**Atmospheric
Chemistry
and Physics
Discussions**

# Technical Note: Trend estimation from irregularly sampled, correlated data

**T. von Clarmann, G. Stiller, U. Grabowski, and J. Orphal**

Karlsruhe Institute of Technology, Institute for Meteorology and Climate Research, Karlsruhe, Germany

Correspondence to: T. von Clarmann (thomas.clarmann@kit.edu)

Full Screen / Esc

Printer-friendly Version

Interactive Discussion

**Abstract**

Estimation of a trend of an atmospheric state variable is often performed by fitting a linear regression line to a set of data of this variable sampled at different times. Often these data are irregularly sampled in space and time and clustered in a sense that error correlations among data points cause a similar error of data points sampled at similar times. Since this can affect the estimated trend, we suggest to take the full error covariance matrix of the data into account. Superimposed periodic variations can be jointly fitted in a straight forward manner, even if the shape of the periodic function is not known. Global data sets, particularly satellite data, can form the basis to estimate the error correlations.

## 1  Introduction

Correct trend estimation is a key question in the discussion of climate change (IPCC, 2007). While fitting a straight line to a sample of data is an almost trivial task, errors in the data set and non-representativeness of the sample add some difficulty to the problem. Assuming normally distributed errors which are uncorrelated over the sample, each data point is simply weighted by the inverse of its variance to obtain a best linear unbiased estimated of the trend (Aitken, 1935). Methods applicable to least squares fitting of data where both the dependent and the independent variables are affected by errors have recently been reviewed by Cantrell (2008).

If the assumption of normal error distribution is questionable, robust linear regression methods help to reduce the sensitivity of the trend to outliers in the sample (Muhlbauer et al. 2009 and references therein). Another cure against non-normality of distributions of residuals are bootstrap methods, introduced by Efron (1979) as a variant to Jackknife methods and applied to atmospheric trend analysis by, e.g. Cox et al. (2002), Gardiner et al. (2008) or Vigouroux et al. (2008).

**Trend estimation from clustered data**

T. von Clarmann et al.

Besides non-normality of the distribution of residuals, correlations between the sampled data are another class of problems. When using multisite means to infer a trend, the standard errors of the means $\sigma_{\text{mean}}$ which determine the weight of each mean in the regression analysis are not the standard deviation $\sigma$ of the sample over the sites divided by the square root of the number of sites n but

$$\sigma_{\text{mean}} = \sqrt{\sigma^2 \left( \frac{1 + (n-1)\bar{r}}{n} \right)}, \tag{1}$$

where $\bar{r}$ is the average intersite correlation coefficient (Jones et al., 1997). This can easily be verified by multiplication of the averaging operator from the left and right to the intersite covariance matrix $\mathbf{S}_i$ according to multivariate Gaussian error propagation:

$$\sigma^2_{\text{mean}} = (\frac{1}{n}, \ldots, \frac{1}{n})\mathbf{S}_i \begin{pmatrix} \frac{1}{n} \\ \vdots \\ \frac{1}{n} \end{pmatrix} \tag{2}$$

This approach solves the problem of intersite correlations and is applicable, e.g., if measurements of the same set of sites are used over the whole period. $\sigma_{\text{mean}}$ calculated under consideration of $\bar{r}$ accounts for the fact that the available sites do not fully represent the population, i.e., the sample mean at a given time is not necessarily identical to the global mean. Since the same set of stations is used over the whole period, the measurements at the given sites are not a random sample.

Weatherhead et al. (1998) discuss how autocorrelations of noise in the data affect the precision of the estimated trend, and they provide a practicable method to consider these autocorrelations to avoid over-optimistic confidence estimated with respect to inferred linear trends. Further, these authors present a tool to estimate the required length of the time series to significantly detect a trend.
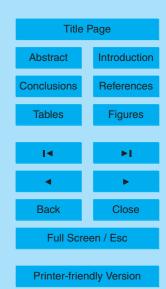
None of these papers, however, tackles the problem how to derive trends on the basis of inhomogeneous data sets. In this technical note, we investigate the problem that

the sampled data are clustered in a sense that the data are groupwise correlated in the time domain, i.e. that data inhomogeneities cause systematic deviations between subsets of the data of time series which, through irregular sampling in time, map onto the time series as errors correlated in the time domain. Since these errors are systematic and are not autoregressive, the autocorrelation concept discussed above does not help. A typical example would be the estimation of a trend of one atmospheric state variable from measurements at two different latitudes, where one measurement site dominates the early part and the other measurement site the later part of the time series. The neglected latitudinal dependence of the observed quantity maps onto the time domain if the atmosphere is irregularly sampled at the different observation sites. Such data sets, where the target variable depends on further variables except the independent variable of the regression analysis, we call inhomogeneous. Irregular sampling of inhomogeneous data leads to clustering, because certain values of the independent variable may go along with certain values of the hidden variable. This dependence can be formulated as correlations, typically the larger, the more similar the value of the hidden variable is. These correlations will, if neglected, not only render the significance analysis of the trend insignificant but can actually change the slope of the regression line, i.e. lead to different trends.

