# Modélisation Numérique
# de l'Écoulement Atmosphérique
# et Assimilation de Données

Olivier Talagrand
Cours 5

12 Mai 2025

*Kalman Filter*

Two questions

- *How to propagate information backwards in time ?* (useful for reassimilation of past data)

- *How to take into account possible dependence in time ?*

Kalman Filter, whether in its standard linear form or in its Ensemble form, does neither.

Rudolf Kálmán (1930-2016)

- Lissage de Kalman.

- Assimilation variationnelle. Principe

- Méthode adjointe. Principe.

- Assimilation variationnelle. Résultats

- La Méthode incrémentale

- Compléments sur l'Estimation Statistique (*BLUE*)

# *Kalman smoother*

Propagates information both forward and backward in time, as does 4DVar, but uses Kalman-type formulæ

Various possibilities

- Define new state vector $X^{\mathrm{T}} \equiv (x_0^{\mathrm{T}}, \ldots, x_K^{\mathrm{T}})$

  and use general *BLUE* formula from a background $X^b$ and associated covariance matrix $\Pi^b$.

  Model equations, which bring information on the $x_i$ 's, must be included in the observation vector and the associated observation operator.

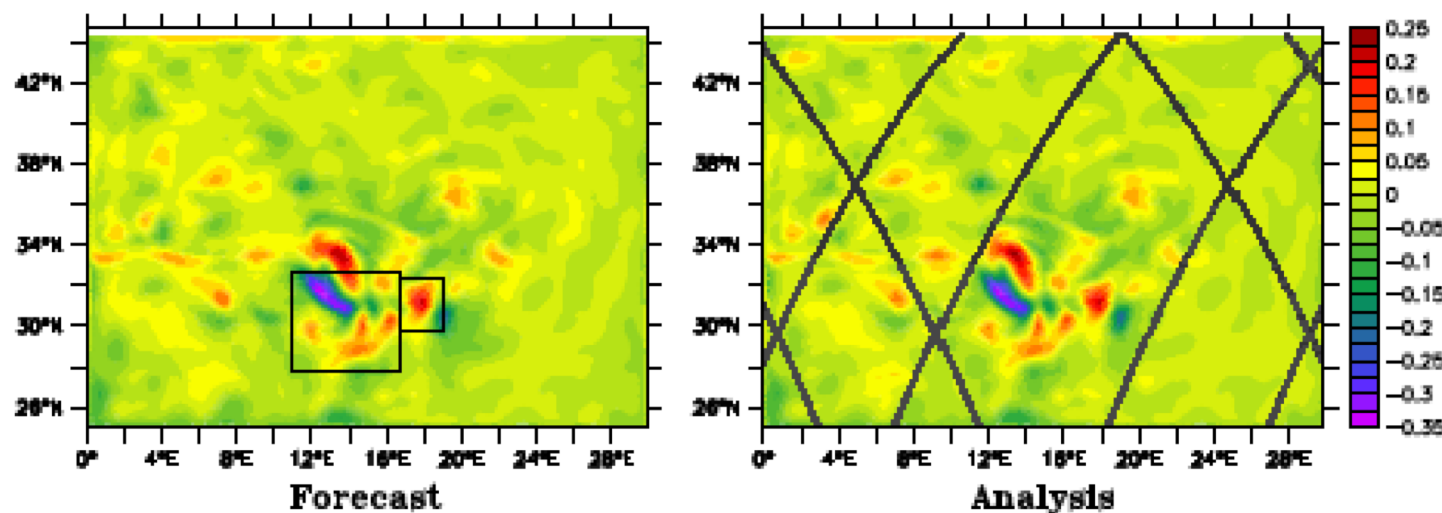  Can take into account temporal correlations

- Update sequentially vector $(x_0^{\mathrm{T}}, \ldots, x_k^{\mathrm{T}})^{\mathrm{T}}$ for increasing $k$

  Cannot take into account temporal correlations

  Algorithms exist in ensemble form

E. Cosme (2015)

Ensemble smoother based on *Singular Evolutive Extended Kalman Filter* (*SEEK*)

Of second type above. Retropropagates corrections on fields backwards in time, but without modifying relative weights given to previous data, *i.e.* cannot be optimal in case of temporal dependence between errors.

E. Cosme,
HDR,
2015,
Lissage
d'ensemble
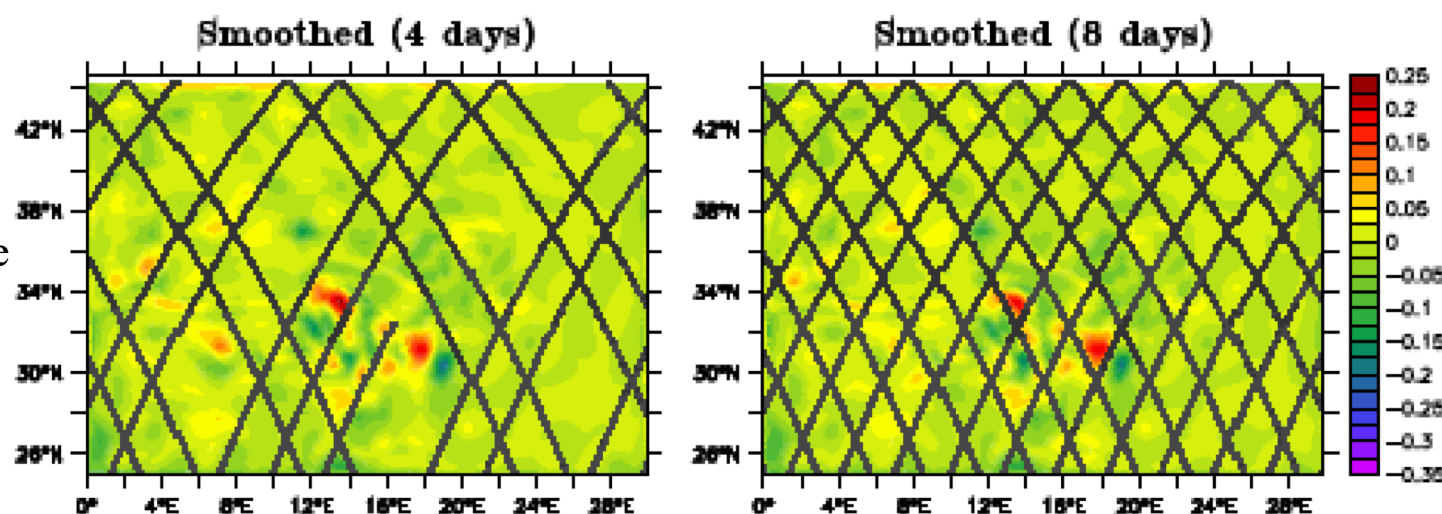SEEK

Données
synthétiques

FIGURE 3.6 – Evolution du champ d'erreur en SSH du jour 38, au cours des étapes d'analyse successives. En haut à gauche : prévision du filtre ; en haut à droite : analyse du filtre. Les observations utilisées pour cette analyse sont distribuées le long des traces grises. En bas à gauche : analyse du lisseur après introduction des observations des jours 40 et 42 ; En bas à droite : analyse du lisseur après introduction des observations des jours 40 à 46.
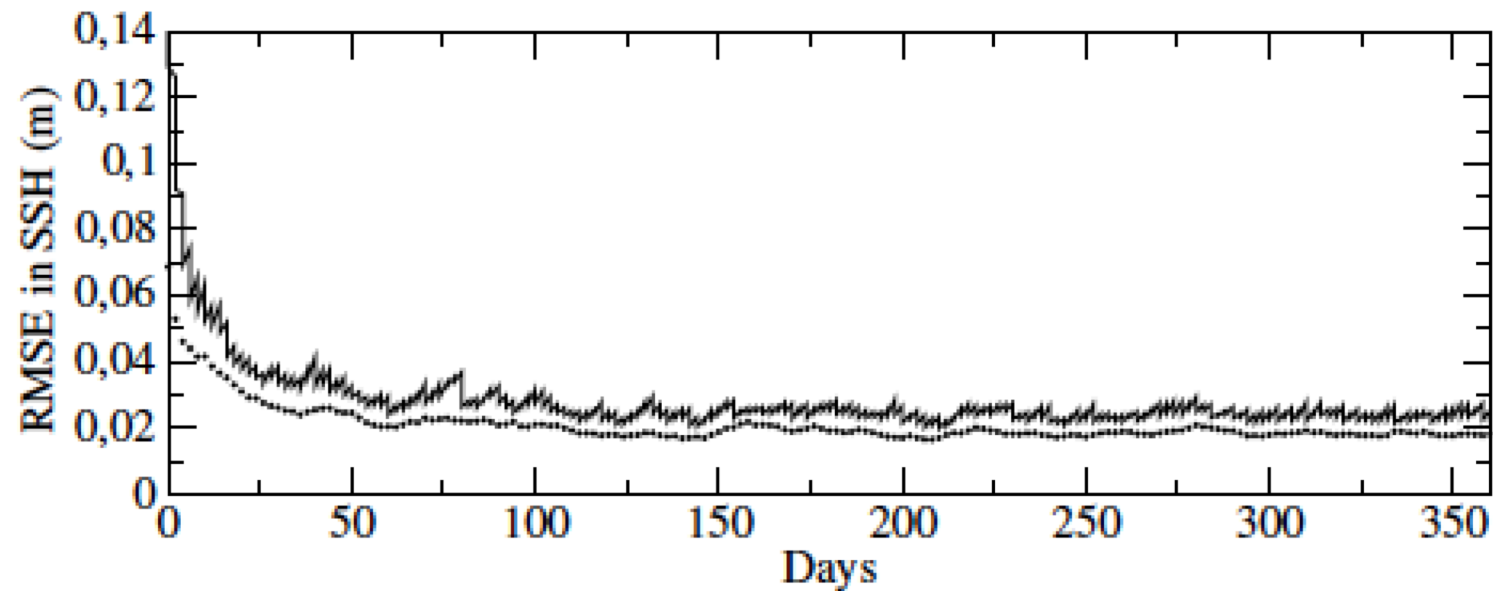
FIGURE 3.7 – Evolution de l'erreur RMS de SSH au cours du temps. Ligne continue : Résultat du filtre (les dents de scie reflètent l'alternance des étapes de prévision et d'analyse) ; Points : lisseur à retard fixe de 8 jours.

