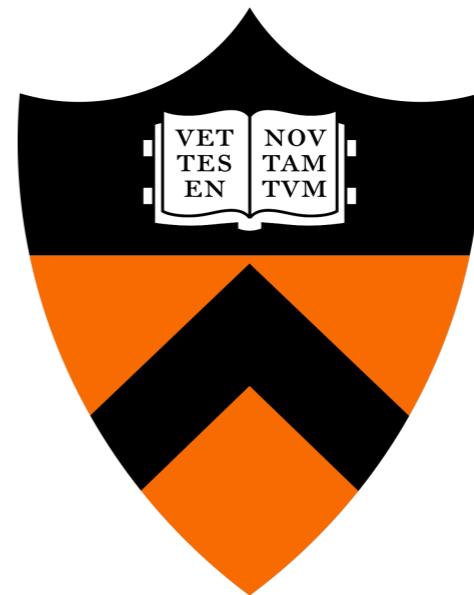


Spectral Methods for Latent Variable Models

Kaizheng Wang

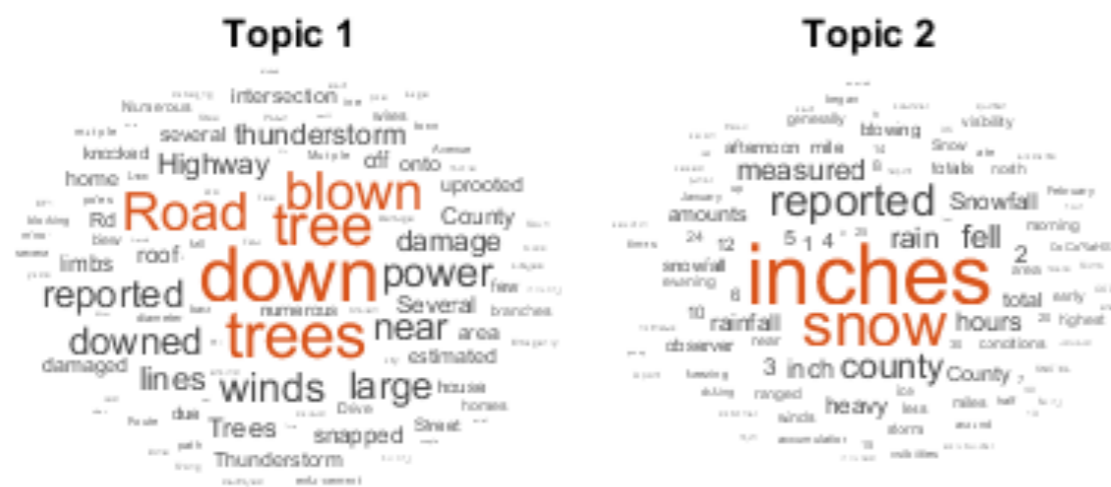
Department of ORFE
Princeton University

March 20th 2020



Data Diversity

Unstructured, heterogeneous and incomplete information:



Matrix Representations

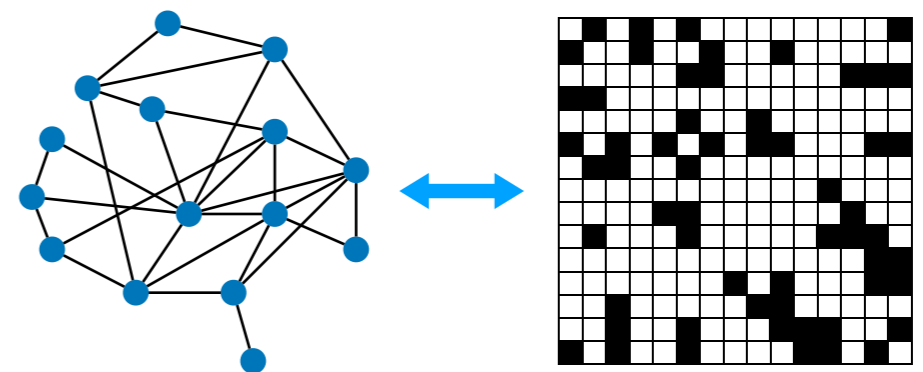
Object-by-feature ($n \times d$):

- **Texts:** document-term;
- **Genetics:** individual-marker;
- **Recomm. systems:** user-item.



Object-by-object ($n \times n$):

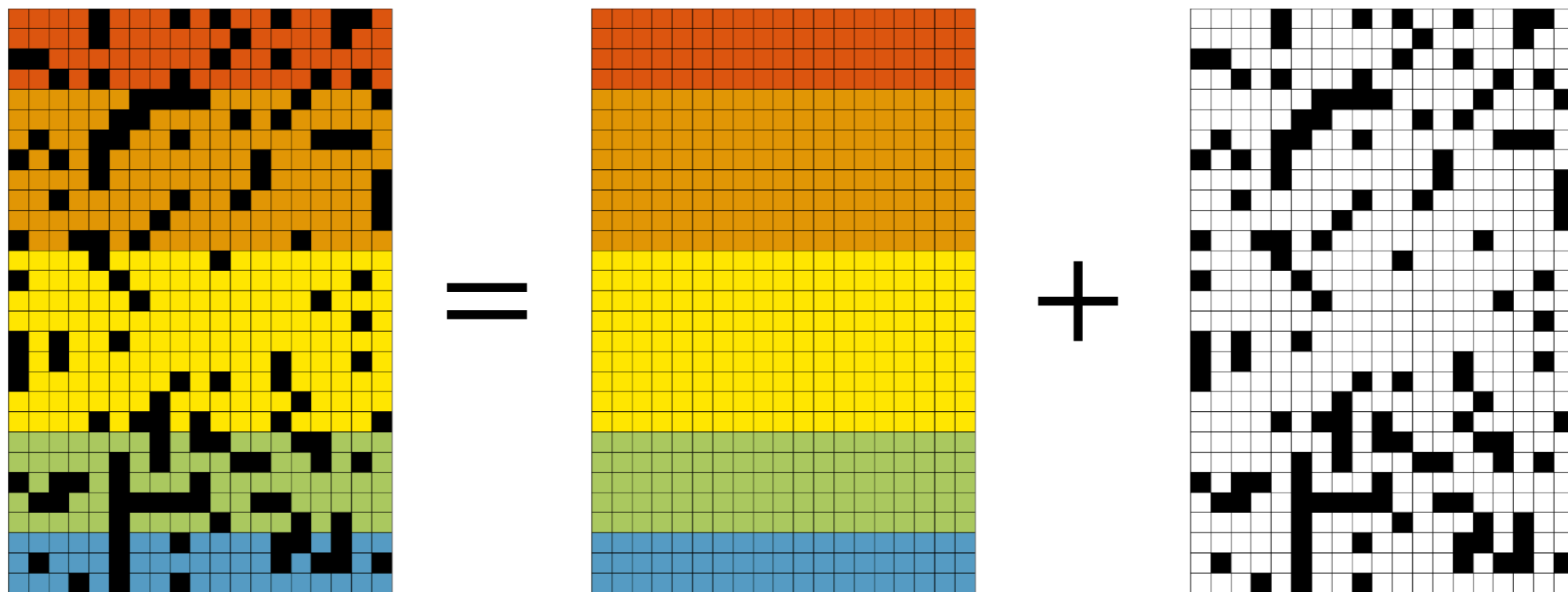
- **Networks:** adjacency matrices.



Matrix Representations

Common belief: **high** ambient dim. but **low** intrinsic dim.

Low-rank approximation:



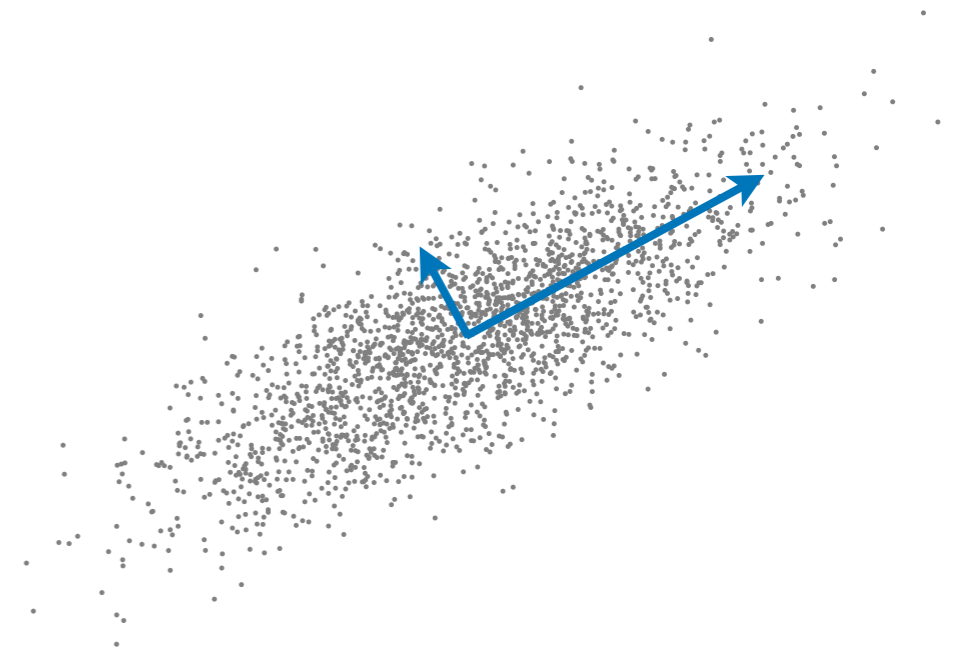
Matrix Representations

Low-dimensional embedding via **latent variables**:

$$\begin{pmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_n \end{pmatrix} = \begin{pmatrix} f_1 \\ \vdots \\ f_n \end{pmatrix} \mathbf{B} + \begin{pmatrix} \mathbf{E}_1 \\ \vdots \\ \mathbf{E}_n \end{pmatrix} .$$