While the proposed concept is quite straightforward rather than novel, we hope that it may be useful to the climate research community where currently error covariances in irregularly sampled data often seem to be ignored, even when inhomogeneous datasets are analyzed.

## 2   Linear trends of clustered data

Assuming a linear trend, we can approximate the temporal development of an atmospheric state variable $y$ as a straight line. A straight line is defined as

$$\hat{y}(x; a, b) = a + bx, \tag{3}$$

where the ˆ symbol indicates a modeled or estimated rather than a measured state variable. In our application $x$ is the time of the measurement, but this concept of regression of clustered data is applicable to a wider context.

For normally distributed but possibly interdependent errors of $y_i$, $i = 1\ldots n$, of which the ex ante[1] estimates are represented by the $n \times n$ covariance matrix $\mathbf{S}_y$, this straight line is the optimal regression line, when the cost function

$$\chi^2 = (\mathbf{y} - (a\mathbf{e} + b\mathbf{x}))^T \mathbf{S}_y^{-1}(\mathbf{y} - (a\mathbf{e} + b\mathbf{x})) \tag{4}$$

is minimum, where $\mathbf{e} = (1,\ldots,1)^T$ and $\mathbf{x} = (x_1,\ldots,x_n)^T$, $\mathbf{y} = (y_1,\ldots,y_n)^T$, and $^T$ denotes the transpose of a matrix. Coefficients $a$ and $b$ are inferred in a well established manner by setting the derivatives $\partial\chi^2/\partial a$ and $\partial\chi^2/\partial b$ to zero. This gives

$$\frac{\partial\chi^2}{\partial a} = -2\mathbf{e}^T\mathbf{S}_y^{-1}(\mathbf{y} - a\mathbf{e} - b\mathbf{x}) = 0; \tag{5}$$

$$\mathbf{e}^T\mathbf{S}_y^{-1}\mathbf{y} = \mathbf{e}^T\mathbf{S}_y^{-1}a\mathbf{e} + \mathbf{e}^T\mathbf{S}_y^{-1}b\mathbf{x};$$

$$a = \frac{\mathbf{e}^T\mathbf{S}_y^{-1}\mathbf{y} - \mathbf{e}^T\mathbf{S}_y^{-1}b\mathbf{x}}{\mathbf{e}^T\mathbf{S}_y^{-1}\mathbf{e}}$$

and

$$\frac{\partial\chi^2}{\partial b} = -2\mathbf{x}^T\mathbf{S}_y^{-1}(\mathbf{y} - a\mathbf{e} - b\mathbf{x}) \tag{6}$$

$$= \mathbf{x}^T\mathbf{S}_y^{-1}\mathbf{y} - \mathbf{x}^T\mathbf{S}_y^{-1}a\mathbf{e} - \mathbf{x}^T\mathbf{S}_y^{-1}b\mathbf{x} = 0;$$

$$\mathbf{x}^T\mathbf{S}_y^{-1}\mathbf{y} = \mathbf{x}^T\mathbf{S}_y^{-1}b\mathbf{x} + \mathbf{x}^T\mathbf{S}_y^{-1}a\mathbf{e}.$$

---

[1]Ex ante error estimates we call error estimates based on propagation of assumed primary errors through the system and can be calculated before the measurement actually has been made, as opposed to ex post error estimates which are based on the standard deviation of a sample of measurements (von Clarmann, 2006).

T. von Clarmann et al.

Printer-friendly Version

Interactive Discussion

Combining Eqs. 5 and 6 gives

$$x^T \mathbf{S}_y^{-1} y = x^T \mathbf{S}_y^{-1} bx + x^T \mathbf{S}_y^{-1} ae \tag{7}$$

$$= x^T \mathbf{S}_y^{-1} bx + x^T \mathbf{S}_y^{-1} e \frac{e^T \mathbf{S}_y^{-1} y - e^T \mathbf{S}_y^{-1} bx}{e^T \mathbf{S}_y^{-1} e}.$$

