



**HAL**  
open science

# An efficient algorithm for the large-scale smoothing of scattered data retrieved from remote sounding experiments

D. Fussen

► **To cite this version:**

D. Fussen. An efficient algorithm for the large-scale smoothing of scattered data retrieved from remote sounding experiments. *Annales Geophysicae*, 2003, 21 (7), pp.1645-1652. hal-00317131

**HAL Id: hal-00317131**

**<https://hal.science/hal-00317131>**

Submitted on 18 Jun 2008

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# An efficient algorithm for the large-scale smoothing of scattered data retrieved from remote sounding experiments

D. Fussen

Institut d'Aéronomie Spatiale de Belgique, Brussels, Belgium

Received: 23 August 2002 – Revised: 22 December 2002 – Accepted: 9 February 2003

**Abstract.** We present a new algorithm for smoothing/interpolation of two-dimensional fields applicable to noisy data observed at scattered sites. The technique is based on a special statistics allowing one to simultaneously minimize the fit residual and the correlation between residuals of adjacent points. The principle of the method is first explained in the 1-D case and then extended to the 2-D case by adjunction of a regularization operator. The method is compared with different algorithms (Loess-Renka, Biharmonic Spline and kriging) in three test cases related to remote sounding of the Earth's atmosphere by space-borne experiments.

**Key words.** Atmospheric composition and structure (evolution of the atmosphere; instruments and techniques; general or miscellaneous)

## 1 Introduction

Data retrieved from remote sounding space experiments generally provide atmospheric information with good coverage in space and time. These data can be validated by comparing results obtained by independent instruments and/or by ground-based observations. However, satellite data from several orbiting platforms correspond to different sounding geometries, such as solar or stellar occultation, nadir or limb viewing, etc. Geolocation of sounded regions is also clearly influenced by the orbital state vector and possible hardware constraints. From these synoptic and composite data sets, it is, therefore, highly desirable to use efficient numerical tools capable of producing regular (gridded) fields on which validation and statistical analysis can be performed (Lait, 2000).

Sophisticated interpolation techniques, such as 4-D variational assimilation (Fisher and Lary, 1995), are presently developed and will be routinely used for future processing of stratospheric data. However, these procedures are computationally very expensive and their success strongly depend

upon the quality of the underlying atmospheric model. An extensive review of atmospheric data analysis can be found in Daley (1991) and a considerable number of methods for approximating two-dimensional (“2-D”) fields from scattered data have been published (see Foley and Hagen, 1994).

This work presents a new algorithm for 2-D smoothing which is based on an expansion in orthogonal polynomials combined with a regularization driven by statistics different from the usual Chi-squared. After disappointing experiences with existing schemes, the method has been designed to be robust (it should always converge), automatic (not requiring interactive inspection of the smoothing level), reasonably fast and not very sensitive to the inaccuracy of experimental error bars. As the main idea of the algorithm is to Minimize the Correlation between Residuals, it will be referred to as the MCR method hereafter.

For the sake of clarity, we will first explain the principles of the algorithm for one-dimensional (“1-D”) problems without addressing a rigorous mathematical framework. In the second part, we will extend it to 2-D cases and we will investigate its performance (with respect to alternative methods) for realistic data sets coming from atmospheric remote sensing experiments.

It is well known that the generalized least-squares techniques require the use of a full covariance matrix to avoid biases induced by possible correlation between data point errors (notice that the data points themselves are highly correlated by the underlying physics). Also, it is clear that the data point processing of instrumental data through inversion algorithms may introduce correlation between values retrieved at adjacent altitudes or between different species retrieved simultaneously. However, it is reasonable to assume that the estimated random experimental errors at different geolocations and sampling times are uncorrelated because they correspond to independent realisations of the instrumental/inversion noise. Clearly, there does not exist a smoothing algorithm capable of removing a bias in a given measurement and if such a bias would exist, it could only be corrected at the level of the combined “measurement-data processing”

chain. The removing of the bias would also require the availability of many measurements at close geolocations and time or the cross-validation with respect to several independent experiments. Also, experimental data from spaceborne instruments are mostly available as independent measurements with estimated variances. If a bias is suspected in the data and if it can be estimated but not exactly computed, it can be quadratically added to the random error variance before applying the smoothing algorithm. The MCR method explained hereafter presupposes that the experimental errors associated with the measurement points have a zero mean and are uncorrelated in space and time.

## 2 1-D MCR

We deal with the problem of approximating  $m$  bivariate data  $(x_i, f_i)$ ,  $i = 1 \dots m$  by a smooth curve  $f(x)$  that should approximate the unknown reality  $g(x)$ . The measured data are supposed to contain a noise component  $\delta_i$  so that

$$f_i = g(x_i) + \delta_i. \quad (1)$$

The standard deviation  $\sigma_i$  of  $\delta_i$  and the experimental data  $x_i$  may have been measured by different instruments with specific random error distributions and possible moderate bias. The classical fitting strategy (Press et al., 1992) consists of minimizing the merit function  $\chi^2$  defined by

$$\chi^2 = \sum_{i=1}^m d_i^2, \quad (2)$$

where the normalized residual  $d_i$  is

$$d_i = \frac{f(x_i; c_0, c_1, \dots, c_n) - f_i}{\sigma_i} \quad (3)$$

and the approximating function  $f(x; c_0, c_1, \dots, c_n)$  depends on  $(n + 1)$  parameters.

