

Sparse Principal Components Analysis

Iain Johnstone, Statistics, Stanford & UC Berkeley

Arthur Yu Lu, Stanford and Renaissance Technologies

FDA Workshop, Gainesville, Jan 10-11, 2003

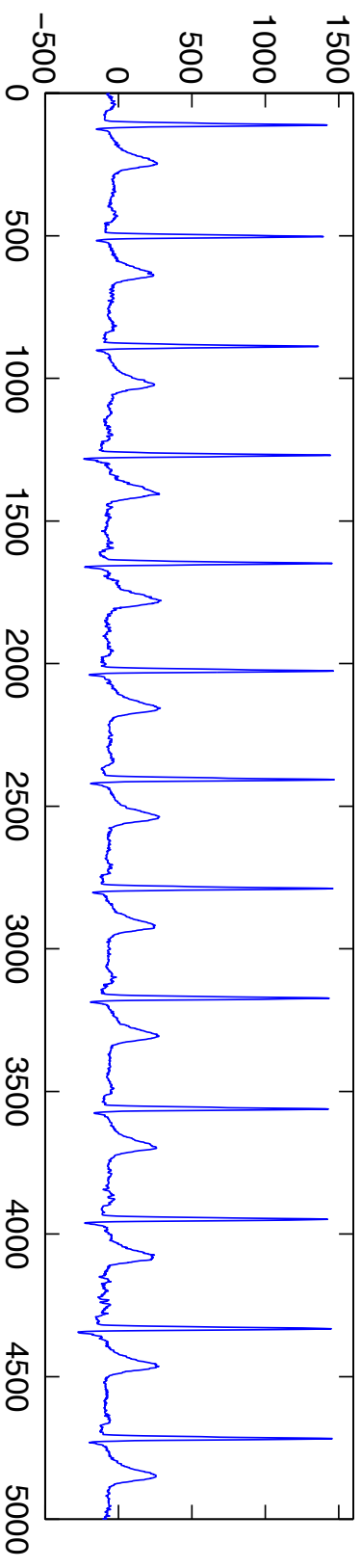
Support thanks: NIH

“Multiscale” Functional Data

Functional data: $x_i = x_i(t_j)$ $j = 1, \dots, p$, time points,
 $i = 1, \dots, n$, cases.

- Focus on PCA: \rightarrow principal modes of variation
- ‘high-dimensional’: $p = O(n)$ or larger
- Signals x_i contain localized features, perhaps on different scales—visible in p.c.’s also?
- Example: ECG signals

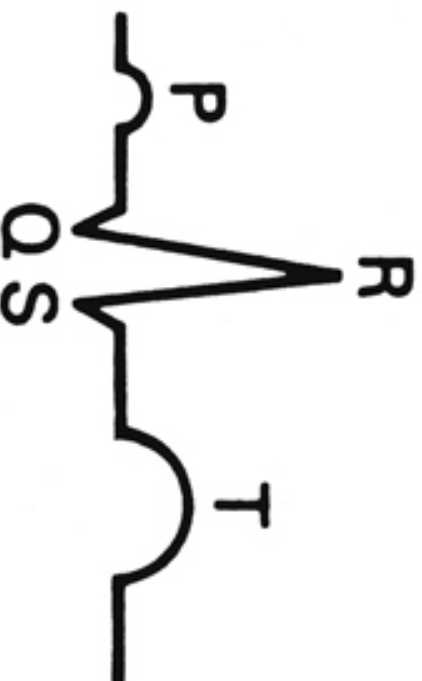
ElectroCardioGram traces



Average beat

vs.

beat to beat variation



Outline

⇒ Multiscale Functional Data $p \asymp n$

- Need for Dimension Reduction: Inconsistency
- Sparse PCA
 - Some sampling properties, incl. consistency
- Examples, incl. ECG

Main Themes

- initial dimension reduction before PCA
 - otherwise, inconsistency!
- use basis with sparse representation
 - so that little is lost in initial dim reduction

Background Theme: role of “Random Matrices”

- Small perturbations of symmetric matrices
- a.s. bounds for extreme eigenvalues of large matrices
- other talks here: large stochastic data matrices: (e.g. Marron)
 - ⇒ more generally, broader role for RMT tools in analysis of FDA methods?

Outline

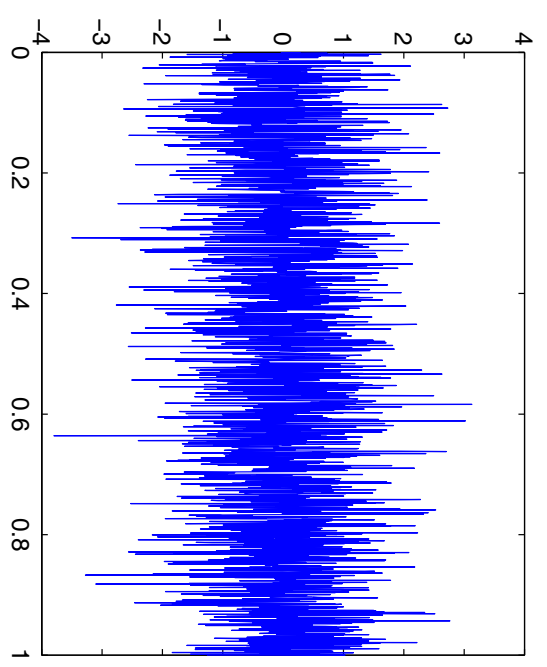
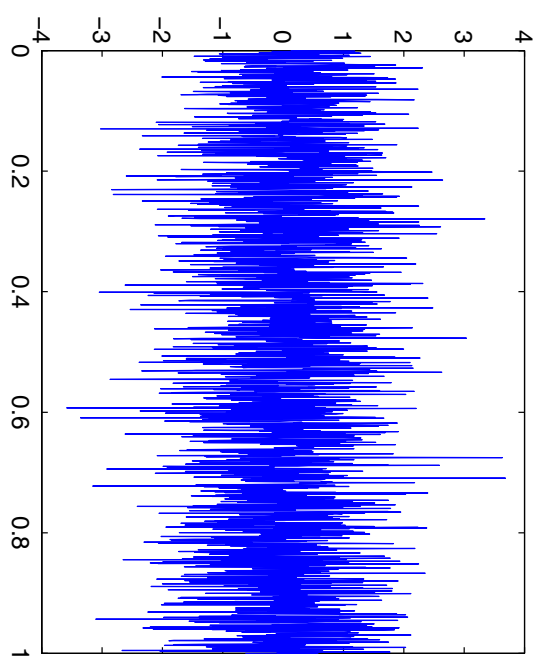
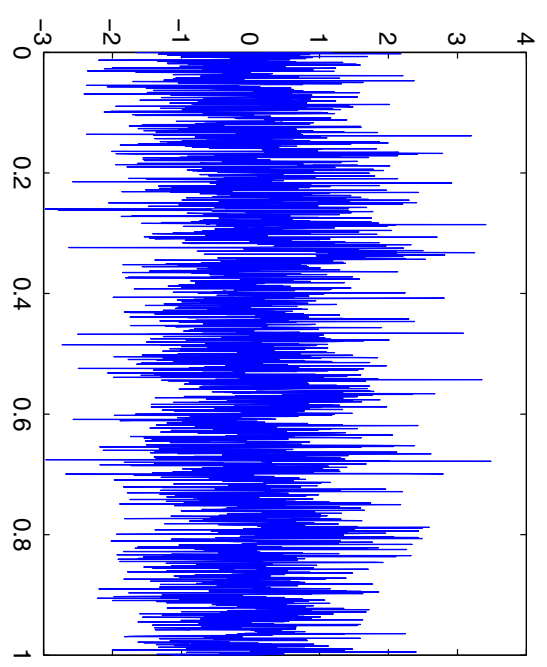
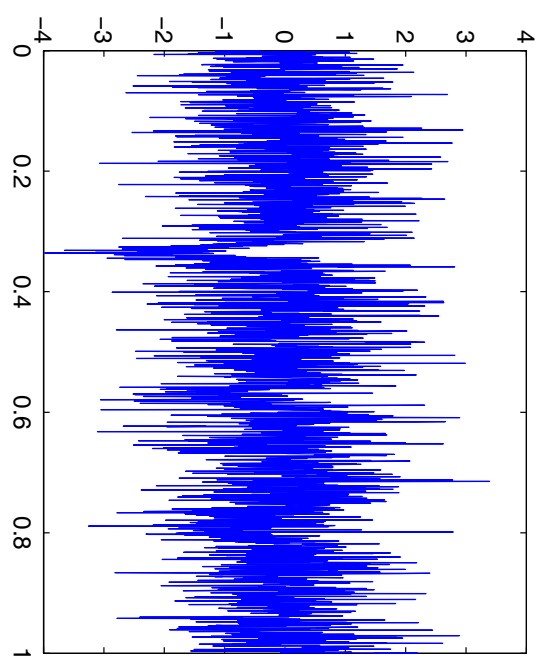
- Multiscale Functional Data $p \asymp n$
- ⇒ Need for Dimension Reduction: Inconsistency
 - Sparse PCA
 - Some sampling properties, incl. consistency
 - Examples, incl. ECG