This can be rearranged as

$$5 \quad x^T \mathbf{S}_y^{-1} bx - \frac{x^T \mathbf{S}_y^{-1} e e^T \mathbf{S}_y^{-1} bx}{e^T \mathbf{S}_y^{-1} e} = \tag{8}$$

$$x^T \mathbf{S}_y^{-1} y - \frac{x^T \mathbf{S}_y^{-1} e e^T \mathbf{S}_y^{-1} y}{e^T \mathbf{S}_y^{-1} e}$$

and finally solved to give $b$:

$$b = \frac{x^T \mathbf{S}_y^{-1} y - \frac{x^T \mathbf{S}_y^{-1} e e^T \mathbf{S}_y^{-1} y}{e^T \mathbf{S}_y^{-1} e}}{x^T \mathbf{S}_y^{-1} x - \frac{x^T \mathbf{S}_y^{-1} e e^T \mathbf{S}_y^{-1} x}{e^T \mathbf{S}_y^{-1} e}} \tag{9}$$

$$= \frac{x^T \mathbf{S}_y^{-1} y e^T \mathbf{S}_y^{-1} e - x^T \mathbf{S}_y^{-1} e e^T \mathbf{S}_y^{-1} y}{x^T \mathbf{S}_y^{-1} x e^T \mathbf{S}_y^{-1} e - x^T \mathbf{S}_y^{-1} e e^T \mathbf{S}_y^{-1} x}.$$

10 Inserting this into Eq. 5 allows to calculate $a$:

$$a = \frac{x^T \mathbf{S}_y^{-1} y - e^T \mathbf{S}_y^{-1} x b}{e^T \mathbf{S}_y^{-1} e} \tag{10}$$

$$= \frac{e^T \mathbf{S}_y^{-1} y - e^T \mathbf{S}_y^{-1} x \frac{x^T \mathbf{S}_y^{-1} y e^T \mathbf{S}_y^{-1} e - x^T \mathbf{S}_y^{-1} e e^T \mathbf{S}_y^{-1} y}{x^T \mathbf{S}_y^{-1} x e^T \mathbf{S}_y^{-1} e - x^T \mathbf{S}_y^{-1} e e^T \mathbf{S}_y^{-1} x}}{e^T \mathbf{S}_y^{-1} e}$$

For unity $\mathbf{S}_y$ this reduces to the widely used parameters $\tilde{a}$ and $\tilde{b}$ of a regression line for data points of uncorrelated errors of equal variance:

$$\tilde{a} = \frac{\sum y_i}{n} - \tilde{b}\frac{\sum x_i}{n}, \tag{11}$$

where

$$\tilde{b} = \frac{n\sum x_i y_i - \sum x_i \sum y_i}{n\sum x_i^2 - (\sum x_i)^2} \tag{12}$$

The uncertainty of the slope $b$ is:

$$\sigma_b^2 = \left(\frac{x^T\mathbf{S}_y^{-1}e^T\mathbf{S}_y^{-1}e - x^T\mathbf{S}_y^{-1}ee^T\mathbf{S}_y^{-1}}{x^T\mathbf{S}_y^{-1}xe^T\mathbf{S}_y^{-1}e - x^T\mathbf{S}_y^{-1}ee^T\mathbf{S}_y^{-1}x}\right)^T \cdot$$
$$\mathbf{S}_y\left(\frac{x^T\mathbf{S}_y^{-1}e^T\mathbf{S}_y^{-1}e - x^T\mathbf{S}_y^{-1}ee^T\mathbf{S}_y^{-1}}{x^T\mathbf{S}_y^{-1}xe^T\mathbf{S}_y^{-1}e - x^T\mathbf{S}_y^{-1}ee^T\mathbf{S}_y^{-1}x}\right) \tag{13}$$

From comparison of Eqs. (9) and (12) we see that the error correlations do not only change the estimated error of the trend but also affect the trend itself, e.g. rotate the regression line.

## 3  Estimation of measurement covariances

Evaluation of Eq. (9) requires knowledge of the covariance matrix $\mathbf{S}_y$. While for some error sources such error assumptions are available and reasonable assumptions on correlations within a class of measurements can be made (e.g. perfect correlation, i.e., $r = 1$ for the calibration error component within all measurements based on the same calibration standard may be reasonable), for evaluation of error covariances representing other error sources, external data may be needed. Typical error correlations in a

time series can be caused by the fact that the sample is composed of measurements at various locations. If the mean measurement times at two locations differ, any difference in the expectation value of the state variable with, e.g., latitude, will map onto the trend. If the latitudinal dependence is too complicated for a simple correction, the related error correlation should at least be included in the covariance matrix $\mathbf{S}_y$.

A source for correlation information are satellite measurements. While satellite data sets often cover only a few years and thus are often not suitable to infer trends, they, in contrast to station measurements or balloon measurements, in some cases cover the globe densely enough for assessment of spatial variability.

