

How far can the statistical error estimation problem be closed by collocated data?

Annika Vogel^{1,2} and Richard Ménard¹

¹Air Quality Research Division, Environment and Climate Change Canada (ECCC), Dorval - QC, Canada

²Rhenish Institute for Environmental Research (RIU) at the University of Cologne, Cologne, Germany

Correspondence: Annika Vogel (annika.vogel@ec.gc.ca)

Abstract. Accurate specification of error statistics required for data assimilation remains an ongoing challenge, partly because their estimation is an ill-posed problem that requires statistical assumptions. Even with the common assumption that background and observation errors are uncorrelated, the problem remains underdetermined. One natural question that could arise is: Can the increasing amount of overlapping observations or other datasets help to reduce the total number of statistical assumptions, or do they introduce more statistical unknowns? In order to answer this question, this paper provides a conceptual view on the statistical error estimation problem for multiple collocated datasets, including a generalized mathematical formulation, an exemplary demonstration with synthetic data as well as a formulation of the minimal and optimal conditions to solve the problem. It is demonstrated that the required number of statistical assumptions increases linearly with the number of datasets. However the number of error statistics that can be estimated increases quadratically, allowing for an estimation of an increasing number of error cross-statistics between datasets for more than three datasets. The presented optimal estimation of full error covariance- and cross-covariance matrices between dataset does not accumulate uncertainties of assumptions among estimations of subsequent error statistics.

1 Introduction

Accurate specification of error statistics used for data assimilation has been an ongoing challenge. It is known that the accuracy of both, background and observation error covariances have a strong impact on the performance of atmospheric data assimilation (e.g., Daley, 1992a, b; Mitchell and Houtekamer, 2000; Desroziers et al., 2005; Li et al., 2009). A number of approaches to estimate optimal error statistics make use of residuals, i.e. the innovations between observation and background states in observation space (Tandeo et al., 2020), but the error estimation problem remains underdetermined. Different approaches exist which aim at closing the error estimation problem, all of which rely on various assumptions. For example, error variances and correlations were estimated a-posteriori by Tangborn et al. (2002); Ménard and Deshaies-Jacques (2018); Voshtani et al. (2022) based on cross-validation of the analysis with independent observations withheld from the assimilation. However, these a-posteriori methods require an iterative calculation of the analysis and the global minimization criterion provides only spatial-mean estimates of optimal error statistics. In recent years, the amount of available datasets has increased rapidly, including overlapping

or collocated observations from several measurements systems. This arises the question whether multiple overlapping datasets
25 can be used to estimate full spatial fields of optimal error statistics a-priori.

Outside of the field of data assimilation, two different methods were developed that allow for a statistically optimal estimation
of scalar error variances for fully collocated datasets. Although being similar, these two methods were developed independently
from each other in different scientific fields. One method, called the three cornered hat (3CH) method, is based on Grubbs
(1948) and Gray and Allan (1974) who developed an estimation method for error variances of three datasets based on their
30 residuals. This method is widely used in physics for decades, but and has only recently been exploited in meteorology (e.g.,
Anthes and Rieckh, 2018; Rieckh et al., 2021; Kren and Anthes, 2021; Xu and Zou, 2021). Nielsen et al. (2022) and Todling
et al. (2022) were the first to independently use the generalized 3CH (G3CH) method to estimate full error covariances matrices.
Todling et al. (2022) used a modification of the G3CH method to estimate the observation error covariance matrix in a data
assimilation framework. They show that when the corners of G3CH are identified with the observation, background and analysis
35 of variational assimilation procedures, this particular error estimation problem can only be closed under the assumption of
optimally.

Independent from these developments, Stoffelen (1998) used three collocated datasets for multiplicative calibration w.r.t.
each other. Following this idea, the triple collocation (TC) method became a well-known tool to estimate scalar error variances
from residual statistics in the fields of hydrology and oceanography (e.g., Scipal et al., 2008; McColl et al., 2014; Sjoberg et al.,
40 2021). Up to now, only scalar error variances estimated with the TC method are rarely used in data assimilation frameworks
(e.g., Crow and van den Berg, 2010; Crow and Yilmaz, 2014). The 3CH and TC methods use different error models leading to
slightly different assumptions and formulations of error statistics. A detailed description, comparison and evaluation of the two
methods is given in Sjoberg et al. (2021). Both methods have in common that they require fully spatio-temporally collocated
datasets with random errors. These errors are assumed to be independent among the realizations of each dataset with common
45 error statistics across all realizations (e.g., Zwieback et al., 2012; Su et al., 2014). In addition, error statistics of the three
datasets are assumed to be pairwise independent, which is the most critical assumption of these methods (Pan et al., 2015;
Sjoberg et al., 2021).

While the estimation of three error variances is well-established for decades, recent developments propose different ap-
proaches to extend the method to a larger number of datasets. As observed e.g. by Su et al. (2014); Pan et al. (2015); Vogelzang
50 and Stoffelen (2021), the problem of error variance estimation from pairwise residuals becomes overdetermined for more than
three datasets. Su et al. (2014); Anthes and Rieckh (2018); Rieckh et al. (2021) averaged all possible solutions of each error
variance which reduces the sensitivity of the error estimates to inaccurate assumptions. Pan et al. (2015) clustered their datasets
into structural groups and performed a two-step estimation of the in-group errors and the mean errors of each group, which
were assumed to be independent. Zwieback et al. (2012) were the first proposing the additional estimation of the error cross-
55 variances between two selected datasets instead of solving an overdetermined system. This extended collocation (EC) method
was applied to scalar soil moisture datasets by Gruber et al. (2016) who estimated one cross-variance in addition to the error
variances of four datasets. Also for four datasets, Vogelzang and Stoffelen (2021) demonstrated the ability to estimate two
cross-variances in addition to the error variances. They observed that the problem can not be solved for all possible combina-

tions of cross-variances to be estimated. However, their approach failed for five dataset due to a missing generalized condition
60 which is required to solve the problem.

This demonstrates that the different approaches available for more than three datasets provide only an incomplete picture of the problem, where each approach is tailored to the specific conditions of the respective application. Aiming for a more general analysis, this paper approaches the problem from a conceptual point-of-view. The main questions to be answered are: How many error statistics can be extracted from residual statistics between multiple collocated datasets? How many statistics
65 remain to be assumed? How do inaccuracies in assumed error statistics affect different formulations of the estimated error statistics? And what are the minimal and optimal conditions to solve the problem?

In order to answer these questions, the general framework of the estimation problem which builds the basis for the remaining sections is introduced in Sect. 2. It provides a conceptual analysis of the general problem w.r.t. the number of knowns and unknowns and the minimum number of assumptions required. Based on this, the mathematical formulation for non-scalar error
70 matrices is derived in Sect. 3 and Sect. 4, respectively. The derivation is based on the formulation of residual statistics as function of error statistics which is introduced in Sect. 3.2. While the exact formulations of error statistics in Sect. 3.3 remain underdetermined in real applications, approximate formulations which provide a closed system of equations are derived in Sect. 4. Some relations presented in these two sections were already formulated previously for scalar problems dealing with error variances only. However, we present formulations for full covariance matrices including off-diagonal covariances between
75 single elements of the state-vector of the respective dataset, as well as cross-covariance matrices between different datasets. Overlap to previous studies is mainly restricted to the formulation for three datasets in Sect. 4.1 and noted accordingly. Based on this, Sect. 4.2 provides a new approach for the estimation of error statistics of all additional datasets which uses a minimal number of assumptions. The theoretical formulations are applied to four synthetic datasets in Sect. 5. It demonstrates the general ability to estimate full error covariances and cross-statistics as well as effects of inaccurate assumptions w.r.t different
80 setups. The theoretical concept proposed in this study is summarized in Sect. 6. This summary aims at providing the most important results in a general context; answering the main research questions of this study without requiring knowledge of the full mathematical theory. It includes the formulation and illustration of minimal requirements to solve the problem for an arbitrary number of datasets and provides criteria for an optimal setup of those. Finally, Sect. 7 concludes the findings and discusses consequences of using the proposed method in the context of high-dimensional data assimilation.

85 2 General framework

Suppose a system of I spatio-temporally collocated datasets, which may include various model forecasts, observations, analyses and any other datasets available in the same state space. The 2nd moment statistics of the random errors of this system (with respect to the truth) can be described by I error covariances w.r.t. each dataset and N_I error cross-covariances w.r.t. each pair of different datasets. In a discrete state space, (cross-)covariances are matrices and the cross-covariance of dataset A and
90 B is the transposed of the cross-covariance of B and A (see Sect. 3.1 for an explicit definition). Considering this equivalence,

the number N_I of error cross-covariances between all different pairs of datasets is:

$$N_I = \sum_{i=1}^{I-1} i = \frac{1}{2} \cdot I \cdot (I - 1) \quad (1)$$

Thus, the total number U_I of error statistics (error covariances and cross-covariances) is:

$$U_I = N_I + I = \frac{1}{2} \cdot I \cdot (I + 1) \quad (2)$$

95 While error statistics w.r.t. the truth are usually unknown in real applications, residual covariances can be calculated from the residuals between each pair of different datasets. The main idea now is to express the known residual statistics as function of unknown error statistics (Sect. 3.2) and combine these equations to eliminate single error statistics (Sect. 3.3, Sect. 4). Because of $j \neq i$ for residuals, each of the I datasets can be combined with all other $I - 1$ datasets. As residual statistics also do not change with the order of datasets in the residual (see Sect. 3.1), the number of known statistics of the system is also given by
100 N_I as defined in Eq. (1). It will be shown in Sect. 3.2.3 that residual cross-covariances contain generally the same information as residual covariances; thus the N_I residual statistics can be given in form of residual covariances or cross-covariances.

Because N_I residual statistics are known, N_I of the U_I error statistics can be estimated and the remaining I have to be assumed in order to close the problem. The set of error statistics to be estimated can generally be chosen according to the specific application, but it will be shown that there are some constraints. Based on the mathematical theory provided in the
105 following sections, the actual minimal conditions to solve the problem will be discussed in Sect. 6.1.

In most applications of geophysical datasets like in data assimilation, the estimation of error covariances is highly crucial while their error cross-covariances are usually assumed to be negligible. Given the greater need to estimate the I error covariances, the remaining number of error cross-covariances which can be additionally estimated D_I is:

$$D_I = N_I - I = \frac{1}{2} \cdot I \cdot (I - 3) \quad (3)$$

110 The relation between the number of datasets, residual covariances, assumed and estimated error statistics is visualized in Fig. 1. $I = 0$ represents the mathematical extension of the problem, were no error- and residual statistics are required when no dataset is considered. For less than three datasets ($0 < I < 3$), D_I is negative because the number of (known) residual covariances is smaller than the number of (unknown) error covariances ($N_I < I$) and thus the problem is underdetermined even when all datasets are assumed to be independent (=zero error cross-covariances). As it is the case in data assimilation of
115 two datasets ($I = 2$), additional assumptions on error statistics are required. The same holds when only one dataset is available ($I = 1$), were the error covariance of this dataset remains unknown because no residual covariance can be formed. For three datasets ($I = 3$), D_I is zero meaning that the problem is fully determined under the assumption of independent errors ($N_I = I$, formulated in Sect. 4.1).

For more than three datasets ($I > 3$), the number of (known) residual covariances exceeds the number of error covariances
120 which would lead to an overdetermined problem assuming independence among all datasets. Instead of solving an overdetermined problem, the additional information can be used to calculate some error cross-covariances (formulated in Sect. 4.2).

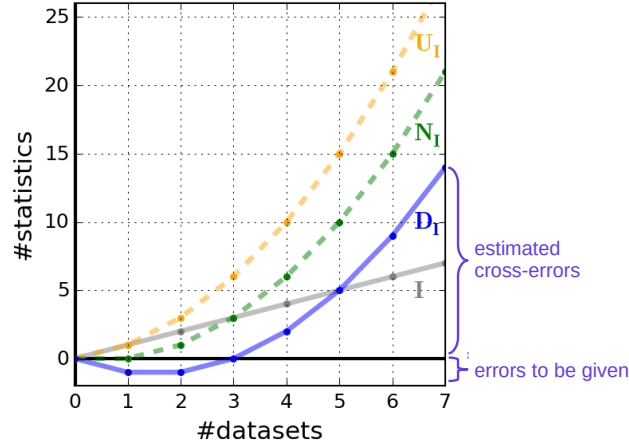


Figure 1. Relation between different numbers of statistics (covariances and cross-covariances) as function of the number of datasets. Shown are I in solid gray (#datasets, #error covariances, #required assumptions), U_I in dashed orange (#error statistics), N_I in dashed green (#residual covariances, #error dependencies, #estimated error statistics), and D_I in solid blue (#estimated error dependencies).

In other words, for $I > 3$ not all datasets need to be assumed to be independent; were D_I gives the number of error cross-covariances which can be estimated in addition to the error covariances from all datasets. For example, half of the error cross-covariances can be estimated for $I = 5$ ($\frac{D_5}{N_5} = \frac{5}{10}$), while two-thirds of them can be estimated for $I = 7$ ($\frac{D_7}{N_7} = \frac{14}{21}$).

125 Although the relative amount of error cross-covariances which can be estimated increases with the number of datasets, an increasing number of $U_I - N_I = I$ assumptions – equal to the number of datasets – is required in order to close the problem because of $U_I > N_I, \forall I > 0$.

Note that almost all numbers presented above apply to the general case were any combination of error covariances and cross-covariances may be given or assumed. While the interpretation of the numbers I , N_I and U_I remains the same in all
130 cases, the only difference is the interpretation of D_I which is less meaningful when also error covariances are assumed.

