

Section 6

Principal Component and Maximum
Covariance Analyses

Maximum Covariance Analysis (MCA)

Purpose

- **Purpose**
 - Find patterns in 2 datasets which are highly correlated (i.e. are frequently met simultaneously). E.g. patterns of SST that go along with patterns of SLP.
- **Applications**
 - Study the coupling between parameters to understand physical mechanisms of climate variations. E.g. how do SST and SLP mutually influence each other?
 - Statistical downscaling. (Translate GCM derived climate change scenarios to the local/regional scale.)
 - Reconstruction / Forecasting

Paired Multivariate Datasets

- **Two paired multivariate datasets**

$$\mathbf{x}(i) = (x_1(i), \dots, x_n(i))$$

$$\mathbf{y}(i) = (y_1(i), \dots, y_m(i))$$

- $i=1, \dots, N$, joint observations of \mathbf{x} , \mathbf{y}
- Two data clouds in two phase spaces
- Centered (means subtracted!), (just for more simple notation!)
- E.g. \mathbf{x} : anomalies of SLP Europe/Atlantic (n gridpoints),
 \mathbf{y} : Temperature anomalies in CH (m stations)
 all Januaries of the 20th century ($N=100$)

- **Data Matrices**

$$\mathbf{X} = \begin{bmatrix} x_1(1) & \cdots & x_n(1) \\ \vdots & \ddots & \vdots \\ x_1(N) & \cdots & x_n(N) \end{bmatrix} \quad \mathbf{Y} = \dots$$

Cross-Covariance



Call me
cross-covariofant

- **Cross-covariance matrix**

$$\mathbf{S}_{xy} := \begin{bmatrix} \text{cov}(x_1, y_1) & \text{cov}(x_1, y_j) & \text{cov}(x_1, y_m) \\ \text{cov}(x_i, y_1) & \text{cov}(x_i, y_j) & \vdots \\ \text{cov}(x_n, y_1) & \dots & \text{cov}(x_n, y_m) \end{bmatrix} = \frac{1}{(N-1)} \mathbf{X}^T \cdot \mathbf{Y}$$

- Univariate cross-covariances between all pairs of components
- $n \times m$ matrix
- In general not square, not symmetric
- diagonal elements do not have special meaning

Cross-Correlation

- **Cross-Correlation matrix**

$$\mathbf{C}_{xy} := \begin{bmatrix} \text{cor}(x_1, y_1) & \text{cor}(x_1, y_j) & \text{cor}(x_1, y_m) \\ \text{cor}(x_i, y_1) & \text{cor}(x_i, y_j) & \vdots \\ \text{cor}(x_n, y_1) & \dots & \text{cor}(x_n, y_m) \end{bmatrix}$$

- Analogous to \mathbf{S}_{xy} but divided by standard deviations
- All matrix elements in $\{-1, +1\}$

$$\mathbf{C}_{xy} = \frac{1}{N-1} \mathbf{D}_x^{-\frac{1}{2}} \cdot \mathbf{X}^T \cdot \mathbf{Y} \cdot \mathbf{D}_y^{-\frac{1}{2}} = \mathbf{D}_x^{-\frac{1}{2}} \cdot \mathbf{S}_{xy} \cdot \mathbf{D}_y^{-\frac{1}{2}}$$

$$\mathbf{D}_x^{-\frac{1}{2}} = \begin{bmatrix} 1/\sigma_{x_k} \end{bmatrix} \quad \mathbf{D}_y^{-\frac{1}{2}} = \begin{bmatrix} 1/\sigma_{y_k} \end{bmatrix}$$

MCA Mathematical Procedure

- **Singular Value Decomposition (SVD)**

>> Appendix A

$\mathbf{S}_{xy} = \text{cov}(\mathbf{X}, \mathbf{Y})$ the cross-covariance matrix ($n \times m$)

- There are r real numbrs $\{\omega_1, \omega_2, \dots, \omega_r\}$, $\omega_k > 0$, *singular values* and r vectors $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}$, n -dim., unit-length, orthogonal, and r vectors $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r\}$, m -dim., unit-length, orthogonal, called *left* (\mathbf{u}_k) and *right* (\mathbf{v}_k) *singular vectors*, such that:

$$\mathbf{S}_{xy} = \mathbf{U}^T \cdot \mathbf{\Omega} \cdot \mathbf{V} \quad \mathbf{\Omega} = [\omega_k], \text{ diagonal, } r \times r$$

$$\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r], \quad \mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r], \quad \text{vectors in columns}$$

R: svd

MCA Interpretation

- **Singular Vectors, Coefficients**

- Consider singular vectors anchored in center of data clouds, spanning sub-spaces of the phasespaces of each dataset (new variables). Not necessarily a basis system ($r \leq n, r \leq m$!)
- Projections of data $\mathbf{x}(i)$ onto singular vectors are new data coordinates:

$$a_j(i) = \mathbf{x}(i)^T \cdot \mathbf{u}_j, \quad b_j(i) = \underbrace{\mathbf{y}(i)^T \cdot \mathbf{v}_j}_{\text{projection of data vector on singular vector}}, \quad j = 1, \dots, r$$

projection of data vector on singular vector

- New coordinates are linear combinations of original variables.
- $a_j(i)$: left coefficients, (also left SVD scores)
 $b_j(i)$: right coefficients, (also right SVD scores)

MCA Interpretation

- **Singular values, cross-covariance**

- Ω is the cross-covariance matrix of the new coordinates $\{a_k\}$, $\{b_k\}$:

>> Appendix B

$$\text{cov}(a_i(.), b_j(.)) = 0 \quad \text{for } i \neq j, \quad \text{cov}(a_i(.), b_i(.)) = \omega_i$$

- Coordinates corresponding to different indices of singular vectors are mutually uncorrelated.
- The first pair of singular vectors $\{\mathbf{u}_1, \mathbf{v}_1\}$ are the phase-space directions, for which the projections have the largest possible cross-covariance. *First coupled mode*.

Subsequent vector pairs $\{\mathbf{u}_k, \mathbf{v}_k\}$ maximise cross-covariance subject to orthogonality on previous pairs. *kth coupled mode*.

