

Validation of remotely sensed profiles of atmospheric state variables: strategies and terminology

T. von Clarmann

Forschungszentrum Karlsruhe, Institut für Meteorologie und Klimaforschung, Karlsruhe, Germany

Received: 18 April 2006 – Published in Atmos. Chem. Phys. Discuss.: 20 June 2006

Revised: 16 August 2006 – Accepted: 15 September 2006 – Published: 25 September 2006

Abstract. This paper summarizes and classifies the various approaches to validation of remote measurements of atmospheric state variables, and tries to recommend a clear and unambiguous terminology. The following approaches have been identified: Intercomparison of individual profiles for accuracy validation; statistical comparison of matched pairs of measurements with respect to bias determination and precision validation; statistical intercomparison of randomly sampled measurements by two instruments, and comparison of a single measurement to an ensemble of measurements. Applicable statistics are shortly reviewed, and recipes for evaluation of the co-incidence error due to less than perfect coincidences are presented. An approach is suggested to quantitatively validate profile measurements when full covariance matrices are unavailable.

1 Introduction

Validation of a data product we understand means to confirm the predicted accuracy of the data product. After a series of self-consistence texts (e.g. Rodgers, 2000), the key element of validation is a statistical analysis of the differences between measurements of a new instrument to be validated, and of a reference instrument already validated. The purpose is to detect and remove any potential bias of the new measurement, and to verify that the estimated precision of the new measurement characterizes the measurements correctly.

Without any validated and bias-corrected reference measurement available, it may also be helpful to intercompare measurements by two or more non-validated instruments. This approach we call “cross validation”. While this approach certainly is no validation in its rigorous sense, it still may help to better characterize the data products.

Correspondence to: T. von Clarmann
(thomas.clarmann@imk.fzk.de)

2 Terminology

Let $\mathbf{x}=(x_1, \dots, x_N)^T$ be a vertical profile of an atmospheric state variable, sampled on a discrete vertical grid of N altitude gridpoints, describing the true atmospheric state at the altitude resolution of the measurement to be validated. Let further $\hat{\mathbf{x}}=(\hat{x}_1, \dots, \hat{x}_N)^T$ be a measurement of \mathbf{x} . The accuracy¹ \mathbf{a} of the measurement $\hat{\mathbf{x}}$ we understand is the square root of the expectation value of the squared differences of the true quantities x_n and their measurements \hat{x}_n :

$$\mathbf{a} = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_N \end{pmatrix} = \begin{pmatrix} \sqrt{\langle (\hat{x}_1 - x_1)^2 \rangle} \\ \sqrt{\langle (\hat{x}_2 - x_2)^2 \rangle} \\ \vdots \\ \sqrt{\langle (\hat{x}_N - x_N)^2 \rangle} \end{pmatrix} \quad (1)$$

The bias \mathbf{b} of a measurement is the expectation value of the deviation of the measured and the true quantity:

$$\mathbf{b} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_N \end{pmatrix} = \langle \hat{\mathbf{x}} - \mathbf{x} \rangle. \quad (2)$$

Depending on the nature of the bias, it can also be multiplicative rather than additive and then is better reported as a relative quantity:

$$\mathbf{b}_{\text{mult.}} = \begin{pmatrix} b_{\text{mult.};1} \\ b_{\text{mult.};2} \\ \vdots \\ b_{\text{mult.};N} \end{pmatrix} = \left\langle \begin{pmatrix} \frac{\hat{x}_1}{x_1} - 1 \\ \frac{\hat{x}_2}{x_2} - 1 \\ \vdots \\ \frac{\hat{x}_N}{x_N} - 1 \end{pmatrix} \right\rangle \quad (3)$$

¹In the statistics as well as remote sensing literature there are at least two different definitions of “accuracy”. The definition in this paper is consistent with the one used by, e.g. Rodgers (2000), Haseloff and Hoffmann (1970) or Walther and Moore (2005), while it is in conflict with Bevington (1969), whose “accuracy” corresponds to the quantity which is called “bias” here.

The precision p of the measurement characterizes the reproducibility of the measurement,

$$p = \begin{pmatrix} p_1 \\ p_2 \\ \vdots \\ p_N \end{pmatrix} = \begin{pmatrix} \sqrt{\langle(\hat{x}_1 - \langle\hat{x}_1\rangle)^2\rangle - \langle(x_1 - \langle x_1\rangle)^2\rangle} \\ \sqrt{\langle(\hat{x}_2 - \langle\hat{x}_2\rangle)^2\rangle - \langle(x_2 - \langle x_2\rangle)^2\rangle} \\ \vdots \\ \sqrt{\langle(\hat{x}_N - \langle\hat{x}_N\rangle)^2\rangle - \langle(x_N - \langle x_N\rangle)^2\rangle} \end{pmatrix}, \quad (4)$$

where the second term under the square root corrects for the natural variability of x . Accuracy, bias, and precision are related by

$$\begin{pmatrix} a_1^2 \\ a_2^2 \\ \vdots \\ a_N^2 \end{pmatrix} = \begin{pmatrix} b_1^2 \\ b_2^2 \\ \vdots \\ b_N^2 \end{pmatrix} + \begin{pmatrix} p_1^2 \\ p_2^2 \\ \vdots \\ p_N^2 \end{pmatrix} \quad (5)$$

Usually, remotely sensed data are provided along with careful data characterization, which includes estimates of the random error covariance matrix $\mathbf{S}_{\text{random}}$ and the systematic error covariance matrix \mathbf{S}_{sys} . In the case of remote measurements, these error estimates typically are the linear mapping of known uncertainties (measurement noise, model parameter uncertainties etc.) onto the retrieved quantities \hat{x}_n (Rodgers, 1990). This means, that these variances and covariances are ex ante estimates, i.e., they do not rely on any statistical comparison with external data. If the measurement includes a priori information, also the smoothing error, characterized by the covariance matrix $\mathbf{S}_{\text{smooth}}$, representing the mapping of the difference between the a priori assumption and the true state onto the measurement, contributes to the total error budget (Rodgers, 2000), and we get for the total error covariance matrix $\mathbf{S}_{\text{total}}$:

$$\mathbf{S}_{\text{total}} = \mathbf{S}_{\text{s+r}} + \mathbf{S}_{\text{smooth}} = \mathbf{S}_{\text{sys}} + \mathbf{S}_{\text{random}} + \mathbf{S}_{\text{smooth}}, \quad (6)$$

where $\mathbf{S}_{\text{s+r}}$ is the sum of $\mathbf{S}_{\text{random}}$ and \mathbf{S}_{sys} . The diagonal elements of these matrices are the related variances $\sigma_{\text{total};n}^2 = \sigma_{\text{total};n,n}^2$, $\sigma_{\text{sys};n}^2 = \sigma_{\text{sys};n,n}^2$, $\sigma_{\text{random};n}^2 = \sigma_{\text{random};n,n}^2$, and $\sigma_{\text{smooth};n}^2 = \sigma_{\text{smooth};n,n}^2$, respectively. The smoothing error can be composed of components constant with time ($\sigma_{\text{smooth,sys};n}^2$) and components randomly varying with time ($\sigma_{\text{smooth,random};n}^2$). Validation then means to confirm the ex ante error estimates by verification that for all n from 1 to N

$$a_n^2 = \sigma_{\text{s+r};n}^2 + \sigma_{\text{smooth};n}^2 = \sigma_{\text{total};n}^2 \quad (7)$$

$$b_n^2 = \sigma_{\text{sys};n}^2 + \sigma_{\text{smooth,sys};n}^2 \quad (8)$$

$$p_n^2 = \sigma_{\text{random};n}^2 + \sigma_{\text{smooth,random};n}^2 \quad (9)$$

This is not as easy as it might seem, because the true atmospheric state x , which is needed to evaluate precision, bias, and accuracy with the formalism outlined above is not available. Instead, we use independent reference measurements, which allow to infer ex post estimates of bias, precision, and

accuracy. A useful strategy in validation is to first search for a possible bias, to quantify the bias in order to allow its correction, and to finally validate the estimated precision. Optimally, also the causes of the bias will be understood and removed. The scope of this paper, however, is restricted to the detection and quantification of the bias, and the validation of the precision estimates.

