

*Regularized Estimation of High Dimensional  
Covariance Matrices*

Peter Bickel

Cambridge

January, 2008

With Thanks to

- E. Levina (Joint collaboration, slides)
- I. M. Johnstone (Slides)
- Choongsoon Bae (Slides)

## *Outline*

1. Why estimate covariance matrices?
2. Some examples
3. Pathologies of empirical covariance matrix
4. Sparsity in various guises
5. A brief review of methods
6. Asymptotic theory for banding, thresholding, SPICE (LASSO penalized Gaussian log likelihood)
7. Simulations and data
8. Some future directions

## Model

$\mathbf{X}_1, \dots, \mathbf{X}_n$ :  $p$  vectors, I.I.D.  $F$ .

- Base model:  $F = N_p(\boldsymbol{\mu}, \Sigma)$

- More generally:  $\mathbb{E}_F |\mathbf{X}|^2 < \infty$

$$\boldsymbol{\mu}(F) = \mathbb{E}_F \mathbf{X}$$

$$\Sigma(F) = \mathbb{E}_F (\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T$$

- Asymptotic theory:  $n, p \rightarrow \infty$ ,  $\Sigma \in \mathcal{T}_{p,n}$
- Extensions to Stationary  $\mathbf{X}_1, \dots, \mathbf{X}_n$ , Wei Biao Wu (2007).

## *Why Estimate Covariance Matrix?*

- Principal component analysis (PCA)
- Linear or quadratic discriminant analysis (LDA/QDA)
- Inferring independence and conditional independence in Gaussian models (graphical models)
- Inference about the mean (e.g. longitudinal mean response curve)

Covariance itself is usually not the end goal:

PCA: Estimation of the **eigenstructure**

LDA/QDA, Graphical models: Estimation of **inverses**

# *Covariance Matrices in Climate Research*

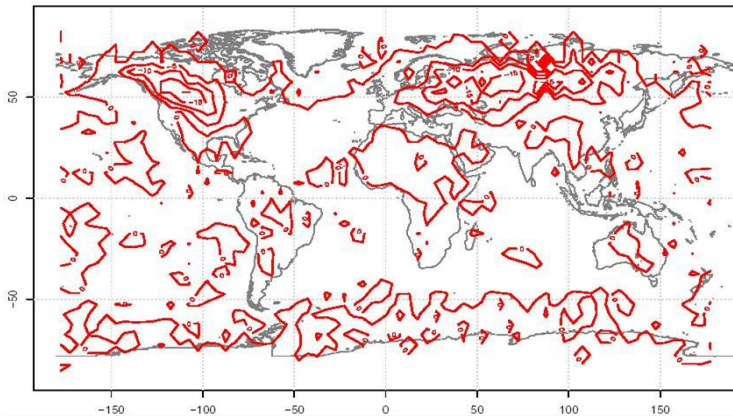
- Example 1

Data courtesy of Serge Guillas

- Empirical Orthogonal Functions(EOF) aka principal components.
- $\mathbf{X}_1, \dots, \mathbf{X}_n$  :  $p$  dimensional stationary vector time series.
- Monthly Jan Temp 1850 – 2006
- Latitude longitude calls  $5^\circ \times 5^\circ$  grid
- $p=2592$
- $n=157$

# Example 1

EOF pattern #1(field)



EOF PATTERN #1 (CLIM.PACT)

## Example 1

- $\mathbf{X}_k = \{X(i, j) : i = \text{latitude}, j = \text{longitude}\}$
- $\mathbb{X}_{n \times p} = (\mathbf{x}_1^T, \dots, \mathbf{x}_n^T)^T$
- $\mathbb{E}(\mathbf{X}) = \mathbf{0}$
- $\Sigma = \text{Var}(\mathbf{X}) = \mathbb{E}(\mathbf{X}\mathbf{X}^T) = \sum_{j=1}^p \lambda_j \mathbf{e}_j \mathbf{e}_j^T$
- $\mathbf{e}_1, \dots, \mathbf{e}_p$  : Principal components.
- $\lambda_1 > \dots > \lambda_p$  : Eigenvalues.
- Goal : Estimate, interpret  $\mathbf{e}_j, j = 1, \dots, K$   
such that  $\frac{\sum_{j=1}^K \lambda_j}{\sum_{j=1}^p \lambda_j}$  large.

