Université Pierre et Marie Curie Master de Sciences et Technologies Spécialité Océan, Atmosphère, Climat et Observations Spatiales

Année 2013-2014

Course Introduction to data assimilation

From numerical modelling to data assimilation

Olivier Talagrand

7 January 2014



.

Fig. 1: Members of day 7 forecast of 500 hPa geopotential height for the ensemble originated from 25 January 1993.



Figure 6 Hurricane Katrina mean-sea-level-pressure (MSLP) analysis for 12 UTC of 29 August 2005 and t+84h high-resolution and EPS forecasts started at 00 UTC of 26 August:

1st row: 1st panel: MSLP analysis for 12 UTC of 29 Aug 2nd panel: MSLP t+84h T_L511L60 forecast started at 00 UTC of 26 Aug 3rd panel: MSLP t+84h EPS-control T_L255L40 forecast started at 00 UTC of 26 Aug Other rows: 50 EPS-perturbed T_L255L40 forecast started at 00 UTC of 26 Aug.

The contour interval is 5 hPa, with shading patters for MSLP values lower than 990 hPa.

ECMWF, Technical Report 499, 2006

Why have meteorologists such difficulties in predicting the weather with any certainty? Why is it that showers and even storms seem to come by chance, so that many people think it is quite natural to pray for them, though they would consider it ridiculous to ask for an eclipse by prayer ? [...] a tenth of a degree more or less at any given point, and the cyclone will burst here and not there, and extend its ravages over districts that it would otherwise have spared. If they had been aware of this tenth of a degree, they could have known it beforehand, but the observations were neither sufficiently comprehensive nor sufficiently precise, and that is the reason why it all seems due to the intervention of chance.

> H. Poincaré, *Science et Méthode*, Paris, 1908 (translated Dover Publ., 1952)

ECMWF Data Coverage (All obs DA) - Synop-Ship-Metar 13/Nov/2011; 00 UTC Total number of obs = 31583



ECMWF Data Coverage (All obs DA) - Temp 13/Nov/2011; 00 UTC Total number of obs = 649



ECMWF Data Coverage (All obs DA) - Pilot-Profiler 13/Nov/2011; 00 UTC Total number of obs = 1592



ECMWF Data Coverage (All obs DA) - Aircraft 13/Nov/2011; 00 UTC Total number of obs = 50106



ECMWF Data Coverage (All obs DA) - AMSU-A 13/Nov/2011; 00 UTC Total number of obs = 607377



ECMWF Data Coverage (All obs DA) - AMV WV 13/Nov/2011; 00 UTC Total number of obs = 175647



ECMWF Data Coverage (All obs DA) - SCAT 13/Nov/2011; 00 UTC Total number of obs = 289170



ECMWF Data Coverage (All obs DA) - GPSRO 13/Nov/2011; 00 UTC Total number of obs = 48559



ECMWF Data Coverage (All obs DA) - Buoy 13/Nov/2011; 00 UTC Total number of obs = 8540



ECMWF Data Coverage (All obs DA) - OZONE 13/Nov/2011; 00 UTC Total number of obs = 81811



December 2007: Satellite data volumes used: around 18 millions per day



quantity of satellite data used per day at ECMWF

Value as of early 2013 : around 25 millions per day

- *Synoptic* observations (ground observations, radiosonde observations), performed simultaneously, by international agreement, in all meteorological stations around the world (00:00, 06:00, 12:00, 18:00 TU), and are in practice concentrated over continents.
- *Asynoptic* observations (satellites, aircraft), performed more or less continuously in time.
- *Direct* observations (temperature, pressure, horizontal components of the wind, moisture), which are local and bear on the variables used for for describing the flow in numerical models.
- *Indirect* observations (radiometric observations, ...), which bear on some more or less complex combination (most often, a one-dimensional spatial integral) of variables used for for describing the flow

$y = H(\mathbf{x})$

H : *observation operator* (for instance, radiative transfer equation)

Échantillonnage de la circulation océanique par les missions altimétriques sur 10 jours : combinaison Topex-Poséidon/ERS-1



S. Louvel, Doctoral Dissertation, 1999



FIG. 1 – Bassin méditerranéen occidental: réseau d'observation tomographique de l'expérience Thétis 2 et limites du domaine spatial utilisé pour les expériences numériques d'assimilation.

E. Rémy, Doctoral Dissertation, 1999

18

Physical laws governing the flow

Conservation of mass

 $D\rho/Dt + \rho \operatorname{div} U = 0$

- Conservation of energy $De/Dt - (p/\rho^2) D\rho/Dt = Q$
- Conservation of momentum $D\underline{U}/Dt + (1/\rho) \operatorname{grad} p - \underline{g} + 2 \underline{\Omega} \wedge \underline{U} = \underline{F}$
- Equation of state $f(p, \rho, e) = 0$ (for a perfect gas $p/\rho = rT$, $e = C_v T$)
- Conservation of mass of secondary components (water in the atmosphere, salt in the ocean, chemical species, ...)
 Dq/Dt + q div<u>U</u> = S

These physical laws must be expressed in practice in discretized (and necessarily imperfect) form, both in space and time \Rightarrow *numerical model* ¹⁹

Parlance of the trade :

- Adiabatic and inviscid, and therefore thermodynamically reversible, processes (everything except Q, <u>F</u> and <u>S</u>) make up 'dynamics'
- Processes described by terms Q, <u>F</u> and <u>S</u> make up 'physics'

- All presently existing numerical models are built on simplified forms of the general physical laws. Global numerical models, used either for large-scale meteorological prediction or for climate simulation, are at present built on the so-called *primitive equations*. Those equations rely on several approximations, the most important of which being the *hydrostatic approximation*, which expresses balance, in the vertical direction, of the gravity and pressure gradient forces. This forbids explicit description of thermal convection, which must be parameterized in some appropriate way.
- More and more *limited-area models* have been progressively developed. They require appropriate definition of lateral boundary conditions (not a simple problem). Most of them are non-hydrostatic, and therefore allow description of convection.

