# The N-Cornered Hat method for Estimating Error Variances between Multiple Data Sets:

## Theoretical considerations and comparisons with the two-cornered hat method

Jeremiah P Sjoberg*, Richard A Anthes, and Therese Rieckh

UCAR/COSMIC

## Abstract

Variations of the three-cornered hat (3CH) method for estimating random error variances associated with three or more data sets have been reported in the literature for studying geophysical data sets, including sea-surface winds, sea-surface temperatures, precipitation, leaf area index and soil moisture. Anthes and Rieckh (2018) and Rieckh and Anthes (2018) used the 3CH method to estimate the error variances of multiple atmospheric sounding data sets. However the methods used often contain subtle variations and make different assumptions. Here we derive the full 3CH equations that relate the error variance of three or more observations to the variance of differences between the data sets and the error covariances among the different data sets.

## Three-cornered hat relations

Consider the data sets $X_1$, $X_2$, and $X_3$ with individual elements $i = 1, 2, \ldots$. We assume that these may be cast as

$$X_{1,i} = T_i + b_{X1} + \varepsilon_{X1,i} \tag{1a}$$
$$X_{2,i} = T_i + b_{X2} + \varepsilon_{X2,i} \tag{1b}$$
$$X_{3,i} = T_i + b_{X3} + \varepsilon_{X3,i}, \tag{1c}$$

where $T_i$ is a set of reference values; the $b$ terms are mean differences of the individual data sets from the reference data set; and the $\varepsilon$ terms are sets of zero mean, not necessarily Gaussian random variations.

To remove the mean difference terms $b$, we subtract the mean $E[\cdot]$ of each data set:

$$X'_{1,i} = T'_i + \varepsilon_{X1,i} \tag{2a}$$
$$X'_{2,i} = T'_i + \varepsilon_{X2,i} \tag{2b}$$
$$X'_{3,i} = T'_i + \varepsilon_{X3,i}, \tag{2c}$$

where primes denote difference from the mean.

The unique set of the variance of differences between data sets can be written

$$\text{Var}\,[X_{1,i} - X_{2,i}] = \text{Var}\,[\varepsilon_{X1,i}] + \text{Var}\,[\varepsilon_{X2,i}] - 2\text{Cov}\,[\varepsilon_{X1,i}, \varepsilon_{X2,i}] \tag{3a}$$
$$\text{Var}\,[X_{1,i} - X_{3,i}] = \text{Var}\,[\varepsilon_{X1,i}] + \text{Var}\,[\varepsilon_{X3,i}] - 2\text{Cov}\,[\varepsilon_{X1,i}, \varepsilon_{X3,i}] \tag{3b}$$
$$\text{Var}\,[X_{2,i} - X_{3,i}] = \text{Var}\,[\varepsilon_{X2,i}] + \text{Var}\,[\varepsilon_{X3,i}] - 2\text{Cov}\,[\varepsilon_{X2,i}, \varepsilon_{X3,i}], \tag{3c}$$

where Cov $[\cdot]$ is the covariance between two quantities.

The relations for error variance can be derived by linearly combining Eqs. (3a-c):

$$\text{Var}\,[\varepsilon_{X1,i}] = \frac{1}{2}\left(\text{Var}\,[X_{1,i} - X_{2,i}] + \text{Var}\,[X_{1,i} - X_{3,i}] - \text{Var}\,[X_{2,i} - X_{3,i}]\right)$$
$$+ \text{Cov}\,[\varepsilon_{X1,i}, \varepsilon_{X2,i}] + \text{Cov}\,[\varepsilon_{X1,i}, \varepsilon_{X3,i}] - \text{Cov}\,[\varepsilon_{X2,i}, \varepsilon_{X3,i}] \tag{4a}$$

$$\text{Var}\,[\varepsilon_{X2,i}] = \frac{1}{2}\left(\text{Var}\,[X_{1,i} - X_{2,i}] + \text{Var}\,[X_{2,i} - X_{3,i}] - \text{Var}\,[X_{1,i} - X_{3,i}]\right)$$
$$+ \text{Cov}\,[\varepsilon_{X1,i}, \varepsilon_{X2,i}] + \text{Cov}\,[\varepsilon_{X2,i}, \varepsilon_{X3,i}] - \text{Cov}\,[\varepsilon_{X1,i}, \varepsilon_{X3,i}] \tag{4b}$$

$$\text{Var}\,[\varepsilon_{X3,i}] = \frac{1}{2}\left(\text{Var}\,[X_{1,i} - X_{3,i}] + \text{Var}\,[X_{2,i} - X_{3,i}] - \text{Var}\,[X_{1,i} - X_{2,i}]\right)$$
$$+ \text{Cov}\,[\varepsilon_{X1,i}, \varepsilon_{X3,i}] + \text{Cov}\,[\varepsilon_{X2,i}, \varepsilon_{X3,i}] - \text{Cov}\,[\varepsilon_{X1,i}, \varepsilon_{X2,i}]. \tag{4c}$$

