

École Doctorale des Sciences de l'Environnement d'Île-de-France

Année Universitaire 2019-2020

Modélisation Numérique
de l'Écoulement Atmosphérique
et Assimilation de Données

Olivier Talagrand

Cours 3

2 Avril 2020

- Bayesian estimation. Continuation.
- Reminder on elementary probability theory. Random vectors and covariance matrices, random functions and covariance functions
- ‘Optimal Interpolation’

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; see, *e.g.* Jaynes, 2007, *Probability Theory: The Logic of Science*, Cambridge University Press).

Assimilation is a problem in bayesian estimation.

Determine the conditional probability distribution for the state of the system, knowing everything we know (see Tarantola, A., 2005, *Inverse Problem Theory and Methods for Model Parameter Estimation*, SIAM).

Bayesian estimation

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

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

$$z = F(x, \zeta) \quad (1)$$

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 ?

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

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

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 (the *curse of dimensionality*).
- 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 $N \approx O(10-100)$).

Coût des différentes composantes de la chaîne de prévision opérationnelle du CEPMMT (septembre 2015, J.-N. Thépaut) :

4DVAR: 9.5%

HRES FC: 4.5%

EDA: 30%

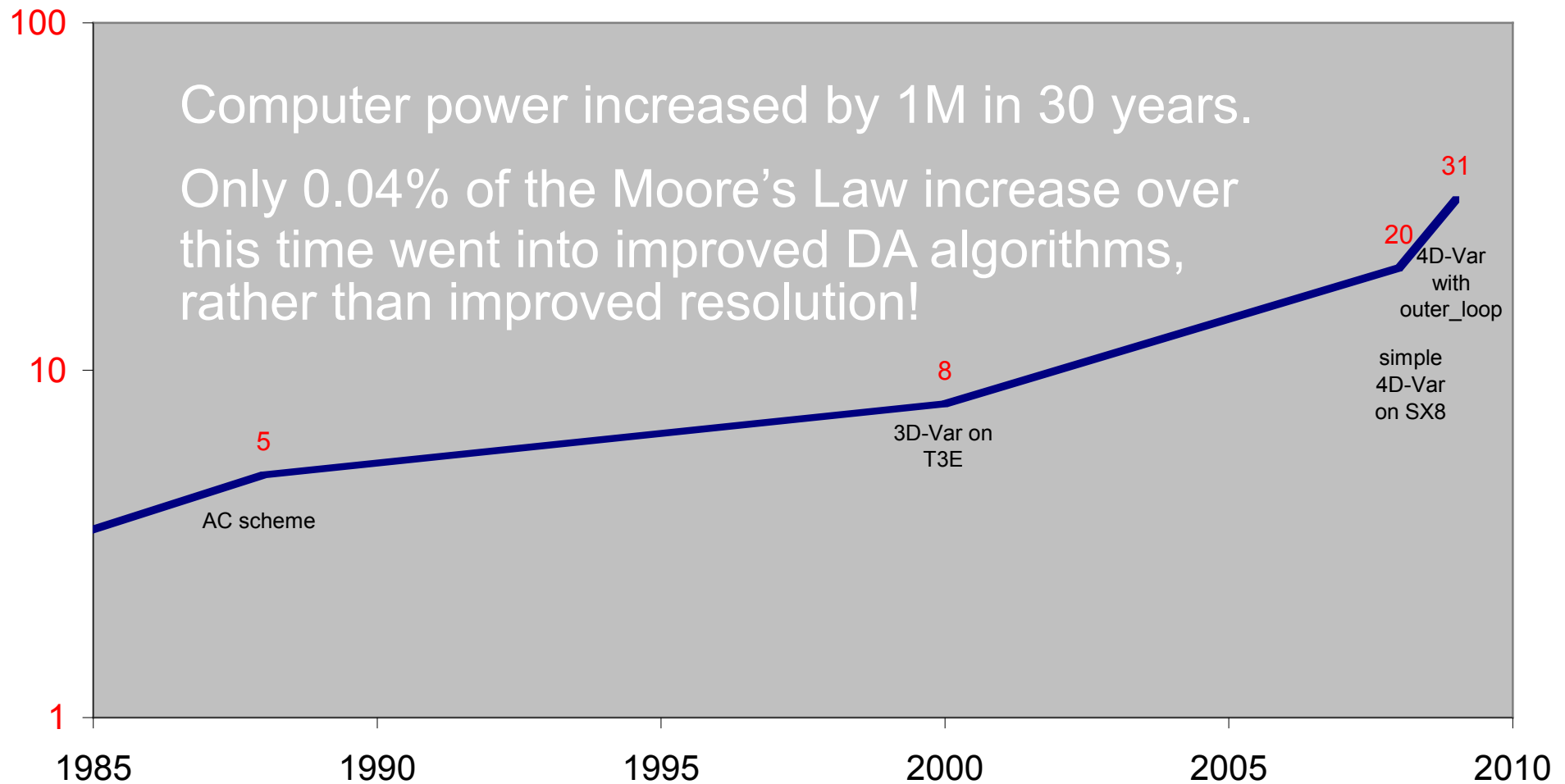
ENS: 22%

ENS: hindcasts 14%

Other: 20% of which BC AN: 3.5% BC FC: 4% BC ENS: 9.5%

L'EDA fournit à la fois les variances d'erreur d'ébauche du 4D-Var, et les perturbations initiales (en complément des vecteurs singuliers) de l'EPS.

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



Courtesy A. Lorenc

Scalar random variable x

Observed outcome of ‘realizations’ of a process that is repeated a large number of times. And also, *a priori* uncertainty on that result.

For any interval $[a, b]$, the probability $P(a < x < b)$ is known (whether inequalities are strict or not may matter).

Probability density function (pdf). Function $p(\xi)$ such that, for any interval $[a, b]$

$$P[a < x < b] = \int_a^b p(\xi) d\xi$$

($p(\xi)$ may contain diracs)

Expectation. Mean of a large number of realizations of x

$$E(x) = \int_{-\infty}^{+\infty} \xi p(\xi) d\xi$$

(may not exist)

Scalar random variable x (continued)

Variance

$$\text{Var}(x) \equiv E\{[x - E(x)]^2\} = E(x^2) - [E(x)]^2$$

Standard deviation

$$\sigma(x) \equiv \sqrt{\text{Var}(x)}$$

Centred variable $x' \equiv x - E(x)$

Couple of random variables $\mathbf{x} = (x_1, x_2)^T$

For any intervals $[a_1, b_1]$, $[a_2, b_2]$, probability $P(a_1 < x_1 < b_1 \text{ and } a_2 < x_2 < b_2)$ is known

Extends to any measurable domain $\mathcal{D} \subset \mathbb{R}^2$

$$P[(x_1, x_2) \in D] = \int_D p(\xi_1, \xi_2) d\xi_1 d\xi_2$$

where $p(\xi_1, \xi_2)$ is probability density function

Covariance

$$\text{Cov}(x_1, x_2) \equiv E(x_1' x_2')$$

$$\text{Corr}(x_1, x_2) \equiv \text{Cov}(x_1, x_2) / (\sigma(x_1) \sigma(x_2)) = \cos \varphi$$

Covariance is a scalar product, and defines Euclidean geometry (on space of finite-variance random variables on a given trial space)

Modulus = standard deviation σ , angle = $\cos^{-1}(\text{Corr})$, orthogonality = decorrelation

Couple of random variables $\mathbf{x} = (x_1, x_2)^T$ (continued)

Independence

x_1 and x_2 independent : knowledge about either one of the variables brings no knowledge about the other one.

