
An Introduction to Nonlinear Principal Component Analysis

Adam Monahan

`monahana@uvic.ca`

School of Earth and Ocean Sciences
University of Victoria

Overview

- Dimensionality reduction

Overview

- Dimensionality reduction
- Principal Component Analysis

Overview

- Dimensionality reduction
- Principal Component Analysis
- Nonlinear PCA

Overview

- Dimensionality reduction
- Principal Component Analysis
- Nonlinear PCA
 - theory

Overview

- Dimensionality reduction
- Principal Component Analysis
- Nonlinear PCA
 - theory
 - implementation

Overview

- Dimensionality reduction
- Principal Component Analysis
- Nonlinear PCA
 - theory
 - implementation
- Applications of NLPCA

Overview

- Dimensionality reduction
- Principal Component Analysis
- Nonlinear PCA
 - theory
 - implementation
- Applications of NLPCA
 - Lorenz attractor

Overview

- Dimensionality reduction
- Principal Component Analysis
- Nonlinear PCA
 - theory
 - implementation
- Applications of NLPCA
 - Lorenz attractor
 - NH Tropospheric LFV

Overview

- Dimensionality reduction
- Principal Component Analysis
- Nonlinear PCA
 - theory
 - implementation
- Applications of NLPCA
 - Lorenz attractor
 - NH Tropospheric LFV
- Conclusions

Dimensionality Reduction

- Climate datasets made up of time series at individual stations/geographical locations

Dimensionality Reduction

- Climate datasets made up of time series at individual stations/geographical locations
- Typical dataset has $P \sim O(10^3)$ time series

Dimensionality Reduction

- Climate datasets made up of time series at individual stations/geographical locations
- Typical dataset has $P \sim O(10^3)$ time series
- Organised structure in atmosphere/ocean flows

Dimensionality Reduction

- Climate datasets made up of time series at individual stations/geographical locations
 - Typical dataset has $P \sim O(10^3)$ time series
 - Organised structure in atmosphere/ocean flows
- ⇒ time series at different locations not independent

Dimensionality Reduction

- Climate datasets made up of time series at individual stations/geographical locations
 - Typical dataset has $P \sim O(10^3)$ time series
 - Organised structure in atmosphere/ocean flows
- ⇒ time series at different locations not independent
- ⇒ data does not fill out isotropic cloud of points in \mathbf{R}^P , but clusters around lower-dimensional surface (reflecting the “attractor”)

Dimensionality Reduction

- Climate datasets made up of time series at individual stations/geographical locations
 - Typical dataset has $P \sim O(10^3)$ time series
 - Organised structure in atmosphere/ocean flows
- ⇒ time series at different locations not independent
- ⇒ data does not fill out isotropic cloud of points in \mathbf{R}^P , but clusters around lower-dimensional surface (reflecting the “attractor”)
- Goal of *dimensionality reduction* in climate diagnostics is to characterise such structures in climate datasets



Dimensionality Reduction

- Realising this goal has both theoretical and practical difficulties:

Dimensionality Reduction

- Realising this goal has both theoretical and practical difficulties:
- **Theoretical:**

Dimensionality Reduction

- Realising this goal has both theoretical and practical difficulties:
- **Theoretical:**
 - what is the precise definition of “structure”?

Dimensionality Reduction

- Realising this goal has both theoretical and practical difficulties:
- **Theoretical:**
 - what is the precise definition of “structure”?
 - how to formulate appropriate statistical model?

Dimensionality Reduction

- Realising this goal has both theoretical and practical difficulties:
- **Theoretical:**
 - what is the precise definition of “structure”?
 - how to formulate appropriate statistical model?
- **Practical:**

Dimensionality Reduction

- Realising this goal has both theoretical and practical difficulties:
- **Theoretical:**
 - what is the precise definition of “structure”?
 - how to formulate appropriate statistical model?
- **Practical:**
 - many important observational climate datasets quite short, with $O(10) - O(1000)$ statistical degrees of freedom

Dimensionality Reduction

- Realising this goal has both theoretical and practical difficulties:
- **Theoretical:**
 - what is the precise definition of “structure”?
 - how to formulate appropriate statistical model?
- **Practical:**
 - many important observational climate datasets quite short, with $O(10) - O(1000)$ statistical degrees of freedom
 - what degree of “structure” can be robustly diagnosed with existing data?

Principal Component Analysis

- A classical approach to dimensionality **principal component analysis (PCA)**

Principal Component Analysis

- A classical approach to dimensionality **principal component analysis (PCA)**
- Look for M -dimensional hyperplane approximation, optimal in least-squares sense

$$\mathbf{X}(t) = \sum_{k=1}^M \langle \mathbf{X}(t), \mathbf{e}_k \rangle \mathbf{e}_k + \epsilon(t)$$

Principal Component Analysis

- A classical approach to dimensionality **principal component analysis (PCA)**
- Look for M -dimensional hyperplane approximation, optimal in least-squares sense

$$\mathbf{X}(t) = \sum_{k=1}^M \langle \mathbf{X}(t), \mathbf{e}_k \rangle \mathbf{e}_k + \epsilon(t)$$

- minimising $E \{ \|\epsilon^2\| \}$

Principal Component Analysis

- A classical approach to dimensionality **principal component analysis (PCA)**
- Look for M -dimensional hyperplane approximation, optimal in least-squares sense

$$\mathbf{X}(t) = \sum_{k=1}^M \langle \mathbf{X}(t), \mathbf{e}_k \rangle \mathbf{e}_k + \epsilon(t)$$

- minimising $E \{ \|\epsilon^2\| \}$
- inner product often (not always) simple dot product

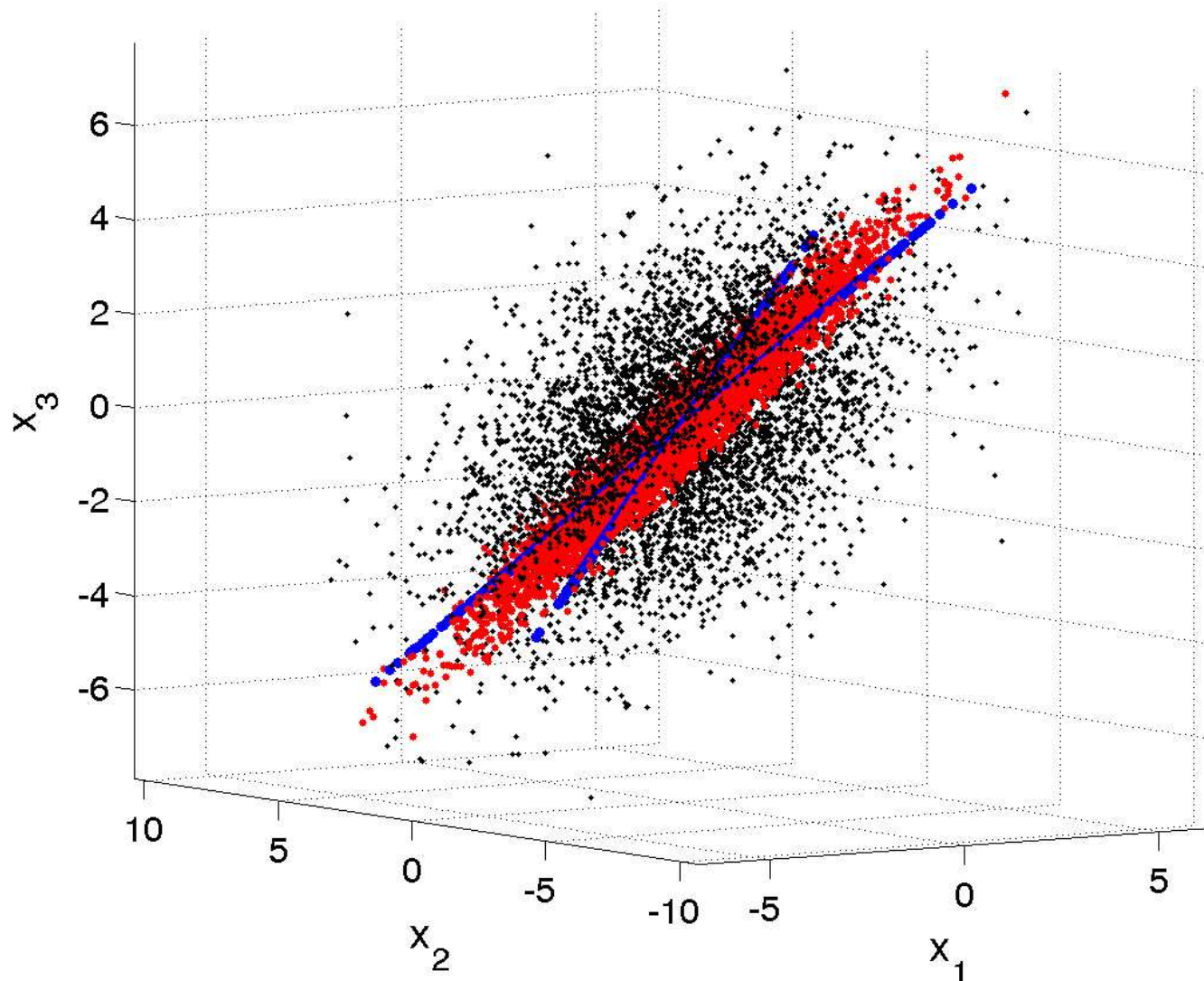
Principal Component Analysis

- A classical approach to dimensionality **principal component analysis (PCA)**
- Look for M -dimensional hyperplane approximation, optimal in least-squares sense

$$\mathbf{X}(t) = \sum_{k=1}^M \langle \mathbf{X}(t), \mathbf{e}_k \rangle \mathbf{e}_k + \epsilon(t)$$

- minimising $E \{ \|\epsilon\|^2 \}$
- inner product often (not always) simple dot product
- Vectors \mathbf{e}_k are the **empirical orthogonal functions (EOFs)**

Principal Component Analysis



Principal Component Analysis

- Operationally, EOFs are found as eigenvectors of covariance matrix (in appropriate norm)

Principal Component Analysis

- Operationally, EOFs are found as eigenvectors of covariance matrix (in appropriate norm)
- PCA optimally efficient characterisation of Gaussian data