If no regression model is available to correct inhomogeneous station data, or if the use of such a regression model seems not justified because the inferred model parameters may not be sufficiently representative, the satellite data still can be used to estimate non-representativeness of the station data in terms of a covariance matrix. Let I be the number of global satellite measurements available, $\boldsymbol{u}_i$ be the global (or at least multi-site[2]) field associated with the $i$th measurement, then the component of the covariance matrix accounting for the representativeness error is

$$\frac{1}{I-1}\sum_{i=1}^{I}(\boldsymbol{u}_i - \bar{\boldsymbol{u}})(\boldsymbol{u}_i - \bar{\boldsymbol{u}})^T - \mathbf{S}_{\text{sat}}, \tag{14}$$

where $\bar{\boldsymbol{u}}$ is the average global field as measured by the satellite,

$$\bar{\boldsymbol{u}} = \frac{1}{I}\sum_{i=1}^{I}\boldsymbol{u}_i, \tag{15}$$

and $\mathbf{S}_{\text{sat}}$ is the covariance matrix characterizing the measurement error of the satellite data. A precondition to this approach is that the time window covered by the satellite data used to infer the representativeness error covariance matrix is small enough that neglect of trends within this time window is justified.

---

[2]Instead of full global fields it is actually sufficient to include only locations for which station measurements are available; each component of the vector represents one geolocation.

**Trend estimation from clustered data**

T. von Clarmann et al.

◄◄ | ►►

◄ | ►

Back | Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion

## 4   Consideration of the annual cycle and related problems

A linear trend may be superimposed with a periodic function of known periodic time, e.g. diurnal or seasonal variation, etc. There are several options to tackle this problem.

If the sample is large enough, the linear trend can be evaluated for subsets of data recorded at the same phase of the periodical variation, and the overall trend is calculated as an optimal mean of the individual trends. This requires binning of data; in the case of a seasonal cycle, the linear trend can be estimated as the mean of the trend over all Januaries, Februaries, etc. Problems occur when the amplitude of the seasonal cycle has a trend in itself and the whole observation period is not an integer multiple of the time of one cycle.

Another option is subtraction of the periodic signal prior to trend estimation. The periodic component of the signal can either be imported from an external source (model, independent data) or from the sample itself. The latter approach is not quite trivial, because the inferred mean periodical signal will, in turn, depend on the trend (periodic analysis usually is defined only for stationary time series, i.e. zero trend), such that either an iterative approach or a multivariate optimization (see below) is required. Care has to be taken to consider the reduction of degrees of freedom implied by inferring the correction from the data themselves.

The problem of the non-stationary nature of time series, which is by definition inherent in trend analysis, can be solved by retrieving the trend, the amplitude of the periodic variation, and possibly the phase and the shape of the oscillation in one step. In the case of a known function of unknown amplitude (e.g. sine) the amplitude can be fitted along with the trend. In the case of unknown phase, it is usually more appropriate to fit amplitudes of a sine and a cosine of the same period length rather than the amplitude and the phase, in order to keep the fit linear. A regression model involving a linear trend superimposed with a single harmonic variation of unknown phase but known period length $l$ is written as

$$\hat{y}(x; a, b, c, d) = a + bx + c\sin\frac{2\pi x}{l} + d\cos\frac{2\pi x}{l}. \tag{16}$$

Setting the partial derivatives of

$$\chi^2 = (\boldsymbol{y}(\boldsymbol{x}) - \hat{\boldsymbol{y}}(\boldsymbol{x}))^T \mathbf{S}_y^{-1} (\boldsymbol{y}(\boldsymbol{x}) - \hat{\boldsymbol{y}}(\boldsymbol{x})) \tag{17}$$

with respect to the parameters of the regression model to zero gives

$$\frac{\partial \chi^2}{\partial a} = -2\boldsymbol{e}^T \mathbf{S}_y^{-1}. \tag{18}$$

$$\left( \boldsymbol{y}(\boldsymbol{x}) - a\boldsymbol{e} - b\boldsymbol{x} - c\sin\frac{2\pi\boldsymbol{x}}{l} - d\cos\frac{2\pi\boldsymbol{x}}{l} \right) = 0$$

$$\frac{\partial \chi^2}{\partial b} = -2\boldsymbol{x}^T \mathbf{S}_y^{-1}. \tag{19}$$

$$\left( \boldsymbol{y}(\boldsymbol{x}) - a\boldsymbol{e} - b\boldsymbol{x} - c\sin\frac{2\pi\boldsymbol{x}}{l} - d\cos\frac{2\pi\boldsymbol{x}}{l} \right) = 0$$

$$\frac{\partial \chi^2}{\partial c} = -2(\sin\frac{2\pi\boldsymbol{x}}{l})^T \mathbf{S}_y^{-1}. \tag{20}$$

$$\left( \boldsymbol{y}(\boldsymbol{x}) - a\boldsymbol{e} - b\boldsymbol{x} - c\sin\frac{2\pi\boldsymbol{x}}{l} - d\cos\frac{2\pi\boldsymbol{x}}{l} \right) = 0$$

$$\frac{\partial \chi^2}{\partial d} = -2(\cos\frac{2\pi\boldsymbol{x}}{l})^T \mathbf{S}_y^{-1}. \tag{21}$$

$$\left( \boldsymbol{y}(\boldsymbol{x}) - a\boldsymbol{e} - b\boldsymbol{x} - c\sin\frac{2\pi\boldsymbol{x}}{l} - d\cos\frac{2\pi\boldsymbol{x}}{l} \right) = 0$$

