

Transelliptical Component Analysis

Fang Han

Han Liu

Johns Hopkins University

Princeton University

Neural Information Processing Systems (NIPS)

Lake Tahoe, Nevada, 2012

General Framework

- Using **semiparametric model**
- Obtaining **nonparametric modeling flexibility**
- Achieving **nearly optimal parametric rates**
- Simple and computational efficient

Outline

- **Sparse Principal Component Analysis**
- **Transelliptical Component Analysis**
- **Concluding Remarks**

Sparse Principal Component Analysis

Leading Eigenvectors Estimation

Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be n i.i.d d -variate random vectors with covariance matrix Σ .

Goal: Estimate the top m leading eigenvectors $\theta_1, \dots, \theta_m$ of Σ based on the samples $\{\mathbf{X}_i\}$.

Applications

- Large-scale genomic data
- Brain imaging data
- Equity data
- ...

Sparse Leading Eigenvectors

Let $\{\mathbf{X}_i : i = 1, \dots, n\}$ be n random samples from a d -variate **Gaussian** or **sub-Gaussian** distribution with covariance matrix Σ .

Goal: Estimate θ_1 based on $\{\mathbf{X}_i\}$, while θ_1 is sparse.

Assumption and Estimation

Assumption 1: Assume that $\|\boldsymbol{\theta}_1\|_0 = s < n$.

Estimation:

$$\hat{\boldsymbol{\theta}}_1^* = \arg \max_{\boldsymbol{v} \in \mathbb{R}^d} \boldsymbol{v}^T \boldsymbol{\Sigma}_n \boldsymbol{v}, \quad \text{subject to } \boldsymbol{v} \in \mathbb{S}^{d-1} \cap \mathbb{B}_0(s),$$

where $\boldsymbol{\Sigma}_n$ is the sample covariance matrix, \mathbb{S}^{d-1} is the unit sphere in \mathbb{R}^d and $\mathbb{B}_0(s) := \{\boldsymbol{v} \in \mathbb{R}^d : \|\boldsymbol{v}\|_0 \leq s\}$.

Rates of Convergence

Assumption 2: Assume that \mathbf{X} is Gaussian or sub-Gaussian distributed.

Theorem. Let $\|\boldsymbol{\theta}_1\|_0 = s$. Then the parametric rate of convergence under the ℓ_2 norm is