Ideally, the value of  $\chi^2$  should decrease when the number of fitting parameters  $c_0, c_1, \dots, c_n$  is increased until a plateau (whose value is about  $m \pm \sqrt{2m}$  for large  $m$ ) is reached. An insufficient number of parameters causes underfitting and some information content of the data set is lost, while too many parameters means overfitting and the fit starts to represent the noise. A simple example is given in the upper part of Fig. 1. The optimal fit is clearly at the border between both extreme regimes. On the other hand, the values of  $\sigma_i$  are often poorly known or even unknown, and the optimal fit has to rely to the appearance of such a plateau. Furthermore, it is not always easy to recognize the existence of the plateau, which can be obscured by multiple extrema.

Quite recently, the Durbin-Watson (Durbin and Watson, 1950) statistics was proposed as an alternative statistics for the least-squares spline approximation of noisy data (Thijsse et al., 1998). Basically, the merit function of the fit is given by:

$$Q = \frac{\sum_{i=1}^{m-1} (d_{i+1} - d_i)^2}{\sum_{i=1}^m d_i^2}, \quad (4)$$

i.e. the ratio of two variances where the denominator is the usual aggregate magnitude of residuals, while the numerator measures the serial correlation between them. For reasonably smooth underlying functions, the evolution of  $Q$  from underfitting to overfitting regimes is easy to understand (see Fig. 1). For underfitting cases, residuals are strongly correlated because two adjacent points tend to have the same residual with respect to the fit ( $d_{i+1} \simeq d_i$  and  $Q$  is closer to 0). In the case of overfitting, residuals tend to be anticorrelated ( $d_{i+1} \simeq -d_i$  and  $Q$  is closer to 4), as the fit rapidly oscillates between neighbouring points. The optimal fit, for which  $Q$  is about  $2 \pm \frac{2}{\sqrt{m}}$ , corresponds to a minimal correlation between residuals ( $\sum_i d_{i+1}d_i \simeq 0$ ) because they only contain noise. A considerable advantage of  $Q$  arises from its relative insensitivity to the knowledge of the  $\sigma_i$ , since these appear both in the numerator and the denominator.

Keeping in mind that our goal is to globally approximate smooth geophysical fields retrieved from space experiments, we decided to develop the unknown function over a basis of orthogonal functions. Excellent (but not unique) candidates are Chebyshev polynomials  $T_k(u) = \cos(k \arccos(u))$  (Press et al., 1992) and, within a domain  $[x_{\min}, x_{\max}]$ ,  $f(x)$  may be expanded as:

$$f(x) \approx \sum_{k=0}^n c_k T_k(u(x)), \quad (5)$$

where

$$u(x) = \frac{2x - (x_{\max} + x_{\min})}{(x_{\max} - x_{\min})}, \quad u \in [-1, 1]. \quad (6)$$

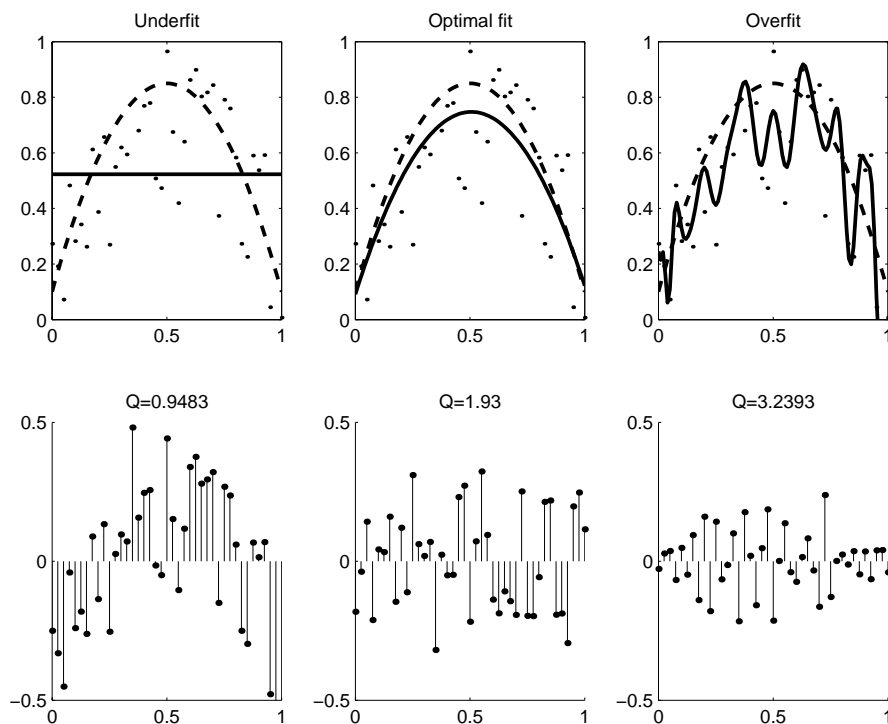
In Fig. 2, we simulated a synthetic atmospheric transmittance measured by an occultation experiment on which detector shot noise (varying as the square root of the signal) has been added. The least-squares linear problem associated with Eq. (5) can be solved for increasing values of  $n$  until the value  $Q = 2$  (nominal) or  $Q_+ = 2 + \frac{2}{\sqrt{m}}$  (conservative) is reached (for  $n \simeq 7$ ). The resulting fitted curve has the desired smoothness and it is able to capture the slope change around 25 km. Considering larger  $n$  values would only result in overfitting (wavy structures) and possible ill-conditioning of the linear system.