## Four-cornered hat relations

Following the derivation of the 3CH relations, it may be shown that the four-cornered hat error variance relation for data set $X_{1,i}$ is

$$\text{Var}\,[\varepsilon_{X1,i}] = \frac{1}{3}\left(\text{Var}\,[X_{1,i} - X_{2,i}] + \text{Var}\,[X_{1,i} - X_{3,i}] + \text{Var}\,[X_{1,i} - X_{4,i}]\right)$$
$$- \frac{1}{6}\left(\text{Var}\,[X_{2,i} - X_{3,i}] + \text{Var}\,[X_{2,i} - X_{4,i}] + \text{Var}\,[X_{3,i} - X_{4,i}]\right)$$
$$+ \frac{2}{3}\left(\text{Cov}\,[\varepsilon_{X1,i}, \varepsilon_{X2,i}] + \text{Cov}\,[\varepsilon_{X1,i}, \varepsilon_{X3,i}] + \text{Cov}\,[\varepsilon_{X1,i}, \varepsilon_{X4,i}]\right)$$
$$- \frac{1}{3}\left(\text{Cov}\,[\varepsilon_{X2,i}, \varepsilon_{X3,i}] + \text{Cov}\,[\varepsilon_{X2,i}, \varepsilon_{X4,i}] + \text{Cov}\,[\varepsilon_{X3,i}, \varepsilon_{X4,i}]\right). \tag{5}$$

## Contact*

sjoberg@ucar.edu

## N-cornered hat relations

The error variance relation derivation may be generalized for $N$ data sets such that we may write the error variance for $X_{1,i}$ as

$$\text{Var}\,[\varepsilon_{1,i}] = \frac{1}{N-1}\sum_{n=2}^{N}\text{Var}\,[X_{1,i} - X_{n,i}]$$
$$- \frac{1}{(N-1)(N-2)}\sum_{n=2}^{N-1}\sum_{m=n+1}^{N}\text{Var}\,[X_{n,i} - X_{m,i}]$$
$$+ \frac{2}{N-1}\sum_{n=2}^{N}\text{Cov}\,[X_{1,i}, X_{n,i}]$$
$$- \frac{2}{(N-1)(N-2)}\sum_{n=2}^{N-1}\sum_{m=j+1}^{N}\text{Cov}\,[X_{n,i}, X_{m,i}],$$

and so on for the additional $N-1$ data sets.

## 3CH relations using 4 data sets

Note that we may use four data sets with the 3CH method to get three separate relations for the error variance of each data set. E.g., consider $X_{1,i}$ and assume all error covariance terms are 0 – a standard assumption for applying the 3CH method to real data – we find

$$\text{Var}\,[\varepsilon_{X1,i}] = \frac{1}{2}\left(\text{Var}\,[X_{1,i} - X_{2,i}] + \text{Var}\,[X_{1,i} - X_{3,i}] - \text{Var}\,[X_{2,i} - X_{3,i}]\right)$$

$$\text{Var}\,[\varepsilon_{X1,i}] = \frac{1}{2}\left(\text{Var}\,[X_{1,i} - X_{2,i}] + \text{Var}\,[X_{1,i} - X_{4,i}] - \text{Var}\,[X_{2,i} - X_{4,i}]\right)$$

$$\text{Var}\,[\varepsilon_{X1,i}] = \frac{1}{2}\left(\text{Var}\,[X_{1,i} - X_{3,i}] + \text{Var}\,[X_{1,i} - X_{4,i}] - \text{Var}\,[X_{3,i} - X_{4,i}]\right).$$