$$O_P \left(\frac{1}{\lambda_1 - \lambda_2} \sqrt{\frac{s \log d}{n}} \right),$$

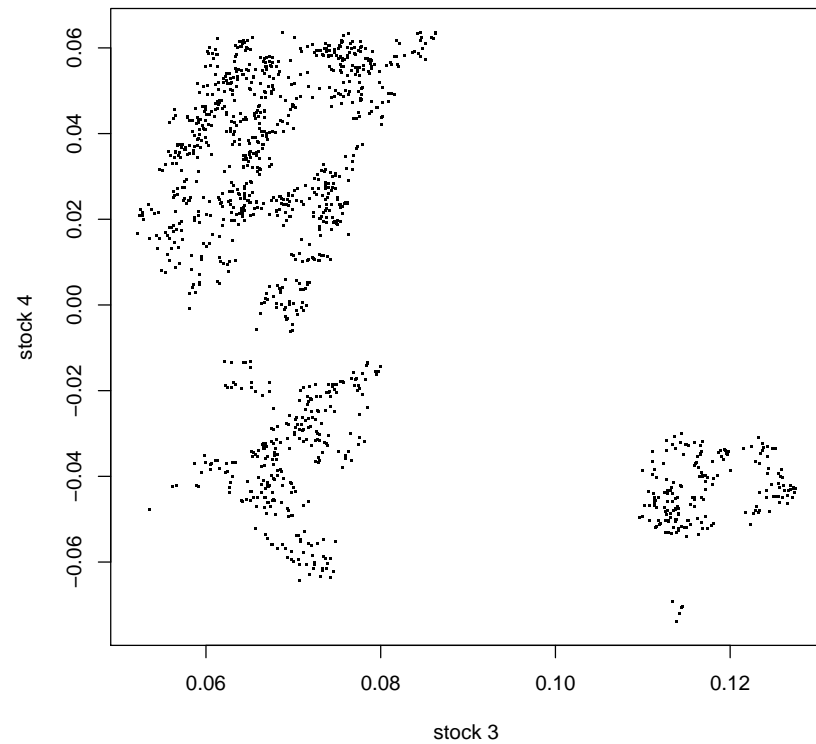
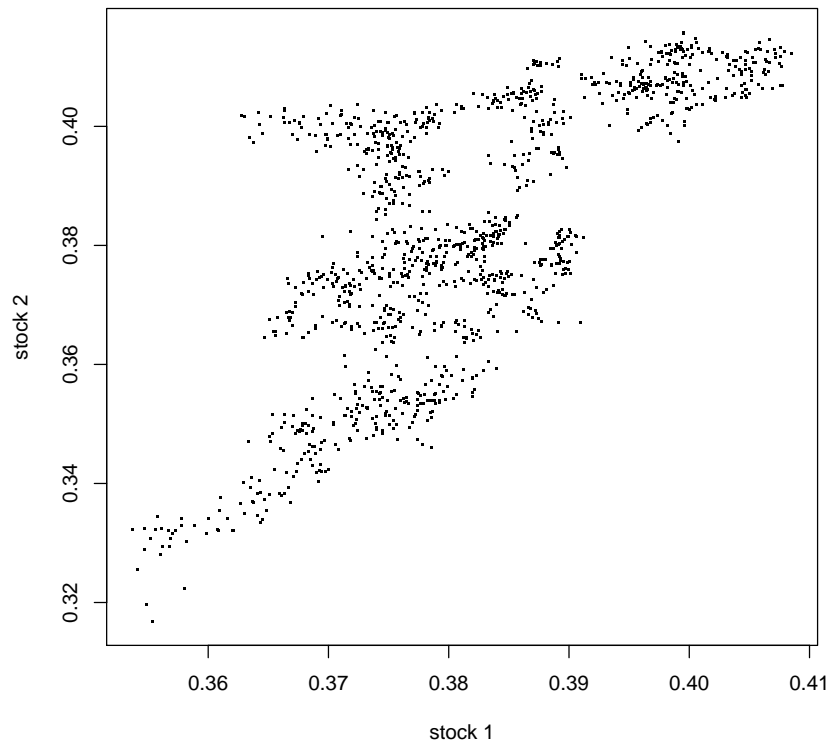
where λ_1 and λ_2 are the top two largest eigenvalues of $\boldsymbol{\Sigma}$.

Constraints

The Gaussian or sub-Gaussian assumption is too **constraint**, appearing to be **rare** in applications.

Equity Data

452 stocks that were consistently in the S&P 500 index between January 1, 2003 though January 1, 2008.