3 Mathematical theory: Exact formulation

This section gives the theoretical formulation for exact statistical formulations of complete error covariance- and cross-covariance matrices from fully spatio-temporally collocated datasets. Similar to the 3CH method, the errors are assumed to be random, independent among different realizations, but with common error statistics for each dataset. The notation is introduced
135 in Sect. 3.1. While the true state and thus error statistics w.r.t. the truth are usually unknown, residual statistics can be calculated from residuals between each pair of datasets. At the same time, residual statistics contain information about error statistics of the datasets involved. The expression of residual statistics as function of error covariances and cross-covariances in Sect. 3.2 provides the basis for the subsequent mathematical theory. Based on these forward relations, inverse relations describe error statistics as function of residual statistics. The general equations of inverse relations are given in Sect. 3.3 which result in a

140 highly underdetermined system of equations. Closed formulations of error statistics for three and more datasets under certain assumptions will be formulated in the subsequent Sect. 4.

This first part of the mathematical theory includes the following new elements: (i) the separation of cross-statistics into a symmetric error dependency and an error asymmetry (Sect. 3.1), (ii) the general formulation of residual statistics as function of error statistics (Sect. 3.2.1 and 3.2.2), (iii) the demonstration of equivalence between residual covariances and cross-covariances
 145 (Sect. 3.2.3), (iv) the general formulation of exact relations between residual- and error statistics (Sect. 3.3).

3.1 Notation

Suppose I datasets, each containing R realizations of spatio-temporally collocated state vectors $\mathbf{x}_i, \forall i \in [1, I]$. Without loss of generality, the following formulation uses unbiased state vectors with zero mean. In practice, each index i, j, k, l may represent any geophysical dataset like model forecasts, climatologies, in-situ- or remote sensing observations, or other datasets.

150 Let $\mathbf{\Gamma}_{i-j; k-l}$ be the *residual cross-covariance matrix* between dataset residuals $i-j$ and $k-l$ with $j \neq i$ and $l \neq k$, where each element (p, q) is given by the expectation over all realizations:

$$\mathbf{\Gamma}_{i-j; k-l}(p, q) := \overline{[\mathbf{x}_i(p) - \mathbf{x}_j(p)]} [\mathbf{x}_k(q) - \mathbf{x}_l(q)] \quad (4)$$

and the *error cross-covariance matrix* $\mathbf{X}_{\tilde{i}; \tilde{j}}$ between the errors of two datasets i and j w.r.t. the true state \mathbf{x}_T :

$$\mathbf{X}_{\tilde{i}; \tilde{j}}(p, q) := \overline{[\mathbf{x}_i(p) - \mathbf{x}_T(p)]} [\mathbf{x}_j(q) - \mathbf{x}_T(q)] \quad (5)$$

155 where the tilde above a dataset index indicates its deviation from the truth and the overbar denotes the expectation over all R realizations. Note that $x_i(p)$ is a scalar element of the dataset vector.

In the symmetric case, each element (p, q) of the *residual covariance matrix* of $i-j$ with $j \neq i$, is given by:

$$\mathbf{\Gamma}_{i-j}(p, q) := \mathbf{\Gamma}_{i-j; i-j}(p, q) \stackrel{(4)}{=} \overline{[\mathbf{x}_i(p) - \mathbf{x}_j(p)]} [\mathbf{x}_i(q) - \mathbf{x}_j(q)] \quad (6)$$

and the *error covariance matrix* $\mathbf{C}_{\tilde{i}}$ of a dataset i w.r.t. the true state \mathbf{x}_T :

$$160 \quad \mathbf{C}_{\tilde{i}}(p, q) := \mathbf{X}_{\tilde{i}; \tilde{i}}(p, q) \stackrel{(5)}{=} \overline{[\mathbf{x}_i(p) - \mathbf{x}_T(p)]} [\mathbf{x}_i(q) - \mathbf{x}_T(q)] \quad (7)$$

where numbers in parenthesis above an equal sign indicate other equations that were used to retrieve the right hand side.

Note that residual- and error cross-covariance matrices are generally asymmetric in the non-scalar formulation presented here, but the following relations hold for residual as well as similarly for error cross-covariance matrices:

$$\mathbf{\Gamma}_{i-j; k-l} \stackrel{(4)}{=} -\mathbf{\Gamma}_{j-i; k-l} \stackrel{(4)}{=} -\mathbf{\Gamma}_{i-j; l-k} \stackrel{(4)}{=} \mathbf{\Gamma}_{j-i; l-k} \quad (8)$$

$$165 \quad \mathbf{\Gamma}_{i-j; k-l} \stackrel{(4)}{=} \left[\mathbf{\Gamma}_{k-l; i-j} \right]^T \quad (9)$$

$$\mathbf{X}_{\tilde{i}; \tilde{j}} \stackrel{(5)}{=} \left[\mathbf{X}_{\tilde{j}; \tilde{i}} \right]^T \quad (10)$$

The symmetric properties of residual- and error covariances follow directly from their definition:

$$\mathbf{\Gamma}_{i-j} \stackrel{(6)}{=} \mathbf{\Gamma}_{j-i} \quad (11)$$

$$\left[\mathbf{\Gamma}_{i-j}\right]^T \stackrel{(6)}{=} \mathbf{\Gamma}_{i-j} \quad (12)$$

170 The sum of an (asymmetric) cross-covariance matrix and its transposed is denoted as *dependency*. For example, the sum of error cross-covariance matrices between i and j is denoted as *error dependency matrix* $\mathbf{D}_{i,j}^{\sim}$:

$$\mathbf{D}_{i,j}^{\sim} := \mathbf{X}_{i,j}^{\sim} + \mathbf{X}_{j,i}^{\sim} \quad (13)$$

Although error cross-covariances may be asymmetric, the error dependency matrix is symmetric by definition:

$$\mathbf{D}_{i,j}^{\sim} \stackrel{(13)}{=} \mathbf{X}_{i,j}^{\sim} + \mathbf{X}_{j,i}^{\sim} \stackrel{(13)}{=} \mathbf{D}_{j,i}^{\sim} \quad (14)$$

$$175 \quad \mathbf{D}_{i,j}^{\sim} \stackrel{(13)}{=} \mathbf{X}_{i,j}^{\sim} + \mathbf{X}_{j,i}^{\sim} \stackrel{(10)}{=} \left[\mathbf{X}_{j,i}^{\sim}\right]^T + \left[\mathbf{X}_{i,j}^{\sim}\right]^T \stackrel{(13)}{=} \left[\mathbf{D}_{i,j}^{\sim}\right]^T \quad (15)$$

Likewise, the sum of the residual cross-covariance matrices between $i-j$ and $k-l$ with $j \neq i$ and $l \neq k$, is denoted as *residual dependency matrix* $\mathbf{D}_{i,j}^{\sim}$:

$$\mathbf{D}_{i-j;k-l} := \mathbf{\Gamma}_{i-j;k-l} + \mathbf{\Gamma}_{k-l;i-j} \quad (16)$$

The difference between a cross-covariance matrix and its transposed is a measure of asymmetry in the cross-covariances
180 and is therefore denoted as *asymmetry*. For example, difference between the error cross-covariance matrices between i and j is denoted as *error asymmetry matrix* $\mathbf{Y}_{i,j}^{\sim}$:

$$\mathbf{Y}_{i,j}^{\sim} := \mathbf{X}_{i,j}^{\sim} - \mathbf{X}_{j,i}^{\sim} \quad (17)$$

Likewise, the difference between the residual cross-covariance matrices between $i-j$ and $k-l$ with $j \neq i$ and $l \neq k$, is denoted as *residual asymmetry matrix* $\mathbf{D}_{i,j}^{\sim}$:

$$185 \quad \mathbf{Y}_{i-j;k-l} := \mathbf{\Gamma}_{i-j;k-l} - \mathbf{\Gamma}_{k-l;i-j} \quad (18)$$

3.2 Residual statistics

For real geophysical problems, the available statistical information are (i) residual covariance matrices of each pair of datasets and (ii) residual cross-covariance matrices between different residuals of datasets. The forward relations of residual covariances and residual cross-covariances as function of error statistics are formulated in the following. For the estimation of error
190 statistics, it is important to quantify the number of independent input statistics which determines the number of possible error estimations. Therefore, this section also includes an evaluation of the relation between residual cross-covariances and residual covariances in order to specify the additional information content of residual cross-covariances.

3.2.1 Residual covariances

Each element (p, q) of the residual covariance matrix between two input datasets i and j can be written as function of their error statistics as follows:

$$\begin{aligned}\mathbf{\Gamma}_{i-j}(p, q) &\stackrel{(6)}{=} \overline{\left\{ [\mathbf{x}_i(p) - \mathbf{x}_T(p)] - [\mathbf{x}_j(p) - \mathbf{x}_T(p)] \right\} \left\{ [\mathbf{x}_i(q) - \mathbf{x}_T(q)] - [\mathbf{x}_j(q) - \mathbf{x}_T(q)] \right\}} \\ &\stackrel{(5),(7)}{=} \mathbf{C}_{\tilde{i}}(p, q) - \mathbf{X}_{\tilde{i};\tilde{j}}(p, q) - \mathbf{X}_{\tilde{j};\tilde{i}}(p, q) + \mathbf{C}_{\tilde{j}}(p, q)\end{aligned}\quad (19)$$

Thus the complete residual covariance matrix of $i - j$ is expressed as:

$$\mathbf{\Gamma}_{i-j} \stackrel{(19)}{=} \underbrace{\mathbf{C}_{\tilde{i}} + \mathbf{C}_{\tilde{j}}}_{\text{"independent residual"}} - \underbrace{\left[\mathbf{X}_{\tilde{i};\tilde{j}} + \mathbf{X}_{\tilde{j};\tilde{i}} \right]}_{\text{"error dependency"} =: \mathbf{D}_{\tilde{i};\tilde{j}}} \quad (20)$$

Equation (20) is an exact formulation of the complete residual covariance matrix of any pair of datasets $i - j$. It holds for all combinations of datasets without any further assumption like independent- or unbiased error statistics. Thus, the residual covariance of any dataset pair consists of (i) the independent residual associated with sum of the error covariances of each dataset, minus (ii) the error dependency corresponding to the sum of their error cross-covariances.

Note that although the error dependency matrix is symmetric by definition, it is the sum of two error cross-covariances which are generally asymmetric and thus differ in the non-scalar formulation. In the scalar case, the two error cross-covariances reduce to their common error cross-variance and the residual covariance reduces to the scalar formulation of the variance as e.g. in Anthes and Rieckh (2018); Sjoberg et al. (2021).

3.2.2 Residual cross-covariances

Each element (p, q) of the residual cross-covariance matrix between two input datasets $i - j$ and $k - l$ can be written as function of their error cross-covariances:

$$\begin{aligned}\mathbf{\Gamma}_{i-j;k-l}(p, q) &\stackrel{(4)}{=} \overline{\left\{ [\mathbf{x}_i(p) - \mathbf{x}_T(p)] - [\mathbf{x}_j(p) - \mathbf{x}_T(p)] \right\} \left\{ [\mathbf{x}_k(q) - \mathbf{x}_T(q)] - [\mathbf{x}_l(q) - \mathbf{x}_T(q)] \right\}} \\ &\stackrel{(5)}{=} \mathbf{X}_{\tilde{i};\tilde{k}}(p, q) - \mathbf{X}_{\tilde{i};\tilde{l}}(p, q) - \mathbf{X}_{\tilde{j};\tilde{k}}(p, q) + \mathbf{X}_{\tilde{j};\tilde{l}}(p, q)\end{aligned}\quad (21)$$

And thus the complete residual cross-covariance matrix between $i - j$ and $k - l$:

$$\mathbf{\Gamma}_{i-j;k-l} \stackrel{(21)}{=} \mathbf{X}_{\tilde{i};\tilde{k}} - \mathbf{X}_{\tilde{i};\tilde{l}} - \mathbf{X}_{\tilde{j};\tilde{k}} + \mathbf{X}_{\tilde{j};\tilde{l}} \quad (22)$$

Equation (22) is a generalized form of Eq. (20) with residuals between different datasets ($i - j$; $k - l$). It consists of four error cross-covariances of the datasets involved. This formulation of residual statistics as function of error statistics provides the basis for the complete theoretical derivation of error estimates in this study. In contrast to the symmetric residual covariance matrix, the residual cross-covariance matrix may be asymmetric for asymmetric error cross-covariances.

3.2.3 Relation of residual statistics

220 In the following, it is demonstrated that combinations of residual cross-covariances contain the same statistical information as residual covariance matrices.