For any intervals $[a_1, b_1], [a_2, b_2]$

$$P(a_1 < x_1 < b_1 \text{ and } a_2 < x_2 < b_2) = P(a_1 < x_1 < b_1) P(a_2 < x_2 < b_2)$$

Equivalently, pdf's verify

$$p(\xi_1, \xi_2) = p_1(\xi_1) p_2(\xi_2)$$

Independence implies decorrelation. Converse is not true

(consider $S = \sin \alpha$, $C = \cos \alpha$, where α is uniformly distributed over $[0, 2\pi]$)

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

- Expectation $E(\mathbf{x}) \equiv [E(x_i)]$; centred vector $\mathbf{x}' \equiv \mathbf{x} - E(\mathbf{x})$
- Covariance matrix

$$E(\mathbf{x}'\mathbf{x}'^T) = [E(x_i'x_j')]$$

dimension $n \times n$

Non-random vector $\boldsymbol{\lambda} = (\lambda_i)_{i=1, \dots, n}$

$$G \equiv \sum_i \lambda_i x_i'$$

$$G^2 = \sum_{i,j} \lambda_i \lambda_j x_i' x_j'$$

$$E(G^2) = \sum_{i,j} \lambda_i \lambda_j E(x_i' x_j') = \boldsymbol{\lambda}^T E(\mathbf{x}'\mathbf{x}'^T) \boldsymbol{\lambda} \geq 0$$

Covariance matrix $E(\mathbf{x}'\mathbf{x}'^T)$ is symmetric non negative (strictly definite positive except if linear relationship holds between the x_i' 's with probability 1).

Change

$$\mathbf{x} \rightarrow \mathbf{y} \equiv P\mathbf{x}$$

$$\mathbf{y}'\mathbf{y}'^T = P\mathbf{x}'(P\mathbf{x}')^T = P\mathbf{x}\mathbf{x}'^T P^T$$

$$E(\mathbf{y}'\mathbf{y}'^T) = P E(\mathbf{x}'\mathbf{x}'^T) P^T$$

In change $\mathbf{x} \rightarrow \mathbf{y}$, eigenvalues of covariance matrix remain > 0 , but can be modified (conserved if $P^T = P^{-1}$, orthogonal matrix).

Eigenvalues can actually take any positive values. In particular, covariance matrix can be made equal to the unit matrix, for instance in the basis of *principal components*.

- Two random vectors

$$\mathbf{x} = (x_1, x_2, \dots, x_n)^T$$

$$\mathbf{z} = (z_1, z_2, \dots, z_p)^T$$

$$E(\mathbf{x}'\mathbf{z}'^T) = E(x_i'z_j')$$

dimension $n \times p$

Change

$$\mathbf{x} \rightarrow \mathbf{u} \equiv A\mathbf{x} \qquad \mathbf{z} \rightarrow \mathbf{v} \equiv B\mathbf{z}$$

$$E(\mathbf{u}'\mathbf{v}'^T) = A E(\mathbf{x}'\mathbf{z}'^T) B^T$$

Covariance matrices will be denoted

$$C_{xx} \equiv E(\mathbf{x}'\mathbf{x}'^T)$$

$$C_{xy} \equiv E(\mathbf{x}'\mathbf{y}'^T)$$

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

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

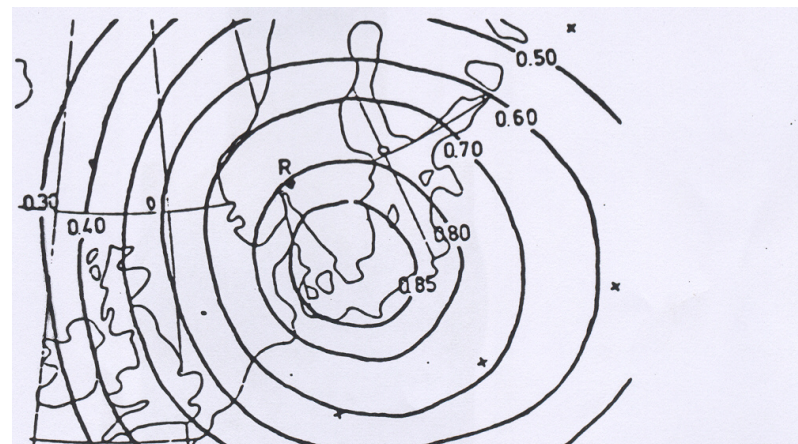
$$(\xi_1, \xi_2) \rightarrow C_{\phi}(\xi_1, \xi_2) \equiv E[\Phi'(\xi_1) \Phi'(\xi_2)]$$

- Correlation function

$$Cor_{\phi}(\xi_1, \xi_2) \equiv E[\Phi'(\xi_1) \Phi'(\xi_2)] / \{Var[\Phi(\xi_1)] Var[\Phi(\xi_2)]\}^{1/2}$$



.: 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

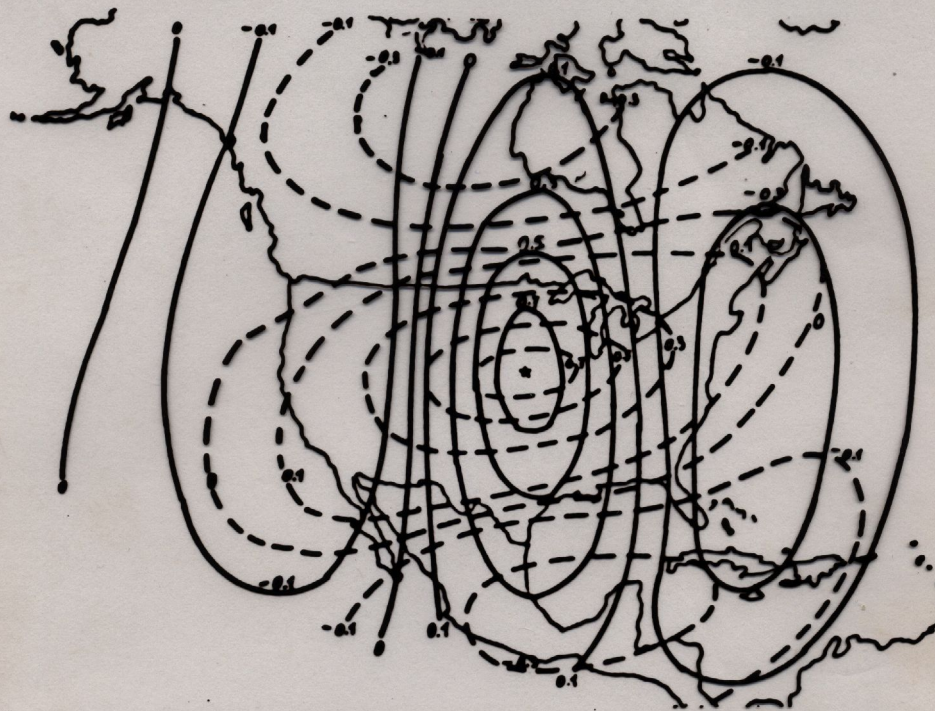
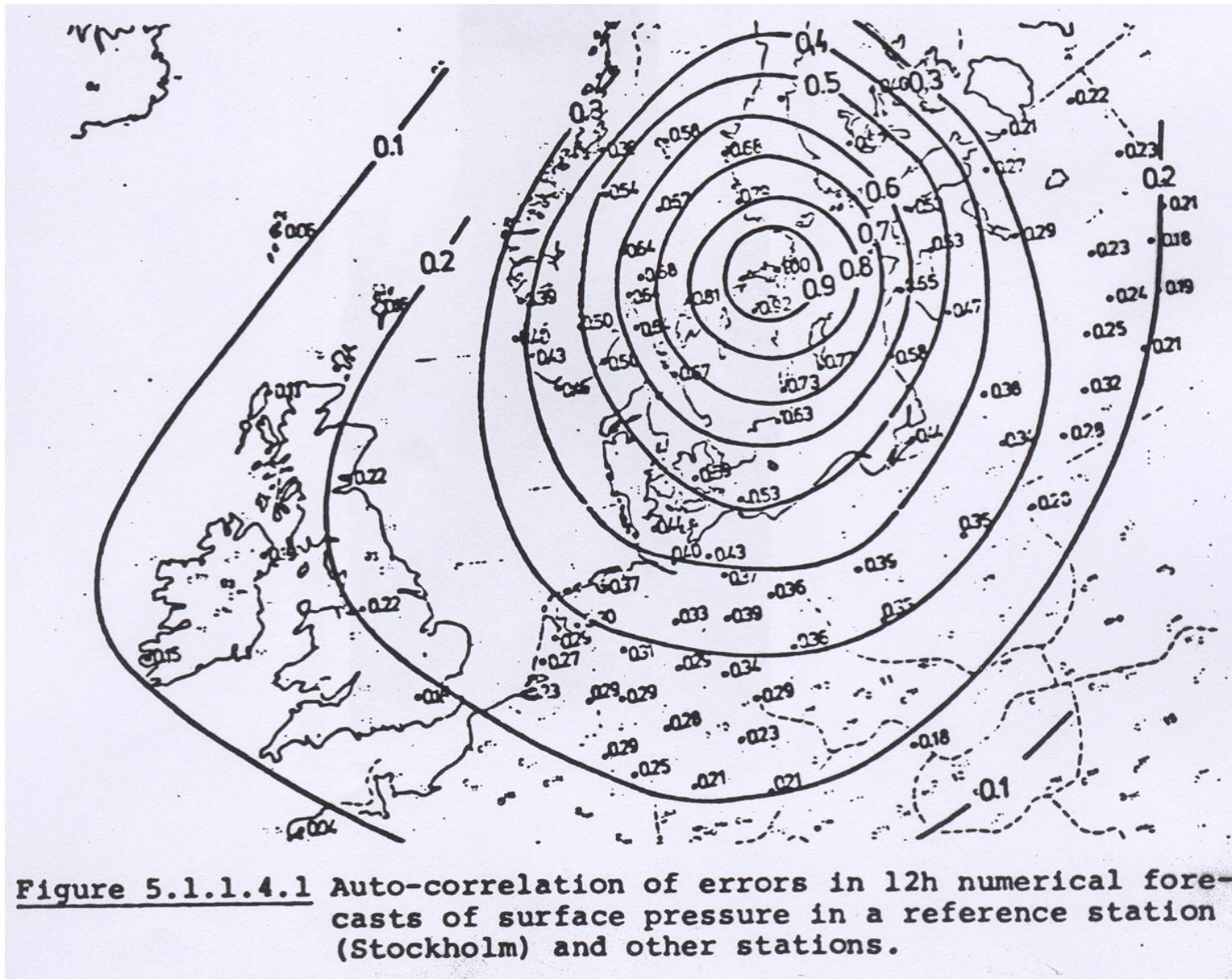