Combining,

$$3\text{Var}\,[\varepsilon_{X1,i}] = \frac{1}{2}(\text{Var}\,[X_{1,i} - X_{2,i}] + \text{Var}\,[X_{1,i} - X_{3,i}] + \text{Var}\,[X_{1,i} - X_{2,i}]$$
$$+ \text{Var}\,[X_{1,i} - X_{4,i}] + \text{Var}\,[X_{1,i} - X_{3,i}] + \text{Var}\,[X_{1,i} - X_{4,i}]$$
$$- \text{Var}\,[X_{2,i} - X_{3,i}] - \text{Var}\,[X_{2,i} - X_{4,i}] - \text{Var}\,[X_{3,i} - X_{4,i}]).$$

This can be reduced to

$$\text{Var}\,[\varepsilon_{X1,i}] = \frac{1}{3}\left(\text{Var}\,[X_{1,i} - X_{2,i}] + \text{Var}\,[X_{1,i} - X_{3,i}] + \text{Var}\,[X_{1,i} - X_{4,i}]\right)$$
$$- \frac{1}{6}\left(\text{Var}\,[X_{2,i} - X_{3,i}] + \text{Var}\,[X_{2,i} - X_{4,i}] + \text{Var}\,[X_{3,i} - X_{4,i}]\right),$$

which is identical to the four-cornered hat error variance relation for $X_{1,i}$ in Eq. (5) with 0-valued error covariance terms.

## Number of relations for the 3CH method

For the 3CH method with $N$ data sets, we can write the system of variance of differences as

$$\text{Var}\,[X_{1,i} - X_{2,i}] = \text{Var}[\varepsilon_{X1,i}] + \text{Var}[\varepsilon_{X2,i}] - 2\text{Cov}[\varepsilon_{X1,i}, \varepsilon_{X2,i}]$$
$$\text{Var}\,[X_{1,i} - X_{3,i}] = \text{Var}[\varepsilon_{X1,i}] + \text{Var}[\varepsilon_{X3,i}] - 2\text{Cov}[\varepsilon_{X1,i}, \varepsilon_{X3,i}]$$
$$\vdots$$
$$\text{Var}\,[X_{1,i} - X_{N,i}] = \text{Var}[\varepsilon_{X1,i}] + \text{Var}[\varepsilon_{XN,i}] - 2\text{Cov}[\varepsilon_{X1,i}, \varepsilon_{XN,i}]$$
$$\vdots$$
$$\text{Var}\,[X_{N-1,i} - X_{N,i}] = \text{Var}[\varepsilon_{X(N-1),i}] + \text{Var}[\varepsilon_{XN,i}]$$
$$- 2\text{Cov}[\varepsilon_{X(N-1),i}, \varepsilon_{XN,i}].$$

Note that in the above, there are $N-1$ relations containing instances of $\text{Var}[\varepsilon]$ for each data set $X$. The method of solving for error variance is to combine any two of these relations along with a third that includes the relevant two data sets that are carried along with those two relations. The total number of error variance relations $\mathcal{N}_\varepsilon$ is thus given by "$(N-1)$ choose 2," or

$$\mathcal{N}_\varepsilon = \sum_{i=1}^{N-2} i = \frac{(N-1)(N-2)}{2}. \tag{6}$$

## Other variance relationships

The set of variances for $X_{1,i}$ can be written

$$\text{Var}\,[X'_{1,i}] = \text{Var}\,[T'_i] + \text{Var}\,[\varepsilon_{X1,i}] + 2E\,[\varepsilon_{X1,i} T'_i]. \tag{7a}$$

For data sets $X'_{1,i}$ and $X'_{2,i}$, the variance of the sum is

$$\text{Var}\,[X'_{1,i} + X'_{2,i}] = E\left[(T'_i + \varepsilon_{X1,i} + T'_i + \varepsilon_{X2,i})^2\right]$$
$$= E\left[4T'^2_i + \varepsilon^2_{X1,i} + \varepsilon^2_{X2,i} + 2\varepsilon_{X1,i}\varepsilon_{X2,i}\right.$$
$$\left. + 4\varepsilon_{X1,i}T'_i + 4\varepsilon_{X2,i}T'_i\right]$$
$$= 4\text{Var}\,[T'_i] + \text{Var}\,[\varepsilon_{X1,i}] + \text{Var}\,[\varepsilon_{X2,i}]$$
$$+ 2\text{Cov}\,[\varepsilon_{X1,i}\varepsilon_{X2,i}] + 4E\,[(\varepsilon_{X1,i} + \varepsilon_{X2,i})T'_i]. \tag{9}$$