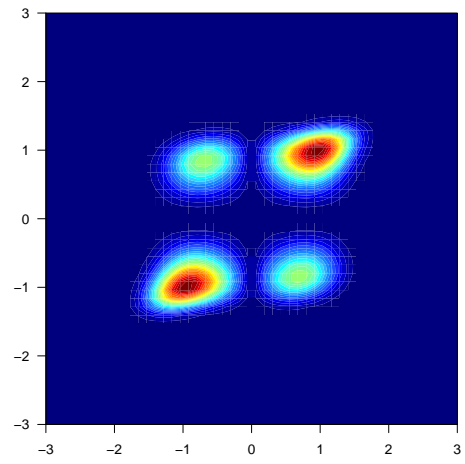
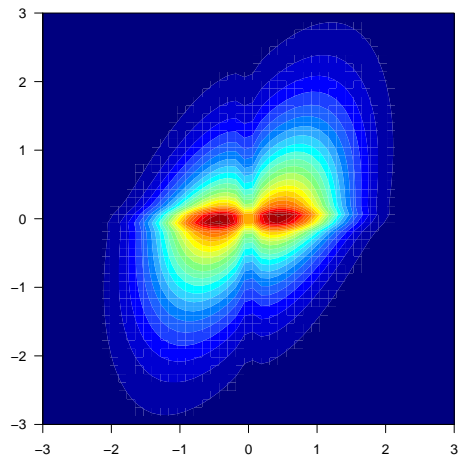
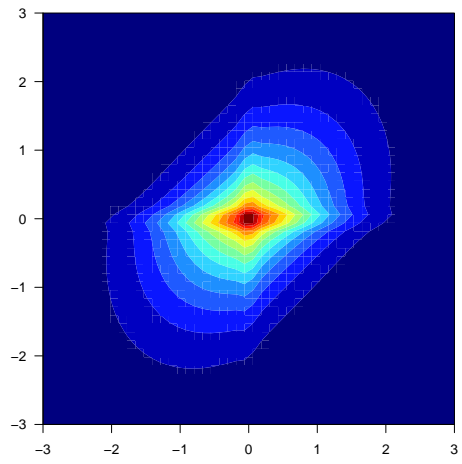
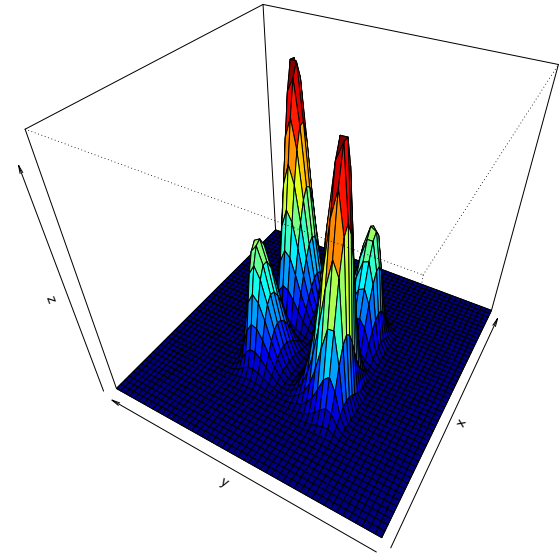
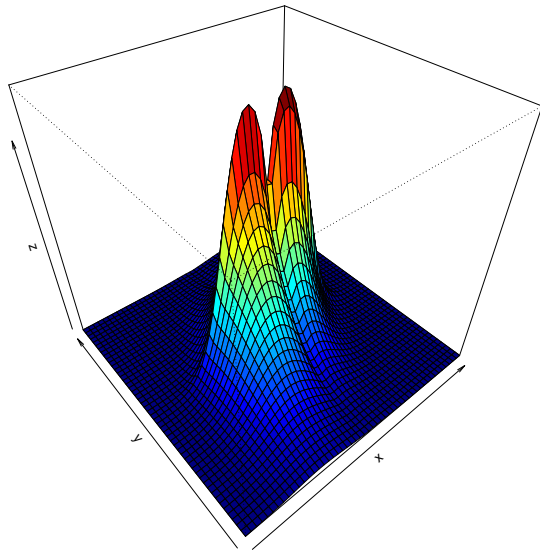
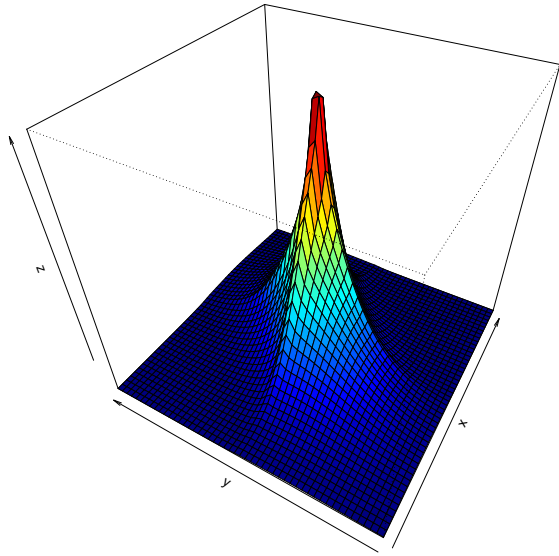


