

# Covariance matrix estimation and error model in D&A

Aurélien Ribes

CNRM - GAME, Météo France - CNRS

Buenos Aires, 16 October 2012

- 1 Introduction
- 2 Covariance matrix estimation
- 3 Model Error in D&A
- 4 Conclusion

# Statistical model

- 1 Introduction
- 2 Covariance matrix estimation
- 3 Model Error in D&A
- 4 Conclusion

# Statistical model

$$Y_\ell = \sum_{i=1}^I \beta_i X_\ell^{(i)} + \varepsilon_\ell \quad (1)$$

- $\ell$  : the location (spatio-temporal),
- $Y$  : observations (spatio-temporal vector),
- $\beta_i$  : scaling factors (real number), unknown,
- $X^{(i)}$  : expected response to the  $i^{\text{th}}$  forcing (spatio-temporal vector), known,
- $\varepsilon$  : internal variability (spatio-temporal vector).

# Statistical model

$$Y_\ell = \sum_{i=1}^I \beta_i X_\ell^{(i)} + \varepsilon_\ell \quad (1)$$

- $\ell$  : the location (spatio-temporal),
- $Y$  : observations (spatio-temporal vector),
- $\beta_i$  : scaling factors (real number), unknown,
- $X^{(i)}$  : expected response to the  $i^{\text{th}}$  forcing (spatio-temporal vector), known,
- $\varepsilon$  : internal variability (spatio-temporal vector).

## 2 key assumptions

- (H1) : The distribution of  $\varepsilon$  is known (from climate models) [D&A],
- (H2) :  $X^{(i)}$  are known (from climate models) [A].

# Optimal fingerprint : estimation

## 2 key assumptions

- (H1) : The distribution of  $\varepsilon$  is known (from climate models) [D&A],
- (H2) :  $X^{(i)}$  are known (from climate models) [A].

Let  $C = \text{Cov}(\varepsilon)$  ;  $C$  is known according to (H2).

$$\hat{\beta} = (X' C^{-1} X)^{-1} X' C^{-1} Y, \quad (2)$$

## Discussed issues

- How can we estimate  $C$  to approximate  $\hat{\beta}$  ?
- Can we deal with uncertainties in  $X$  ?

- 1 Introduction
- 2 Covariance matrix estimation**
- 3 Model Error in D&A
- 4 Conclusion

# High dimension in climate datasets

Let us assume that climate models are able to perfectly simulate the climate system dynamics.

## Typical climate dataset (e.g. near-surface temperature)

- Spatial dimension :  $5^\circ \times 5^\circ \sim 2600$  grid-points,
- Temporal dimension : 50 - 100 ans (instrumental period),
- Dimension of  $Y \sim 10^5$ .
- Internal variability is described by  $C \sim 10^5 \times 10^5$ .
- The estimation of  $C$  requires *at least*  $10^5$  realisations of  $\varepsilon$ , i.e.  $10^7$  yrs of control simulations (vs about  $\sim 10^4$  yrs available).

Some options :

- Decrease the dimension of  $Y$ ,
- Look for an estimator of  $C$  *accurate* in large dimension.



# Decreasing the dimension

Statistical investigation of climate at the global scale requires to reduce the spatio-temporal dimension of datasets.

- Decadal means,
- Projection on principal components,
- Projection on spherical harmonics (e.g. truncation T4, ~ spatial scales  $> 5000$  kms),
- Use of simple climate indices (globale mean, land-sea contrast, inter-hemispheric contrast, annual cycle, etc).

# Estimation of $C$ in large dimension

Let us assume that  $x_1, \dots, x_n \sim N(0, C)$  are available for estimating  $C$  ( $p \times p$ ).

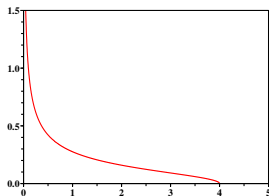
# Estimation of $C$ in large dimension

Let us assume that  $x_1, \dots, x_n \sim N(0, C)$  are available for estimating  $C$  ( $p \times p$ ).

What about  $\hat{C}$  ?

The sample estimate  $\hat{C}$  is a poor estimator of  $C$  in large dimension ( $n$  close to  $p$ ).

Illustration : case  $C = I$ , distribution of the eigenvalues of  $\hat{C}$  when  $n, p \rightarrow \infty$  (Marčenko-Pastur distribution).