MCA Interpretation

- **Within space variance**

- Singular vectors do not maximize variance in individual spaces.
- Singular vectors are not necessarily aligned along directions of large data spread or cloud symmetry.
- Left and right coefficients are in general not uncorrelated between themselves:

$$\text{var}(a_i, a_j) \neq 0 \quad \text{var}(b_i, b_j) \neq 0 \quad \text{for } i \neq j$$

- *Cumulative Explained Variance Fraction* of first l modes

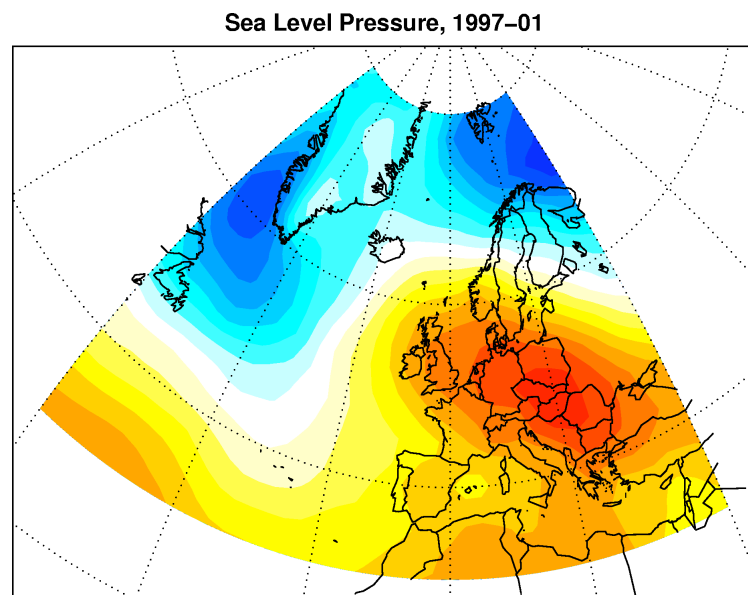
$$CEVF_x^l = \sum_k^l \text{var}(a_k) / \text{tr}(\mathbf{S}_{xx}) \quad CEVF_y^l = \sum_k^l \text{var}(b_k) / \text{tr}(\mathbf{S}_{yy})$$

MCA Interpretation

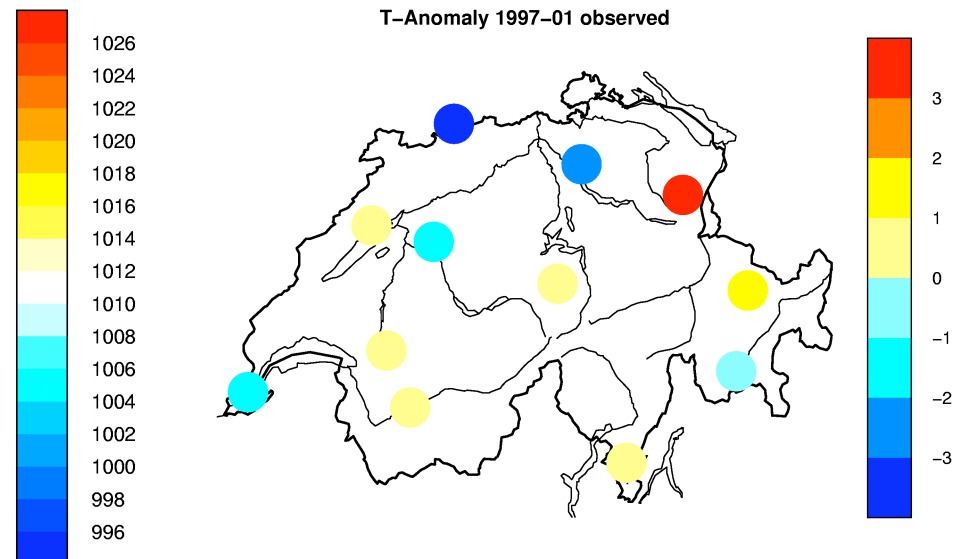
- Singular vectors describe patterns of anomalies in each dataset that tend to occur simultaneously (are linearly correlated). *Modes of co-variability, coupled modes, canonical pairs.*
- Coefficients represent amplitude (emphasis) of the respective patterns in each sample.
- The first few coefficient pairs have large co-variance and often show high correlations. *Dominant modes.*

Example: Swiss T <=> SLP

How does sea level pressure influence winter temperature-distributions in Switzerland?



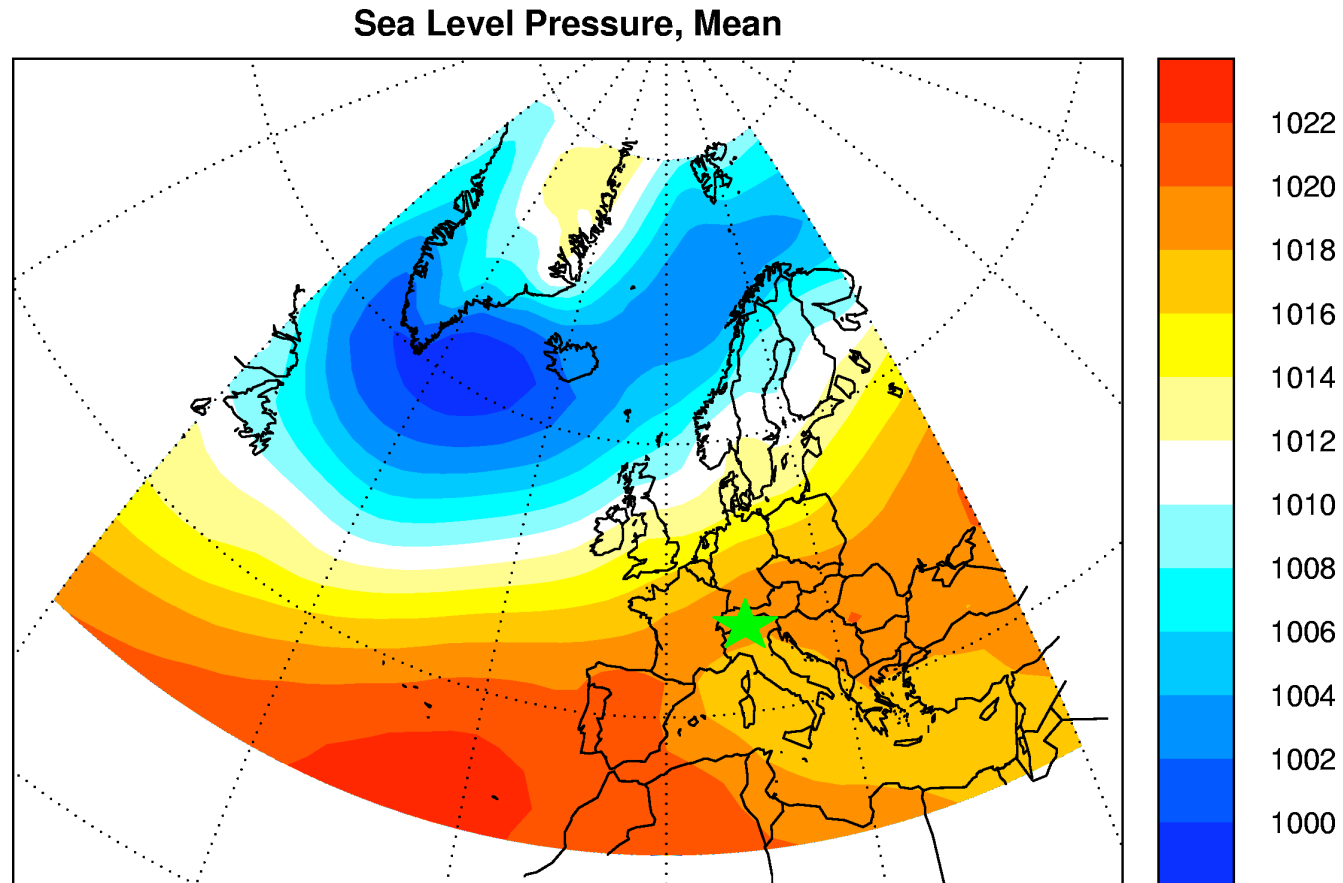
SLP: DJF, 1957 –1994
(60°W–40°E, 30°N–80°N)
41 x 21 grid points, 2.5 degrees
ECMWF, ERA40