## 3 2-D MCR

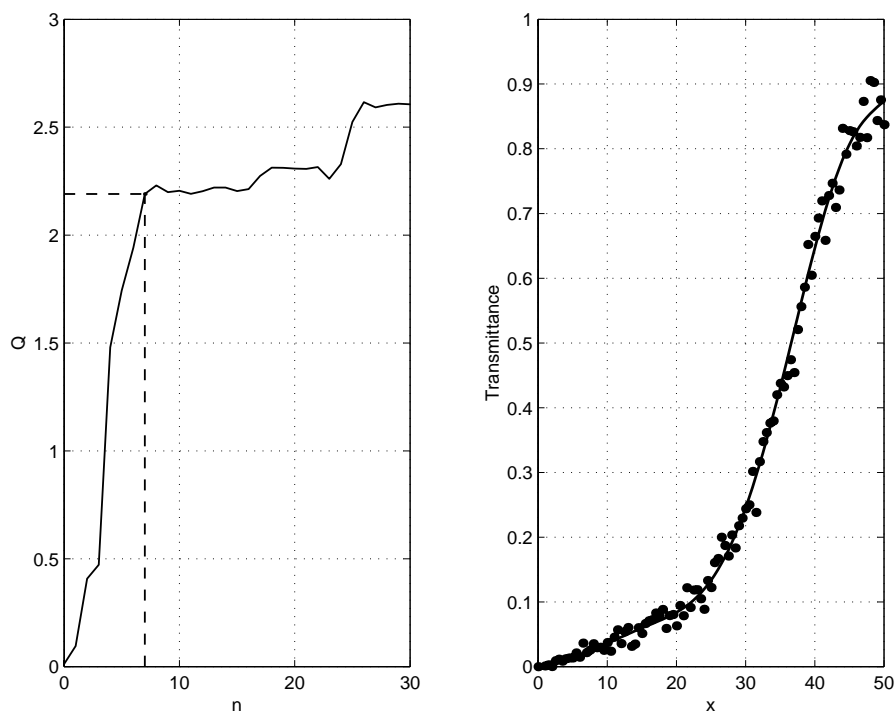
It is possible to generalize the orthogonal function expansion to the 2-D case by:

$$f(x, y) \approx \sum_{k=0}^n \sum_{l=0}^{n-k} c_{kl} T_k(u(x)) T_l(v(y)) \quad (7)$$

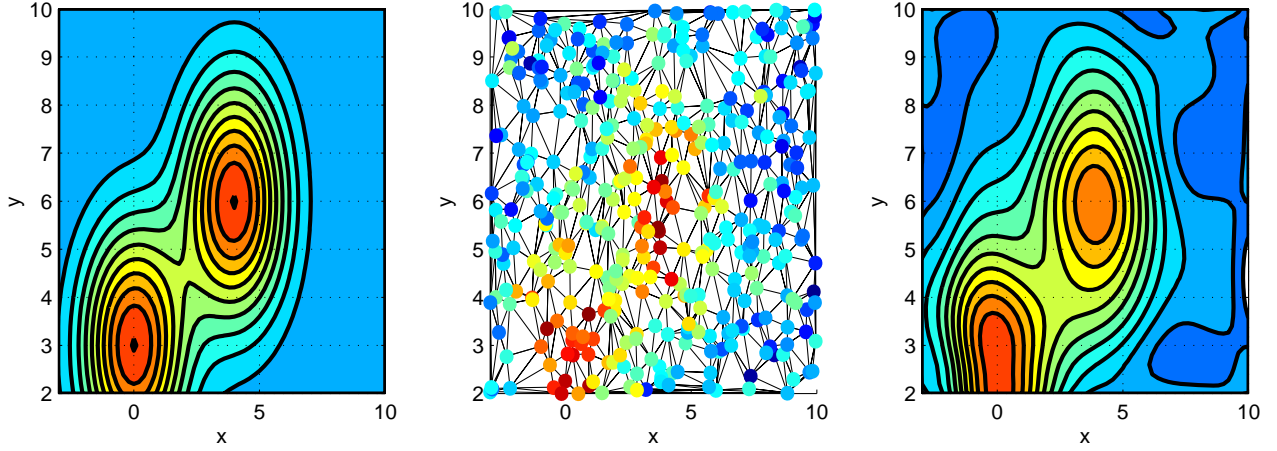
(Chebyshev polynomials could be replaced by spherical harmonics if a full coverage of the globe is desired). A natural way to consider the correlation between adjacent measurement points is to perform a Delaunay triangulation of the considered domain. For a given data point  $i$ , we consider



**Fig. 1.** Top: Transition from underfitting to overfitting for the 1-D case. Dots: experimental data. Dashed line: reality from which experimental data were generated by adding a Gaussian noise of 0.2 amplitude. Full line: polynomial fit. Bottom: respective residuals and associated  $Q$  value (Eq. 4). From left to right, notice the evolution in the correlation of adjacent residuals.



**Fig. 2.** A typical example of 1-D MCR fitting. Left: Evolution of  $Q$  (Eq. 4) versus  $n$  (Eq. 5), the polynomial order of the fit. Right: A synthetic transmittance signal with shot noise added (dots) and the optimal MCR fit ( $n = 7$ ).