3 Comparison of co-incident measurements

3.1 General aspects

Let \hat{x}_{val} and \hat{x}_{ref} be two vertical profiles of the same quantity, measured by instruments “val” (instrument to be validated) and “ref” (independent reference instrument), respectively. The profiles and related diagnostic data have to be represented on a common grid, which usually implies regridding of one or both profiles (Calisesi et al., 2005). Further, if the measurements include a priori information, both profiles have to be transformed to the same a priori profile, and the smoothing error of the difference, $\mathbf{S}_{\text{smooth,diff}}$, has to be estimated (Rodgers and Connor, 2003). This smoothing error difference can be minimized or nullified using a method proposed by Ridolfi et al. (2006a), or it can be restricted to sub-scale differences by transformation of the data to a dedicated representation (von Clarmann and Grabowski, 2006).

Rodgers and Connor (2003) suggest to quantify profile intercomparison by application of a χ^2 test.

$$\chi^2 = (\hat{x}_{\text{val}} - \hat{x}_{\text{ref}})^T \mathbf{S}_{\text{diff}}^{-1} (\hat{x}_{\text{val}} - \hat{x}_{\text{ref}}), \quad (10)$$

where \mathbf{S}_{diff} is the ex ante estimate of the error covariance matrix of the difference $\hat{x}_{\text{val}} - \hat{x}_{\text{ref}}$. The actual value of χ^2 allows to conclude if the differences $\hat{x}_{\text{val}} - \hat{x}_{\text{ref}}$ are consistent with the ex ante estimates of the uncertainties of the difference, represented by its covariance matrix \mathbf{S}_{diff} , or if there is a significant inconsistency. The integral of the χ^2 probability density function from the actual χ^2 to infinity yields the probability $P_{\text{acc}}(\chi^2)$ that the actual χ^2 value occurs accidentally, i.e. that the differences can be explained by the error estimates; the integral from zero to the actual χ^2 yields the probability of a substantial difference. In particular, underestimation of the variances of the differences would lead to a χ^2 value larger than its expectation value, which, in the case of a regular \mathbf{S}_{diff} , is the number of comparison pairs. If the probability of a substantial difference is above a certain threshold, e.g. 95%, then the difference is significant at 5% confidence level and there is statistical evidence that there is something wrong with the data or related error estimates. If this probability is below 95%, then the falsification has failed because the difference is not significant at the chosen (5%) confidence level.

Usually, more than one pair of co-incident profiles is available, and Eq. (10) can be applied to a larger ensemble of comparison pairs. For a large ensemble of K independent

intercomparison pairs, the χ^2 test can either be performed for L subsets of data or as one single χ^2 test involving the χ^2 probability density function of K degrees of freedom. The division into subsets allows to check if the χ^2 values found follow the expected χ^2 distribution. E.g., Migliorini et al. (2004) have detected suspicious ozone profiles in their comparison ensemble by comparison of the expected and the found χ^2 distribution. The probability of disagreement for the complete ensemble, $P_{\text{dis}}(L, \chi_{\text{max}}^2) = 1 - P_{\text{acc}}(L, \chi_{\text{max}}^2)$, can be estimated according to the multiplication axiom,

$$P_{\text{dis}}(L, \chi_{\text{max}}^2) \leq P_{\text{dis}}(1, \chi_{\text{max}}^2)^L, \quad (11)$$

where χ_{max}^2 is the largest χ^2 value found in the ensemble of comparison subsets, and $P_{\text{dis}}(1, \chi_{\text{max}}^2)$ is the related probability of substantial disagreement in the subset where χ_{max}^2 is found. The drawback of this approach is that it is not sufficient in a sense that the probability estimate is based on the maximum χ^2 value only and thus does not use all available information. This implies that this test is not very robust because it is very sensitive to outliers. For determination of $P_{\text{acc}}(\chi^2)$ the single χ^2 test involving the complete ensemble of comparison pairs is superior because of its inherent sufficiency. The safest is to combine both approaches. Discrepancies can then point at non-representative outliers in the comparison ensemble.

Large probabilities of substantial differences can have three different causes: 1. The ex ante error estimates may have been too optimistic. 2. The initial ensemble size was chosen too small and not representative. In this case, a larger comparison ensemble may help to achieve a larger $P_{\text{acc}}(\chi^2)$. If, however, the initial sample was representative, even larger χ^2 values will most probably occur in a larger ensemble, and $P_{\text{acc}}(\chi^2)$ will not improve. It is, of course, important to work with pre-defined random samples and not to adjust the sample or the sample size to the maximum $P_{\text{acc}}(\chi^2)$. 3. Large χ^2 can also be associated with a particular subset of the sample which can be characterized by some objective criterion. Migliorini et al. (2004), e.g. have found problems in O_3 data from spectra suspected to be cloud contaminated. In such a case it may be appropriate to define a kind of data filter and to validate only the subset of the data which passes the filter. There are, however, two traps in this approach: First, the filter should not use the quantity to be validated itself as a filter criterion. Second, the new analysis system, of which the newly defined filter is a part, has to be validated using an independent comparison ensemble. When the original sample is used, it will always be possible to tune the data filter such that good agreement between the intercompared data is achieved.

While, in the case of good ex ante error estimates, we certainly can do better than just reject the hypothesis of inconsistency at 5% confidence level, which still allows a probability

$P_{\text{dis}}(L, \chi_{\text{max}}^2)$ of up to 95%, we cannot prove equivalence of $\hat{\mathbf{x}}_{\text{val}}$ and $\hat{\mathbf{x}}_{\text{ref}}$ at a confidence level better than 50%, because

$$\lim_{L \rightarrow \infty} P_{\text{acc}}(L, \langle \chi_{\text{max}}^2 \rangle) = 0.5. \quad (12)$$

Knowledge of all error terms contributing to \mathbf{S}_{diff} is crucial to allow a meaningful estimate of the significance of the differences $\hat{\mathbf{x}}_{\text{val}} - \hat{\mathbf{x}}_{\text{ref}}$. The ex ante estimate of the covariance matrix of the difference with elements $s_{\text{diff};m,n}$ is usually calculated as

$$\mathbf{S}_{\text{diff}} = \mathbf{S}_{\text{s+r,val}} + \mathbf{S}_{\text{s+r,ref}} + \mathbf{S}_{\text{coinc.}} + \mathbf{S}_{\text{smooth,diff}}, \quad (13)$$

where $\mathbf{S}_{\text{coinc.}}$ describes the spatial and temporal co-incidence error in terms of variances and co-variances, which are important to be quantified and considered in the case of less than perfect co-incidences of the two measurements (see Sect. 3.2). If both the validation and the reference measurement have a common error source, this introduces correlations. This applies e.g. when the same or correlated temperature profiles or spectroscopic data are used to derive both $\hat{\mathbf{x}}_{\text{val}}$ and $\hat{\mathbf{x}}_{\text{ref}}$. In the case of such correlations, \mathbf{S}_{diff} can be evaluated as

$$\begin{aligned} \mathbf{S}_{\text{diff}} = & (\mathbf{I}, -\mathbf{I}) \begin{pmatrix} \mathbf{S}_{\text{s+r,val}}, \mathbf{C}_{\text{s+r,val,ref}} \\ \mathbf{C}_{\text{s+r,val,ref}}^T, \mathbf{S}_{\text{s+r,ref}} \end{pmatrix} (\mathbf{I}, -\mathbf{I})^T \\ & + \mathbf{S}_{\text{coinc.}} + \mathbf{S}_{\text{smooth,diff}}, \end{aligned} \quad (14)$$

where \mathbf{I} is $N \times N$ unity and where matrix $\mathbf{C}_{\text{s+r}}$ contains the related covariance elements $r_{\text{s+r,val,ref};m,n} \sigma_{\text{s+r,val},m} \sigma_{\text{s+r,ref},n}$ between the new measurement “val” and the reference measurements “ref”, where $r_{\text{s+r,val,ref};m,n}$ is the correlation coefficient of the combined systematic and random errors of the validation measurement at altitude m and the reference measurement at altitude n . Comparison of two individual profiles does not allow to distinguish between precision and bias validation.