E. Cosme, HDR, 2015, Lissage d'ensemble SEEK

Other variants of Ensemble Kalman Smoothers


*An iterative ensemble Kalman smoother* (Bocquet and Sakov, 2014. *Q. J. R. Meteorol. Soc.*)


*An Iterative Ensemble Kalman Smoother in Presence of Additive Model Error* (Fillion *et al.*, 2019, *SIAM/ASA J. Uncertainty Quantification)*

*Case of data that are distributed over time*

Suppose for instance available data consist of

- Background estimate at time $0$

$$x_0^b = x_0 + \zeta_0^b \qquad E(\zeta_0^b \zeta_0^{bT}) = P_0^b$$

- Observations at times $k = 0, \ldots, K$

$$y_k = H_k x_k + \varepsilon_k \qquad E(\varepsilon_k \varepsilon_j^T) = R_k \, \delta_{kj}$$

- Model (supposed for the time being to be exact)

$$x_{k+1} = M_k x_k \qquad k = 0, \ldots, K\text{-}1$$

Errors assumed to be unbiased and uncorrelated in time, $H_k$ and $M_k$ linear

Then objective function

$$\xi_0 \in S \rightarrow$$

$$\mathcal{J}(\xi_0) \equiv (1/2) (x_0^b - \xi_0)^T [P_0^b]^{-1} (x_0^b - \xi_0) + (1/2) \, \Sigma_k [y_k - H_k \xi_k]^T R_k^{-1} [y_k - H_k \xi_k]$$

$$\equiv \qquad\qquad\qquad \mathcal{J}_b \qquad\qquad + \qquad\qquad \mathcal{J}_o$$

subject to $\xi_{k+1} = M_k \xi_k, \qquad k = 0, \ldots, K\text{-}1$

# Principle of 4D-VAR assimilation

$$\mathcal{J}(\boldsymbol{\xi}_0) = (1/2)\,(\boldsymbol{x}_0^b - \boldsymbol{\xi}_0)^T\,[\boldsymbol{P}_0^b]^{-1}\,(\boldsymbol{x}_0^b - \boldsymbol{\xi}_0) + (1/2)\,\Sigma_k[\boldsymbol{y}_k - \boldsymbol{H}_k\boldsymbol{\xi}_k]^T\,\boldsymbol{R}_k^{-1}\,[\boldsymbol{y}_k - \boldsymbol{H}_k\boldsymbol{\xi}_k]$$

subject to  $\boldsymbol{\xi}_{k+1} = \boldsymbol{M}_k\boldsymbol{\xi}_k\,,\quad k = 0, \ldots, K\text{-}1$

Background is not necessary, if observations are in sufficient number to overdetermine the problem. Nor is strict linearity.

*Four-Dimensional Variational Assimilation*

**'4D-Var'**

How to minimize objective function with respect to initial state $u = \xi_0$ ($u$ is called the *control variable* of the problem) ?

Use iterative minimization algorithm, each step of which requires the explicit knowledge of the local gradient $\nabla_u \mathcal{J} \equiv (\partial \mathcal{J}/\partial u_i)$ of $\mathcal{J}$ with respect to $u$.

How to numerically compute the gradient $\nabla_u \mathcal{J}$ ?

Direct perturbation, in order to obtain partial derivatives $\partial \mathcal{J}/\partial u_i$ by finite differences ? That would require as many explicit computations of the objective function $\mathcal{J}$ as there are components in $\boldsymbol{u}$. Practically impossible.

Gradient computed by *adjoint method*.

-  Méthode adjointe. Principe.

**Adjoint Method**

*Input vector* $\boldsymbol{u} = (u_i)$, $\dim\boldsymbol{u} = n$

Numerical process, implemented on computer (*e. g.* integration of numerical model)

$$\boldsymbol{u} \rightarrow \boldsymbol{v} = \boldsymbol{G}(\boldsymbol{u})$$

$\boldsymbol{v} = (v_j)$ is *output vector*, $\dim\boldsymbol{v} = m$

Perturbation $\delta\boldsymbol{u} = (\delta u_i)$ of input. Resulting first-order perturbation on $\boldsymbol{v}$

$$\delta v_j = \Sigma_i \, (\partial v_j / \partial u_i) \, \delta u_i$$

or, in matrix form

$$\delta\boldsymbol{v} \; = \; \boldsymbol{G'} \, \delta\boldsymbol{u}$$

where $\boldsymbol{G'} \equiv (\partial v_j / \partial u_i)$ is local matrix of partial derivatives, or *jacobian matrix*, of $\boldsymbol{G}$.

**Adjoint Method (continued 1)**

$$\delta v \;=\; G' \, \delta u \qquad\qquad (D)$$

Scalar function of output

$$\mathcal{J}(v) \;=\; \mathcal{J}[G(u)]$$

Gradient $\nabla_u \mathcal{J}$ of $\mathcal{J}$ with respect to input $u$?

'Chain rule'

$$\partial \mathcal{J}/\partial u_i = \Sigma_j \, \partial \mathcal{J}/\partial v_j \, (\partial v_j/\partial u_i)$$

or

$$\nabla_u \mathcal{J} \;=\; G'^{\mathrm{T}} \, \nabla_v \mathcal{J} \qquad\qquad (A)$$

## Adjoint Method (continued 2)

$G$ is the composition of a number of successive steps

$$G = G_N \circ \ \ldots \circ \ G_2 \circ \ G_1$$

'Chain rule'

$$G' = G_N' \ \ldots \ G_2' \ G_1'$$

Transpose

$$G'^{\mathrm{T}} = G_1'^{\mathrm{T}} \ G_2'^{\mathrm{T}} \ \ldots \ G_N'^{\mathrm{T}}$$

Transpose, or *adjoint*, computations are performed in reversed order of direct computations.

If $G$ is nonlinear, local jacobian $G'$ depends on local value of input $u$. Any quantity which is an argument of a nonlinear operation in the direct computation will be used again in the adjoint computation. It must be kept in memory from the direct computation (or else be recomputed again in the course of the adjoint computation).

If everything is kept in memory, total operation count of adjoint computation is at most 4 times operation count of direct computation (in practice about 2).
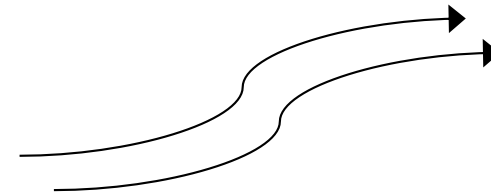
## A few basics

- Basic (nonlinear) model

$$\boldsymbol{x}_{k+1} = \boldsymbol{M}_k(\boldsymbol{x}_k)$$

- Perturbation $\delta\boldsymbol{x}_0$ at time $0$. Resulting perturbation $\delta\boldsymbol{x}_k$ evolves in time according to

$$\delta x_{k+1} = \boldsymbol{M}_k(\boldsymbol{x}_k + \delta\boldsymbol{x}_k) - \boldsymbol{M}_k(\boldsymbol{x}_k)$$

$$= \boldsymbol{M}_k{'}(\boldsymbol{x}_k)\ \delta\boldsymbol{x}_k + o(\delta\boldsymbol{x}_0)$$

where $\boldsymbol{M}_k{'}(\boldsymbol{x}_k)$ is jacobian of $\boldsymbol{M}_k$ at point $x_k$.

$$\delta\boldsymbol{\xi}_{k+1} = \boldsymbol{M}_k{'}(\boldsymbol{x}_k)\ \delta\boldsymbol{\xi}_k$$

is *tangent linear model* along solution $x_k$.

**A few basics** (continuation)

*Tangent linear model*

$$\delta\boldsymbol{\xi}_{k+1} = \mathbf{M}_k{}'(x_k)\,\delta\boldsymbol{\xi}_k$$

*Adjoint model*

$$\lambda_k = [\mathbf{M}_k{}'(\boldsymbol{x}_k)]^{\mathrm{T}}\,\lambda_{k+1}$$

Describes evolution with respect to $k$ of gradient of a scalar function $\mathcal{J}$ with respect to $\boldsymbol{x}_k$.