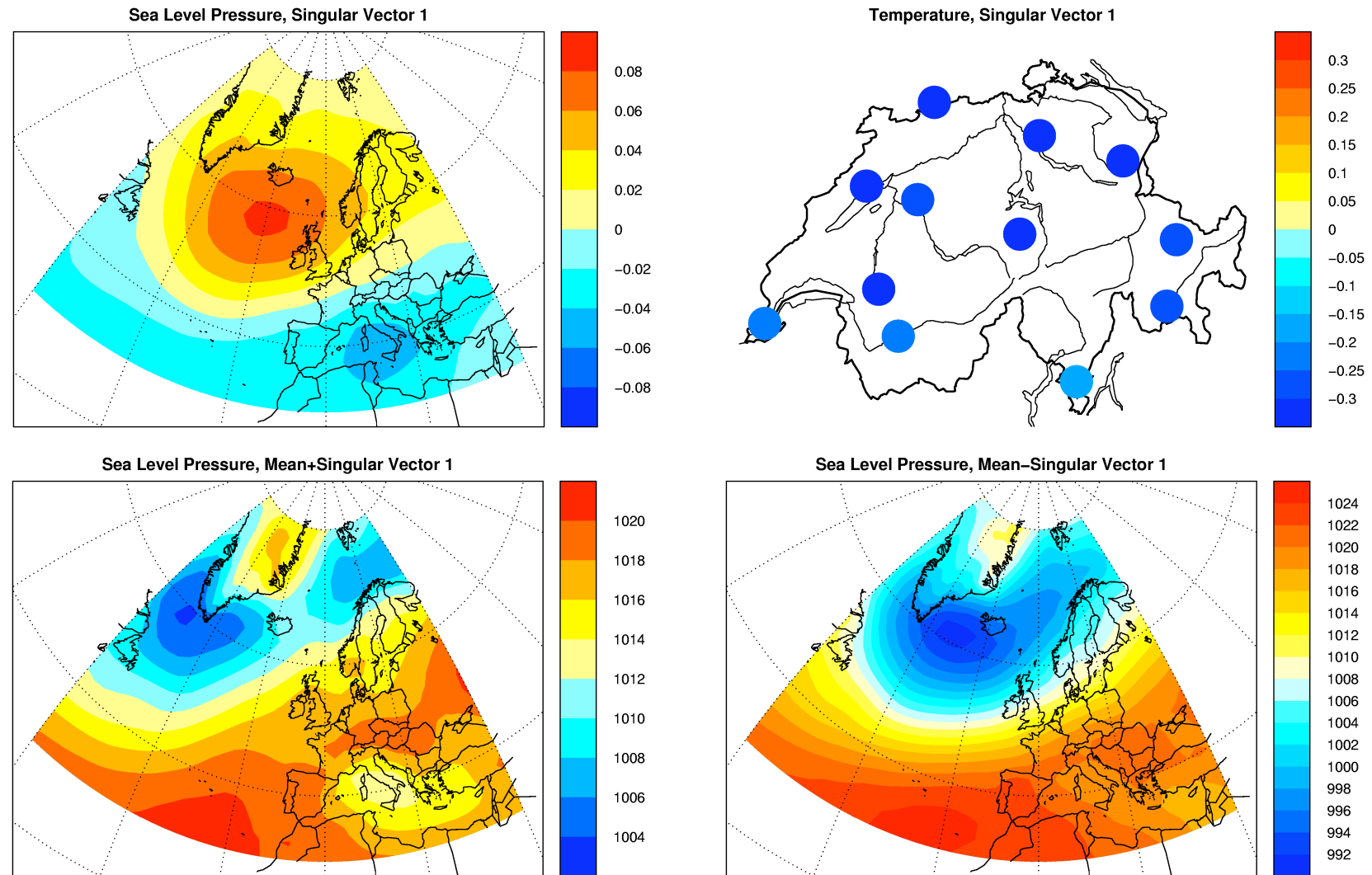
T: DJF, 1957 –1994
12 stations
MeteoSwiss