3.2 Determination of co-incidence error in time and space

Usually, only profiles are selected for comparison which meet a certain co-incidence criterion in time and space or any other adequate co-ordinates d like solar zenith angle, potential vorticity, equivalent latitude etc. The actual difference Δd in this quantity is the mismatch, and the maximum allowed mismatch is the co-incidence criterion Δd_{max} .

Variability of most atmospheric state variables is composed of a functional term and a random term. The abundance of a certain species, for example, may have a typical latitudinal dependence or a typical diurnal variation, which are superimposed by random fluctuations caused by the actual small-scale atmospheric situation. Whenever applicable, the functional term should be corrected first by a correction function \mathbf{M} , which can either be some appropriate parametrization, or alternatively a tabulated data set. With d_{val} and d_{ref} being the co-ordinates of the validation and

reference measurements, respectively, the uncorrected reference measurement $\hat{x}_{\text{ref,uncorrected}}$ is corrected as

$$\hat{x}_{\text{ref}} = \hat{x}_{\text{ref,uncorrected}} + \mathbf{M}(d_{\text{val}}) - \mathbf{M}(d_{\text{ref}}) \quad (15)$$

and only the residual random part of the co-incidence error with respect to the corrected reference measurement \hat{x}_{ref} should be characterized by the covariance matrix $\mathbf{S}_{\text{coinc.}}$. Otherwise the co-incidence error may not follow a Gaussian distribution, and errors due to systematic sampling differences in d may inadvertently be treated as random co-incidence errors. This may happen, e.g., if the abundance of an atmospheric constituent which is characterized by a strong diurnal change is observed by two instruments at instrument-specific local times. An example of application of a correction function \mathbf{M} is found in Ridolfi et al. (2006b)² who use ECMWF (European Centre for Medium-Range Weather Forecasts) temperature analyses to estimate the component of the differences between MIPAS and radiosonde temperatures which are explained by mismatch in space and time. A similar approach was chosen by Cortesi et al. (2006)³ for ozone.

To quantify the residual co-incidence error caused by finer structures in d than those accounted for by the correction function \mathbf{M} , a sufficiently fine resolved typical reference data set \hat{x}_r of state variable $x(d)$ is needed. Let the reference data set contain $K(\Delta d)$ independent pairs of data points separated by the mismatch $\Delta d = d_{\text{val}} - d_{\text{ref}}$. Then, the co-incidence error $\mathbf{S}_{\text{coinc.}}$ can be evaluated as a function of Δd as

$$s_{\text{coinc.};m,n}(\Delta d) = \frac{\sum_{k=1}^{K(\Delta d)} \left(\Delta \hat{x}_{r;m}(\Delta d) \right)_k \left(\Delta \hat{x}_{r;n}(\Delta d) \right)_k}{K} - s_{\text{err,diff};m,n} \quad (16)$$

where

$$\left(\Delta \hat{x}_{r;m}(\Delta d) \right)_k = \left(\hat{x}_{r;m}(d) - \hat{x}_{r;m}(d + \Delta d) - M_m(d) + M_m(d + \Delta d) \right)_k \quad (17)$$

and

$$\left(\Delta \hat{x}_{r;n}(\Delta d) \right)_k = \left(\hat{x}_{r;n}(d) - \hat{x}_{r;n}(d + \Delta d) - M_n(d) + M_n(d + \Delta d) \right)_k \quad (18)$$

m and n identify the profile gridpoints, and

$$s_{\text{err,diff};m,n} = 2s_{\text{err};m,n}, \quad (19)$$

²Ridolfi, M., Blum, U., Carli, B., et al.: Geophysical Validation of temperature retrieved from MIPAS/ENVISAT atmospheric Limb-emission measurements, Atmos. Chem. Phys. Discuss., in preparation, 2006b.

³Cortesi, U., Blom, C., Blumenstock, Th., et al.: Co-ordinated validation activity and quality assessment of MIPAS-ENVISAT Ozone data, Atmos. Chem. Phys. Discuss., in preparation, 2006.

where $s_{\text{err};m,n}$ is an element of the random error profile covariance matrix \mathbf{S}_{err} of the reference data set of the state variable \hat{x}_r . The factor of two accounts for error propagation through the calculation of the difference. \mathbf{S}_{err} , which is assumed constant and uncorrelated with geolocation, cannot be obtained from the scatter of the reference sample, because the latter contains the natural variability we are trying to isolate. Instead, a true ex ante estimate is needed, e.g. by error propagation calculation, sensitivity studies or similar means. The M terms account for the difference already explained by the functional mismatch correction.

In order to get $K(\Delta d)$ large enough for meaningful statistics, binning of $\mathbf{S}_{\text{coinc.}}$ is recommended, i.e. evaluation of $\mathbf{S}_{\text{coinc.}}([\Delta d_1, \Delta d_2])$ for all mismatches in a range from Δd_1 to Δd_2 , where $\mathbf{S}_{\text{coinc.}}$ is sufficiently linear in Δd . If such a bin $[\Delta d_1, \Delta d_2]$ covers the entire co-incidence criterion, i.e. $\Delta d_1 = 0$ and Δd_2 equals the co-incidence criterion, it is no longer necessary to care about the Δd -dependence of $\mathbf{S}_{\text{coinc.}}$, but the mean co-incidence error $\overline{\mathbf{S}_{\text{coinc.}}} \approx \mathbf{S}_{\text{coinc.}}(\overline{\Delta d})$ can be used for the entire ensemble of co-incidences.

Meteorological analyses, satellite measurements or modeled atmospheric fields can be used as reference data sets to evaluate the co-incidence error on a larger scale. It is important to carefully assess any possible reduction of the horizontal variability in these datasets through application of background or a priori knowledge in the sense of variational data assimilation (e.g. Ide et al., 1997) or optimal estimation retrievals (Rodgers, 1976), respectively. For determination of small-scale temporal fluctuations stationary in situ measurements or ground-based remote sensing measurements are better suited, while for small-scale spatial fluctuations aircraft measurements are the first choice.

Multi-dimensional co-incidence can be assessed component-wise by evaluation of Eq. (16) for each co-incidence direction (e.g. latitude, longitude and time) and summing up the respective co-incidence error covariance matrices. In the case where the variation of the state variable under assessment is correlated between two of these dimensions, the summation has to be replaced by the following scheme:

$$\mathbf{S}_{\text{coinc.}} = (\mathbf{I}, \mathbf{I}) \begin{pmatrix} \mathbf{S}_{\text{coinc.};1}, \mathbf{C}_{\text{coinc.};1,2} \\ \mathbf{C}_{\text{coinc.};1,2}^T, \mathbf{S}_{\text{coinc.};2} \end{pmatrix} (\mathbf{I}, \mathbf{I})^T, \quad (20)$$

where the subscripts of the covariance matrices $\mathbf{S}_{\text{coinc.};l}$ and the cross-dimension covariances $\mathbf{C}_{\text{coinc.};k,l}$ denote the dimensions along which the variabilities are analyzed. Such correlations may apply, e.g., to the mixing ratio of an inert trace gas the abundance of which is ruled by transport processes. The existence of a prevailing direction of wind in combination with a prevailing gradient in the field of the state variable then introduces such correlations.