$n \times d$ $n \times r$ $r \times d$ $n \times d$
samples **latent** latent noises
 coordinates bases

Principal Component Analysis (PCA)
— truncated SVD

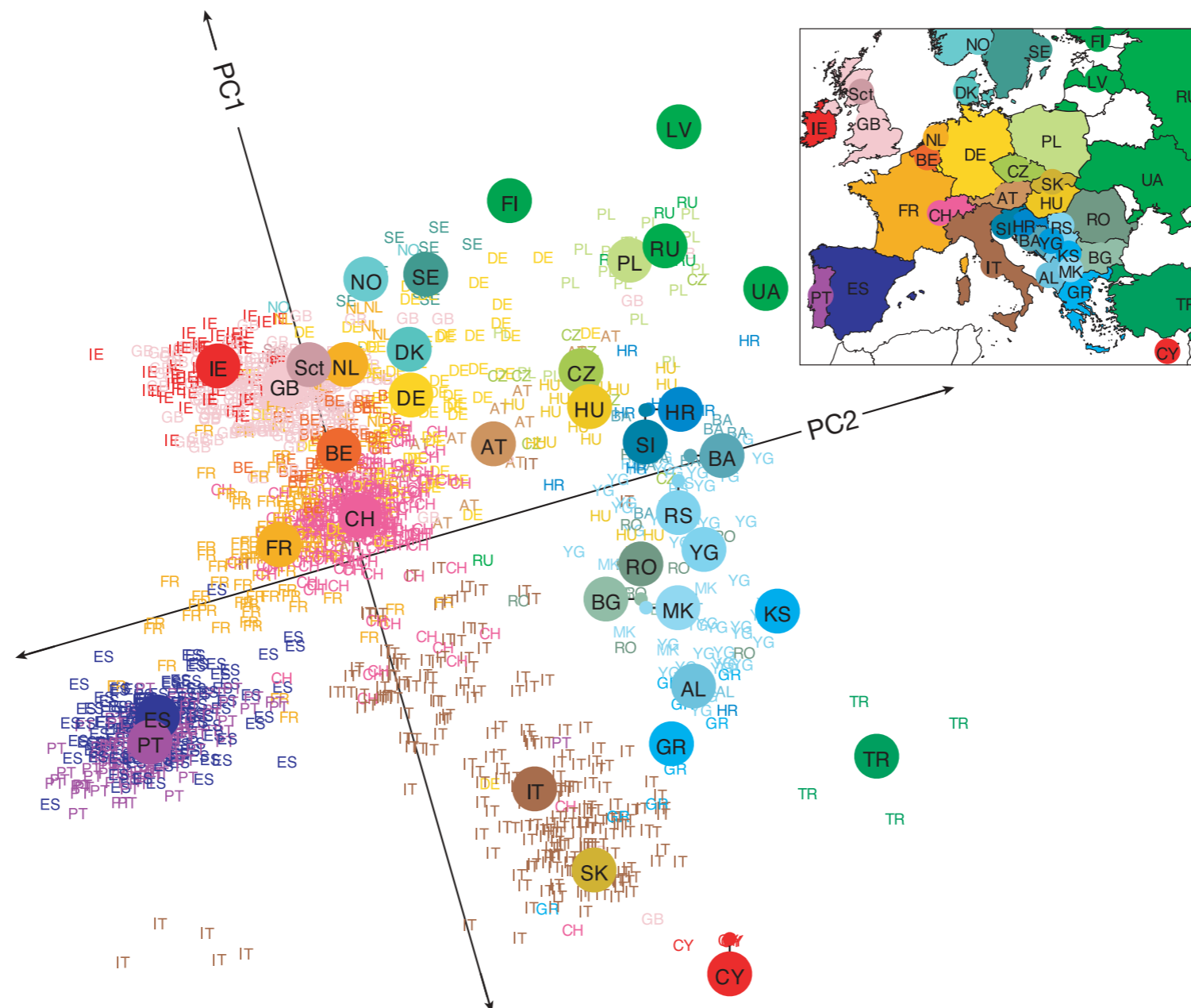


Example: Genes Mirror Geography within Europe

Novembre et al. (2008), Nature.

$n = 1387$ individuals and $d = 197146$ SNPs;

Figure 1a: 2-dim. embedding vs. labels.



Outline

- Distributed PCA and linearization of eigenvectors
- An ℓ_p theory for spectral methods
- Summary and future directions



Distributed PCA and linearization of eigenvectors

Principal Component Analysis

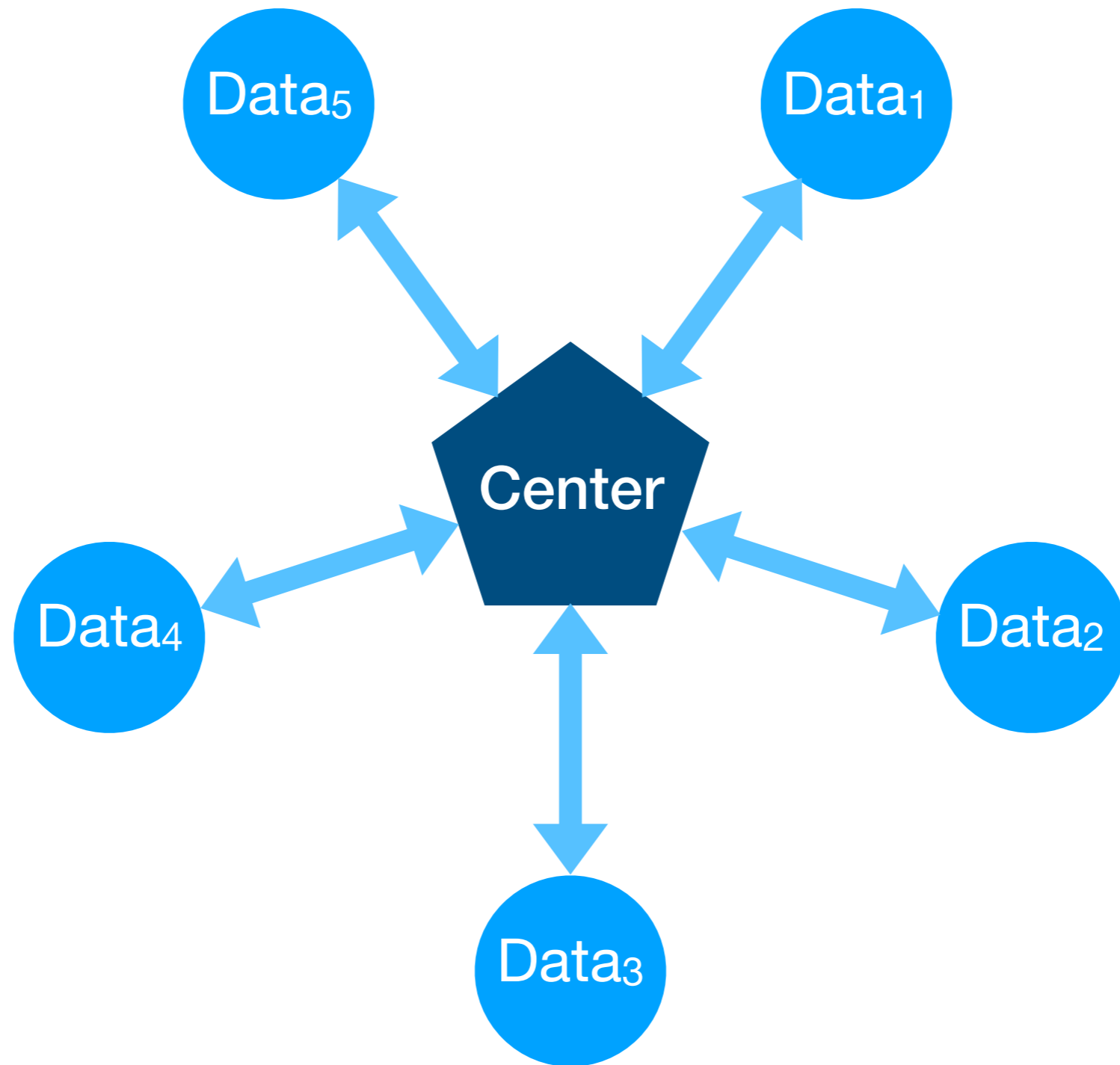
Data: $\{\mathbf{X}_i\}_{i=1}^n \subseteq \mathbb{R}^d$ i.i.d., $\mathbb{E}\mathbf{X}_i = \mathbf{0}$, $\mathbb{E}(\mathbf{X}_i\mathbf{X}_i^\top) = \Sigma$.