One might be too lazy to calculate a general algebraic solution for such a system of equations but prefer to solve this numerically for parameters $a \ldots d$ by any appropriate solver provided by the program library at hand. The advantage of this approach is that it does not require a stationary time series to evaluate the amplitudes of the oscillations.

**Trend estimation from clustered data**

T. von Clarmann et al.

Printer-friendly Version

Interactive Discussion

If need be, this type of analysis can also involve seasonal "overtones", i.e. additional periodic functions of periods which are an integer fraction of $l$, which may be made subject to lowpass filtering. This generalization of the schemes presented here to applications with more than one pair of periodic functions is straightforward, and the relationship to harmonic analysis is obvious.

If the shape of the periodic variation is not known a priori, it can be inferred from the data themselves in one step with the trend estimation. For binned data (e.g. when monthly means are used to infer a trend with superimposed seasonal variation) monthly corrections are fitted along with slope and axis intercept. The regression model then is

$$\hat{y} = a + bx + c_{\text{month}}(x), \tag{22}$$

where $c_{\text{month}}$ is the monthly correction applicable to time $x$. The cost function to be minimized for this application is

$$\chi^2 = (\boldsymbol{y}(\boldsymbol{x}) - (a + b\boldsymbol{x} + \boldsymbol{c}^T \mathbf{U}))^T \mathbf{S}_y^{-1}(\boldsymbol{y}(\boldsymbol{x}) - (a + b\boldsymbol{x} + \boldsymbol{c}^T \mathbf{U})) \tag{23}$$

where $\mathbf{U}$ is a selection matrix with all elements in the $i$th row zero except for column $j$, where $j$ represents the month when measurement $x_i$ was made, where the matrix element is one.

If binning or averaging is to be avoided, the periodic correction terms $c$ can also be defined points on the $x$ axis rather than bins. The actual correction for a given $y(x)$ can then be estimated by interpolation, leading to the regression model

$$\hat{y} = a + bx + c_{\text{month}}(x) + \frac{c_{\text{month}+1} - c_{\text{month}}}{d(\text{month})} d(x), \tag{24}$$

where for clarity but without sacrificing a more general applicability of the concept we assume $c_{\text{month}}(x)$ is the periodic correction of the first day of month, $d(\text{month})$ is the number of days of the month, and $d(x)$ is the day of the month. Periodicity is assumed in a sense that month+12=month. The cost function to be minimized for this application has the same structure as Eq. (23):

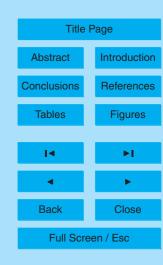$$\chi^2 = (\boldsymbol{y}(\boldsymbol{x}) - (a + b\boldsymbol{x} + \boldsymbol{c}^T \mathbf{V}))^T \mathbf{S}_y^{-1}(\boldsymbol{y}(\boldsymbol{x}) - (a + b\boldsymbol{x} + \boldsymbol{c}^T \mathbf{V})) \tag{25}$$

**Trend estimation from clustered data**

T. von Clarmann et al.

$\mathbf{V}$ is a matrix with all elements in the $i$th row zero except for column $j$ and $j + 1$ (or 1, if $j$ denotes the last column), where $j$ and $j + 1$ represent the month when measurement $x(i)$ was made, and the subsequent month, respectively. The respective matrix elements are the weights of the monthly correction factors $c_j$ and $c_{j+1}$:

$$v_{i,j} = \frac{d(x)}{d(\text{month})} \tag{26}$$

$$v_{i,j+1} = \frac{d(\text{month}) - d(x) + 1}{d(\text{month})} \tag{27}$$

The minimization of the cost functions of Eqs. (23) and (25) follows the same scheme as outlined for the cost function in Eq. (17).