# Regularisation of $C$ (1)

## Principle

### Principle

We use an estimator of  $C$  such as

$$\tilde{C} = \gamma \hat{C} + \rho I.$$

# Regularisation of $C$ (1)

## Principle

### Principle

We use an estimator of  $C$  such as

$$\tilde{C} = \gamma \hat{C} + \rho I.$$

### LW estimate (Ledoit & Wolf, 2004)

- Introduction of estimators  $\hat{\gamma}, \hat{\rho}$  of  $\gamma, \rho$  to minimise the mean square error

$$E \left( \|\tilde{C} - C\|_{\mathcal{M}}^2 \right).$$

- 

$$\hat{C}_l = \hat{\gamma} \hat{C} + \hat{\rho} I.$$

# Regularisation of $C$ (1)

## Principle

### Principle

### LW estimate (Ledoit & Wolf, 2004)

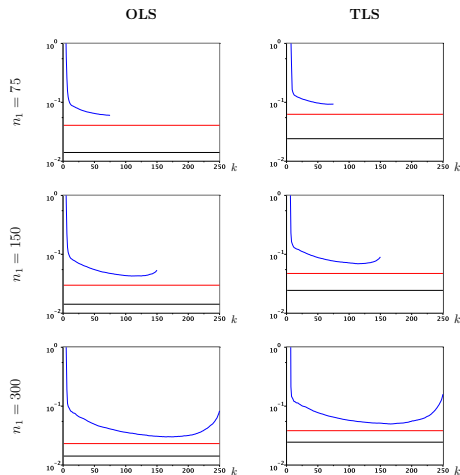
$$\hat{C}_l = \hat{\gamma} \hat{C} + \hat{\rho} l.$$

### New estimator

$$\hat{\beta}_l = (X' \hat{C}_l^{-1} X)^{-1} X' \hat{C}_l^{-1} Y.$$

# Regularisation of $C$ (2)

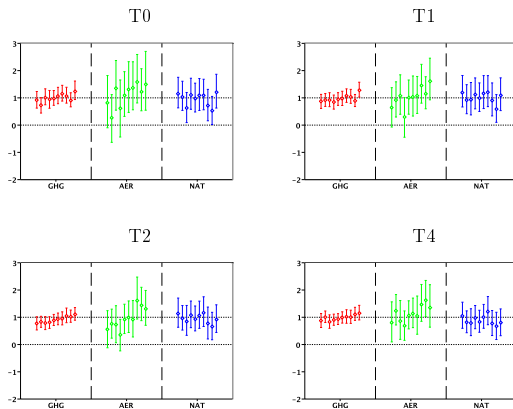
Results : Mean square error



**FIG. :** Mean square error of  $\beta$ -estimates based on  $\widehat{C}_l$  (red, regularisation),  $\widehat{C}_q^+$ ,  $q = 1 \dots p$  (blue,  $q$ -truncation) and  $C$  (black, perfect estimation). Estimation based on Monte-Carlo simulations, for three values of  $n$  ( $p = 250$  here), and under OLS and TLS models.

# Regularisation of $C$ (3)

Resultats : Illustration

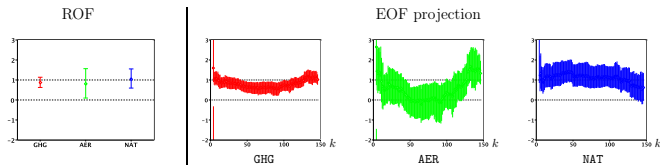


**FIG. :** Results obtained by regularising  $\hat{C}$  on pseudo-observations (historical simulations with the CNRM-CM5 climate model), based on different spatial resolutions in the pre-processing step (spherical harmonics T0, T1, T2 and T4).



# Regularisation of $C$ (3)

Resultats : Illustration



**FIGURE 1:** Sensitivity of the results to the choice  $k$  in the case of a projection on  $k$  principal components.

# Questions

- How the LW estimate performs if  $p \gg n$ ?
- Can we propose other way of regularising (eg with other priors than  $I$ ) ?
- Can we use more recent results from random matrix theory ?

- 1 Introduction
- 2 Covariance matrix estimation
- 3 Model Error in D&A**
- 4 Conclusion

# Key assumptions (recall)

## 2 key assumptions

- (H1) : The distribution of  $\varepsilon$  is known (from climate models) [D&A],
- (H2) :  $X^{(i)}$  are known (from climate models) [A].

Model error in D&A primarily refers to the climate model error, i.e. error in  $X$ .

# Internal variability within the model

## 1 Problem :

$X^{(i)}$  is estimated from climate model simulations, which do contain internal variability.

What is actually observed is  $X^{(i)} + \varepsilon_X^{(i)}$ .

# Internal variability within the model

## 1 Problem :

$X^{(i)}$  is estimated from climate model simulations, which do contain internal variability.

What is actually observed is  $X^{(i)} + \varepsilon_X^{(i)}$ .

## 2 Solution :

One may use a **Total Least Squares** algorithm instead of ordinary least squares (Allen & Stott, 2003).