There exist at present two forms of discretization

- Gridpoint discretization
- (Semi-)spectral discretization (mostly for global models, and most often only in the horizontal direction)

Finite element discretization, which is very common in many forms of numerical modelling, is rarely used for modelling of the atmosphere. It is more frequently used for oceanic modelling, where it allows to take into account the complicated geometry of coast-lines.

Schematic of a gridpoint atmospheric model (L. Fairhead /LMD-CNRS)





The grids of two of the models of Météo-France (La Météorologie)

- In gridpoint models, meteorological fields are defined by values at the nodes of a the grid. Spatial and temporal derivatives are expressed by finite differences.
- In spectral models, fields are defined by the coefficients of their expansion along a prescribed set of basic functions. In the case of global meteorological models, those basic functions are the spherical harmonics (eigenfunctions of the laplacian at the surface of the sphere).



Linear operations, and in particular differentiation, are performed in spectral space, while nonlinear operations and 'physical' computations (advection, diabatic heating and cooling, ...) are performed in gridpoint physical space. This requires constant transformations from one space to the other, which are made possible at an acceptable cost through the systematic use of Fast Fourier Transforms.

For that reason, those models are called *semi-spectral*.

- In the parlance of the trade, one distinguishes two different parts in models. The 'dynamics' deals with the physically reversible processes (pressure forces, Coriolis force, advection, ...), while the 'physics' deals with physically irreversible processes, in particular the diabatic heating term Q in the energy equation, and also the parameterization of subgrid scales effects.
- Numerical schemes have been progressively developed and validated for the 'dynamics' component of models, which are by and large considered now to work satisfactorily (although regular improvements are still being made).

The situation is different as concerns 'physics', where many problems remain (as concerns for instance subgrid scales parameterization, the water cycle and the associated exchanges of energy, or the exchanges between the atmosphere and the underlying medium). 'Physics' as a whole remains the weaker point of models, and is still the object of active research.



European Centre for Medium-range Weather Forecasts (ECMWF, Reading, UK)

Since 26 January 2010

Horizontal spherical harmonics triangular truncation T1279 (horizontal resolution \approx 16 kilometres, but still hydrostatic)

91 levels on the vertical (0 - 80 km)

Dimension of state vector $n \approx 1.5 \ 10^9$

Timestep = 10 minutes

Sunday 23 September 2012 00UTC ©ECMWF Forecast t+144 VT: Saturday 29 September 2012 00UTC Surface: Mean sea level pressure / 850-hPa wind speed



Saturday 29 September 2012 00UTC ©ECMWF Analysis t+000 VT: Saturday 29 September 2012 00UTC Surface: Mean sea level pressure / 850-hPa wind speed





Saturday 29 September 2012 00UTC ©ECMWF Analysis t+000 VT: Saturday 29 September 2012 00UTC Surface: Mean sea level pressure / 850-hPa wind speed





Figure 2: 500 hPa geopotential height 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 2011–July 2012.

Persistence = 0; climatology = 50 at long range



Figure 11: WMO-exchanged scores from global forecast centres. RMS error over northern extratropics for 500 hPa geopotential height (top) and mean sea level pressure (bottom). In each

http://www.ecmwf.int/

publications/library/ ecpublications/_pdf/ tm/601-700/tm688.pdf







Figure 12: As Figure 11 for the southern hemisphere.
Anomaly correlation of ECMWF 500hPa height forecasts





Figure 10: Forecast performance in the tropics. Curves show the monthly average RMS vector wind errors at 200 hPa (top) and 850 hPa (bottom) for one-day (blue) and five-day (red) forecasts. 12-month moving average scores are also shown (in bold).



Figure 19: Verification of 2 m temperature forecasts against European SYNOP data on the GTS for 60-hour (night-time) and 72-hour (daytime) forecasts. Lower pair of curves shows bias, upper curves are standard deviation of error.

http://www.ecmwf.int/publications/library/ecpublications/_pdf/tm/601-700/tm688.pdf



Figure 8: Ensemble spread (standard deviation, dashed lines) and RMS error of ensemble-mean (solid lines) for winter 2011–2012 (upper figure in each panel), and differences of ensemble spread and RMS error of ensemble mean for last three winter seasons (lower figure in each panel, negative values indicate spread is too small); plots are for 500 hPa geopotential (top) and 850 hPa temperature (bottom) over the extratropical northern hemisphere for forecast days 1 to 15.

publications/library/ ecpublications/_pdf/ tm/601-700/tm688.pdf



Figure 9: CPRSS for 500 hPa height (top) and 850 hPa temperature (bottom) ensemble forecasts for winter (December-February) over the extratropical northern hemisphere. Skill from the ensemble day 1–15 forecasts is shown for winters 2011–12 (red), 2010–11 (blue), 2009–10 (green), 2008–09 (magenta) and 2007–08 (cyan). The ensemble only ran to ten days in 2005–06 (orange).

41





FIG. 3. Evolution of forecast errors from 1981 to 2012 for N.Hem (a and c) and S.Hem (b and d). Operational forecasts (blue) and ERA Interim (green). Note that before 1986 the operational analysis is used to verify the operational forecasts, after 1986 ERA Interim is used for the verification (with an overlap of 6 months present).

Remaining Problems

Mostly in the 'physics' of models (Q and F terms in basic equations)

- Water cycle (evaporation, condensation, influence on radiation absorbed or emitted by the atmosphere)

- Exchanges with ocean or continental surface (heat, water, momentum, ...)

- ...

