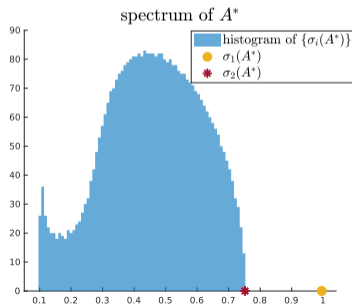
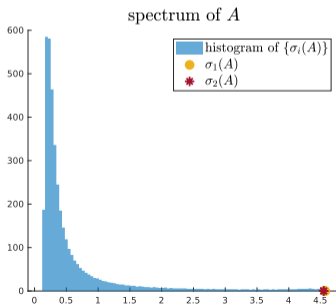
# A spectral estimator on steroids

## Theorem (Efficient algorithm [ZM24])

The estimator

$$\hat{u}^{\text{spec}} = G_2(\Xi)u_1\left(\underbrace{G_1(\Xi)AF_1(\Sigma)}_{A^*}\right)$$

- achieves the Bayes-optimal weak recovery threshold;
- achieves the Bayes-optimal overlap if  $\Xi = I_n$ .



**THANK YOU!**

arXiv:2405.13912

zephyr.z798@gmail.com

-  Romina Ahumada, Carlos Allende Prieto, Andrés Almeida, Friedrich Anders, et al.  
The 16th data release of the sloan digital sky surveys: First release from the apogee-2 southern survey and full release of eboss spectra.  
*The Astrophysical Journal Supplement Series*, 249(1):3, jun 2020.
-  Jinho Baik, Gérard Ben Arous, and Sandrine Péché.  
Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices.  
*Ann. Probab.*, 33(5):1643–1697, 2005.
-  David Hong, Kyle Gilman, Laura Balzano, and Jeffrey A. Fessler.  
HePPCAT: probabilistic PCA for data with heteroscedastic noise.  
*IEEE Trans. Signal Process.*, 69:4819–4834, 2021.
-  David Hong, Fan Yang, Jeffrey A. Fessler, and Laura Balzano.  
Optimally weighted PCA for high-dimensional heteroscedastic data.  
*SIAM J. Math. Data Sci.*, 5(1):222–250, 2023.

 Brad W. Lyke, Alexandra N. Higley, J. N. McLane, Danielle P. Schurhammer, et al.

The sloan digital sky survey quasar catalog: Sixteenth data release.  
*The Astrophysical Journal Supplement Series*, 250(1):8, aug 2020.

 Yihan Zhang and Marco Mondelli.

Matrix denoising with doubly heteroscedastic noise: Fundamental limits and optimal spectral methods.  
*arXiv preprint arXiv:2405.13912*, 2024.