# Covariance Matrices in Climate Research

- Example 2

## Data assimilation

- $X_j$  = ave. “pressure”, “temperature”, ... in  $50\text{km} \times 50\text{km} \times$  variable block of atmosphere,  $|J| \asymp 10^7$ , computer model.
- $\mathbf{X}_i = \mathbf{X}(t_i)$   $i = 1, \dots, T$ .
- $\mathbf{Y}(t_i)$  : Data vectors.
- Ensemble:  $\mathbf{X}_j^U(t)$ ,  $1 \leq j \leq n$ .
- Data assimilated :  $\mathbf{X}_j^F(t)$ ,  $1 \leq j \leq n$  uses  $\hat{\Sigma}^{-1}$  as estimate of true  $[\text{Var}(\mathbf{X}_1^U)]^{-1}$  for Kalman gain.



## *Pathologies of Empirical Covariance Matrix*

Observe  $\mathbf{X}_1, \dots, \mathbf{X}_n$ , i.i.d.  $p$ -variate random variables

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T$$

- MLE, for Gaussian unbiased (almost), well-behaved (and well studied) for fixed  $p$ ,  $n \rightarrow \infty$ . But very **noisy** if  $p$  is large.
- **Singular** if  $p > n$ , so  $\hat{\Sigma}^{-1}$  is not uniquely defined.
- Computational issues with  $\hat{\Sigma}^{-1}$  for large  $p$ .
- LDA completely breaks down if  $p/n \rightarrow \infty$
- Eigenstructure can be inconsistent if  $p/n \rightarrow \gamma > 0$ .

# Eigenvalues

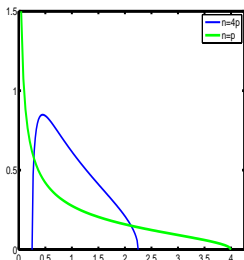
Description of spreading phenomenon:

Empirical distribution function: for eigenvalues  $\{\hat{\ell}_i\}_{j=1}^p$

$$G_p(t) = p^{-1} \#\{\hat{\ell}_j \leq t\} \rightarrow G(t) \leftrightarrow g(t)dt.$$

Marčenko-Pastur, (67), For  $F = N_p(\mathbf{0}, I)$ ,  $p/n \rightarrow \gamma$

$$g^{MP}(t) = \frac{\sqrt{(b_+ - t)(t - b_-)}}{2\pi\gamma t},$$
$$b_{\pm} = (1 \pm \sqrt{\gamma})^2.$$



## Eigenvectors

D.Paul (2006) For  $N_p(\mathbf{0}, \Sigma)$

$$\Sigma = \text{diag}(\lambda, 1, \dots, 1) = \begin{pmatrix} \lambda & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ 0 & & \ddots & \\ 0 & \dots & & 1 \end{pmatrix}, \quad 1 < \lambda < 1 + \sqrt{\gamma},$$

Eigenvector  $\hat{\mathbf{e}} \leftrightarrow \hat{\lambda}$ ,  $\mathbf{e} \leftrightarrow \lambda$ ,

$$|\mathbf{e} - \hat{\mathbf{e}}|^2 \xrightarrow{\text{a.s.}} 2$$

## Fisher's LDA

- For classification into two populations using:  
 $\mathbf{X}_1, \dots, \mathbf{X}_n \sim N_p(-\Delta, \Sigma), \mathbf{X}_{n+1}, \dots, \mathbf{X}_{2n} \sim N_p(\Delta, \Sigma)$
- Minimax behaviour over  $\mathcal{T}_{p,n} = \{\Sigma : \Delta^T \Sigma^{-1} \Delta \geq c > 0\}$  of LDA based on MLE  $\hat{\Sigma}^{-1}$  (Moore-Penrose for  $p \geq 2^{n-2}$ ) is equivalent to coin tossing if  $\frac{p}{n} \rightarrow \infty$  (Bickel - Levina (2004))
- True for any rule!. What if  $\mathcal{T}$  consists of “sparse” matrices?