Single Component Model

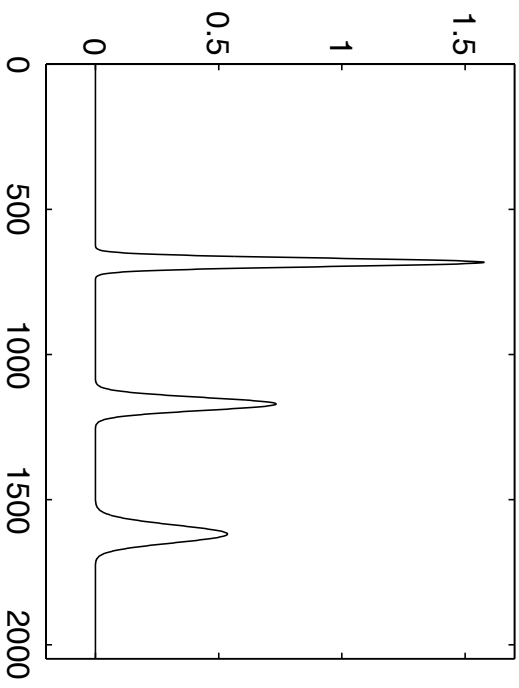
[to illustrate, e.g., need for dimension reduction]

$$x_i = v_i \rho + \sigma z_i \quad i = 1, \dots, n$$

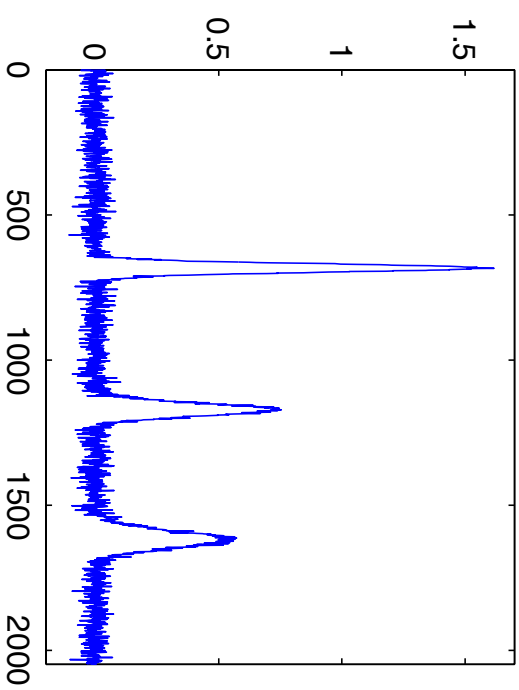
- $\rho \in \mathbb{R}^p$, single component to be estimated
(e.g. $\|\rho\| = 10, p = 2048$)
- $v_i \stackrel{i.i.d.}{\sim} N(0, 1)$ random effects (e.g. $n = 1024$)
- $z_i \stackrel{i.i.d.}{\sim} N_p(0, I)$ Gaussian noise (e.g. $\sigma = 1$)



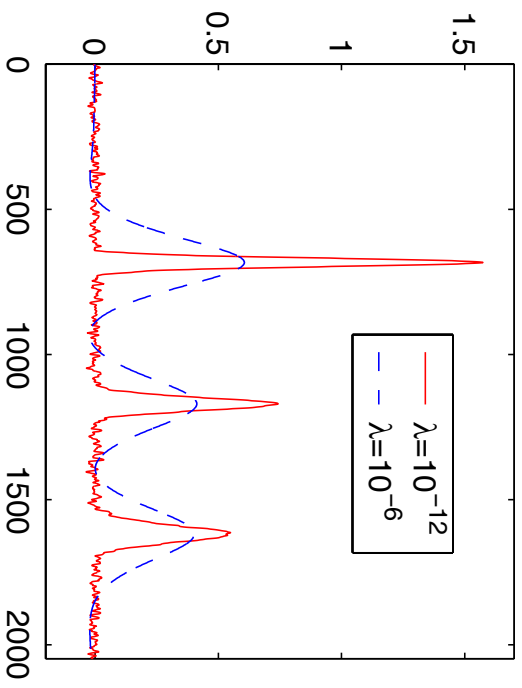
True PC, $p = 2048$, $n = 1024$, $\|p\| = 10$.



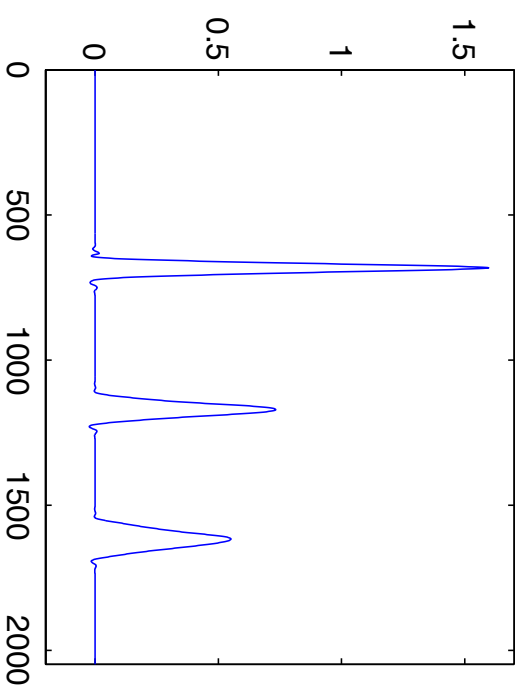
standard PCA



smoothed PCA



ASPCA, $w = 99.5\%$, $K = 372$.