For $k = i$, the residual dependency between $i - j$ and $i - l$ becomes:

$$\begin{aligned}
 \mathbf{\Gamma}_{i-j;i-l} + \mathbf{\Gamma}_{i-l;i-j} &\stackrel{(21)}{=} \mathbf{C}_{\tilde{i}} - \mathbf{X}_{\tilde{i};\tilde{l}} - \mathbf{X}_{\tilde{j};\tilde{i}} + \mathbf{X}_{\tilde{j};\tilde{l}} + \mathbf{C}_{\tilde{i}} - \mathbf{X}_{\tilde{i};\tilde{j}} - \mathbf{X}_{\tilde{l};\tilde{i}} + \mathbf{X}_{\tilde{l};\tilde{j}} \\
 &\stackrel{(13)}{=} 2 \mathbf{C}_{\tilde{i}} - \mathbf{D}_{\tilde{i};\tilde{l}} - \mathbf{D}_{\tilde{j};\tilde{i}} + \mathbf{D}_{\tilde{j};\tilde{l}} + [\mathbf{C}_{\tilde{j}} + \mathbf{C}_{\tilde{j}}] - [\mathbf{C}_{\tilde{j}} + \mathbf{C}_{\tilde{j}}] \\
 225 \quad &= \mathbf{C}_{\tilde{i}} + \mathbf{C}_{\tilde{l}} - \mathbf{D}_{\tilde{i};\tilde{l}} + \mathbf{C}_{\tilde{j}} + \mathbf{C}_{\tilde{i}} - \mathbf{D}_{\tilde{j};\tilde{i}} - \mathbf{C}_{\tilde{j}} - \mathbf{C}_{\tilde{l}} + \mathbf{D}_{\tilde{j};\tilde{l}} \\
 &\stackrel{(20)}{=} \mathbf{\Gamma}_{i-l} + \mathbf{\Gamma}_{j-i} - \mathbf{\Gamma}_{j-l}
 \end{aligned} \tag{23}$$

The relation between residual covariances and residual cross-covariances in Eq. (23) is exact and holds for all datasets without any further assumption. In case of symmetric residual cross-covariances $(\mathbf{\Gamma}_{i-j;i-l} = \mathbf{\Gamma}_{i-l;i-j} \stackrel{(23)}{=} \frac{1}{2} [\mathbf{\Gamma}_{i-l} + \mathbf{\Gamma}_{j-i} - \mathbf{\Gamma}_{j-l}])$, the residual cross-covariance matrices are fully determined by the symmetric residual covariances.

230 In the general asymmetric case, Eq. (23) can be rewritten as:

$$\begin{aligned}
 \mathbf{\Gamma}_{i-l} + \mathbf{\Gamma}_{j-i} - \mathbf{\Gamma}_{j-l} &\stackrel{(23)}{=} \mathbf{\Gamma}_{i-j;i-l} + \mathbf{\Gamma}_{i-l;i-j} \stackrel{(18)}{=} \mathbf{\Gamma}_{i-j;i-l} + [\mathbf{\Gamma}_{i-j;i-l} - \mathbf{Y}_{i-j;i-l}] \\
 \iff \mathbf{\Gamma}_{i-j;i-l} &= \frac{1}{2} [\mathbf{\Gamma}_{i-l} + \mathbf{\Gamma}_{j-i} - \mathbf{\Gamma}_{j-l}] + \frac{1}{2} \mathbf{Y}_{i-j;i-l}
 \end{aligned} \tag{24}$$

Equation (24) shows that each individual residual cross-covariance consists of a symmetric contribution including residual covariances between the datasets and an asymmetric contribution being half of the related residual asymmetry matrix. Thus, 235 residual cross-covariances may only provide additional information on asymmetries of error statistics, but not on symmetric statistics (like error covariances).

3.3 Exact error statistics

As an extension to previous work, this section provides generalized formulations of error covariances, cross-covariances, and dependencies in matrix form. These formulations are based on the relations between residual- and error statistics in Eq. (20) and Eq. (22). Note that the general formulations presented here do not provide a closed system of equations which can be 240 solved in real applications. They serve as basis for the approximate solutions which are formulated in the subsequent section.

3.3.1 Error statistics from residual covariances

Equation (20) shows that each residual covariance matrix can be expressed by the error covariances of the two datasets involved and their error dependency. The goal is to find an inverse formulation of an error covariance matrix as function of residual 245 covariances which does not include other (unknown) error covariances matrices. By combining the formulations of three residuals $\mathbf{\Gamma}_{i-j}$, $\mathbf{\Gamma}_{j-k}$, and $\mathbf{\Gamma}_{k-i}$ between the same three datasets i , j , and k and expressing each using Eq. (20), a single error

covariance can be eliminated:

$$\mathbf{C}_{\tilde{i}}^{(20)ij} \equiv \mathbf{\Gamma}_{i-j} + \mathbf{D}_{\tilde{i};\tilde{j}} - \mathbf{C}_{\tilde{j}} \quad (25)$$

$$\begin{aligned} & \stackrel{(20)jk}{=} \mathbf{\Gamma}_{i-j} + \mathbf{D}_{\tilde{i};\tilde{j}} - \mathbf{\Gamma}_{j-k} - \mathbf{D}_{\tilde{j};\tilde{k}} + \mathbf{C}_{\tilde{k}} \\ 250 \quad & \stackrel{(20)ki}{=} \mathbf{\Gamma}_{i-j} + \mathbf{D}_{\tilde{i};\tilde{j}} - \mathbf{\Gamma}_{j-k} - \mathbf{D}_{\tilde{j};\tilde{k}} + \mathbf{\Gamma}_{k-i} + \mathbf{D}_{\tilde{k};\tilde{i}} - \mathbf{C}_{\tilde{i}} \\ \iff \quad & \mathbf{C}_{\tilde{i}} = \frac{1}{2} \left[\underbrace{\mathbf{\Gamma}_{i-j} + \mathbf{\Gamma}_{k-i} - \mathbf{\Gamma}_{j-k}}_{\text{"independent contribution"}} + \underbrace{\mathbf{D}_{\tilde{i};\tilde{j}} + \mathbf{D}_{\tilde{k};\tilde{i}} - \mathbf{D}_{\tilde{j};\tilde{k}}}_{\text{"dependent contribution"}} \right] \end{aligned} \quad (26)$$

were the indication of used equations above the equal signs are extended by indices which denote to which datasets this equation has been applied. For example, " $\stackrel{(20)ki}{=}$ " indicates that the relation in Eq. (20) was applied to datasets k and i to achieve the right hand side.

255 Equation (26) provides a general formulation of error covariances as function of residual covariances and error dependencies. It holds for all combinations of datasets without any further assumption (e.g. independence). Thus, each error covariance can be formulated as a sum of an independent contribution of three residual covariances w.r.t. any pair of other datasets and an dependent contribution of the three related error dependencies. While the independent contribution can be calculated from residual statistics between input datasets, the dependent contribution is generally unknown in real applications.

260 Given I datasets, the total number of different formulations of each error covariance in Eq. (26) is determined by the number of different pairs of the other datasets which is $\sum_{i=1}^{I-2} i = \frac{1}{2}(I-1)(I-2)$ (see also Sjoberg et al., 2021). The scalar equivalent of Eq. (26) where the dependency matrices reduce to twice the error cross-variances has been formulated previously in the 3CH method e.g. in Anthes and Rieckh (2018); Sjoberg et al. (2021). Very recently, the full matrix form was used by Nielsen et al. (2022); Todling et al. (2022). Note that in the literature, the dependent contribution in Eq. (26) is denoted as cross-covariances
265 between the errors.

A formulation of each individual error dependency matrix as function of the error covariances of the two datasets and their residual covariance results directly from Eq. (20):

$$\mathbf{D}_{\tilde{i};\tilde{j}} \stackrel{(20)}{=} \mathbf{C}_{\tilde{i}} + \mathbf{C}_{\tilde{j}} - \mathbf{\Gamma}_{i-j} \quad (27)$$

270 Being a symmetric matrix, residual covariances cannot provide information on error asymmetries and thus on asymmetric components of error cross-covariances. Only the symmetric component of error cross-covariances could be estimated from half the error dependency which is equivalent to a zero error asymmetry matrix:

$$\mathbf{D}_{\tilde{i};\tilde{j}} + \mathbf{Y}_{\tilde{i};\tilde{j}} \stackrel{(13),(17)}{=} \left[\mathbf{X}_{\tilde{i};\tilde{j}} + \cancel{\mathbf{X}_{\tilde{j};\tilde{i}}} \right] + \left[\mathbf{X}_{\tilde{i};\tilde{j}} - \cancel{\mathbf{X}_{\tilde{j};\tilde{i}}} \right] \iff \mathbf{X}_{\tilde{i};\tilde{j}} = \frac{1}{2} \left[\mathbf{D}_{\tilde{i};\tilde{j}} + \mathbf{Y}_{\tilde{i};\tilde{j}} \right] \quad (28)$$

3.3.2 Error statistics from residual cross-covariances

275 The general forward formulation of residual cross-covariances in Eq. (22) consists of error cross-covariances of the four datasets involved. Setting for example $k = i$, provides an inverse formulation of error covariances of i :

$$\Gamma_{i-j;i-l} \stackrel{(22)}{=} \mathbf{C}_{\tilde{i}} - \mathbf{X}_{\tilde{i};l}^{\sim} - \mathbf{X}_{\tilde{j};i}^{\sim} + \mathbf{X}_{\tilde{j};l}^{\sim} \iff \mathbf{C}_{\tilde{i}} = \Gamma_{i-j;i-l} + \mathbf{X}_{\tilde{i};l}^{\sim} + \mathbf{X}_{\tilde{j};i}^{\sim} - \mathbf{X}_{\tilde{j};l}^{\sim} \quad (29)$$

The scalar formulation of Eq. (29) was previously given in Zwieback et al. (2012).

Similarly to Eq. (26) from residual covariances, the number of formulations of each error covariance from different pairs of other datasets in Eq. (29) is $\sum_{i=1}^{I-2} i = \frac{1}{2}(I-1)(I-2)$. In addition, there are four possibilities to write each error covariances from the same pairs of other datasets using the relations of residual cross-covariances in Eq. (8). Each of the four possibilities results from setting one pair of datasets in definition of residual cross-covariances in Eq. (22) equal.

Two of the error cross-covariances in Eq. (29) can be rewritten by applying Eq. (29) to the error covariance of dataset j :

$$\mathbf{C}_{\tilde{j}} \stackrel{(29)_j}{=} \Gamma_{j-i;j-l} + \mathbf{X}_{\tilde{j};l}^{\sim} + \mathbf{X}_{\tilde{i};j}^{\sim} - \mathbf{X}_{\tilde{i};l}^{\sim} \iff \mathbf{X}_{\tilde{i};l}^{\sim} - \mathbf{X}_{\tilde{j};l}^{\sim} = \Gamma_{j-i;j-l} + \mathbf{X}_{\tilde{i};j}^{\sim} - \mathbf{C}_{\tilde{j}} \quad (30)$$

285 With this, Eq. (29) becomes:

$$\mathbf{C}_{\tilde{i}} \stackrel{(30)}{=} \Gamma_{i-j;i-l} + \Gamma_{j-i;j-l} - \mathbf{C}_{\tilde{j}} + \mathbf{X}_{\tilde{i};j}^{\sim} + \mathbf{X}_{\tilde{j};i}^{\sim} \stackrel{(13)}{=} \Gamma_{i-j;i-l} + \Gamma_{j-i;j-l} - \mathbf{C}_{\tilde{j}} + \mathbf{D}_{\tilde{i};j}^{\sim} \quad (31)$$

Because the residual cross-covariances can be rewritten as:

$$\Gamma_{i-j;i-l} + \Gamma_{j-i;j-l} \stackrel{(22)}{=} \mathbf{C}_{\tilde{i}} - \cancel{\mathbf{X}_{\tilde{i};l}^{\sim}} - \cancel{\mathbf{X}_{\tilde{j};i}^{\sim}} + \cancel{\mathbf{X}_{\tilde{j};l}^{\sim}} + \mathbf{C}_{\tilde{j}} - \cancel{\mathbf{X}_{\tilde{j};l}^{\sim}} - \cancel{\mathbf{X}_{\tilde{i};j}^{\sim}} + \mathbf{X}_{\tilde{i};l}^{\sim} \stackrel{(13)}{=} \mathbf{C}_{\tilde{i}} + \mathbf{C}_{\tilde{j}} - \mathbf{D}_{\tilde{i};j}^{\sim} \stackrel{(20)}{=} \Gamma_{i-j} \quad (32)$$

the formulation of error covariances based on residual cross-covariances in Eq. (31) is symmetric and equivalent to the formulation based on residual covariances from Eq. (25).

The forward formulation of residual cross-covariances does not allow for an elimination of one single error cross-covariance even when multiple equations are combined. One formulation of an error cross-covariance matrix as function of residual cross-covariances results directly from the forward relation:

$$295 \quad \mathbf{X}_{\tilde{j};l}^{\sim} \stackrel{(29)}{=} \Gamma_{i-j;i-l} - \mathbf{C}_{\tilde{i}} + \mathbf{X}_{\tilde{i};l}^{\sim} + \mathbf{X}_{\tilde{j};i}^{\sim} \quad (33)$$

Note that the third dataset i on the right hand side of Eq. (33) can be any other dataset ($i \neq j, i \neq l$). Thus for any $I > 2$, there are $I - 2$ formulations of each error cross-covariance $\mathbf{X}_{\tilde{j};l}^{\sim}$, which are all equivalent in the exact formulation.