Begert et al. 2005, Simmons&Gibbson 2000

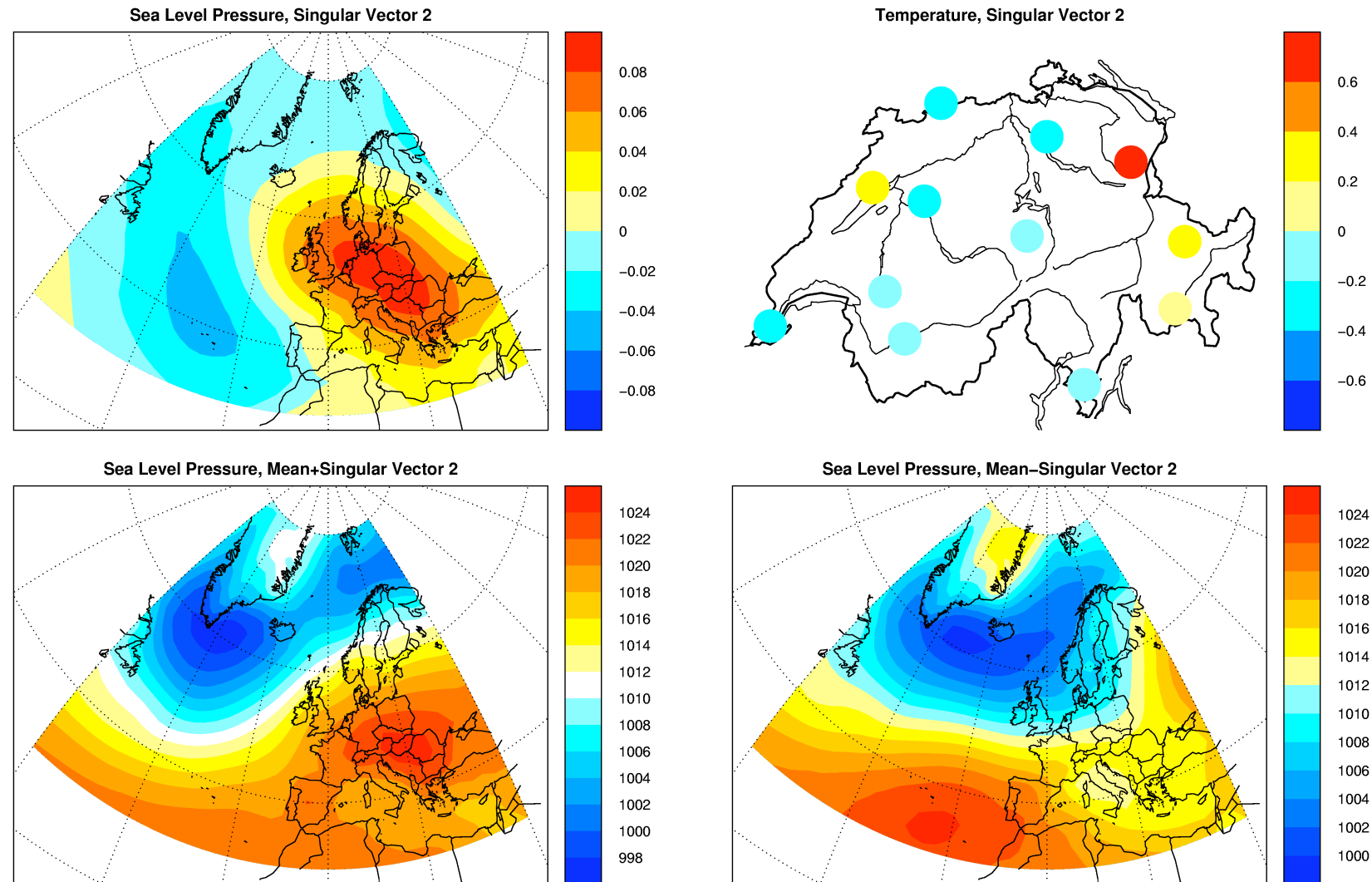
Example: Mean winter SLP



Example: Singular Vector Pair 1

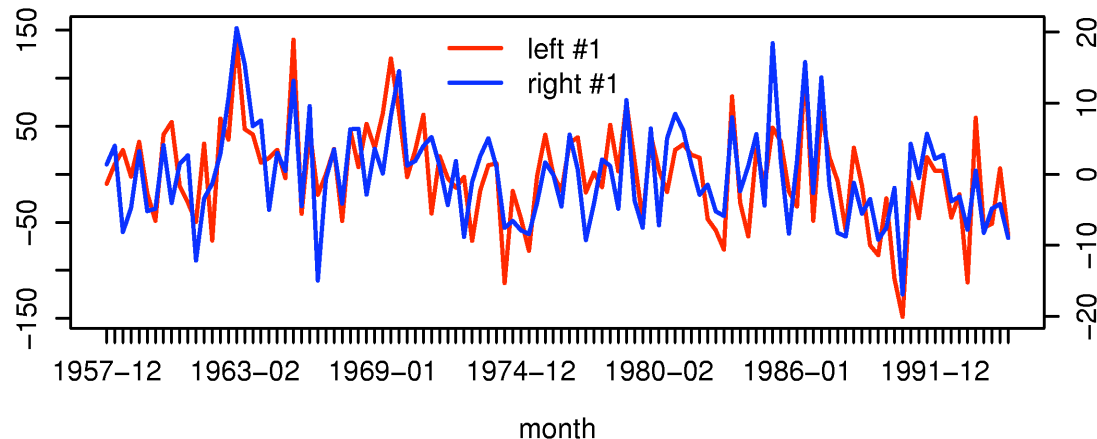


Example: Singular Vector Pair 2

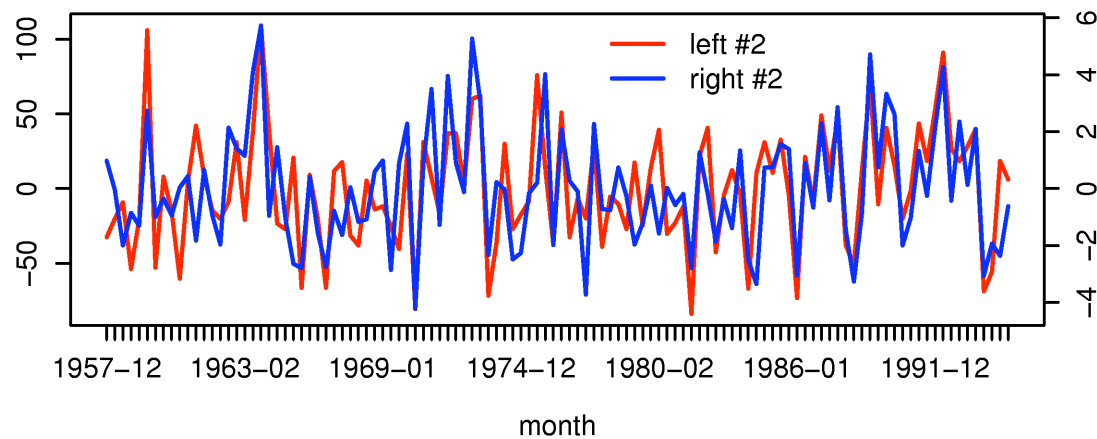


Example: Coefficients

Coefficients 1st mode



Coefficients 2nd mode



Mode	Correlation
1	0.70
2	0.67
3	0.60
4	0.39
5	0.28
6	0.37
7	0.29
8	0.22
9	0.33
10	0.20
11	0.34
12	0.34

Measuring Cross-Covariance

- ***Total squared covariance:***

$$\|S_{xy}\|_F := \sum_i^n \sum_j^m s_{ij}^2 = \sum_i^r \omega_i^2$$

Frobenius Norm

- ***Squared covariance fraction:***

- of singular vect. pair k :

$$SCF_k = \omega_k^2 / \sum_i^r \omega_i^2$$

- Note *squared* quantities compared to PCA!