**Fig. 3.** Test case TC1. Left: The exact field consisting of the sum of 2 Gaussian functions of amplitude 1. Middle: The same field, with addition of random Gaussian noise (standard deviation  $s_0 = 0.2$ ), sampled at 400 scattered locations. Right: reconstructed field by MCR method ( $n = 10$ ).

the list  $\{1, 2, \dots, p(i)\}$  of all neighbours, i.e. all points belonging to a Delaunay triangle whose point  $i$  is a vertex. For the MCR statistics, we propose

$$Q = \frac{\sum_{i=1}^m \sum_{j=1}^{p(i)} (d_{ij} - d_i)^2}{\sum_{i=1}^m p(i) d_i^2}, \quad (8)$$

where  $d_{ij}$  refers to the residual of the  $j$ -th neighbour of point  $i$ . This statistics has the same asymptotic properties as those present in Eq. (4) after normalization with respect to the number  $p$  of neighbours surrounding the point  $i$ . In case of underfitting,  $d_{ij} \simeq d_i$  and  $Q$  almost vanish, while overfitting corresponds to anticorrelated residuals ( $d_{ij} \simeq -d_i$ ) and makes  $Q$  closer to 4 if the number of data points is large enough. The optimal fit is achieved when the sum of products of neighbour residuals times the point residual, summed over all points, approaches zero, i.e.  $Q \simeq 2 \pm \frac{2}{\sqrt{m}}$ .

An additional problem of the 2-D case is the large increase in the number of coefficients  $c_{kl}$  to compute, which may cause severe ill-conditioning (Dierckx, 1993). It means that the solution becomes extremely sensitive to the noise amplitude even if the statistical criterion is fulfilled. Therefore, it is necessary to introduce a regularization technique (Hansen, 1992). If  $\mathbf{A}$  represents the design matrix of the problem (hereafter, we will consider  $\sigma_i = 1$  in order to simplify the notation)

$$A = \begin{pmatrix} T_0(u_1)T_0(v_1) & \dots & T_0(u_1)T_n(v_1) & \dots & T_n(u_1)T_0(v_1) \\ T_0(u_2)T_0(v_2) & \dots & & & \vdots \\ \vdots & & & & \vdots \\ T_0(u_m)T_0(v_m) & \dots & T_0(u_m)T_n(v_m) & \dots & T_n(u_m)T_0(v_m) \end{pmatrix}, \quad (9)$$

the constrained linear inversion consists of minimizing the Lagrangian merit function consisting of two quadratic forms

$$\| \mathbf{A} \mathbf{c} - \mathbf{b} \|^2 + \lambda \mathbf{c}^T \mathbf{U} \mathbf{c}, \quad (10)$$

where  $\mathbf{c}$  and  $\mathbf{b}$ , respectively stand for the vector of unknown coefficients  $c_{kl}$  and the vector of  $m$  data points. The regularization operator  $\mathbf{U}$  is aimed at measuring the total surface smoothness by using direct partial differentiation of Eq. (7). It is constructed from the smoothness measure  $\eta = \mathbf{c}^T \mathbf{U} \mathbf{c}$  given by Dierckx (1993)

$$\eta = \int_{-1}^1 \int_{-1}^1 \left( \left( \frac{\partial f}{\partial u} \right)^2 + \left( \frac{\partial f}{\partial v} \right)^2 \right) du dv \quad (11)$$

and the minimization problem defined by expression (10) can be algebraically solved:

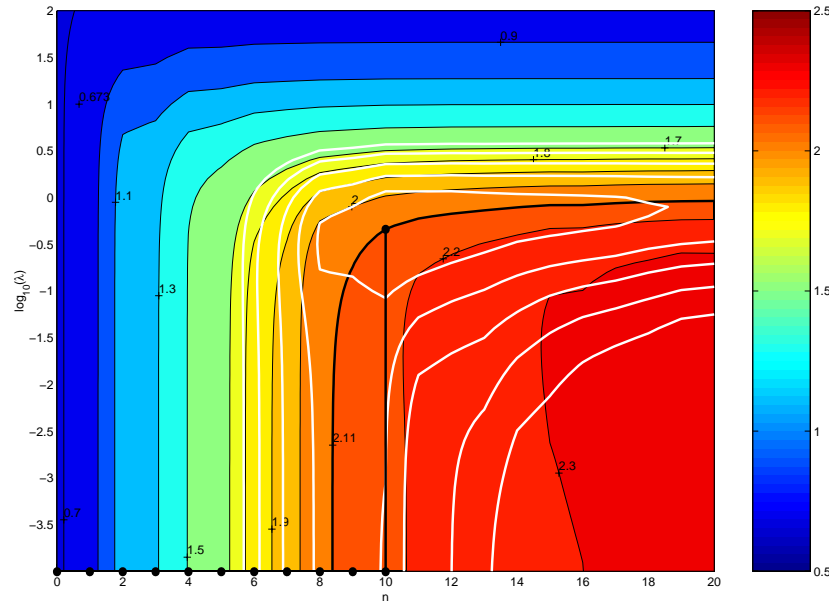
$$\mathbf{c}(\lambda) = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{U})^{-1} \mathbf{A}^T \mathbf{b}. \quad (12)$$

Before investigating examples related to non-uniform sampling in remote sensing, it is worth illustrating the full algorithm by a simple test case (test case 1 = ‘‘TC1’’) that is represented in Fig. 3. A synthetic 2-D field  $f_0(x, y)$  consisting of 2 Gaussian functions of amplitude 1 has been perturbed by random Gaussian noise with a standard deviation of  $s_0 = 0.2$ . The noisy field  $f_i(x, y)$  was sampled by a set of 400 measurements at random locations in the considered domain. The goal of the algorithm is twofold: to reconstruct the unperturbed field at the sampling points and to construct the best possible gridded field everywhere in the domain. If we call  $f_1$  and  $f_g$  the MCR estimated values at the  $m$  sampling locations and at the  $m_g$  grid nodes, respectively, the performance can be estimated by means of the following quantities:

$$s = \left[ \frac{\sum_{i=1}^m (f_1(i) - f_i)^2}{m} \right]^{1/2} \quad (13)$$

$$s_1 = \left[ \frac{\sum_{i=1}^m (f_1(i) - f_0(i))^2}{m} \right]^{1/2} \quad (14)$$

$$s_g = \left[ \frac{\sum_{g=1}^{m_g} (f_g(g) - f_0(g))^2}{m_g} \right]^{1/2}, \quad (15)$$