Any of the formulations of error cross-covariances can also be used for a formulation of the error dependency matrix

$$\begin{aligned} & \mathbf{D}_{\tilde{j};l}^{\sim} \Big|_{\text{cross}} \text{ which is equivalent to the formulation based on residual covariances } \mathbf{D}_{\tilde{j};l}^{\sim} \Big|_{\text{covar}} : \\ 300 \quad & \mathbf{D}_{\tilde{j};l}^{\sim} \Big|_{\text{cross}} \stackrel{(13)}{=} \mathbf{X}_{\tilde{j};l}^{\sim} + \mathbf{X}_{\tilde{l};j}^{\sim} \stackrel{(33)}{=} \Gamma_{j-i;l-i} - \mathbf{C}_{\tilde{i}} + \mathbf{X}_{\tilde{j};i}^{\sim} + \mathbf{X}_{\tilde{i};l}^{\sim} + \Gamma_{l-i;j-i} - \mathbf{C}_{\tilde{i}} + \mathbf{X}_{\tilde{i};j}^{\sim} + \mathbf{X}_{\tilde{l};i}^{\sim} \\ & \stackrel{(13)}{=} \Gamma_{j-i;l-i} + \Gamma_{l-i;j-i} - 2 \mathbf{C}_{\tilde{i}} + \mathbf{D}_{\tilde{i};j}^{\sim} + \mathbf{D}_{\tilde{i};l}^{\sim} \stackrel{(23)}{=} \Gamma_{i-j} + \mathbf{D}_{\tilde{i};j}^{\sim} + \Gamma_{i-l} + \mathbf{D}_{\tilde{i};l}^{\sim} - \Gamma_{j-l} - 2 \mathbf{C}_{\tilde{i}} \\ & \stackrel{(20)}{=} \cancel{\mathbf{C}_{\tilde{i}}} + \mathbf{C}_{\tilde{j}} + \cancel{\mathbf{C}_{\tilde{i}}} + \mathbf{C}_{\tilde{l}} - \Gamma_{j-l} - 2 \mathbf{C}_{\tilde{i}} = \mathbf{C}_{\tilde{j}} + \mathbf{C}_{\tilde{l}} - \Gamma_{j-l} \stackrel{(27)}{=} \mathbf{D}_{\tilde{j};l}^{\sim} \Big|_{\text{covar}} \end{aligned} \quad (34)$$

The equivalence demonstrates that the exact formulations of error statistics from residual covariances and cross-covariances are consistent to each other. This consistency applies to the exact formulations of all symmetric error statistics (error covari-
 305 ances and dependencies) and results directly from the fact that the basic formulation of residual covariances in Eq. (20) is a special case of the formulation of residual cross-covariances in Eq. (22).

4 Mathematical theory: Approximate formulation

Based on the exact formulations in Sect. 3 which remain underdetermined in real applications, this section provides approx-
 310 imate formulations for three and more datasets which provide a closed system of equations. Section 4.1 describes the long-
 known closure of the system for three datasets, but for full covariance matrices. An extension for any additional dataset $I > 3$
 using a minimal number of assumptions is introduced in Sect. 4.2. It includes the estimation of additional error covariances
 and some error cross-statistics to estimate a maximum amount of error statistics.

In addition to the optimal extension to more than three datasets, this second part of the mathematical theory includes the
 following new elements: (i) the analysis of differences from residual covariance- and cross-covariance estimates (Sect. 4.1.2),
 315 (ii) the determination of uncertainties caused by assumed error statistics (Sect. 4.1.3 and 4.2.3), and (iii) the comparison of the
 approximation from three- ("triangular estimation") and more ("sequential estimation") datasets (Sect. 4.2.4).

4.1 Approximation for three datasets

As demonstrated in Sect. 2, at least three collocated datasets are required to estimate all error covariances ($U_I \geq 0$). For three
 datasets ($I = 3$), three residual covariances ($N_3 = 3$) can be calculated between each pair of datasets. At the same time, there
 320 are six unknown error statistics ($U_3 = 6$): three error covariances and three error cross-statistics (cross-covariances or depen-
 dencies). Thus, the problem is under-determined and three error statistics ($U_3 - N_3 = 3$) have to be assumed in order to close
 the system. The most common approach, which is also used in 3CH and TC methods, is to assume zero error cross-statistics
 between all pairs of datasets: $\mathbf{X}_{\tilde{i}\tilde{j}} = 0 \Leftrightarrow \mathbf{D}_{\tilde{i}\tilde{j}} = 0, \forall i, j \in [1, 3], j \neq i$. The approximation of the three error covariances
 can also be formulated in a Hilbert space which allows for an illustrative geometric interpretation as in Pan et al. (2015)
 325 (their Fig. 1). Because the assumption of zero cross-covariance equals zero error dependency, it is denoted as "assumption of
 independence" or "independence assumption" thereafter.

The independence assumption resembles the innovation covariance consistency of data assimilation, where the residual co-
 variance between background and observation datasets - denoted as innovation covariance - is assumed to be equal to the sum
 of their error covariances in the formulation of the analysis (e.g., Daley, 1992b; Ménard, 2016):

$$330 \quad \Gamma_{i-j}^{(20)} \approx_{\{in\}} \mathbf{C}_{\tilde{i}} + \mathbf{C}_{\tilde{j}} \quad (35)$$

where " $\approx_{\{in\}}$ " indicates the assumption of independence between the two datasets, i.e. $\mathbf{X}_{\tilde{i}\tilde{j}} = 0$.

Because all error cross-statistics need to be assumed in this setup, approximations of these cross-covariances and dependen-
 cies only reproduces the initially assumed statistics and do not provide any new information.

4.1.1 Error covariance estimates

335 Assuming independent error statistics among all three datasets, or similarly that error dependencies are negligible compared to residual covariances $\mathbf{D}_{i;\tilde{j}} \ll \mathbf{\Gamma}_{i-j}, \forall j \neq i$ gives an estimate of each error covariance matrix as function of three residual covariances:

$$\mathbf{C}_{i_{\{in3\}}}^{(26)} \approx \frac{1}{2} [\mathbf{\Gamma}_{i-j} + \mathbf{\Gamma}_{k-i} - \mathbf{\Gamma}_{j-k}] \quad (36)$$

Were " \approx " indicates the assumption of independence among all three datasets involved.

340 In the scalar case, Eq. (36) reduces to the equivalent formulation for error variances known from the TC and 3CH method (e.g., Pan et al., 2015; Sjoberg et al., 2021). Thus, the long-known 3CH estimation of error variances from residual variances among three datasets holds similarly for complete error covariance matrices from residual covariances under the independence assumption. In fact, the approximation in Eq. (36) requires only the assumption that the dependent contribution of Eq. (26) vanishes. However combining this condition for the error covariance estimates of all three datasets results in the need for each error dependency to be zero.

Under the assumption of independence among all three datasets $\mathbf{X}_{i;\tilde{j}} = 0, \forall i, j$, their error covariance matrices can also be directly estimated from residual cross-covariances:

$$\mathbf{C}_{i_{\{in3\}}}^{(29)} \approx \mathbf{\Gamma}_{i-j;i-l} \quad (37)$$

And likewise:

$$350 \quad \mathbf{C}_{i_{\{in3\}}}^{(29)} \approx \mathbf{\Gamma}_{i-l;i-j} \quad (38)$$

As described in Sect. 3.3.2 on exact cross-covariance statistics, every error covariance from residual cross-covariances has four equivalent formulations for each pair of other datasets which provide the same result in the exact case, but might differ in the approximate formulation. Equation (37) and (38) provide two different approximations of each error covariance matrix from residual cross-covariances based on each pair of other datasets. In the simplified case of scalar statistics, the two different formulations in Eq. (37) and (38) reduce to the same residual cross-variance which was previously formulated by e.g. Crow and van den Berg (2010); Zwieback et al. (2012); Pan et al. (2015).

4.1.2 Differences

Equations (36) to (38) provide three different estimates of an error covariance matrix for each pair of other datasets. Using the relation between residual covariances and cross-covariances from Sect. 3.2.3 and the symmetric properties of residual statistics allow for a comparison of the three estimates:

$$\mathbf{C}_{i_{\{in3\}}}^{(37)} \approx \mathbf{\Gamma}_{i-j;i-l} \stackrel{(24),(36)}{=} \mathbf{C}_{i_{(36)}} + \frac{1}{2} \mathbf{Y}_{i-j;i-l} \quad (39)$$

$$\mathbf{C}_{i_{\{in3\}}}^{(38)} \approx \mathbf{\Gamma}_{i-l;i-j} \stackrel{(24),(36)}{=} \mathbf{C}_{i_{(36)}} - \frac{1}{2} \mathbf{Y}_{i-j;i-l} \quad (40)$$

The three independent estimates of a error covariance matrix from the same pair of other datasets differ only by their residual asymmetry. Thus, differences between the estimates from Eq. (36) to (38) provide no additional information about symmetric error statistics.

While the estimation from residual covariances remains symmetric by definition, the estimates of error covariances from residual cross-covariances may become asymmetric. This asymmetry can be eliminated using the residual asymmetry matrix which is also equivalent to averaging both formulations of error covariances from residual cross-covariances:

$$\mathbf{C}_{i\{in3\}}^{(36)} \approx \frac{1}{2} [\mathbf{\Gamma}_{i-j} + \mathbf{\Gamma}_{l-i} - \mathbf{\Gamma}_{j-l}] \stackrel{(39)}{=} \mathbf{\Gamma}_{i-j;i-l} - \frac{1}{2} \mathbf{Y}_{i-j;i-l} \stackrel{(40)}{=} \mathbf{\Gamma}_{i-l;i-j} + \frac{1}{2} \mathbf{Y}_{i-j;i-l} \quad (41)$$

All three estimates become equivalent if the residual cross-covariances and thus, error cross-covariances are symmetric ($\rightarrow \mathbf{X}_{i;\tilde{j}} = \frac{1}{2} \mathbf{D}_{i;\tilde{j}} = \mathbf{X}_{\tilde{j};i}, \forall i, j$). This is also the case for scalar statistics, were the equivalence between scalar error variance estimates from residual variances and cross-variances was previously shown by Pan et al. (2015).

4.1.3 Uncertainties of approximation

The independence assumption introduces the following absolute uncertainties $\Delta \mathbf{C}_{\tilde{i}}$ of the three different estimates for each dataset i :

$$\Delta \mathbf{C}_{\tilde{i}} \Big|_{(36)} := \mathbf{C}_{\tilde{i}} \Big|_{\text{true}} - \mathbf{C}_{\tilde{i}} \Big|_{(36)} \stackrel{(26),(36)}{=} \frac{1}{2} [\Delta \mathbf{D}_{\tilde{i};\tilde{j}} + \Delta \mathbf{D}_{\tilde{i};\tilde{k}} - \Delta \mathbf{D}_{\tilde{j};\tilde{k}}] \quad (42)$$

$$\Delta \mathbf{C}_{\tilde{i}} \Big|_{(37)} := \mathbf{C}_{\tilde{i}} \Big|_{\text{true}} - \mathbf{C}_{\tilde{i}} \Big|_{(37)} \stackrel{(29),(37)}{=} \Delta \mathbf{X}_{\tilde{j};\tilde{i}} + \Delta \mathbf{X}_{\tilde{i};\tilde{k}} - \Delta \mathbf{X}_{\tilde{j};\tilde{k}} \quad (43)$$

$$\Delta \mathbf{C}_{\tilde{i}} \Big|_{(38)} := \mathbf{C}_{\tilde{i}} \Big|_{\text{true}} - \mathbf{C}_{\tilde{i}} \Big|_{(38)} \stackrel{(29),(38)}{=} \Delta \mathbf{X}_{\tilde{i};\tilde{j}} + \Delta \mathbf{X}_{\tilde{k};\tilde{i}} - \Delta \mathbf{X}_{\tilde{k};\tilde{j}} \quad (44)$$

where $\Delta \mathbf{D}_{\tilde{i};\tilde{j}}$ and $\Delta \mathbf{X}_{\tilde{i};\tilde{j}}$ are the uncertainties of the estimated error dependencies and cross-covariances, respectively.

The absolute uncertainty of the estimates depends similarly on the (neglected) error cross-covariances or dependencies among the three datasets. While the error dependencies to the two other datasets contribute positively, the dependency between the two others is subtracted. If these dependencies cancel out ($\Delta \mathbf{D}_{\tilde{i};\tilde{j}} + \Delta \mathbf{D}_{\tilde{i};\tilde{k}} = \Delta \mathbf{D}_{\tilde{j};\tilde{k}}$), the estimate of one dataset might be exact even if all three dependencies are non-zero. However, two exact estimates can only be achieved if one (e.g. $\Delta \mathbf{D}_{\tilde{i};\tilde{j}} = 0 \wedge \Delta \mathbf{D}_{\tilde{i};\tilde{k}} = \Delta \mathbf{D}_{\tilde{j};\tilde{k}}$) or all three dependencies are zero. A special case was observed by Todling et al. (2022) who showed that the estimations of background, observation and analysis errors in a variational data assimilation system become exact if the analysis is optimal. In this particular case, no assumptions on dependencies are required because the optimality of the analysis induces vanishing dependencies.

Estimated error covariances might even contain negative values if error dependencies are large compared to the true error covariance of a dataset. If the true error covariances differ significantly among highly correlated datasets, the neglected error dependency between two datasets might become much larger than the smaller error covariance, e.g. $\Delta \mathbf{D}_{\tilde{k};\tilde{i}} - \Delta \mathbf{D}_{\tilde{j};\tilde{k}} \approx 0$, $\frac{1}{2} \Delta \mathbf{D}_{\tilde{i};\tilde{j}} > \mathbf{C}_{\tilde{i}} \Big|_{\text{true}}$. Thus, the estimated error covariance matrices might not be positive definite if the independence assumption among three datasets is not fulfilled. This phenomena was also described and demonstrated by Sjoberg et al. (2021) for scalar

problems. However, the generalization to covariances matrices is expected to increase the occurrence of negative values were correlations between two entries of the state are low, thus relative differences and sampling errors become large.

395 4.2 Approximation for multiple datasets

While independence among all datasets is required to estimate the error covariances of three datasets ($I = 3$), the use of more than three datasets ($I > 3$) enables the additional estimation of some error dependencies or cross-covariances (compare Sect. 2). Although this potential of cross-statistic estimation was previously indicated by Gruber et al. (2016); Vogelzang and Stoffelen (2021) for scalar problems, a generalized formulation exploiting its full potential by minimizing the number of assumptions is yet missing.