## Adjoint Method (continued 3)

$$\mathcal{J}(\xi_0) = (1/2)(x_0^b - \xi_0)^T [P_0^b]^{-1} (x_0^b - \xi_0) + (1/2) \Sigma_k [y_k - H_k \xi_k]^T R_k^{-1} [y_k - H_k \xi_k]$$

$$\text{subject to } \xi_{k+1} = M_k \xi_k, \qquad k = 0, \ldots, K\text{-}1$$

Control variable $\qquad \xi_0 = u$

Adjoint equation

$$\lambda_K = \qquad\qquad H_K^T R_K^{-1} [H_K \xi_K - y_K]$$

....

$$\lambda_k = M_k^T \lambda_{k+1} + H_k^T R_k^{-1} [H_k \xi_k - y_k] \qquad\qquad k = K\text{-}1, \ldots, 1$$

....

$$\lambda_0 = M_0^T \lambda_1 + H_0^T R_0^{-1} [H_0 \xi_0 - y_0] + [P_0^b]^{-1} (\xi_0 - x_0^b)$$

$$\nabla_u \mathcal{J} = \lambda_0$$

Result of direct integration ($\xi_k$), which appears in quadratic terms in expression of objective function, must be kept in memory from direct integration.

# Principle of 4D-VAR assimilation

**Adjoint Method (continued 4)**

**Nonlinearities ?**

$$\mathcal{J}(\boldsymbol{\xi}_0) = (1/2)\,(\boldsymbol{x}_0^b - \boldsymbol{\xi}_0)^T\,[\boldsymbol{P}_0^b]^{-1}\,(\boldsymbol{x}_0^b - \boldsymbol{\xi}_0) + (1/2)\,\Sigma_k[\boldsymbol{y}_k - \boldsymbol{H}_k(\boldsymbol{\xi}_k)]^T\,\boldsymbol{R}_k^{-1}\,[\boldsymbol{y}_k - \boldsymbol{H}_k(\boldsymbol{\xi}_k)]$$

subject to $\boldsymbol{\xi}_{k+1} = \boldsymbol{M}_k(\boldsymbol{\xi}_k)$, $\qquad k = 0, \dots, K\text{-}1$

Control variable $\qquad\qquad \boldsymbol{\xi}_0 = \boldsymbol{u}$

Adjoint equation

$$\boldsymbol{\lambda}_K = \qquad\qquad\quad \boldsymbol{H}_K{}'^T\,\boldsymbol{R}_K^{-1}\,[\boldsymbol{H}_K(\boldsymbol{\xi}_K) - \boldsymbol{y}_K]$$

....

$$\boldsymbol{\lambda}_k = \boldsymbol{M}_k{}'^T\boldsymbol{\lambda}_{k+1} + \boldsymbol{H}_k{}'^T\,\boldsymbol{R}_k^{-1}\,[\boldsymbol{H}_k(\boldsymbol{\xi}_k) - \boldsymbol{y}_k] \qquad\qquad\qquad k = K\text{-}1, \dots, 1$$

....

$$\boldsymbol{\lambda}_0 = \boldsymbol{M}_0{}'^T\boldsymbol{\lambda}_1 + \boldsymbol{H}_0{}'^T\,\boldsymbol{R}_0^{-1}\,[\boldsymbol{H}_0(\boldsymbol{\xi}_0) - \boldsymbol{y}_0] + [\boldsymbol{P}_0^b]^{-1}\,(\boldsymbol{\xi}_0 - \boldsymbol{x}_0^b)$$

$$\nabla_u\mathcal{J} = \boldsymbol{\lambda}_0$$

Not approximate (it gives the exact gradient $\nabla_u\mathcal{J}$), and really used as described here.