**Fig. 4.**  $Q$  (Eq. 8) as a function of  $n$  and  $\log(\lambda)$  for test case TC1. Black isopleths and color scale refer to the value of  $Q$ , with optimal value for  $Q_+ = 2.11$  (black bold isopleth). White isopleths: fit error  $s_1(n, \lambda)$  (Eq. 14) showing the clear existence of a minimum. Connected black dots: the path followed by the MCR algorithm.

where  $s$  is the estimation of the noise RMS (to be compared with the true value  $s_0$ ),  $s_1$  is the RMS error with respect to the true field at the sampling locations (referred to by the summation index  $i$ ) and  $s_g$  is the RMS error with respect to the true field at the grid nodes (index  $g$ ).

Like in the 1-D case, the algorithm proceeds by increasing the value of  $n$  until the conservative value  $Q_+$  is reached (Fig. 4). At this stage,  $n$  may be safely incremented by one or two, in order to ensure a maximal sensitivity to data because the regularization in Eq. (12) is not yet working ( $\lambda = 0$ ). The isopleths associated with  $Q(c(n, \lambda))$  exhibit a characteristic corner-shaped form that expresses the progressive transition from a regime of high sensitivity to data and low smoothness, to a regime of low sensitivity and high smoothness. The optimal estimation of  $f_0$ , if it exists, should lie somewhere in the corner region. Therefore, keeping  $n$  constant, the algorithm iteratively increases the value of  $\lambda$  by means of a standard zero finding numerical scheme and computes  $Q(c(n, \lambda))$  at each step until the  $Q_+$  isopleth is reached again, which will always be possible due to the corner-shaped form of the curve. It should be noticed that all  $(n, \lambda)$  solutions lying in the corner region can be considered as equivalent and the present algorithm is a simple way to discover a cheap solution (with a low  $n$  value). Indeed, augmenting  $n$  along the  $Q_+$  isopleth would increase the computational cost without reducing  $s_1$  and, eventually, the ill-conditioning would become redhibitory. Also, the technique to discover a quasi-optimal  $(n, \lambda)$  doublet is not unique. For instance, it would also be possible to detect the maximal curvature along the  $Q_+$  isopleth when both  $n$  and its associated  $\lambda$  value are varied, as done in many regularization problems (Hansen, 1992). However, this procedure turns out to be more expensive.

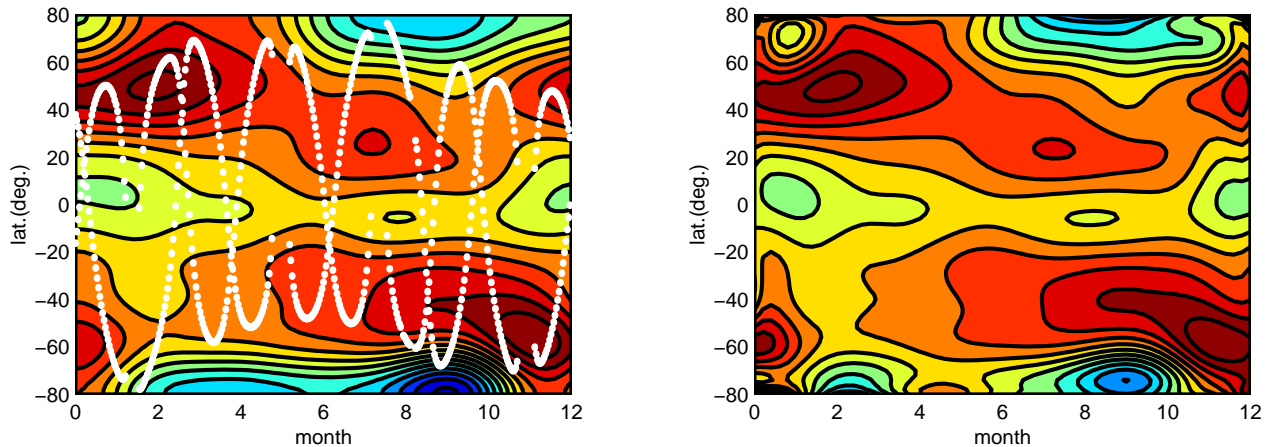
In Fig. 3, it can be seen that the proposed MCR algorithm is quite effective ( $s_1 \simeq s_0/4$ ) in recovering the unperturbed field in presence of important noise.

#### 4 Three competing methods

In order to evaluate the efficiency of the MCR algorithm, it is worth comparing its performance with three other published algorithms applied to the same test cases.

The first one is the Loess (Cleveland and Devlin, 1988) method (referred to as LR,) in which the smoothing is locally performed by using a bivariate cubic polynomial fitted to an adjustable number of nearest neighbours. In the surface fitting commercial software (TableCurve 3D<sup>TM</sup>) that we have used, interpolation on the regular grid is performed by locally-weighted fitting with the Renka triangulation-based procedure (Renka, 1996) applied to the smoothed values. When the number of neighbours is interactively increased by the user, the degree of smoothing varies from overfitting to underfitting. In the comparisons hereafter, we have selected the optimal number that produces an exact value for the estimation of  $s_0$ . This gives a lower bound of the smoothing error caused by the LR method.