## Eigenstructure Sparsity

- Write,

$$\Sigma = \Sigma_0 + \sigma^2 I_p, \quad \text{where } \Sigma_0 = \sum_{j=1}^M \lambda_j \boldsymbol{\theta}_j \boldsymbol{\theta}_j^t$$

and  $\lambda_1 \geq \dots \geq \lambda_M > 0$ ,  $\{\boldsymbol{\theta}_j\}$  orthonormal.

- **Equivalent :**

$$\mathbf{x}_i = \mu + \sum_{j=1}^M \sqrt{\lambda_j} v_{ji} \boldsymbol{\theta}_j + \sigma \mathbf{z}_i, \quad i = 1, \dots, n, \quad v_{ji} \sim i.i.d. \mathcal{N}(0, 1)$$

## *Eigenstructure Sparsity defined*

- $p, n$  large but,
  - (i)  $M$  small fixed.
  - (ii)  $\theta_j$  “sparse”
- “well approximated” by  $\theta_{js}$ ,  $\|\theta_{js}\|_0 \leq s$ ,  $s$  “small”.

$\|\mathbf{v}_{js}\|_0 \equiv \#$  of nonzero coordinates of  $\mathbf{v}$ .

## *Label Dependent Sparsity of $\Sigma$*

- Write  $\Sigma = \|\sigma_{ij}\|$ ,
  - $|\sigma_{ij}|$  small (effectively 0) for  $|i - j|$  large
  - More generally, given metric  $m$  on  $J$ ,  $|\sigma_{ij}|$  small if  $m(i, j)$  large
- Note:  $\sigma_{ij} = 0 \Leftrightarrow X_i \perp X_j$  under Gaussianity.

## *Label Dependent Sparsity of $\Sigma^{-1}$*

If  $\Sigma^{-1} = \|\sigma^{ij}\|$ ,

- $|\sigma^{ij}|$  small for "many"  $(i, j)$  pairs if  $|i - j|$  is large.

Note:

$\sigma^{ij} = 0 \Rightarrow X_i \perp X_j \mid \{X_k : k \neq i, j\}$  for Gaussian case.



## Examples of Label Dependent Sparsity

(i)  $\mathbf{X}$  stationary  $\sigma(i, j) = \sigma(|i - j|)$

$\leftrightarrow$  spectral density  $f$ ,  $0 < \varepsilon \leq f \leq \frac{1}{\varepsilon} < \infty$

- Ergodic ARMA processes satisfy

(ii)  $T = S + K \leftrightarrow \mathbf{X} = \mathbf{Y} + \mathbf{Z}$

$\mathbf{Y}, \mathbf{Z}$  independent,  $S \leftrightarrow \mathbf{Y}$ ,  $K \leftrightarrow \mathbf{Z}$ ,  $S = [s(i - j)]$  as in (i),

$K$  Hilbert-Schmidt:

- $\sum_{i,j} K^2(i, j) < \infty$  ( $Z_m \xrightarrow{p} 0$ , non stationary)
- (i)  $\Rightarrow \sum_i s^2(i) < \infty$
- $S = 0 \Rightarrow \Sigma$  singular but correspond to functional data analysis

A typical goal: Estimate  $\Sigma^{-1}$  for prediction.

## Permutation Invariant Sparsity of $\Sigma$ or $\Sigma^{-1}$

a) Each row of  $\Sigma$  sparse or sparsely approximable

e.g. If  $\sigma_i = \{\sigma_{i,j} : 1 \leq j \leq p\}$ ,  $\|\sigma_i\|_0 \leq s$  for all  $i$ .

b) Each row of  $\Sigma^{-1}$  sparsely approximable.

- a) roughly implies b) if  $\lambda_{\min}(\Sigma) \geq \delta > 0$ .
- A more general notion (El Karoui (2007) AS to appear).

*(Roughly)* The number of loops of length  $k$  in the undirected graph:  
 $i \sim j$  iff  $\sigma_{ij} \neq 0$  grows more slowly than  $p^{k/2}$  as  $k \rightarrow \infty$ . The  
corresponding definition for  $\Sigma^{-1}$  is **not** equivalent in general.

