

One major question in data assimilation is to determine the statistics of the errors affecting the data to be assimilated. It is those statistics that define the weights to be given to the data in the assimilation. However, they can never be fully determined without external hypotheses, *i. e.* hypotheses that cannot be objectively validated on the basis of the data alone.

The authors present and discuss an approach that is appropriate for the situations in which a number of sets of collocated data are available. They consider only second-order statistical moments (first-order moments, *i. e.* biases, are also required, but their identification is an independent problem). Covariances and cross-covariances of the differences ('innovations') between those different sets of data are known from the data, and are linearly related to the statistics of the underlying data errors. By appropriate *a priori* specification of a number of those data errors statistics (the external hypotheses), all, or part, of the remaining error statistics are solution of a system of linear matrix equations. That approach, which originated from the so-called *three cornered hat* (3CH) method, has been used in a number of applications, but not much so far in assimilation of geophysical data.

Given  $I$  sets of collocated data, the unknowns (covariances and cross-covariances of data errors) are in number  $U_I = (1/2) I (I+1)$  (Eq. 2 of the paper). Concerning the innovations, their second order moments are not independent, and they are combinations of only  $N_I = (1/2) I (I-1)$  of them (l. 95). This leads to a linear system of  $N_I$  matrix equations with  $U_I$  unknowns (that system is basically expressed, although in what is to me a cursory passing remark, by Eq. 22). The degree of underdeterminacy of the system is  $U_I - N_I = I$ . The view that suffices to choose *a priori*  $I$  of the unknowns to close algebraically the system is correct for  $I = 3$ , but not necessarily for larger values (at least if, as the authors want, no error covariance is specified *a priori*). The purpose of the authors, in addition to stating precisely and discussing the problem, is to determine minimal conditions for its solution (... *what are the minimal and optimal conditions to solve the problem?*, l. 65). They also present numerical results obtained from synthetic data.

The article is instructive, and certainly contains material that is worth publishing. But it needs in my opinion substantial improvement.

1. My main comment is that I have found it very difficult to understand the very logic of the paper (and I am actually still not even sure I have fully understood). A succinct analysis shows that, for  $I > 1$ , system (22) (strictly speaking, a system of  $N_I$  equations which is equivalent to 22) is of rank  $N_I$ , which shows that by appropriately choosing  $I$  of the unknown error covariances and cross-covariances, one can obtain the values of all the other unknowns. My understanding is that the authors show that these  $I$  *a priori* chosen error covariances and cross-covariances cannot be chosen arbitrarily, and that there are constraints in that choice (especially in the case considered by the authors, in which only cross-covariances are to be chosen *a priori*). If it is so, I think it must be stated more explicitly.

2. Subsection 6.1 (*Minimal conditions*) contains what I understand are the authors' main conclusions. That Subsection states two conditions (ll. 527-529) that are presented as the minimal conditions ensuring existence and uniqueness of the solution of system (22) (at least, that is my understanding)

- (i) *all three error dependencies between one triple of datasets are needed (this triple of independent datasets is called "basic triangle")*
- (ii) *at least one error dependency of each additional dataset to any prior datasets is needed*

I is not clear to me whether these two conditions are mathematically exact (if yes, explain more clearly where they are proven in the paper, or give a reference ; if not, say clearly they are only reasonable conjectures).

3. I find that Sections 3 and 4, although they boil down to elementary algebraic manipulations, are intricate and difficult to follow.

*a.* Eq. (22) expresses the basic links between innovation and error statistics (denoted respectively  $\Gamma$  and  $\mathbf{X}$ ). Although algebraically obvious, it is the crux of the method, and should be stressed more strongly as such.

*b.* The derivation of Eq. 23 (ll. 211-213) is strange, since it suggests (l. 211) that one must go through the error statistics  $\mathbf{X}$  to obtain the equation, while the latter expresses necessarily links between the innovation statistics  $\Gamma$ , and can be easily be proved directly.

*c.* Eq. (34) is also strange in that it purports to show the ‘equivalence’ between two expressions for the error dependencies  $\mathbf{D}$ . Those two expressions are basically obtained from Eq. (22), and the reader would think they must necessarily be the same. I presume the authors want to stress that inappropriate choice of the *a priori* chosen error cross-covariances can lead to inconsistencies. But, rather than demonstrating consistency, it would be preferable to show an explicit example of inconsistency. Actually, my understanding is that Eqs (39-40) precisely show an example of inconsistency. If I am mistaken about the significance of Eq. (34), say more explicitly what that significance is.

*d.* The authors, for some unspecified reason, consider only the ‘error-dependencies’, *i.e.* the symmetric part of the error cross-correlations matrices (Eq. 20), and ignore the anti-symmetric part. Why so ?

*e.* It is not clearly said why the number of independent innovation covariances and cross-covariances is equal to  $N_I = (1/2) I (I-1)$  (that is rather simple, but must be said more clearly). The mutual dependence between those quantities is expressed by Eq. 23, the significance of which (in addition to my remark *b* above) should be stressed more strongly.

These are only examples of places that can cause confusion in the mind of a reader who is a newcomer to the approach described in the paper, as elementary as that approach may fundamentally be. I think Sections 3 and 4 could be rewritten in a clearer and more concise way, with more stress on the logic of the approach and on the two fundamental aspects upon which it is based. First, that the observed innovation covariances and cross-covariances are redundant. Second, the basic link between between the innovation and errors covariances and cross-covariances, expressed by Eq. (22) (or any other equivalent equation for that matter).

4. And, for a final (but I think important) comment, any algebraic solution to system (22) will not be acceptable in then present context. It must also define a proper (symmetric non negative) global error covariance matrix (in particular, the estimated error covariance matrices  $C_i$  of the various individual datasets, in addition to being symmetric, must be non-negative). The authors hardly mention this point. Do the conditions (i-ii) stated in subsection (6.1) lead to a proper global covariance matrix ? Since system (22) expresses necessary conditions between error and innovation variances and covariances, I presume that if the *a priori* specified variances and cross-covariances are compatible with a globally symmetric non-negative matrix that is itself compatible with the  $\Gamma_{i,j;k,l}$  ‘s (Eq. 22), the estimated

variances and cross-covariances will also be. I do not ask the authors to necessarily give a full answer to that question, but it should be clearly mentioned and at least briefly discussed. In particular, if the authors do not have a full answer to that question, it should clearly stated as remaining an open question.

It may that the response to some of the questions I raise above is available in the literature, in particular in the literature the authors mention. If so, please give precise references.

I would have a number of other comments, bearing on both scientific and editing aspects of the paper, but they are of lesser importance, and I will wait for a possible revised version for mentioning them.