Goal: estimate the principal subspace spanned by the K leading eigenvectors of Σ .

PCA: $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)^\top \xrightarrow{\text{SVD}} \hat{\mathbf{U}} = (\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_K) \in \mathbb{R}^{d \times K}$



Distributed PCA



Distributed PCA

m local machines in total, each has n samples.

1. PCA in parallel: the ℓ -th machine conducts

$$\mathbf{X}^{(\ell)} \in \mathbb{R}^{n \times d} \xrightarrow{\text{SVD}} \hat{\mathbf{U}}^{(\ell)} \in \mathbb{R}^{d \times K}$$

and sends $\hat{\mathbf{U}}^{(\ell)}$ to the central server;

2. Aggregation: $\{\hat{\mathbf{U}}^{(\ell)}\}_{\ell=1}^m \longrightarrow \hat{\mathbf{U}} \in \mathbb{R}^{d \times K}$.

Related works: McDonald et al. 2009; Zhang et al. 2013; Lee et al., 2015; Battey et al., 2015; Qu et al. 2002; El Karoui and d'Aspremont 2010; Liang et al. 2014.

Center of Subspaces

How to find $\hat{U} \in \mathcal{O}_{d \times K}$ that best summarizes $\{\hat{U}^{(\ell)}\}_{\ell=1}^m$?

Center of Subspaces

How to find $\hat{U} \in \mathcal{O}_{d \times K}$ that best summarizes $\{\hat{U}^{(\ell)}\}_{\ell=1}^m$?

- **Subspace distance:**

$$\rho(V, W) = \|VV^\top - WW^\top\|_F.$$

- **Least squares:**

$$\hat{U} = \operatorname{argmin}_{V \in \mathcal{O}_{d \times K}} \sum_{\ell=1}^m \rho^2(V, \hat{U}^{(\ell)}).$$

- **Algorithm:**

$$(\hat{U}^{(1)}, \dots, \hat{U}^{(m)}) \in \mathbb{R}^{d \times mK} \xrightarrow{\text{SVD}} \hat{U} \in \mathcal{O}_{d \times K}.$$

Theoretical Results

Assume that \mathbf{X}_i is sub-Gaussian, i.e. $\|\Sigma^{-1/2}\mathbf{X}_i\|_{\psi_2} \lesssim 1$.

Define the effective rank and condition number as

$$r = \text{Tr}(\Sigma)/\lambda_1, \quad \kappa = \lambda_1/(\lambda_K - \lambda_{K+1}).$$

Theoretical Results

Assume that \mathbf{X}_i is sub-Gaussian, i.e. $\|\Sigma^{-1/2}\mathbf{X}_i\|_{\psi_2} \lesssim 1$.

Define the effective rank and condition number as

$$r = \text{Tr}(\Sigma)/\lambda_1, \quad \kappa = \lambda_1/(\lambda_K - \lambda_{K+1}).$$

Theorem (FWWZ, AoS 2019)

There exists constant C such that when $n \geq C\kappa^2\sqrt{K}r$,

$$\left\| \|\hat{\mathbf{U}}\hat{\mathbf{U}}^\top - \mathbf{U}\mathbf{U}^\top\|_{\text{F}} \right\|_{\psi_1} \lesssim \underbrace{\kappa\sqrt{\frac{Kr}{mn}}}_{\text{variance}} + \underbrace{\kappa^2\frac{\sqrt{K}r}{n}}_{\text{bias}}.$$

Theoretical Results

Assume that \mathbf{X}_i is sub-Gaussian, i.e. $\|\Sigma^{-1/2}\mathbf{X}_i\|_{\psi_2} \lesssim 1$.

Define the effective rank and condition number as

$$r = \text{Tr}(\Sigma)/\lambda_1, \quad \kappa = \lambda_1/(\lambda_K - \lambda_{K+1}).$$

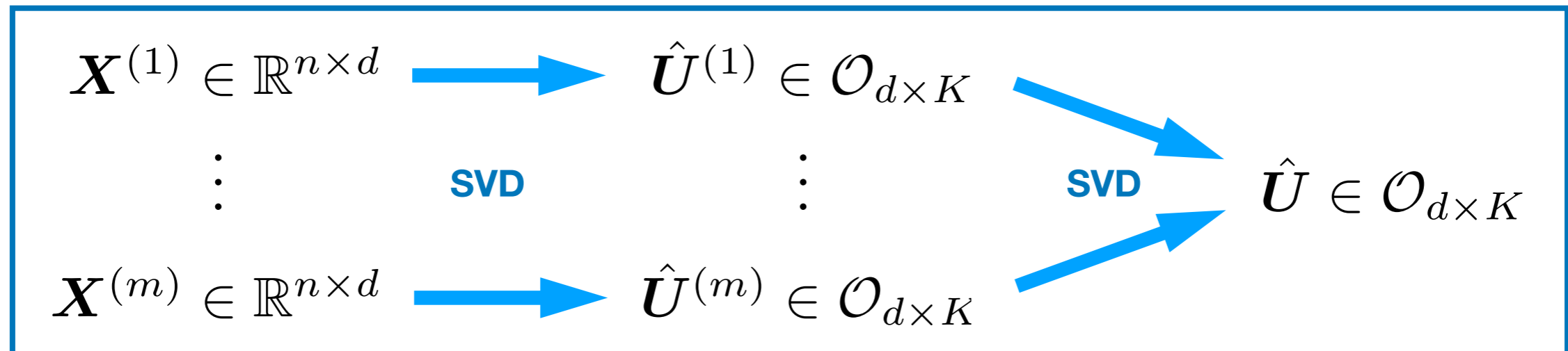
Theorem (FWWZ, AoS 2019)

There exists constant C such that when $n \geq C\kappa^2\sqrt{K}r$,

$$\left\| \|\hat{\mathbf{U}}\hat{\mathbf{U}}^\top - \mathbf{U}\mathbf{U}^\top\|_{\text{F}} \right\|_{\psi_1} \lesssim \underbrace{\kappa\sqrt{\frac{Kr}{mn}}}_{\text{variance}} + \underbrace{\kappa^2\frac{\sqrt{K}r}{n}}_{\text{bias}}.$$

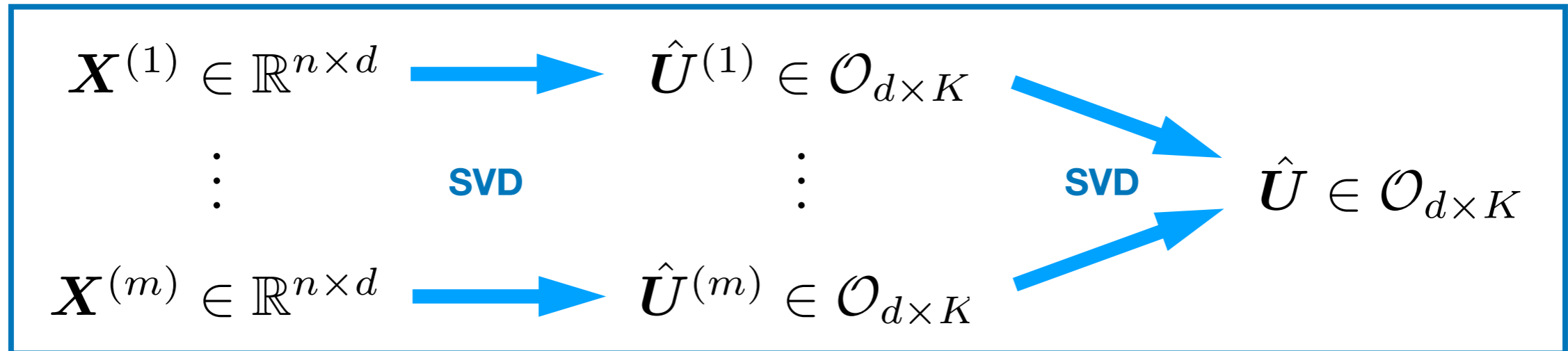
- If $m \lesssim n/(\kappa^2 r)$, distributed PCA is optimal.
- The condition $n \geq C\kappa^2\sqrt{K}r$ cannot be improved.

Analysis of Aggregation



$\hat{\mathbf{U}}$: eigenvectors of $\frac{1}{m} \sum_{\ell=1}^m \hat{\mathbf{U}}^{(\ell)} \hat{\mathbf{U}}^{(\ell)\top}$.

Analysis of Aggregation



\hat{U} : eigenvectors of $\frac{1}{m} \sum_{\ell=1}^m \hat{U}^{(\ell)} \hat{U}^{(\ell)\top}$.