Another option to handle co-incidence errors in L dimensions is to define a norm of the following type which transforms the multi-dimensional mismatch

$\Delta \mathbf{d}=(\Delta d_1, \dots, \Delta d_L)$ to a scalar mismatch distance Δd :

$$\Delta d = \sqrt{\sum_1^L (w_l \Delta d_l)^2}, \quad (21)$$

where w_l are weighting factors reflecting the expected variability of the state variable with the respective direction l . Steck et al. (2006)⁴, e.g., have used

$$\Delta d = \sqrt{\Delta_{\text{long}}^2 + \Delta_{\text{lat}}^2 + (\Delta_t v_w)^2} \quad (22)$$

where Δ_{long} and Δ_{lat} are longitudinal and latitudinal mismatch distances, Δ_t is the mismatch in time, v_w is the typical windspeed. This particular norm holds for analysis of transport-dominated abundances of trace species without prevailing gradients and wind directions.

3.3 Smoothing

Additional complication arises if the measurements to be compared characterize air parcels of non-zero extension in the direction of d . In this case the smoothing error in direction of d and the co-incidence error can no longer be treated as independent. Smoothing error and co-incidence error are errors of the same nature, since both characterize differences caused by the fact that two instruments observe different parts of the atmosphere. The pure co-incidence error describes the error component caused by the variation of the atmospheric state over a distance Δd between the two disjoint air parcels observed by the two instruments to be compared. Contrary to that, the smoothing error difference quantifies the error component caused by the different weights of different parts within one air parcel observed by two instruments. Both error terms can be estimated by one formalism which includes both the smoothing error application and the co-incidence error application as well as the case of partly overlapping air parcels.

Here we first discuss the quite general case that for both the reference measurement and the measurement to be validated smoothing in all three spatial directions has to be considered, and where the observed air parcels may or may not overlap. Inclusion of additional dimensions (e.g. time) is straightforward and will not explicitly be discussed here. Later, some convenient simplifications will be mentioned.

In a first step, we store all relevant elements of the 3-dimensional (3-D) fine-resolved field of the atmospheric state variable under assessment in a vector

⁴Steck, T., Blumenstock, T., Clarmann, T., Glatthor, N., Grabowski, U., Hase, F., Hochschild, G., Höpfner, M., Kellmann, S., Kiefer, M., Kopp, G., Linden, A., Milz, M., Oelhaf, H., Stiller, G. P., Wetzell, G., Zhang, G., Fischer, H., Funke, B., Wand, D. Y., Gathen, P., Hansen, G., Stebel, K., Kyrö, E., Allaart, M., Redondas Marrero, A., Remsberg, E., Russell III, J., Steinbrecht, W., Yela, M., and Raffalski, U.: Validation of ozone measurements from MIPAS-Envisat, in preparation, 2006.

$\mathbf{x}_{3-D}=(x_{3-D;1}, \dots, x_{3-D;L})^T$. Those elements are considered relevant, whose entries in the rows of either of the 3-D averaging kernel matrices of the size $L \times L$, $\mathbf{A}_{3-D,\text{ref}}$ and $\mathbf{A}_{3-D,\text{val}}$, are non-zero, i.e., elements of the 3-D field which are seen by at least one of both instruments. When reshaping the 3-D field to the vector \mathbf{x}_{3-D} , the ordering is arbitrary but unambiguous and we use the notation $l(i, j, n)$ for the l -th element of the vector which represents the element of the n -th altitude, and the i -th and j -th geolocation in the original 3-D field. Following the concept of Rodgers (2000) but disregarding noise and other measurement errors, the atmospheric state as seen by the instrument to be validated and the reference instrument are

$$\hat{\mathbf{x}}_{3-D,\text{val}} = \mathbf{A}_{3-D,\text{val}} \mathbf{x}_{3-D} + (\mathbf{I} - \mathbf{A}_{3-D,\text{val}}) \mathbf{x}_{3-D,\text{a}} \quad (23)$$

and

$$\hat{\mathbf{x}}_{3-D,\text{ref}} = \mathbf{A}_{3-D,\text{ref}} \mathbf{x}_{3-D} + (\mathbf{I} - \mathbf{A}_{3-D,\text{ref}}) \mathbf{x}_{3-D,\text{a}}, \quad (24)$$

respectively, where \mathbf{I} is $L \times L$ unity and where $\mathbf{x}_{3-D,\text{a}}$ is the common a priori 3-D information which may be included in the measurements. The difference $\Delta \hat{\mathbf{x}}_{3-D}$ of the measurements caused both by spatial mismatch and different smoothing characteristics then is

$$\Delta \hat{\mathbf{x}}_{3-D} = (\mathbf{A}_{3-D,\text{val}} - \mathbf{A}_{3-D,\text{ref}}) (\mathbf{x}_{3-D} - \mathbf{x}_{3-D,\text{a}}), \quad (25)$$

where the smoothing component is accounted for by the deviation of each of the relevant sub-matrices of \mathbf{A}_{3-D} from unity, while the co-incidence component caused by Δd is accounted for by the different placement of the relevant sub-matrices in the full \mathbf{A}_{3-D} matrix. The difference profile $\Delta \hat{\mathbf{x}}$ at the nominal geolocation (i, j) then has the elements

$$\Delta \hat{x}_n = \Delta \hat{x}_{3-D;l(i,j;n)} \quad (26)$$

This way to calculate the differences can be applied to the correction scheme suggested in Eq. (15) where the difference $\mathbf{M}(d_{\text{val}}) - \mathbf{M}(d_{\text{ref}})$ can be replaced by the $\Delta \hat{\mathbf{x}}$ values from Eq. (26) applied to the $\Delta \hat{\mathbf{x}}_{3-D}$ field from Eq. (25), where \mathbf{x}_{3-D} is generated by the model \mathbf{M} . Further, Eqs. (25–26) can be used to calculate the $\Delta \hat{\mathbf{x}}(\Delta d)$ terms in Eq. (16) which is used for a statistical estimate of the residual co-incidence error after a possible correction. In the latter application, the elements $s_{\text{err};m,n}$ of the covariance matrix \mathbf{S}_{err} describing the part of the observed differences of pairs of profiles caused by errors in the 3-D dataset itself rather than its variability in d at a given geolocation (i, j) are calculated as

$$s_{\text{err,diff};m,n} = \left((\mathbf{A}_{3-D,\text{val}} - \mathbf{A}_{3-D,\text{ref}}) \mathbf{S}_{3-D} (\mathbf{A}_{3-D,\text{val}} - \mathbf{A}_{3-D,\text{ref}})^T \right)_{l(i,j;m),l(i,j;n)}, \quad (27)$$

where \mathbf{S}_{3-D} is the $L \times L$ covariance matrix representing the uncertainties of $\hat{\mathbf{x}}_{3-D}$. Equation (16) can then be extended to describe the combined smoothing and co-incidence error:

$$(\mathbf{S}_{\text{coinc}} + \mathbf{S}_{\text{smooth}})_{m,n} =$$

$$\frac{1}{K} \left(\sum_{k=1}^K \left(\mathbf{A}_{3-D,\text{val}} - \mathbf{A}_{3-D,\text{ref}} \right) (\mathbf{x}_{3-D,\text{r}} - \mathbf{M}_{3-D}) \right)_{l(i(k),j(k),m)} \quad (28)$$

$$\times \left(\mathbf{A}_{3-D,\text{val}} - \mathbf{A}_{3-D,\text{ref}} \right) (\mathbf{x}_{3-D,\text{r}} - \mathbf{M}_{3-D})_{l(i(k),j(k),n)}$$

$$- (\mathbf{A}_{3-D,\text{val}} - \mathbf{A}_{3-D,\text{ref}}) \mathbf{S}_{3-D} (\mathbf{A}_{3-D,\text{val}} - \mathbf{A}_{3-D,\text{ref}})^T,$$

where $l(i(k), j(k), n)$ denotes the position of the point of the 3 dimensional field in the vector, and where $i(k)$, $j(k)$, and m or n denote the horizontal nominal co-ordinates and the altitude of the k -th sample pair used to evaluate the combined co-incidence and smoothing error covariance matrix. \mathbf{M}_{3-D} is the 3-dimensional correction model tabulated in a vector according to the predefined reshaping convention $(\mathbf{M}_{3-D})_{k(i,j,n)} = M_n(d(i, j))$.