Truncation

- Retain only first l coupled modes
- Projection onto first l modes
 - yields an approximation of original data

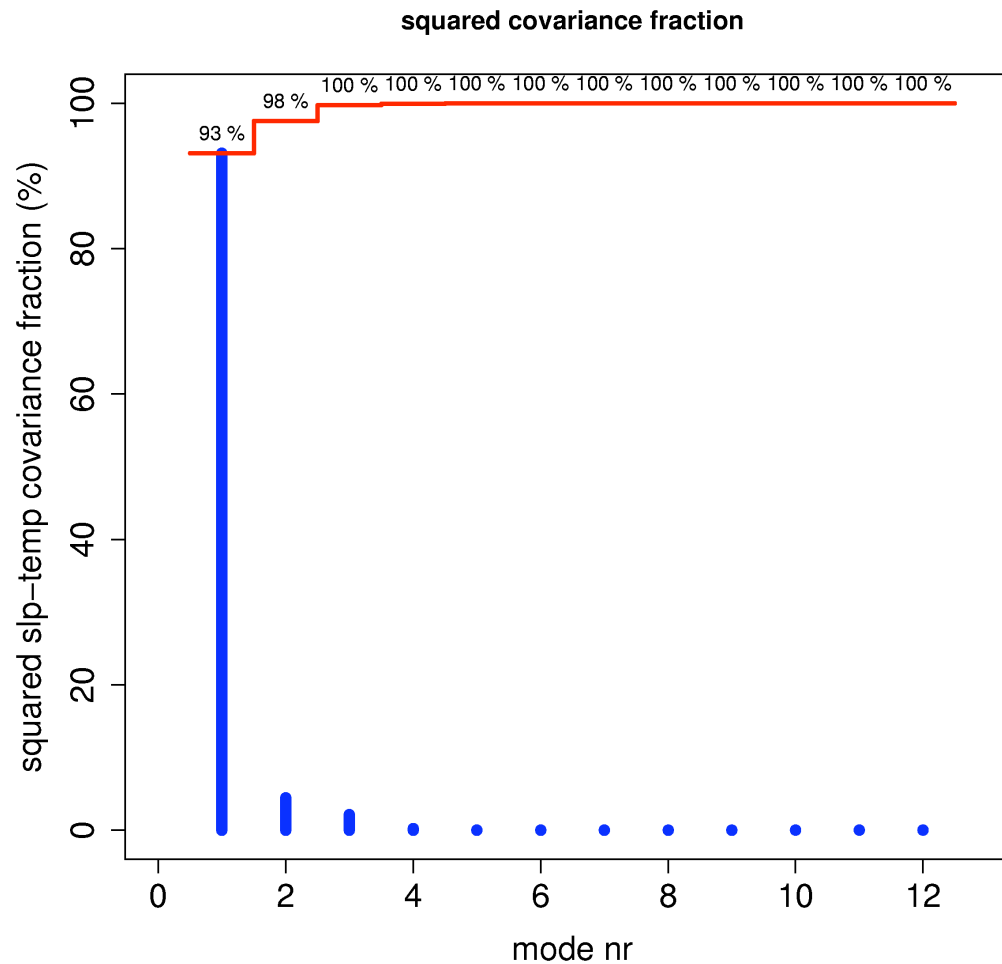
$${}_l \tilde{\mathbf{x}}(i) = \sum_k^l a_k(i) \cdot \mathbf{u}_k \quad {}_l \tilde{\mathbf{y}}(i) = \sum_k^l b_k(i) \cdot \mathbf{v}_k \quad l \leq r$$

- ***Cumulative squared covariance fraction:***

Residual cross-covariance,
i.e. not reproduced by approximation

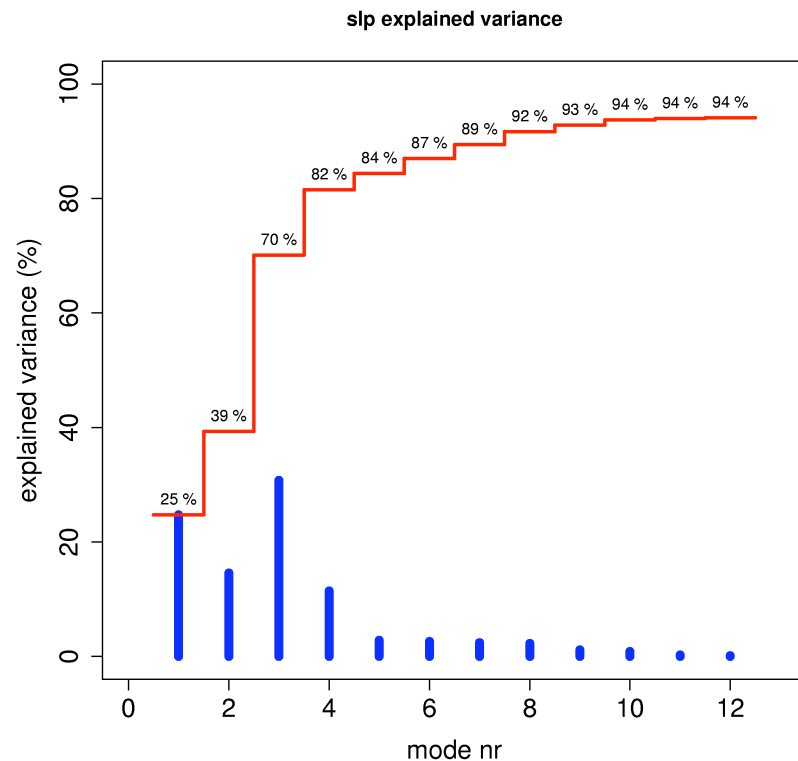
$$CSCF_l = 1 - \frac{\overbrace{\|\mathbf{S}_{xy} - \tilde{\mathbf{S}}_{xy}^l\|_F}}{\|\mathbf{S}_{xy}\|_F} = \sum_k^l \omega_k^2 / \sum_k^r \omega_k^2 = \sum_k^l SCF_k$$

Example: Squared Covariance Fraction

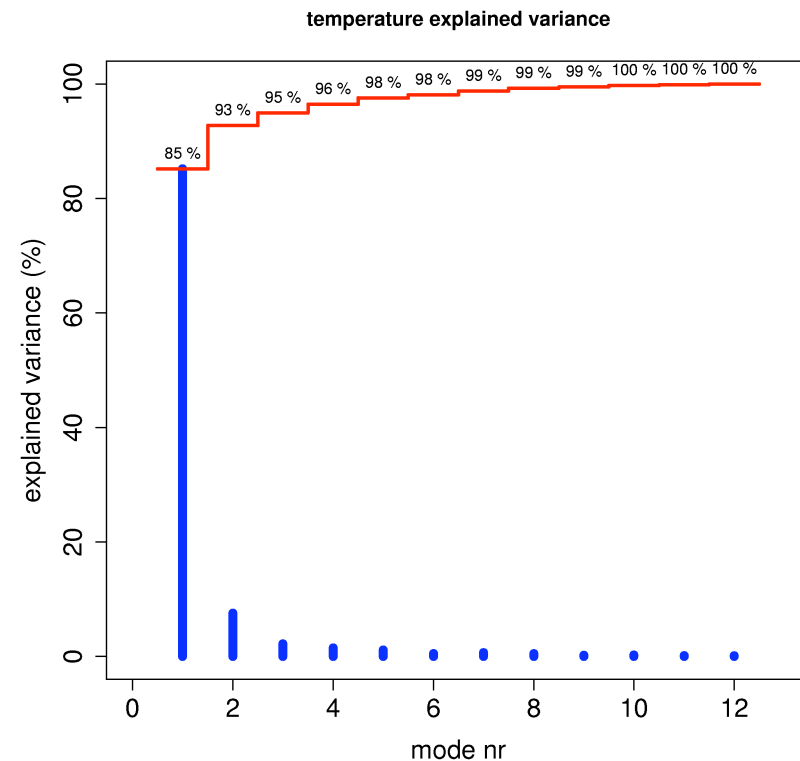


Example: Explained Variance

Sea Level Pressure



Temperature



Reconstruction / Prediction

- **Purpose**
 - Exploit cross-covariance to reconstruct/predict a right data from a left data (or vice versa).
- **Linear model between left and right coefficients**