Purpose of assimilation : reconstruct as accurately as possible the state of the atmospheric or oceanic flow, using all available appropriate information. The latter essentially consists of

- The observations proper, which vary in nature, resolution and accuracy, and are distributed more or less regularly in space and time.
- The physical laws governing the evolution of the flow, available in practice in the form of a discretized, and necessarily approximate, numerical model.
- 'Asymptotic' properties of the flow, such as, *e. g.*, geostrophic balance of middle latitudes. Although they basically are necessary consequences of the physical laws which govern the flow, these properties can usefully be explicitly introduced in the assimilation process.

Both observations and 'model' are affected with some uncertainty \Rightarrow uncertainty on the estimate.

For some reason, uncertainty is conveniently described by probability distributions (don't know too well why, but it works).

Assimilation is a problem in bayesian estimation.

Determine the conditional probability distribution for the state of the system, knowing everything we know

Assimilation is one of many '*inverse problems*' encountered in many fields of science and technology

- solid Earth geophysics
- plasma physics
- 'nondestructive' probing
- navigation (spacecraft, aircraft,)
- ...

Solution most often (if not always) based on Bayesian, or probabilistic, estimation. 'Equations' are fundamentally the same.

Difficulties specific to assimilation of meteorological observations :

- Very large numerical dimensions ($n \approx 10^{6}$ -10⁹ parameters to be estimated, $p \approx 1$ -3.10⁷ observations per 24-hour period). Difficulty aggravated in Numerical Weather Prediction by the need for the forecast to be ready in time.

- Non-trivial, actually chaotic, underlying dynamics

ratio of supercomputer costs: 1 day's assimilation / 1 day forecast



Relative cost of the various components of the operational prediction suite at ECMWF (september 2011, J.-N. Thépaut) :

4DVAR: 17% Ensemble Data Assimilation (EDA) : 15% Deterministic model : 13% Ensemble Prediction System (EPS) : 53% others : 2%

EDA produces both the background error covariances for 4D-Var and the initial perturbations (in addition to Singular Vectors) for EPS.

 $z_1 = x + \zeta_1$ density function $p_1(\zeta) \propto \exp[-(\zeta^2)/2s_1]$ $z_2 = x + \zeta_2$ density function $p_2(\zeta) \propto \exp[-(\zeta^2)/2s_2]$ ζ_1 and ζ_2 mutually independent

 $P(x = \xi \mid z_1, z_2)$?

 $x = \xi \iff \zeta_1 = z_1 - \xi \text{ and } \zeta_2 = z_2 - \xi$

• $P(x = \xi | z_1, z_2) \propto p_1(z_1 - \xi) p_2(z_2 - \xi)$ $\propto \exp[-(\xi - x^a)^2/2p^a]$

where $1/p^a = 1/s_1 + 1/s_2$, $x^a = p^a (z_1/s_1 + z_2/s_2)$

Conditional probability distribution of *x*, given z_1 and $z_2 : \mathcal{N}[x^a, p^a]$ $p^a < (s_1, s_2)$ independent of z_1 and z_2



Fig. 1.1: Prior pdf p(x) (dashed line), posterior pdf $p(x|y^o)$ (solid line), and Gaussian likelihood of observation $p(y^o|x)$ (dotted line), plotted against x for various values of y^o . (Adapted from Lorenc and Hammon 1988.)

$$z_1 = x + \zeta_1$$
$$z_2 = x + \zeta_2$$

Same as before, but ζ_1 and ζ_2 are now distributed according to exponential law with parameter *a*, *i*. *e*.

```
p(\zeta) \propto \exp[-|\zeta|/a]; \operatorname{Var}(\zeta) = 2a^2
```

Conditional probability density function is now uniform over interval $[z_1, z_2]$, exponential with parameter a/2 outside that interval

 $E(x \mid z_1, z_2) = (z_1 + z_2)/2$

Var $(x \mid z_1, z_2) = a^2 (2\delta^3/3 + \delta^2 + \delta + 1/2) / (1 + 2\delta)$, with $\delta = |z_1 - z_2| / (2a)$ Increases from $a^2/2$ to ∞ as δ increases from 0 to ∞ . Can be larger than variance $2a^2$ of original errors (probability 0.08)

(Entropy - *fplnp* always decreases in bayesian estimation)

Bayesian estimation

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

Data vector z, belonging to data space $\mathcal{D}(\dim \mathcal{D} = m)$, available.

 $z = F(x, \zeta) \tag{1}$

where ζ is a random element representing the uncertainty on the data (or, more precisely, on the link between the data and the unknown state vector).

For example

 $z = \Gamma x + \zeta$

Bayesian estimation (continued)

Probability that $x = \xi$ for given ξ ?

 $x = \xi \implies z = F(\xi, \zeta)$

$$P(x = \xi \mid z) = P[z = F(\xi, \zeta)] / \int_{\xi'} P[z = F(\xi', \zeta)]$$

Unambiguously defined iff, for any ζ , there is at most one x such that (1) is verified.

 \Leftrightarrow data contain information, either directly or indirectly, on any component of *x*. *Determinacy* condition.

Bayesian estimation is however impossible in its general theoretical form in meteorological or oceanographical practice because

- It is impossible to explicitly describe a probability distribution in a space with dimension even as low as $n \approx 10^3$, not to speak of the dimension $n \approx 10^{6-9}$ of present Numerical Weather Prediction models.
- Probability distribution of errors on data very poorly known (model errors in particular).

One has to restrict oneself to a much more modest goal. Two approaches exist at present

- Obtain some 'central' estimate of the conditional probability distribution (expectation, mode, ...), plus some estimate of the corresponding spread (standard deviations and a number of correlations).
- Produce an ensemble of estimates which are meant to sample the conditional probability distribution (dimension $N \approx O(10-100)$).



Figure 2. 500 mb height field produced by the operational analysis procedure of Direction de la Météorologie for 00 GMT, 26 April 1984. Units: dam, contour interval: 4 dam. The field has been truncated to the truncation of the model used for the experiments described in the article.