# TLS : Model, estimation

## Revised statistical model (TLS)

Regression equation  $Y_0 = X_0\beta$

One observes  $\begin{cases} X = X_0 + \varepsilon_X \\ Y = Y_0 + \varepsilon_Y \end{cases}$

$$\text{Cov}(\varepsilon_X) = \text{Cov}(\varepsilon_Y) \quad (\text{or } \text{Cov}(\varepsilon_X) = \text{Cov}(\varepsilon_Y)/m)$$

# TLS : Model, estimation

## Revised statistical model (TLS)

Regression equation  $Y_0 = X_0\beta$

One observes  $\begin{cases} X = X_0 + \varepsilon_X \\ Y = Y_0 + \varepsilon_Y \end{cases}$

$$\text{Cov}(\varepsilon_X) = \text{Cov}(\varepsilon_Y) \quad (\text{or } \text{Cov}(\varepsilon_X) = \text{Cov}(\varepsilon_Y)/m)$$

- Alternative writing :

$$Y = (X - \varepsilon_X)\beta + \varepsilon_Y.$$

- $Y$  and  $X$  play symmetric roles,
- $\hat{\beta}$ ,  $\hat{X}_0$  and  $\hat{Y}_0$  provided by the SVD of  $[X, Y]$ ,
- Confidence intervals are not symmetric in terms of  $\beta$ .



# Modeling uncertainty (1)

## 1 Problem :

Even if internal variability within the model is accounted for,  $X^{(i)}$  is model-dependent (different models would provide different  $X^{(i)}$ ).

## 2 Solution :

Take into account the model discrepancies into the analysis : EIV method (Errors In Variables, Huntigford & al., 2006).

# Modeling uncertainty (1)

## TLS model revised

Regression equation  $Y_0 = X_0\beta$

One observes 
$$\begin{cases} Y = Y_0 + \varepsilon_{Y,IV} \\ X = X_0 + \varepsilon_{X,IV} + \nu_{Mod} \end{cases}$$

With 
$$\begin{cases} \text{Cov}(\varepsilon_{Y,IV}) = C, & \text{(I.V.)} \\ \text{Cov}(\varepsilon_{X,IV} + \nu_{Mod}) = C + \Sigma_{Mod}, & \text{(I.V. + Mod. Uncert.)} \end{cases}$$

# Modeling uncertainty (1)

## TLS model revised

Regression equation

$$Y_0 = X_0\beta$$

One observes

$$\begin{cases} Y = Y_0 + \varepsilon_{Y,IV} + \epsilon_{Obs} \\ X = X_0 + \varepsilon_{X,IV} + \nu_{Mod} \end{cases}$$

With

$$\begin{cases} \text{Cov}(\varepsilon_{Y,IV} + \epsilon_{Obs}) = C + \Sigma_{Obs}, & (\text{I.V. + Obs. Uncert.}) \\ \text{Cov}(\varepsilon_{X,IV} + \nu_{Mod}) = C + \Sigma_{Mod}, & (\text{I.V. + Mod. Uncert.}) \end{cases}$$

# Modeling uncertainty (1)

## TLS model revised

Regression equation

$$Y_0 = X_0\beta$$

One observes

$$\begin{cases} Y = Y_0 + \varepsilon_{Y,IV} + \epsilon_{Obs} \\ X = X_0 + \varepsilon_{X,IV} + \nu_{Mod} \end{cases}$$

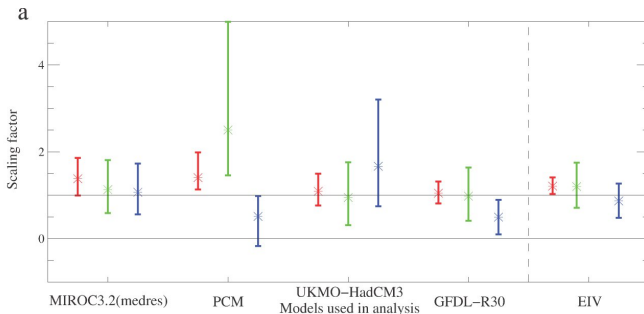
With

$$\begin{cases} \text{Cov}(\varepsilon_{Y,IV} + \epsilon_{Obs}) = C + \Sigma_{Obs}, & (\text{I.V.} + \text{Obs. Uncert.}) \\ \text{Cov}(\varepsilon_{X,IV} + \nu_{Mod}) = C + \Sigma_{Mod}, & (\text{I.V.} + \text{Mod. Uncert.}) \end{cases}$$

- Estimation and hypothesis testing are more complicated than in TLS.
- Only a few studies have used this approach until now.

# Modeling uncertainties (2)

## Illustration



**FIGURE 2:** Illustration of the “EIV” method (IPCC, AR4, 2007)

# Questions

- Can we improve the inference in EIV ?
- How can we estimate  $\Sigma_{Obs}$  and  $\Sigma_{Mod}$  ?
- ...

# Conclusions

## Covariance matrix estimation

- Estimation of large covariance matrix is required in D&A,
- The use of regularised estimates may help,
- Are such estimates potentially useful in DA ?

## Model error

- Model error in D&A primarily refers to uncertainty on the climate response simulated by climate models,
- This error is typically taken into account in Errors In Variable regression models.