If the \mathbf{S}_{err} matrices vary between the data pairs over which the summation in Eqs. (16) or (28) runs, this has to be taken into account by appropriate weighting of the terms.

If an ex ante estimate of the covariance matrix $\mathbf{S}_{3-D,a}$ which describes the variability and correlations of \mathbf{x}_{3-D} is available, the formalism proposed by Rodgers and Connor (2003) for estimation of the smoothing error difference can be directly applied to estimate the combined effect of smoothing error difference and co-incidence error, $\mathbf{S}_{\text{coinc.}} + \mathbf{S}_{\text{smooth.}}$, at a given location (i, j) , without evaluation of a reference data set $\hat{\mathbf{x}}_r$:

$$s_{\text{coinc.};m,n} + s_{\text{smooth.};m,n} = \quad (29)$$

$$\left((\mathbf{A}_{3-D,\text{val}} - \mathbf{A}_{3-D,\text{ref}})^T \mathbf{S}_{3-D,a} (\mathbf{A}_{3-D,\text{val}} - \mathbf{A}_{3-D,\text{ref}}) \right)_{l(i,j;m),l(i,j;n)}$$

For particular cases the relationship of smoothing error and co-incidence error can be simplified due to the nature of measurements to be compared. Further, the rigorous approach often is not possible since the needed data are unavailable. For example, $\hat{\mathbf{x}}_r$ may not be available as sufficiently high resolved 3-D field but only as 1- or 2-dimensional cross-section. Further, the \mathbf{A}_{3-D} averaging kernel matrix may not be provided along with the data product. Sometimes the rigorous approach is inappropriate, because the given data sets are not accurate enough to justify the related effort. In this case, the errors introduced by the approximations may no longer be significant.

Obvious simplifications are: The instruments may sound non-overlapping air parcels; this is reflected by disjoint non-zero entries in the 3-D averaging kernel matrices and, as a consequence, correlation terms between the air parcels in Eq. (27) turn zero. The reference measurement can be a point measurement, e.g. an in situ measurement, where no averaging kernel has to be considered; in this case, the reference data set $\hat{\mathbf{x}}_r$ can first be smoothed by application of the matrix of vertical averaging kernels and the assessment of the horizontal smoothing and co-incidence error then requires only the horizontal components of the 3-D averaging kernel matrix. Further, \mathbf{S}_{3-D} may have some diagonal structure, e.g. because the errors are assumed uncorrelated in geolocations. Often the cross-line-of-sight component of the 3-D averag-

ing kernel matrix is negligible. Various retrieval schemes use $\mathbf{x}_{3-D,a} = 0$.

The entries of \mathbf{A}_{3-D} related to vertical smoothing are contained in the profile averaging kernel matrix, which often is available as part of the diagnostics of the data sets under assessment. For the elements related to horizontal smoothing the situation is different: For nadir sounders and the cross-line-of-sight component of limb sounders, related averaging components just are identical to the field of view function, possibly widened by the horizontal smearing caused by the motion of the instrument while the measurement is made. If the field of view is narrower than the horizontal spacing of gridpoints in the data set $\hat{\mathbf{x}}_r$, this component can be ignored. The horizontal along-line-of-sight component can be obtained e.g. from perturbational analysis or analytically from 2-D radiative transfer modelling and retrieval tools (see, e.g. Steck et al., 2005; Carlotti et al., 2001). If the horizontal components of \mathbf{A}_{3-D} are not available, they can be approximated by $\mathbf{R}\mathbf{A}$, where \mathbf{A} is the vertical profile averaging kernel matrix, and \mathbf{R} the $I \times N$ dimensional ray-tracing operator, which maps altitudes z_1, \dots, z_n to along-track distances d_1, \dots, d_I according to the observation geometry. Elements of \mathbf{A} representing contributions from below the tangent altitude are assigned to the tangent point geolocation. This approximation, however, neglects both the mapping of any horizontal smoothing error onto the retrieved profile, and the asymmetry of the horizontal averaging kernel around the tangent point of a limb viewing measurement. This approach has been chosen by Ridolfi et al. (2006b)² and Cortesi et al. (2006)³ to account for the horizontal smoothing of MIPAS in the co-incidence correction. These authors have used vertically smoothed ECMWF fields as correction model in Eq. (15) and have considered the horizontal smoothing by the formalism of Eqs. (25–26), where $\mathbf{A}_{3-D,\text{ref}}$ was assumed unity and where $\mathbf{x}_{3-D,a}$ was zero.

4 Bias determination

To determine the bias between two measurement systems, a statistical ensemble of measurements is needed. This ensemble can either be composed of K matching pairs of measurements or random samples of K and L measurements of each measuring system, respectively.

4.1 Statistical bias determination with matching pairs of measurements

The mean difference between measurements to be validated and co-incident reference measurements can be compared with its statistical uncertainty in order to determine any bias between the measurement to be validated and the reference measurement and its significance. With K pairs of co-incident measurements available, the bias b_{diff} between these measurements is estimated as (here and henceforth we use

the \checkmark symbol to denote ex post estimates based on a statistical comparison with reference data)

$$\checkmark_{\text{diff}} = \frac{\sum_{k=1}^K (\hat{x}_{\text{val};k} - \hat{x}_{\text{ref};k})}{K} \quad (30)$$

The statistical uncertainty of the bias is characterized by the related covariance matrix \checkmark_{bias} , the elements of which are estimated as

$$\begin{aligned} \checkmark_{\text{bias};m,n} &= \frac{\sum_{k=1}^K (\hat{x}_{\text{val};m,k} - \bar{x}_{\text{val};m})(\hat{x}_{\text{val};n,k} - \bar{x}_{\text{val};n})}{K(K-1)} + \\ &\frac{\sum_{k=1}^K (\hat{x}_{\text{ref};m,k} - \bar{x}_{\text{ref};m})(\hat{x}_{\text{ref};n,k} - \bar{x}_{\text{ref};n})}{K(K-1)} - \\ &\frac{\sum_{k=1}^K (\hat{x}_{\text{val};m,k} - \bar{x}_{\text{val};m})(\hat{x}_{\text{ref};n,k} - \bar{x}_{\text{ref};n})}{K(K-1)} - \\ &\frac{\sum_{k=1}^K (\hat{x}_{\text{val};n,k} - \bar{x}_{\text{val};n})(\hat{x}_{\text{ref};m,k} - \bar{x}_{\text{ref};m})}{K(K-1)} \\ &= \frac{\sum_{k=1}^K (\hat{x}_{\text{val};m,k} - \hat{x}_{\text{ref};m,k} - \checkmark_{\text{diff};m})(\hat{x}_{\text{val};n,k} - \hat{x}_{\text{ref};n,k} - \checkmark_{\text{diff};n})}{K(K-1)}, \quad (31) \end{aligned}$$

where

$$\bar{x} = \frac{\sum_{k=1}^K \hat{x}_k}{K}. \quad (32)$$

This assessment does not need any error estimates of \hat{x}_{val} or \hat{x}_{ref} .