Courtier and Talagrand, QJRMS, 1987



Figure 1. Geographical distribution of the observations used for the assimilation experiments. (a): geopotential observations; (b): wind observations. At most of the points plotted, several observations were made at successive synoptic hours. On each of the two charts, the heavy line delineates the Aleutian depression (see Figure 2).



Figure 2. 500 mb height field produced by the operational analysis procedure of Direction de la Météorologie for 00 GMT, 26 April 1984. Units: dam, contour interval: 4 dam. The field has been truncated to the truncation of the model used for the experiments described in the article.



Figure 3. 500 mb height field produced for 00 GMT, 26 April 1984, by the variational analysis minimizing the distance function defined by Eqs. (1)-(2) over a 24-hour period. Units: dam; contour interval: 4 dam.

500-hPa geopotential field as determined by : (left) operational assimilation system of French Weather Service (3D, primitive equation) and (right) experimental variational system (2D, vorticity equation)

Courtier and Talagrand, QJRMS, 1987

Random vector $\mathbf{x} = (x_1, x_2, ..., x_n)^T = (x_i)$ (e. g. pressure, temperature, abundance of given chemical compound at *n* grid-points of a numerical model)

- Expectation $E(x) = [E(x_i)]$; centred vector x' = x E(x)
- Covariance matrix

$$E(\boldsymbol{x}^{\prime}\boldsymbol{x}^{\prime \mathrm{T}}) = [E(x_{i}^{\prime}x_{j}^{\prime})]$$

dimension nxn, symmetric non-negative (strictly definite positive except if linear relationship holds between the x_i 's with probability 1).

- Two random vectors
 - $\boldsymbol{x} = (x_1, x_2, \dots, x_n)^{\mathrm{T}}$ $\boldsymbol{y} = (y_1, y_2, \dots, y_p)^{\mathrm{T}}$

 $E(\mathbf{x}'\mathbf{y}'^{\mathrm{T}}) = E(x_i'y_i')$

dimension *nxp*

Random function $\Phi(\xi)$ (field of pressure, temperature, abundance of given chemical compound, ...; ξ is now spatial and/or temporal coordinate)

- Expectation $E[\Phi(\xi)]$; $\Phi'(\xi) = \Phi(\xi) E[\Phi(\xi)]$
- Variance $Var[\varphi(\xi)] = E\{[\varphi'(\xi)]^2\}$
- Covariance function

$$(\xi_{1}, \xi_{2}) \rightarrow C_{\Phi}(\xi_{1}, \xi_{2}) \equiv E[\Phi'(\xi_{1}) \Phi'(\xi_{2})]$$

Correlation function

•

 $Cor_{\varphi}(\xi_{1}, \xi_{2}) = E[\Phi'(\xi_{1}) \Phi'(\xi_{2})] / \{Var[\Phi(\xi_{1})] Var[\Phi(\xi_{2})]\}^{1/2}$



.: Isolines for the auto-correlations of the 500 mb geopotential between the station in Hannover and surrounding stations. From Bertoni and Lund (1963)



Isolines of the cross-correlation between the 500 mb geopotential in station 01 384 (R) and the surface pressure in surrounding stations.

After N. Gustafsson



After N. Gustafsson



After N. Gustafsson

Optimal Interpolation

Random field $\Phi(\xi)$ Observation network $\xi_1, \xi_2, ..., \xi_p$ For one particular realization of the field, observations

 $y_j = \Phi(\xi_j) + \varepsilon_j$, j = 1, ..., p, making up vector $\mathbf{y} = (y_j)$

Estimate $x = \Phi(\xi)$ at given point ξ , in the form

 $x^{a} = \alpha + \Sigma_{j} \beta_{j} y_{j} = \alpha + \beta^{T} y$, where $\beta = (\beta_{j})$

 α and the β_j 's being determined so as to minimize the expected quadratic estimation error $E[(x-x^a)^2]$

Optimal Interpolation (continued 1)

Solution

$$x^{a} = E(x) + E(x'\mathbf{y}'^{\mathrm{T}}) \left[E(\mathbf{y}'\mathbf{y}'^{\mathrm{T}})\right]^{-1} \left[\mathbf{y} - E(\mathbf{y})\right]$$

i.e.,
$$\boldsymbol{\beta} = [E(\mathbf{y}'\mathbf{y}'^{\mathrm{T}})]^{-1} E(x'\mathbf{y}')$$

 $\alpha = E(x) - \boldsymbol{\beta}^{\mathrm{T}} E(\mathbf{y})$

Estimate is unbiased $E(x-x^a) = 0$

Minimized quadratic estimation error

 $E[(x-x^{a})^{2}] = E(x'^{2}) - E(x'y'^{T}) [E(y'y'^{T})]^{-1} E(y'x')$

Estimation made in terms of deviations from expectations x' and y'.

Optimal Interpolation (continued 2)

 $x^{a} = E(x) + E(x'y'^{T}) [E(y'y'^{T})]^{-1} [y - E(y)]$ $y_{j} = \Phi(\xi_{j}) + \varepsilon_{j}$ $E(y_{j}'y_{k}') = E[\Phi'(\xi_{j}) + \varepsilon_{j}'][\Phi'(\xi_{k}) + \varepsilon_{k}']$

If observation errors ε_j are mutually uncorrelated, have common variance *s*, and are uncorrelated with field Φ , then

$$E(y_j'y_k') = C_{\Phi}(\xi_j, \xi_k) + s\delta_{jk}$$

and

$$E(x'y_j') = C_{\Phi}(\xi, \xi_j)$$









Optimal Interpolation (continued 3)