Averaging **reduces variance** but **retains bias**.

- Variance: controlled by Davis-Kahan:

$$\|\hat{U}^{(\ell)} \hat{U}^{(\ell)\top} - UU^\top\|_{\text{F}} \lesssim \|(\hat{\Sigma}^{(\ell)} - \Sigma)U\|_{\text{F}} / \Delta.$$

- Bias: how large it is?

Linearization of Eigenvectors

Theorem (FWWZ, AoS 2019)

$$\|\hat{U}^{(\ell)}\hat{U}^{(\ell)\top} - [UU^\top + f(\hat{\Sigma}^{(\ell)} - \Sigma)]\|_F \lesssim [\|(\hat{\Sigma}^{(\ell)} - \Sigma)U\|_F / \Delta]^2,$$

$f : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d \times d}$ is a **linear** functional determined by Σ .

Linearization of Eigenvectors

Theorem (FWWZ, AoS 2019)

$$\|\hat{U}^{(\ell)}\hat{U}^{(\ell)\top} - [UU^\top + f(\hat{\Sigma}^{(\ell)} - \Sigma)]\|_F \lesssim [\|(\hat{\Sigma}^{(\ell)} - \Sigma)U\|_F/\Delta]^2,$$

$f : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d \times d}$ is a **linear** functional determined by Σ .

More precise than **Davis-Kahan**:

$$\|\hat{U}^{(\ell)}\hat{U}^{(\ell)\top} - UU^\top\|_F \lesssim \|(\hat{\Sigma}^{(\ell)} - \Sigma)U\|_F/\Delta.$$

PCA has small bias:

$$\|\mathbb{E}(\hat{U}^{(\ell)}\hat{U}^{(\ell)\top}) - UU^\top\|_F \lesssim [\|(\hat{\Sigma}^{(\ell)} - \Sigma)U\|_F/\Delta]^2.$$

Summary

Theoretical guarantees for **distributed PCA**:

- **Bias** and variance of PCA;
- **Linearization** of eigenvectors, high-order Davis-Kahan.

Paper (alphabetical order):

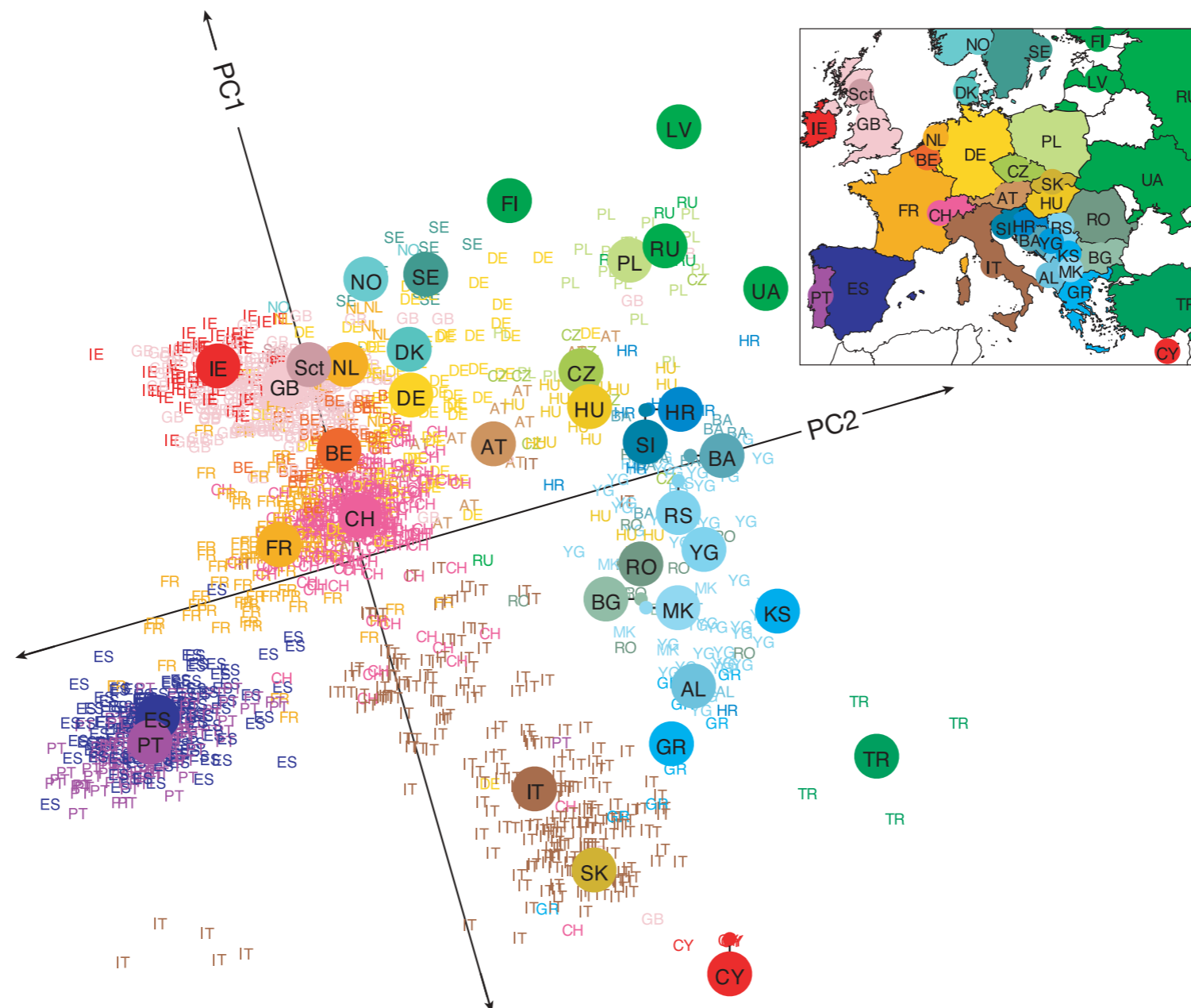
- Fan, Wang, **Wang** and Zhu. Distributed estimation of principal eigenspaces. **The Annals of Statistics**, 2019.

Example: Genes Mirror Geography within Europe

Novembre et al. (2008), Nature.

$n = 1387$ individuals and $d = 197146$ SNPs;

Figure 1a: 2-dim. embedding vs. labels.



A Pipeline for Spectral Methods

1. Similarity matrix construction

e.g. Gram XX^T , adjacency A ;

2. Spectral decomposition

get r eigen-pairs $\{\lambda_j, \mathbf{u}_j\}_{j=1}^r$;

3. r -dim. embedding

e.g. using the rows of $(\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r)$;

4. Downstream tasks

e.g. visualization.

Ext.: { robust, probabilistic, sparse, nonnegative } PCA.

Pearson (1901), Hotelling (1933), Schölkopf (1997), Tipping and Bishop (1999), Shi and Malik (2000), Ng et al. (2002), Belkin and Niyogi (2003), Von Luxburg (2007)

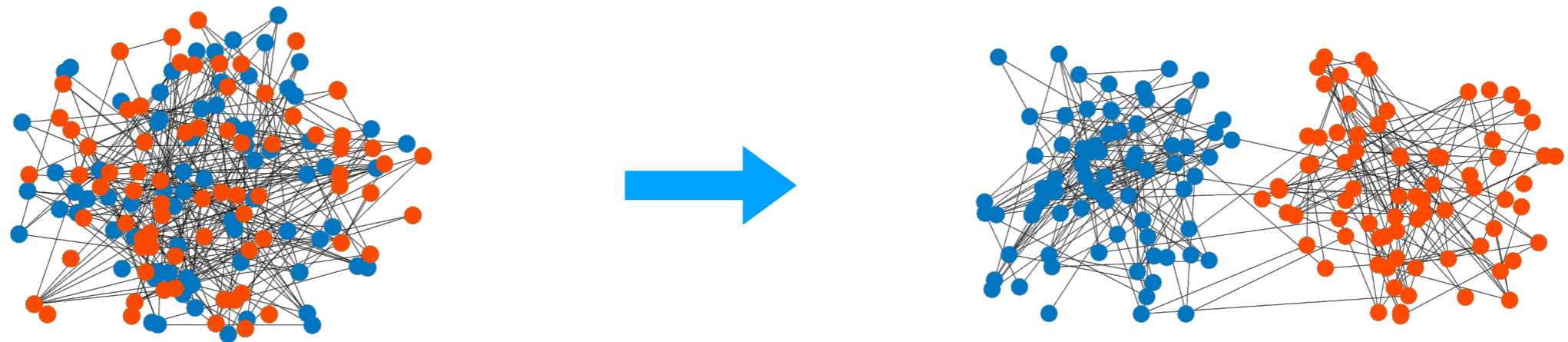


An ℓ_p theory for spectral methods

- **Network analysis and Wigner-type matrices**
- Mixture model and Wishart-type matrices