The consistence of \checkmark_n and the ex ante estimate of the systematic error $\sigma_{\text{sys},n}$ of the retrieved state parameter \hat{x}_n , or of $\checkmark_{\text{diff};n}$ and $\sigma_{\text{diff},\text{sys},n}$, respectively, can easily be checked (see, e.g. Ridolfi et al., 2006b², for application to MIPAS temperature validation, or Cortesi et al., 2006³, for ozone validation). Evaluation of the significance of the bias then requires χ^2 statistics, where

$$\langle \chi_{\text{bias}}^2 \rangle = \langle \checkmark_{\text{diff}}^T \checkmark_{\text{bias}}^{-1} \checkmark_{\text{diff}} \rangle = N. \quad (33)$$

The consistency of the bias with the ex ante estimates of the systematic error components can also be evaluated by a χ^2 test:

$$\langle \chi_{\text{bias}}^2 \rangle = \langle \checkmark_{\text{diff}}^T (\mathbf{S}_{\text{diff},\text{sys}} + \checkmark_{\text{bias}})^{-1} \checkmark_{\text{diff}} \rangle = N. \quad (34)$$

Covariance matrix $\mathbf{S}_{\text{diff},\text{sys}}$ includes the systematic components of \mathbf{S}_{diff} :

$$\mathbf{S}_{\text{diff},\text{sys}} = \mathbf{S}_{\text{val},\text{sys}} + \mathbf{S}_{\text{ref},\text{sys}} + \mathbf{S}_{\text{smooth},\text{sys}} \quad (35)$$

Obviously, neither the root mean squares difference of profiles obtained from two measurement systems nor $1/\sqrt{K}$ of the root mean squares difference are a measure of the significance of the bias.

If ex ante precision estimates of differences

$$\mathbf{S}_{\text{diff},\text{random}} = \mathbf{S}_{\text{val},\text{random}} + \mathbf{S}_{\text{ref},\text{random}} + \mathbf{S}_{\text{coinc.}} + \mathbf{S}_{\text{smooth},\text{random}} \quad (36)$$

are available and the uncertainties are known to vary within the sample, the measurements can be weighted accordingly to determine the weighted bias \checkmark_{diff} :

$$\checkmark_{\text{diff}} = \left(\sum_{k=1}^K \mathbf{S}_{\text{diff},\text{random};k}^{-1} \right)^{-1} \left(\sum_{k=1}^K \mathbf{S}_{\text{diff},\text{random};k}^{-1} (\hat{x}_{\text{val};k} - \hat{x}_{\text{ref};k}) \right) \quad (37)$$

The bias uncertainty in terms of covariance matrix then is

$$\mathbf{S}_{\text{bias}} = \left(\sum_{k=1}^K \mathbf{S}_{\text{diff},\text{random};k}^{-1} \right)^{-1}, \quad (38)$$

which is an ex ante estimate. Thus, the use of \checkmark_{diff} and \mathbf{S}_{bias} as determined from Eqs. (38) and (37) for precision validation (see Sect. 5) is not fully conclusive, because it depends on typically unvalidated ex ante precision estimates.

The most probable estimate of the multiplicative bias from a given sample of K measurement pairs $\hat{x}_{\text{val};n,k}$ and $\hat{x}_{\text{ref};n,k}$, each affected by an additive random error of constant expectation is

$$\checkmark_{\text{diff},\text{mult.};n} = \frac{\checkmark_{\text{diff};n}}{\frac{\sum_{k=1}^K \hat{x}_{\text{ref};n,k}}{K}} = \frac{\checkmark_{\text{diff};n}}{\bar{x}_{\text{ref};n}} = \frac{\sum_{k=1}^K \hat{x}_{\text{val};n,k}}{\sum_{k=1}^K \hat{x}_{\text{ref};n,k}} - 1. \quad (39)$$

This estimator gives larger weight to the ratios determined from large $\hat{x}_{\text{val};n,k}$ and $\hat{x}_{\text{ref};n,k}$, because their ratio is less affected by the measurement error. The covariance matrix $\checkmark_{\text{bias},\text{mult.}}$ has the elements

$$\begin{aligned} \checkmark_{\text{bias},\text{mult.};m,n} &= \frac{\checkmark_{\text{val};m,n} \checkmark_{\text{val};m} \checkmark_{\text{val};n}}{\bar{x}_{\text{ref};m} \bar{x}_{\text{ref};n}} \\ &+ \frac{\bar{x}_{\text{val};m} \bar{x}_{\text{val};n} \checkmark_{\text{ref};m,n} \checkmark_{\text{ref};m} \checkmark_{\text{ref};n}}{\bar{x}_{\text{ref};m}^2 \bar{x}_{\text{ref};n}^2} \\ &- \frac{\bar{x}_{\text{val};m} \checkmark_{\text{val};m} \checkmark_{\text{ref};m} \checkmark_{\text{ref};n}}{\bar{x}_{\text{ref};m}^2 \bar{x}_{\text{ref};n}} \\ &- \frac{\bar{x}_{\text{val};n} \checkmark_{\text{val};m} \checkmark_{\text{ref};n} \checkmark_{\text{ref};m}}{\bar{x}_{\text{ref};m}^2 \bar{x}_{\text{ref};n}}, \quad (40) \end{aligned}$$

where

$$\checkmark_{\text{val};m} = \frac{1}{\sqrt{K}} \checkmark_{\text{val};m} = \sqrt{\frac{\sum_{k=1}^K (\hat{x}_{\text{val};m,k} - \bar{x}_{\text{val};m})^2}{K(K-1)}}, \quad (41)$$

$$\checkmark_{\text{ref};m} = \frac{1}{\sqrt{K}} \checkmark_{\text{ref};m} = \sqrt{\frac{\sum_{k=1}^K (\hat{x}_{\text{ref};m,k} - \bar{x}_{\text{ref};m})^2}{K(K-1)}}, \quad (42)$$

$$\checkmark_{\text{val};m,n} = \frac{\sum_{k=1}^K (\hat{x}_{\text{val};m,k} - \bar{x}_{\text{val};m})(\hat{x}_{\text{val};n,k} - \bar{x}_{\text{val};n})}{(K-1)\checkmark_{\text{val};m}\checkmark_{\text{val};n}}, \quad (43)$$

$$\checkmark_{\text{ref};m,n} = \frac{\sum_{k=1}^K (\hat{x}_{\text{ref};m,k} - \bar{x}_{\text{ref};m})(\hat{x}_{\text{ref};n,k} - \bar{x}_{\text{ref};n})}{(K-1)\checkmark_{\text{ref};m}\checkmark_{\text{ref};n}}, \quad (44)$$

and

$$\check{r}_{m,n} = \frac{\sum_{k=1}^K (\hat{x}_{\text{val};m,k} - \bar{\hat{x}}_{\text{val};m})(\hat{x}_{\text{ref};n,k} - \bar{\hat{x}}_{\text{ref};n})}{(K-1)\check{\sigma}_{\text{val};m}\check{\sigma}_{\text{ref};n}}. \quad (45)$$

This gives for the variances

$$\begin{aligned} \check{\sigma}_{\text{bias,mult};n}^2 &= \check{s}_{\text{bias,mult};n,n} \\ &= \frac{1}{\bar{\hat{x}}_{\text{ref};n}^4} (\check{\sigma}_{\text{val};n}^2 \bar{\hat{x}}_{\text{ref};n}^2 + \check{\sigma}_{\text{ref};n}^2 \hat{x}_{\text{val};n}^2 \\ &\quad - 2\check{r}_{n,n}\check{\sigma}_{\text{val};n}\check{\sigma}_{\text{ref};n}\bar{\hat{x}}_{\text{ref};n}\bar{\hat{x}}_{\text{val};n}), \end{aligned} \quad (46)$$

The simplified expression

$$\check{\sigma}_{\text{b,rel};n} = \frac{\check{\sigma}_{\text{bias};n}}{\bar{\hat{x}}_{\text{ref};n}} \quad (47)$$

ignores the uncertainty of $\bar{\hat{x}}_{\text{ref};n}$.