Extensions to non-Gaussian (Abnormal)

Liu, Lafferty and Wasserman (2009, JMLR)

Definition. A random vector $\mathbf{X} = (X_1, \dots, X_d)^T$ is said to follow a **nonparanormal** distribution if there exists a set of **unspecified** univariate increasing functions $\{f_j\}_{j=1}^d$ such that

$$(f_1(X_1), f_2(X_2), \dots, f_d(X_d))^T \sim N_d(\mathbf{0}, \mathbf{\Sigma}), \quad \text{where } \text{diag}(\mathbf{\Sigma}) = \mathbf{1}.$$



Liu et.al. Annals of Statistics (2012)

For each pair (X_j, X_k) , **Spearman's rho coefficient** of (X_j, X_k) , denoted by $\hat{\rho}_{jk}$, is the correlation of the ranks. Let

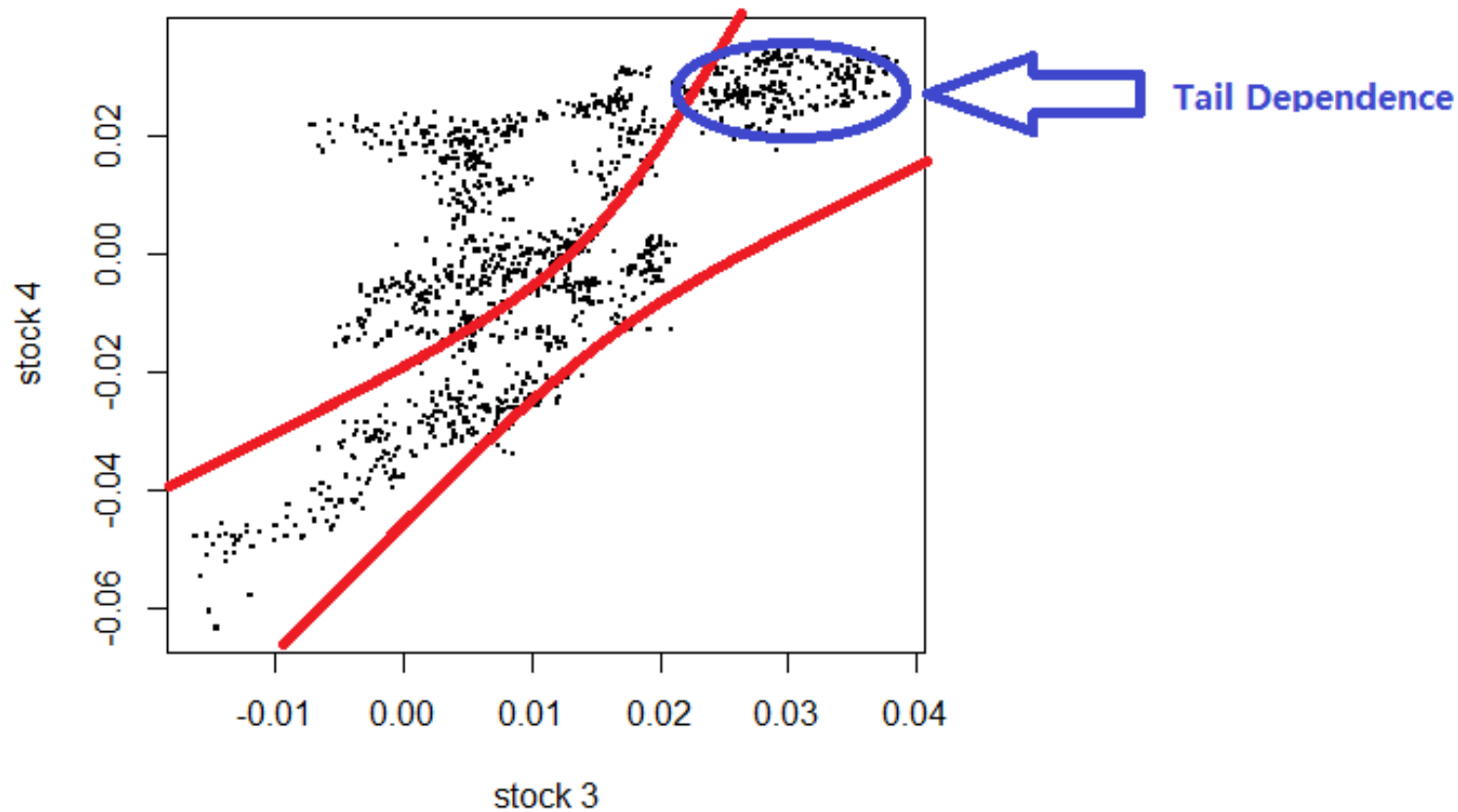
$$\hat{\mathbf{R}}^\rho := \left[2 \sin \left(\frac{\pi}{6} \hat{\rho}_{jk} \right) \right].$$