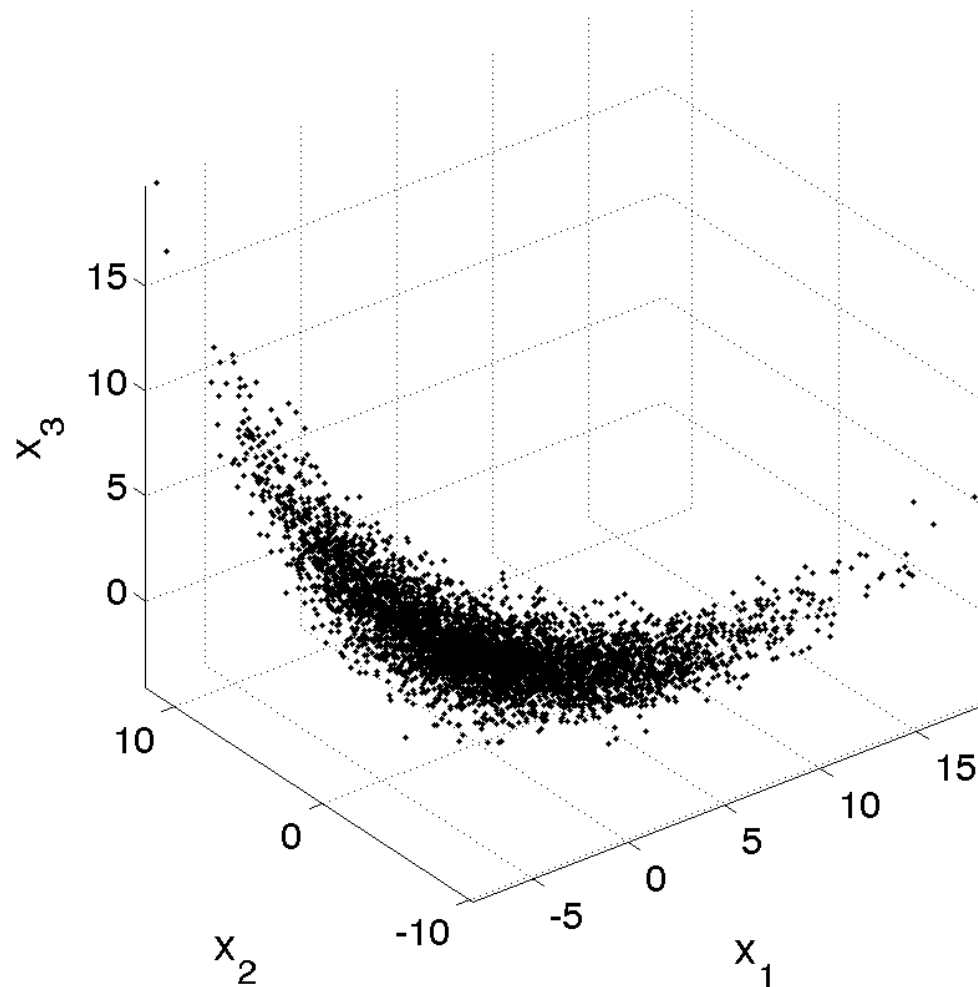
Principal Component Analysis

- Operationally, EOFs are found as eigenvectors of covariance matrix (in appropriate norm)
- PCA optimally efficient characterisation of Gaussian data
- More generally: PCA provides optimally parsimonious data compression for any dataset whose distribution lies along orthogonal axes

Principal Component Analysis

- Operationally, EOFs are found as eigenvectors of covariance matrix (in appropriate norm)
- PCA optimally efficient characterisation of Gaussian data
- More generally: PCA provides optimally parsimonious data compression for any dataset whose distribution lies along orthogonal axes
- But what if the underlying low-dimensional structure is curved rather than straight?
(cigars vs. bananas)

Nonlinear Low-Dimensional Structure



Nonlinear PCA

- An approach to diagnosing nonlinear low-dimensional structure is *Nonlinear PCA (NLPCA)*

Nonlinear PCA

- An approach to diagnosing nonlinear low-dimensional structure is *Nonlinear PCA (NLPCA)*
- Goal: find functions (with $M < P$)

$$\mathbf{s}_f : \mathbb{R}^P \rightarrow \mathbb{R}^M, \quad \mathbf{f} : \mathbb{R}^M \rightarrow \mathbb{R}^P$$

such that

$$\mathbf{X}(t) = (\mathbf{f} \circ \mathbf{s}_f)(\mathbf{X}(t)) + \epsilon(\mathbf{t})$$

where

Nonlinear PCA

- An approach to diagnosing nonlinear low-dimensional structure is *Nonlinear PCA (NLPCA)*
- Goal: find functions (with $M < P$)

$$\mathbf{s}_f : \mathbb{R}^P \rightarrow \mathbb{R}^M, \quad \mathbf{f} : \mathbb{R}^M \rightarrow \mathbb{R}^P$$

such that

$$\mathbf{X}(t) = (\mathbf{f} \circ \mathbf{s}_f)(\mathbf{X}(t)) + \epsilon(\mathbf{t})$$

where

- $E \{ \|\epsilon^2\| \}$ is minimised

Nonlinear PCA

- An approach to diagnosing nonlinear low-dimensional structure is *Nonlinear PCA (NLPCA)*
- Goal: find functions (with $M < P$)

$$\mathbf{s}_f : \mathbb{R}^P \rightarrow \mathbb{R}^M, \quad \mathbf{f} : \mathbb{R}^M \rightarrow \mathbb{R}^P$$

such that

$$\mathbf{X}(t) = (\mathbf{f} \circ \mathbf{s}_f)(\mathbf{X}(t)) + \epsilon(\mathbf{t})$$

where

- $E \{ \|\epsilon^2\| \}$ is minimised
- $\mathbf{f}(\lambda) \sim$ approximation manifold

Nonlinear PCA

- An approach to diagnosing nonlinear low-dimensional structure is *Nonlinear PCA (NLPCA)*
- Goal: find functions (with $M < P$)

$$\mathbf{s}_f : \mathbb{R}^P \rightarrow \mathbb{R}^M, \quad \mathbf{f} : \mathbb{R}^M \rightarrow \mathbb{R}^P$$

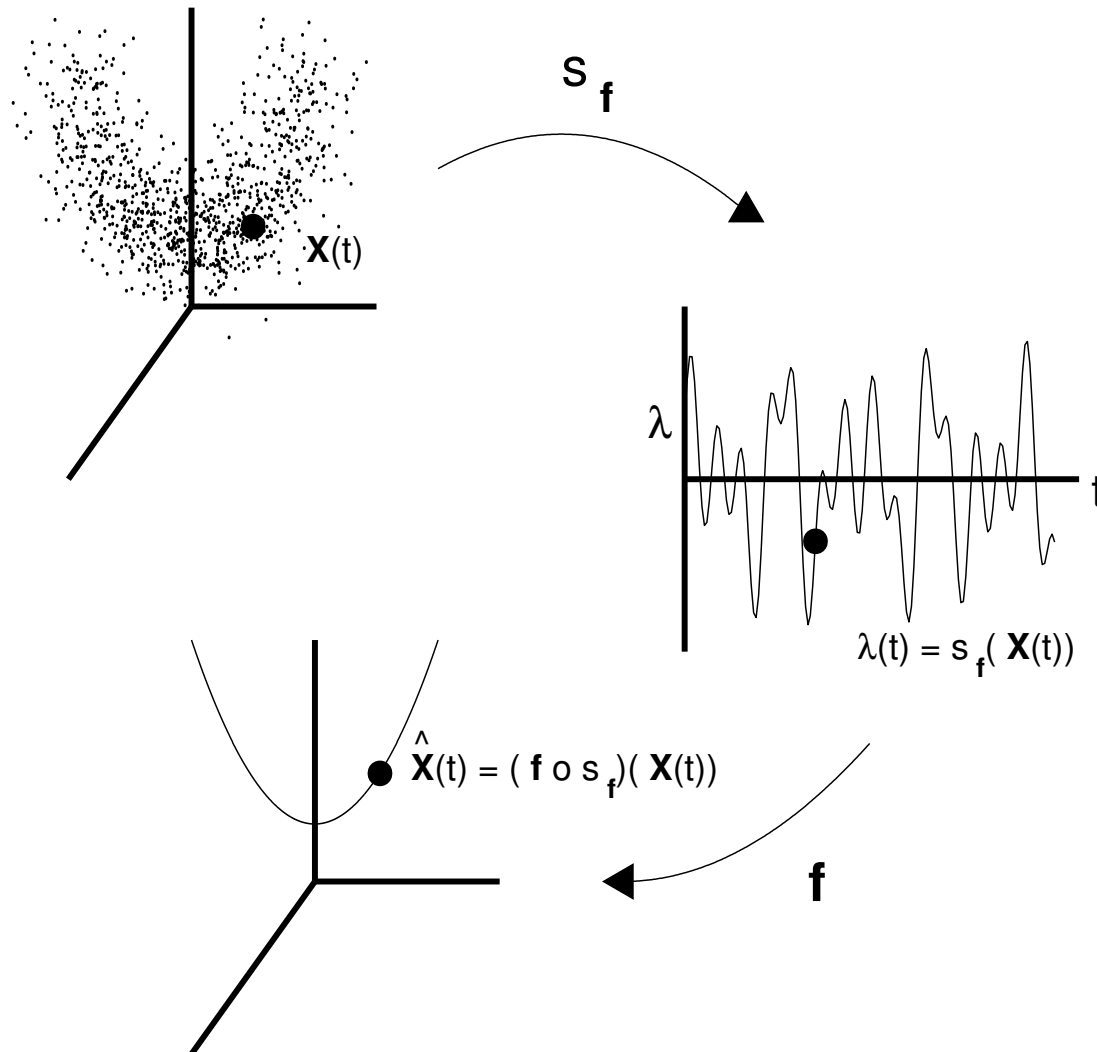
such that

$$\mathbf{X}(t) = (\mathbf{f} \circ \mathbf{s}_f)(\mathbf{X}(t)) + \epsilon(\mathbf{t})$$

where

- $E \{ \|\epsilon^2\| \}$ is minimised
- $\mathbf{f}(\lambda) \sim$ approximation manifold
- $\lambda(t) = \mathbf{s}_f(\mathbf{X}(t)) \sim$ manifold parameterisation (time series)

Nonlinear PCA



From Monahan, Fyfe, and Pandolfo (2003)

Nonlinear PCA

- As with PCA, “fraction of variance explained” is a measure of quality of approximation

Nonlinear PCA

- As with PCA, “fraction of variance explained” is a measure of quality of approximation
- PCA is a special case of NLPCA

Nonlinear PCA

- As with PCA, “fraction of variance explained” is a measure of quality of approximation
- PCA is a special case of NLPCA
- When implemented, NLPCA should reduce to PCA if:

Nonlinear PCA

- As with PCA, “fraction of variance explained” is a measure of quality of approximation
- PCA is a special case of NLPCA
- When implemented, NLPCA should reduce to PCA if:
 - data is Gaussian

Nonlinear PCA

- As with PCA, “fraction of variance explained” is a measure of quality of approximation
- PCA is a special case of NLPCA
- When implemented, NLPCA should reduce to PCA if:
 - data is Gaussian
 - not enough data is available to robustly characterise non-Gaussian structure

NLPCA: Implementation

- Implemented NLPCA using neural networks (convenient, not necessary)

NLPCA: Implementation

- Implemented NLPCA using neural networks (convenient, not necessary)
- Parameter estimation more difficult than for PCA

NLPCA: Implementation

- Implemented NLPCA using neural networks (convenient, not necessary)
- Parameter estimation more difficult than for PCA
- PCA model is linear in statistical parameters:

$$Y = MX$$

so variational problem has unique analytic solution

NLPCA: Implementation

- Implemented NLPCA using neural networks (convenient, not necessary)
- Parameter estimation more difficult than for PCA
- PCA model is linear in statistical parameters:

$$Y = MX$$

so variational problem has unique analytic solution

- NLPCA model nonlinear in model parameters, so solution

NLPCA: Implementation

- Implemented NLPCA using neural networks (convenient, not necessary)
- Parameter estimation more difficult than for PCA
- PCA model is linear in statistical parameters:

$$Y = MX$$

so variational problem has unique analytic solution

- NLPCA model nonlinear in model parameters, so solution
 - may not be unique

NLPCA: Implementation

- Implemented NLPCA using neural networks (convenient, not necessary)
- Parameter estimation more difficult than for PCA
- PCA model is linear in statistical parameters:

$$Y = MX$$

so variational problem has unique analytic solution

- NLPCA model nonlinear in model parameters, so solution
 - may not be unique
 - must be found through numerical minimisation

NLPCA: Parameter Estimation

- Two fundamental issues regarding parameter estimation common to *all* statistical models:

NLPCA: Parameter Estimation

- Two fundamental issues regarding parameter estimation common to *all* statistical models:

Reproducibility

NLPCA: Parameter Estimation

- Two fundamental issues regarding parameter estimation common to *all* statistical models:

Reproducibility

- model must be robust to the introduction of new data

NLPCA: Parameter Estimation

- Two fundamental issues regarding parameter estimation common to *all* statistical models:

Reproducibility

- model must be robust to the introduction of new data
- new observations shouldn't fundamentally change model

NLPCA: Parameter Estimation

- Two fundamental issues regarding parameter estimation common to *all* statistical models:

Reproducibility

- model must be robust to the introduction of new data
- new observations shouldn't fundamentally change model

Classifiability:

NLPCA: Parameter Estimation

- Two fundamental issues regarding parameter estimation common to *all* statistical models:

Reproducibility

- model must be robust to the introduction of new data
- new observations shouldn't fundamentally change model

Classifiability:

- model must be robust to details of optimisation procedure

NLPCA: Parameter Estimation

- Two fundamental issues regarding parameter estimation common to *all* statistical models:

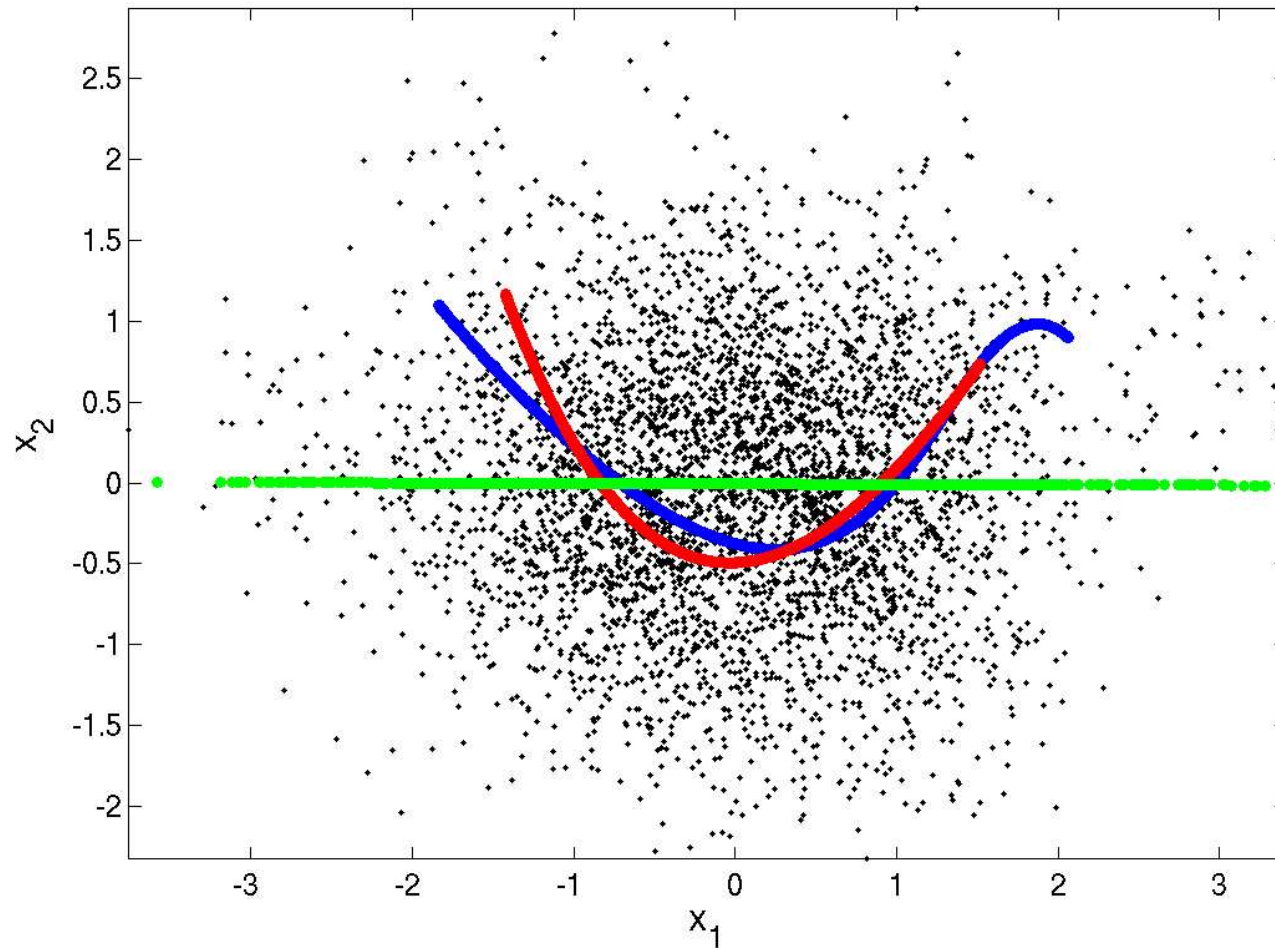
Reproducibility

- model must be robust to the introduction of new data
- new observations shouldn't fundamentally change model

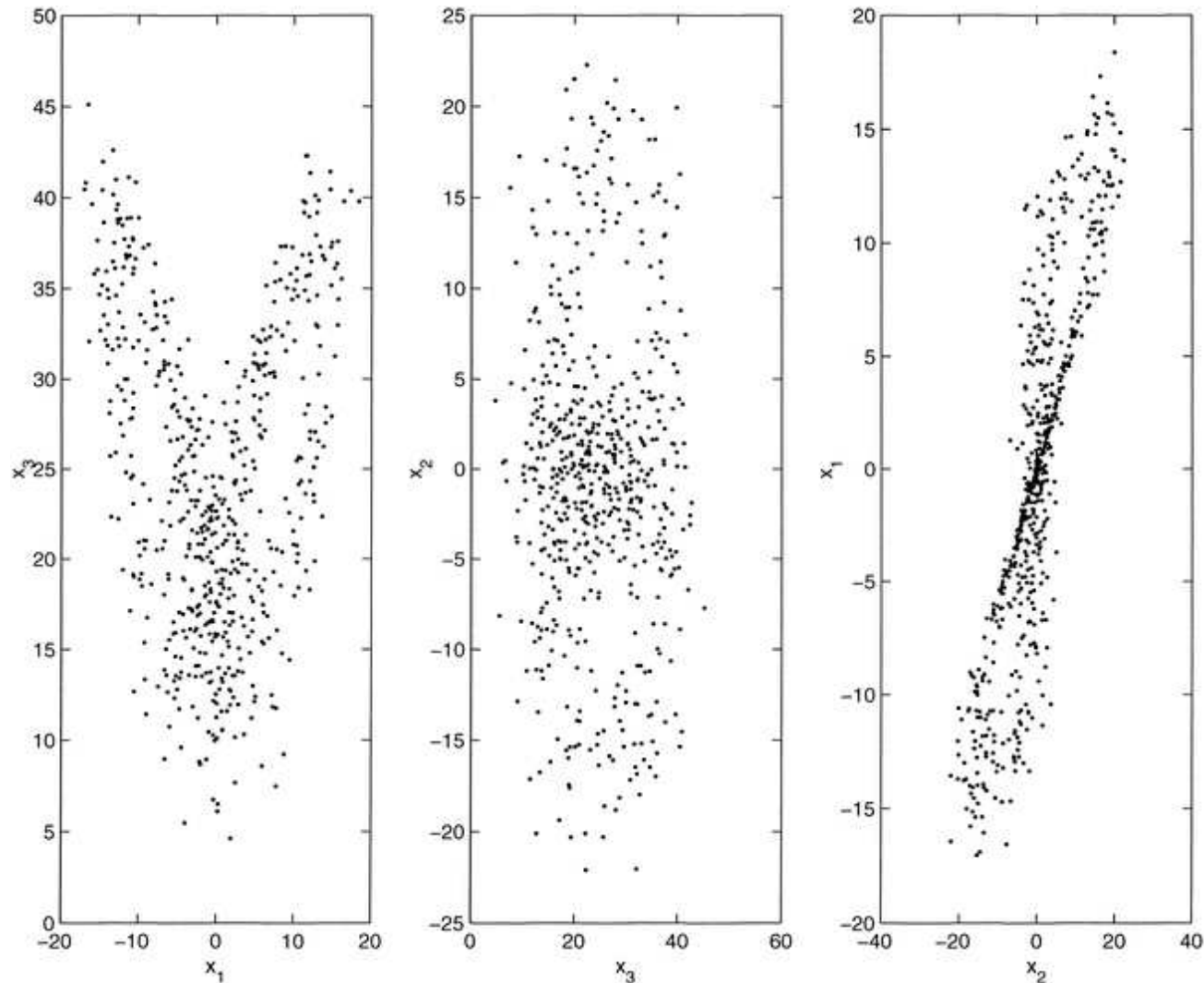
Classifiability:

- model must be robust to details of optimisation procedure
- model shouldn't depend on initial parameter values

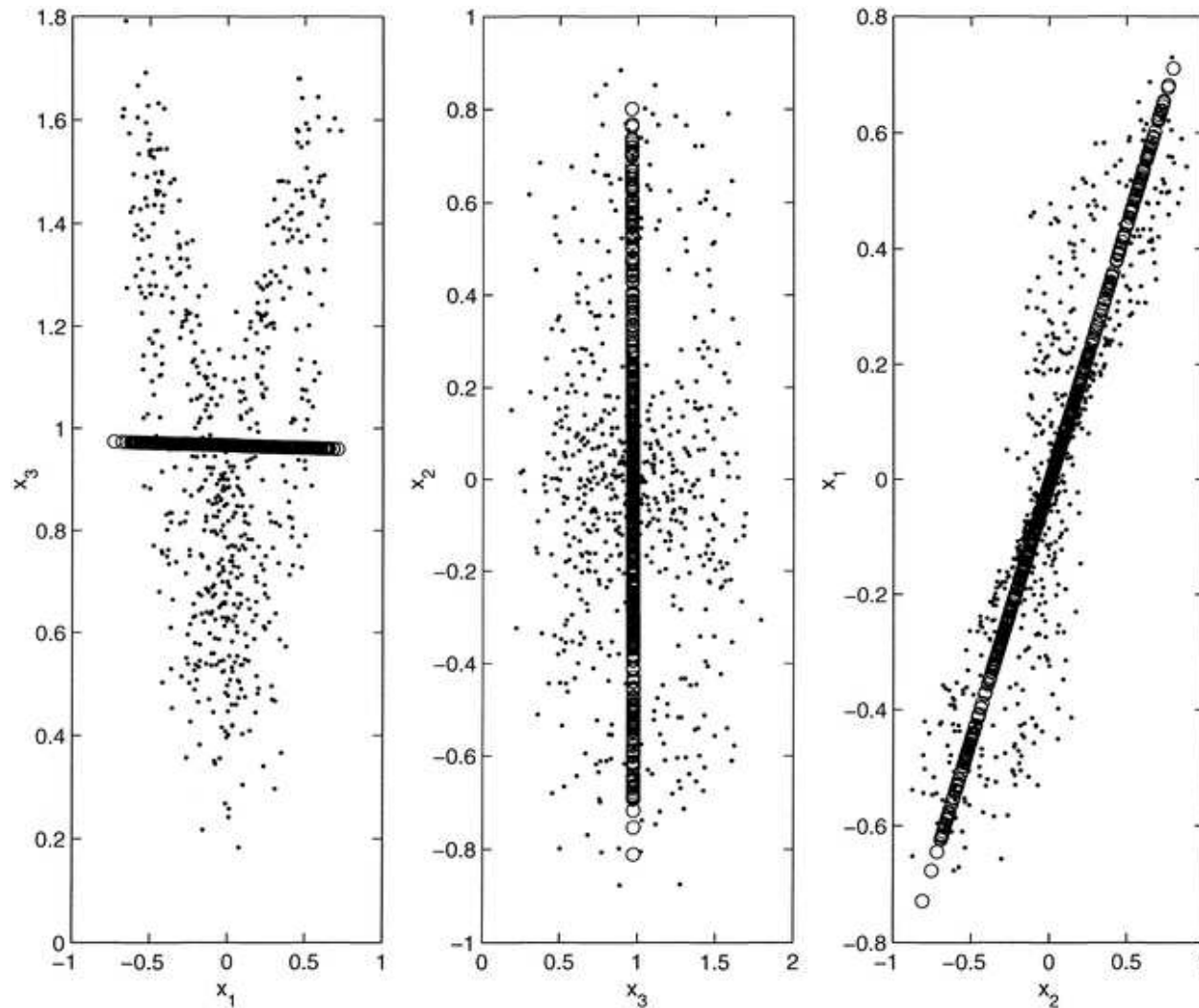
NLPCA: Synthetic Gaussian Data



Applications of NLPCA: Lorenz Attractor



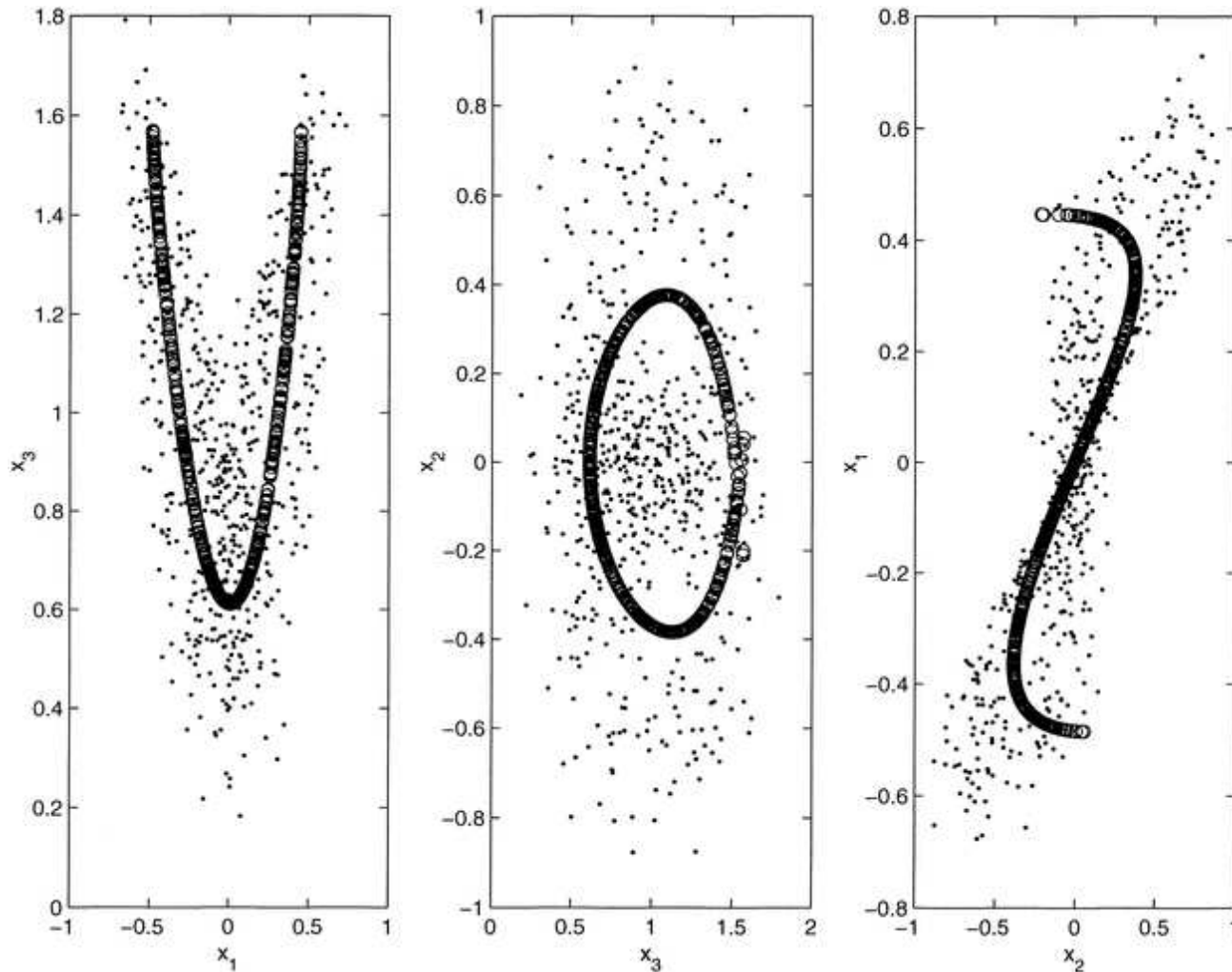
Applications of NLPCA: Lorenz Attractor



UVic

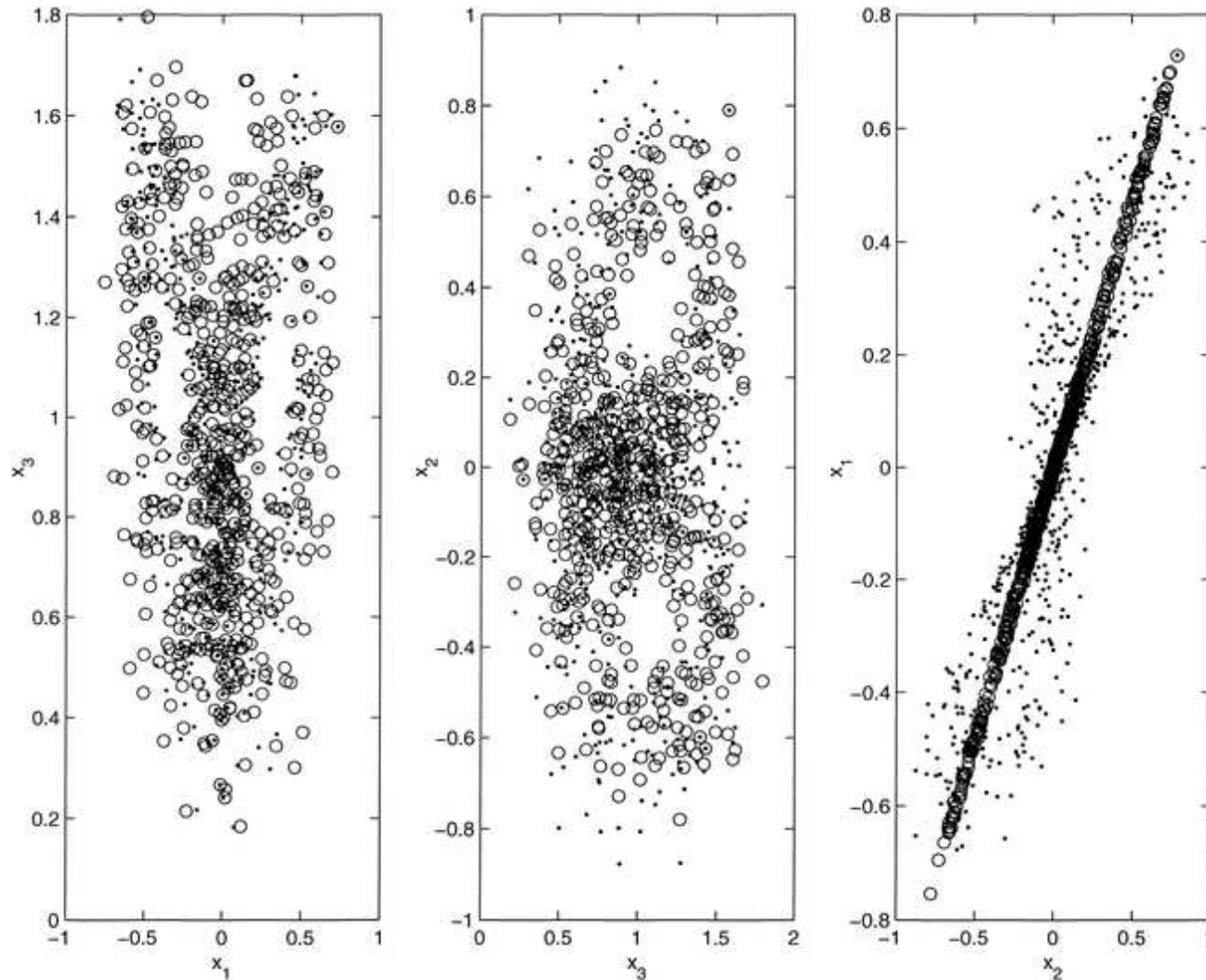
1D PCA approximation (60%)