## Permutation Invariant Sparsity of $\Sigma^{-1}$

- $\Sigma^{-1}$  corresponds to a graphical model.
- $\mathcal{N}(i) = \{j : \sigma^{ij} \neq 0, j \neq i\}$
- Sparsity means neighbourhood size  $\ll p$ .

Example: Gene networks

- Goal : Determine  $\mathcal{N}(i)$   $i = 1, \dots, p$ .

*Construction of Estimates of  $\Sigma$  or Parameters of  $\Sigma$   
which work well for  $\mathcal{T}$  having sparse members*

- **Eigenstructure sparsity**

*a)* Sparse PCA d'Aspremont *et al.*, SIAM Review (2007), Hastie, Tibshirani, Zhou (2007).

- PCA with Lasso penalty on coordinates of eigenvectors
- Focus on computation, no asymptotics

*Construction of Estimates of  $\Sigma$  or Parameters of  $\Sigma$   
which work well for  $\mathcal{T}$  having sparse members*

- **Eigenstructure sparsity continued**

*b)* Johnstone and Lu, JASA (2006).

- Sparse PCA for  $\Sigma$  corresponding to signals  
 $\mathbf{X} = (X(t_1), \dots, X(t_p))^T$
- Represent in wavelet basis
- Use only  $k \ll p$  coefficients to represent data to whose covariance matrix PCA is applied
- Asymptotics for large  $p$  and examples

*Construction of Estimates of  $\Sigma$  or Parameters of  $\Sigma$   
which work well for  $\mathcal{T}$  having sparse members*

- **Label Dependent Sparsity**

- Banding

Bickel and Levina (2004), Bernoulli

Bickel and Levina (2007), Annals of Statistics (To appear)

a) "Banding"  $\hat{\Sigma}$

For any  $M \equiv \|m_{ij}\|$ ,

$$B_k(M) \equiv \|m_{ij} \mathbf{1}(|i-j| \leq k)\|$$
$$\hat{\Sigma} \equiv \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

Estimate:  $B_k(\hat{\Sigma})$ ,  $k = o(p)$ .

*Construction of Estimates of  $\Sigma$  or Parameters of  $\Sigma$   
which work well for  $\mathcal{T}$  having sparse members*

- **Label Dependent Sparsity continued**

b) “Banding”  $\Sigma^{-1}$

$$\Sigma^{-1} = AA^T, \text{ } A \text{ Lower triangular}$$

Estimate:  $\hat{\Sigma}^{-1} \equiv \hat{A}_k \hat{A}_k^T$ , where  $\hat{A}_k$  is obtained by regression.

Pourahmadi and Wu (2003), Biometrika.

Huang, Liu, Pourahmadi (2007)

*Construction of Estimates of  $\Sigma$  or Parameters of  $\Sigma$   
which work well for  $\mathcal{T}$  having sparse members*

- **Permutation Invariant Sparsity**

(a) Thresholding

For any  $M = \|m_{ij}\|$ ,  $T_t(M) \equiv \|m_{ij}\mathbf{1}(|m_{ij}| \geq t)\|$

Estimate:  $T_t(\widehat{\Sigma})$ ,  $|t| = o(1)$

El Karoui (2007)

Bickel and Levina (2007).



# Construction of Estimates of $\Sigma$ or Parameters of $\Sigma$ which work well for $\mathcal{T}$ having sparse members

- **Permutation Invariant Sparsity continued**

- b) SPICE

Estimate of  $\Sigma^{-1}$ :  $\Sigma^{-1}$  sparse

$$F(P, M) \equiv \text{Trace}(MP) - \log \det(P) + \lambda |P|_1,$$

where  $|M|_1 = \sum_{i,j} |m_{ij}|$

- a)  $\hat{P}^{(1)} \equiv \arg \min \{ F(P, \hat{\Sigma}) : P \succeq 0 \}$

- Yuan, Li (2007), *Biometrika*. Fixed  $p$  asymptotics

- Banerjee *et al.* (2006), ICML computation

- b)  $\hat{R}^{(2)} \equiv \arg \min \{ F(P, \hat{R}) : R \succeq 0, r_{ii} = 1, \text{ all } i \}$