- Assimilation variationnelle. Résultats

Temporal evolution of the 500-hPa geopotential autocorrelation with respect to point located at 45N, 35W. From top to bottom: initial time, 6- and 24-hour range. Contour interval 0.1. After F. Bouttier.
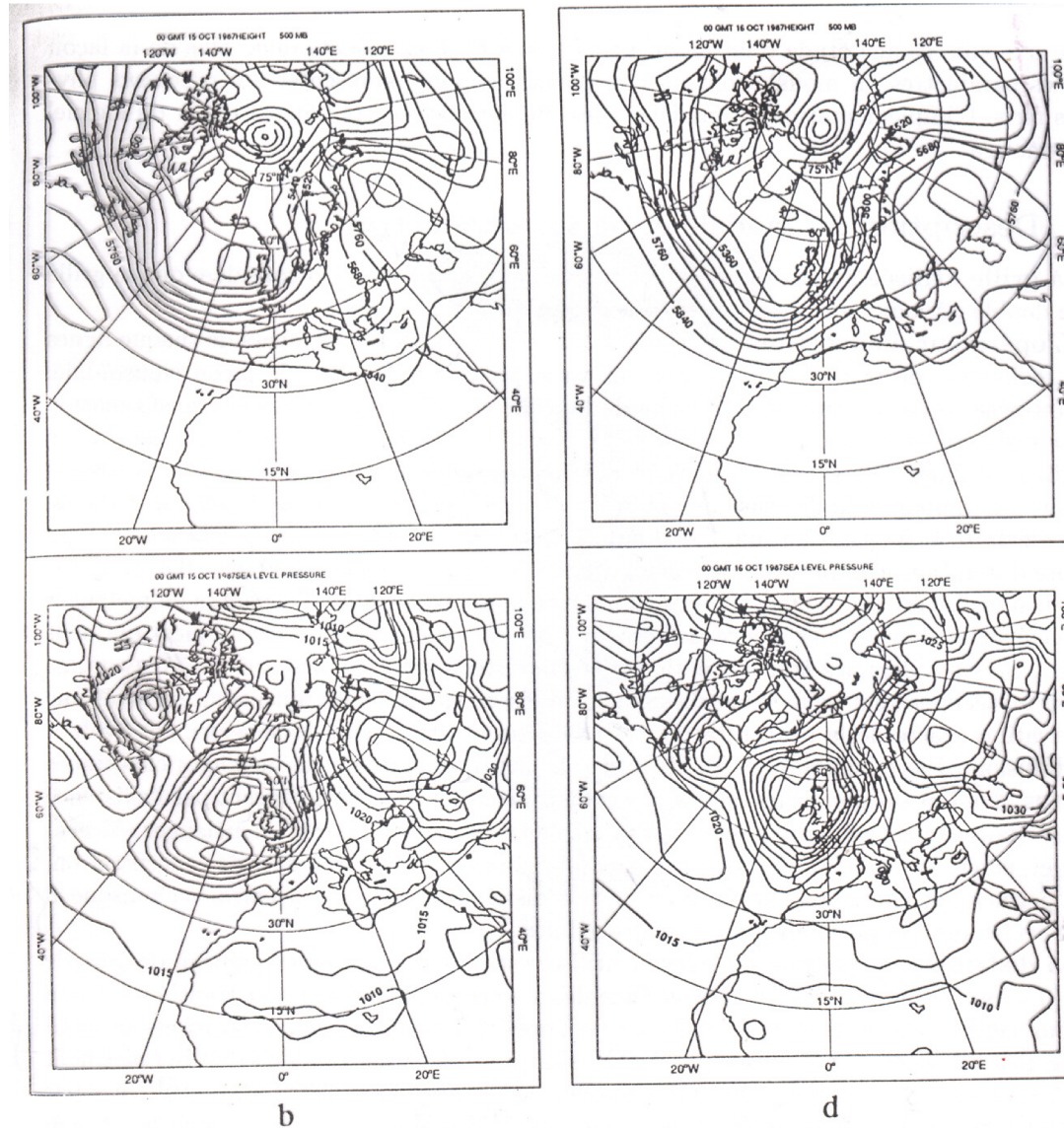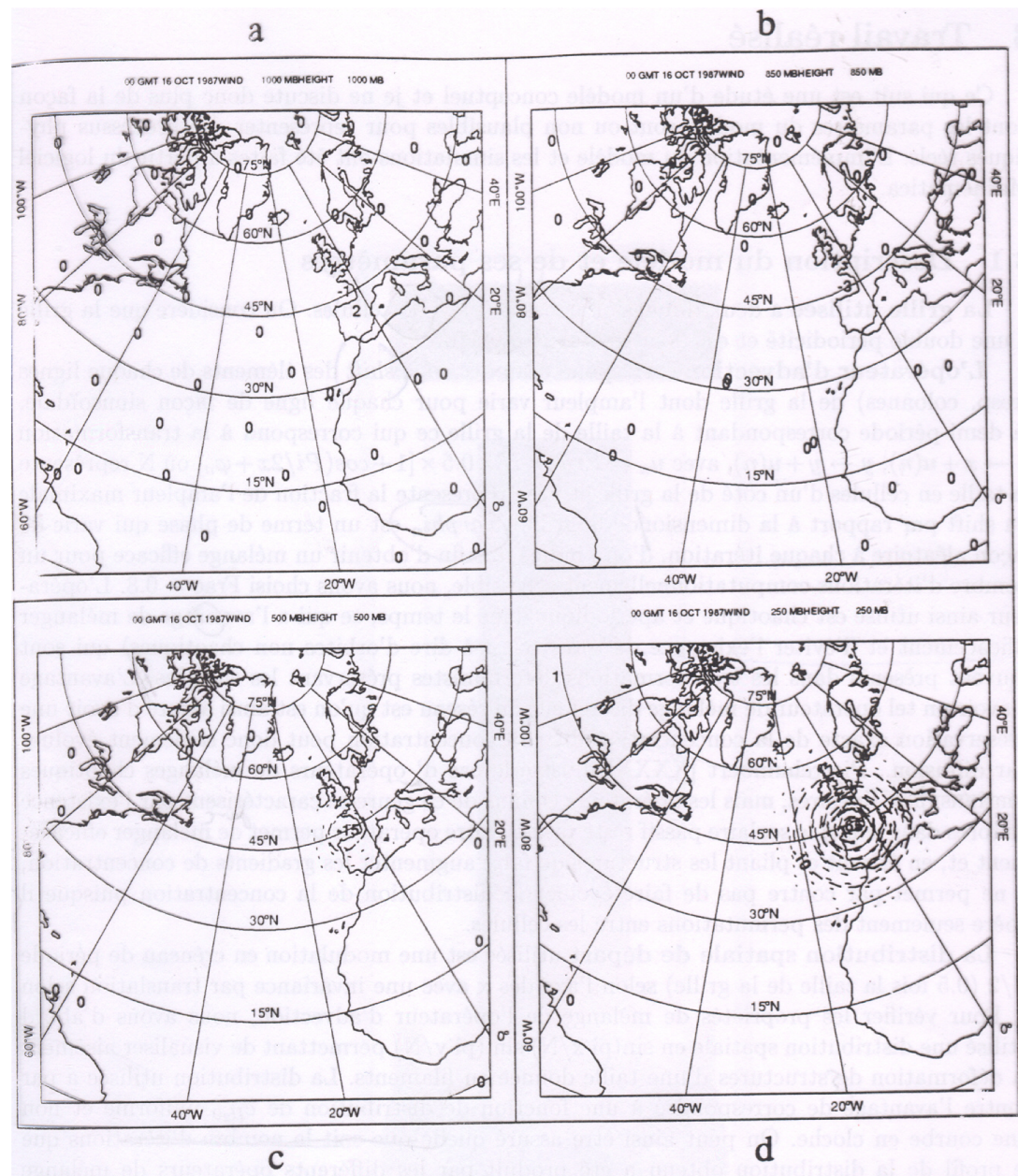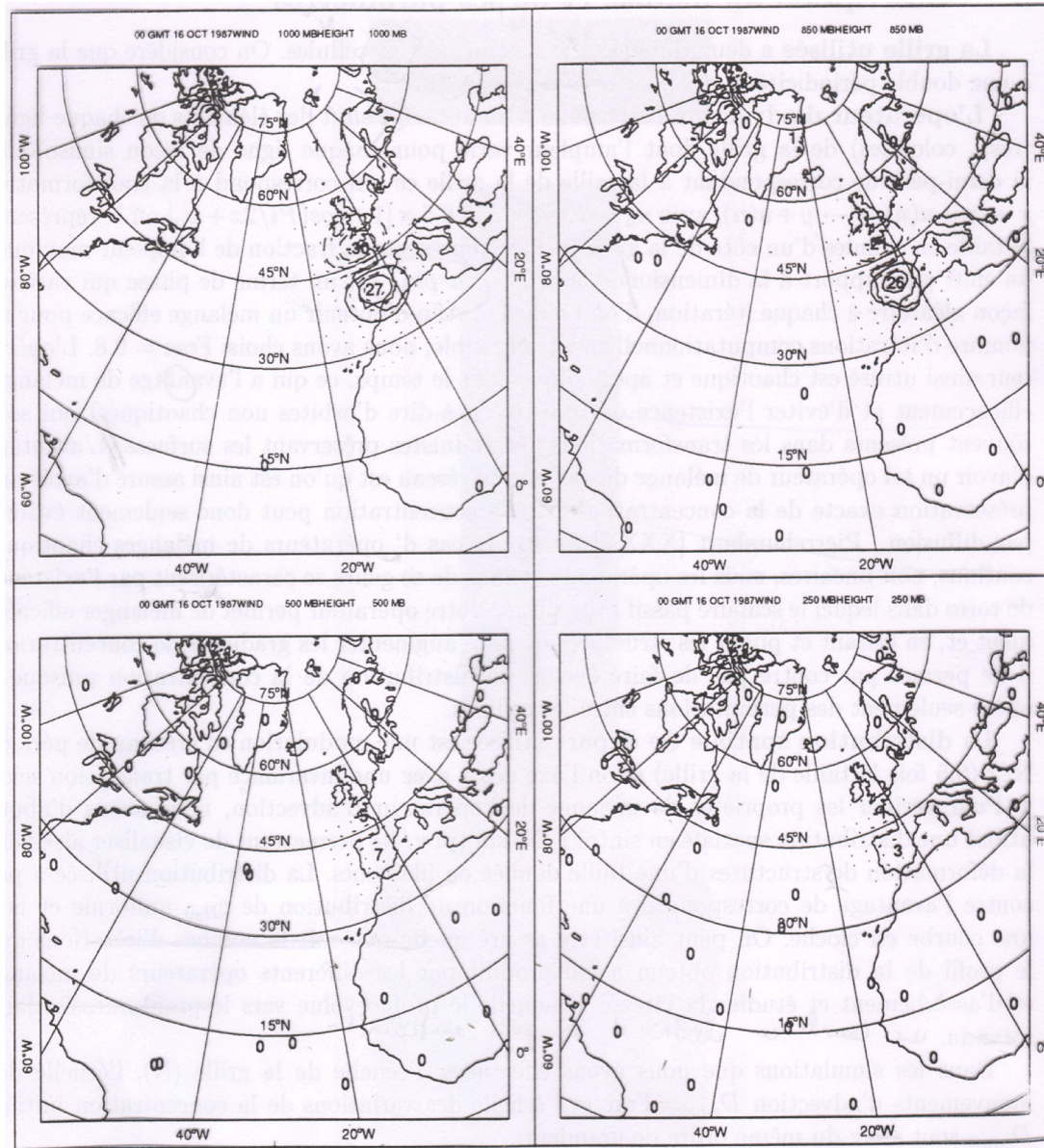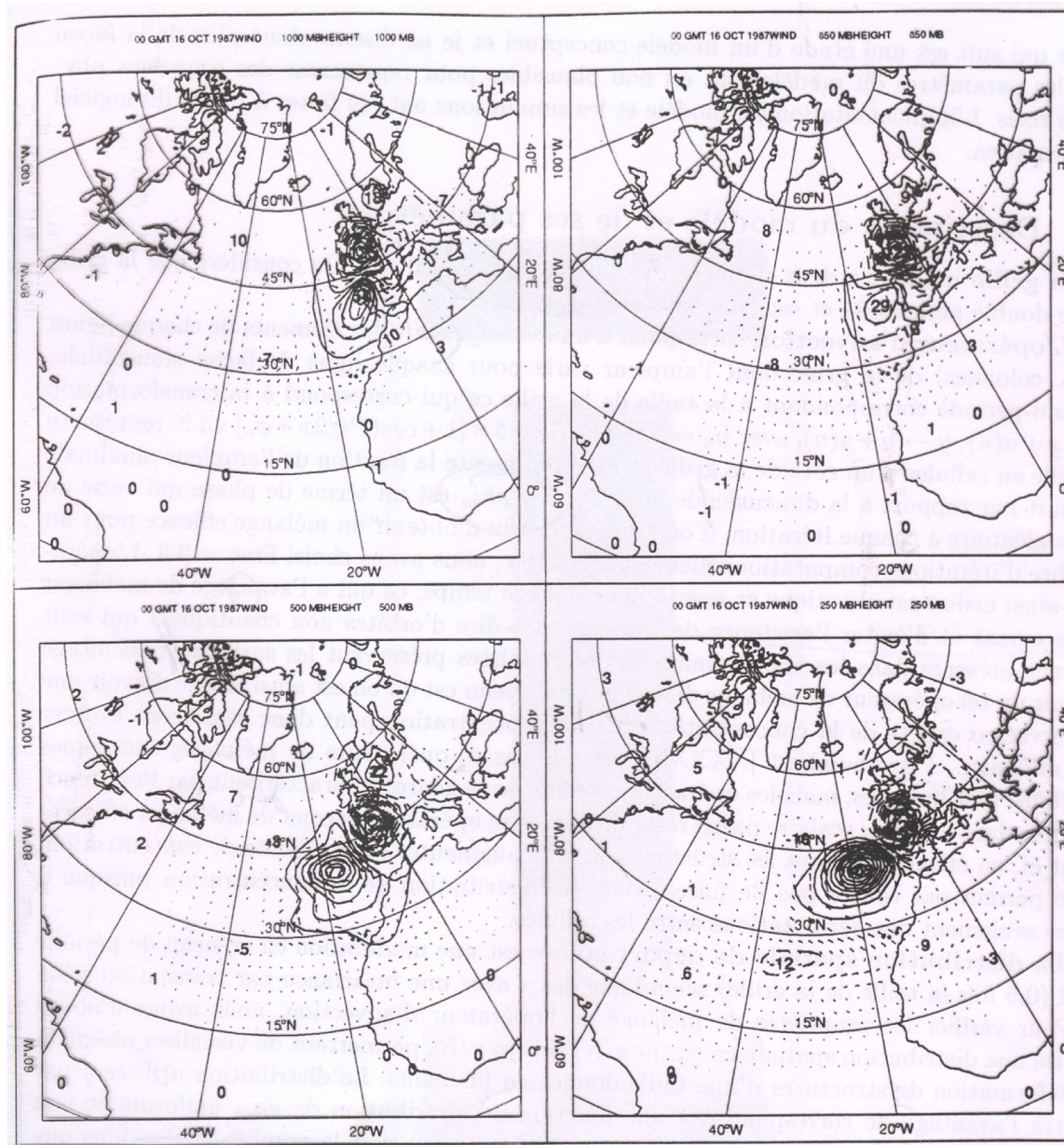
25

FIG. 1. Background fields for 0000 UTC 15 October–0000 UTC 16 October 1987. Shown here are the Northern Hemisphere (a) 500-hPa geopotential height and (b) mean sea level pressure for 15 October and the (c) 500-hPa geopotential height and (d) mean sea level pressure for 16 October. The fields for 15 October are from the initial estimate of the initial conditions for the 4DVAR minimization. The fields for 16 October are from the 24-h T63 adiabatic model forecast from the initial conditions. Contour intervals are 80 m and 5 hPa.