Applications of NLPCA: Lorenz Attractor



1D NLPCA approximation (76%)

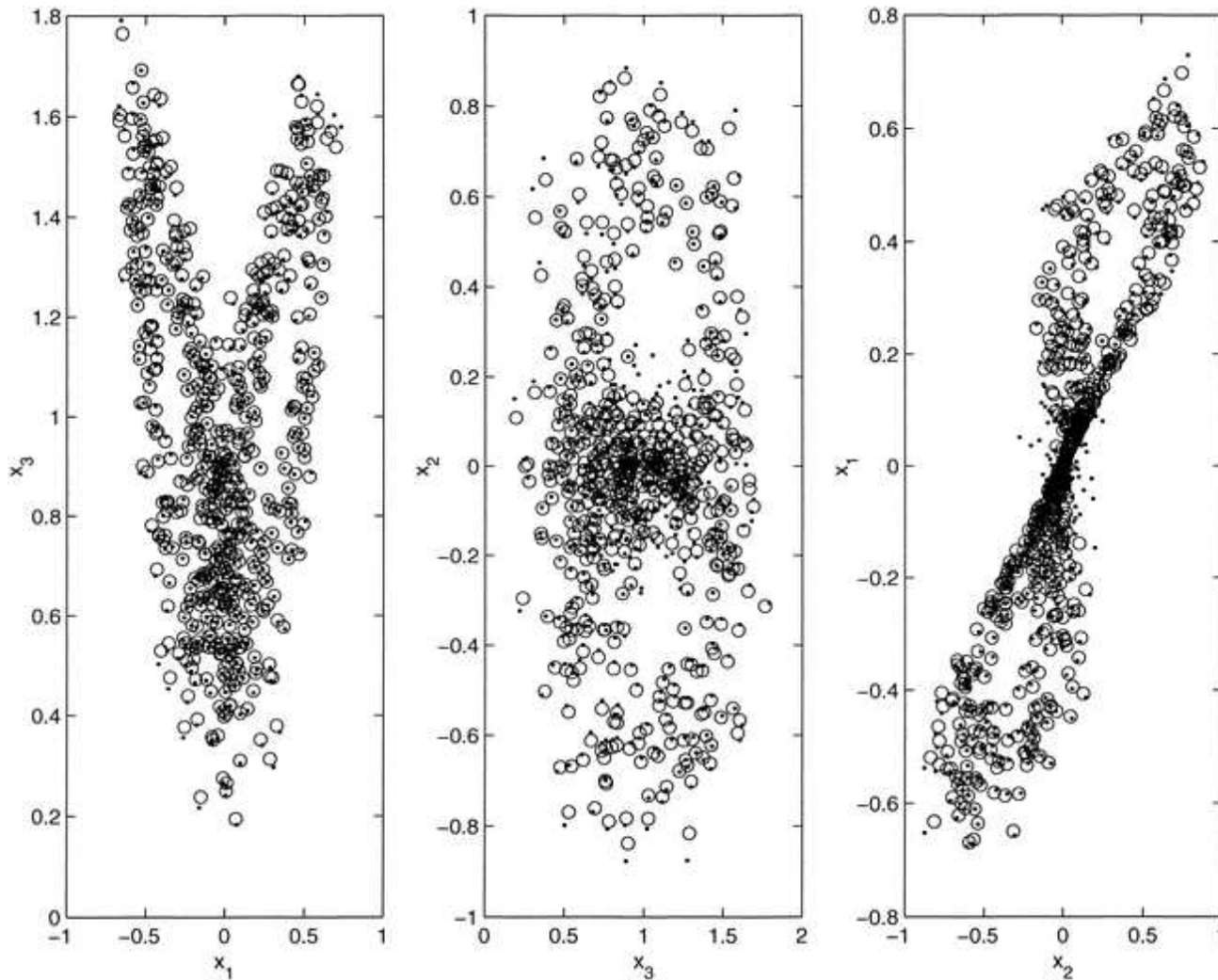
Applications of NLPCA: Lorenz Attractor



UVic

2D PCA approximation (94%)

Applications of NLPCA: Lorenz Attractor

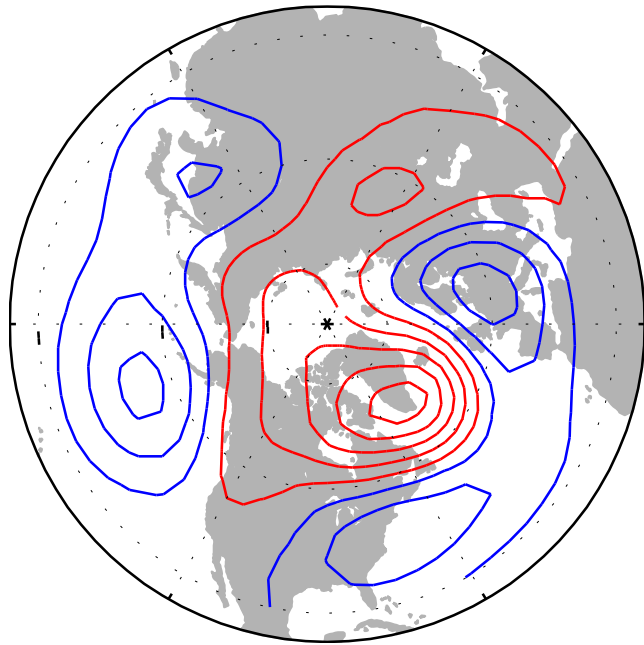


UVic

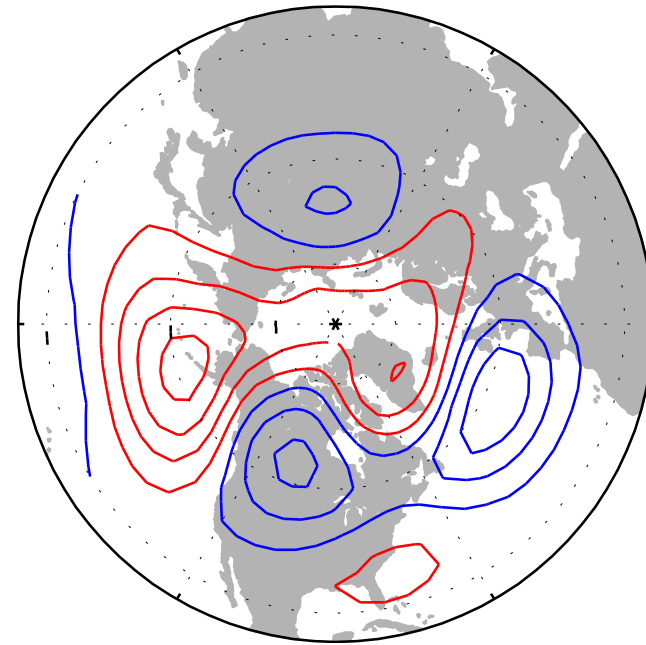
2D NLPCA approximation (97%)