The number of fit variables may be unreasonably large compared to the sample size, leading to poorly determined regression parameters. Since the reduction of seasonal correction parameters may be undesirable because it would imply quite crude discretization, the correction function $c$ can be smoothed by reduction of the month-to-month differences. This can be achieved by a constraint as proposed by Tikhonov (1963) or, in a context different from ours, by Phillips (1962) or Twomey (1963). This leads to a modified cost function

$$\chi^2 = \left(\mathbf{y}(\mathbf{x}) - (a + b\mathbf{x} + \mathbf{c}^T \mathbf{V})\right)^T \mathbf{S}_y^{-1} \cdot$$
$$\left(\mathbf{y}(\mathbf{x}) - (a + b\mathbf{x} + \mathbf{c}^T \mathbf{V})\right) + $$
$$\gamma \mathbf{c}^T \mathbf{L}^T \mathbf{L} \mathbf{c}$$
$$, \tag{28}$$

where $\gamma$ is a scalar to adjust the strength of the smoothing, and $\mathbf{L}$ is a $k \times k$ ($k$ is the number of seasonal correction parameters) cyclic first order differences matrix,

enhanced by two zero rows to cope for the parameters representing axis intercept and slope which shall not be subject to regularization[3].

It may be desirable to allow for a seasonal variation of the strength of the constraint. In this case we use a $k \times k$-dimensional diagonal matrix $\mathbf{D}$ to control the regularization strength and and use the following cost function:

$$\chi^2 = \left( \boldsymbol{y}(\boldsymbol{x}) - (a + b\boldsymbol{x} + \boldsymbol{c}^T \mathbf{V}) \right)^T \mathbf{S}_y^{-1} \cdot$$
$$\left( \boldsymbol{y}(\boldsymbol{x}) - (a + b\boldsymbol{x} + \boldsymbol{c}^T \mathbf{V}) \right) +$$
$$\boldsymbol{c}^T \mathbf{L}^T \mathbf{D} \mathbf{L} \boldsymbol{c}$$

(29)

Again, it is the use of the full covariance matrix which makes the method applicable to inhomogeneous data irregularly sampled in space and time. In any case the uncertainty of the slope in terms of variance is

$$\sigma_b^2 = \left( \frac{\partial b}{\partial \boldsymbol{y}} \right)^T \mathbf{S}_y \left( \frac{\partial b}{\partial \boldsymbol{y}} \right)$$

(30)

and the uncertainty of the intersect is

$$\sigma_a^2 = \left( \frac{\partial a}{\partial \boldsymbol{y}} \right)^T \mathbf{S}_y \left( \frac{\partial a}{\partial \boldsymbol{y}} \right).$$

(31)

The off-diagonal elements of $\mathbf{S}_y$ will determine whether the data errors map either predominantly onto the slope or onto the intersect of the regression curve. For example, large positive correlations throughout the data lead to large intersect errors while the slope remains quite well determined with sometimes surprisingly small uncertainties.

---

[3]Usually first order difference matrices as used in, e.g., vertical profile retrieval from remotely sensed data (c.f., e.g. Steck and von Clarmann 2001) are of the size $(k-1) \times k$. The difference with respect to that application is that we have a cyclic application here, i.e. the 12th and the 1st month are also constrained to each other.

**Trend estimation from clustered data**

T. von Clarmann et al.

Full Screen / Esc

Printer-friendly Version

Interactive Discussion

The extreme case would be fully correlated data errors. It is well known that such a constant bias in the data does not affect the trend at all. Consideration of full covariance matrices allows the correct treatment of realistic cases, where the errors are neither purely random nor purely systematic but may include correlations within subsets of the data.

## 5  Testing

If correct uncertainties and error covariances have been assumed and the errors are normally distributed, we expect $\chi^2$ to equal the number of degrees of freedom. If $\chi^2$ is larger (or smaller) than the respective upper (or lower) percentile of the $\chi^2$ function, the probability that the error characteristics disagree with the actual data and their uncertainties is larger than this percentile (1 minus percentile). The number of degrees of freedom is $n - i$, where n is the number of data points and $i$ is the number of fitted parameters. If a cyclic first order differences smoothness constraint is applied as proposed in Eq. (28), the number of degrees of freedom is the rank of the regularization matrix $\mathbf{L}^T \mathbf{L}$ which is $n - i_u$, where $i_u$ is the number of unregularized regression parameters. The number of degrees of freedom $n - i_u$ is a consequence of the cyclic application of the constraint. For non-cyclic applications of the first order differences smoothness constraint as applied in profile retrieval, the applicable number of degrees of freedom would be $n - i_u - 1$. $\chi^2$-statistics will lead to meaningful results only if the regularization term $\gamma \mathbf{L}^T \mathbf{DL}$ represents the true statistics of the differences of the state variables between adjacent sampling points but not for ad hoc choices.