Thépaut *et al*., 1993, *Mon. Wea. Rev.*, **121**, 3393-3414

Analysis increments in a 3D-Var corresponding to a height observation at the 250-hPa pressure level (no temporal evolution of background error covariance matrix)
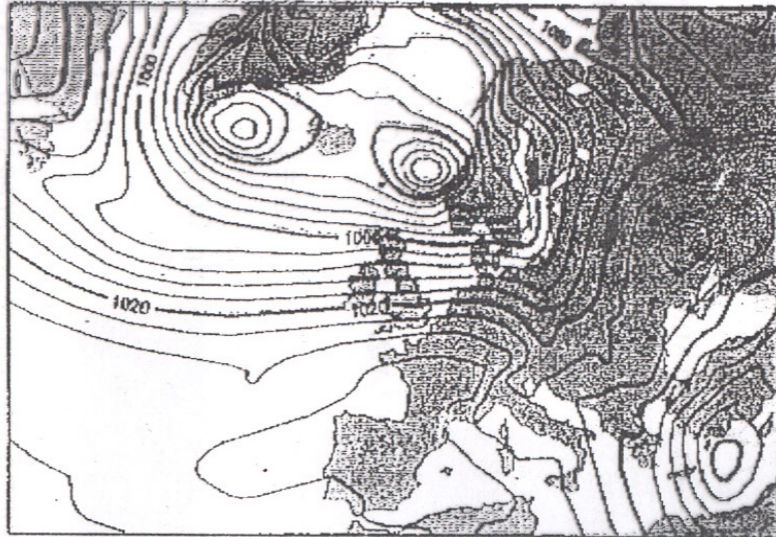
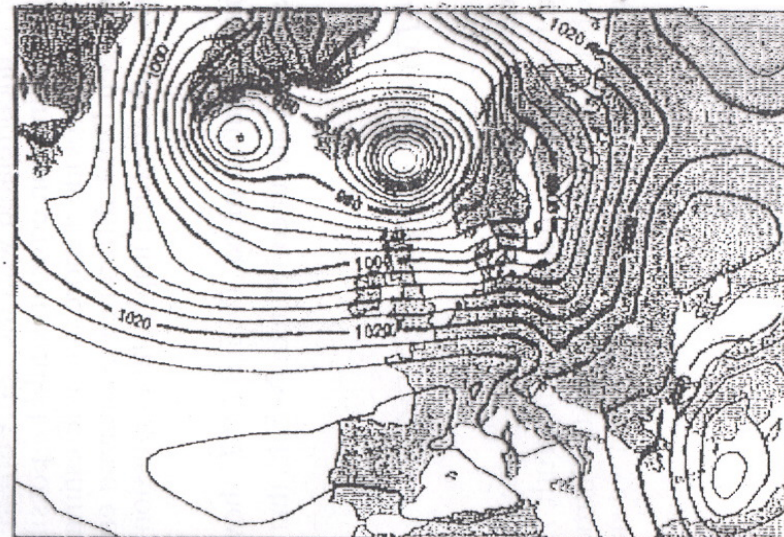Thépaut *et al*., 1993, *Mon. Wea. Rev.*, **121**, 3393-3414

Same as before, but at the end of a 24-hr 4D-Var

Thépaut *et al.*, 1993, *Mon. Wea. Rev.*, **121**, 3393-3414

Analysis increments in a 3D-Var corresponding to a *u*-component wind observation at the 1000-hPa pressure level (no temporal evolution of background error covariance matrix)

Thépaut *et al*., 1993, *Mon. Wea. Rev.*, **121**, 3393-3414

Same as before, but at the end of a 24-hr 4D-Var

Thépaut *et al*., 1993, *Mon. Wea. Rev.*, **121**, 3393-3414

3-day forecast from 3D-Var analysis

3-day forecast from 4D-Var analysis

3D-Var verifying analysis

4D-Var verifying analysis
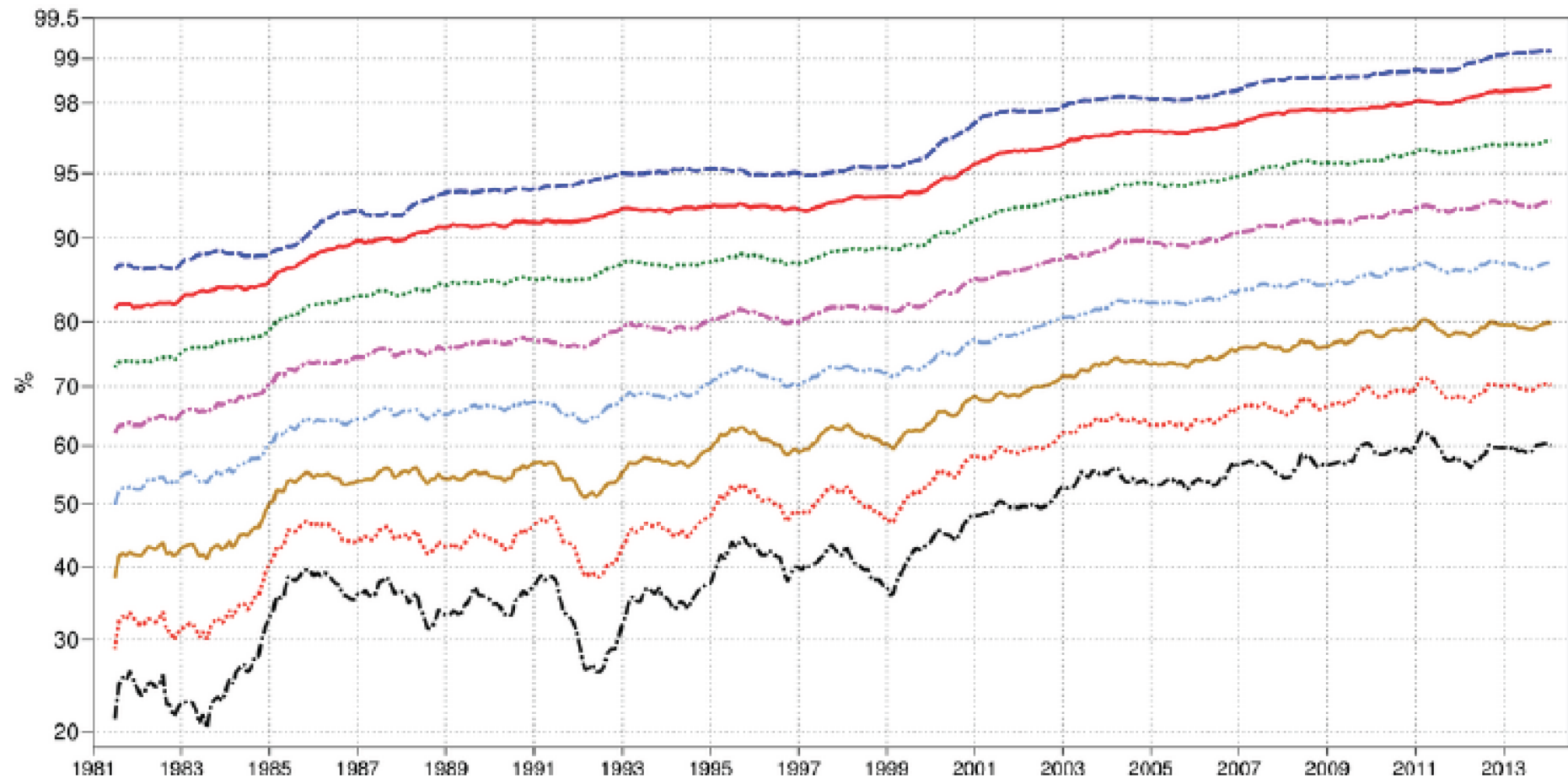
ECMWF, Results on one FASTEX case (1997)