Applications of NLPCA: NH Tropospheric LFV

EOF₁

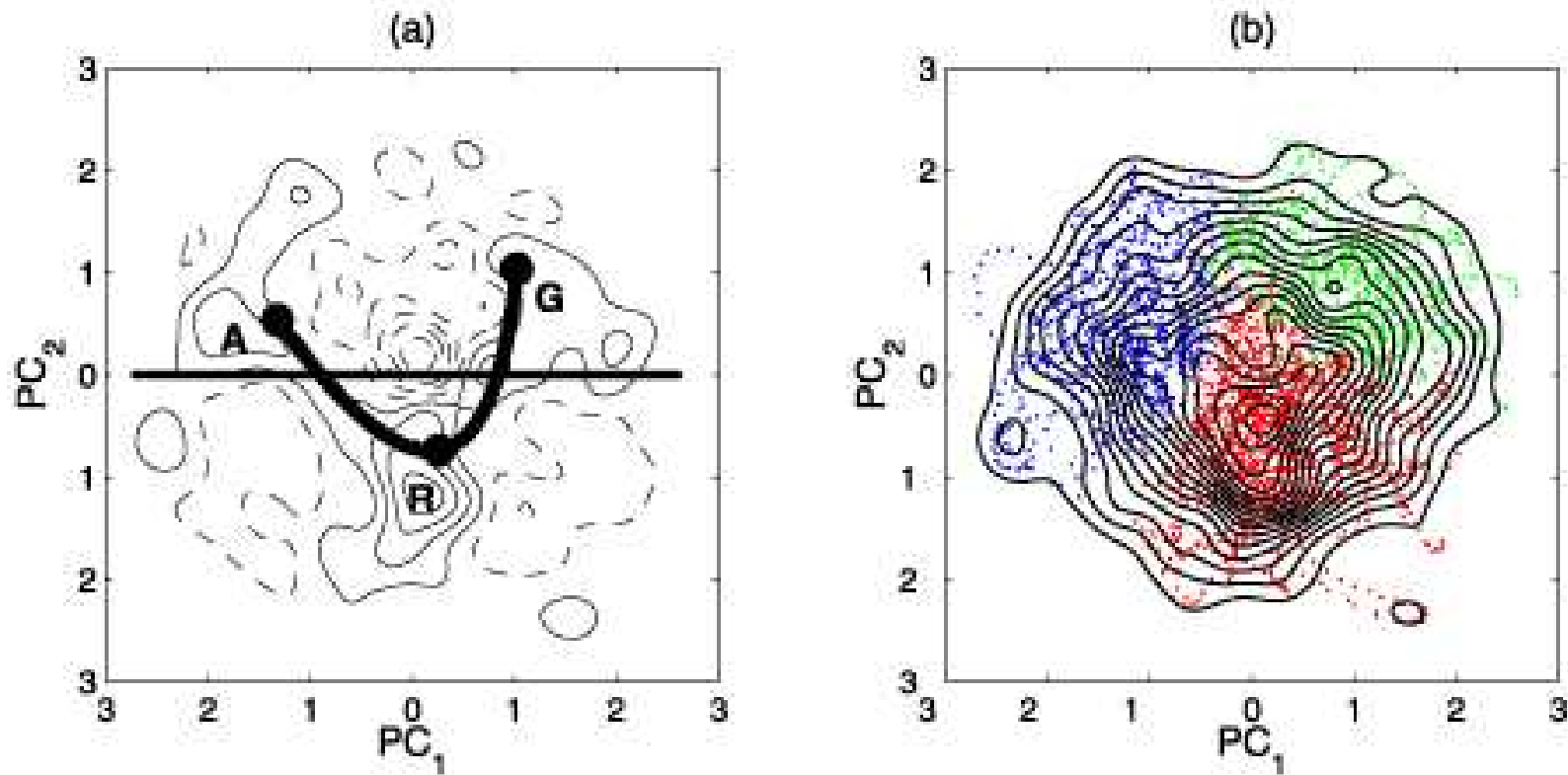


EOF₂



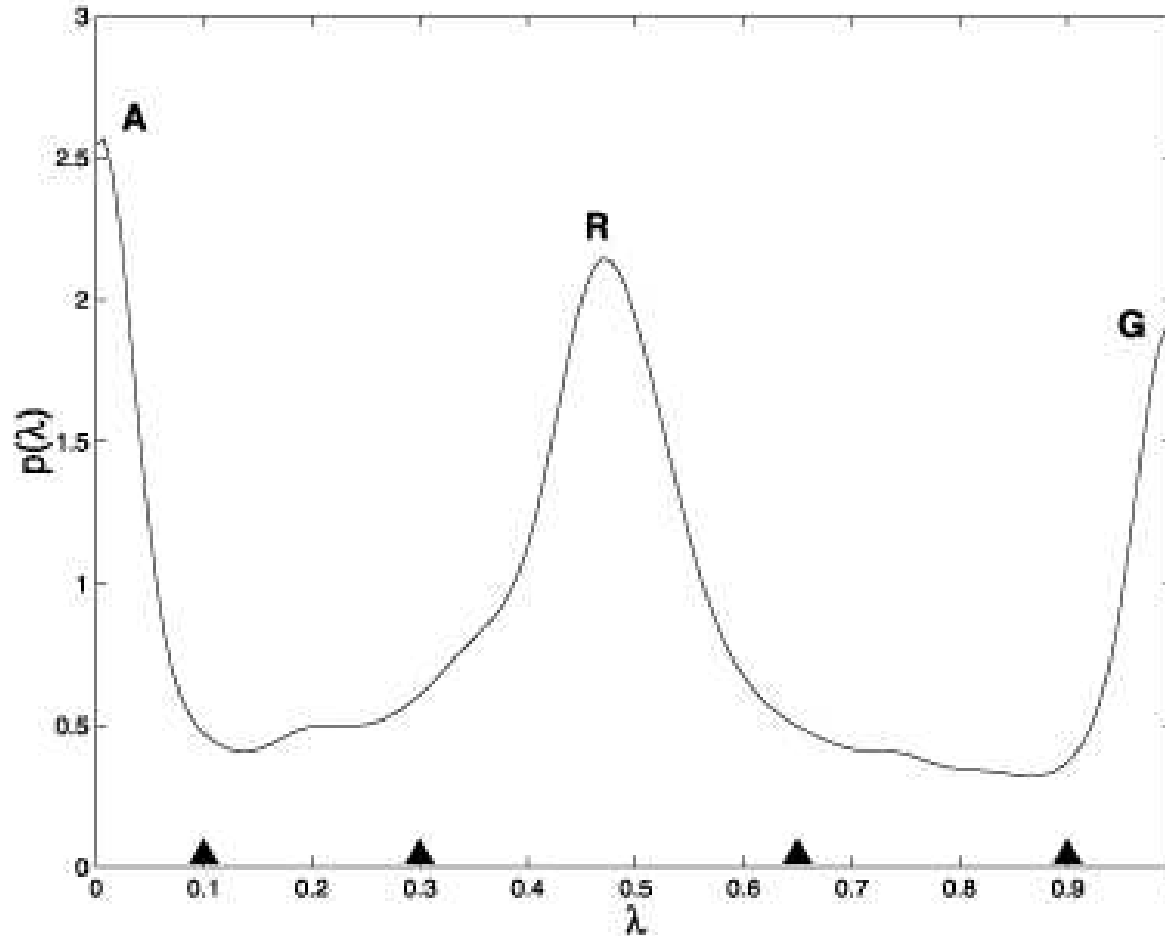
10-day lowpass-filtered 500 hPa geopotential height EOFs

Applications of NLPCA: NH Tropospheric LFV

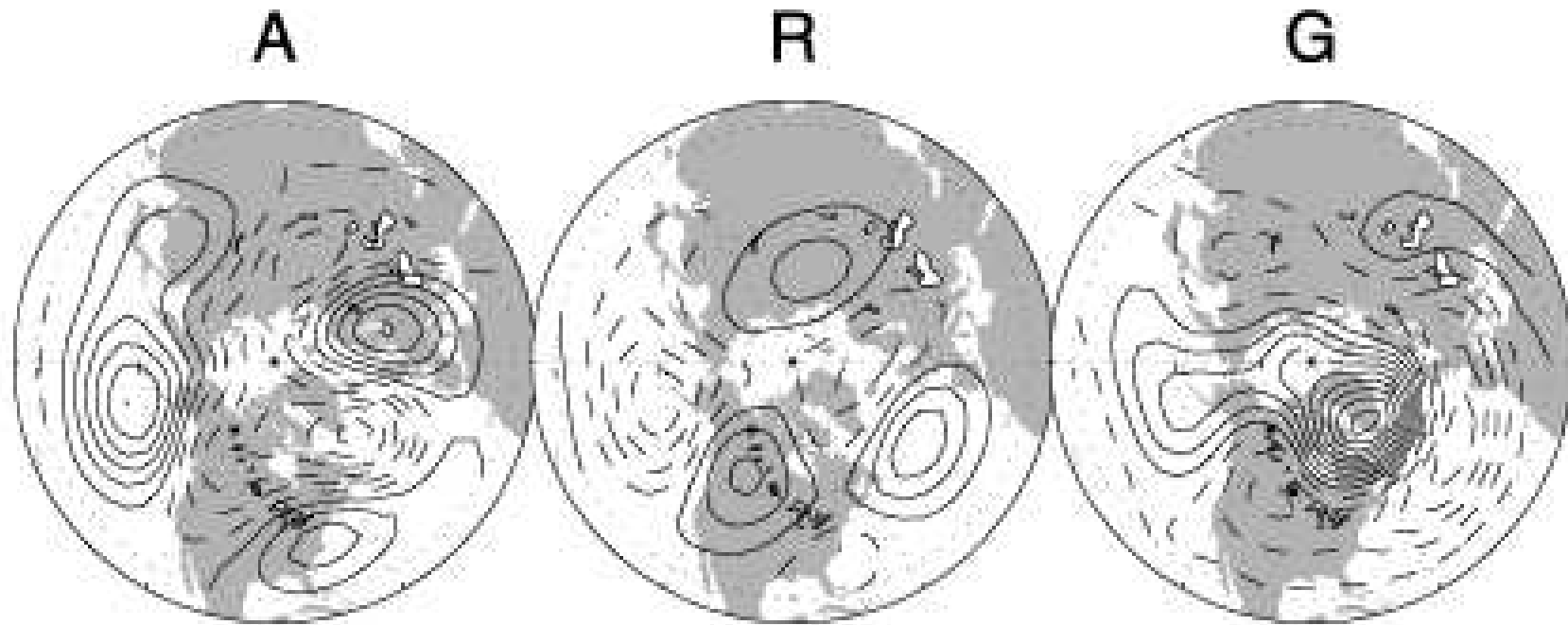


1D NLPCA Approximation: spatial structure
(PCA: 14.8%; NLPCA 18.4%)