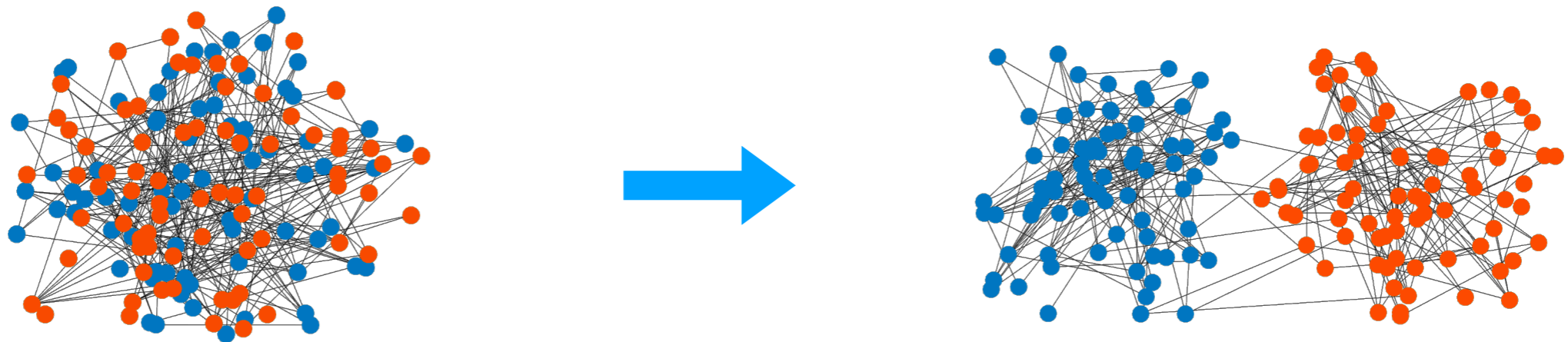
Community Detection and SBM

Community detection in networks:



Community Detection and SBM

Community detection in networks:



Stochastic Block Model (Holland et al., 1983)

Symmetric adjacency matrix $A \in \{0, 1\}^{n \times n}$, $|J| = |J^c| = \frac{n}{2}$:

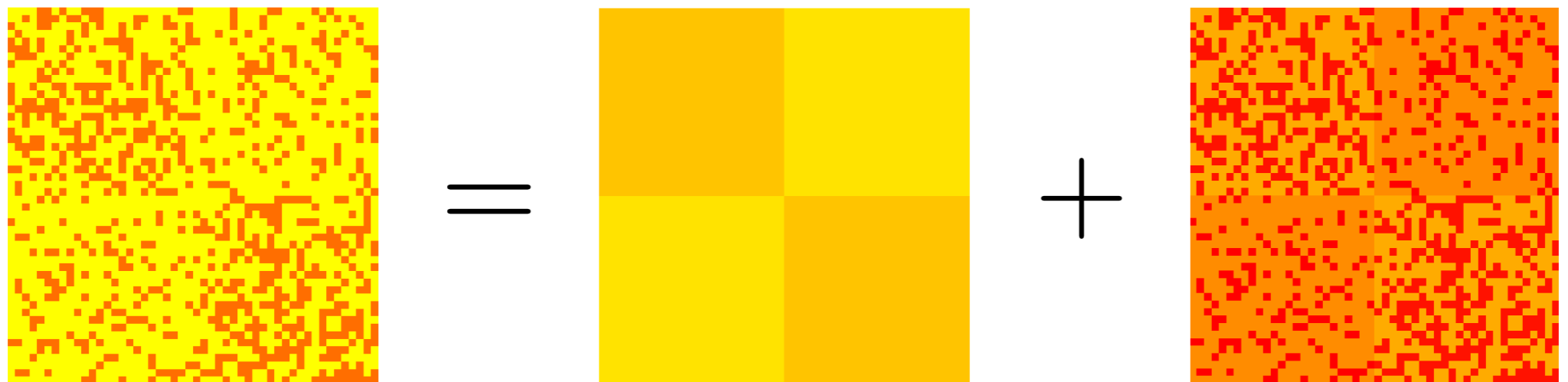
$$\mathbb{P}(A_{ij} = 1) = \begin{cases} p, & \text{if } i, j \text{ or } i, j \\ q, & \text{if } i, j \text{ or } i, j \end{cases}.$$

McSherry (2001), Coja-Oghlan (2006), Rohe et al. (2011), Mossel et al. (2013), Massoulié (2014), Lelarge et al. (2015), Chin et al. (2015), Abbe et al. (2016), Zhang and Zhou (2016).

Community Detection and SBM

$$\mathbb{E}A = \begin{pmatrix} p\mathbf{1}_{J,J} & q\mathbf{1}_{J,J^c} \\ q\mathbf{1}_{J^c,J} & p\mathbf{1}_{J^c,J^c} \end{pmatrix} = \frac{p+q}{2n}\mathbf{1}\mathbf{1}^\top + \frac{p-q}{2n} \begin{pmatrix} \mathbf{1}_J \\ -\mathbf{1}_{J^c} \end{pmatrix} (\mathbf{1}_J^\top \quad -\mathbf{1}_{J^c}^\top).$$

The 2nd eigenvector $\bar{u} = \frac{1}{\sqrt{n}}(\mathbf{1}_J - \mathbf{1}_{J^c})$ reveals (J, J^c) .



$$A = \mathbb{E}A + A - \mathbb{E}A$$

Community Detection and SBM

$$\mathbb{E}A = \begin{pmatrix} p\mathbf{1}_{J,J} & q\mathbf{1}_{J,J^c} \\ q\mathbf{1}_{J^c,J} & p\mathbf{1}_{J^c,J^c} \end{pmatrix} = \frac{p+q}{2n}\mathbf{1}\mathbf{1}^\top + \frac{p-q}{2n} \begin{pmatrix} \mathbf{1}_J \\ -\mathbf{1}_{J^c} \end{pmatrix} \begin{pmatrix} \mathbf{1}_J^\top & -\mathbf{1}_{J^c}^\top \end{pmatrix}.$$

The 2nd eigenvector $\bar{u} = \frac{1}{\sqrt{n}}(\mathbf{1}_J - \mathbf{1}_{J^c})$ reveals (J, J^c) .

Spectral method:

$$A \xrightarrow{\text{SVD}} \text{the 2nd eigenvector } u \longrightarrow \text{sgn}(u)$$

To recover (J, J^c) , we need $u \approx \bar{u}$ in a uniform way.

Classical ℓ_2 bounds (Davis and Kahan, 1970) are too loose!

Optimality of Spectral Method

Theorem (AFWZ, AoS 2020+)

Let $a \neq b$ and $\mathbb{P}(A_{ij} = 1) = \begin{cases} \frac{a \log n}{n}, & \text{if } i, j \text{ or } i, j \\ \frac{b \log n}{n}, & \text{if } i, j \text{ or } i, j \end{cases}$.

- Exact recovery w.h.p. when $(\sqrt{a} - \sqrt{b})^2 > 2$;
- Error rate $n^{-(\sqrt{a} - \sqrt{b})^2/2}$ when $(\sqrt{a} - \sqrt{b})^2 \leq 2$.
- Optimality (Abbe et al., 2016; Zhang and Zhou, 2016).

Optimality of Spectral Method

Theorem (AFWZ, AoS 2020+)

Let $a \neq b$ and $\mathbb{P}(A_{ij} = 1) = \begin{cases} \frac{a \log n}{n}, & \text{if } i, j \text{ or } i, j \\ \frac{b \log n}{n}, & \text{if } i, j \text{ or } i, j \end{cases}$.

- Exact recovery w.h.p. when $(\sqrt{a} - \sqrt{b})^2 > 2$;
- Error rate $n^{-(\sqrt{a} - \sqrt{b})^2/2}$ when $(\sqrt{a} - \sqrt{b})^2 \leq 2$.
- Optimality (Abbe et al., 2016; Zhang and Zhou, 2016).

Key ingredients:

- **Entrywise** linear approximation $\mathbf{u} = \mathbf{A}\mathbf{u}/\lambda \approx \mathbf{A}\bar{\mathbf{u}}/\bar{\lambda}$;
- Weighted sum of independent Bernoulli variables.

l_∞ Analysis: Linearization

Theorem (Linear approximation)

$$\mathbb{P}(\sqrt{n} \|\mathbf{u} - A\bar{\mathbf{u}}/\bar{\lambda}\|_\infty < \varepsilon_n) > 1 - n^{-3}.$$

ℓ_∞ Analysis: Linearization

Theorem (Linear approximation)

$$\mathbb{P}(\sqrt{n} \|\mathbf{u} - \mathbf{A}\bar{\mathbf{u}}/\bar{\lambda}\|_\infty < \varepsilon_n) > 1 - n^{-3}.$$

General results (AFWZ, AoS 2020+)

- Singular vectors of **Wigner-type** matrices:
 - ▶ Symmetric, independent entries above the diagonal;
 - ▶ Rectangular, independent entries;
- $\min_{\mathbf{O} \in \mathcal{O}_{r \times r}} \|\mathbf{U}\mathbf{O} - \mathbf{A}\bar{\mathbf{U}}\bar{\Lambda}^{-1}\|_{2,\infty} \ll \|\bar{\mathbf{U}}\|_{2,\infty}$;
- Applications: synchronization, matrix completion, (inference).