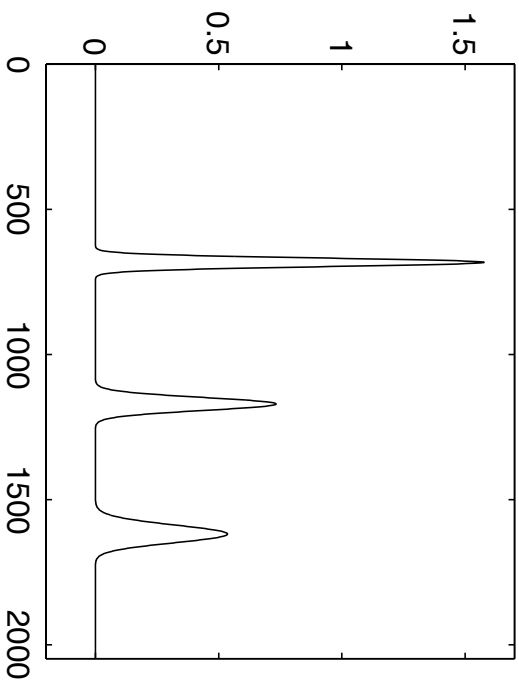
Smooth functional PCA

- Rice & Silverman (1991), Silverman (1996), Ramsay & Silverman (1997)
- Seek smoothed p.c.'s ξ by maximizing

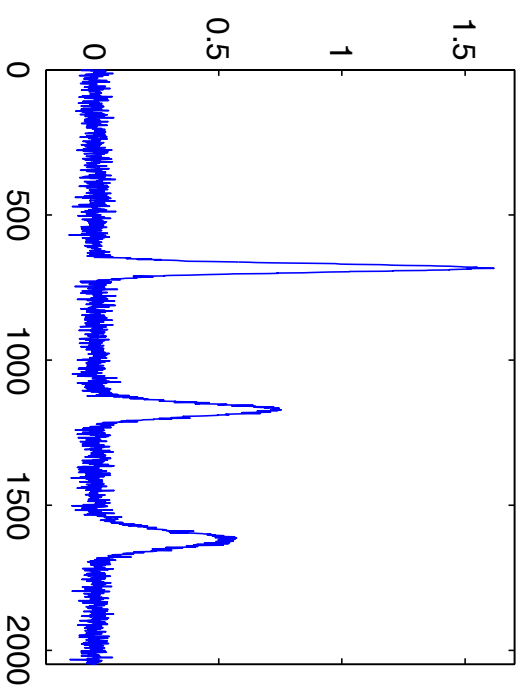
$$\frac{\text{Var}(\xi^T x_i)}{\|\xi\|^2 + \lambda \|D^2 \xi\|^2},$$

- Our example: no choice of λ can *both*
 - preserve peak heights, *and*
 - remove baseline noise

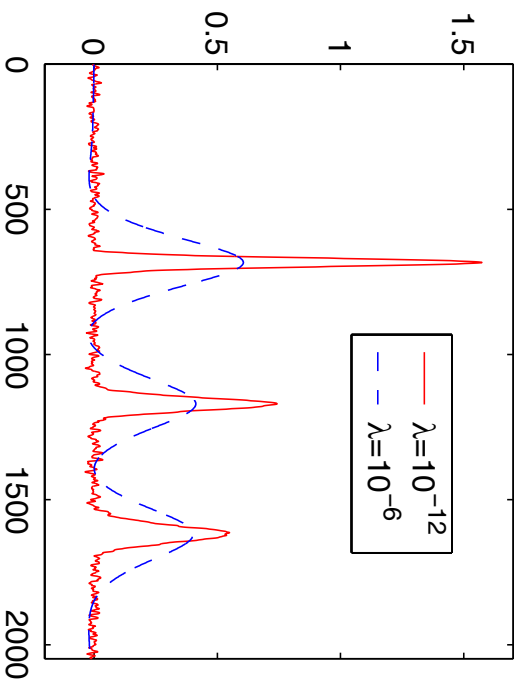
True PC, $p = 2048$, $n = 1024$, $\|p\| = 10$.



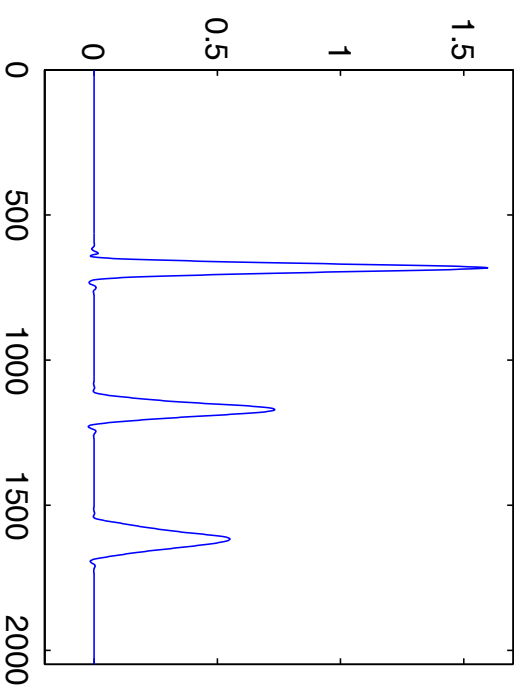
standard PCA



smoothed PCA



ASPCA, $w = 99.5\%$, $K = 372$.



Consistency

Single component model: $x_i = u_i\rho + \sigma z_i$, $i = 1, \dots, n$

Suppose $p(n)/n \rightarrow c$, $\|\rho(n)\| \rightarrow \varrho > 0$.

Then a.s.

$$\overline{\lim}_{n \rightarrow \infty} \sin \angle(\hat{\rho}, \rho) \leq 6\sigma\sqrt{c}/\varrho.$$

In particular:

- **consistent** if $p/n \rightarrow c = 0$
- correct rate if p fixed: $O(\sqrt{c}) = O(1/\sqrt{n})$.
- *but* positive if $c > 0$.

A Matrix Perturbation Theorem

- Suppose
- A, E are symmetric
 - q_1, \hat{q}_1 are the principal eigenvectors of $A, A + E$
 - $\lambda_1(A) - \delta \geq \lambda_2(A) \geq \lambda_\nu(A), \nu \geq 2.$

Then

$$\sin \angle(\hat{q}_1, q_1) \leq (4/\delta) \|E\|_2$$

consequence of more general result for invariant subspaces
(Stewart, Stewart - Sun).

A Multicomponent Model

$$x_i = \sum_{j=1}^m v_i^j \rho^j + \sigma z_i, \quad i = 1, \dots, n$$

- ρ^j are unknown, mutually orthogonal, $\|\rho^1\| \geq \dots \geq \|\rho^m\|$.
- $v_i^j \stackrel{i.i.d.}{\sim} N(0, 1)$ random effects
- $z_i \sim N_p(0, I)$ noise, independent of $\{v_i^j\}$.