Only if the expected errors of the differences are proportional to reference values, the mean relative deviation of a state parameter x_n at altitude gridpoint n is calculated as

$$\check{b}_{\text{diff,mult};n} = \frac{\sum_{k=1}^K \frac{\hat{x}_{\text{val};n,k} - \hat{x}_{\text{ref};n,k}}{\hat{x}_{\text{ref};n,k}}}{K}. \quad (48)$$

is a better estimate of the multiplicative bias, because in this case the uncertainty of each of the ratios is equal, requiring equal weight of each ratio in the average. The elements of its covariance matrix are calculated as

$$\check{s}_{\text{bias,mult};m,n} = \frac{\sum_{k=1}^K \left(\frac{\hat{x}_{\text{val};m,k} - \hat{x}_{\text{ref};m,k}}{\hat{x}_{\text{ref};m,k}} - \check{b}_{\text{diff,mult};m} \right) \left(\frac{\hat{x}_{\text{val};n,k} - \hat{x}_{\text{ref};n,k}}{\hat{x}_{\text{ref};n,k}} - \check{b}_{\text{diff,mult};n} \right)}{K(K-1)}, \quad (49)$$

4.2 Bias determination by statistical comparison of random samples

It is not necessary to use matched pairs for validation. Random samples are sufficient but any sampling artefacts have to be carefully excluded. A parametrization as suggested in Sect. 3.2, Eq. (15) may help to reduce systematic sampling errors.

When two instruments provide large but independent, i.e. unmatched, random samples of measurements, the bias can be determined as the difference of respective mean values:

$$\check{b}_{\text{diff}} = \frac{\sum_{k=1}^K \hat{x}_{\text{val};k}}{K} - \frac{\sum_{l=1}^L \hat{x}_{\text{ref};l}}{L} = \bar{\hat{x}}_{\text{val}} - \bar{\hat{x}}_{\text{ref}}, \quad (50)$$

where K and L are the respective sample sizes. The respective covariance matrix has the elements

$$\begin{aligned} \check{s}_{\text{bias};m,n} &= \frac{\sum_{k=1}^K (\hat{x}_{\text{val};m,k} - \bar{\hat{x}}_{\text{val};m})(\hat{x}_{\text{val};n,k} - \bar{\hat{x}}_{\text{val};n,k})}{K(K-1)} \\ &\quad + \frac{\sum_{l=1}^L (\hat{x}_{\text{ref};m,l} - \bar{\hat{x}}_{\text{ref};m,l})(\hat{x}_{\text{ref};n,l} - \bar{\hat{x}}_{\text{ref};n,l})}{L(L-1)}. \end{aligned} \quad (51)$$

Obviously, any non-randomness of the samples can cause an apparent bias or hide an existing bias.

5 Precision validation

5.1 Precision determination with matching pairs of measurements

The expectation value of the root mean squares difference of a pair of measurements of the same atmospheric state is the accuracy of the difference. In terms of variances and covariances of a profile of differences, this means that

$$\langle (\hat{x}_{\text{val};m} - \hat{x}_{\text{ref};m})(\hat{x}_{\text{val};n} - \hat{x}_{\text{ref};n}) \rangle = s_{\text{diff};m,n}. \quad (52)$$

For accuracy validation, this can be rewritten in terms of χ^2 statistics, where the actual χ^2 is evaluated from a sample of size K of profiles with N altitude gridpoints each:

$$\langle \chi^2 \rangle = \left\langle \sum_{k=1}^K (\hat{x}_{\text{val};k} - \hat{x}_{\text{ref};k})^T \mathbf{S}_{\text{diff}}^{-1} (\hat{x}_{\text{val};k} - \hat{x}_{\text{ref};k}) \right\rangle = K \times N \quad (53)$$

If, however, there is a bias \mathbf{b}_{diff} between the measurement systems, this should be evaluated in a preceding step (see Sect. 4, Eq. 30) and removed in order to validate the precision of the measurement rather than the accuracy. This leads to the following χ^2 statistics:

$$\begin{aligned} \langle \chi^2 \rangle &= \left\langle \sum_{k=1}^K (\hat{x}_{\text{val};k} - \hat{x}_{\text{ref};k} - \check{b}_{\text{diff}})^T \mathbf{S}_{\text{diff,random}}^{-1} \right. \\ &\quad \left. (\hat{x}_{\text{val};k} - \hat{x}_{\text{ref};k} - \check{b}_{\text{diff}}) \right\rangle \\ &= (K-1)N \end{aligned} \quad (54)$$

$\mathbf{S}_{\text{diff,random}}$ is the random component of \mathbf{S}_{diff} according to Eq. (14) and \check{b}_{diff} has been estimated from the same sample (see, e.g. Ridolfi et al., 2006b², for application to MIPAS temperature validation, or Cortesi et al., 2006³, for ozone validation).

While Eq. (53) can be evaluated for single profiles ($K = 1$), Eq. (54) needs a sample of profiles ($K > 1$) in order to distinguish between precision and bias, unless an altitude-independent bias $\check{b}_{\text{diff}} = (\check{b}_{\text{diff}}, \dots, \check{b}_{\text{diff}})^T$ is assumed, where

$$\check{b}_{\text{diff}} = \frac{\sum_{n=1}^N (\hat{x}_{\text{val};n} - \hat{x}_{\text{ref};n})}{N}. \quad (55)$$

Equation (54) then reads

$$\begin{aligned} \langle \chi^2 \rangle &= \langle (\hat{x}_{\text{val}} - \hat{x}_{\text{ref}} - \check{b}_{\text{diff}})^T \mathbf{S}_{\text{diff,random}}^{-1} (\hat{x}_{\text{val}} - \hat{x}_{\text{ref}} - \check{b}_{\text{diff}}) \rangle \\ &= N - 1. \end{aligned} \quad (56)$$

5.2 Precision validation by comparison of random samples

The scatter of a sample of measurements is composed of both the measurement random error (characterized by covariance matrices $\mathbf{S}_{\text{random,val}}$ or $\mathbf{S}_{\text{random,ref}}$, respectively) and the natural variability (characterized by its covariance matrix \mathbf{S}_{nat}). The natural variability of two randomly sampled data sets,

however, is the same, regardless if we observe the atmosphere with the one or the other instrument. Thus, we have to verify

$$\begin{aligned}\check{\mathbf{S}}_{\text{val,nat}} &= \check{\mathbf{S}}_{\text{val,sample}} - \mathbf{S}_{\text{val,random}} \\ &= \check{\mathbf{S}}_{\text{ref,sample}} - \mathbf{S}_{\text{ref,random}} \\ &= \check{\mathbf{S}}_{\text{ref,nat}},\end{aligned}\quad (57)$$

where the elements of $\check{\mathbf{S}}_{\text{val,sample}}$ are

$$\check{s}_{\text{val,sample};m,n} = \frac{\sum_{k=1}^K (\hat{x}_{\text{val};m,k} - \bar{\hat{x}}_{\text{val};m})(\hat{x}_{\text{val};n,k} - \bar{\hat{x}}_{\text{val};n})}{K-1} \quad (58)$$

and where the elements of $\check{\mathbf{S}}_{\text{ref,sample}}$ are

$$\check{s}_{\text{ref,sample};m,n} = \frac{\sum_{l=1}^L (\hat{x}_{\text{ref};m,l} - \bar{\hat{x}}_{\text{ref};m})(\hat{x}_{\text{ref};n,l} - \bar{\hat{x}}_{\text{ref};n})}{L-1}. \quad (59)$$

Testing if related variances are equal can be performed with the *F*-test (see, e.g., Press et al., 1989). Care must be taken that the samples are really random samples of the same population. Further, this strategy discussed here is particularly sensitive to an artificial reduction of the variability of one of the measurement data sets through the use of retrieval schemes involving Bayesian statistics, where each single profile is pushed towards some a priori information (see, e.g. Rodgers, 2000, for application of Bayesian statistics to retrieval theory). Beyond this, the result of the *F*-test is particularly sensitive to deviations of the actual distributions from normal distributions. While the approach proposed here is valid theoretically, these inherent traps make it difficult to use, and the author is not aware of any actual application to atmospheric measurements.