## 6  Application areas

Application areas for trend estimation under consideration of covariances are (a) measurements at multiple sites, when at different times different sites dominate the sample

[ACPD](#)

9, 27675–27692, 2009

**Trend estimation from clustered data**

T. von Clarmann et al.

Title Page

Abstract | Introduction

Conclusions | References

Tables | Figures

◄ | ►

◄ | ►

Back | Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion

(e.g. Engel et al. 2009); (b) measurements with multiple measurement systems with specific errors which are correlated for all measurements done with the same system but independent between the measurement systems. Age of air measurements based on different calibration standards as used by Engel et al. (2009) fall into this category, where Eq. (9) should be used for trend estimation; (c) combination of data from two measurement systems which cover different episodes. The continuation of HALOE water wapour time series (Randel et al., 2004; Rosenlof and Reid, 2008) with MIPAS data (Milz et al., 2005) or any other set of satellite measurements of one trace gas by different sensors would be a typical application. These time series contain significant seasonality, hence trend estimation by minimization of cost functions Eqs. (16)–(17), (23), (25), (28), or (29) is recommended. In cases without overlap in time where the data sets cannot be calibrated one to the other, this approach of treatment of systematic errors is particularly useful. This applies to, e.g., MIPAS $H_2O$ measurements before (Milz et al., 2009) and after (von Clarmann et al., 2009) 2004, when the instrument was operated at different spectral resolutions.
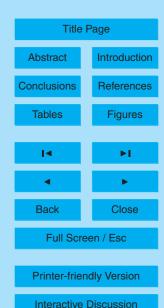
## 7 Conclusions

In case of irregular temporal and spatial sampling and/or multiple measurement systems, intersite and/or intersystem error correlations have to be considered for trend estimatation. To disregard the correlations not only renders the significance analysis meaningless but leads to wrong estimates of the trend itself. Intersite correlations can be estimated from satellite data. The regression model can easily be adapted for periodic corrections of known period length but unknown phase, shape and amplitude. This scheme solves the problem that usual approaches to infer periodic corrections rely on the time series being stationary, which is inherently not true in the case of trend estimation.

# References

Aitken, A. C.: On Least Squares and Linear Combinations of Observations, P. Roy. Soc. Edinb., 55, 42–48, 1935. 27676

Cantrell, C. A.: Technical Note: Review of methods for linear least-squares fitting of data and application to atmospheric chemistry problems, Atmos. Chem. Phys., 8, 5477–5487, 2008, http://www.atmos-chem-phys.net/8/5477/2008/. 27676

Cox, M., Harris, P., Hilton, M., and Woods, P.: Method for evaluating trends in ozone concentration data and its application to data from the UK Rural Ozone Monitoring Network, NPL Report CMSC 15/02, 2002. 27676

Efron, B.: The 1977 Rietz lecture: Bootstrap methods: Another look at the jackknife, Ann. Stat., 7, 1–26, 1979. 27676

Engel, A., Möbius, T., Bönisch, H., Schmidt, U., Heinz, R., Levin, I., Atlas, E., Aoki, S., Nakazawa, T., Sugawara, S., Moore, F., Hurst, D., Elkins, J., Schauffler, S., Andrews, A., and Boering, K.: Age of stratospheric air unchanged within uncertainties over the past 30 years, Nat. Geosci., 2, 28–31, 2009. 27689

Gardiner, T., Forbes, A., de Mazière, M., Vigouroux, C., Mahieu, E., Demoulin, P., Velazco, V., Notholt, J., Blumenstock, T., Hase, F., Kramer, I., Sussmann, R., Stremme, W., Mellqvist, J., Strandberg, A., Ellingsen, K., and Gauss, M.: Trend analysis of greenhouse gases over Europe measured by a network of ground-based remote FTIR instruments, Atmos. Chem. Phys., 8, 6719–6727, 2008, http://www.atmos-chem-phys.net/8/6719/2008/. 27676

IPCC: Climate Change 2007: The Physical Science Basis. Contribution of Working Group I to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change, edited by: Solomon, S., Qin, D., Manning, M., Chen, Z., Marquis, M., Averyt, K. B., Tignor M., and Miller, H. L., Cambridge University Press, Cambridge, UK and New York, NY, USA, 996 pp., 2007. 27676

Jones, P. D., Osborn, T. J., and Briffa, K. R.: Estimating sampling errors in large-scale temperature averages, J. Climate, 10, 2548–2568, 1997. 27677

Milz, M., von Clarmann, T., Fischer, H., Glatthor, N., Grabowski, U., Höpfner, M., Kellmann, S., Kiefer, M., Linden, A., Mengistu Tsidu, G., Steck, T., Stiller, G. P., Funke, B., López-Puertas, M., and Koukouli, M. E.: Water Vapor Distributions Measured with the Michelson Interferometer for Passive Atmospheric Sounding on board Envisat (MIPAS/Envisat), J. Geophys. Res.,