For asymptotics,

$$(\|\rho^1(n)\|, \dots, \|\rho^j(n)\|, \dots) \xrightarrow{\ell^1} (\varrho_1, \dots, \varrho_j, \dots).$$

Inconsistency

In either single or multicomponent model,

Theorem If $p/n \rightarrow c > 0$, then

$$\liminf_{n \rightarrow \infty} E \sin \angle(\hat{\rho}^1, \rho^1) > 0.$$

- Noise does not average out in PCA if too many dimensions p relative to n .
- Suggests: reduce p to $k \ll p$ before starting PCA

What goes wrong if $p/n \rightarrow c > 0$

$$\begin{aligned}
 X &= \rho v^T + \sigma Z & S &= n^{-1} X X^T \\
 & & &= \underbrace{\frac{v^T v}{n} \rho \rho^T + \frac{\sigma^2}{n} Z Z^T}_{D} + \underbrace{u \rho^T + \rho^T u}_{B} \\
 & & &= D + B
 \end{aligned}$$

- eigenvalues of $n^{-1} Z Z^T$ do not approach 1:
 - $\lambda_{max}, \lambda_{min} \xrightarrow{a.s.} (1 \pm \sqrt{c})^2$ [Geman, Silverstein]
- $u = \sigma n^{-1} Z v$ does not vanish:
 - $\|u\|^2 \stackrel{D}{=} \sigma^2 \frac{p}{n} \frac{\chi^2(n)}{n} \frac{\chi^2(p)}{p} \xrightarrow{a.s.} \sigma^2 c > 0.$

Why $c > 0$ forces inconsistency

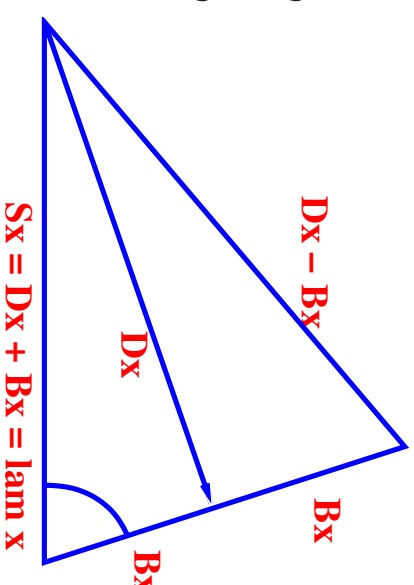
Suppose $S_{\pm} = D \pm B$ have principal e-vectors $\hat{\rho}_{\pm}$

$\hat{\rho}_+$, $\hat{\rho}_-$ have same distribution (symmetry)

But $\|B\rho\| > 0$ and $B\rho \sim \perp (D+B)\rho$

$\Rightarrow \hat{\rho}_+$ and $\hat{\rho}_-$ cannot both be close to ρ

\Rightarrow both are inconsistent.



Outline

- Multiscale Functional Data $p \asymp n$
- Need for Dimension Reduction: Inconsistency
 - ⇒ Sparse PCA
 - Some sampling properties, incl. consistency
- Examples, incl. ECG

Sparse PCA Algorithm

Basis

$$x_i(t) = \sum_{\nu} x_{i\nu} e_{\nu}(t),$$

$$i = 1, \dots, n$$

$$O(np \log p)$$

Subset

$$\hat{\sigma}_{\nu}^2 = \text{Var}\{x_{i\nu}, i = 1, \dots, n\}$$

$\hat{I}_k = \{\nu \leftrightarrow \text{largest } k \text{ variances}\}$

$$O(p \log p)$$

Reduced PCA

on $\{x_{i\nu} : \nu \in \hat{I}_k\}$,

→ eigenvectors $\hat{\rho}_j$

$$O(k^3)$$

Thresholding

$$\hat{\rho}_{j\nu}^* = \eta_H(\hat{\rho}_{j\nu}, \delta)$$

$$O(k)$$

Reconstruct

$$\hat{\rho}_j(t) = \sum_{\nu} \hat{\rho}_{j\nu}^* e_{\nu}(t).$$

$$O(k^2 p)$$

Sparse PCA - Choice of Basis

In basis $\{e_\nu(t)\}$, a population p.c. $\{\rho\}$ has coefficients $\{\rho_\nu\}$:

$$\rho(t) = \sum_{\nu=1}^p \rho_\nu e_\nu(t).$$

Sparsity and weak l_p Say $\rho \in w l_p(C)$ if

$$|\rho_\nu| \leq C\nu^{-1/p}, \quad \nu = 1, 2, \dots$$

- p small \Rightarrow rapid decay of ordered coefficients
- choose basis to exploit sparsity

Wavelet Bases and Sparsity

- Expand ρ in wavelet basis $\{\psi_{jk}\}$: $\rho = \sum_{jk} \rho_{jk} \psi_{jk}$
- order coefficients $\rho_{\nu} = \nu$ -th largest $|\rho_{jk}|$
- **Fact:** smoothness (even non-homogeneous) \Rightarrow sparse wavelet representation:

$$\rho \in B_{p,q}^{\alpha} \quad \Rightarrow \quad (\rho_{\nu}) \in w\ell_p \quad p = 2/(2\alpha + 1).$$

- Hence use wavelet basis here, but algorithm *could* use others..

Sparse PCA Algorithm

Basis

$$x_i(t) = \sum_{\nu} x_{i\nu} e_{\nu}(t),$$

$$i = 1, \dots, n$$

$$O(np \log p)$$

Subset

$$\hat{\sigma}_{\nu}^2 = \text{Var}\{x_{i\nu}, i = 1, \dots, n\}$$

$\hat{I}_k = \{\nu \leftrightarrow \text{largest } k \text{ variances}\}$

$$O(p \log p)$$

Reduced PCA

on $\{x_{i\nu} : \nu \in \hat{I}_k\}$,

→ eigenvectors $\hat{\rho}_j$

$$O(k^3)$$

Thresholding

$$\hat{\rho}_{j\nu}^* = \eta_H(\hat{\rho}_{j\nu}, \delta)$$

$$O(k)$$

Reconstruct

$$\hat{\rho}_j(t) = \sum_{\nu} \hat{\rho}_{j\nu}^* e_{\nu}(t).$$

$$O(k^2 p)$$

Choosing subset size k from data

Aim: Choose \hat{I} to capture most of population p.c.'s ρ variance:

$$\sum_{\nu \in \hat{I}} \rho_{\nu}^2 = w(n) \sum \rho_{\nu}^2, \quad w(n) \nearrow 1.$$

Possibilities: a) $\hat{I} = \{\nu : \hat{\sigma}_{(\nu)}^2 \geq \hat{\sigma}^2(1 + L_n)\}$, or

b) define excess over noise using percentiles of $\chi_{(n)}^2$:

$$\hat{\tau}_{(\nu)}^2 = \hat{\sigma}_{(\nu)}^2 \chi_{(n), \nu/n}^2, \quad \text{and,}$$

$$\hat{I} = \left\{ \nu : \sum_{\nu=1}^{\hat{k}} \hat{\tau}_{(\nu)}^2 \geq w(n) \sum_{\nu} \hat{\tau}_{(\nu)}^2 \right\}.$$

Sparse PCA Algorithm

Basis $x_i(t) = \sum_{\nu} x_{i\nu} e_{\nu}(t),$

$$i = 1, \dots, n$$

$$O(np \log p)$$

Subset $\hat{\sigma}_{\nu}^2 = \text{Var}\{x_{i\nu}, i = 1, \dots, n\}$

$\hat{I}_k = \{\nu \leftrightarrow \text{largest } k \text{ variances}\}$ $O(p \log p)$

Reduced PCA on $\{x_{i\nu} : \nu \in \hat{I}_k\},$

\rightarrow eigenvectors $\hat{\rho}_j$ $O(k^3)$

Thresholding $\hat{\rho}_{j\nu}^* = \eta_H(\hat{\rho}_{j\nu}, \delta)$

$$O(k)$$

Reconstruct $\hat{\rho}_j(t) = \sum_{\nu} \hat{\rho}_{j\nu}^* e_{\nu}(t).$

$$O(k^2 p)$$

Computational Complexity

For standard PCA on $n \times p$ data set X , running time is

$$O((p \wedge n)^3).$$

For $k = o(p) \asymp n$, reduced PCA is $O(k^3)$.