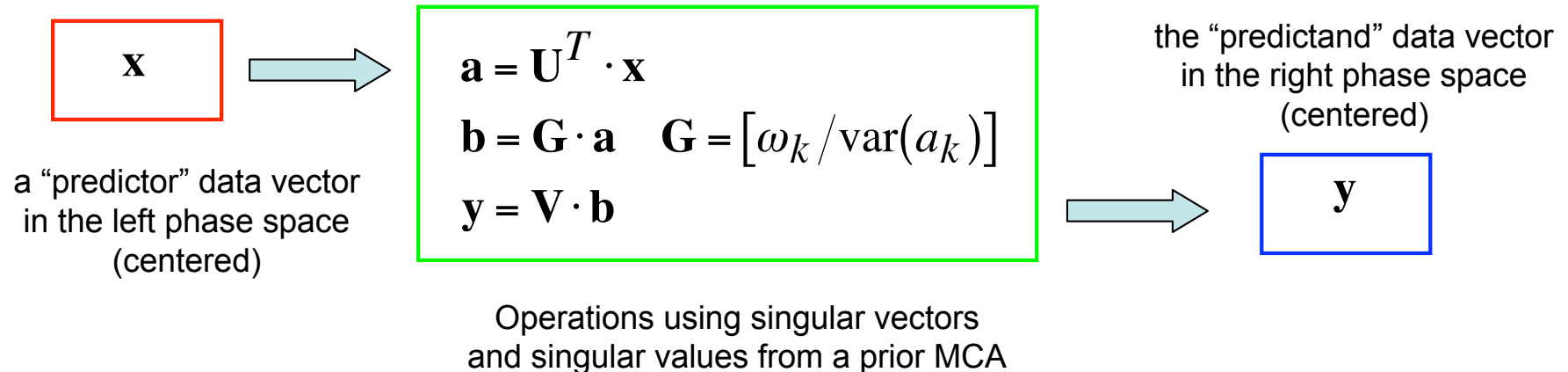
$$b_k(i) = \beta_{k0} + \sum_l^r \beta_{kl} \cdot a_l(i) + \varepsilon_k(i)$$

Annotations:
- β_{k0} : right coefficients
- β_{kl} : left coefficients (predictors)
- $\varepsilon_k(i)$: random noise
- β_{kl} : regression coefficients

- Simplifies to:

$$b_k(i) = \beta_{kk} \cdot a_k(i), \quad \beta_{kk} = \frac{\omega_k}{\text{var}(a_k)} \quad \begin{array}{l} a_k, b_l \text{ centered} \Rightarrow \beta_{k0} = 0 \\ \text{var}(a_k, b_l) = 0 \Rightarrow \beta_{kl} = 0, \quad k \neq l \end{array}$$

Reconstruction / Prediction



- **MCA prediction / reconstruction equation:**

$$\mathbf{y} = \mathbf{V} \cdot \mathbf{G} \cdot \mathbf{U}^T \cdot \mathbf{x}$$

- Possibly only using a few leading coupled modes. I.e. all matrices truncated to the number of desired modes.

Reconstruction / Prediction

- **Accuracy of reconstruction depends on**
 - The degree of covariance between the two fields.
 - The cumulative squared covariance fraction represented by the leading modes. The more modes the better the reconstruction.
 - The variance explained by the right singular vectors in the right phase space. I.e. the component of the right space that is related to the left space.

Example: Reconstruction

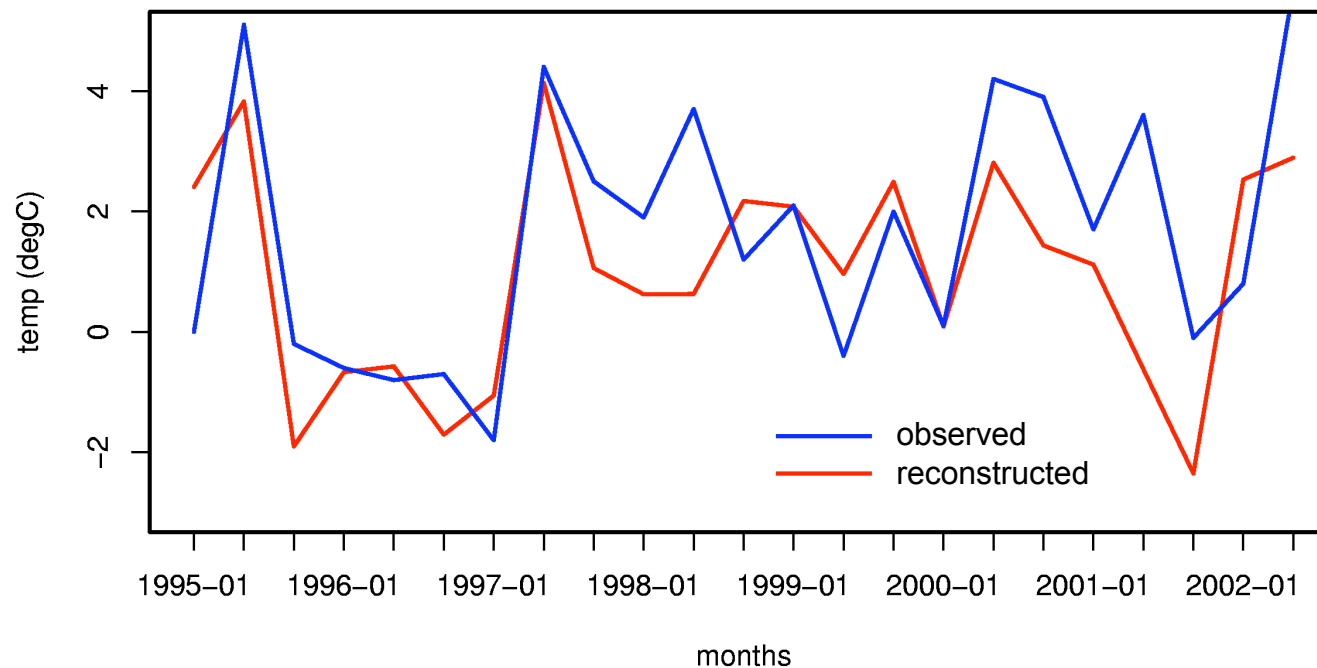
Reconstruction for 1995-2002 (winter months) using:

SLP (left field) as predictor

MCA calibrated for 1957-1994

2 leading coupled modes

Reconstructed (red): Zurich-MeteoSchweiz

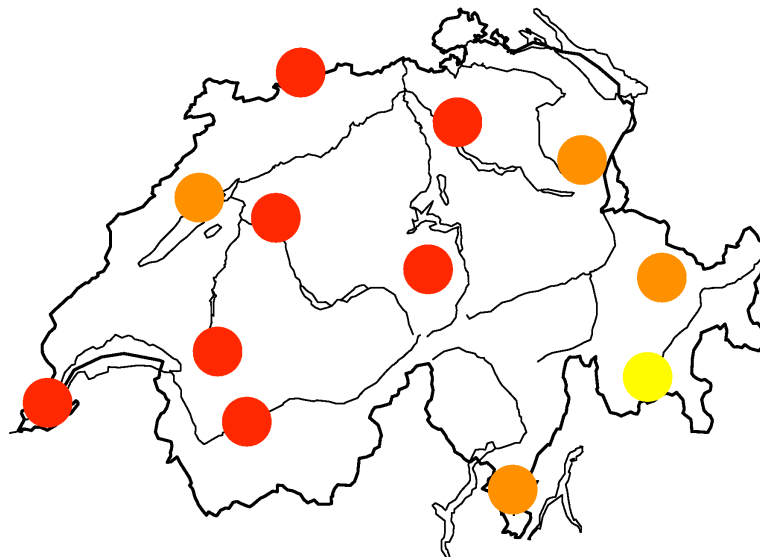


Example: Reconstruction

February 1997

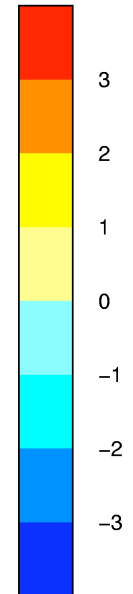
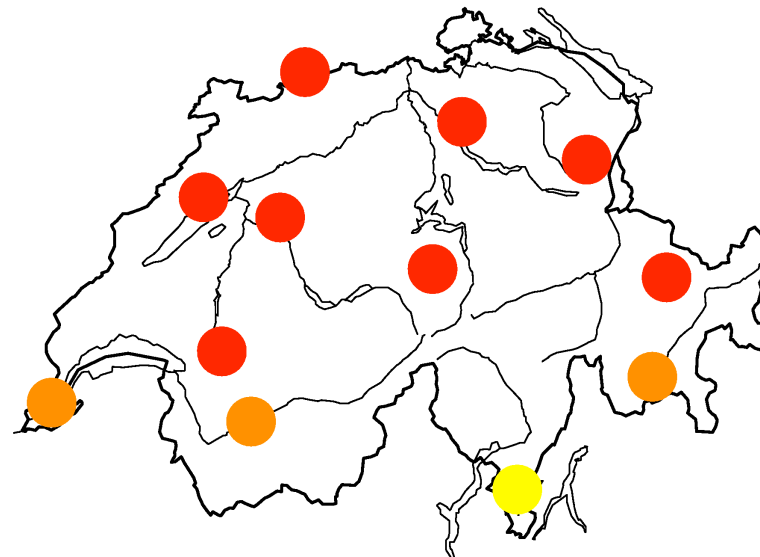
Observed

T-Anomaly 1997-02 observed



Reconstructed

T-Anomaly 1997-02 reconstructed from 2 MCA modes

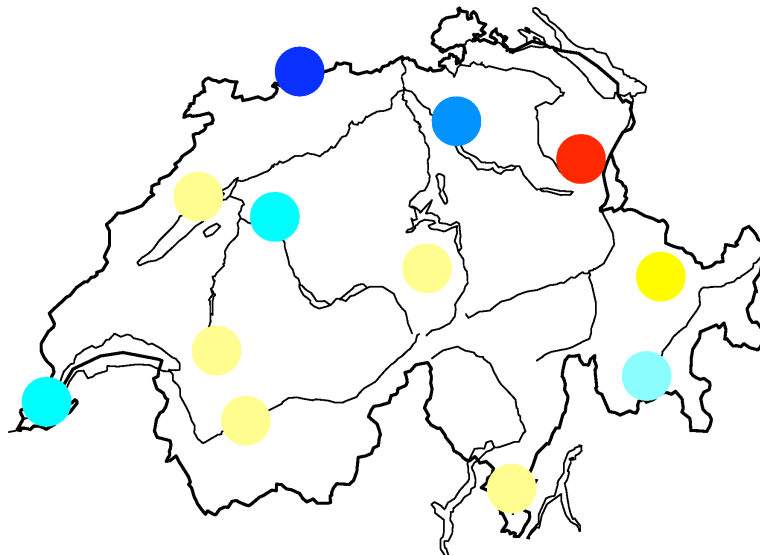


Example: Reconstruction

January 1997

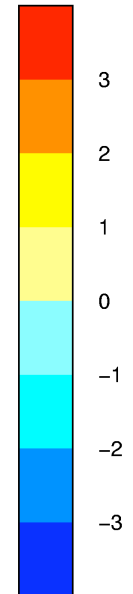
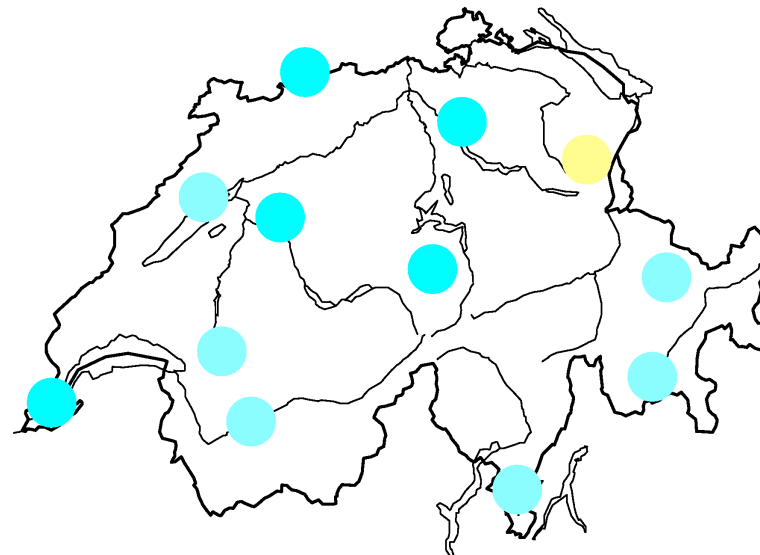
Observed

T-Anomaly 1997-01 observed



Reconstructed

T-Anomaly 1997-01 reconstructed from 2 MCA modes



Canonical Correlation Analysis

- **Similar to MCA but searching for max. correlation**
- **Procedure**
 - Conduct SVD with cross-correlation matrix
- **Interpretation (like for MCA)**
 - Canonical pairs describe coupled modes, emphasis on correlation
- **CCA or MCA?**
 - MCA focuses on *covariance*: Modes tend to be large where the variance is large. Danger that physical modes are confounded by large variance.
 - CCA focuses on *correlation*: Coupling is identified also if associated variance is low. Danger that physical modes are confounded by small (insignificant) variations (sampling problems).