110, D24307, doi:10.1029/2005JD005973, 2005. 27689

Milz, M., v. Clarmann, T., Bernath, P., Boone, C., Buehler, S. A., Chauhan, S., Deuber, B., Feist, D. G., Funke, B., Glatthor, N., Grabowski, U., Griesfeller, A., Haefele, A., Höpfner, M., Kämpfer, N., Kellmann, S., Linden, A., Müller, S., Nakajima, H., Oelhaf, H., Remsberg, E., Rohs, S., Russell III, J. M., Schiller, C., Stiller, G. P., Sugita, T., Tanaka, T., Vömel, H., Walker, K., Wetzel, G., Yokota, T., Yushkov, V., and Zhang, G.: Validation of water vapour profiles (Version 13) retrieved by the IMK/IAA scientific retrieval processor based on full resolution spectra measured by MIPAS on board Envisat, Atmos. Meas. Techn., 2, 379–399, 2009. 27689

Muhlbauer, A., Spichtinger, P., and Lohmann, U.: Application and comparison of robust linear regression methods for trend estimation, J. Appl. Meteorol. Climatol., 48, 1961–1970, doi:10.1175/2009JAMC1851.1, 2009. 27676

Phillips, D.: A Technique for the numerical solution of certain integral equations of first kind, J. Ass. Comput. Mat., 9, 84–97, 1962. 27686

Randel, W. J., Wu, F., Oltmans, S. J., Rosenlof, K., and Nedoluha, G. E.: Interannual changes of stratospheric water vapor and correlations with tropical tropopause temperatures, J. Atmos. Sci., 61, 2133–2148, 2004. 27689

Rosenlof, K. H. and Reid, G. C.: Trends in the temperature and water vapor content of the tropical lower stratosphere: Sea surface connection, J. Geophys. Res., 113, D06107, doi:10.1029/2007/JD009109, 2008. 27689

Steck, T. and von Clarmann, T.: Constrained profile retrieval applied to the observation mode of the Michelson Interferometer for Passive Atmospheric Sounding, Appl. Opt., 40, 3559–3571, 2001. 27687

Tikhonov, A.: On the solution of incorrectly stated problems and method of regularization, Dokl. Akad. Nauk. SSSR, 151, 501–504, 1963. 27686

Twomey, S.: On the Numerical Solution of Fredholm Integral Equations of the First Kind by the Inversion of the Linear System Produced by Quadrature, J. ACM, 10, 97–101, 1963. 27686

Vigouroux, C., De Mazière, M., Demoulin, P., Servais, C., Hase, F., Blumenstock, T., Kramer, I., Schneider, M., Mellqvist, J., Strandberg, A., Velazco, V., Notholt, J., Sussmann, R., Stremme, W., Rockmann, A., Gardiner, T., Coleman, M., and Woods, P.: Evaluation of tropospheric and stratospheric ozone trends over Western Europe from ground-based FTIR network observations, Atmos. Chem. Phys., 8, 6865–6886, 2008, http://www.atmos-chem-phys.net/8/6865/2008/. 27676

**Trend estimation from clustered data**

T. von Clarmann et al.

Title Page

| Abstract | Introduction |
| Conclusions | References |
| Tables | Figures |

| ◄◄ | ►► |
| ◄ | ► |
| Back | Close |

Full Screen / Esc

Printer-friendly Version

Interactive Discussion

von Clarmann, T.: Validation of remotely sensed profiles of atmospheric state variables: strategies and terminology, Atmos. Chem. Phys., 6, 4311–4320, 2006, http://www.atmos-chem-phys.net/6/4311/2006/. 27679

von Clarmann, T., Höpfner, M., Kellmann, S., Linden, A., Chauhan, S., Funke, B., Grabowski, U., Glatthor, N., Kiefer, M., Schieferdecker, T., Stiller, G. P., and Versick, S.: Retrieval of temperature, $H_2O$, $O_3$, $HNO_3$, $CH_4$, $N_2O$, $ClONO_2$ and $ClO$ from MIPAS reduced resolution nominal mode limb emission measurements, Atmos. Meas. Techn., 2, 159–175, 2009. 27689

Weatherhead, E. C., Reinsel, G. C., Tiao, G. C., Meng, X.-L., Choi, D., Cheang, W.-K., Keller, T., DeLuisi, J., Wuebbles, D. J., Kerr, J. B., Miller, A. J., Oltmans, S. J., and Frederick, F. E.: Factors affecting the detection of trends: Statistical considerations and applications to environmental data, J. Geophys. Res., 103, 17149–17161, 1998. 27677

◄ | ►◄

◄ | ►

Back | Close

Full Screen / Esc