Overall, if say $p \geq n$, reduce from

$$O(p^3) \rightarrow O(k^2 p + np \log p).$$

Sparse PCA Algorithm

Basis

$$x_i(t) = \sum_{\nu} x_{i\nu} e_{\nu}(t),$$

$$i = 1, \dots, n$$

$$O(np \log p)$$

Subset

$$\hat{\sigma}_{\nu}^2 = \text{Var}\{x_{i\nu}, i = 1, \dots, n\}$$

$$\hat{I}_k = \{\nu \leftrightarrow \text{largest } k \text{ variances}\} \quad O(p \log p)$$

Reduced PCA

$$\text{on } \{x_{i\nu} : \nu \in \hat{I}_k\},$$

$$\rightarrow \text{eigenvectors } \hat{\rho}_j \quad O(k^3)$$

Thresholding

$$\hat{\rho}_{j\nu}^* = \eta_H(\hat{\rho}_{j\nu}, \delta)$$

$$O(k)$$

Reconstruct

$$\hat{\rho}_j(t) = \sum_{\nu} \hat{\rho}_{j\nu}^* e_{\nu}(t).$$

$$O(k^2 p)$$

Thresholding subset eigenvectors

Reason: estimated e-vectors on reduced variable set still noisy.

By analogy with wavelet shrinkage in regression: *keep large coefficients, kill noise.*

Here, usual *hard thresholding*:

$$\eta_H(y, \delta) = \begin{cases} y & |y| \geq \delta \\ 0 & \text{otherwise} \end{cases}$$

Other variants: soft, SCAD ...

Choice of δ . For now, use usual “universal” threshold

$$\delta = \hat{\tau}_j \sqrt{2 \log k}, \quad \hat{\tau}_j = MAD\{\hat{\rho}_{j,\nu}, \nu = 1, \dots, k\} / 0.6745.$$

Sparse PCA Algorithm

Basis $x_i(t) = \sum_{\nu} x_{i\nu} e_{\nu}(t),$

$$i = 1, \dots, n$$

$$O(np \log p)$$

Subset $\hat{\sigma}_{\nu}^2 = \text{Var}\{x_{i\nu}, i = 1, \dots, n\}$

$\hat{I}_k = \{\nu \leftrightarrow \text{largest } k \text{ variances}\}$ $O(p \log p)$

Reduced PCA on $\{x_{i\nu} : \nu \in \hat{I}_k\},$

\rightarrow eigenvectors $\hat{\rho}_j$ $O(k^3)$

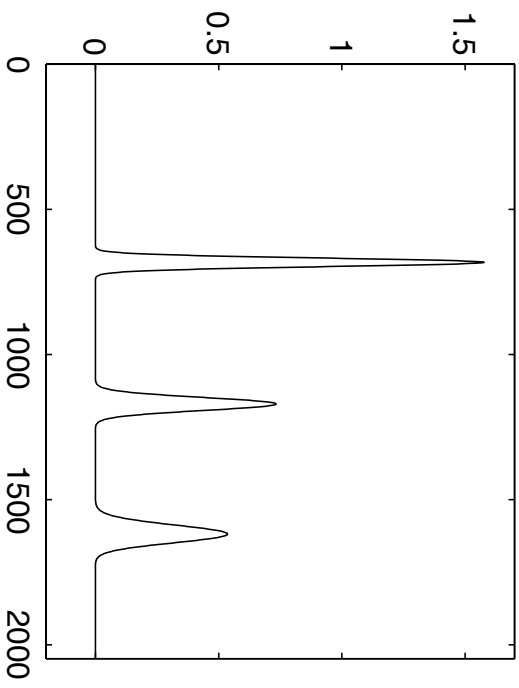
Thresholding $\hat{\rho}_{j\nu}^* = \eta_H(\hat{\rho}_{j\nu}, \delta)$

$$O(k)$$

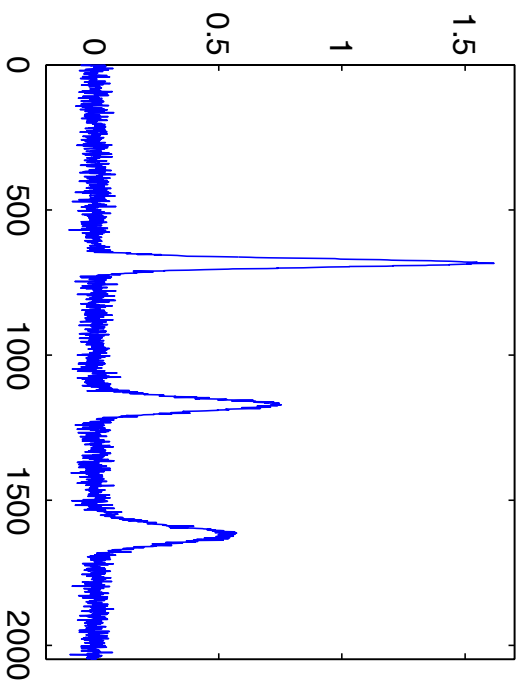
Reconstruct $\hat{\rho}_j(t) = \sum_{\nu} \hat{\rho}_{j\nu}^* e_{\nu}(t).$

$$O(k^2 p)$$

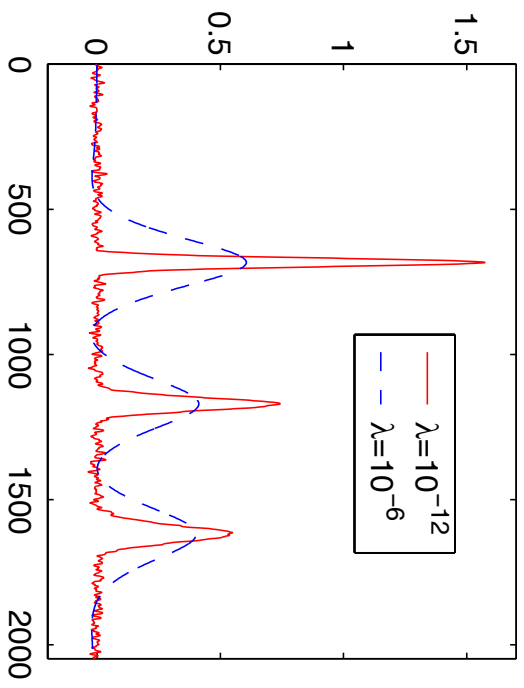
True PC, $p = 2048$, $n = 1024$, $\|p\| = 10$.