A General ℓ_∞ Theory

Merits and demerits of ℓ_∞ analysis

- 👍 Characterizes individual objects precisely
 - ▶ Results and tools apply to ncvx opt. (MWCC, FoCM 2019);
- 👎 Requires strong signals for **uniform** control:
 - ▶ e.g. degree $\gtrsim \log n$ for SBM.

Successor: ℓ_p analysis with $p < \infty$

- Controls a **vast majority** of the entries.

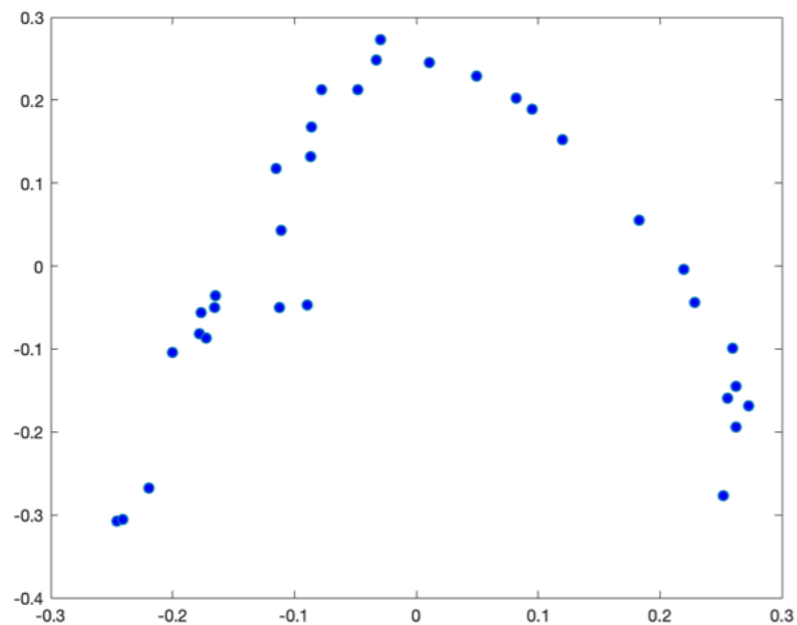
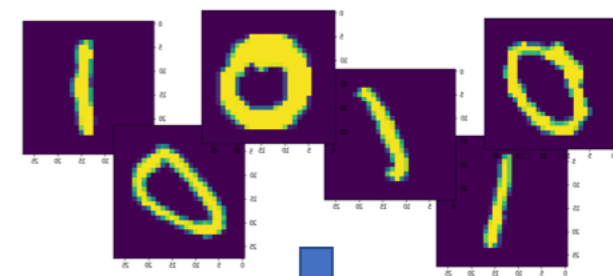
An ℓ_p theory for spectral methods

- Network analysis and Wigner-type matrices
- **Mixture model and Wishart-type matrices**

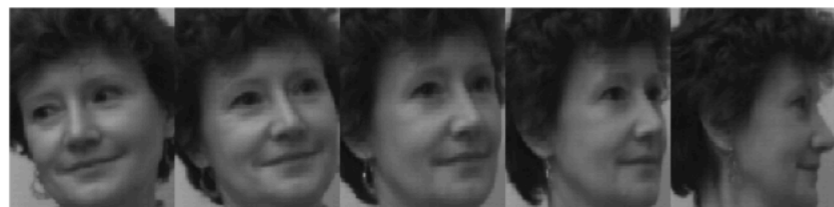
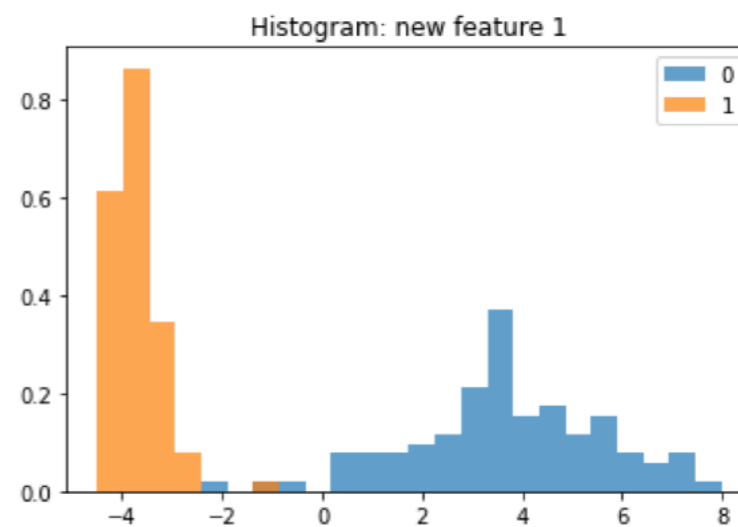
Dimensionality Reduction



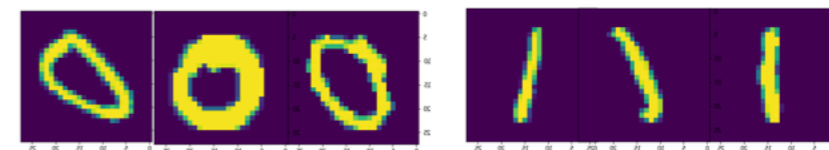
Data



Embedding



Structures



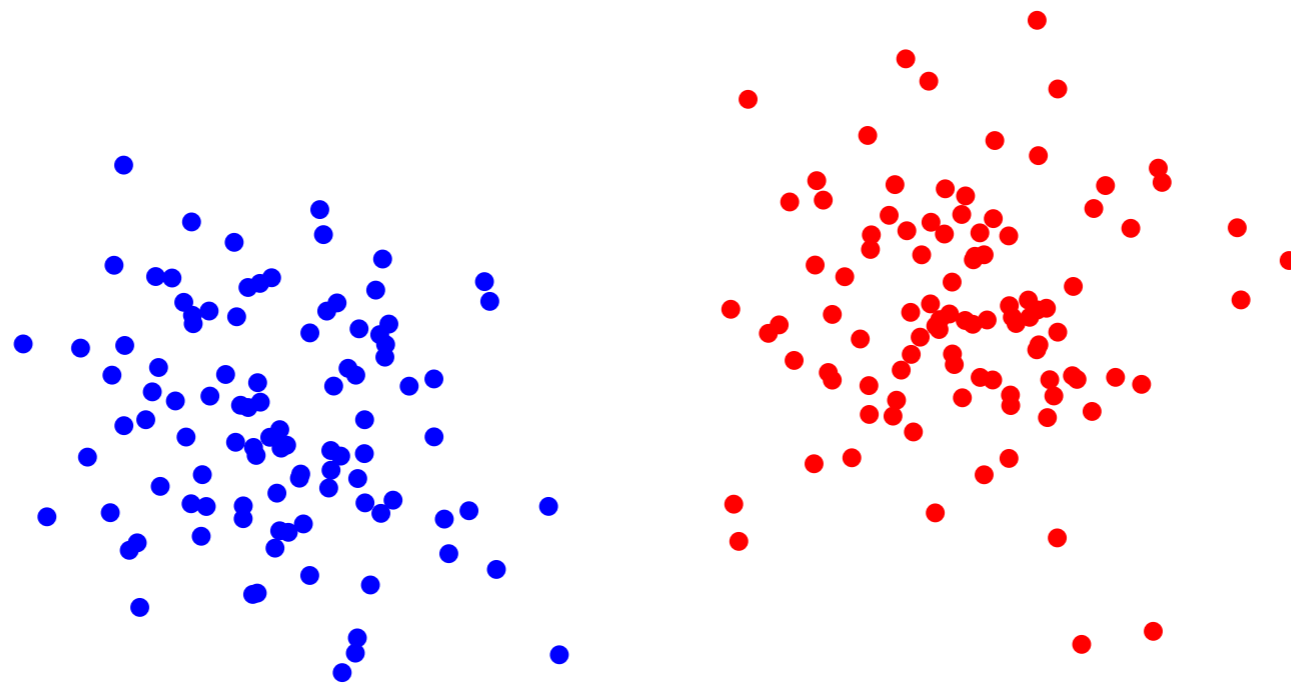
Gaussian Mixture Model

Heteroscedastic model

$$l_i \in \{\pm 1\}, \quad \mathbf{X}_i | l_i \sim \mathcal{N}(l_i \boldsymbol{\mu}, \boldsymbol{\Sigma}_i), \quad \boldsymbol{\Sigma}_i \preceq \boldsymbol{\Sigma}.$$

$$\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)^\top \in \mathbb{R}^{n \times d}$$

$$\text{Low-rank: } \mathbb{E}(\mathbf{X} | \ell) = \ell \boldsymbol{\mu}^\top.$$



Spectral Method

Recall $\mathbb{E}(\mathbf{X} | \ell) = \ell \boldsymbol{\mu}^\top$.