- $\hat{P}^{(2)} \equiv \hat{D} \hat{R} \hat{D}, \hat{D} \equiv \text{diag}(\hat{\sigma}_{ii}).$

- $\hat{R} \equiv \hat{D}^{-1} \hat{\Sigma} \hat{D}^{-1}$

- a),b): Rothman *et al.* (2007), Submitted to JRSS(B) Asymptotics

*Construction of Estimates of  $\Sigma$  or Parameters of  $\Sigma$   
which work well for  $\mathcal{T}$  having sparse members*

- **Permutation Invariant Sparsity continued**

- c)* Methods for neighbourhood estimation

- Meinshausen and Bühlmann (2006), Zhao and Yu (2007):  
Regression methods + Lasso
- Wainwright (2006): Bounds
- Kalisch and Bühlmann (2007): Using the PC algorithm of Spertes  
*et al.* (2000) to roughly hierarchically test whether partial  
correlation coefficients are 0

## Some Asymptotic Theory

- Matrix norms

$$M \equiv \|m_{ij}\|_{p \times p}$$
$$|\mathbf{x}|_r^r \equiv \sum_{j=1}^p |x_j|^r, \quad \mathbf{x} = (x_1, \dots, x_p)$$

- Operator norms

$$\|M\|_{(2,2)} = \lambda_{\max}^{1/2}(MM^T)$$
$$\|M\|_{(1,1)} = \max_j \sum_{i=1}^p |m_{ij}|$$
$$\|M\|_{(\infty, \infty)} = \max_i \sum_{j=1}^p |m_{ij}|$$
$$|M|_{\infty} \equiv \max_{i,j} |m_{ij}|$$
$$\|M\|_F^2 \equiv \sum_{i,j} m_{ij}^2 : \text{Frobenius norm}$$

## Properties

- For any operator norm,

$$\|AB\| \leq \|A\|\|B\|$$

Henceforth,  $\|M\|_{(2,2)} \equiv \|M\|$ .

- If  $M_{p \times p}$  is symmetric,

$$\|M\| = \text{Max} \left\{ \left| \lambda_{\text{Max}}^{(M)} \right|, \left| \lambda_{\text{Min}}^{(M)} \right| \right\}.$$

- $\|M\| \leq [\|M\|_{(1,1)}\|M\|_{(\infty,\infty)}]^{1/2}$ .
- $\|M\| \leq \|M\|_F$ .

## Basic Property of $\|\cdot\|$

Given  $A_n, B_n$  symmetric  $\|A_n - B_n\| \rightarrow 0$ ,

suppose  $\lambda_1(B_n) > \lambda_2(B_n) > \cdots > \lambda_k(B_n) > \lambda_{k+1}(B_n)$

and define  $\lambda_j(A_n)$  analogously.

Suppose  $\lambda_{j+1}(B_n) < \lambda_j(B_n) - \Delta$ ,  $1 \leq j \leq k$

Dimension  $B_n$  arbitrary,  $k, \Delta > 0$  fixed.

Then,

a)  $|\lambda_j(A_n) - \lambda_j(B_n)| = O(\Delta^{-j})\|A_n - B_n\|$

b) If  $E_{jA}$  respectively  $E_{jB}$  is projection operator onto eigenspace corresponding to  $\lambda_j$ , then

$$\|E_{jA} - E_{jB}\| = O(\Delta^{-j}\|A_n - B_n\|)$$

## B-L (2006) Main Result I

Banded estimator :

$$\hat{\Sigma}_{k,p}(i,j) = \hat{\Sigma}_p(i,j) \cdot \mathbf{1}(|i-j| \leq k)$$

Let

$$\begin{aligned} \mathcal{U}(\epsilon_0, \alpha, C) = & \{ \Sigma : 0 < \epsilon_0 \leq \lambda_{\min}(\Sigma) \leq \lambda_{\max}(\Sigma) \leq 1/\epsilon_0, \\ & \max_j \sum_i \{ |\sigma_{ij}| : |i-j| > k \} \leq Ck^{-\alpha} \text{ for all } k \geq 0 \}. \end{aligned}$$