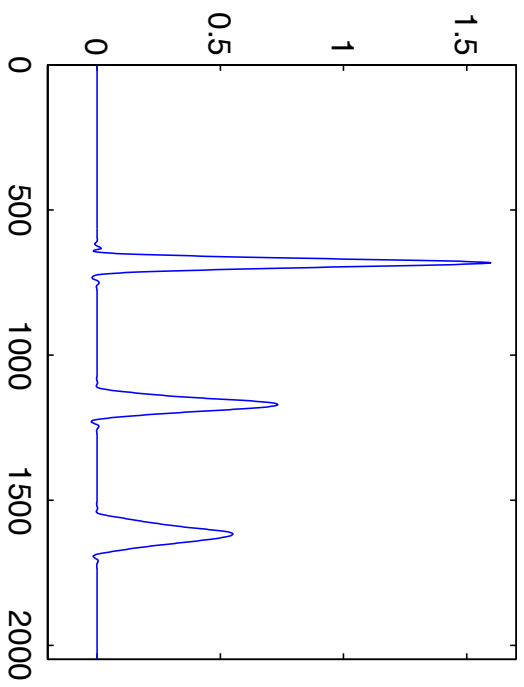
standard PCA

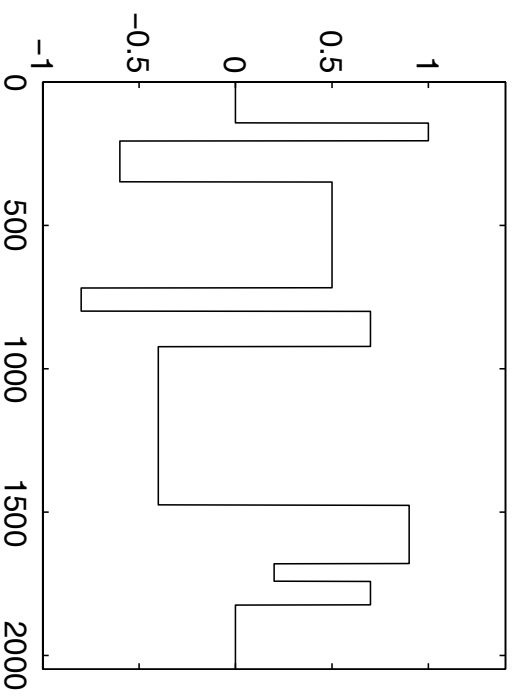


smoothed PCA

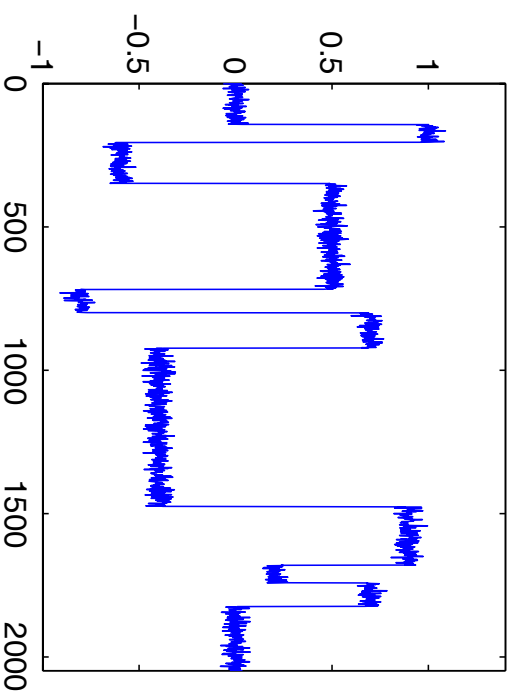
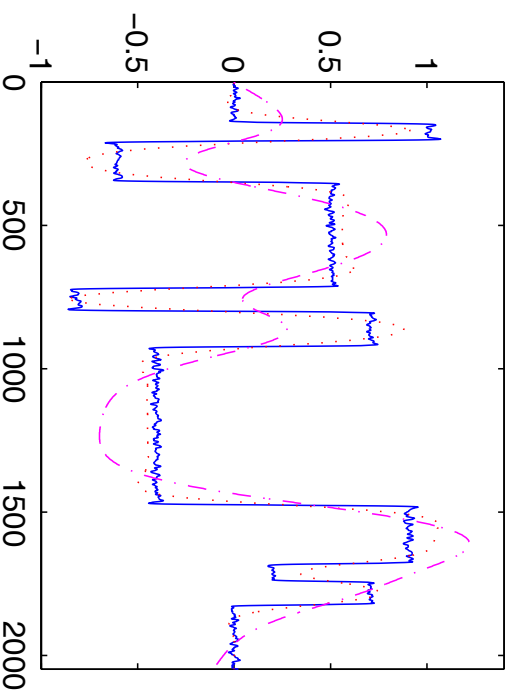
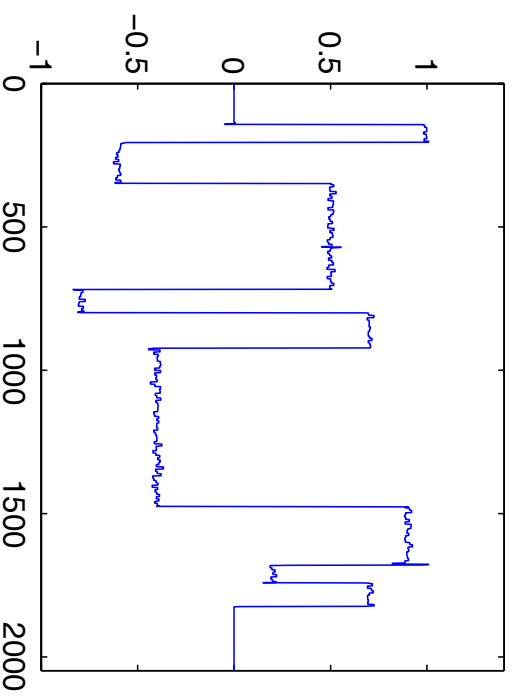