As described in Sect. 2, $D_I > 0$ gives the number of error cross-statistics which can potentially be estimated in addition to all error covariances for $I > 3$ datasets. Consequently, for each additional dataset $i > 3$, its cross-statistics to one prior dataset $\text{ref}(i) < i$ is needed to be assumed in order to close the problem. This prior dataset $\text{ref}(i)$ is denoted as "reference dataset" of dataset i . In the following, the approximate estimation of error covariances and cross-statistics (cross-covariances or dependencies) under the "partly independence assumption" $\mathbf{D}_{i;\text{ref}(i)} \approx 0$ is formulated for all additional dataset ($\forall i > 3$). This estimation procedure of error statistics of additional dataset based on their reference datasets is denoted as "sequential estimation" in contrast to the "triangular estimation" from an independent triple of datasets presented in Sect. 4.1.

4.2.1 Error covariance estimates

As in the estimation for $I = 3$ datasets, the error covariances of the first three datasets can be estimated from residual covariances or cross-covariances using Eq. (36), (37) or (38). This triple of the first three datasets which are assumed to be pairwise independent is denoted as "basic triangle".

Based on this, each additional error covariance can directly be calculated w.r.t. its reference dataset $\text{ref}(i) < i$:

$$\mathbf{C}_{i\{inI\}}^{(25)} \approx \mathbf{\Gamma}_{i-\text{ref}(i)} - \mathbf{C}_{\text{ref}(i)} \quad (45)$$

where " \approx " indicates the assumption of independence to the reference dataset, i.e. $\mathbf{X}_{i;\text{ref}(i)} \approx 0$.

415 Similarly, each additional error covariance can be estimated from two residual cross-covariances w.r.t its reference dataset $\text{ref}(i)$ and any other dataset j :

$$\mathbf{C}_{i\{inI\}}^{(31)} \approx \mathbf{\Gamma}_{i-\text{ref}(i);i-j} + \mathbf{\Gamma}_{\text{ref}(i)-i;\text{ref}(i)-j} - \mathbf{C}_{\text{ref}(i)} \quad (46)$$

From the equivalence of residual statistics in Eq. (32) it follows that the two formulations of error covariances in Eq. (45) and Eq. (46), respectively, are equivalent and produce exactly the same estimates even if the underlying assumptions are not perfectly fulfilled.

4.2.2 Error cross-covariance and dependency estimates

Once the error covariances are estimated, the remaining residual covariances can be used to calculate the error dependencies to all other prior datasets $j \neq \text{ref}(i), j < i$:

$$\mathbf{D}_{i;\tilde{j}}^{(27)} \stackrel{(27)}{=} \mathbf{C}_{\tilde{i}} + \mathbf{C}_{\tilde{j}} - \mathbf{\Gamma}_{i-j} \quad (47)$$

425 In contrast to residual covariances, the asymmetric formulation of residual cross-covariances allows for an estimation of remaining error cross-covariances including their asymmetries. The error cross-covariance to each other prior dataset $j \neq \text{ref}(i), j < i$ can be estimated sequentially using again the reference dataset $\text{ref}(i)$:

$$\mathbf{X}_{\tilde{i};\tilde{j}}^{(33)} \stackrel{(33)}{\approx}_{\{inI\}} \mathbf{\Gamma}_{\text{ref}(i)-i;\text{ref}(i)-j} - \mathbf{C}_{\widetilde{\text{ref}(i)}} + \mathbf{X}_{\widetilde{\text{ref}(i)};\tilde{j}} \quad (48)$$

Based on this, the symmetric error dependencies can be estimated from its definition in Eq. (13). The equivalence between 430 the formulations of error dependencies from residual covariances and cross-covariances was shown in Eq. (34).

Note that the error cross-covariances $\mathbf{X}_{\tilde{j};\tilde{i}}$ and dependencies $\mathbf{D}_{\tilde{j};\tilde{i}}$ of each subsequent datasets $j > i$ to dataset j result directly from their symmetric properties in Eq. (10) and Eq. (14), respectively.

4.2.3 Uncertainties in approximation

The absolute uncertainty $\Delta \mathbf{C}_{\tilde{i}}$ of an additional error covariance estimate of any dataset $3 < i < I$ is formulated recursively 435 w.r.t. its reference dataset $\text{ref}(i)$:

$$\Delta \mathbf{C}_{\tilde{i}} \Big|_{(45)} := \mathbf{C}_{\tilde{i}} \Big|_{\text{true}} - \mathbf{C}_{\tilde{i}} \Big|_{(45)} \stackrel{(25),(45)}{=} \Delta \mathbf{D}_{\tilde{i};\widetilde{\text{ref}(i)}} - \Delta \mathbf{C}_{\widetilde{\text{ref}(i)}} \quad (49)$$

$$\Delta \mathbf{C}_{\tilde{i}} \Big|_{(46)} := \mathbf{C}_{\tilde{i}} \Big|_{\text{true}} - \mathbf{C}_{\tilde{i}} \Big|_{(46)} \stackrel{(31),(46)}{=} \Delta \mathbf{D}_{\tilde{i};\widetilde{\text{ref}(i)}} - \Delta \mathbf{C}_{\widetilde{\text{ref}(i)}} \quad (50)$$

The two estimates of error covariances from residual covariances in Eq. (49) and from cross-covariances in Eq. (50) are equivalent, and the uncertainty of the latter is independent of the selection of the third dataset j in the residual cross-covariances 440 (compare Eq. (46)). Thus, absolute uncertainties of estimations from residual covariances and cross-covariances differ only in the uncertainties w.r.t. the basic triangle given in Eq. (42) to (44).

With this, a series of reference datasets $\{m_f\} = m_1, \dots, m_F$, with m_F being the reference of i , and the m_{f-1} reference of m_F and so on, with $m_{f-1} < m_f < i, \forall f$ and $m_1 = j \leq 3$ are defined from the target dataset to the basic triangle. Then, the absolute uncertainty $\Delta \mathbf{C}_{\tilde{i}}$ of each error covariance estimate is:

$$\begin{aligned} 445 \quad \Delta \mathbf{C}_{\tilde{i}} &\stackrel{(49)}{=} \Delta \mathbf{D}_{\tilde{i};\widetilde{m_F}} - \Delta \mathbf{C}_{\widetilde{m_F}} = \Delta \mathbf{D}_{\tilde{i};\widetilde{m_F}} - \Delta \mathbf{D}_{\widetilde{m_F};\widetilde{m_{F-1}}} + \Delta \mathbf{C}_{\widetilde{m_{F-1}}} = \dots \\ &\stackrel{(42)}{=} \Delta \mathbf{D}_{\tilde{i};\widetilde{m_F}} + \sum_{f=F-1}^1 \left[(-1)^{F-f} \cdot \Delta \mathbf{D}_{\widetilde{m_{f+1}};\widetilde{m_f}} \right] + (-1)^F \cdot \frac{1}{2} \left[\Delta \mathbf{D}_{\tilde{j};\tilde{k}} + \Delta \mathbf{D}_{\tilde{j};\tilde{l}} - \Delta \mathbf{D}_{\tilde{k};\tilde{l}} \right] \end{aligned} \quad (51)$$

Were $k, l \leq 3$ are the other two datasets in the basic triangle.

According to Eq. (51), uncertainties in the estimations of additional error covariances result from the partly independence assumption of the additional datasets in the series of reference datasets and the independence assumption in the basic triangle.

450 Due to the changing sign between the intermediate dependencies as well as within the basic triangle, the individual uncertainties may cancel out. Thus, absolute uncertainties do not necessarily increase with more intermediate reference datasets.

Although Eq. (47) is exact, the dependency estimate of each additional pair of datasets $(i; j)$ is influenced by uncertainties in the estimations of the related error covariances:

$$455 \quad \Delta \mathbf{D}_{\tilde{i}; \tilde{j}} := \mathbf{D}_{\tilde{i}; \tilde{j}} \Big|_{\text{true}} - \mathbf{D}_{\tilde{i}; \tilde{j}} \Big|_{(47)} \stackrel{(27), (47)}{=} \Delta \mathbf{C}_{\tilde{i}} + \Delta \mathbf{C}_{\tilde{j}} \quad (52)$$

were the uncertainties of the two error covariances are given in Eq. (51).

And the absolute uncertainties of estimates of additional error cross-covariances based on residual cross-covariances can be determined recursively using Eq. (51):

$$\Delta \mathbf{X}_{\tilde{i}; \tilde{j}} := \mathbf{X}_{\tilde{i}; \tilde{j}} \Big|_{\text{true}} - \mathbf{X}_{\tilde{i}; \tilde{j}} \Big|_{(48)} \stackrel{(33), (48)}{=} \Delta \mathbf{X}_{\widetilde{\text{ref}(i); \tilde{j}}} + \Delta \mathbf{X}_{\widetilde{\tilde{i}; \text{ref}(i)}} - \Delta \mathbf{C}_{\widetilde{\text{ref}(i)}} \quad (53)$$

460 In contrast to error covariances, the uncertainties of error cross-covariances sum up in the two series of reference datasets. However, this sum is subtracted by the two sums of uncertainties in error covariances of these datasets, whose elements may cancel partly (not shown).

4.2.4 Comparison to approximation from three datasets

It can be shown that the sequential formulation of an error covariance from its reference dataset is consistent with the triangular formulation from three independent datasets in Sect. 4.1 in the basic triangle. Given the triangular estimate of one error covariance $\mathbf{C}_{\tilde{i}}|_{\triangleleft}$ from Eq. (36), the error covariances $\mathbf{C}_{\tilde{j}}|_{\triangleleft}$ of the other two datasets in the basic triangle are equal to their sequential formulation $\mathbf{C}_{\tilde{j}}|_{\vdash}$ from Eq. (45) with reference dataset $\text{ref}(j) = i$:

$$\mathbf{C}_{\tilde{j}}|_{\vdash} \stackrel{(45)}{\approx}_{\{inI\}} \mathbf{\Gamma}_{i-j} - \mathbf{C}_{\tilde{i}}|_{\triangleleft} \stackrel{(36)_i}{\approx}_{\{in3\}} \mathbf{\Gamma}_{j-i} - \frac{1}{2} [\mathbf{\Gamma}_{i-j} + \mathbf{\Gamma}_{k-i} - \mathbf{\Gamma}_{j-k}] = \frac{1}{2} [\mathbf{\Gamma}_{i-j} + \mathbf{\Gamma}_{j-k} - \mathbf{\Gamma}_{i-k}] \stackrel{(36)_j}{=} \mathbf{C}_{\tilde{j}}|_{\triangleleft} \quad (54)$$

Thus, only one error covariance needs to be calculated with Eq. (36) while all other can be estimated from Eq. (45). Note that

470 although even if only $\mathbf{C}_{\tilde{i}}$ is calculated from the fully independent formulation in the basic triangle, the independence assumption among all three pairs of datasets in the basic triangle remains.

Instead of using the sequential estimation for additional datasets $i > 3$, the error covariances could also be estimated by defining another independent triangle $(i; j; k)$, with $k = \text{ref}(j)$, $j = \text{ref}(i)$. Because the definition of another independent triangle requires an additional independence assumption between $(i; k)$ (i.e. $\mathbf{D}_{\tilde{i}; \tilde{k}} = 0$), this triangular estimate $\mathbf{C}_{\tilde{i}}|_{\triangleleft}$ from Eq. (36) differs from the sequential estimate $\mathbf{C}_{\tilde{i}}|_{\vdash}$ from Eq. (45) using its reference dataset ($\mathbf{C}_{\tilde{j}} \rightarrow \mathbf{C}_{\tilde{i}}$), were their absolute errors compare as follows:

$$\left| \Delta \mathbf{C}_{\tilde{i}}|_{\vdash} \right| - \left| \Delta \mathbf{C}_{\tilde{i}}|_{\triangleleft} \right| \stackrel{(42), (49)}{=} \left| \Delta \mathbf{D}_{\tilde{i}; \tilde{j}} - \Delta \mathbf{C}_{\tilde{j}} \right| - \frac{1}{2} \left| \Delta \mathbf{D}_{\tilde{i}; \tilde{j}} + \Delta \mathbf{D}_{\tilde{i}; \tilde{k}} - \Delta \mathbf{D}_{\tilde{j}; \tilde{k}} \right| \quad (55)$$

The sequential estimation of error covariances of an error covariance becomes favourable if the estimation of the error co-
 480 variance of its reference dataset is of similar accuracy as the uncertainty in their dependent assumption $(\Delta \mathbf{C}_{\tilde{j}} \rightarrow \Delta \mathbf{D}_{\tilde{i};\tilde{j}})$.
 And the triangular estimation becomes favourable if the accuracy of the additional independence assumption is of the order
 of the difference between the uncertainties of other two error dependencies $(\Delta \mathbf{D}_{\tilde{i};\tilde{k}} \rightarrow \Delta \mathbf{D}_{\tilde{i};\tilde{j}} - \Delta \mathbf{D}_{\tilde{j};\tilde{k}})$; i.e. if the additional
 independence assumption is of similar accuracy as the other two dependent assumptions.

Note that the absolute uncertainties presented here only account for uncertainties due to the underlying assumptions on error
 485 cross-statistics and not due to imperfect residual statistics occurring e.g. from finite sampling. A discussion of those effects for
 scalar problems can be found in Sjöberg et al. (2021).