6 Comparison of a single measurement with a random sample of measurements

If only a single profile measurement is available which does not co-incide with any of the measurements to be validated, it can be checked if this single profile measurement belongs to the distribution defined sample of size *K* of the measurements to be validated. The applicable χ^2 test then uses

$$\chi^2 = (\hat{\mathbf{x}}_{\text{val}} - \hat{\mathbf{x}}_{\text{ref}})^T (\check{\mathbf{S}}_{\text{val,ensemble}} + \mathbf{S}_{\text{val,sys}} + \mathbf{S}_{\text{ref,s+r}} + \mathbf{S}_{\text{smooth,diff}})^{-1} (\hat{\mathbf{x}}_{\text{val}} - \hat{\mathbf{x}}_{\text{ref}}) \quad (60)$$

where $\mathbf{S}_{\text{smooth,diff}}$ characterizes the applicable smoothing error difference, and $\check{\mathbf{S}}_{\text{val,ensemble}}$ is the ex post ensemble covariance matrix of the measurements to be validated. Its elements are calculated as

$$\check{s}_{\text{val,ensemble};m,n} = \sum_{k=1}^K \frac{(\hat{x}_{\text{val};m,k} - \bar{\hat{x}}_{\text{val};m})(\hat{x}_{\text{val};n,k} - \bar{\hat{x}}_{\text{val};n})}{K-1}. \quad (61)$$

Again, considerations as outlined in the context of Eq. (14) may apply.

7 What if full retrieval error covariance matrices are not available?

Without ex ante estimates of the profile covariances available, we cannot draw any quantitative conclusion on the reliability of the retrieved profiles in the sense of χ^2 statistics. Often, however, after debiasing, there are at least no horizontal error correlations to be considered. Then, state variables can be compared and χ^2 statistics can be set up for a large ensemble of size *K* of scalar measurements to be validated $\hat{x}_{\text{val},n,k}$ and reference measurements $\hat{x}_{\text{ref},n,k}$ at a single selected altitude $z(n)$. This corresponds to “map validation” instead of “profile validation”. All formulation discussed in this paper then is applied to the particular case where $N=1$. χ^2 -testing in this application leads to a valid conclusion on the reliability of a measurement $\hat{x}_{\text{val},n}$ at the selected altitude $z(n)$. Of course, this procedure can be performed for all altitudes of interest independently. We consider a profile measurement system validated if we can validate the values at each altitude. If, after debiasing, correlations in the time domain can be excluded, the rationale outlined above also can be applied to time series validation. Ridolfi et al. (2006b)² have combined the map validation and time series validation approach by statistically analyzing differences between MIPAS temperatures and radiosonde temperatures from two stations measured at various times. The statistical analysis was performed for altitude bins defined such that each MIPAS limb scan (i.e. each profile) was represented only once in each bin, justifying to disregard any error correlations in altitude.

8 Conclusions

Recipes and terminology for statistical validation of a profile measurement system have been suggested which cover both bias and precision validation and which are applicable to both matched pairs of co-incident measurements and random samples of measurements. Further, a recipe has been suggested to validate profile measurements in a statistical rigorous way even if their full profile covariance matrices are not available. While in real life it will not always be possible to apply these approaches at full rigorosity, validation scientists certainly will find workarounds and simplifications. It is hoped that this paper at least supports better communication in the validation community by suggesting a more or less consistent terminology. Further, ad hoc validation approaches may serve their purpose better, once clarified which rigorous approach they are meant to replace.

Acknowledgements. The author would like to thank S. Ceccherini, K.-H. Fricke, S. Mikuteit, M. Ridolfi, G. Stiller, and the reviewers for helpful comments.

Edited by: P. Hartogh

References

- Bevington, B. R.: Data reduction and error analysis for physical sciences, MacGraw-Hill Book Company, New York, 1969.
- Calisesi, Y., Soebijanta, V. T., and van Oss, R.: Regridding of remote soundings: Formulation and application to ozone profile comparison, *J. Geophys. Res.*, 110, D23306, doi:10.1029/2005JD006122, 2005.
- Carlotti, M., Dinelli, B. M., Raspollini, P., and Ridolfi, M.: Geofit approach to the analysis of limb-scanning satellite measurements, *Appl. Opt.*, 40, 1872–1885, 2001.
- Haseloff, O. W. and Hoffmann, H.-J.: Kleines Lehrbuch der Statistik, de Gruyter, Berlin, 1970.
- Ide, K., Courtier, P., Ghil, M., and Lorenc, A. C.: Unified Notation for Data Assimilation: Operational, Sequential and Variational, *J. Meteorol. Soc. Japan*, 75, 1B, 1997.
- Migliorini, S., Piccolo, C., and Rodgers, C. D.: Intercomparison of direct and indirect measurements: Michelson Interferometer for Passive Atmospheric Sounding (MIPAS) versus sonde ozone profiles, *J. Geophys. Res.*, 109, D19316, doi:10.1029/2004JD004988, 2004.
- Press, W. H., Flannery, B. P., Teukolsky, S. A., and Vetterling, W. T.: Numerical Recipes, Cambridge University Press, Cambridge, 1989.
- Ridolfi, M., Ceccherini, S., and Carli, B.: Optimal interpolation method for intercomparison of atmospheric measurements, *Opt. Lett.*, 31, 855–857, 2006.
- Rodgers, C. D.: Retrieval of Atmospheric Temperature and Composition From Remote Measurements of Thermal Radiation, *Rev. Geophys. Space Phys.*, 14, 609–624, 1976.
- Rodgers, C. D.: Characterization and error analysis of profiles retrieved from remote sounding measurements, *J. Geophys. Res.*, 95, 5587–5595, 1990.
- Rodgers, C. D.: Inverse Methods for Atmospheric Sounding: Theory and Practice, vol. 2 of Series on Atmospheric, Oceanic and Planetary Physics, edited by: Taylor, F. W., World Scientific, 2000.
- Rodgers, C. D. and Connor, B. J.: Intercomparison of remote sounding instruments, *J. Geophys. Res.*, 108, 4116, doi:10.1029/2002JD002299, 2003.
- Steck, T., Höpfner, M., von Clarmann, T., and Grabowski, U.: Tomographic retrieval of atmospheric parameters from infrared limb emission observations, *Appl. Opt.*, 44, 3291–3301, 2005.
- von Clarmann, T. and Grabowski, U.: Elimination of hidden a priori information from remotely sensed profile data, *Atmos. Chem. Phys. Discuss.*, 6, 6723–6751, 2006, <http://www.atmos-chem-phys-discuss.net/6/6723/2006/>.
- Walther, B. A. and Moore, J. L.: The concepts of bias, precision and accuracy, and their use in testing the performance of species richness estimators, with a literature review of estimator performance, *Ecography*, 28, 815–829, 2005.