Figure 4.2.4.3: Isolines for the auto-correlation of the 500 mb u-wind component (dashed line) and the auto-correlation of the 500 mb v-wind component (full line). The "star" indicates the position of the reference station. (From Buel (1972).

After N. Gustafsson



After N. Gustafsson

Covariance function can be

homogeneous $C_{\phi}(\xi_1, \xi_2) = H(\xi_1 - \xi_2)$

or *isotropic* $C_{\phi}(\xi_1, \xi_2) = K(|\xi_1 - \xi_2|)$
(on the sphere, no difference)

N points $\xi_1, \xi_2, \dots, \xi_N$ in state space

N non-random coefficients $\lambda_1, \lambda_2, \dots, \lambda_N$

$$G \equiv \sum_i \lambda_i \Phi'(\xi_i)$$

$$E(G^2) = \sum_{i,j} \lambda_i \lambda_j C_{\phi}(\xi_i, \xi_j) \geq 0$$

$$E(G^2) = \sum_{i,j} \lambda_i \lambda_j C_{\Phi}(\xi_i, \xi_j) \geq 0$$

covariance functions are of *positive type* (or *definite positive*). Conversely, a function of positive type can be shown to be the covariance function of a random function.

Examples

On a circle, function $C(\xi_1, \xi_2) = \cos(\xi_1 - \xi_2)$ is covariance function of random function $\Phi(\xi) = 2 \cos(\xi + \alpha)$, where α is uniformly distributed over $[0, 2\pi]$.

In R^n , squared exponential

$$C(\xi_1, \xi_2) = \exp[- (\xi_1 - \xi_2)^T B^{-1} (\xi_1 - \xi_2)]$$

Bochner-Khintchin theorem. Homogeneous function C
 $(\xi_1, \xi_2) = H(\xi_1 - \xi_2)$ over R^n of positive type \Leftrightarrow Fourier
Transform of H is real ≥ 0 .

- ‘Optimal Interpolation’. Basic theory and basic properties. A simple example.

Optimal Interpolation

Random field $\Phi(\xi)$

Observation network $\xi_1, \xi_2, \dots, \xi_p$

For one particular realization of the field, observations

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

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

$$x^a = \alpha + \sum_j \beta_j y_j = \alpha + \boldsymbol{\beta}^T \mathbf{y}, \quad \text{where } \boldsymbol{\beta} = (\beta_j)$$

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)

$E[(x-x^a)^2]$ minimum $\Rightarrow E(x-x^a) = 0$ Estimate x^a is unbiased.

$$x^a = \alpha + \sum_j \beta_j y_j$$

$$E(x^a) = \alpha + \sum_j \beta_j E(y_j)$$

$$x^a - E(x) = \sum_j \beta_j [y_j - E(y_j)]$$

Computations are to be made on centred variables

$x'^a \equiv x^a - E(x)$ is the linear combination of the $y_j' = y_j - E(y_j)$ that minimizes the distance to $x' = x - E(x)$. It is the orthogonal projection, in the sense of covariance, of x' onto the space spanned by the y_j' 's.

Optimal Interpolation (continued 2)

$x' - x'^a$ uncorrelated with y_j'

$$E[(x' - x'^a) y_j'] = 0$$

$$x'^a = \sum_k \beta_k y_k'$$

$$\Rightarrow \sum_k \beta_k E(y_k' y_j') = E(x' y_j')$$

in matrix form $C_{yy} \boldsymbol{\beta} = C_{yx}$

Optimal Interpolation (continued 3)

Solution

$$\begin{aligned}x^a &= E(x) + E(x'y'^T) [E(y'y'^T)]^{-1} [y - E(y)] \\ &= E(x) + C_{xy} [C_{yy}]^{-1} [y - E(y)]\end{aligned}$$

$$\begin{aligned}i. e., \quad \beta^T &= C_{xy} [C_{yy}]^{-1} \\ \alpha &= E(x) - \beta^T E(y)\end{aligned}$$

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

Minimized quadratic estimation error

$$\begin{aligned}E[(x-x^a)^2] &= E(x'^2) - E[(x'^a)^2] \\ &= C_{xx} - C_{xy} [C_{yy}]^{-1} C_{yx}\end{aligned}$$

Estimation made in terms of deviations x' and y' from expectations $E(x)$ and $E(y)$.

Optimal Interpolation (continued 4)

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

$$y_j = \Phi(\xi_j) + \varepsilon_j$$

$$E(y_j'y_k') = E[\Phi'(\xi_j) + \varepsilon_j'] [\Phi'(\xi_k) + \varepsilon_k']$$

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

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

and

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

Optimal Interpolation (continued 5)

Unique observation ($p=1$) $y_1 = \Phi(\xi_1) + \varepsilon_1$

Value $x = \Phi(\xi)$ at some point ξ to be estimated
(all values assumed to be centred)