Applications of NLPCA: NH Tropospheric LFV

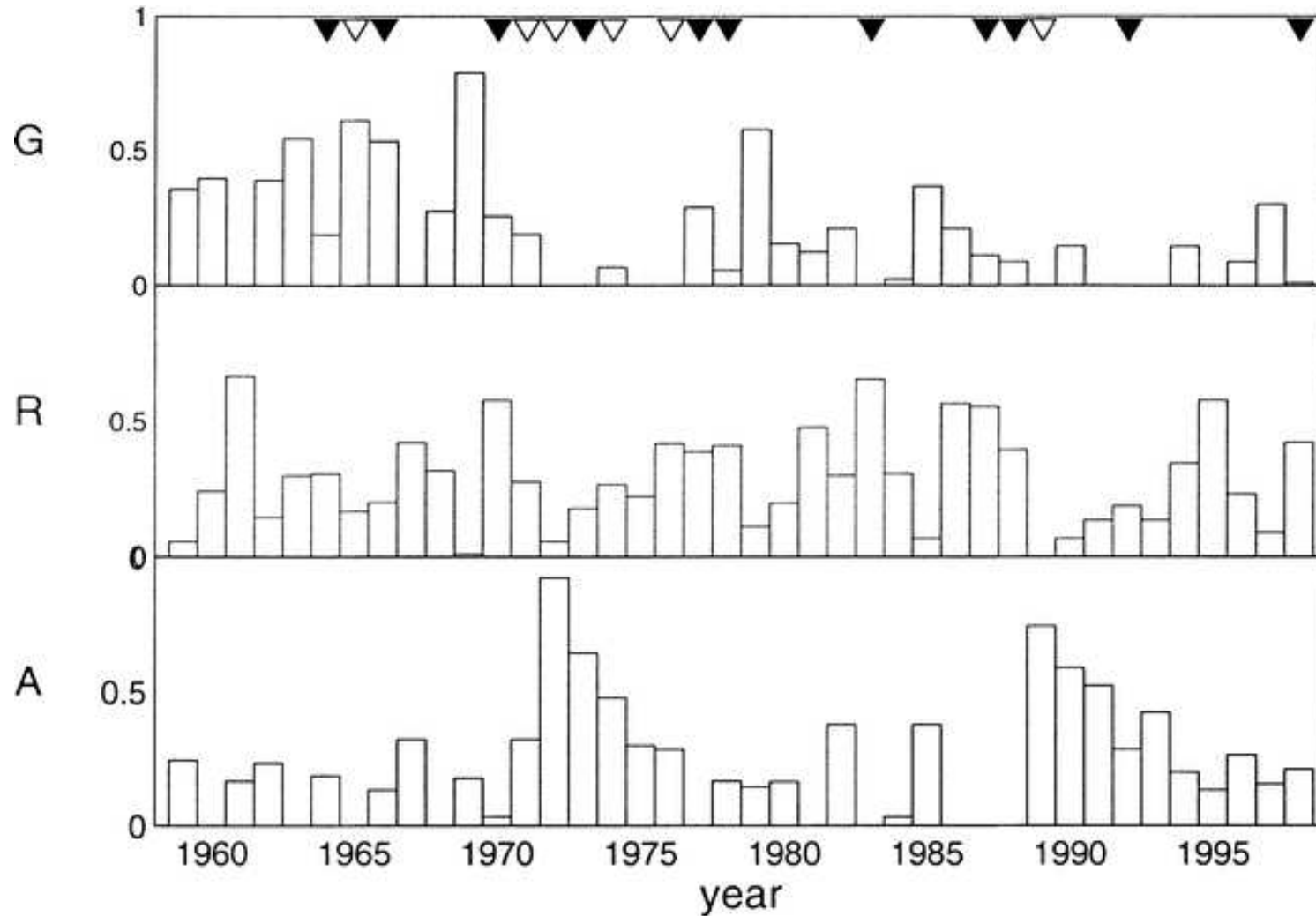


Applications of NLPCA: NH Tropospheric LFV



1D NLPCA Approximation: regime maps

Applications of NLPCA: NH Tropospheric LFV



UVic 1D NLPCA Approximation: interannual variability

NLPCA: Limitations and Drawbacks

- Parameter estimation in NLPCA (as in any nonlinear statistical model) must be done very carefully to ensure robust approximation

NLPCA: Limitations and Drawbacks

- Parameter estimation in NLPCA (as in any nonlinear statistical model) must be done very carefully to ensure robust approximation
- ⇒ analysis time-consuming, data hungry

NLPCA: Limitations and Drawbacks

- Parameter estimation in NLPCA (as in any nonlinear statistical model) must be done very carefully to ensure robust approximation
 - ⇒ analysis time-consuming, data hungry
 - ⇒ insufficiently careful analysis leads to spurious results (e.g. Christiansen, 2005)

NLPCA: Limitations and Drawbacks

- Parameter estimation in NLPCA (as in any nonlinear statistical model) must be done very carefully to ensure robust approximation
 - ⇒ analysis time-consuming, data hungry
 - ⇒ insufficiently careful analysis leads to spurious results (e.g. Christiansen, 2005)
- Theoretical underpinning of NLPCA is weak

NLPCA: Limitations and Drawbacks

- Parameter estimation in NLPCA (as in any nonlinear statistical model) must be done very carefully to ensure robust approximation
 - ⇒ analysis time-consuming, data hungry
 - ⇒ insufficiently careful analysis leads to spurious results (e.g. Christiansen, 2005)
- Theoretical underpinning of NLPCA is weak
 - ⇒ no “rigorous” theory of sampling variability

NLPCA: Limitations and Drawbacks

- Parameter estimation in NLPCA (as in any nonlinear statistical model) must be done very carefully to ensure robust approximation
 - ⇒ analysis time-consuming, data hungry
 - ⇒ insufficiently careful analysis leads to spurious results (e.g. Christiansen, 2005)
- Theoretical underpinning of NLPCA is weak
 - ⇒ no “rigorous” theory of sampling variability
 - Information theory may provide new tools with

NLPCA: Limitations and Drawbacks

- Parameter estimation in NLPCA (as in any nonlinear statistical model) must be done very carefully to ensure robust approximation
 - ⇒ analysis time-consuming, data hungry
 - ⇒ insufficiently careful analysis leads to spurious results (e.g. Christiansen, 2005)
- Theoretical underpinning of NLPCA is weak
 - ⇒ no “rigorous” theory of sampling variability
- Information theory may provide new tools with
 - better sampling properties

NLPCA: Limitations and Drawbacks

- Parameter estimation in NLPCA (as in any nonlinear statistical model) must be done very carefully to ensure robust approximation
 - ⇒ analysis time-consuming, data hungry
 - ⇒ insufficiently careful analysis leads to spurious results (e.g. Christiansen, 2005)
- Theoretical underpinning of NLPCA is weak
 - ⇒ no “rigorous” theory of sampling variability
- Information theory may provide new tools with
 - better sampling properties
 - better theoretical basis

Conclusions

- Traditional PCA optimal for dimensionality reduction only if data distribution falls along orthogonal axes

Conclusions

- Traditional PCA optimal for dimensionality reduction only if data distribution falls along orthogonal axes
- Can define nonlinear generalisation, NLPCA, which can robustly characterise nonlinear low-dimensional structure in datasets

Conclusions

- Traditional PCA optimal for dimensionality reduction only if data distribution falls along orthogonal axes
- Can define nonlinear generalisation, NLPCA, which can robustly characterise nonlinear low-dimensional structure in datasets
- NLPCA approximations can provide a fundamentally different characterisation of data than PCA approximations

Conclusions

- Traditional PCA optimal for dimensionality reduction only if data distribution falls along orthogonal axes
- Can define nonlinear generalisation, NLPCA, which can robustly characterise nonlinear low-dimensional structure in datasets
- NLPCA approximations can provide a fundamentally different characterisation of data than PCA approximations
- Implementation of NLPCA difficult and lacking in underlying theory; represents a first attempt at a big (and challenging) problem



Acknowledgements

- William Hsieh (UBC)
- Lionel Pandolfo (UBC)
- John Fyfe (CCCma)
- Qiaobin Teng (CCCma)
- Benyang Tang (JPL)

Parameter Estimation in NLPCA

- An ensemble approach was taken

Parameter Estimation in NLPCA

- An ensemble approach was taken
- For a large number N (~ 50) of trials:

Parameter Estimation in NLPCA

- An ensemble approach was taken
- For a large number N (~ 50) of trials:
 - data was randomly split into *training* and *validation* sets (taking autocorrelation into account)

Parameter Estimation in NLPCA

- An ensemble approach was taken
- For a large number N (~ 50) of trials:
 - data was randomly split into *training* and *validation* sets (taking autocorrelation into account)
 - a random initial parameter set was selected

Parameter Estimation in NLPCA

- An ensemble approach was taken
- For a large number N (~ 50) of trials:
 - data was randomly split into *training* and *validation* sets (taking autocorrelation into account)
 - a random initial parameter set was selected
- For each ensemble member, iterative minimisation procedure carried out until either:

Parameter Estimation in NLPCA

- An ensemble approach was taken
- For a large number N (~ 50) of trials:
 - data was randomly split into *training* and *validation* sets (taking autocorrelation into account)
 - a random initial parameter set was selected
- For each ensemble member, iterative minimisation procedure carried out until either:
 - error over training data stopped changing

Parameter Estimation in NLPKA

- An ensemble approach was taken
- For a large number N (~ 50) of trials:
 - data was randomly split into *training* and *validation* sets (taking autocorrelation into account)
 - a random initial parameter set was selected
- For each ensemble member, iterative minimisation procedure carried out until either:
 - error over training data stopped changing
 - error over validation data started increasing

Parameter Estimation in NLPCA

- An ensemble approach was taken
- For a large number N (~ 50) of trials:
 - data was randomly split into *training* and *validation* sets (taking autocorrelation into account)
 - a random initial parameter set was selected
- For each ensemble member, iterative minimisation procedure carried out until either:
 - error over training data stopped changing
 - error over validation data started increasing
- Method **does not** look for global error minimum

Parameter Estimation in NLPCA

- Ensemble member becomes **candidate model** if

$$\langle \|\epsilon\|^2 \rangle_{validation} \leq \langle \|\epsilon\|^2 \rangle_{training}$$

Parameter Estimation in NLPCA

- Ensemble member becomes **candidate model** if

$$\langle \|\epsilon\|^2 \rangle_{validation} \leq \langle \|\epsilon\|^2 \rangle_{training}$$

- Candidate models compared

Parameter Estimation in NLPKA

- Ensemble member becomes **candidate model** if

$$\langle \|\epsilon\|^2 \rangle_{validation} \leq \langle \|\epsilon\|^2 \rangle_{training}$$

- Candidate models compared
 - if they share same shape and orientation

Parameter Estimation in NLPKA

- Ensemble member becomes **candidate model** if

$$\langle \|\epsilon\|^2 \rangle_{validation} \leq \langle \|\epsilon\|^2 \rangle_{training}$$

- Candidate models compared
 - if they share same shape and orientation
- ⇒ approximation is robust

Parameter Estimation in NLPKA

- Ensemble member becomes **candidate model** if

$$\langle \|\epsilon\|^2 \rangle_{validation} \leq \langle \|\epsilon\|^2 \rangle_{training}$$

- Candidate models compared
 - if they share same shape and orientation

⇒ approximation is robust

 - if they differ in shape and orientation

Parameter Estimation in NLPKA

- Ensemble member becomes **candidate model** if

$$\langle \|\epsilon\|^2 \rangle_{validation} \leq \langle \|\epsilon\|^2 \rangle_{training}$$

- Candidate models compared
 - if they share same shape and orientation
 - ⇒ approximation is robust
 - if they differ in shape and orientation
 - ⇒ approximation is not robust

Parameter Estimation in NLPKA

- Ensemble member becomes **candidate model** if

$$\langle \|\epsilon\|^2 \rangle_{validation} \leq \langle \|\epsilon\|^2 \rangle_{training}$$

- Candidate models compared
 - if they share same shape and orientation
 - ⇒ approximation is robust
 - if they differ in shape and orientation
 - ⇒ approximation is not robust
- If approximation not robust, model simplified & procedure repeated until robust model found