5 Experiments

This section illustrates the capabilities to estimate full error covariance matrices of all datasets and some error dependencies.
 Four collocated datasets ($I = 4$) are generated synthetically on a 1D domain with 25 grid-points. Each dataset consists of 20.000
 490 realizations at each grid-point which are randomly sampled around the true value of 5.0. The spatial variation of prescribed
 error variances and spatial error correlations differ for each dataset. Datasets (1;2;3) span the basic triangle and dataset 1 is
 the reference of dataset 4 ($\text{ref}(4) = 1$). This allows the estimation of all four error covariances of each dataset and two error
 dependencies between the datasets (2;4) and (3;4) (compare Sect. 2). All other error dependencies need to be assumed and
 are set to zero for this experiment (in accordance to the formulation in Sect. 4). Note that the change of error dependencies
 495 between the different experiments affects their true statistics. The experiments presented in this section are based on the
 symmetric estimations from residual covariances derived in Sect. 4 which are summarized in Algorithm A1. Similar results
 would be obtained from estimations from cross-covariances given in Algorithm A2, but this short illustration is restricted to a
 general demonstration using symmetric statistics only.

The plots are structured as follows: Each subplot combines two covariance matrices; one shown in the upper-left part and
 500 the other in the lower-right part. Because all matrices involved are symmetric, it is sufficient to show only one half of each
 matrix. The two matrices are separated by a thick gray diagonal bar and shifted off-diagonal so that diagonal variances are right
 above/below the gray bar, respectively. Statistics that might become negative are shown as absolute quantities in order to show
 them with the same color-code. In each row, the upper-left parts are matrices which are usually unknown in real applications
 (as they require the knowledge of the truth) and the lower-right parts are known/estimated matrices. The first row contains the
 505 error dependencies and residual covariances of each dataset pair. Here, gray asterisks in the upper-left subplot indicate that
 these error dependency matrices are assumed to be zero in the estimation. The second row contains the true and estimated error
 covariances and dependencies. The third row gives the absolute difference between the true and estimates matrices.

5.1 Uncertainties in additional dependencies

The experiment shown in Fig. 2a contains only true error dependencies between datasets (2;4) and (3;4). This is consistent
 510 to the selected estimation setup in which (1;2;3) build the basic triangle and dataset 4 is sequentially estimated w.r.t. dataset

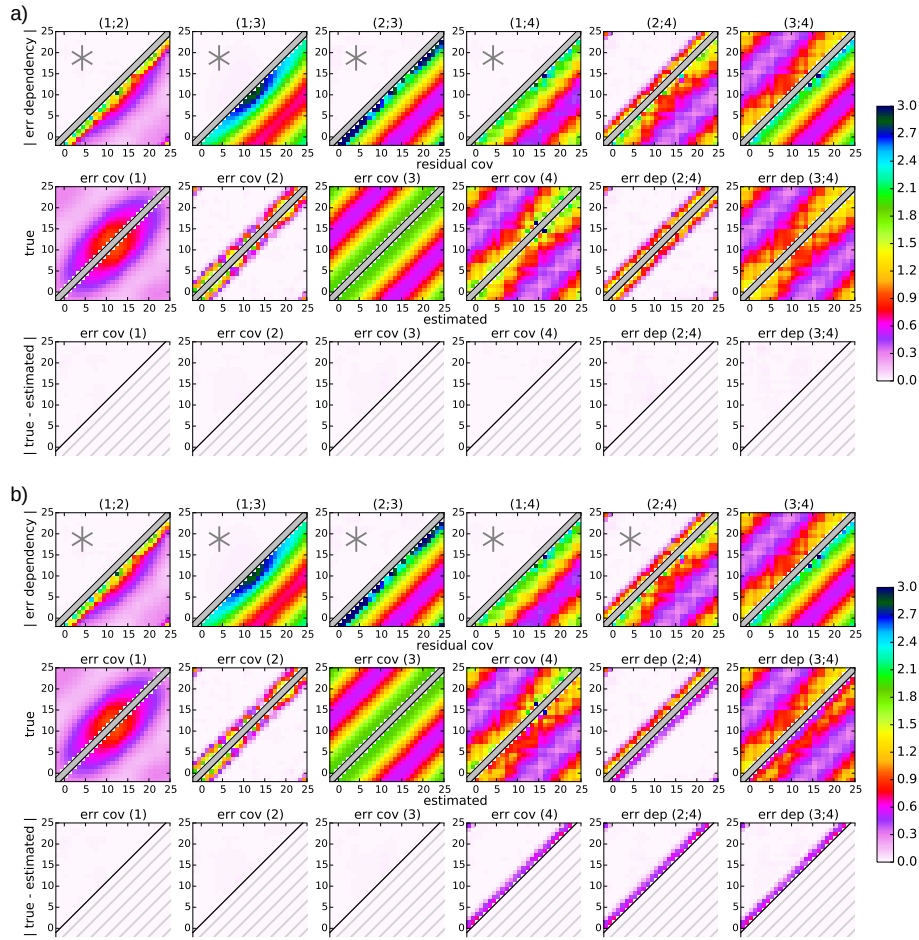


Figure 2. Covariance matrices for 4 datasets ($I = 4$) with true dependencies of datasets (2;4) and (3;4). Datasets (1;2;3) build the basic triangle. Dataset 4 is estimated (a) from its reference dataset 1 ("sequential estimation") and (b) from an additional independent triangle (1;2;4) ("triangular estimation").

1. Consequently, all error covariances and the two remaining dependencies are estimated accurately in accordance to Eq. (42), (51), and (52). The estimation method is able to reproduce true error covariance matrices of all datasets and error dependency matrices between some datasets independent of the complexity of the statistics if the assumptions are sufficiently fulfilled. For comparison, the error covariance matrix of dataset 4 is estimated from an additional independent triangle (1;2;4) in Fig. 2b.

515 The triangular estimation requires the additional independence assumption between datasets (2;4) which is not fulfilled in this experiment. The positive true error dependency has equivalent impact on the estimated error covariance of dataset 4 and its dependencies to datasets 2 and 3. All three matrices are underestimated w.r.t. the true statistics by the half of the neglected error dependency in accordance to Eq. (42) and (52) applied to the triangle (1;2;4).

5.2 Uncertainties in basic triangle

520 The effects of neglected dependencies in the basic triangle is illustrated in Fig. 3 were a true positive error dependency appears between datasets (2;3). In accordance to Eq. (42), (51), and (52), the neglected dependency in the basic triangle affects all estimated statistics. While the error covariance of dataset 1 is overestimated, all other statistics are underestimated. For the sequential estimation in Fig. 3a, uncertainties in the estimated error dependencies are equal to the neglected error dependency and uncertainties in error covariances are halved (compare Eq. (51)). For the triangular estimation of dataset 4 in Fig. 3b, the
525 effects of the two neglected dependencies between (2;3) and between (2;4) are combined. As shown in Fig. 2b, the error covariance of dataset 4 is underestimated by half the neglected dependency between (2;4). The uncertainties two estimated error dependencies (2;4) and (3;4) are the sum of the uncertainties of the error covariances of the two datasets involved in accordance to Eq. (52).

In this setup, the sequential estimation of the additional dataset (here 4) from its reference dataset (here 1) is more accurate
530 because the neglected dependency in the basic triangle (here (2;3)) is small compared to the neglected additional dependency in the triangular estimation (here (2;4)), which is in accordance to Eq. (55). This changes in Fig. 4, where the neglected dependency (here (2;3)) in the basic triangle is larger than the one additional one in the triangular estimation (here (2;4)). Because the sequential estimation is more sensitive to uncertainties in the basic triangle (Fig. 4a), the triangular estimation (Fig. 4b) becomes more accurate. This holds for the error covariance estimate of dataset 4 as well as the two estimated error
535 dependencies (2;4) and (3;4). This particular setup also demonstrates that uncertainties due to neglected dependencies can become larger than the actual true statistics (here e.g. in the error dependency of datasets (2;4) for both estimation methods) which creates negative values in the estimate.

Note that the choice of the estimation method affects only the uncertainty of subsequent estimates which are directly or indirectly referring to the uncertain assumption. In this case, the estimations in the basic triangle are not affected by the
540 estimation method.

6 Conceptual summary

This section provides a summary of the statistical error estimation method proposed in this study focusing on its technical application. Section 6.1 summarises the general assumptions and requirements including an exemplary visualisation, and Sect. 6.2 formulates rules for an optimal setup of datasets w.r.t. imperfect assumptions. An algorithmic summary the calculation of error
545 statistics from residual covariances and cross-covariances, respectively, is given in Appendix A.

6.1 Minimal conditions

For error statistics that need to be assumed, their specific formulation may have different forms. The easiest and most common assumption is to set their error correlations and thus the error cross-covariances and dependencies to zero. This assumption used in Sect. 4.1 and 4.2 is equivalent to the 3CH and TC methods. However, any non-zero error statistics can be defined and

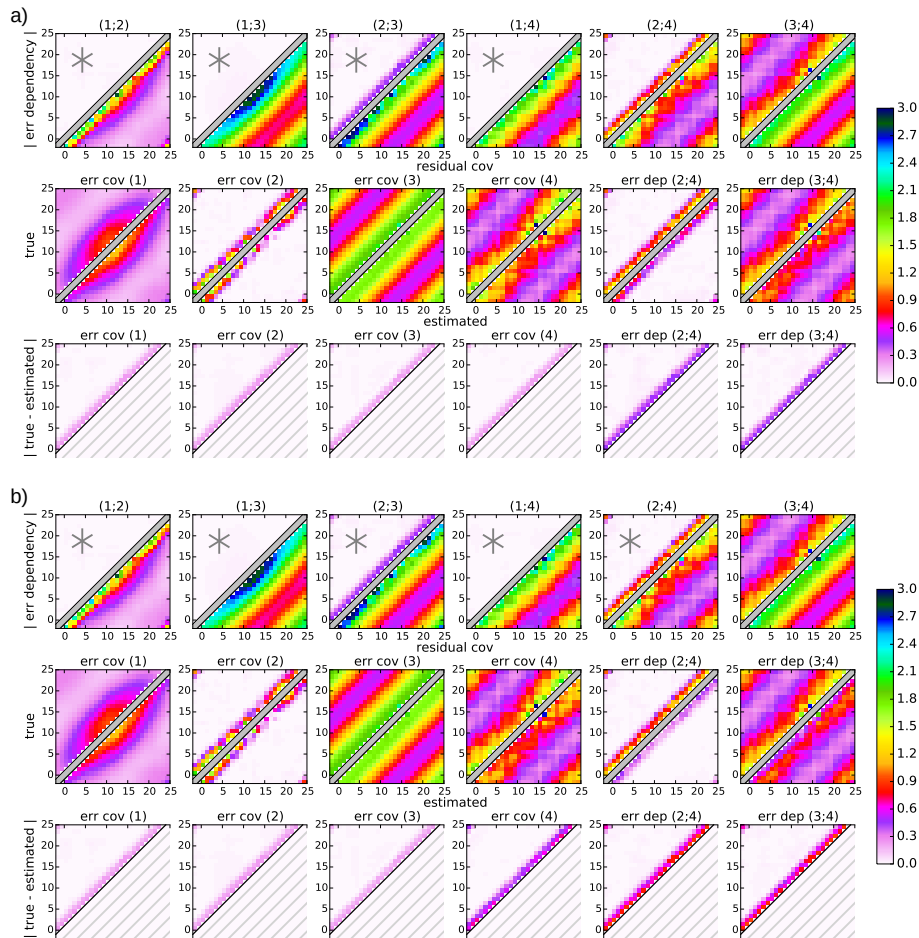


Figure 3. Covariance matrices for 4 datasets ($I = 4$) with true dependencies of datasets (2;3), (2;4) and (3;4). As in Fig. 2, but with an neglected dependency in the basic triangle between datasets (2;3).

550 used in the general form which is summarized in Appendix. A. This also includes assuming error statistics as function of other error statistics including the ones estimated during the calculation. The only restriction is that all assumed error statistics must be fully determined by other error statistics or predefined values without introducing additional degrees of freedom.

In the common case were all error covariances and some error dependencies (or cross-covariances) are estimated, there are two requirements for the setup of datasets: (i) all three error dependencies among one triple of datasets are needed (this triple of independent datasets is called "basic triangle"), and (ii) at least one error dependency of each additional dataset to any prior
555 datasets is needed (this prior dataset is called "reference dataset" of the referring additional datasets).