 $x^{a} = E(x) + E(x'\mathbf{y}'^{\mathrm{T}}) \left[E(\mathbf{y}'\mathbf{y}'^{\mathrm{T}})\right]^{-1} \left[\mathbf{y} - E(\mathbf{y})\right]$

Vector

 $\boldsymbol{\mu} = (\boldsymbol{\mu}_j) \equiv [E(\mathbf{y}'\mathbf{y}'^{\mathrm{T}})]^{-1} [\mathbf{y} - E(\mathbf{y})]$

is independent of variable to be estimated

 $x^{a} = E(x) + \sum_{i} \mu_{i} E(x'y_{i}')$

 $\Phi^{a}(\xi) = E[\Phi(\xi)] + \sum_{j} \mu_{j} E[\Phi'(\xi) y_{j}']$ $= E[\Phi(\xi)] + \sum_{j} \mu_{j} C_{\Phi}(\xi, \xi_{j})$

Correction made on background expectation is a linear combination of the p functions

 $E[\Phi'(\xi) y_i'] \cdot E[\Phi'(\xi) y_i'] [= C_{\Phi}(\xi, \xi_i)]$

considered as a function of estimation position ξ , is the *representer* associated with observation y_i .
Optimal Interpolation (continued 4)

Univariate interpolation. Each physical field (*e. g.* temperature) determined from observations of that field only.

Multivariate interpolation. Observations of different physical fields are used simultaneously. Requires specification of cross-covariances between various fields.

Cross-covariances between mass and velocity fields can simply be modelled on the basis of geostrophic balance.

Cross-covariances between humidity and temperature (and other) fields still a problem.

Schlatter's (1975) multivariate covariances

Specified as multivariate 2-point functions.

Not easy to ensure that specified functions are actually valid covariances.

Used in OI and related observation-space methods.

Courtesy A. Lorenc

© Crown copyright Met Office Andre



FIG. 3. Correlations among the variables h, u, and v based upon the expression $\mu = 0.95 \exp(-1.24s^2)$ for height-height correlation and the geostrophic relations. Diagrams centered at 110° W, 35° N. Tick marks 500 km apart.

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 (ζ^b, ε) . Assume $E(\zeta^b) = 0$, $E(\varepsilon) = 0$, $E(\zeta^b \varepsilon^T) = 0$ (not restrictive) Set $E(\zeta^b \zeta^{b_T}) = P^b$ (also often denoted *B*), $E(\varepsilon \varepsilon^T) = R$ Best Linear Unbiased Estimate (continuation 1)

$$\mathbf{x}^b = \mathbf{x} + \boldsymbol{\zeta}^b \tag{1}$$

$$\mathbf{y} = H\mathbf{x} + \boldsymbol{\varepsilon} \tag{2}$$

A probability distribution being known for the couple $(\boldsymbol{\zeta}^b, \boldsymbol{\varepsilon})$, eqs (1-2) define probability distribution for the couple $(\boldsymbol{x}, \boldsymbol{y})$, with

 $E(\mathbf{x}) = \mathbf{x}^b$, $\mathbf{x}' = \mathbf{x} - E(\mathbf{x}) = -\boldsymbol{\zeta}^b$

 $E(\mathbf{y}) = H\mathbf{x}^b$, $\mathbf{y}' = \mathbf{y} - E(\mathbf{y}) = \mathbf{y} - H\mathbf{x}^b = \boldsymbol{\varepsilon} - H\boldsymbol{\zeta}^b$

 $d = y - Hx^b$ is called the *innovation vector*.

Best Linear Unbiased Estimate (continuation 2)

Apply formulæ for Optimal Interpolation

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

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

Equivalent set of formulæ

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

Matrix $K = P^b H^T [HP^b H^T + R]^{-1} = P^a H^T R^{-1}$ is gain matrix.

If probability distributions are *globally* gaussian, *BLUE* achieves bayesian estimation, in the sense that $P(x | x^b, y) = \mathcal{N}[x^a, P^a]$.





After A. Lorenc

Best Linear Unbiased Estimate (continuation 4)

Variational form of the *BLUE*

BLUE x^a minimizes following scalar objective function, defined on state space

$$\begin{split} \boldsymbol{\xi} \in \boldsymbol{S} \rightarrow \\ \boldsymbol{\mathcal{J}}(\boldsymbol{\xi}) &= (1/2) \left(\boldsymbol{x}^{b} - \boldsymbol{\xi} \right)^{\mathrm{T}} [P^{b}]^{-1} \left(\boldsymbol{x}^{b} - \boldsymbol{\xi} \right) + (1/2) \left(\boldsymbol{y} - H\boldsymbol{\xi} \right)^{\mathrm{T}} R^{-1} \left(\boldsymbol{y} - H\boldsymbol{\xi} \right) \\ &= \boldsymbol{\mathcal{J}}_{b} \qquad + \boldsymbol{\mathcal{J}}_{o} \end{split}$$
$$\boldsymbol{\mathcal{S}} \boldsymbol{D} - \boldsymbol{V} \boldsymbol{a} \boldsymbol{r}^{\prime} \end{split}$$

Can easily, and heuristically, be extended to the case of a nonlinear observation operator H.

Used operationally in USA, Australia, China, ...

Question. How to introduce temporal dimension in estimation process ?

- Logic of Optimal Interpolation can be extended to time dimension.
- But we know much more than just temporal correlations. We know explicit dynamics.

Real (unknown) state vector at time k (in format of assimilating model) x_k . Belongs to state space $S(\dim S = n)$

Evolution equation

 $x_{k+1} = M_k(x_k) + \eta_k$

 M_k is (known) model, η_k is (unknown) model error

Sequential Assimilation

• Assimilating model is integrated over period of time over which observations are available. Whenever model time reaches an instant at which observations are available, state predicted by the model is updated with new observations.