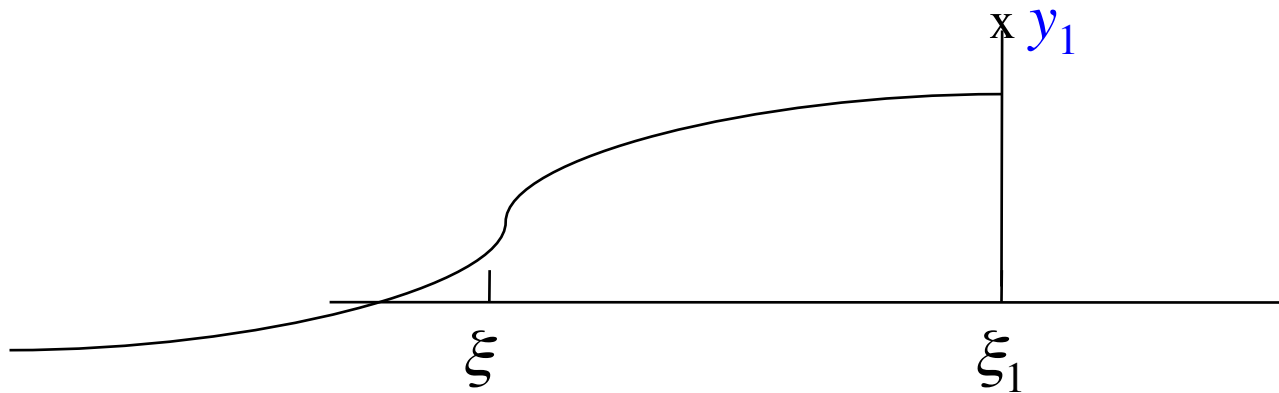
$$C_{yy} \beta = C_{yx}$$

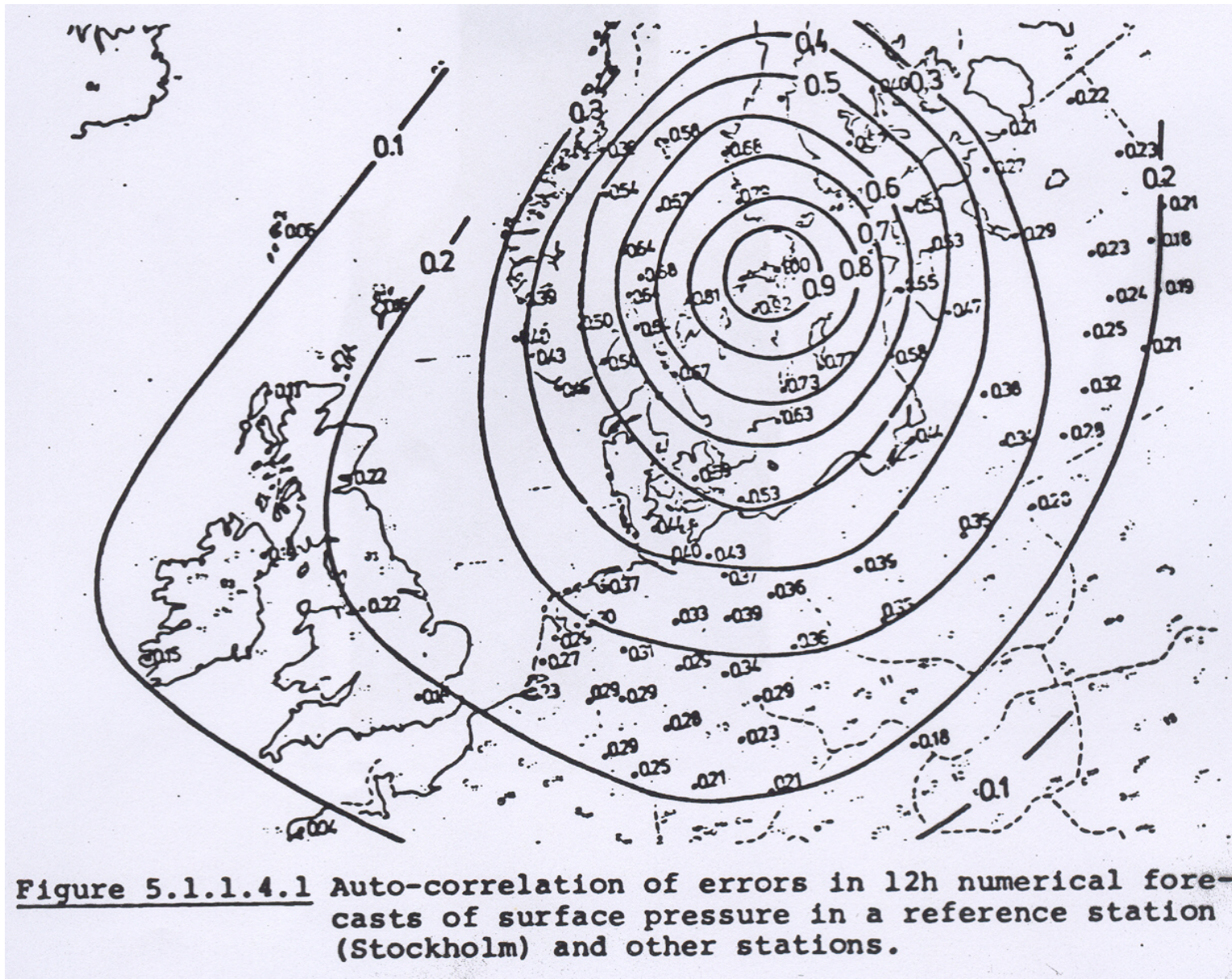
$$C_{yy} = E(y_1^2) = C_{\Phi}(\xi_1, \xi_1) + r \quad C_{yx} = C_{\Phi}(\xi, \xi_1)$$

$$x^a = \Phi^a(\xi) = \frac{C_{\Phi}(\xi, \xi_1)}{C_{\Phi}(\xi_1, \xi_1) + r} y_1$$

Optimal Interpolation (continued 6)

$$x^a = \Phi^a(\xi) = \frac{C_\Phi(\xi, \xi_1)}{C_\Phi(\xi_1, \xi_1) + r} y_1$$

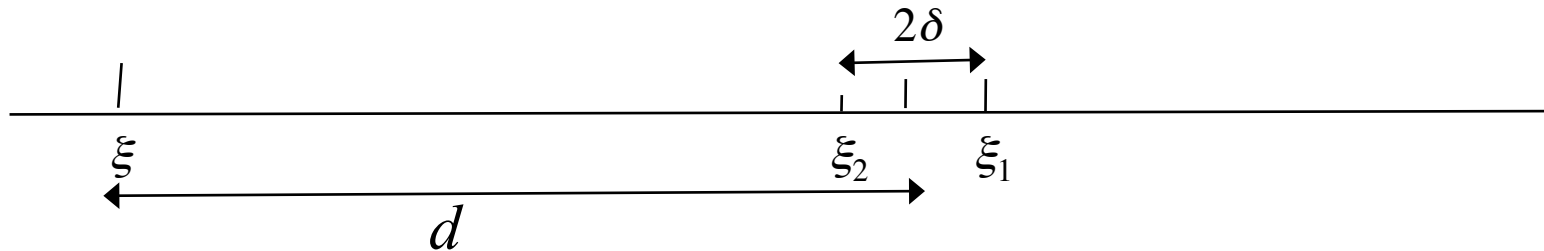




After N. Gustafsson

Optimal Interpolation (continued 7)

Two mutually close observations ($p=2$) $y_j = \Phi(\xi_j) + \varepsilon_j$, $j = 1, 2$



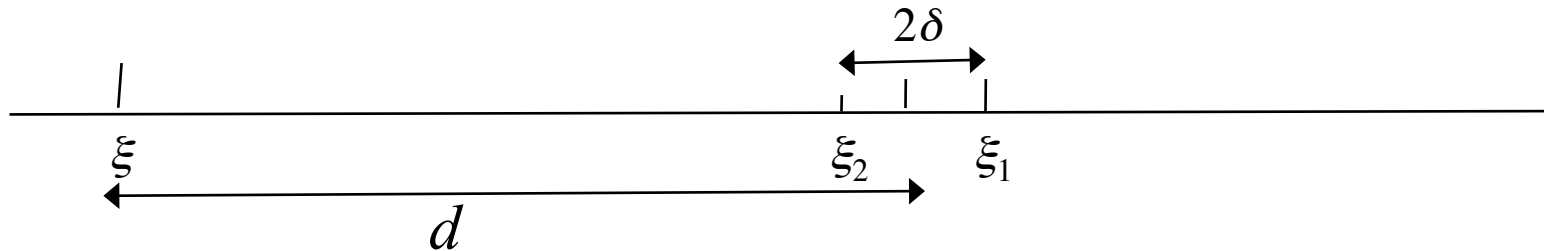
Homogeneous covariance function $C_\phi(\chi_1, \chi_2) = \Gamma(\chi_1 - \chi_2)$

Linear system for weights β_j 's

$$\begin{pmatrix} \Gamma(0) + r & \Gamma(2\delta) \\ \Gamma(2\delta) & \Gamma(0) + r \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} = \begin{pmatrix} \Gamma(d + \delta) \\ \Gamma(d - \delta) \end{pmatrix}$$

Optimal Interpolation (continued 8)

Two mutually close observations ($p=2$) $y_j = \Phi(\xi_j) + \varepsilon_j$, $j = 1, 2$