The second method is the biharmonic spline scheme (BS) proposed by Sandwell (1987), in which Green functions of the biharmonic operator are used for minimum curvature interpolation of irregularly spaced data points. The interpolating surface is a linear combination of Green functions centered at each data point. By reducing the number of model parameters, noisy data can be fitted in a least-squares sense due to the minimal curvature property of any solution of the biharmonic equation. The interpolation scheme



**Fig. 5.** Test case TC2: Left: Exact ozone VMR field at 30 mbar (color scale maps linearly from 2 to 5 ppm) with superimposed SAGE II geolocations. Right: Reconstructed field by MCR ( $n = 14$ ).

**Table 1.** Results of test cases TC1, TC2, TC3 for the LR, BS, K and MCR methods. The value of  $s$  is normalized by the exact value for the considered test. Value of  $s_1$  and  $s_g$  are normalized to the best value of the respective test

	$s$	$s$	$s$	$s$	$s_1$	$s_1$	$s_1$	$s_1$	$s_g$	$s_g$	$s_g$	$s_g$
	LR	BS	K	MCR	LR	BS	K	MCR	LR	BS	K	MCR
TC1	-	-	<b>.988</b>	.970	1.270	-	1.083	<b>1.000</b>	1.460	3.679	1.048	<b>1.000</b>
TC2	-	-	.906	<b>.922</b>	1.284	-	1.055	<b>1.000</b>	1.235	1.782	1.029	<b>1.000</b>
TC3	-	-	.984	<b>.994</b>	1.095	-	<b>1.000</b>	1.082	<b>1.000</b>	2.359	1.082	1.033

reproduces the input data at the sampling points so that the smoothing error has only been computed at the regular grid nodes. For this comparison test, we have used the routine “grid data” as implemented in the numerical software package MATLAB<sup>TM</sup> (release 12).

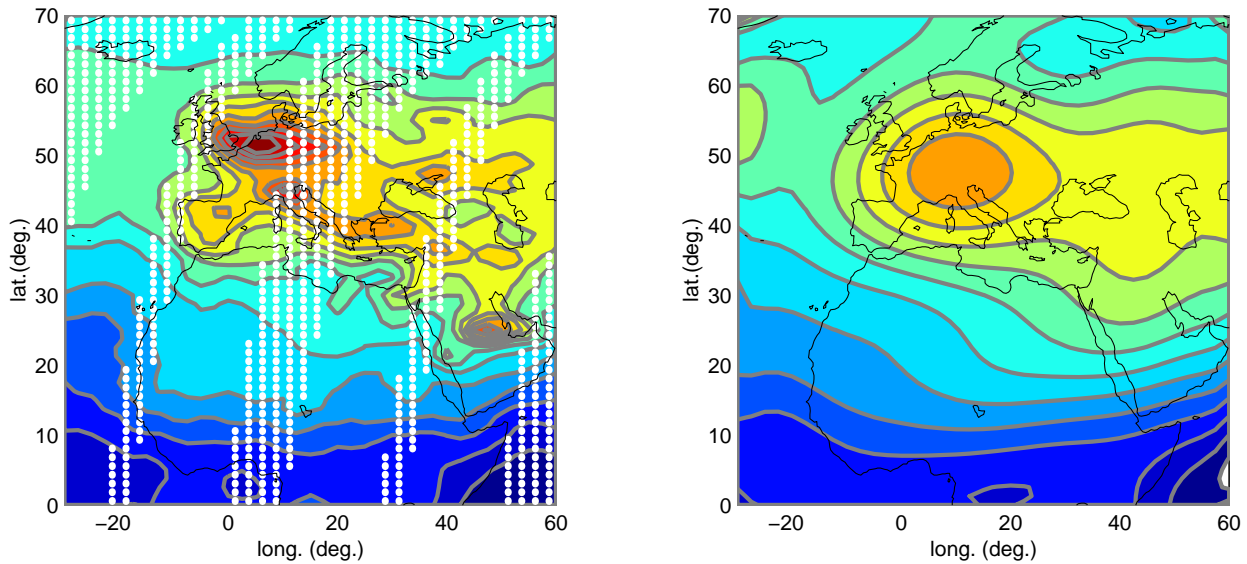
The third competing method is “kriging” which has been extensively used in geophysics (see, e.g. Tranchant and Vincent, 2000). Briefly, kriging is a minimum variance method closely connected with statistical interpolation techniques (Daley, 1991). The analyzed field is considered as being reducible to the sum of a mean value, the non-stationary expectation of the field, and a random fluctuation with a zero expectation. A predictor is then constructed from a linear combination of data values whose coefficients minimize the estimation variance, with a constraint ensuring that the estimator is unbiased. The kriging system of equations is found to relate the optimal coefficient set and the covariance between the sample points, which is assumed to depend only on the interdistance. The covariance is then computed by means of a “semi-variogram” and fitted by a predefined theoretical model. In the comparison exercise, we have used the EasyKrig code (version 2.1) developed by Chu (2000). We also made the assumption that the sampled geophysical field was quasi-stationary, which allows one to use ordinary kriging. In case of suspected non-stationarity, universal kriging should be preferred. Kriging is a powerful statistical method, and robust techniques have been proposed to fit

the semi-variograms (Journel and Huijbregt, 1989), including the possible use of anisotropic forms if privileged directions are expected (see Tranchant and Vincent, 2000, for an extended discussion). Here, however, we do assume that the use of an Gaussian-linear isotropic semi-variogram is accurate enough. In Table 1, we present the results of the comparison of  $s$ ,  $s_1$  and  $s_g$  (Eqs. 13–15) obtained by MCR and by the competing methods for test cases TC1, TC2, TC3 (the latter two are described hereafter). For  $s_1$  and  $s_g$ , reported values have been normalized to the best score for each test case, whereas for  $s$  the normalization was performed with respect to the exact value. For TC1, MCR is clearly superior to the LR and BS methods. It is also slightly superior to kriging, except for the estimation of noise standard deviation.