### Theorem 1

If  $\mathbf{X}$  is Gaussian and  $k_n \asymp (n^{-1} \log p)^{-\frac{1}{2(\alpha+1)}}$ , then, uniformly on  $\Sigma \in \mathcal{U}(\epsilon_0, \alpha, C)$ ,

$$\|\hat{\Sigma}_{k_n,p} - \Sigma_p\| = O_P\left((n^{-1} \log p)^{\frac{\alpha}{2(\alpha+1)}}\right) = \|\hat{\Sigma}_{k_n,p}^{-1} - \Sigma_p^{-1}\|$$

The banded estimator and its inverse are consistent if  $\frac{\log p}{n} \rightarrow 0$

### Remark

$\lambda_{\min} \geq \epsilon_0$  not needed if only convergence in  $\|\cdot\|$  to  $\Sigma$  is needed.

# *An Approximation Theorem (Bickel and Lindner (2007))*

If  $m\|B_k - A\| \leq \delta$ , ( $B_k$  banded of width  $2k$ ),  $\|A\| = M < \infty$ ,  
 $\|A^{-1}\| = m^{-1} < \infty$ ,  $K \equiv \frac{M}{m}$ ,

There exists  $N(m, M, \delta)$ ,  $B_{kN}$  such that

$$\|B_{kN} - A^{-1}\| \leq \frac{1}{M} c(m, M, \delta) \delta,$$

where  $c(m, M, \delta)$  tends to  $c(K)$  as  $\delta \rightarrow 0$ .

$B_{kN}$  has a simply computable form.

## Thresholding

**Theorem** (*Thresholding*) (B-L Submitted to AS)

$$\text{If } \mathcal{U}_t(q, c_0(p), M) = \left\{ \Sigma : \sigma_{ii} \leq M_0, \sum_{j=1}^p |\sigma_{ij}|^q \leq c_0(p), \text{ all } i \right\},$$
$$0 \leq q < 1, t_n = M \sqrt{\frac{\log p}{n}},$$

$$\left\| T_{t_n}(\hat{\Sigma}) - \Sigma \right\| = O_p \left( c_0(p) \left( \frac{\log p}{n} \right)^{(1-q)/2} \right)$$

*Note*

$q = 0 \leftrightarrow c_0(p)$  : non zero entries per row

**More subtle consistency results** : N. El Karoui (2007).



## *SPICE (Rothman et al.)*

If  $S \equiv \{(i, j) : \sigma^{ij} \neq 0\}$ ,  $|S| \leq s$ ,  $\|\Sigma^{-1}\| \leq m^{-1}$ ,  $\|\Sigma\| \leq M$  and  $\mathbf{X}$  is Gaussian,

$$\begin{aligned}\left\|\hat{\rho}^{(1)} - \Sigma\right\|_F^2 &= O_p\left(\left(p + s\right)\frac{\log p}{n}\right) \\ \left\|\hat{\rho}^{(2)} - \Sigma\right\|^2 &= O_p\left(\frac{s \log p}{n}\right)\end{aligned}$$

# Choosing the “*Banding*” (or *Thresholding* or *SPICE*) Parameter

Ideally want to minimize risk

$$R(k) = E\|\hat{\Sigma}_k - \Sigma\|$$

Estimate via a **resampling scheme**:

- Split the data into two samples of size  $n_1, n_2, N$  times at random
- Let  $\hat{\Sigma}_1^{(\nu)}, \hat{\Sigma}_2^{(\nu)}$  be the two sample covariance matrices from the  $\nu$ -th split. The risk can be estimated by

$$\hat{R}(k) = \frac{1}{N} \sum_{\nu=1}^N \|(\hat{\Sigma}_1^{(\nu)})_k - \hat{\Sigma}_2^{(\nu)}\|$$

- We used  $n_1 = n/3, N = 50$ , and the  $L_1$  matrix norm instead of  $L_2$ .
- Oracle properties for Frobenius norm (Bickel and Levina (2007), Annals of Statistics submitted).