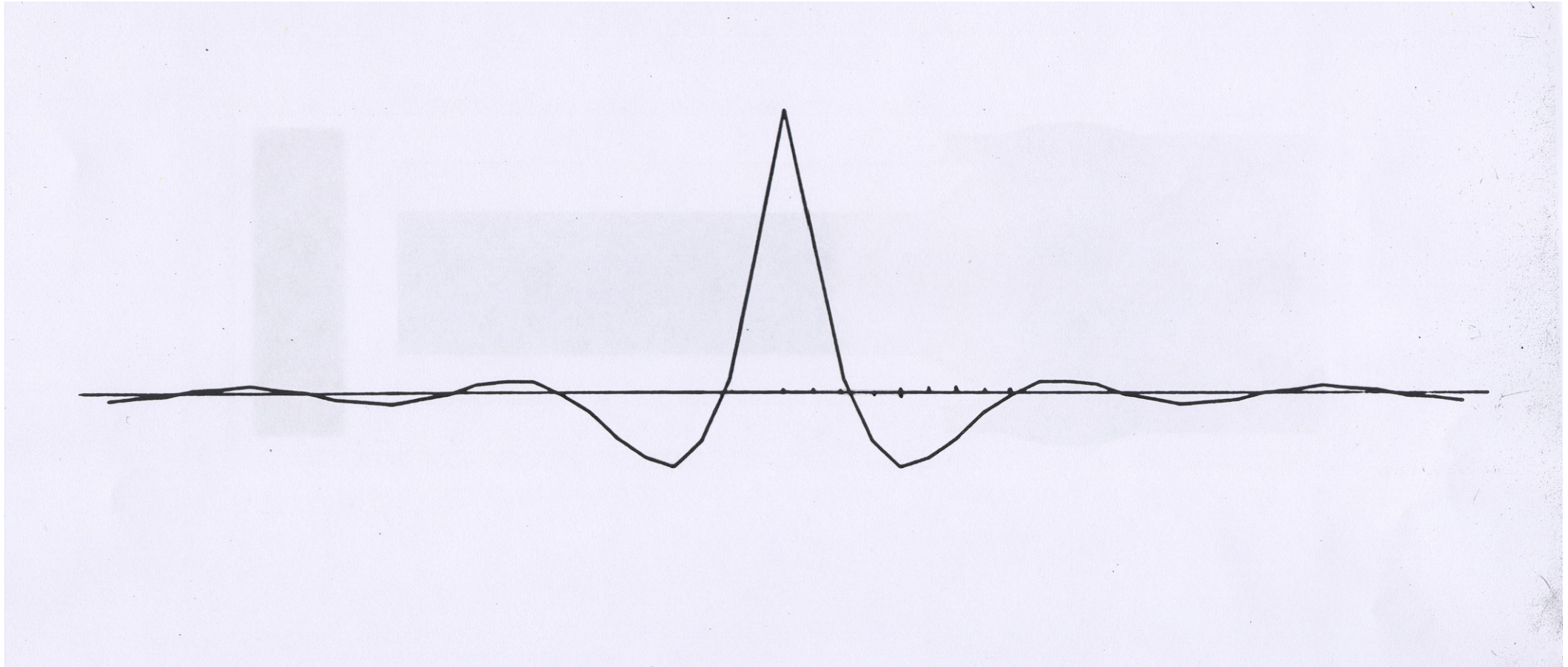


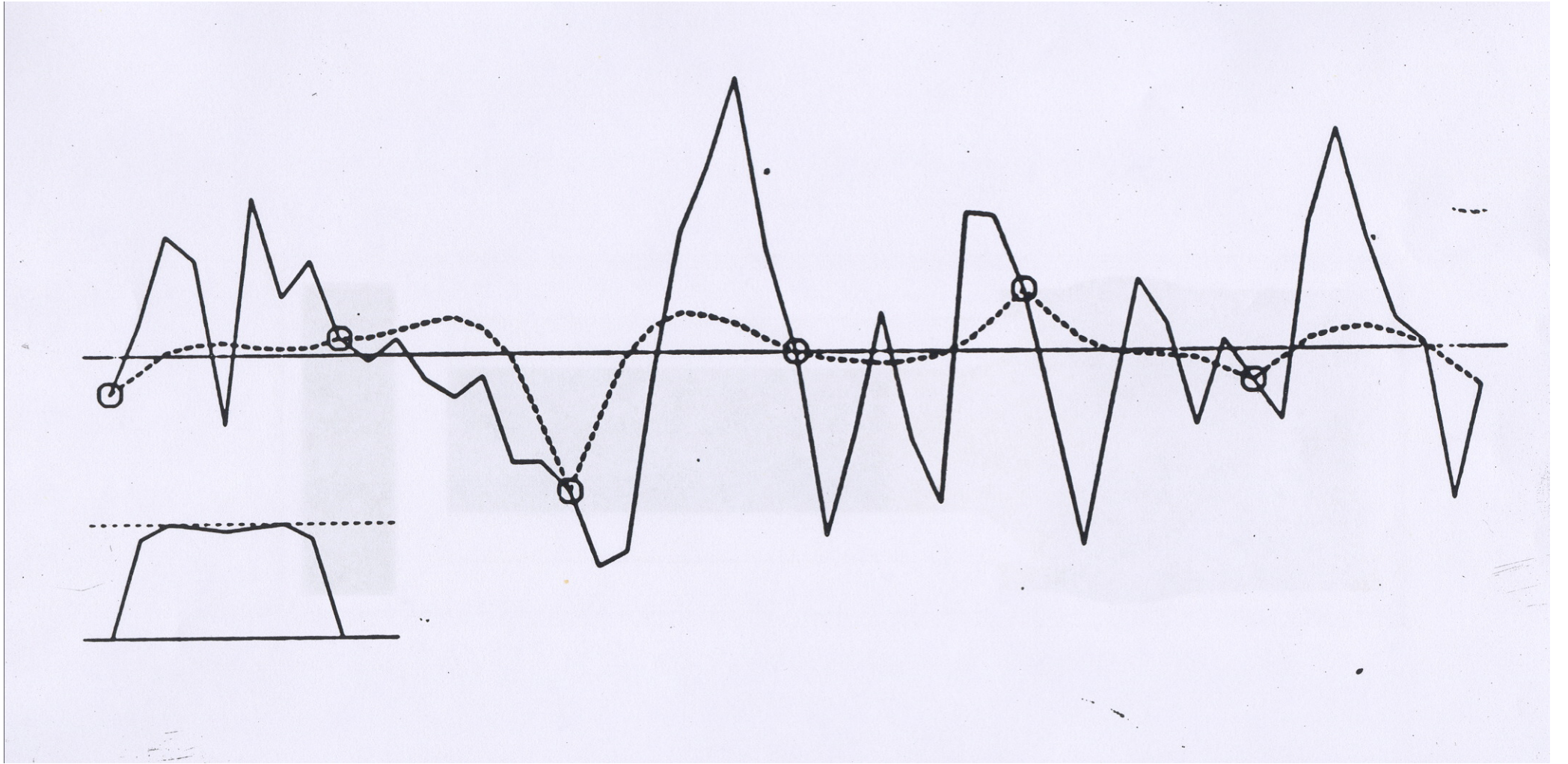
$$\beta_1 + \beta_2 = \frac{\Gamma(d + \delta) + \Gamma(d - \delta)}{\Gamma(0) + \Gamma(2\delta) + r}$$