(see Wilks 2005, Chap 12 for details)

Summary

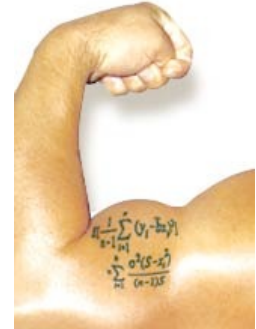
- **PCA & MCA are interesting techniques for characterising the (co-)variability in spatial datasets.**
- **Results require careful interpretation, corroboration in sensitivity experiments.**
- **Well established instruments in climate science:**
 - Technique of data reduction
 - Parsimonious reconstruction models (e.g. in historical climatology, paleo climatology).
 - Evaluation of physical mechanisms in climate models.
 - Statistical (empirical) forecasting, climate change downscaling.
 - Source of hypothesis building for later modelling exercises.

Section 6

Principal Component and Maximum
Covariance Analyses

Appendix MCA

Appendix A



- **Singular Value Decomposition (SVD)**

Q a real-valued $n \times m$ matrix (e.g. a cross-covariance matrix)
 $r = \text{rank}(\mathbf{Q}) \leq \min(n, m)$

► There exist real-valued matrices \mathbf{U} , \mathbf{V} , $\mathbf{\Omega}$ such that:

$$\mathbf{Q} = \mathbf{U} \cdot \mathbf{\Omega} \cdot \mathbf{V}^T$$

$\mathbf{\Omega}$: a diagonal $r \times r$ matrix:

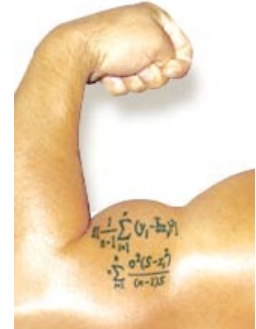
$$\omega_1 \geq \omega_2 \geq \dots \geq \omega_r > 0 \quad \text{the singular values}$$

\mathbf{U} : a $n \times r$ matrix $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r]$, orthonormal: $\mathbf{U}^T \cdot \mathbf{U} = \mathbf{1}$
 \mathbf{u}_k $k=1 \dots r$, the *left singular vectors*.

\mathbf{V} : a $m \times r$ matrix $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r]$, orthonormal: $\mathbf{V}^T \cdot \mathbf{V} = \mathbf{1}$
 \mathbf{v}_k $k=1 \dots r$, the *right singular vectors*.

See Wilks 2005, Sec. 9.3.5

Appendix A



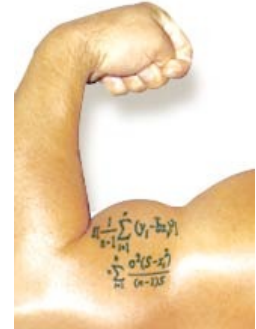
- Singular vectors $\{\mathbf{u}_k\}$ and $\{\mathbf{v}_k\}$ ($k=1..r$) constitute orthonormal coordinate system in r -dim subspace of n -dim and m -dim phase-spaces (i.e. not necessarily complete)
- Sum of squared singular values

$$\|Q\|_F^2 := \sum_i^n \sum_j^m q_{ij}^2 = \sum_i^r \omega_i^2$$

Frobenius norm = sum of squared matrix elements
= sum of squared singular values

see e.g. Bretherton et al. 1992

Appendix A



- **Calculation of SVD**

- Singular vectors $\{\mathbf{u}_k\}$ are eigenvectors of $\mathbf{Q}\mathbf{Q}^T$ (symm, $n \times n$) and $\{\mathbf{v}_k\}$ of $\mathbf{Q}^T\mathbf{Q}$ (symm, $m \times m$)

$$(\mathbf{Q} \cdot \mathbf{Q}^T) \cdot \mathbf{L} = \mathbf{U} \mathbf{\Omega} \mathbf{V}^T \cdot \mathbf{V} \mathbf{\Omega}^T \mathbf{U}^T \mathbf{U} = \mathbf{L} \mathbf{\Omega}^2, \quad \text{dito for } \mathbf{V}$$

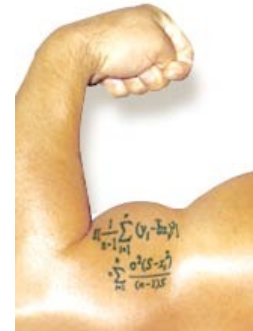
- Singular left vectors can be obtained directly from singular right vectors (and vice versa):

$$\mathbf{Q}\mathbf{V} = \mathbf{U}\mathbf{\Omega}, \quad \mathbf{Q}^T \mathbf{U} = \mathbf{V}\mathbf{\Omega}^T$$

$$\mathbf{Q}\mathbf{v}_k = \omega_k \mathbf{u}_k, \quad \mathbf{Q}^T \mathbf{u}_k = \omega_k \mathbf{v}_k, \quad k = 1, \dots, r$$

- In practice: dedicated software to svd

Appendix B



- **Transformation of Cross-Covariance Matrix**

- If \mathbf{X} , \mathbf{Y} are centered data matrices, the cross-covariance matrix is:

$$\mathbf{S}_{xy} = \frac{1}{N-1} \mathbf{X}^T \cdot \mathbf{Y}$$

- Let \mathbf{U} , \mathbf{V} be transformation matrices (with projection vectors in columns), the data matrices in transformed coordinates are:

$$\mathbf{A} = \mathbf{X} \cdot \mathbf{U}, \quad \mathbf{B} = \mathbf{Y} \cdot \mathbf{V}$$

- The cross-covariance matrix of transformed variables is:

$$\mathbf{S}_{ab} = \frac{1}{N-1} \mathbf{A}^T \cdot \mathbf{B} = \frac{1}{N-1} \mathbf{U}^T \mathbf{X}^T \cdot \mathbf{Y} \mathbf{V} = \mathbf{U}^T \mathbf{S}_{xy} \mathbf{V}$$

- When \mathbf{U} , \mathbf{V} are singular vectors systems of \mathbf{S}_{xy} then:

$$\begin{aligned} \mathbf{S}_{xy} &= \mathbf{U} \mathbf{\Omega} \mathbf{V}^T \\ \mathbf{U}^T \cdot \mathbf{U} &= \mathbf{1}, \mathbf{V}^T \cdot \mathbf{V} = \mathbf{1} \end{aligned} \Rightarrow \mathbf{S}_{ab} = \mathbf{U}^T \mathbf{S}_{xy} \mathbf{V} = \mathbf{U}^T \mathbf{U} \mathbf{\Omega} \mathbf{V}^T \mathbf{V} = \mathbf{\Omega}$$