31

**Figure 3:** 500 hPa geopotential height mean square error skill score for Europe (top) and the northern hemisphere extratropics (bottom), showing 12-month moving averages for forecast ranges from 24 to 192 hours. The last point on each curve is for the 12-month period August 2013–July 2014.
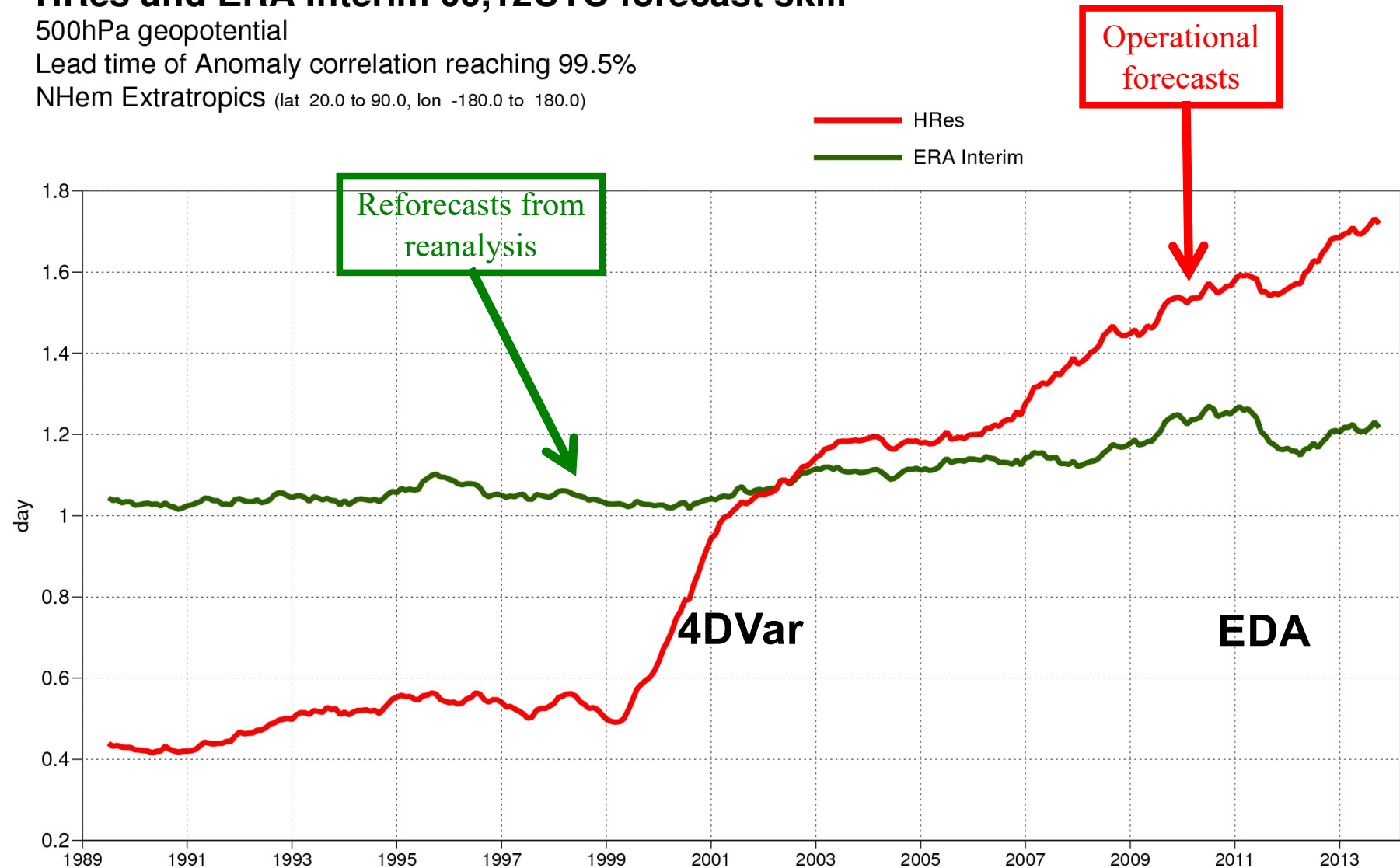
Persistence = 0 ; climatology = 50 at long range

# Initial state error reduction



**HRes and ERA Interim 00,12UTC forecast skill**
500hPa geopotential
Lead time of Anomaly correlation reaching 99.5%
NHem Extratropics (lat 20.0 to 90.0, lon -180.0 to 180.0)

Credit E. Källén, ECMWF

*Strong Constraint 4D-Var* is now used operationally at several meteorological centres (Météo-France, UK Meteorological Office, Canadian Meteorological Centre, Japan Meteorological Agency, …) and, for a number of years, at ECMWF. The latter now has a 'weak constraint' component in its operational system.

**Time-correlated Errors** (*continuation from course 4*)

If data errors are correlated in time, it is not possible to discard observations as they are used. In particular, if model error is correlated in time, all observations are liable to be reweighted as assimilation proceeds.

Variational assimilation can take time-correlated errors into account.

Example of time-correlated observation errors. Global covariance matrix

$$\mathcal{R} = (\boldsymbol{R}_{kk'} = E(\boldsymbol{\varepsilon}_k \boldsymbol{\varepsilon}_{k'}^{\mathrm{T}}))$$

Objective function

$\boldsymbol{\xi}_0 \in S \rightarrow$

$\mathcal{J}(\boldsymbol{\xi}_0) = (1/2) (\boldsymbol{x}_0^b - \boldsymbol{\xi}_0)^{\mathrm{T}} [\boldsymbol{P}_0^b]^{-1} (\boldsymbol{x}_0^b - \boldsymbol{\xi}_0) + (1/2) \Sigma_{kk'} [\boldsymbol{y}_k - \boldsymbol{H}_k \boldsymbol{\xi}_k]^{\mathrm{T}} [\mathcal{R}^{-1}]_{kk'} [\boldsymbol{y}_{k'} - \boldsymbol{H}_{k'} \boldsymbol{\xi}_{k'}]$

where $[\mathcal{R}^{-1}]_{kk'}$ is the $kk'$-sub-block of global inverse matrix $\mathcal{R}^{-1}$.

Similar approach for time-correlated model error.

## Time-correlated Errors (continuation 4)

Temporal correlation of observational error has been introduced by ECMWF (Järvinen *et al*., 1999) in variational assimilation of high-frequency surface pressure observations (correlation originates in that case in representativeness error).

Identification and quantification of time correlation of errors, especially model errors ?

In the linear case, and if errors are uncorrelated in time, Kalman Smoother and Variational Assimilation are algorithmically equivalent, and produce the *BLUE* of the state of the system from all data available over the assimilation window (Kalman Filter produces the *BLUE* only at the end of the final time of the window). If in addition errors are globally Gaussian, both algorithms achieve Bayesian estimation.

If errors are correlated in time, only some Kalman Smoothers can take into account the corresponding correlations, and be equivalent with Variational Assimilation.

- La Méthode incrémentale

**Incremental Method for Variational Assimilation**

Variational assimilation, as it has been described, requires the use of the adjoint of the full model.

Simplifying the adjoint as such can be very dangerous. The computed gradient would not be exact, and experience shows that optimization algorithms (and especially efficient ones) are very sensitive to even slight misspecification of the gradient.

Principle of *Incremental Method* (Courtier *et al.*, 1994, *Q. J. R. Meteorol. Soc.*) : simplify simultaneously the (local tangent linear) dynamics and the corresponding adjoint.

**Incremental Method** (continuation 1)

- Basic (nonlinear) model

$$\boldsymbol{\xi}_{k+1} = \boldsymbol{M}_k(\boldsymbol{\xi}_k)$$

- Tangent linear model

$$\delta\boldsymbol{\xi}_{k+1} = \boldsymbol{M}_k{}' \, \delta\boldsymbol{\xi}_k$$

where $\boldsymbol{M}_k{}'$ is jacobian of $\boldsymbol{M}_k$ at point $\boldsymbol{\xi}_k$.

- Adjoint model

$$\lambda_k = \boldsymbol{M}_k{}'^{\mathrm{T}} \, \lambda_{k+1} + \ldots$$

*Incremental Method*. Simplify both $\boldsymbol{M}_k{}'$ and $\boldsymbol{M}_k{}'^{\mathrm{T}}$ consistently.

## Incremental Method (continuation 2)

More precisely, for given solution $\xi_k^{(0)}$ of nonlinear model, replace tangent linear and adjoint models respectively by

$$\delta\xi_{k+1} = L_k\,\delta\xi_k \qquad\qquad (2)$$

and

$$\lambda_k = L_k^{\mathrm{T}}\,\lambda_{k+1} + \ldots$$

where $L_k$ is an appropriate simplification of jacobian $M_k'$.

It is then necessary, in order to ensure that the result of the adjoint integration is the exact gradient of the objective function, to modify the basic model in such a way that the solution emanating from $\xi_0^{(0)} + \delta\xi_0$ is equal to $\xi_k^{(0)} + \delta\xi_k$, where $\delta\xi_k$ evolves according to (2). This makes the basic dynamics exactly linear.

## Incremental Method (continuation 3)

As concerns the observation operators in the objective function, a similar procedure can be implemented if those operators are nonlinear. This leads to replacing $H_k(\xi_k)$ by $H_k(\xi_k^{(0)}) + N_k \delta\xi_k$, where $N_k$ is an appropriate 'simple' linear operator (possibly, but not necessarily, the jacobian of $H_k$ at point $\xi_k^{(0)}$). The objective function depends only on the initial $\delta\xi_0$ deviation from $\xi_0^{(0)}$, and reads

$$\mathcal{J}_I(\delta\xi_0) \; = \; (1/2)\,(x_0^b - \xi_0^{(0)} - \delta\xi_0)^T\,[P_0^b]^{-1}\,(x_0^b - \xi_0^{(0)} - \delta\xi_0)$$
$$+ (1/2)\,\Sigma_k[d_k - N_k\delta\xi_k]^T\,R_k^{-1}\,[d_k - N_k\delta\xi_k]$$

where $d_k \equiv y_k - H_k(\xi_k^{(0)})$ is the innovation at time $k$, and the $\delta\xi_k$ evolve according to

$$\delta\xi_{k+1} = L_k\,\delta\xi_k \qquad\qquad (2)$$

With the choices made here, $\mathcal{J}_I(\delta\xi_0)$ is an exactly quadratic function of $\delta\xi_0$. The minimizing perturbation $\delta\xi_{0,m}$ defines a new initial state $\xi_0^{(1)} \equiv \xi_0^{(0)} + \delta\xi_{0,m}$, from which a new solution $\xi_k^{(1)}$ of the basic nonlinear equation is determined. The process is restarted in the vicinity of that new solution.