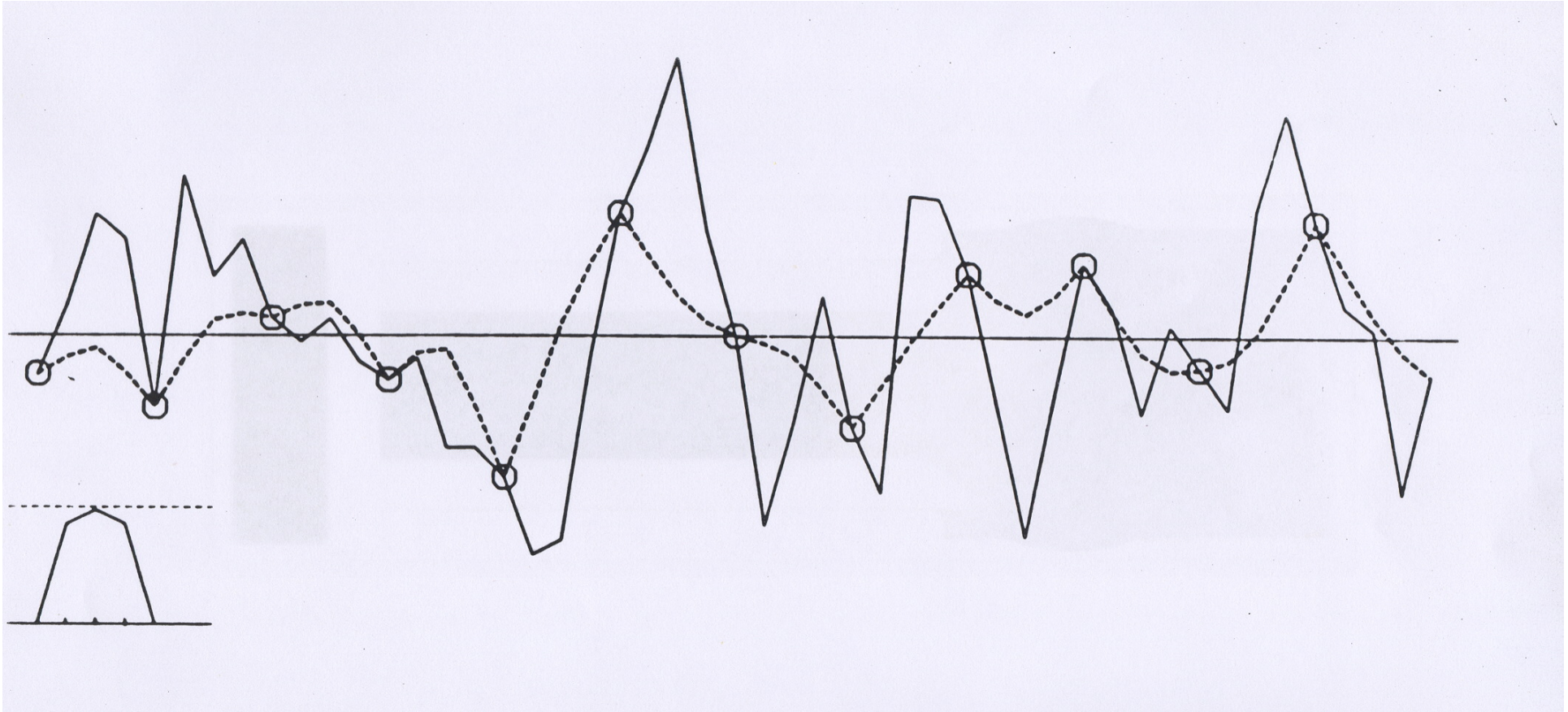
For small δ ,

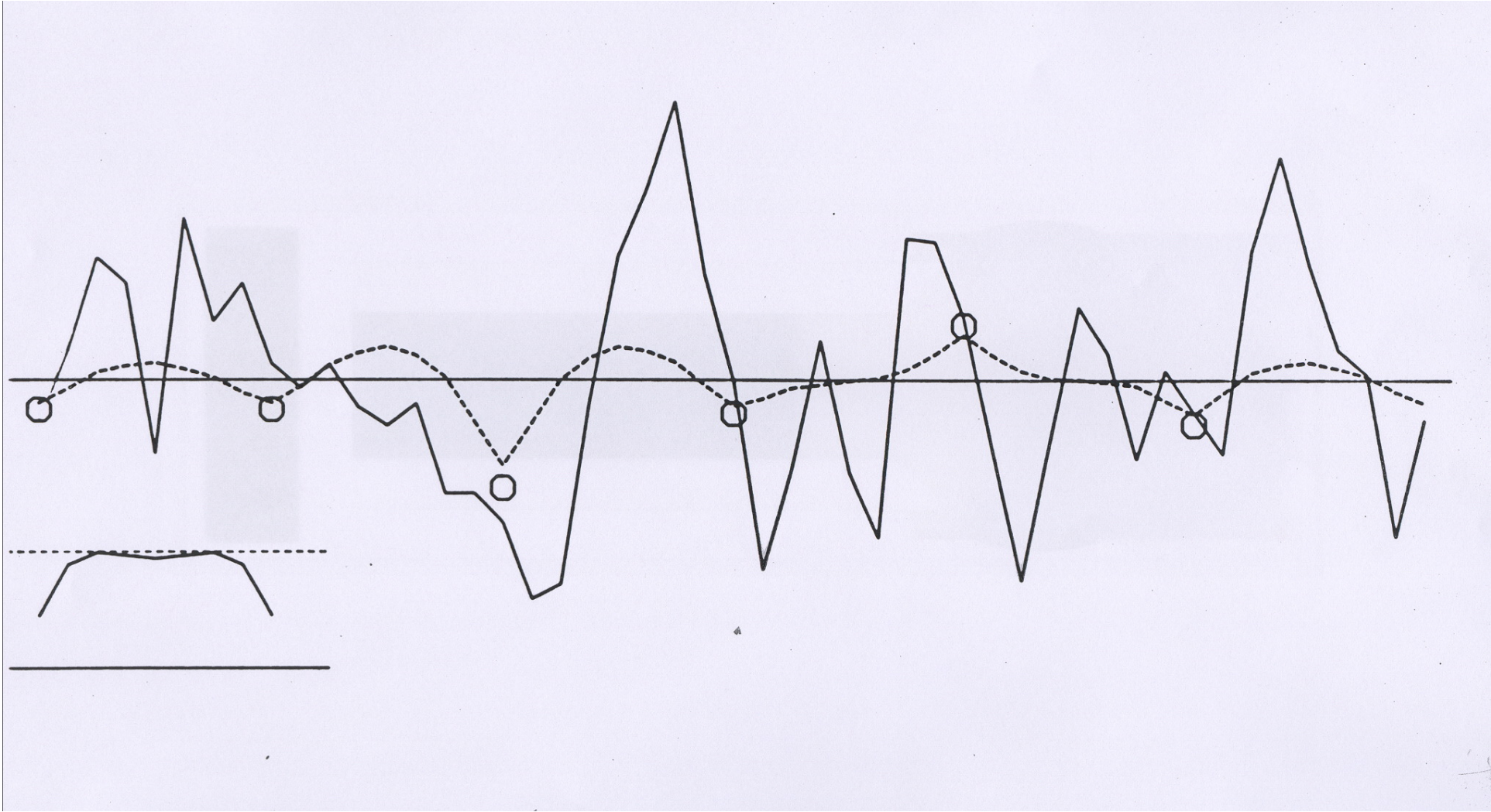
$$\beta_1 + \beta_2 = \frac{\Gamma(d)}{\Gamma(0) + r/2}$$

Sum equal weight that would be given to a unique observation located at position d , with error $r/2$









Optimal Interpolation (continued 10)

$$x^a = E(x) + C_{xy} [C_{yy}]^{-1} [y - E(y)]$$

Vector

$$\boldsymbol{\mu} = (\mu_j) \equiv [C_{yy}]^{-1} [y - E(y)]$$

is independent of variable to be estimated

$$x^a = E(x) + \sum_j \mu_j E(x'y_j')$$

Optimal Interpolation (continued 11)

$$x^a = E(x) + \sum_j \mu_j E(x'y_j')$$

$$\Phi^a(\xi) = E[\Phi(\xi)] + \sum_j \mu_j E[\Phi'(\xi) y_j']$$

Under hypotheses made above, $E[\Phi'(\xi) y_j'] = C_\phi(\xi, \xi_j)$

$$\Phi^a(\xi) = E[\Phi(\xi)] + \sum_j \mu_j C_\phi(\xi, \xi_j)$$

Correction made on background expectation is a linear combination of the p functions $C_\phi(\xi, \xi_j)$

$C_\phi(\xi, \xi_j)$, considered as a function of estimation position ξ , is the *representer* associated with observation y_j .