Theorem. Confined in the **nonparanormal**, $\|\hat{\mathbf{R}}^\rho - \mathbf{\Sigma}\|_{\max} = O_P \left(\sqrt{\frac{\log d}{n}} \right)$.

Parametric (optimal) rates in graph recovery and **near-parametric rate** in principal component analysis (Han and Liu, NIPS 2012).

Good Enough?

Theorem. The tail dependence is **zero** in Nonparanormal (Gaussian Copula).





Elliptical Distribution

When the density exists, the elliptical distribution has the **density**:

$$f(\mathbf{x}) = k \cdot g((\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})),$$

where g is an **unspecified** univariate positive function. In this case, we represent it by $EC_d(\boldsymbol{\mu}, \boldsymbol{\Sigma}, g)$.

Note: Elliptical distributions can be **very heavy tailed** and even have **infinite first or second moments**.

Transelliptical Distribution

Definition. A random vector $\mathbf{X} = (X_1, \dots, X_d)^T$ is said to follow a **transelliptical** distribution if there exists a set of **unspecified** univariate increasing functions $\{f_j\}_{j=1}^d$ such that

$$(f_1(X_1), f_2(X_2), \dots, f_d(X_d))^T \sim EC_d(\mathbf{0}, \mathbf{\Sigma}, g), \quad \text{where } \text{diag}(\mathbf{\Sigma}) = \mathbf{1}.$$

In this case, we represent $\mathbf{X} \sim TE_d(\mathbf{\Sigma}, g; f_1, \dots, f_d)$.

Kendall's tau and Its Invariance Property

Population Kendall's tau:

$$\tau(X_j, X_k) = \mathbb{P}((X_j - \tilde{X}_j)(X_k - \tilde{X}_k) > 0) - \mathbb{P}((X_j - \tilde{X}_j)(X_k - \tilde{X}_k) < 0),$$

where $(\tilde{X}_j, \tilde{X}_k)$ is a independent copy of (X_j, X_k) .

Theorem. Given $\mathbf{X} \sim TE_d(\boldsymbol{\Sigma}, g; f_1, \dots, f_d)$,

$$\boldsymbol{\Sigma}_{jk} = \sin\left(\frac{\pi}{2}\tau(X_j, X_k)\right)$$

.

Invariant to both g and f_j .

Transelliptical Component Analysis

Let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be n independent realizations of \mathbf{X} . Using the Kendall's tau correlation coefficient estimate:

$$\hat{\tau}_{jk} = \frac{2}{n(n-1)} \sum_{1 \leq i < i' \leq n} \text{sign}(x_{ij} - x_{i'j})(x_{ik} - x_{i'k}).$$

Let

$$\hat{\mathbf{R}}^\tau := \left[\sin \left(\frac{\pi}{2} \hat{\tau}_{jk} \right) \right].$$

Plugging $\hat{\mathbf{R}}^\tau$ into any sparse principal component algorithm.