These two requirements are a logical summary of the mathematical derivations in Sect. 3 and 4 and are valid for all number of datasets $I \leq 3$. They provide the necessary conditions for the existence of a solution under the given assumptions (compare Sect. 4.1.1, 4.2.1, and 4.2.2). Optimality and uniqueness of this solution w.r.t. different formulations and setups are achieved

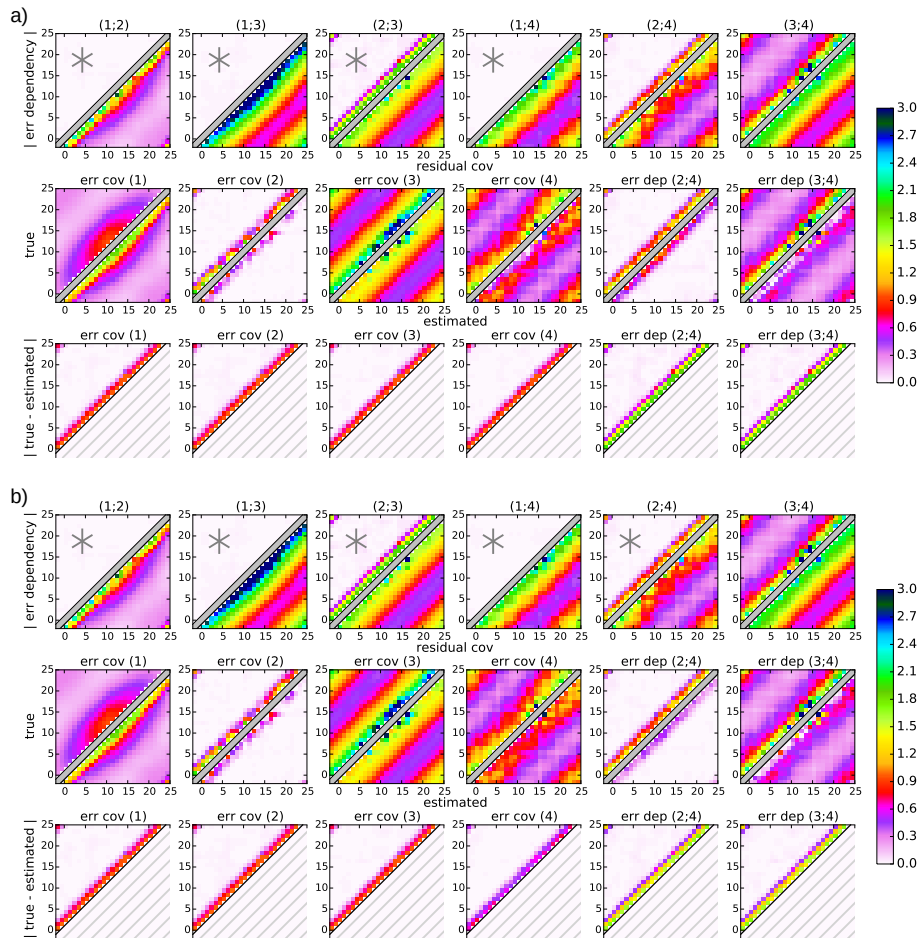


Figure 4. Covariance matrices for 4 datasets ($I = 4$) with true dependencies of datasets (2;3), (2;4) and (3;4). As in Fig. 2, but with an increased dependency in the basic triangle between datasets (2;3).

560 when - and exactly when - the required assumptions are accurate (i.e. vanishing uncertainties of assumed error statistics in Sect. 4.1.2, 4.1.3, 4.2.3, and 4.2.4).

Previously, Vogelzang and Stoffelen (2021) observed that some setups for four and five datasets do not produce a solution for the problem, but without discussing the general requirements. The limited solveability was also found by Gruber et al. (2016) for four datasets, who came up with an unnecessarily strong requirement that each dataset has to be part of an independent
565 triangle.

An exemplary setup of assumed dependencies for $I = 10$ datasets is visualized in Fig. 5. The dependencies among three datasets (1;2;3) is needed to be assumed ("basic triangle"). Then, one dependency of each additional dataset $i > 3$ to any prior dataset j (with $j < i$) is assumed ("sequential estimation"). In general, there is no further restriction on the selection of

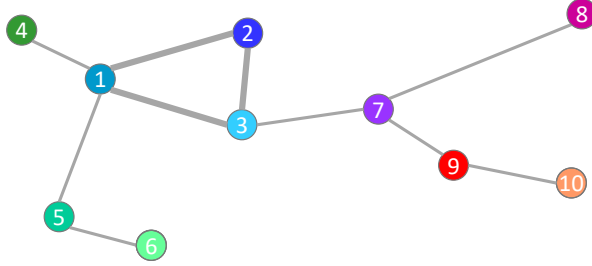


Figure 5. Independence tree: Exemplary visualization of assumed dependencies (gray lines) between 10 datasets (colored dots). The assumed dependencies in the basic triangle (1;2;3) are indicated by thicker lines.

reference datasets in order to close the estimation problem. Note that similar conditions can be derived for cases where also
570 error covariances are given or assumed, which is not part of this paper.

6.2 Optimal setup

In real applications, there might be significant differences in estimated error statistics from different setups as observed e.g. by Vogelzang and Stoffelen (2021) in the scalar case. The relative accuracy of an error covariance estimate is proportional to the ratio between the residual covariance Γ_{i-j} and the absolute uncertainty $\Delta \mathbf{D}_{i;j}^{\sim}$ of the assumed error dependency, which can
575 be interpreted similar to a signal-to-noise ratio. In other words, the larger the residual covariance and the better the absolute estimate of the error dependency to the reference dataset, the more accurate is the estimated error covariance. Because uncertainties in error estimate do not necessarily sum up along a branch of the independence tree (compare Sect. 4.2.3, Eq. (51)), a large residual-to-dependency ratio w.r.t. to the reference is more important than a low number of intermediate reference datasets. In order to achieve optimal estimates, the setup of datasets should be selected according to the expected accuracy of
580 estimated dependencies which minimize the residual-to-dependency ratio for each dataset:

$$\max_j \left(\frac{\Gamma_{i-j}}{\Delta \mathbf{D}_{i;j}^{\sim}} \right) : j \rightarrow \text{ref}(i) \iff \min_j \left(\Delta \rho_{i;j}^{\sim} \right) : j \rightarrow \text{ref}(i) \quad , \quad \forall i \quad (56)$$

The maximal residual-to-dependency ratio is equivalent to the minimal uncertainty in normalized error correlations $\Delta \rho_{i;j}^{\sim} := \frac{\Delta \mathbf{D}_{i;j}^{\sim}}{\sqrt{\mathbf{C}_i \mathbf{C}_j}}$. For example, if the error correlation of one dataset to another is comparably well known, this dataset is best suited as reference dataset. If estimated error dependencies are set to zero, the dataset to which the independence assumption is most
585 certain should be selected as reference dataset. Supposing that distances between datasets indicate their expected degree of independence in the independence tree, the setup visualized in Fig. 5 is not optimal. An example for an improved setup is shown in Fig. 6, which is expected to provide more accurate error estimates.

While uncertainties in the basic triangle only contribute half, they effect the estimations of error statistics of all datasets (compare Sect. 4.1.3, 4.2.3, and 5.2). This has two implications: Firstly, the basic triangle which is defined as the triple of
590 dataset that has the smallest error correlations produces the smallest overall uncertainty w.r.t. all error estimates. Ideally, the

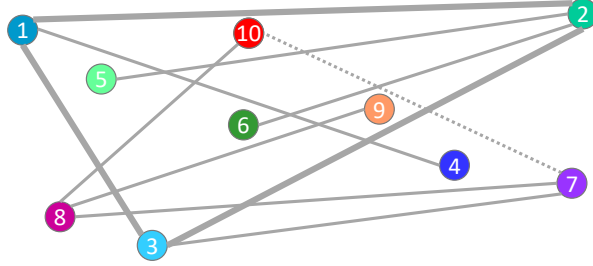


Figure 6. Improved independence tree: As Fig. 5, but with modified setup for more accurate error estimates. Distances between datasets represent the accuracy of assumed dependencies between the error statistics. While locations are the same, the numbers and colors of the datasets has been changed according to the modified setup. An example for an alternative formulation of an additional independent triangle is indicated as dotted line.

basic triangle should be set as a triple of datasets which are highly independent – or at least with reasonably small dependencies among each other.

Secondly, if another independent triangle can be assumed for an additional dataset with similar accuracy as the dependency to its reference dataset, the additional error estimate may be more accurate using the triangular estimation from this additional independent triangle rather than the sequential estimation (compare Sect. 4.2.4 and 5). The additional independent triangle does not need to be connected to the basic triangle and may also have multiple independence branches, thus acting as additional basic triangle. For example, in the setup shown in Fig. 6, the estimation of dataset 7 is sensitive to the dependency(3;7) to its reference dataset, and less sensitive to dependencies in the basic triangle (1;2;3). If the dependency(7;10) could be assumed with higher accuracy than these dependencies, the error covariances of dataset 7 can alternatively be calculated from the independent triangle (7;8;10) and the independence assumption between (3;7) can be dropped. Thus, multiple independence trees can be defined around multiple separated basic triangles.

Furthermore, it is also possible to average the estimated error statistics of a dataset from multiple independent triangles similar to an application of the N-cornered hat method (e.g., Sjöberg et al., 2021) for an arbitrary subset of datasets. This setup builds an overestimated problem which requires the assumption of more error dependencies than the minimal requirements. However, it might be beneficial if multiple independent triangles containing the same dataset can be estimated with similar accuracy. In this case, potential uncertainties in the assumptions are expected to be reduced by the average over similar accurate estimates. Also an extension to weighted averages of different estimations is possible were the weights reflect the expected accuracy of each estimation formulation w.r.t. the others.

7 Conclusions

Despite the generalized matrix-formulation, the main features of the presented approach are (i) its generality defining the flexible setup for any number of datasets according to the specific application, (ii) its optimality w.r.t. a minimal number of

assumptions required, and (iii) its suitability to include expected non-zero dependencies between any pair of datasets. In contrast, the scalar N-CH method averages all estimates of each dataset which is equivalent to assuming that the independence assumption among each dataset triple is fulfilled with the same accuracy. However, this is not the case for most applications to geophysical datasets. For example, Rieckh et al. (2021) applied the N-CH method to multiple atmospheric model and observational datasets and discussed neglected levels of independence between different datasets, which are expected to vary significantly. Pan et al. (2015) tried to account for such variations by clustering the datasets into structural groups; which however requires more assumptions than necessary and makes the result highly sensitive to the selected grouping. In contrast, the method presented here provides an optimal and flexible approach to handle multiple datasets with different levels of expected independence. Depending on the specific application, the estimation may be based on the minimal number of assumptions required or a (weighted) average over any number of estimations with similar expected accuracies.

An important application of the presented method is expected to be numerical weather prediction (NWP) where short-term forecasts from multiple national centers can be used to estimate error statistics required for data assimilation. In contrast to previous statistical methods, potential dependencies among the forecasts, i.e. due to the assimilation of similar observations, can be considered in the error estimation and even explicitly quantified. Future work will show how this statistical approach compares to state-of-the-art background error estimates based on computation-expensive Monte-Carlo- or ensemble-methods. While the presented method ensures symmetry of error covariances, positive definiteness might not be fulfilled in real applications due to inaccurate assumptions or sampling uncertainties.

In comparison to a-posteriori methods which statistically estimate optimal error covariances for data assimilation, an a-priori error estimation of collocated datasets has three main advantages: (i) optimal error statistics are calculated analytically without requiring an iterative minimization including multiple executions of the assimilation, (ii) complete covariance matrices provide spatially-resolved fields of error statistics at each collocated location including spatial- and cross-species correlations, and (iii) error statistics of all datasets are estimated without selecting one dataset as reference. This enables the consideration of more than two datasets in the assimilation. Given sufficiently estimated error statistics, the final analysis w.r.t. to all datasets will be closer to the truth than any analysis between two datasets only. Thus, the rapidly increasing number of geophysical observations and model forecast enables improved analyses through increasingly overlapping datasets, where optimal error statistics can be calculated for example with the method presented here. Especially the possibility to estimate optimal error cross-covariances between datasets provides important information for data assimilation where the violation of the independence assumption remains a major challenge (Tandeo et al., 2020).

However, current data assimilation schemes are not suited for multiple overlapping datasets and cross-errors between datasets are assumed to be negligible. In contrast, the statistical error estimation method presented in this study is explicitly tailored to multiple datasets which cannot be assumed to be independent. Thus, the estimated error covariances are not consistent with assimilation algorithms assuming (two) independent datasets. If the estimated error dependencies among all assimilated datasets are small, the independence assumption may be regarded as sufficiently fulfilled. The error estimation method then provides optimal error covariances for assimilation and information on the accuracy of the independence assumption. Otherwise, generalized assimilation schemes are needed to be developed for a proper use of this additional statistical information in data

assimilation. Although increasing their complexity, such generalized assimilation schemes enable fundamental improvements in terms of an optimal analysis from multiple datasets w.r.t. their error covariances and cross-statistics.

Appendix A: Algorithm

650 The general estimation procedure of error statistics for $I \geq 3$ datasets is summarized in Algorithm A1 and A2. The algorithms require respectively, residual covariances or cross-covariances among all I datasets (calculated from residual statistics) and I assumed error dependencies or cross-covariances. Based on this, the first error covariance matrix in the basic triangle is calculated. Then, error statistics of the remaining datasets are calculated sequentially in an iterative procedure; introducing a new dataset i with given residual statistics (covariances or cross-covariances) to dataset $\text{ref}(i)$ for each $i \in [2, I]$ with $\text{ref}(i) < i$.
 655 Note that this is equivalent to estimating the independent estimations of all three datasets in basic triangle and sequentially estimate all additional estimates for datasets $i > 3$ (compare Sect. 4.2.4).

Algorithm A1 is formulated for symmetric statistic matrices, where error covariances $\text{errcov}(i;:::)$ of each dataset i and error dependency matrices $\text{errdep}(i;j;:::)$ between each pair (i,j) are estimated from symmetric residual covariances $\text{rescov}(i-j;:::)$. In Algorithm A2, the error covariance- and cross-covariance matrices $\text{errcross}(i;j;:::)$ of each pair
 660 (i,j) are estimated from residual cross-covariances $\text{rescross}(i-j;i-k;:::)$ between $(i-j;i-k)$. Here, the third dataset k in the residual cross-covariances can be freely selected and does not affect the accuracy of the estimates (compare Sect. 4.2.3). Each operation applies element-wise to each matrix-element indicated by the last two indices $(:::)$, where matrices may contain different locations of the same quantity as well as different fields for multiple quantities of any dimension (=multivariate covariances). Transposed matrices w.r.t. the two location indices are indicated by $[]^T$.