## 5 Two geophysical examples in the Earth’s remote sounding by space experiments

The second test case (“TC2”) concerns the reconstruction of the ozone field at 30 mbar in a time-latitude domain. The example has been conceived from geolocations inspired by the well-known SAGE II experiment (Chu et al., 1989; Mauldin III L. E. et al., 1985), which has furnished an invaluable series of atmospheric altitude profiles of ozone, nitrogen dioxide, aerosols and water vapour. SAGE is a typical solar occultation experiment that was launched in October 1984 on the Earth Radiation Budget Satellite into a 56 degree





**Fig. 6.** Test case TC3: Left: Exact  $\text{NO}_2$  columns (color scale maps linearly from 2 to  $5 \cdot 10^{15} \text{ mol cm}^{-2}$ ) with superimposed GOME pixels for 15 Sep 1999. Right: Reconstructed field by MCR ( $n = 10$ ).

inclination orbit. It makes measurements primarily at middle and low latitudes, but reaches the polar regions in both hemispheres several times a year due to the combination of the orbital plane precession and of the seasonal effect (see Fig. 5). About 15 sunset and 15 sunrise measurements are performed per day at a quasi-constant latitude. Over the period of about 1 month, the latter slowly varies along sinusoidal tracks and extends over a seasonally dependent altitude of approximately  $70^\circ \text{ S}$  to  $70^\circ \text{ N}$ . In test case TC2, we arbitrarily selected the SAGE geolocations of year 1992.

The assumed geophysical field of interest is the ozone volume mixing ratio (VMR) and, for the sake of realism, we have constructed the “exact” field from the ozone climatology published by Fortuin and Kelder (1998). Daily ozone VMR have been averaged and they were assigned a measurement random error of 2%, based on the published values of Cunnold et al. (1989). From a data set of 638 measurements, the test consisted in the reconstruction of the ozone field on a regular grid having a  $0.2 \text{ month} \times 2 \text{ degree}$  resolution. It can be seen in Fig. 5 that the MCR method performs well (the optimal order was found to be  $n = 14$  with a total of 120 coefficients) and that fine structures are correctly extracted within domains of missing data (e.g. around months 5 and 7 in tropical regions). Also, at high latitudes, partial information about well-known ozone minima is extracted by the algorithm, although data are scarce or even inexistent. Even by using a basis of spherical harmonics (or geodesic distances for kriging), it is very important to keep in mind that an interpolating/smoothing algorithm is not able to produce information not measured by the experiment. An important feature of ozone field is located near the edge of the south polar vortex (see, for instance, Wauben et al., 1997; McIntyre, 1989; Vincent and Tranchant, 1999) which is exterior to the SAGE sampling domain. Nevertheless, MCR

gives valuable, although incomplete information, over this region by smoothly extrapolating the data points. In Table 1, MCR turns out to be the best estimator of  $s$ ,  $s_1$  and  $s_g$ .

The third test case (“TC3”) deals with  $\text{NO}_2$  columns retrieved from the Global Ozone Monitoring Experiment (GOME) on board the ESA ERS-2 satellite (Burrows et al., 1999) that was launched into a polar heliosynchronous orbit. The instrument is a nadir-viewing UV-visible spectrometer and allows for a global coverage in about 3 days. The test case has been built by assuming that the true geophysical field is the monthly averaged  $\text{NO}_2$  columns for September 1999 above Europe and North Africa (longitude  $30^\circ \text{ W}$ – $60^\circ \text{ E}$ , latitude  $0^\circ$ – $70^\circ \text{ N}$ ) (see Fig. 6). According to Lambert et al. (1999), a random retrieval error of 15% has been added to produce the measured data for one observation day (15 September 1999).

Test case TC3 consists of recovering the true field from data (820 values) representing about one-third of the full coverage, with clustered pixels along the satellite tracks. This was achieved by using MCR with an optimal  $n = 10$  (66 coefficients) and reconstruction was performed on a grid having a  $2.5 \text{ deg.}$  (in longitude)  $\times$   $1.25 \text{ deg.}$  (in latitude) resolution. By inspecting Fig. 6, it is clear that the reconstruction is not as accurate as for TC2 case. For instance, the pollution peak above western Europe is well identified but not separated into local maxima above Germany and North Italy. Fine structures above the Black Sea have also been smoothed out. Actually, the lower spatial resolution obtained by the reconstruction algorithm is a consequence of the quite high noise level and not of the sampling geometry. In comparison with alternative methods (Table 1), it is interesting to observe that MCR performs slightly worse than KR at the sampling locations, due to the locally very high correlation between GOME adjacent pixels, whereas it is better



for interpolating/extrapolating on the regular grid. For test case TC3, LR gives the best results on the grid, if we remember that it should be fed with the exact noise estimation that could be delivered by MCR.

## 6 Conclusions

Surface fitting of noisy data from scattered locations can be performed by various methods. The final accuracy of the reconstruction obviously depends on the algorithm intrinsic error, the geometrical distribution of samples and the noise level of the observations.

In this work, we have developed a new algorithm capable of producing an optimal fit in the sense that a trade-off is achieved between underfitting where information is lost and overfitting where spurious structures are induced by experimental noise. The 2-D MCR method is based on minimal correlation of minimal residuals, associated to a standard regularization technique. The method is fast, quite robust and easy to implement.