Variational Assimilation

• Assimilating model is globally adjusted to observations distributed over observation period. Achieved by minimization of an appropriate scalar *objective function* measuring misfit between data and sequence of model states to be estimated.

• Observation vector at time *k*

$$y_{k} = H_{k}x_{k} + \varepsilon_{k}$$

$$E(\varepsilon_{k}) = 0 \quad ; \quad E(\varepsilon_{k}\varepsilon_{j}^{T}) = R_{k} \delta_{kj}$$

$$H_{k} \text{ linear}$$

Evolution equation

 $x_{k+1} = M_k x_k + \eta_k$ $E(\eta_k) = 0 \quad ; \quad E(\eta_k \eta_j^{\mathrm{T}}) = Q_k \delta_{kj}$ $M_k \text{ linear}$

$$k=0,\ldots,K-1$$

 $k = 0, \ldots, K$

• $E(\eta_k \varepsilon_j^{\mathrm{T}}) = 0$ (errors uncorrelated in time)

At time k, background x_k^b and associated error covariance matrix P_k^b known

Analysis step

$$x_{k}^{a} = x_{k}^{b} + P_{k}^{b} H_{k}^{T} [H_{k} P_{k}^{b} H_{k}^{T} + R_{k}]^{-1} (y_{k} - H_{k} x_{k}^{b})$$

$$P_{k}^{a} = P_{k}^{b} - P_{k}^{b} H_{k}^{T} [H_{k} P_{k}^{b} H_{k}^{T} + R_{k}]^{-1} H_{k} P_{k}^{b}$$

• Forecast step

$$\begin{aligned} x^{b}_{k+1} &= M_{k} x^{a}_{k} \\ P^{b}_{k+1} &= E[(x^{b}_{k+1} - x_{k+1})(x^{b}_{k+1} - x_{k+1})^{\mathrm{T}}] = E[(M_{k} x^{a}_{k} - M_{k} x_{k} - \eta_{k})(M_{k} x^{a}_{k} - M_{k} x_{k} - \eta_{k})^{\mathrm{T}}] \\ &= M_{k} E[(x^{a}_{k} - x_{k})(x^{a}_{k} - x_{k})^{\mathrm{T}}]M_{k}^{\mathrm{T}} - E[\eta_{k} (x^{a}_{k} - x_{k})^{\mathrm{T}}] - E[(x^{a}_{k} - x_{k})\eta_{k}^{\mathrm{T}}] + E[\eta_{k} \eta_{k}^{\mathrm{T}}] \\ &= M_{k} P^{a}_{k} M_{k}^{\mathrm{T}} + Q_{k} \end{aligned}$$

At time k, background x_k^b and associated error covariance matrix P_k^b known

Analysis step

 $x^{a}_{\ k} = x^{b}_{\ k} + P^{b}_{\ k} H^{T}_{k} [H_{k} P^{b}_{\ k} H^{T}_{k} + R_{k}]^{-1} (y_{k} - H_{k} x^{b}_{\ k})$ $P^{a}_{\ k} = P^{b}_{\ k} - P^{b}_{\ k} H^{T}_{k} [H_{k} P^{b}_{\ k} H^{T}_{k} + R_{k}]^{-1} H_{k} P^{b}_{\ k}$

Forecast step

$$x^{b}_{k+1} = M_k x^{a}_k$$
$$P^{b}_{k+1} = M_k P^{a}_k M_k^{\mathrm{T}} + Q_k$$

Kalman filter (KF, Kalman, 1960)

Must be started from some initial estimate (x_0^b, P_0^b)

If all operators are linear, and if errors are uncorrelated in time, Kalman filter produces at time k the *BLUE* x_k^b (resp. x_k^a) of the real state x_k from all data prior to (resp. up to) time k, plus the associated estimation error covariance matrix P_k^b (resp. P_k^a).

If in addition errors are gaussian, the corresponding conditional probability distributions are the respective gaussian distributions $\mathcal{N}[x^b_k, P^b_k]$ and $\mathcal{N}[x^a_k, P^a_k]$.

A didactie example (Ghiletal)
Basotropic model
$$\begin{cases} \frac{\partial q}{\partial t} + div(q \xi) = c\\ \frac{\partial y}{\partial t} + gad(q + iv) + kx(g + s) \xi = c \end{cases}$$

One dimension, periodic
Une dimension, periodic
Lineari Continent





Fig. 2

The components of the total expected rms error (Erms), (trace: $P_{\rm b}$)^{1/2}, in the estimation of solutions to the stochastic-dynamic system (Y,H), with Y given by (3.6) and H = (I 0). System noise is absent, Q = 0. The filter used is the standard K-B filter (2.11) for the model.

a) Erms over land; b) Erms over the ocean; c) Erms over the entire L-domain

In each one of the figures, each curve represents one component of the total Erms error. The curves labelled U, V, and P represent the u component, v component and \$ component, respectively. They are found by summing the diagonal elements of Pk which correspond to u, v, and \$, respectively, dividing by the number of terms in the sum, and then taking the square root. In a) the summation extends over land points only, in b) over ocean points only, and in c) over the entire L-domain. The vertical axis is scaled in such a way that 1.0 corresponds to an Erms error of vmax for the U and V curves, and of \$0 for the P curve. The observational error level is 0.089 for the U and V curves, and 0.080 for the P curve. The curves labelled T represent the total Erms error over each region. Each T curve is a weighted average of the corresponding U, V, and P curves, with the weights chosen in such a way that the T curve measures the error in the total energy $u^2 + v^2 + \frac{1}{\sqrt{2}}$, conserved by the system (3.1). The observational noise level for the T curve is then 0.088. Notice the immediate error decrease over land and the gradual decrease over the ocean. The total estimation error tends to zero.