## Incremental Method (continuation 4)

This defines a system of two-level nested loops for minimization. Advantage is that many degrees of freedom are available for defining the simplified operators $L_k$ and $N_k$, and for defining an appropriate trade-off between practical implementability and physical usefulness and accuracy. It is the incremental method which, together with the adjoint method, makes variational assimilation possible.

*First-Guess-At-the-right-Time 3D-Var* (*FGAT 3D-Var*). Corresponds to $L_k = I_n$. Assimilation is four-dimensional in that observations are compared to a first-guess which evolves in time, but is three-dimensional in that no dynamics other than the trivial dynamics expressed by the unit operator is present in the minimization.

Buehner *et al.* (*Mon. Wea. Rev.*, 2010)

For the same numerical cost, and in meteorologically realistic situations, Ensemble Kalman Filter and Variational Assimilation produce results of similar quality.

- Compléments sur l'Estimation Statistique (*BLUE*)

## *Best Linear Unbiased Estimate*

*State vector $x$*, belonging to *state space $S$* ($\dim S = n$), to be estimated.

Available data in the form of

- A '*background*' estimate (*e. g.* forecast from the past), belonging to *state space*, with dimension $n$

$$x^b = x + \zeta^b$$

- An additional set of data (*e. g.* observations), belonging to *observation space*, with dimension $p$

$$y = Hx + \varepsilon$$

$H$ is known linear *observation operator*.

Assume probability distribution is known for the couple $(\zeta^b, \varepsilon)$.

Assume $E(\zeta^b) = 0$, $E(\varepsilon) = 0$, $E(\zeta^b \varepsilon^T) = 0$ (not restrictive)

Set $E(\zeta^b \zeta^{bT}) \equiv P^b$ (also often denoted $B$), $E(\varepsilon \varepsilon^T) \equiv R$

*From course 3*

**Best Linear Unbiased Estimate**

$$x^a = x^b + P^b\,H^T\,[HP^bH^T + R]^{-1}\,(y - Hx^b)$$
$$P^a = P^b - P^b\,H^T\,[HP^bH^T + R]^{-1}\,HP^b$$

$x^a$ is the *Best Linear Unbiased Estimate* (*BLUE*) of $x$ from $x^b$ and $y$.

Equivalent set of formulæ

$$x^a = x^b + P^a\,H^T\,R^{-1}\,(y - Hx^b)$$
$$[P^a]^{-1} = [P^b]^{-1} + H^T\,R^{-1}H$$

Vector $d \equiv y - Hx^b$ is *innovation vector*
Matrix $K \equiv P^b\,H^T\,[HP^bH^T + R]^{-1} = P^a\,H^T\,R^{-1}$ is *gain matrix*.

If couple $(\zeta^b, \varepsilon)$ is Gaussian, *BLUE* achieves bayesian estimation, in the sense that $P(x \mid x^b, y) = \mathcal{N}[x^a, P^a]$.

Condition $E(\varepsilon\zeta^{bT}) = 0$ is not mathematically restrictive. Setting $E(\varepsilon\zeta^{bT}) \equiv D$ (possibly $\neq 0$), and coming back to the general formula

$$x^a = E(x) + C_{xy}\,[C_{yy}]^{-1}\,[y - E(y)]$$
$$P^a = C_{xx} - C_{xy}\,[C_{yy}]^{-1}\,C_{yx}$$

with again $x' = x - E(x) = -\zeta^b$, $E(\zeta^b\zeta^{bT}) = P^b$

$$y' = y - E(y) = y - Hx^b = \varepsilon - H\zeta^b$$

$$C_{xy} = E(x'y'^{T}) \;=\; E[-\zeta^b(\varepsilon - H\zeta^b)^{T}] \;=\; - \underset{D^{T}}{E(\zeta^b\varepsilon^{T})} + \underset{P^b}{E(\zeta^b\zeta^{bT})}H^{T} = -D^{T} + P^bH^{T}$$

$$C_{xy} = -D^{T} + P^bH^{T}$$

$$C_{yy} = E(y'y'^{T}) \;=\; E[(\varepsilon - H\zeta^b)(\varepsilon - H\zeta^b)^{T}]$$
$$= \underset{R}{E(\varepsilon\varepsilon^{T})} - \underset{D}{E(\varepsilon\zeta^{bT})}H^{T} - H\underset{D^{T.}}{E(\varepsilon\zeta^{bT})} + H\underset{P^b}{E(\zeta^b\zeta^{bT})}H^{T}$$

$$C_{yy} = R - DH^{T} - HD^{T} + HP^bH^{T}$$

Leading to expressions

$$x^a = x^b + [P^bH^T - D^T] [HP^bH^T - DH^T - HD^T + R]^{-1} (y - Hx^b)$$

$$P^a = P^b - [P^bH^T - D^T] [HP^bH^T - DH^T - HD^T + R]^{-1} [HP^b - D]$$

This is equivalent to replacing the observation vector $y$ with the vector $v \equiv y - D[P^b]^{-1}x^b$, the error of which is uncorrelated with $\zeta^b$, and then using the formulæ for the case of no correlation between background and observation errors.

But the hypothesis of no correlation is almost always made in practice, although it is certainly not always verified (observations performed by a same satellite instrument, which have been through a same post-processing, are very likely to have correlated errors).

Now, taking into account correlations between backgound and observation errors does not render, as shown in course 4, the corresponding estimate optimal. That would require to modify the weights that have been given to previous data.

**Bayesian Estimation**

**Data of the form**

$$z = \Gamma x + \zeta, \qquad\qquad \zeta \sim \mathcal{N}[0, S]$$

Known data vector $z$ belongs to *data space* $\mathcal{D}$, $\dim \mathcal{D} = m$,
Unknown state vector $x$ belongs to *state space* $X$, $\dim X = n$
$\Gamma$ known ($m\mathrm{x}n$)-matrix, $\zeta$ unknown 'error'

Probability that $x = \xi$ given in $X$ ?     $x = \xi \Rightarrow \zeta = z - \Gamma\xi$

$$P(\zeta = z - \Gamma\xi) \propto \exp[ -(z - \Gamma\xi)^{\mathrm{T}} S^{-1} (z - \Gamma\xi)/2 ] \propto \exp[ -(\xi - x^a)^{\mathrm{T}} (P^a)^{-1} (\xi - x^a)/2 ]$$

where

$$x^a = (\Gamma^{\mathrm{T}} S^{-1} \Gamma)^{-1} \Gamma^{\mathrm{T}} S^{-1} z$$
$$P^a = (\Gamma^{\mathrm{T}} S^{-1} \Gamma)^{-1}$$

Then conditional probability distribution is

$$P(x \mid z) = \mathcal{N}[x^a, P^a]$$

**Bayesian Estimation** (continuation 1)

$$z = \varGamma x + \zeta, \qquad \zeta \sim \mathcal{N}[0, S]$$

Then

$$P(x \mid z) = \mathcal{N}[x^a, P^a]$$

with

$$x^a = (\varGamma^T S^{-1} \varGamma)^{-1} \varGamma^T S^{-1} z$$
$$P^a = (\varGamma^T S^{-1} \varGamma)^{-1}$$

*Determinacy condition* : rank $\varGamma = n$. Data contain information, directly or indirectly, on every component of state vector $x$. Requires $m \geq n$.

**Variational form**

$$P(\boldsymbol{x} \mid \boldsymbol{z}) \propto \exp[\, -(\boldsymbol{z} - \boldsymbol{\Gamma}\boldsymbol{\xi})^{\mathrm{T}} \boldsymbol{S}^{-1} (\boldsymbol{z} - \boldsymbol{\Gamma}\boldsymbol{\xi})/2 \,] \propto \exp[\, -(\boldsymbol{\xi} - \boldsymbol{x}^a)^{\mathrm{T}} (\boldsymbol{P}^a)^{-1} (\boldsymbol{\xi} - \boldsymbol{x}^a)/2 \,]$$

Conditional expectation $\boldsymbol{x}^a$ minimizes following scalar *objective function*, defined on state space $\mathcal{X}$

$$\boldsymbol{\xi} \in \mathcal{X} \to \mathcal{J}(\boldsymbol{\xi}) \equiv (1/2) \, [\boldsymbol{\Gamma}\boldsymbol{\xi} - \boldsymbol{z})]^{\mathrm{T}} \boldsymbol{S}^{-1} [\boldsymbol{\Gamma}\boldsymbol{\xi} - \boldsymbol{z}]$$

$$\boldsymbol{P}^a = [\partial^2 \mathcal{J} / \partial \boldsymbol{\xi}^2]^{-1}$$

If data still of the form

$$z = \Gamma x + \zeta$$

but 'error' $\zeta$, which still has expectation $0$ and covariance $S$, is not Gaussian, expressions

$$x^a = (\Gamma^T S^{-1} \Gamma)^{-1} \Gamma^T S^{-1} z$$

$$P^a = (\Gamma^T S^{-1} \Gamma)^{-1}$$

do not achieve Bayesian estimation, but define least-variance linear estimate of $x$ from $z$ (*Best Linear Unbiased Estimator, BLUE*), and associated estimation error covariance matrix.