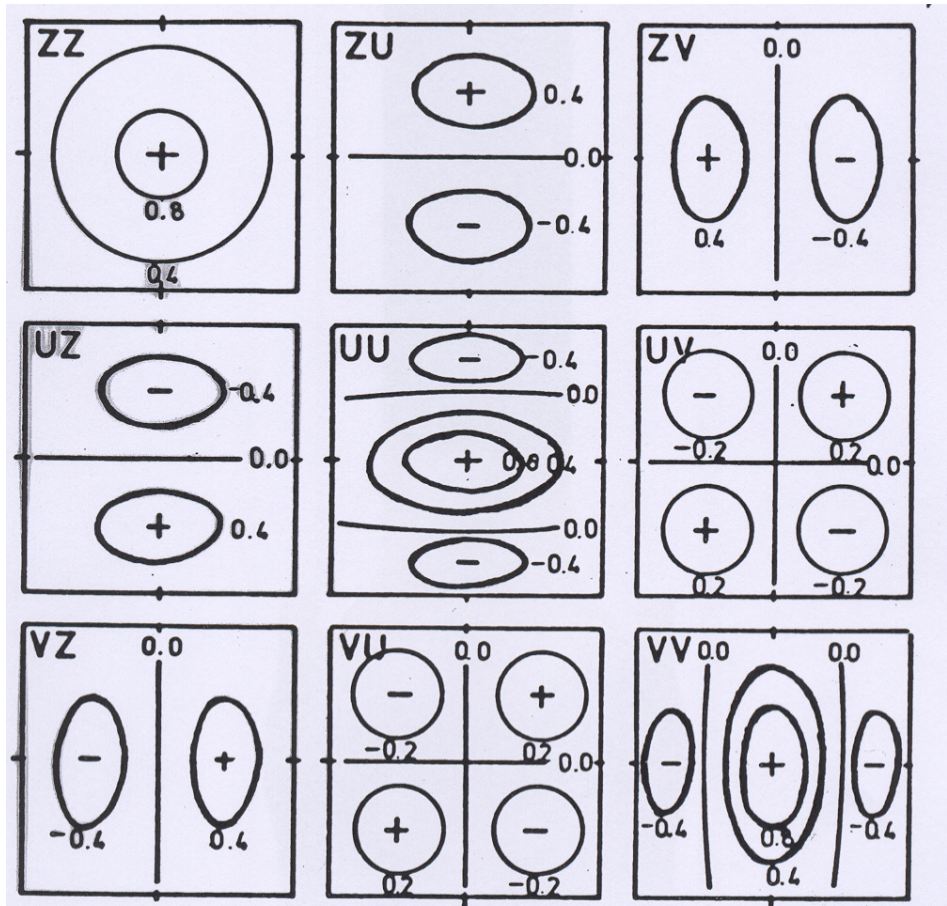
Optimal Interpolation (continued 12)

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.



4.: Schematic illustration of correlation functions and cross-correlation functions for multi-variate analysis derived by the geostrophic assumption.

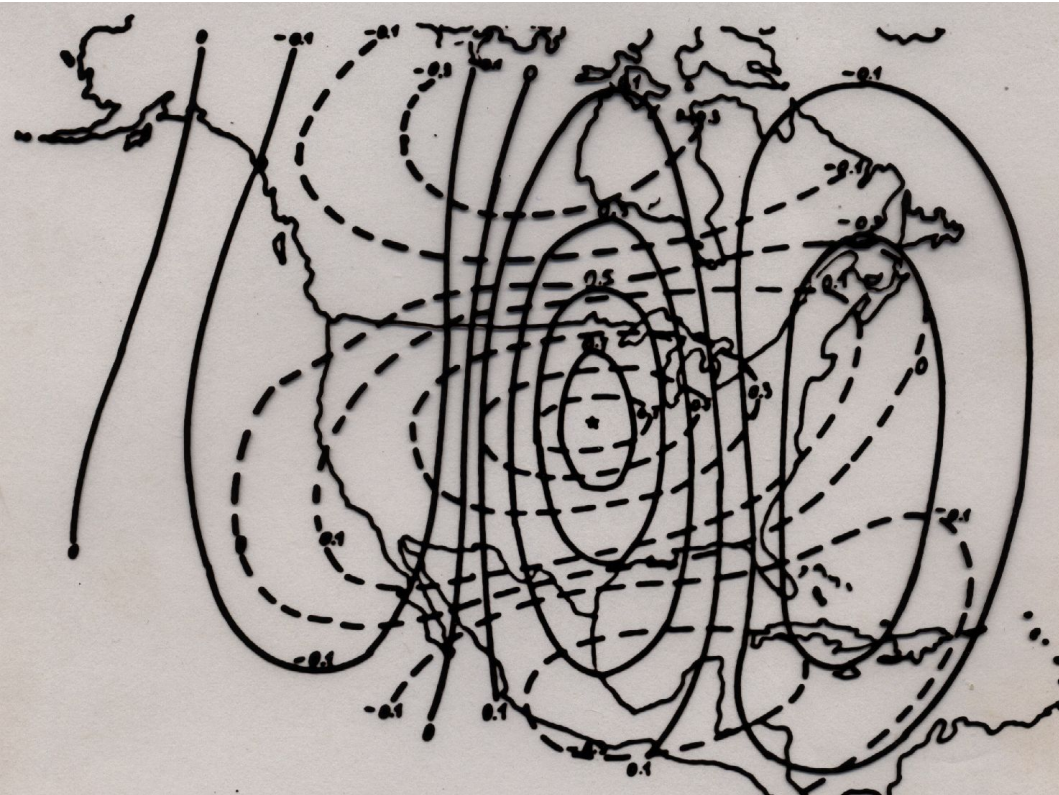
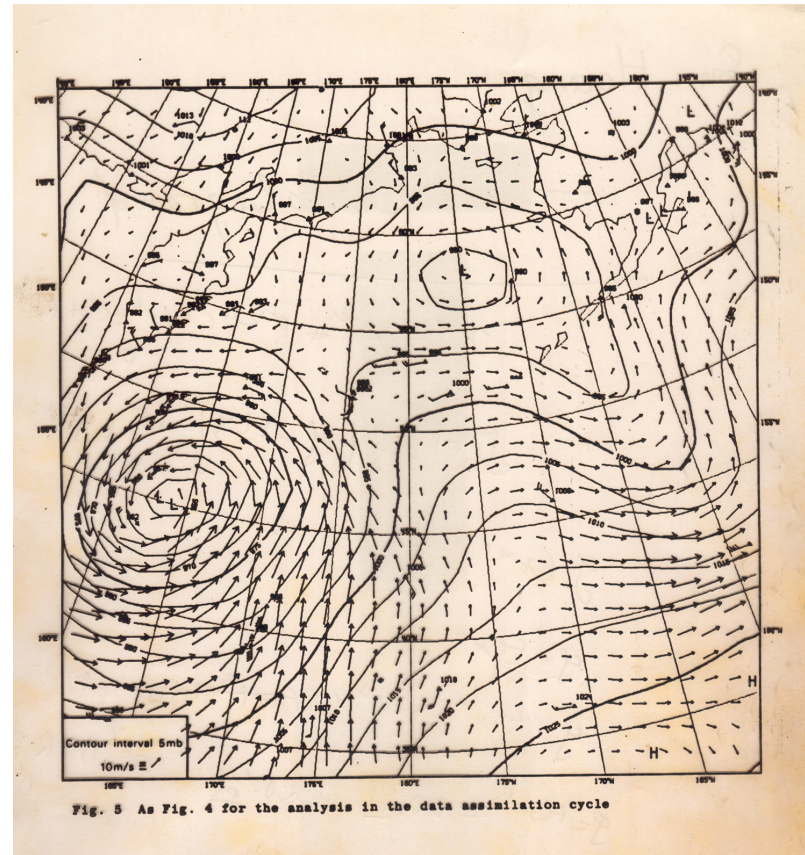
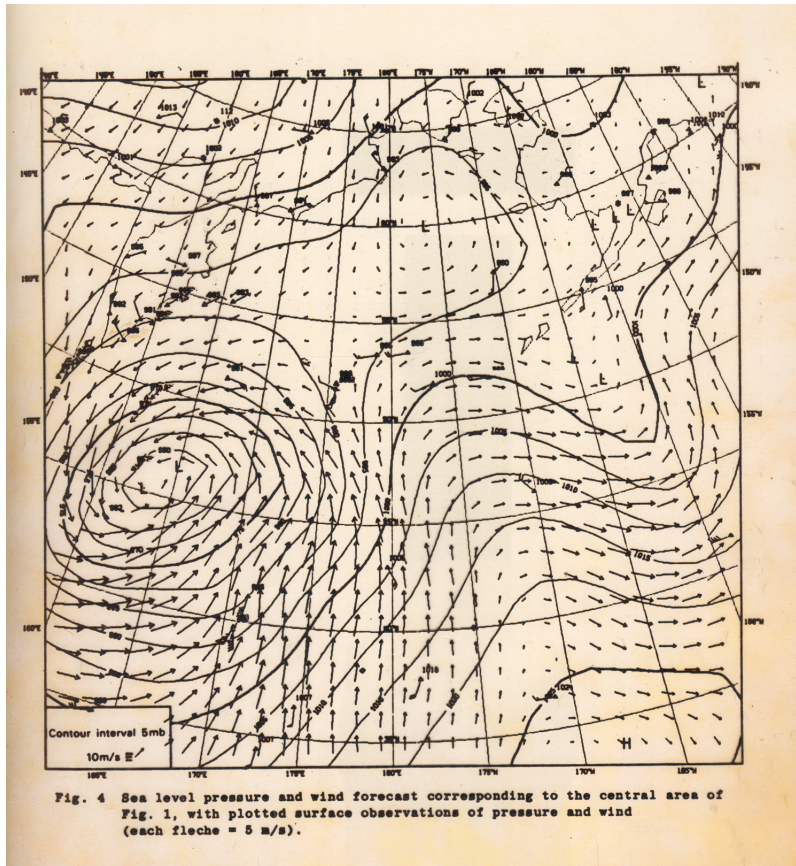