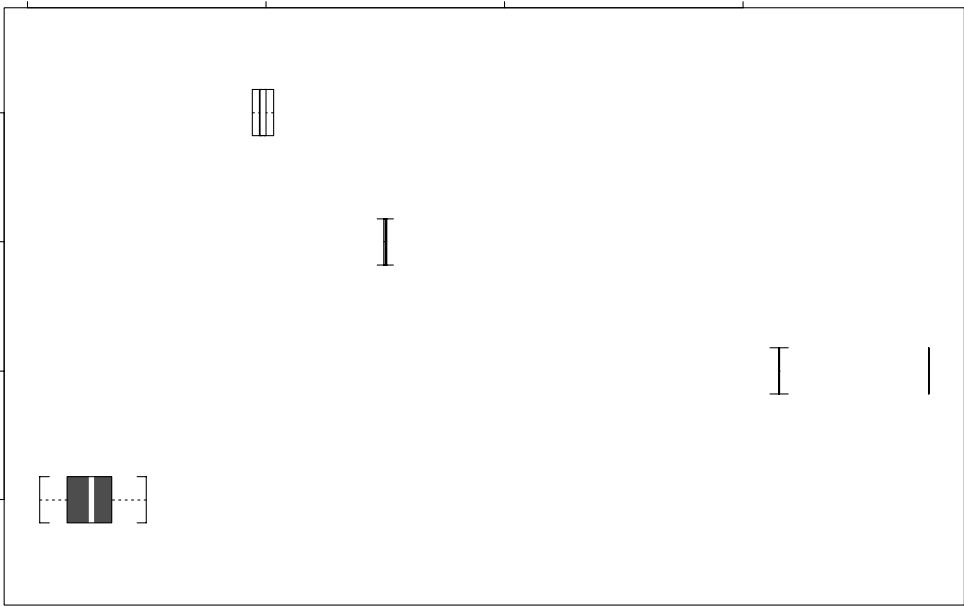
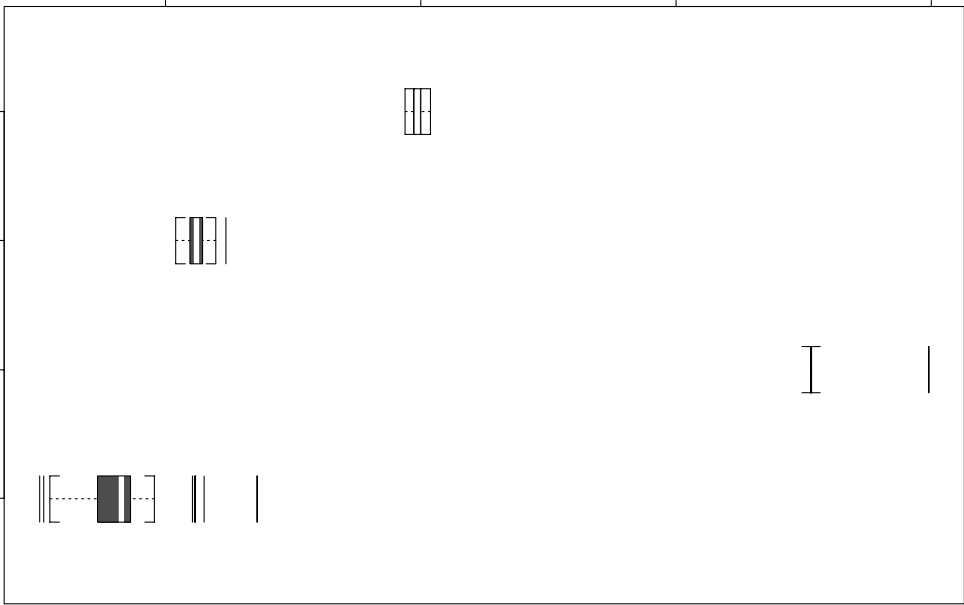
ASPCA, $w = 99.5\%$, $K = 372$.



True PC, $p = 2048$, $n = 1024$, $\|p\| \approx 25$.

Standard PCA

Smoothed PCA, $b:10^{-12}$, $r:10^{-8}$, $m:10^{-6}$.ASPCA, $w = 99.5\%$, $k = 438$.



Speed and Accuracy Comparison

	Standard PCA	Smoothed $\lambda : 10^{-12}$	Smoothed $\lambda : 10^{-6}$	Sparse PCA
ASE (3-peak)	9.681e-04	1.327e-04	3.627e-2	7.500e-05
Time (3-peak)	~ 12min	~ 47 min	~ 43 min	1 min 15 s
ASE (step)	9.715e-04	3.174e-3	1.694e-2	1.947e-04
Time (step)	~ 12min	~ 47 min	~ 46 min	1 min 31 s

Outline

- Multiscale Functional Data $p \asymp n$
- Need for Dimension Reduction: Inconsistency
- Sparse PCA
 - ⇒ Some sampling properties, incl. consistency
- Examples, incl. ECG

Correct Selection Properties

Does $\hat{I}_k \leftrightarrow \hat{\sigma}_{(1)}^2, \dots, \hat{\sigma}_{(k)}^2$ include the “right” variables?

- i.e. include all “large”: $I_{in} = \{\nu : \sigma_\nu^2 \geq \sigma_{(k)}^2 (1 + \alpha_n)\}$,
- and exclude all “small”: $I_{out} = \{\nu : \sigma_\nu^2 \leq \sigma_{(k)}^2 (1 - \alpha_n)\}$.

Define: $FE = \text{False Exclusion} = \{I_{in} \subset \hat{I}_k\}^c$

$FI = \text{False Inclusion} = \{I_{out} \subset \hat{I}_k^c\}^c$

Theorem Suppose $\hat{\sigma}_\nu^2 \sim \sigma_\nu^2 \chi_{(n)}^2 / n$. If $\alpha = \gamma \sqrt{\log n / n}$, then

$$P\{FE \cup FI\} \leq 3pk n^{-b\gamma^2},$$

E.g.: $(k, p, n) = (50, 1000, 1000)$, $1 + \alpha_n = 1.25^2$, $P \leq .02$

Consistency of Sparse PCA

Single component model. Suppose (i) $p/n \rightarrow c > 0$,

$$(ii) \quad \|\rho(n)\| \rightarrow \varrho > 0.$$

Assume *Sparsity*: $\rho(n) \in w_{k_p}(C)$ uniformly in n

Subset selection rule: $\hat{I} = \{v : \hat{\sigma}_v^2 > \sigma^2(1 + c\sqrt{2\log p}\sqrt{2/n})\}$

Let $\hat{\rho}$ denote sparse PCA estimate based on \hat{I} .

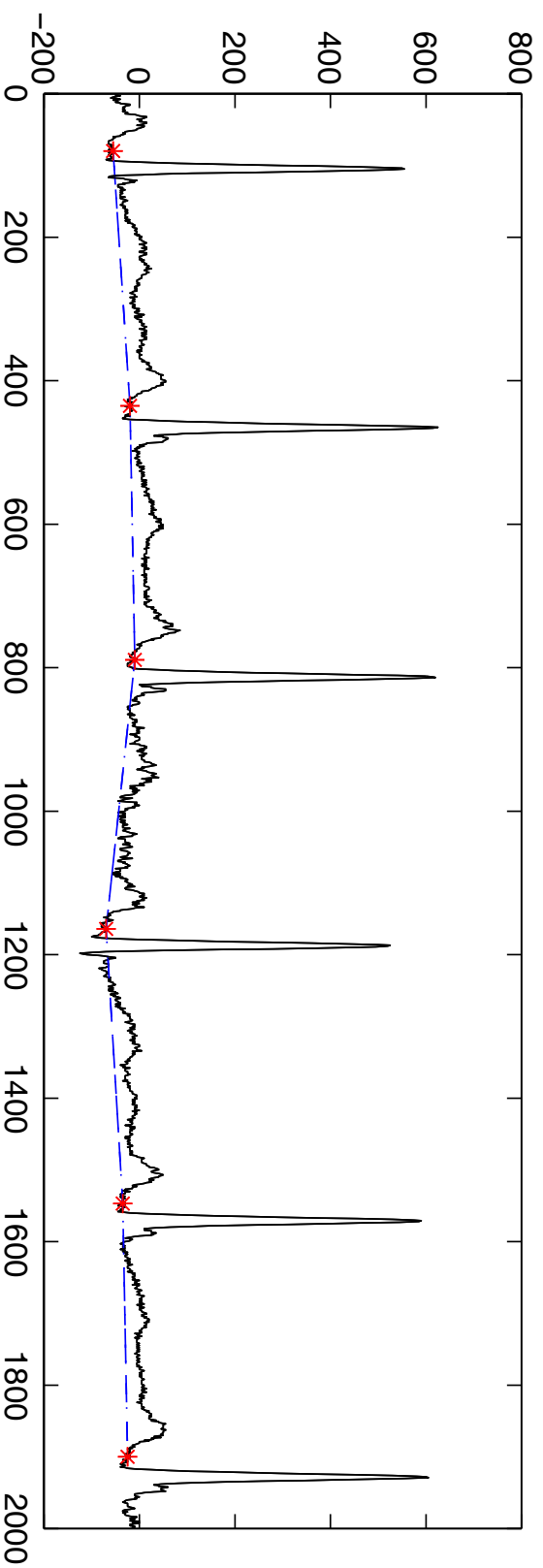
Theorem $\angle(\hat{\rho}, \rho) \xrightarrow{a.s.} 0.$