Spectral method

1. Get the hollowed Gram matrix

$$G = \mathcal{H}(\mathbf{X} \mathbf{X}^\top) \in \mathbb{R}^{n \times n};$$

2. $G \xrightarrow{\text{SVD}}$ the 1st eigenvector $\mathbf{u} \longrightarrow \text{sgn}(\mathbf{u})$.

Hollowing

- improves concentration;
- helps tackle heteroscedasticity.

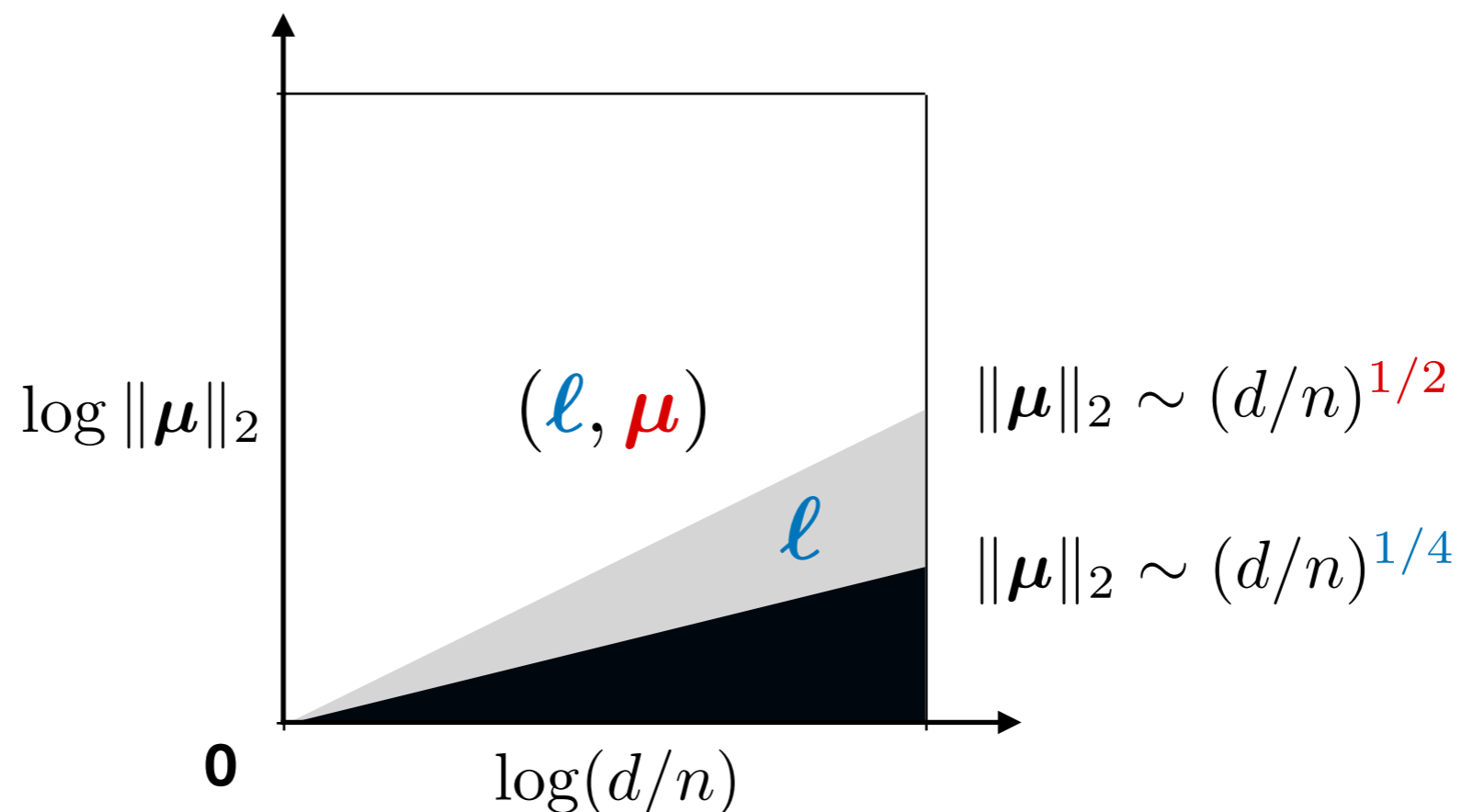
Related works: Montanari and Sun 2018, Ndaoud 2018, Cai et al. 2019.

Theoretical Challenges

- **Dependency:** $\mathcal{H}(X X^\top)$ has **Wishart**-type distribution.
 - ▶ SBM: **Wigner**-type adjacency matrix.
- **High dimensionality:** most existing results require $d \lesssim n$.
 - ▶ **Clustering** vs. **parameter est.:** $\frac{1}{2}N(\boldsymbol{\mu}, \mathbf{I}_d) + \frac{1}{2}N(-\boldsymbol{\mu}, \mathbf{I}_d)$

Theoretical Challenges

- **Dependency:** $\mathcal{H}(X X^\top)$ has **Wishart**-type distribution.
 - ▶ SBM: **Wigner**-type adjacency matrix.
- **High dimensionality:** most existing results require $d \lesssim n$.
 - ▶ **Clustering** vs. **parameter est.:** $\frac{1}{2}N(\boldsymbol{\mu}, \mathbf{I}_d) + \frac{1}{2}N(-\boldsymbol{\mu}, \mathbf{I}_d)$



ℓ_p Analysis: Linear Approximation

\mathbf{u} is the 1st eigenvector of $\mathbf{G} = \mathcal{H}(\mathbf{X}\mathbf{X}^\top) \in \mathbb{R}^{n \times n}$,

$$\text{SNR} = \frac{\|\boldsymbol{\mu}\|_2^4}{\|\boldsymbol{\Sigma}\|_2 \|\boldsymbol{\mu}\|_2^2 + \|\boldsymbol{\Sigma}\|_F^2/n} \gg 1,$$

$\bar{\mathbf{X}} = \ell \boldsymbol{\mu}^\top$, $(\bar{\lambda}, \bar{\mathbf{u}})$ is the 1st eigen-pair of $\bar{\mathbf{X}}\bar{\mathbf{X}}^\top$.

ℓ_p Analysis: Linear Approximation

\mathbf{u} is the 1st eigenvector of $\mathbf{G} = \mathcal{H}(\mathbf{X}\mathbf{X}^\top) \in \mathbb{R}^{n \times n}$,

$$\text{SNR} = \frac{\|\boldsymbol{\mu}\|_2^4}{\|\boldsymbol{\Sigma}\|_2 \|\boldsymbol{\mu}\|_2^2 + \|\boldsymbol{\Sigma}\|_F^2/n} \gg 1,$$

$\bar{\mathbf{X}} = \ell\boldsymbol{\mu}^\top$, $(\bar{\lambda}, \bar{\mathbf{u}})$ is the 1st eigen-pair of $\bar{\mathbf{X}}\bar{\mathbf{X}}^\top$.

Theorem

If $2 \leq p \lesssim \text{SNR}$, there exists $\varepsilon_n \rightarrow 0$ s.t.

$$\mathbb{P} \left(\|\mathbf{u} - \mathbf{G}\bar{\mathbf{u}}/\bar{\lambda}\|_p < \varepsilon_n \|\bar{\mathbf{u}}\|_p \right) > 1 - e^{-p}.$$

- Applies to **RKHS (dim. = ∞)**: kernel PCA;
- Adaptivity: large SNR \Rightarrow large $p \Rightarrow$ strong result.

ℓ_p Analysis: Corollaries

When $p \asymp \text{SNR} \gtrsim \log n$: $\|\cdot\|_p \asymp \|\cdot\|_\infty$ and

$$\mathbb{P}\left(\|\mathbf{u} - \mathbf{G}\bar{\mathbf{u}}/\bar{\lambda}\|_\infty < \frac{\varepsilon_n}{\sqrt{n}}\right) > 1 - \frac{1}{n^3}.$$

ℓ_p Analysis: Corollaries

When $p \asymp \text{SNR} \gtrsim \log n$: $\|\cdot\|_p \asymp \|\cdot\|_\infty$ and

$$\mathbb{P}\left(\|\mathbf{u} - \mathbf{G}\bar{\mathbf{u}}/\bar{\lambda}\|_\infty < \frac{\varepsilon_n}{\sqrt{n}}\right) > 1 - \frac{1}{n^3}.$$

Corollary (optimal clustering, $\Sigma = I_d$)

- Exact recovery w.h.p. when $\text{SNR} > 2 \log n$;
- Error rate $e^{-\text{SNR}/[2+o(1)]}$ when $1 \ll \text{SNR} \leq 2 \log n$.

- Kmeans (Lu and Zhou 2016);
- Spectral (Vempala and Wang 2004, Jin et al. 2017, Ndaoud 2018, Löffler et al. 2019);
- SDP (Mixon et al. 2016, Royer 2017, Fei and Chen 2018, Giraud and Verzelen 2018, Chen and Yang 2018).

Summary

Sharpness of **spectral methods**.

- Abbe, Fan, **Wang** and Zhong. Entrywise eigenvector analysis of random matrices with low expected rank. **The Annals of Statistics**, 2020+.
- Abbe, Fan and **Wang**. An ℓ_p analysis of kernel PCA and contextual network analysis. Manuscript.