## Simulation Examples

- Covariance of AR(1):  $\Sigma \in \mathcal{U}$

$$\sigma_{ij} = \rho^{|i-j|}$$

$n = 100, p = 10, 100, 200, \rho = 0.1, 0.5, 0.9.$

- Fractional Gaussian noise (FGN): long-range dependence, not in  $\mathcal{U}$

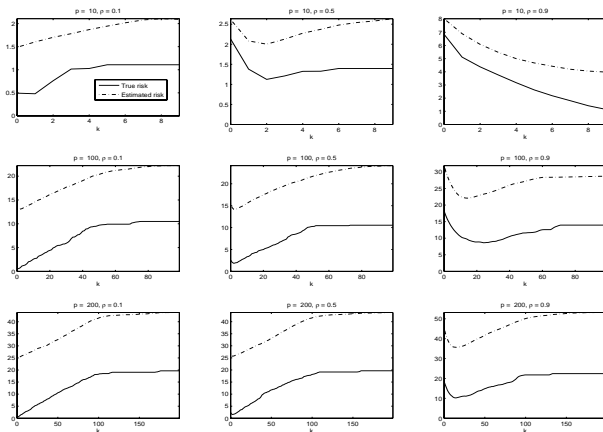
$$\sigma_{ij} = \frac{1}{2} \left[ (|i-j| + 1)^{2H} - 2|i-j|^{2H} + (|i-j| - 1)^{2H} \right]$$

$H \in [0.5, 1]$  is the Hurst parameter

$H = 0.5$  is white noise;  $H = 1$  is perfect dependence

$n = 100, p = 10, 100, 200, H = 0.5, 0.6, 0.7, 0.8, 0.9.$

# True and Estimated Risk for AR(1)



## (1, 1) Loss and Choice of $k$ for AR(1)

p	$\rho$	Mean(SD)		Loss		
		$k_1$	$\hat{k}$	$\hat{\Sigma}_{\hat{k}}$	$\hat{\Sigma}_{k_1}$	$\hat{\Sigma}$
10	0.1	0.5(0.5)	0.0(0.2)	0.5	0.4	1.1
10	0.5	3.3(0.8)	2.0(0.6)	1.1	1.0	1.3
10	0.9	8.6(0.7)	8.9(0.3)	1.5	1.5	1.5
100	0.1	0.2(0.4)	0.1(0.3)	0.6	0.6	10.2
100	0.5	2.7(0.7)	2.3(0.5)	1.6	1.5	10.6
100	0.9	21.3(4.5)	15.9(2.6)	9.2	8.5	13.5
200	0.1	0.2(0.4)	0.2(0.4)	0.7	0.6	20.4
200	0.5	2.4(0.7)	2.7(0.5)	1.8	1.7	20.8
200	0.9	20.2(4.5)	16.6(2.4)	9.9	9.5	24.5

*Table:* AR(1): Oracle and estimated  $k$  and the corresponding loss values.

$$k_1 = \operatorname{argmin}_k \|\hat{\Sigma}_k - \Sigma\|_{(1,1)}.$$

# Operator Norm Results for Various Estimates and Situations

$p$	$\rho$	Sample	Band	Band-P	Thresh	SPICE
10	0.1	0.60(0.01)	0.35(0.01)	0.35(0.01)	0.40(0.01)	0.35(0.01)
10	0.5	0.73(0.02)	0.64(0.02)	0.74(0.02)	0.80(0.02)	0.72(0.02)
10	0.9	1.02(0.05)	1.02(0.05)	1.02(0.05)	1.02(0.05)	1.01(0.05)
100	0.1	2.86(0.02)	0.45(0.01)	0.45(0.01)	0.50(0.01)	0.45(0.01)
100	0.5	3.39(0.03)	0.95(0.01)	2.04(0.00)	1.96(0.01)	1.52(0.01)
100	0.9	6.03(0.12)	4.45(0.09)	6.03(0.12)	6.18(0.12)	5.57(0.12)
200	0.1	4.63(0.02)	0.49(0.01)	0.49(0.01)	0.51(0.01)	0.50(0.01)
200	0.5	5.47(0.03)	1.00(0.01)	2.06(0.00)	2.04(0.00)	1.60(0.01)
200	0.9	9.77(0.16)	5.24(0.09)	9.76(0.16)	8.80(0.10)	7.50(0.10)