Outline

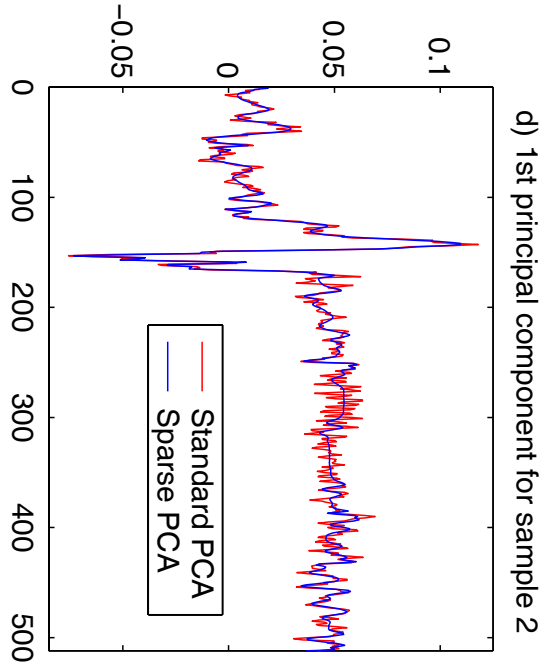
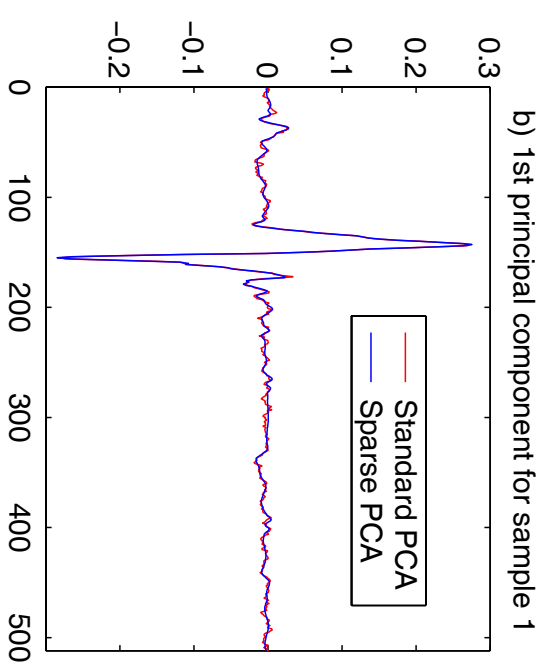
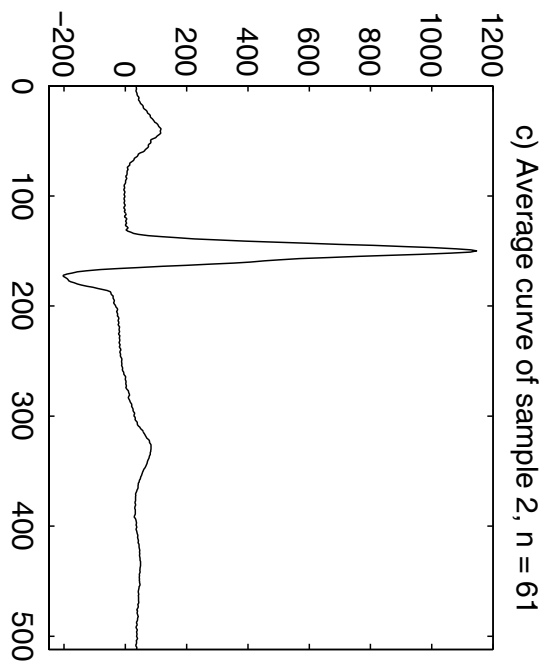
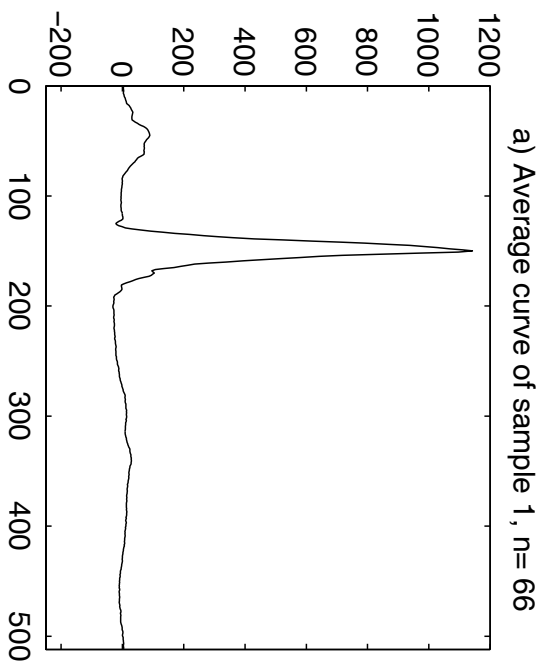
- Multiscale Functional Data $p \asymp n$
 - Need for Dimension Reduction: Inconsistency
 - Sparse PCA
 - Some sampling properties, incl. consistency
- ⇒ Examples, incl. ECG

ECG example

Preprocessing: piecewise linear baseline wander removal,
registration at R —wave maximum,
interpolation to 512 samples per cycle



Data: V. Froelicher, J. Froning, Palo Alto VA Hospital.



Remarks

- both p.c.s correspond to change in shape of R-wave peak
- noise is larger in second case: $\hat{\sigma}_1^2 = 24.97$, $\hat{\sigma}_2^2 = 82.12$.
- Sparse PCA uses $< 10\%$ of computing time for standard PCA.

Much more to do here:

- more work on interpretation with cardiologists
- effect of registration
- thresholding, multiple leads, ...

Conclusion: Main Themes Again

- initial dimension reduction before PCA
 - otherwise, inconsistency!
- use basis with sparse representation
 - so that little is lost in initial dim reduction
- Background role for large random matrices

SAMSI, 'Random Matrices' and FDA?

- Statistical and Applied Math Sci Institute, at NISS, NC.
- RMT active area in math, physics and probability
- Possible semester program, Spring 2005: bring together statisticians and applied math people
- Aim: formulate methodologic & theory questions from statistical areas to profit from RMT tools
- Possible statistical areas: climatology (EOFs..), document retrieval, Functional Data Analysis
- Array of problems where p is not fixed, or not small....

References

- Ramsay, J. O. & Silverman, B. W. (1997), *Functional Data Analysis*, Springer.
- Rice, J. A. & Silverman, B. W. (1991), 'Estimating the mean and covariance structure nonparametrically when the data are curves', *Journal of the Royal Statistical Society, Series B (Methodological)* **53**, 233–243.
- Silverman, B. W. (1996), 'Smoothed functional principal components analysis by choice of norm', *Annals of Statistics* **24**(1), 1–24.