M. Ghil *et al*.



M. Ghil *et al*.

Fig. 6 This figure and the following ones show the properties of the estimated algorithm (2.11) in the presence of system noise, Q ≠ 0. This figure gives the Erms estimation error, and is homologous to Fig. 2. Notice the sharper increase of error over land between synoptic times, and the convergence of each curve to a periodic, nonzero function.

Nonlinearities ?

Model is usually nonlinear, and observation operators (satellite observations) tend more and more to be nonlinear.

Analysis step

 $x_{k}^{a} = x_{k}^{b} + P_{k}^{b} H_{k}^{'T} [H_{k}^{'} P_{k}^{b} H_{k}^{'T} + R_{k}]^{-1} [y_{k} - H_{k}(x_{k}^{b})]$ $P_{k}^{a} = P_{k}^{b} - P_{k}^{b} H_{k}^{'T} [H_{k}^{'} P_{k}^{b} H_{k}^{'T} + R_{k}]^{-1} H_{k}^{'} P_{k}^{b}$

Forecast step

Extended Kalman Filter (EKF, heuristic !)

Costliest part of computation

 $P^b_{k+1} = M_k P^a_{\ k} M_k^{\ \mathrm{T}} + Q_k$

Multiplication by M_k = one integration of the model between times k and k+1. Computation of $M_k P^a_{\ k} M_k^{\ T} \approx 2n$ integrations of the model

Need for determining the temporal evolution of the uncertainty on the state of the system is the major difficulty in assimilation of meteorological and oceanographical observations



Analysis of 500-hPa geopotential for 1 December 1989, 00:00 UTC (ECMWF, spectral truncation T21, unit *m*. After F. Bouttier)



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.

Two solutions :

• Low-rank filters (Heemink, Pham, ...)

Reduced Rank Square Root Filters, Singular Evolutive Extended Kalman Filter,

• Ensemble filters (Evensen, Anderson, Kalnay, ...)

Uncertainty is represented, not by a covariance matrix, but by an ensemble of point estimates in state space which are meant to sample the conditional probability distribution for the state of the system (dimension $N \approx O(10-100)$).

Ensemble is evolved in time through the full model, which eliminates any need for linear hypothesis as to the temporal evolution. How to update predicted ensemble with new observations ?

Predicted ensemble at time $t : \{x_i^b\}, \quad i = 1, ..., N$ Observation vector at same time : $y = Hx + \varepsilon$

• Gaussian approach

Produce sample of probability distribution for real observed quantity Hx $y_i = y - \varepsilon_i$ where ε_i is distributed according to probability distribution for observation error ε .

Then use Kalman formula to produce sample of 'analysed' states

 $x_{i}^{a} = x_{i}^{b} + P^{b} H^{T} [HP^{b}H^{T} + R]^{-1} (y_{i} - Hx_{i}^{b}), \qquad i = 1, ..., N$ (2)

where P^b is covariance matrix of predicted ensemble $\{x_i^b\}$.

Remark. If P^b was exact covariance matrix of background error, (2) would achieve Bayesian estimation, in the sense that $\{x_i^a\}$ would be a sample of conditional probability distribution for x, given all data up to time t.

Called *Ensemble Kalman Filter*. Has now become one of the two major powerful algorithms for assimilation of meteorological and oceanographical observations.

Local Ensemble Transform Kalman Filter (LETKF, Kalnay and colleagues)

Month-long Performance of EnKF vs. 3Dvar with WRF



Better performance of EnKF than 3DVar also seen in both 12-h forecast and posterior analysis in terms of root-mean square difference averaged over the entire month

(Meng and Zhang 2007c, MWR, in review)

Variational Assimilation

• Observation vector at time *k*

 $y_k = H_k x_k + \varepsilon_k$ $E(\varepsilon_k) = 0 \quad ; \quad E(\varepsilon_k \varepsilon_j^{\mathrm{T}}) = R_k \, \delta_{kj}$

• Evolution equation

 $\begin{aligned} x_{k+1} &= M_k x_k + \eta_k \\ E(\eta_k) &= 0 \quad ; \quad E(\eta_k \eta_j^{\mathrm{T}}) = Q_k \, \delta_{kj} \\ E(\eta_k \varepsilon_j^{\mathrm{T}}) &= 0 \end{aligned}$

• Background estimate at time 0

 $\begin{aligned} x^{b}{}_{0} &= x_{0} + \zeta^{b}{}_{0} \\ E(\zeta^{b}{}_{0}) &= 0 \quad ; \ E(\zeta^{b}{}_{0} \zeta^{b}{}_{0}{}^{\mathrm{T}}) \equiv P^{b}{}_{0} \\ E(\zeta^{b}{}_{0}\varepsilon_{k}{}^{\mathrm{T}}) &= 0 \quad ; \ E(\zeta^{b}{}_{0}\eta_{k}{}^{\mathrm{T}}) = 0 \end{aligned}$

• Errors uncorrelated in time

k = 0, ..., K - 1

 $k = 0, \ldots, K$

Variational assimilation leads to the following weak constraint objective function

$$\begin{aligned} (\xi_0, \, \xi_1, \, \dots, \, \xi_K) &\to \\ \mathcal{J}(\xi_0, \, \xi_1, \, \dots, \, \xi_K) \\ &= (1/2) \, (x^b_{\ 0} - \, \xi_0)^{\mathrm{T}} \, [P^b_{\ 0}]^{-1} \, (x^b_{\ 0} - \, \xi_0) \\ &+ (1/2) \, \Sigma_{k=0, \, \dots, \, K} \, (y_k - H_k \xi_k)^{\mathrm{T}} \, R_k^{-1} \, (y_k - H_k \xi_k) \\ &+ (1/2) \, \Sigma_{k=0, \, \dots, \, K-1} \, (\xi_{k+1} - M_k \xi_k)^{\mathrm{T}} \, Q_k^{-1} \, (\xi_{k+1} - M_k \xi_k) \end{aligned}$$