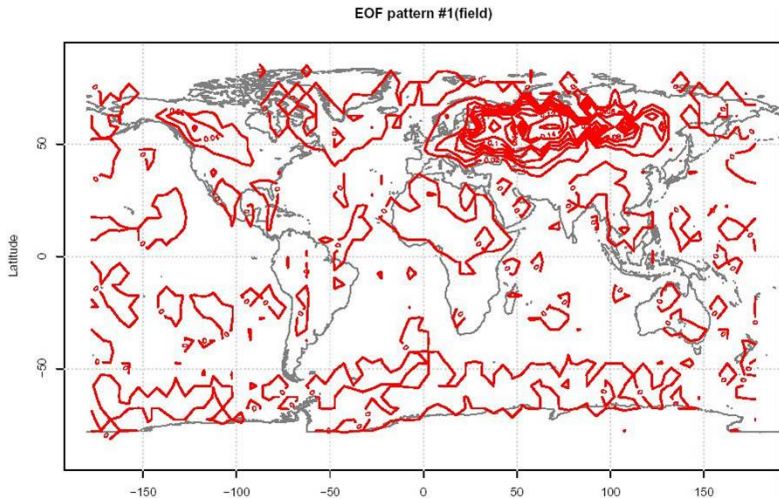
Table: Operator Norm Loss, Averages and SEs over 100 replications. AR(1)

# *Operator Norm Results for Various Estimates and Situations*

$p$	H	Sample	Band	Band-P	Thresh	SPICE
10	0.5	0.60(0.01)	0.27(0.01)	0.27(0.01)	0.35(0.01)	0.28(0.01)
10	0.7	0.67(0.02)	0.88(0.02)	1.03(0.03)	0.88(0.04)	1.05(0.05)
10	0.9	0.96(0.04)	0.96(0.04)	0.96(0.04)	0.96(0.04)	1.13(0.05)
100	0.5	2.83(0.01)	0.38(0.01)	0.38(0.01)	0.46(0.01)	0.39(0.01)
100	0.7	3.43(0.05)	4.46(0.02)	5.37(0.00)	5.36(0.00)	5.30(0.02)
100	0.9	8.02(0.28)	8.03(0.28)	8.02(0.28)	8.02(0.28)	14.18(0.70)
200	0.5	4.60(0.02)	0.44(0.01)	0.44(0.01)	0.46(0.01)	0.45(0.01)
200	0.7	5.81(0.06)	6.46(0.01)	7.40(0.00)	7.40(0.00)	7.39(0.00)
200	0.9	14.34(0.43)	14.65(0.46)	14.34(0.43)	14.34(0.43)	30.16(0.60)

*Table:* Operator Norm Loss, Averages and SEs over 100 replications. FGN

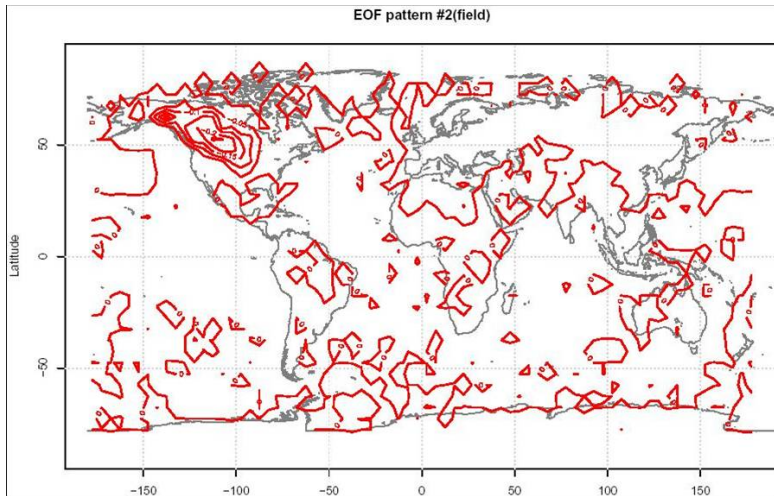
## Example 1 continued



EOF PATTERN #1 (THRESHOLDING)



## Example 1 continued



## *Some Important Directions*

- Choice of  $k$  or other regularization parameters.
- Canonical correlation regularized estimation.
- Independent component analysis regularization.
- Estimation of parameters governing independence and conditional independence in graphical models.