MCR performs well with respect to existing methods for geophysical applications. In particular, we have tested its efficiency for three kinds of sampling geometries: a randomly dispersed data set (TC1), a line type data set (TC2) and a clustered data set (TC3) for different levels of noise. The method is, therefore, well suited for interpolation in validation campaigns or for the construction of climatological fields.

In future work, we will investigate the generalization of the algorithm to higher dimensions. Clear geophysical objectives will include interpolation between vertical profiles observed at different locations and temporal evolution of 2-D fields.

*Acknowledgements.* This work was partly performed within projects “SADE (Prodex 6)” and “Measurement, understanding and climatology of stratospheric aerosols” funded by the SSTC/DWTC service of the Belgian Government. The author also thanks his colleagues J.-C. Lambert, J. Granville and P. Gérard for having provided him with GOME data.

Topical Editor O. Boucher thanks two referees for their help in evaluating this paper.

## References

- Burrows, J. P., Weber, M., Buchwitz, M., Rozanov, V., Weissenmayer, A., Richter, A., DeBeek, R., Hoogen, R., Bramsted, K., Eichmann, K., and Eisinger, M.: The global ozone monitoring experiment (GOME): Mission concept and first scientific results, *J. Atmos. Sci.*, 56, 151–175, 1999.
- Chu, D.: The globec kriging software package – easykrig 2.1, Tech. rep., Woods Hole Oceanographic Institution, [ftp://globec.whoi.edu/pub/software/kriging/easy\\_krig/](ftp://globec.whoi.edu/pub/software/kriging/easy_krig/), 2000.
- Chu, W. P., McCormick, M. P., Lenoble, J., Brogniez, C., and Pruvost, P.: SAGE II Inversion Algorithm, *J. Geophys. Res.*, 94, 8839–8851, 1989.
- Cleveland, W. S. and Devlin, S. J.: Locally weighted regression: An approach to regression analysis by local fitting, *J. Am. Stat. Assoc.*, 83, 596–610, 1988.
- Cunnold, D. M., Chu, W. P., Barnes, R. A., McCormick, M. P., and Veiga, R. E.: Validations of SAGE II ozone measurements, *J. Geophys. Res.*, 94, 8447–8460, 1989.
- Daley, R.: *Atmospheric data analysis*, Cambridge University Press, 1991.
- Dierckx, P.: *Curve and Surface Fitting with Splines*, Oxford Science Publications, 1993.
- Durbin, J. and Watson, G. S.: Testing for serial correlation in least-squares regression, *Biometrika*, 37, 409–428, 1950.
- Fisher, M. and Lary, D. J.: Lagrangian 4-dimensional variational data assimilation of chemical species, *Q. J. R. Meteorol. Soc.*, 121, 1681–1704, 1995.
- Foley, T. A. and Hagen, H.: Advances in scattered data interpolation, *Surv. Math. Ind.*, 4, 71–84, 1994.
- Fortuin, J. P. F. and Kelder, H.: An ozone climatology based on ozonesonde and satellite measurements, *J. Geophys. Res.*, 103, 31 709–31 734, 1998.
- Hansen, P. C.: Numerical tools for analysis and solution of Fredholm integral equations of the first kind, *Inverse problems*, 8, 849–872, 1992.
- Journel, A. G. and Huijbregt, J.: *Mining geostatistics*, Academic Press, London, 1989.
- Lait, L.: Effects of satellite scanning configurations on derived gridded fields, *J. Geophys. Res.*, 105, 9063–9074, 2000.
- Lambert, J.-C., Granville, J., Roozendaal, M., Müller, J.-F., Pommereau, J.-P., Goutail, F., and Sarkissian, A.: A pseudo-global correlative study of ERS-2 GOME NO<sub>2</sub> data with ground-, balloon-, and space-based observations, in *Proc. European Symposium on Atmospheric Measurements from Space (ESAMS) / ESA WPP-161*, 1, pp. 217–224, ESA/ESTEC, The Netherlands, 1999.
- Mauldin III, L. E., Zaun, N. H., McCormick, M. P., Guy, J. H., and Vaughn, W. R.: Stratospheric Aerosol and Gas Experiment II Instrument: A Functional Description, *Opt. Eng.*, 24, 307–312, 1985.
- McIntyre, M. E.: On the antarctic ozone hole, *J. Atmos. Terr. Phys.*, 51, 29–43, 1989.
- Press, W. H., Teukolsky, S. A., Vetterling, W. T., and Flannery, B. P.: *Numerical Recipes in FORTRAN*, Second Edition, Cambridge University Press, Cambridge, 1992.
- Renka, R.: Algorithm 752: Software for scattered data fitting with a constrained surface under tension, *ACM transactions on mathematical software*, 22, 9–17, 1996.
- Sandwell, D. T.: Biharmonic spline interpolation of GEOS-3 and SEASAT altimeter data, *Geophys. Res. Lett.*, 14, 139–142, 1987.
- Thijsse, B. J., Hollanders, M. A., and Hendrikse, J.: A practical algorithm for least-squares spline approximation of data containing noise, *Computers in physics*, 12, 393–399, 1998.
- Tranchant, B. J. S. and Vincent, A. P.: Statistical interpolation of ozone measurements from satellite data (TOMS, SBUV and SAGE II) using the kriging method, *Ann. Geophysicae*, 18, 666–678, 2000.
- Vincent, A. P. and Tranchant, B. J.: Anisotropic turbulent diffusion for ozone transport at 520 K, *J. Geophys. Res.*, 104, 27 209–27 215, 1999.
- Wauben, W. M., Bintanja, R., van Velthoven, P. F., and Kelder, H.: On the magnitude of transport out of the antarctic vortex, *J. Geophys. Res.*, 102, 1229–1238, 1997.