If model error is ignored ($Q_k=0$), problem reduces to minimizing

- $\xi_0 \rightarrow \mathcal{J}(\xi_0) = (1/2) (x_0^{b_0} \xi_0)^{\mathrm{T}} [P_0^{b_0}]^{-1} (x_0^{b_0} \xi_0)$ + $(1/2) \sum_{k=0, ..., K} (y_k - H_k \xi_k)^{\mathrm{T}} R_k^{-1} (y_k - H_k \xi_k)$ subject to
 - $\xi_{k+1} = M_k \xi_k \qquad , \qquad k = 0, \ldots, K-1$

Strong constraint four-dimensional variational assimilation, or strong constraint *4D-Var*

Used operationally in several meteorological centres (Météo-France, UK Meteorological Office, Canadian Meteorological Centre (maybe not any more ?), Japan Meteorological Agency, ...) and, until recently, at ECMWF. The latter now has a 'weak constraint' component in its operational system.

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

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

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

Gradient computed by *adjoint method*, which proceeds, in the space of partial derivatives, in reverse order of direct computations.

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

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

Adjoint Method

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

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

$$u \rightarrow v = G(u)$$

- $\mathbf{v} = (v_i)$ is output vector, $\dim \mathbf{v} = m$
- Perturbation $\delta u = (\delta u_i)$ of input. Resulting first-order perturbation on v

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

• or, in matrix form

- $\delta v = G' \delta u$
- where $G' = (\partial v_j / \partial u_i)$ is local matrix of partial derivatives, or jacobian matrix, of G.

Adjoint Method (continued 1)

$$\delta v = G' \delta u \tag{D}$$

• Scalar function of output

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

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

'Chain rule'

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

or

•
$$\nabla_{\boldsymbol{u}} \mathcal{J} = \boldsymbol{G}^{\mathsf{T}} \nabla_{\boldsymbol{v}} \mathcal{J}$$
 (A)

Adjoint Method (continued 2)

• *G* is the composition of a number of successive steps

•

$$\boldsymbol{G} = \boldsymbol{G}_N \circ \ldots \circ \boldsymbol{G}_2 \circ \boldsymbol{G}_1$$

'Chain rule'

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

Transpose

 $G'^{\rm T} = G_1'^{\rm T} G_2'^{\rm T} \dots G_N'^{\rm T}$

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).

Adjoint Approach

 $\mathcal{J}(\xi_0) = (1/2) (x_0^b - \xi_0)^T [P_0^b]^{-1} (x_0^b - \xi_0) + (1/2) \Sigma_k [y_k - H_k \xi_k]^T R_k^{-1} [y_k - H_k \xi_k]$ subject to $\xi_{k+1} = M_k \xi_k$, k = 0, ..., K-1

Control variable $\xi_0 = u$

Adjoint equation

 $\lambda_{K} = H_{K}^{T} R_{K}^{-1} [H_{K} \xi_{K} - y_{K}]$ $\lambda_{k} = M_{k}^{T} \lambda_{k+1} + H_{k}^{T} R_{k}^{-1} [H_{k} \xi_{k} - y_{k}]$ k = K-1, ..., 1 $\lambda_{0} = M_{0}^{T} \lambda_{1} + H_{0}^{T} R_{0}^{-1} [H_{0} \xi_{0} - y_{0}] + [P_{0}^{b}]^{-1} (\xi_{0} - x_{0}^{b})$ $\nabla_{\mu} \mathcal{J} = \lambda_{0}$

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

Adjoint Approach (continued 2)

Nonlinearities ?

 $\mathcal{J}(\xi_0) = (1/2) (x_0^{\ b} - \xi_0)^{\mathrm{T}} [P_0^{\ b}]^{-1} (x_0^{\ b} - \xi_0) + (1/2) \sum_k [y_k - H_k(\xi_k)]^{\mathrm{T}} R_k^{-1} [y_k - H_k(\xi_k)]$ subject to $\xi_{k+1} = M_k(\xi_k)$, $k = 0, \dots, K-1$

Control variable $\xi_0 = u$

Adjoint equation

 $\lambda_{K} = H_{K}^{T} R_{K}^{-1} [H_{K}(\xi_{K}) - y_{K}]$ $\lambda_{k} = M_{k}^{T} \lambda_{k+1} + H_{k}^{T} R_{k}^{-1} [H_{k}(\xi_{k}) - y_{k}]$ k = K-1, ..., 1 $\lambda_{0} = M_{0}^{T} \lambda_{1} + H_{0}^{T} R_{0}^{-1} [H_{0}(\xi_{0}) - y_{0}] + [P_{0}^{b}]^{-1} (\xi_{0} - x_{0}^{b})$

$$\nabla_{u}\mathcal{J} = \lambda_{0}$$

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



Buehner (2008)

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

Assimilation, which originated from the need of defining initial conditions for numerical weather forecasts, has progressively extended to many diverse applications

- Oceanography
- Atmospheric chemistry (both troposphere and stratosphere)
- Oceanic biogeochemistry
- Ground hydrology
- Terrestrial biosphere and vegetation cover
- Glaciology
- Magnetism (both planetary and stellar)
- Plate tectonics
- Planetary atmospheres (Mars, ...)
- Reassimilation of past observations (mostly for climatological purposes, ECMWF, NCEP/NCAR)
- Identification of source of tracers
- Parameter identification
- A priori evaluation of anticipated new instruments
- Definition of observing systems (Observing Systems Simulation Experiments)
- Validation of models
- Sensitivity studies (adjoints)
- ...

It has now become a major tool of numerical environmental science

Assimilation is related to

- Estimation theory
- Probability theory
- Atmospheric and oceanic dynamics
- Atmospheric and oceanic predictability
- Instrumental physics
- Optimisation theory
- Control theory
- Algorithmics and computer science
- ...