Extensions

- ▶ ranking (CFMW, AoS 2019), topic models, ncvx optimization;
- ▶ statistical inference based on linear representation.

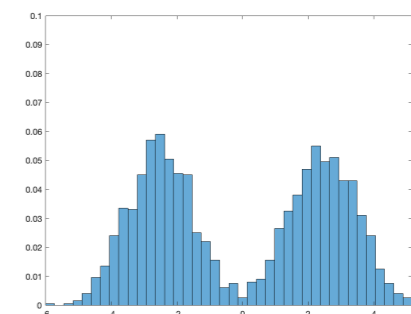
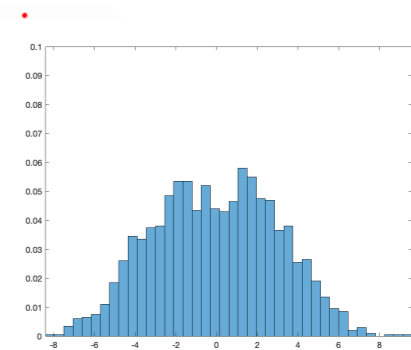
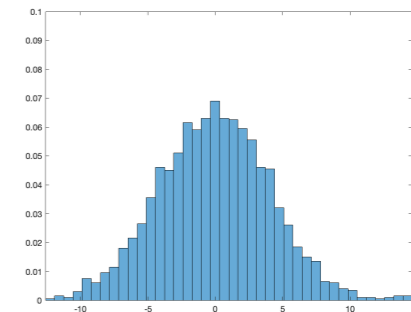
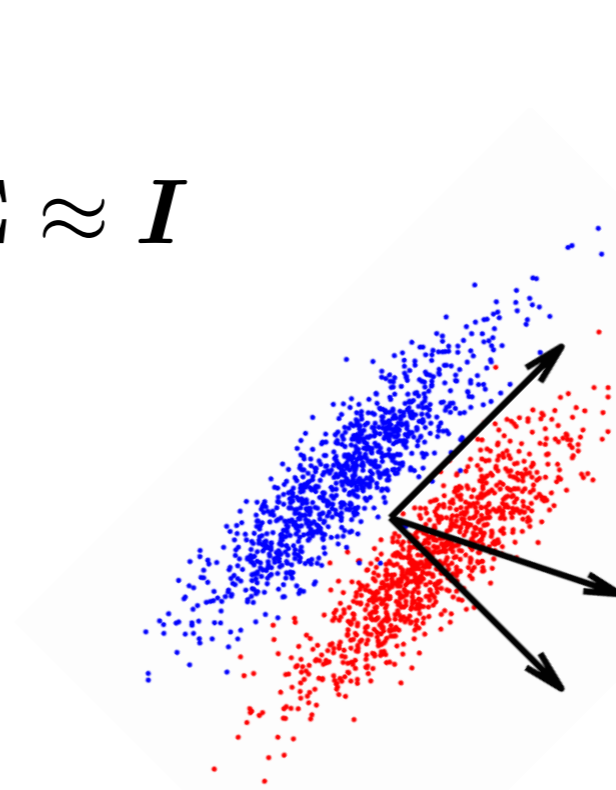
(alphabetical orders)

Finding a Needle in a Haystack

Spectral methods are **powerful** but **not omnipotent**.

$\frac{1}{2}\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + \frac{1}{2}\mathcal{N}(-\boldsymbol{\mu}, \boldsymbol{\Sigma})$: covariance $\boldsymbol{\mu}\boldsymbol{\mu}^\top + \boldsymbol{\Sigma}$

- Max variance \neq useful
- PCA: $\|\boldsymbol{\mu}\|_2^2 / \|\boldsymbol{\Sigma}\|_2 \gg 1$ or $\boldsymbol{\Sigma} \approx \mathbf{I}$

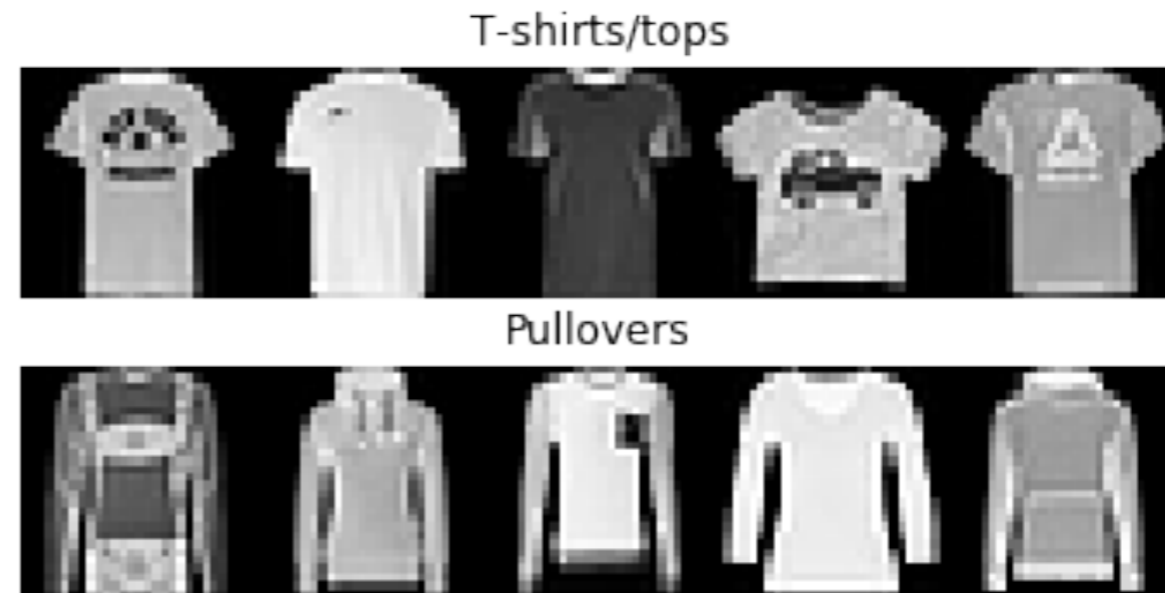


Seek for clustering-friendly projections!

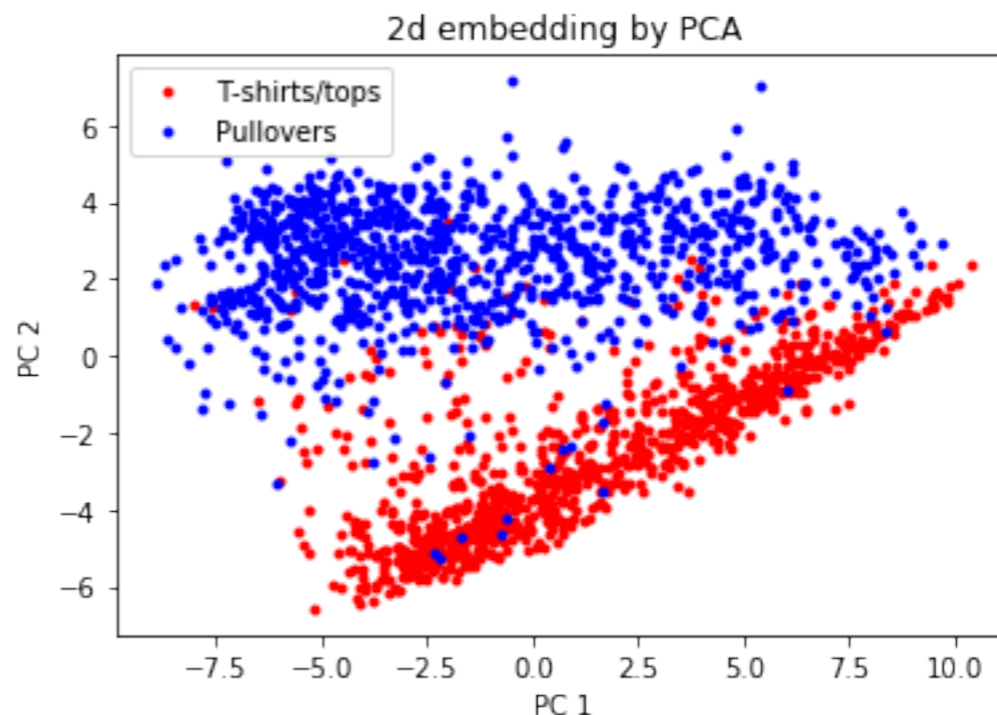
Example: Fashion-MNIST

70000 fashion products, 10 categories (Xiao et al. 2017).

- T-shirts/tops
- Pullovers



Visualization by PCA



A CURE for Clustering Problems

Clustering via Uncoupled REgression (CURE):

Wang, Yan and Diaz. Efficient clustering for stretched mixtures: landscape and optimality. Submitted.

- ▶ **Clustering** -> **classification**;
- ▶ **Stat.** and **comp.** guarantees under mixture models.

Q & A

Thank you!