Figure 4.2.4.3: Isolines for the auto-correlation of the 500 mb u-wind component (dashed line) and the auto-correlation of the 500 mb v-wind component (full line). The "star" indicates the position of the reference station. (From Buel (1972)).

After N. Gustafsson



After A. Lorenc, MWR, 1981

1200 GMT 19 January 1979

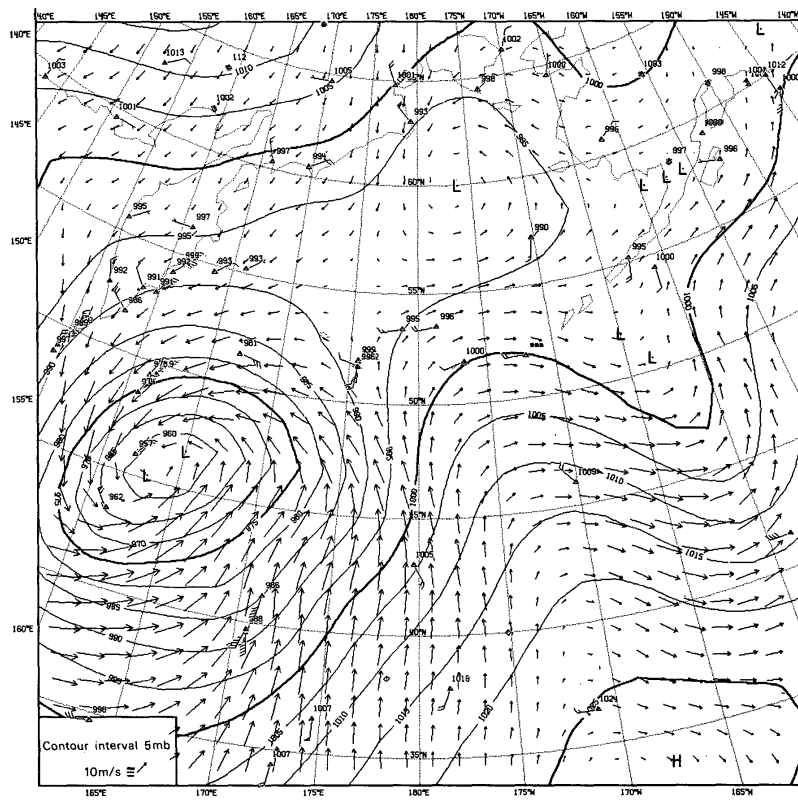


FIG. 14. Sea level pressure and wind forecast corresponding to the central area of Fig. 11, with plotted surface observations of sea level pressure and wind (each barb = 5 m s⁻¹).

1200 GMT 19 January 1979

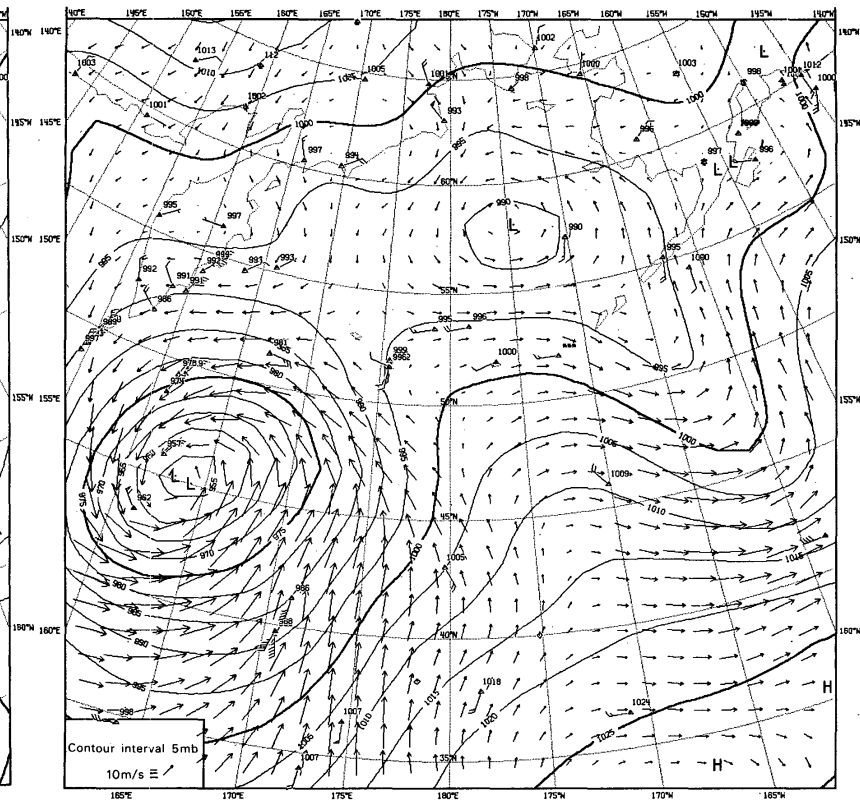


FIG. 15. As in Fig. 14 for the analysis in the data-assimilation cycle.

After A. Lorenc, MWR, 1981

Optimal Interpolation (continued 5)

Observation vector \mathbf{y}

Estimation of a scalar x

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

$$\begin{aligned} p^a &\equiv E[(x-x^a)^2] = E(x'^2) - E[(x'^a)^2] \\ &= C_{xx} - C_{xy} [C_{yy}]^{-1} C_{yx} \end{aligned}$$

Estimation of a vector \mathbf{x}

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

$$\begin{aligned} \mathbf{P}^a &\equiv E[(\mathbf{x}-\mathbf{x}^a) (\mathbf{x}-\mathbf{x}^a)^T] = E(\mathbf{x}'\mathbf{x}'^T) - E(\mathbf{x}'^a \mathbf{x}'^{aT}) \\ &= C_{xx} - C_{xy} [C_{yy}]^{-1} C_{yx} \end{aligned}$$

Optimal Interpolation (continued 6)

$$\mathbf{x}^a = E(\mathbf{x}) + \mathbf{C}_{xy} [\mathbf{C}_{yy}]^{-1} [\mathbf{y} - E(\mathbf{y})]$$

$$\mathbf{P}^a = \mathbf{C}_{xx} - \mathbf{C}_{xy} [\mathbf{C}_{yy}]^{-1} \mathbf{C}_{yx}$$

If probability distribution for couple (\mathbf{x}, \mathbf{y}) is Gaussian (with, in particular, covariance matrix

$$\mathbf{C} \equiv \begin{pmatrix} \mathbf{C}_{xx} & \mathbf{C}_{xy} \\ \mathbf{C}_{yx} & \mathbf{C}_{yy} \end{pmatrix}$$

then Optimal Interpolation achieves Bayesian estimation, in the sense that

$$P(\mathbf{x} | \mathbf{y}) = \mathcal{N}[\mathbf{x}^a, \mathbf{P}^a]$$

Cours à venir

~~Jeudi 19 Mars~~

~~Jeudi 26 mars~~

~~Jeudi 02 avril~~

Jeudi 09 avril

Jeudi 16 avril

Jeudi 23 avril

Jeudi 30 avril

Jeudi 14 mai