665 The equations relate to the general exact formulations which requires some error dependencies or cross-covariances to be given (compare Sect. 3). The explicit calculation of the error cross-statistics (dependencies or cross-covariances) is not needed if only error covariances are of interest. In theory, both algorithms provide the same error estimations (compare Sect.3.2.3). The decision to estimate error statistics from residual covariances (Algorithm A1) or cross-covariances (Algorithm A2) depends on the availability of residual statistics and the need for asymmetric error cross-covariances, which can only be estimated with
 670 Algorithm A2 (compare Sect. 3.3.1).

Algorithm A1 Iterative calculation of error covariances and dependencies for I datasets from residual covariances.

Require: $\text{innocov}(i-\text{ref}(i);:::) \forall i \in [2, I], \text{innocov}(1-3;:::)$

Require: $\text{errdep}(i;\text{ref}(i);:::) \forall i \in [2, I], \text{errdep}(1;3;:::)$

$$\text{errcov}(1;:::) \leftarrow 0.5 \cdot \left[\text{innocov}(2-1;:::) + \text{innocov}(1-3;:::) - \text{innocov}(3-2;:::) + \text{errdep}(2;1;:::) + \text{errdep}(1;3;:::) - \text{errdep}(3;2;:::) \right] \quad \triangleright \sim \text{Eq. (26)}$$

for $i = 2, I$ do

$$\text{errcov}(i;:::) \leftarrow \text{innocov}(i-\text{ref}(i);:::) + \text{errdep}(i;\text{ref}(i);:::) - \text{errcov}(\text{ref}(i);:::) \quad \triangleright \sim \text{Eq. (25)}$$

for $j = 1, i-1$ do

if $j \neq \text{ref}(i)$ then

$$\text{errdep}(i;j;:::) \leftarrow \text{errcov}(i;:::) + \text{errcov}(j;:::) - \text{innocov}(i-j;:::) \quad \triangleright \sim \text{Eq. (27)}$$

end if

$$\text{errdep}(j;i;:::) \leftarrow \text{errdep}(i;j;:::) \quad \triangleright \sim \text{Eq. (14)}$$

end for

end for

Algorithm A2 Iterative calculation of error covariances and cross-covariances for I datasets from residual cross-covariances.

Require: $\text{innocross}(i\text{-ref}(i); i\text{-}j; :::), \text{innocross}(\text{ref}(i)\text{-}i; \text{ref}(i)\text{-}j; :::) \forall i \in [2, I], j \neq \text{ref}(i), j \neq i, \text{innocross}(1\text{-}2; 1\text{-}3; :::)$
Require: $\text{errcross}(i; \text{ref}(i); :::) \forall i \in [2, I], \text{errcross}(1; 3; :::)$

```

for  $i = 2, I$  do
   $\text{errcross}(\text{ref}(i); i; :::) \leftarrow \text{errcross}(i; \text{ref}(i); :::)^T$   $\triangleright \sim \text{Eq. (10)}$ 
end for
 $\text{errcov}(1; :::) \leftarrow \text{innocross}(1\text{-}2; 1\text{-}3; :::) + \text{errcross}(1; 3; :::) + \text{errcross}(2; 1; :::) - \text{errcross}(2; 3; :::)$   $\triangleright \sim \text{Eq. (29)}$ 
for  $i = 2, I$  do
   $\text{errcov}(i; :::) \leftarrow \text{innocross}(i\text{-ref}(i); i\text{-}j; :::) + \text{innocross}(\text{ref}(i)\text{-}i; \text{ref}(i)\text{-}j; :::) - \text{errcov}(\text{ref}(i); :::)$ 
     $+ \text{errcross}(i; \text{ref}(i); :::) + \text{errcross}(\text{ref}(i); i; :::)$   $\triangleright \sim \text{Eq. (31)}$ 
  for  $j = 1, i - 1$  do
    if  $j \neq \text{ref}(i)$  then
       $\text{errcross}(i; j; :::) \leftarrow \text{innocross}(\text{ref}(i)\text{-}i; \text{ref}(i)\text{-}j; :::) - \text{errcov}(\text{ref}(i); :::)$ 
         $+ \text{errcross}(\text{ref}(i); j; :::) + \text{errcross}(i; \text{ref}(i); :::)$   $\triangleright \sim \text{Eq. (33)}$ 
       $\text{errcross}(j; i; :::) \leftarrow \text{errcross}(i; j; :::)^T$   $\triangleright \sim \text{Eq. (10)}$ 
    end if
  end for
end for
end for

```

Author contributions. AV developed the approach, derived the theory, performed the experiments, and wrote the manuscript. RM supervised the work, and revised the manuscript.

Competing interests. The authors declare that they have no conflict of interest.

Acknowledgements. The authors thank Ricardo Todling and one anonymous reviewer for their thoughtful and valuable feedback on the manuscript.

References

- Anthes, R. and Rieckh, T.: Estimating observation and model error variances using multiple data sets, *Atmospheric Measurement Techniques*, 11, 4239–4260, <https://doi.org/10.5194/amt-11-4239-2018>, 2018.
- Crow, W. T. and van den Berg, M. J.: An improved approach for estimating observation and model error parameters in soil moisture data
680 assimilation, *Water Resources Research*, 46, <https://doi.org/https://doi.org/10.1029/2010WR009402>, 2010.
- Crow, W. T. and Yilmaz, M. T.: The Auto-Tuned Land Data Assimilation System (ATLAS), *Water Resources Research*, 50, 371–385, <https://doi.org/https://doi.org/10.1002/2013WR014550>, 2014.
- Daley, R.: The Effect of Serially Correlated Observation and Model Error on Atmospheric Data Assimilation, *Monthly Weather Review*, 120, 164 – 177, [https://doi.org/10.1175/1520-0493\(1992\)120<0164:TEOSCO>2.0.CO;2](https://doi.org/10.1175/1520-0493(1992)120<0164:TEOSCO>2.0.CO;2), 1992a.
- 685 Daley, R.: The Lagged Innovation Covariance: A Performance Diagnostic for Atmospheric Data Assimilation, *Monthly Weather Review*, 120, 178 – 196, [https://doi.org/10.1175/1520-0493\(1992\)120<0178:TLICAP>2.0.CO;2](https://doi.org/10.1175/1520-0493(1992)120<0178:TLICAP>2.0.CO;2), 1992b.
- Desroziers, G., Berre, L., Chapnik, B., and Poli, P.: Diagnosis of observation, background and analysis-error statistics in observation space, *Quarterly Journal of the Royal Meteorological Society*, 131, 3385–3396, <https://doi.org/https://doi.org/10.1256/qj.05.108>, 2005.
- Gray, J. and Allan, D.: A Method for Estimating the Frequency Stability of an Individual Oscillator, in: 28th Annual Symposium on Frequency
690 Control, pp. 243–246, <https://doi.org/10.1109/FREQ.1974.200027>, 1974.
- Grubbs, F. E.: On Estimating Precision of Measuring Instruments and Product Variability, *Journal of the American Statistical Association*, 43, 243–264, <https://doi.org/10.1080/01621459.1948.10483261>, 1948.
- Gruber, A., Su, C.-H., Crow, W. T., Zwieback, S., Dorigo, W. A., and Wagner, W.: Estimating error cross-correlations in soil moisture data sets using extended collocation analysis, *Journal of Geophysical Research: Atmospheres*, 121, 1208–1219,
695 <https://doi.org/https://doi.org/10.1002/2015JD024027>, 2016.
- Kren, A. C. and Anthes, R. A.: Estimating Error Variances of a Microwave Sensor and Dropsondes aboard the Global Hawk in Hurricanes Using the Three-Cornered Hat Method, *Journal of Atmospheric and Oceanic Technology*, 38, 197 – 208, <https://doi.org/10.1175/JTECH-D-20-0044.1>, 2021.
- Li, H., Kalnay, E., and Miyoshi, T.: Simultaneous estimation of covariance inflation and observation errors within an ensemble Kalman filter,
700 *Quarterly Journal of the Royal Meteorological Society*, 135, 523–533, <https://doi.org/https://doi.org/10.1002/qj.371>, 2009.
- McColl, K. A., Vogelzang, J., Konings, A. G., Entekhabi, D., Piles, M., and Stoffelen, A.: Extended triple collocation: Estimating errors and correlation coefficients with respect to an unknown target, *Geophysical Research Letters*, 41, 6229–6236, <https://doi.org/https://doi.org/10.1002/2014GL061322>, 2014.
- Ménard, R.: Error covariance estimation methods based on analysis residuals: theoretical foundation and convergence properties derived from simplified observation networks, *Quarterly Journal of the Royal Meteorological Society*, 142, 257–273, <https://doi.org/https://doi.org/10.1002/qj.2650>, 2016.
- Ménard, R. and Deshaies-Jacques, M.: Evaluation of Analysis by Cross-Validation. Part I: Using Verification Metrics, *Atmosphere*, 9, <https://doi.org/10.3390/atmos9030086>, 2018.
- Mitchell, H. L. and Houtekamer, P. L.: An Adaptive Ensemble Kalman Filter, *Monthly Weather Review*, 128, 416 – 433,
710 [https://doi.org/10.1175/1520-0493\(2000\)128<0416:AAEKF>2.0.CO;2](https://doi.org/10.1175/1520-0493(2000)128<0416:AAEKF>2.0.CO;2), 2000.

- Nielsen, J. K., Gleisner, H., Syndergaard, S., and Lauritsen, K. B.: Estimation of refractivity uncertainties and vertical error correlations in collocated radio occultations, radiosondes and model forecasts, *Atmospheric Measurement Techniques Discussions*, 2022, 1–28, <https://doi.org/10.5194/amt-2022-121>, 2022.
- 715 Pan, M., Fisher, C. K., Chaney, N. W., Zhan, W., Crow, W. T., Aires, F., Entekhabi, D., and Wood, E. F.: Triple collocation: Beyond three estimates and separation of structural/non-structural errors, *Remote Sensing of Environment*, 171, 299–310, <https://doi.org/https://doi.org/10.1016/j.rse.2015.10.028>, 2015.
- Rieckh, T., Sjöberg, J. P., and Anthes, R. A.: The Three-Cornered Hat Method for Estimating Error Variances of Three or More Atmospheric Datasets. Part II: Evaluating Radio Occultation and Radiosonde Observations, Global Model Forecasts, and Reanalyses, *Journal of Atmospheric and Oceanic Technology*, 38, 1777 – 1796, <https://doi.org/10.1175/JTECH-D-20-0209.1>, 2021.
- 720 Scipal, K., Holmes, T., de Jeu, R., Naeimi, V., and Wagner, W.: A possible solution for the problem of estimating the error structure of global soil moisture data sets, *Geophysical Research Letters*, 35, <https://doi.org/https://doi.org/10.1029/2008GL035599>, 2008.
- Sjöberg, J. P., Anthes, R. A., and Rieckh, T.: The Three-Cornered Hat Method for Estimating Error Variances of Three or More Atmospheric Datasets. Part I: Overview and Evaluation, *Journal of Atmospheric and Oceanic Technology*, 38, 555 – 572, <https://doi.org/10.1175/JTECH-D-19-0217.1>, 2021.
- 725 Stoffelen, A.: Toward the true near-surface wind speed: Error modeling and calibration using triple collocation, *Journal of Geophysical Research: Oceans*, 103, 7755–7766, <https://doi.org/https://doi.org/10.1029/97JC03180>, 1998.
- Su, C.-H., Ryu, D., Crow, W. T., and Western, A. W.: Beyond triple collocation: Applications to soil moisture monitoring, *Journal of Geophysical Research: Atmospheres*, 119, 6419–6439, <https://doi.org/https://doi.org/10.1002/2013JD021043>, 2014.
- Tandeo, P., Ailliot, P., Bocquet, M., Carrassi, A., Miyoshi, T., Pulido, M., and Zhen, Y.: A Review of Innovation-Based Methods to Jointly
- 730 Estimate Model and Observation Error Covariance Matrices in Ensemble Data Assimilation, *Monthly Weather Review*, 148, 3973 – 3994, <https://doi.org/10.1175/MWR-D-19-0240.1>, 2020.
- Tangborn, A., Ménard, R., and Ortland, D.: Bias correction and random error characterization for the assimilation of high-resolution Doppler imager line-of-sight velocity measurements, *Journal of Geophysical Research: Atmospheres*, 107, ACL 5–1–ACL 5–15, <https://doi.org/https://doi.org/10.1029/2001JD000397>, 2002.
- 735 Todling, R., Semane, N., Anthes, R., and Healy, S.: The relationship between two methods for estimating uncertainties in data assimilation, *Quarterly Journal of the Royal Meteorological Society*, <https://doi.org/https://doi.org/10.1002/qj.4343>, 2022.
- Vogelzang, J. and Stoffelen, A.: Quadruple Collocation Analysis of In-Situ, Scatterometer, and NWP Winds, *Journal of Geophysical Research: Oceans*, 126, e2021JC017 189, <https://doi.org/https://doi.org/10.1029/2021JC017189>, 2021.
- Voshtani, S., Ménard, R., Walker, T. W., and Hakami, A.: Assimilation of GOSAT Methane in the Hemispheric CMAQ; Part I: Design of the
- 740 Assimilation System, *Remote Sensing*, 14, <https://doi.org/10.3390/rs14020371>, 2022.
- Xu, X. and Zou, X.: Global 3D Features of Error Variances of GPS Radio Occultation and Radiosonde Observations, *Remote Sensing*, 13, <https://doi.org/10.3390/rs13010001>, 2021.
- Zwieback, S., Scipal, K., Dorigo, W., and Wagner, W.: Structural and statistical properties of the collocation technique for error characterization, *Nonlinear Processes in Geophysics*, 19, 69–80, <https://doi.org/10.5194/npg-19-69-2012>, 2012.