Expressions

$$x^a = (\boldsymbol{\Gamma}^T \boldsymbol{S}^{-1} \boldsymbol{\Gamma})^{-1} \boldsymbol{\Gamma}^T \boldsymbol{S}^{-1} \boldsymbol{z}$$

$$\boldsymbol{P}^a = (\boldsymbol{\Gamma}^T \boldsymbol{S}^{-1} \boldsymbol{\Gamma})^{-1}$$

are valid in both the Gaussian case and the general linear (*BLUE*) case. But, although, they are algebraically identical, they do not have the same significance. In the Gaussian case, as said, they solve entirely the problem of Bayesian estimation. For any data vector $\boldsymbol{z}$, $\boldsymbol{x}^a$ and $\boldsymbol{P}^a$ are respectively the expectation and covariance of the conditional (Gaussian) probability distribution $P(\boldsymbol{x} \mid \boldsymbol{z})$. In the general linear case, $\boldsymbol{x}^a$ and $\boldsymbol{P}^a$ have no necessary Bayesian meaning. In particular, for a given $\boldsymbol{z}$, $\boldsymbol{P}^a$, which is the covariance matrix of the estimation error over all possible realizations of $\boldsymbol{z}$ *(i.e.* of the error $\zeta$), can be very different from the corresponding Bayesian covariance matrix.

Expressions

$$x^a = (\Gamma^T S^{-1} \Gamma)^{-1} \Gamma^T S^{-1} z$$

$$P^a = (\Gamma^T S^{-1} \Gamma)^{-1}$$

are invariant in linear invertible change of coordinates, in either data or state space. If determinacy condition is verified, data vector $z$ can be transformed, through linear invertible change of coordinates in data space, into

$$x^b = x + \zeta^b$$

$$y = Hx + \varepsilon$$

from which the formulæ derived previously can be obtained (in both cases $E(\varepsilon\zeta^{bT}) = 0$ and $\neq 0$) .

Three sets of equivalent equations

$$z = \Gamma x + \zeta, \; E(\zeta) = 0, \; E(\zeta\zeta^{\mathrm{T}}) \equiv S^b$$

$$x^a = (\Gamma^{\mathrm{T}} S^{-1} \Gamma)^{-1} \, \Gamma^{\mathrm{T}} S^{-1} \, z$$

$$P^a = (\Gamma^{\mathrm{T}} S^{-1} \Gamma)^{-1}$$

$$x^b = x + \zeta^b, \; E(\zeta^b) = 0, \; E(\zeta^b\zeta^{b\mathrm{T}}) \equiv P^b,$$

$$y = Hx + \varepsilon, \qquad E(\varepsilon) = 0, \; E(\varepsilon\varepsilon^{\mathrm{T}}) \equiv R$$

$E(\zeta^b\varepsilon^{\mathrm{T}}) = 0$ (not restrictive)

$$x^a = x^b + P^b H^{\mathrm{T}} [HP^b H^{\mathrm{T}} + R]^{-1} (y - Hx^b)$$

$$P^a = P^b - P^b H^{\mathrm{T}} [HP^b H^{\mathrm{T}} + R]^{-1} HP^b$$

$$x^a = x^b + P^a H^{\mathrm{T}} R^{-1} (y - Hx^b)$$

$$[P^a]^{-1} = [P^b]^{-1} + H^{\mathrm{T}} R^{-1} H$$

The three equations have the same algebraic structure

$$z = \Gamma x + \zeta$$

$$x^a = (\Gamma^T S^{-1} \Gamma)^{-1} \Gamma^T S^{-1} z$$
$$P^a = (\Gamma^T S^{-1} \Gamma)^{-1}$$

$$x^b = x + \zeta^b, \quad y = Hx + \varepsilon$$

$$x^a = x^b + P^b H^T [HP^b H^T + R]^{-1} (y - Hx^b)$$
$$P^a = P^b - P^b H^T [HP^b H^T + R]^{-1} HP^b$$

$$x^a = x^b + P^a H^T R^{-1} (y - Hx^b)$$
$$[P^a]^{-1} = [P^b]^{-1} + H^T R^{-1} H$$

Data vector

$$z = \Gamma x + \zeta$$

Analysis $x^a$ (whatever its exact meaning) minimizes following scalar *objective function*, defined on state space $\mathcal{X}$

$$\xi \in \mathcal{X} \rightarrow \mathcal{J}(\xi) \equiv (1/2) [\Gamma\xi - z)]^T S^{-1} [\Gamma\xi - z]$$

where $S = E(\zeta\zeta^T)$ is covariance matrix of data error $\zeta$

(for example 4D-Var

$$\mathcal{J}(\xi_0) = (1/2) (x_0^b - \xi_0)^T [P_0^b]^{-1} (x_0^b - \xi_0) + (1/2) \Sigma_k [y_k - H_k\xi_k]^T R_k^{-1} [y_k - H_k\xi_k])$$

Consider quantity $\quad D = z_1^{\mathrm{T}} S^{-1} z_2 = z_1^{\mathrm{T}} [\mathrm{E}(\zeta\zeta^{\mathrm{T}})]^{-1} z_2$

where $z_1$ and $z_2$ are any two vectors in data space

Change of coordinates $\; z \equiv Tw$

$$\zeta = T\chi \quad \Rightarrow S = \mathrm{E}(\zeta\zeta^{\mathrm{T}}) = \mathrm{E}[T\chi(T\chi)^{\mathrm{T}}] = T\,\mathrm{E}(\chi\chi^{\mathrm{T}})T^{\mathrm{T}}$$

$$D = w_1^{\mathrm{T}} T^{\mathrm{T}} [T\,\mathrm{E}(\chi\chi^{\mathrm{T}})T^{\mathrm{T}}]^{-1} Tw_2 = w_1^{\mathrm{T}} T^{\mathrm{T}} T^{-\mathrm{T}} [\mathrm{E}(\chi\chi^{\mathrm{T}})]^{-1} T^{-1} Tw_2$$

$$D = w_1^{\mathrm{T}} [\mathrm{E}(\chi\chi^{\mathrm{T}})]^{-1} w_2$$

Expression $\quad\quad\quad D = z_1{}^\mathrm{T}\, S^{-1}\, z_2$

defines proper scalar product, and associated norm, on data space

Called **_Mahalanobis norm_**

Prasanta Chandra Mahalanobis (1893 -1972)

Minimizing objective function

$$\mathcal{J}(\xi) \equiv (1/2) \, [\boldsymbol{\Gamma}\boldsymbol{\xi} - \boldsymbol{z})]^{\mathrm{T}} \, \boldsymbol{S}^{-1} \, [\boldsymbol{\Gamma}\boldsymbol{\xi} - \boldsymbol{z}]$$

amounts to orthogonal projection onto space $\Gamma(\mathcal{X})$, followed by
inversion through $\boldsymbol{\Gamma}$ (generalized inverse)

*Gaussian variables*

Unidimensional

$$\mathcal{N}[m, a] \sim (2\pi a)^{-1/2} \exp\left[-(1/2a)(\xi - m)^2\right]$$

Dimension *n*

$$\mathcal{N}[\boldsymbol{m}, \boldsymbol{A}] \sim \left[(2\pi)^n \det\boldsymbol{A}\right]^{-1/2} \exp\left[-(1/2) \underbrace{(\boldsymbol{\xi}-\boldsymbol{m})^T \boldsymbol{A}^{-1}(\boldsymbol{\xi}-\boldsymbol{m})}_{\text{Mahalanobis norm}}\right]$$

*Mahalanobis norm*

## Entropy of a probability distribution

Probability distribution over domain described by coordinate $\xi$, with probability density $p(\xi)$. *Entropy*

$$S \equiv -\int p \ln p \, d\xi$$

Entropy of a probability distribution is a measure of the associated uncertainty. The larger the entropy, the larger the uncertainty. A uniform probability distribution over an interval of length $a$ has entropy $\ln a$, which tends to $-\infty$ as $a$ tends to zero. A one-dimensional Gaussian probability distribution with variance $s$ has entropy $\ln\sqrt{(2\pi e s)}$.

For given variance $s$, entropy is largest for the Gaussian distribution.

**Entropy of a probability distribution** (continuation)

Data of the form (see slide 51)

$$z = \mathit{\Gamma} x + \zeta$$

The knowledge of a probability distribution for $\zeta$ defines a conditional probability distribution $P(x|z)$ for $x$. Assuming that only the expectation and covariance matrix $S$ of $\zeta$ are known, for which distribution of $\zeta$ is the entropy of $P(x|z)$ largest ?

Response. The entropy of $P(x|z)$ is largest when $\zeta$ is Gaussian.

If the probability distribution for $\zeta$ is unknown, assuming that it is Gaussian is in a sense the 'least committing' choice.

# Cours à venir

~~Mercredi 2 avril~~

~~Vendredi 11 avril~~

~~Vendredi 18 avril~~

~~Mercredi 23 avril~~

~~Lundi 12 mai~~

Mercredi 28 mai

Mercredi 11 juin

Mercredi 18 juin