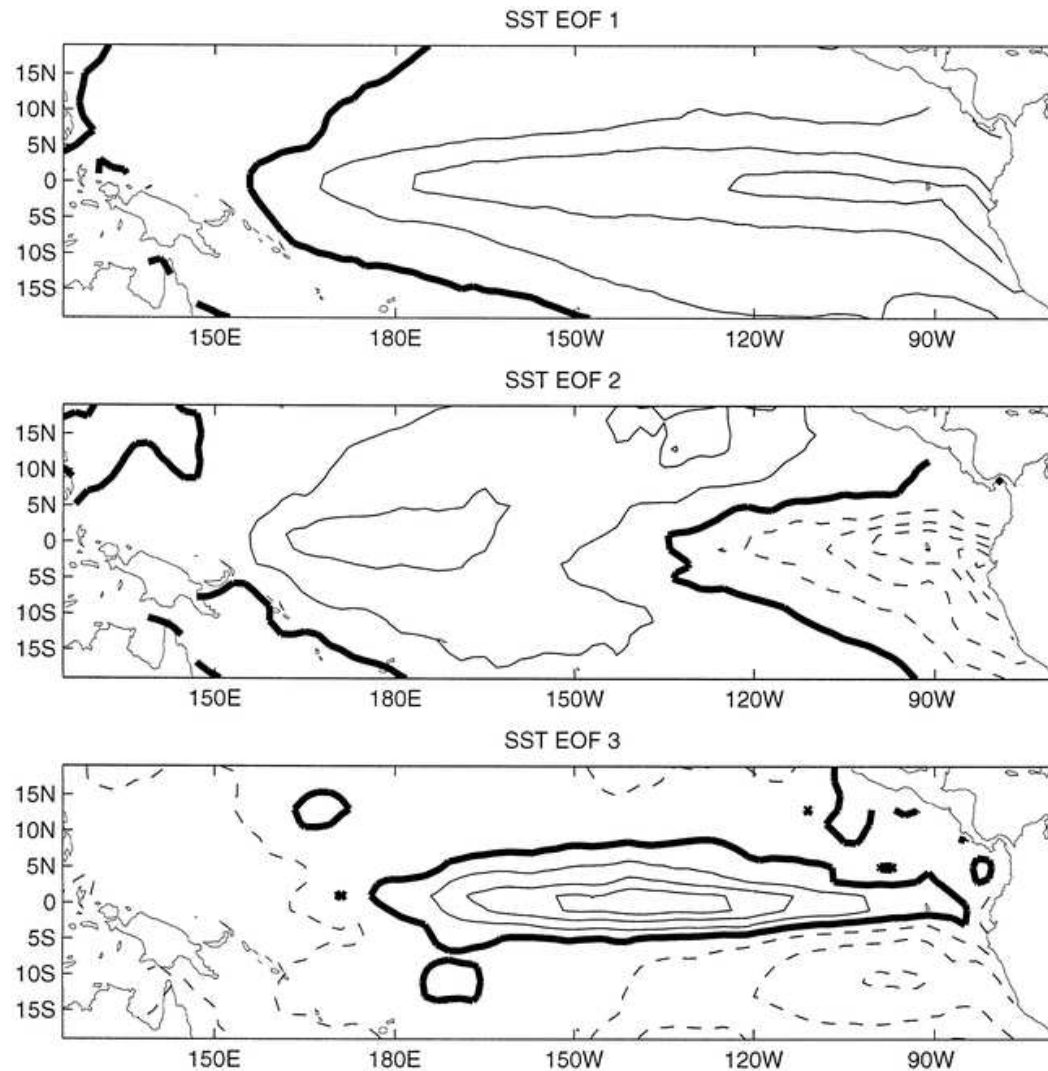
Parameter Estimation in NLPCA

- Procedure will ultimately yield PCA solution if no robust non-Gaussian structure present

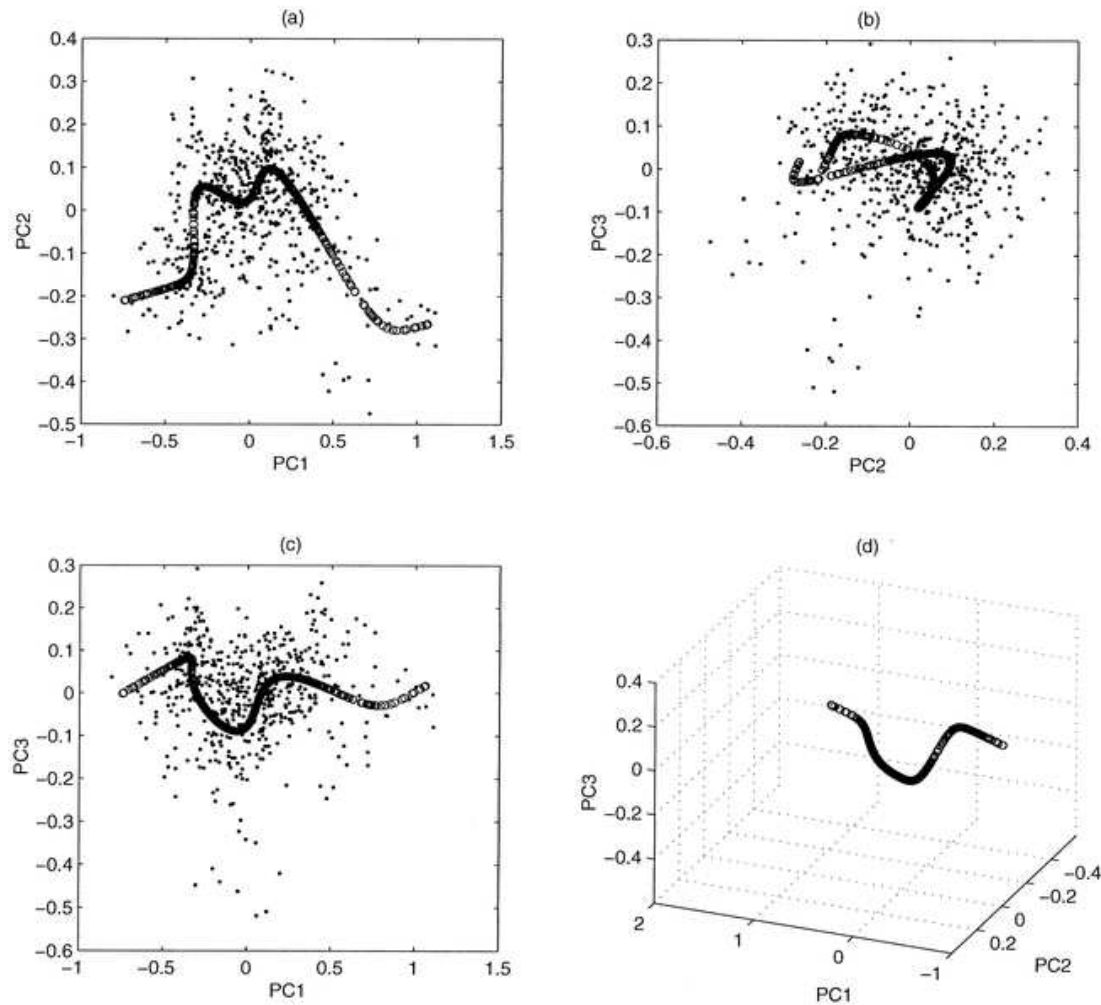
Parameter Estimation in NLPCA

- Procedure will ultimately yield PCA solution if no robust non-Gaussian structure present
- Such a careful procedure necessary to avoid finding spurious non-Gaussian structure

Applications of NLPCA: Tropical Pacific SST

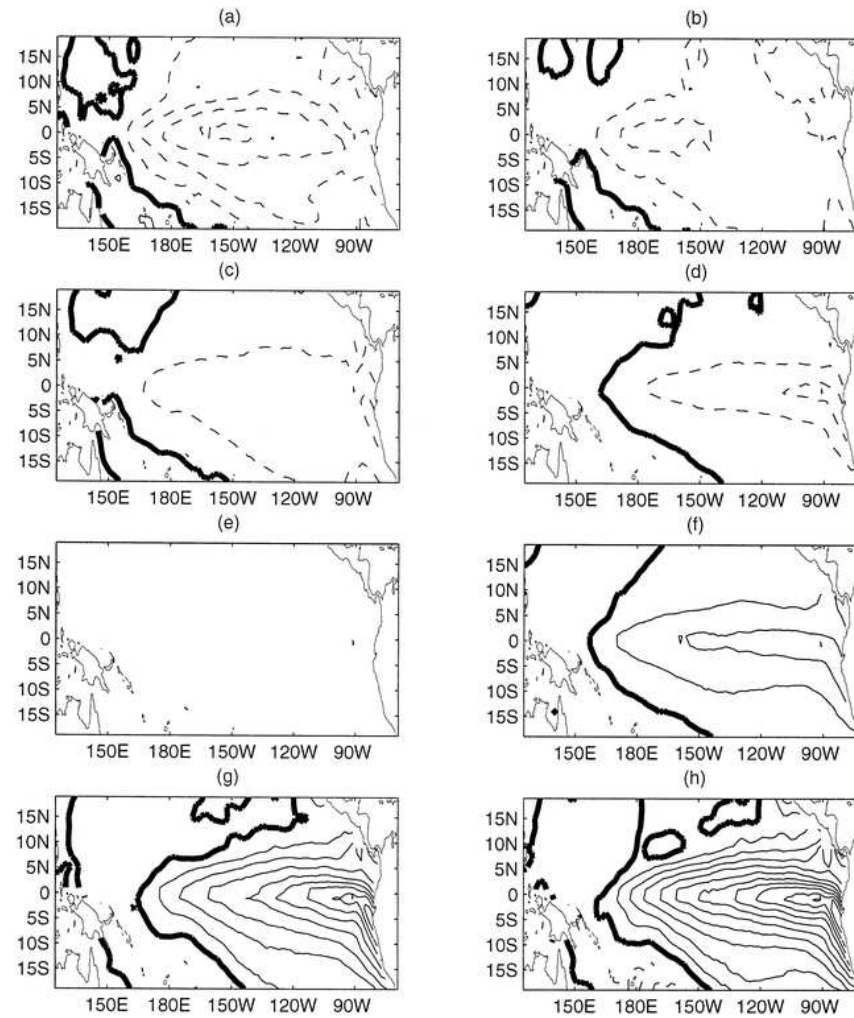


Applications of NLPCA: Tropical Pacific SST



1D NLPCA Approximation: spatial structure

Applications of NLPCA: Tropical Pacific SST



1D NLPCA Approximation: spatial structure