Theoretical Results

In estimating θ_1 and its support, we have

- the near-parametric rate of convergence under the ℓ_2 norm:

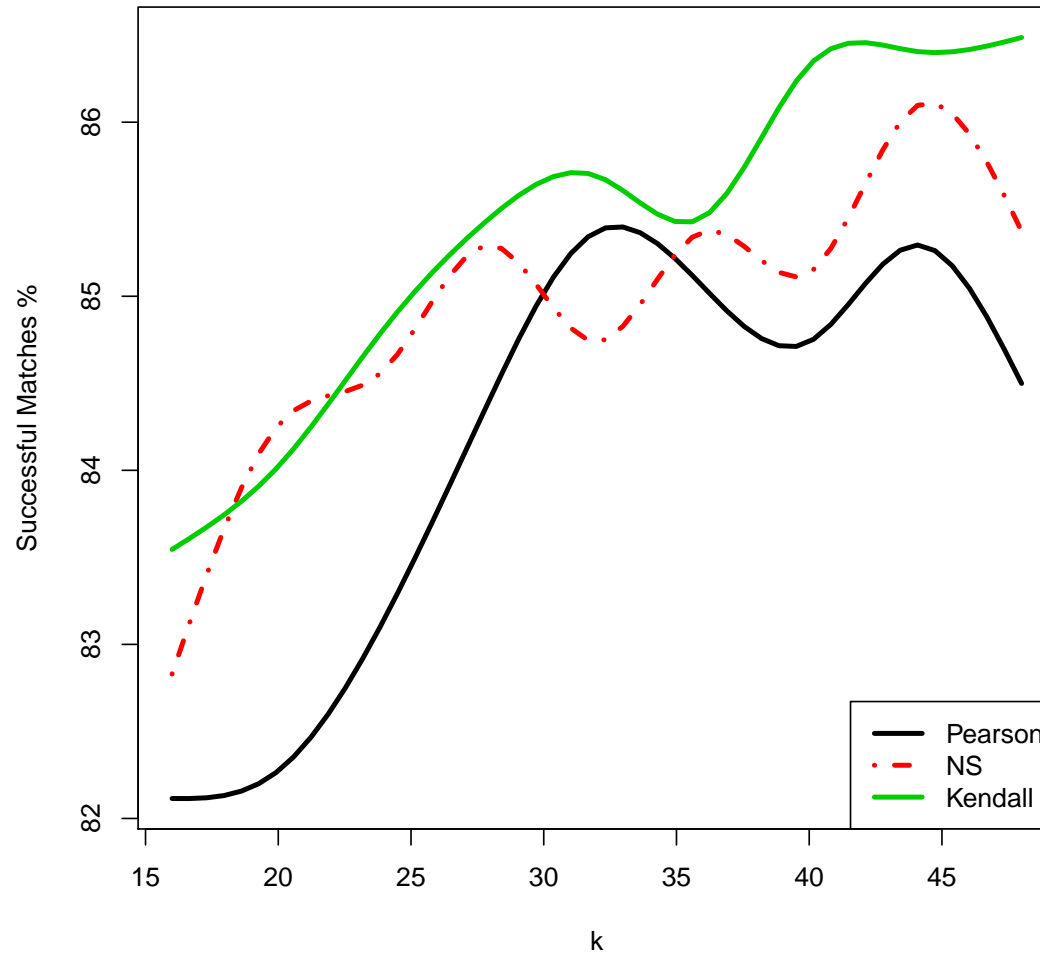
$$O_P \left(\frac{s}{\lambda_1 - \lambda_2} \sqrt{\frac{\log d}{n}} \right)$$

- a support recovery threshold at an order of

$$O_P \left(\frac{s}{\lambda_1 - \lambda_2} \sqrt{\frac{\log d}{n}} \right)$$

Equity Data Again

Prediction of the Market Trend:



Beyond PCA

- Directed Graphs Estimation
- Un-directed Graphs Estimation
- Discriminant Analysis
- Independent Component Analysis
- Canonical Component Analysis
- ...

Remarks

- Model Flexibility
- Parametric or near-parametric rate
- Procedure simple and computation efficiency

Thanks!

for more information, please come to the poster session (Tu66).