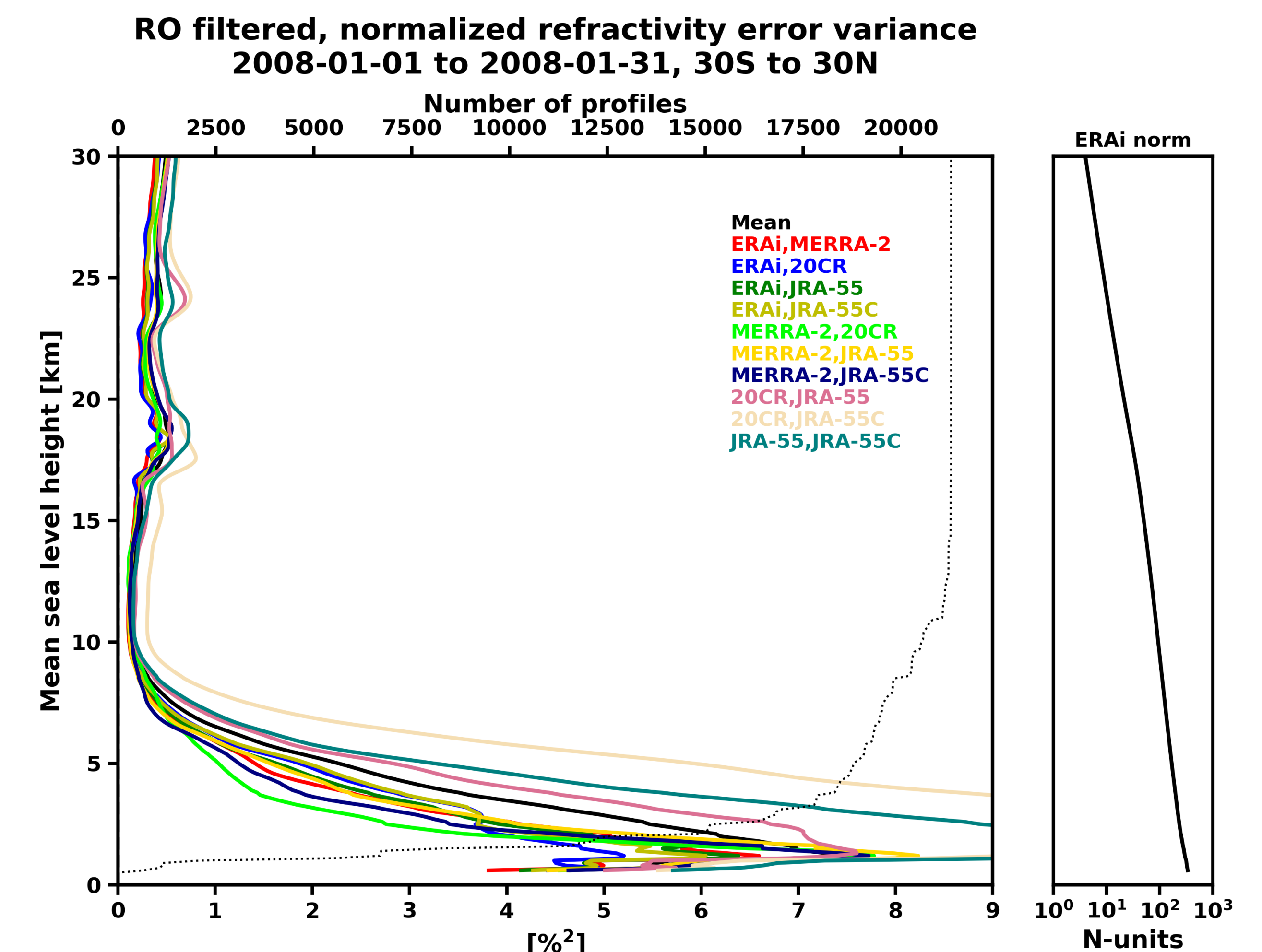## Example: 3CH with six data sets



Figure 1: Normalized 3CH error variance estimates for refractivity from COSMIC radio occultation data spanning 2008-07-01 through 2008-07-31. The respective triplets are labeled by color and the mean of the ten estimates is shown in black. The number of input data points are shown by the black dotted curve. The normalization is shown on the right.

For six data sets – COSMIC RO, ERAi, MERRA-2, 20th Century Reanalysis, JRA-55, and JRA-55C – as in Figure 1, we can produce ten estimates of the error variance for a given data set. This allows us to evaluate the spread between the estimates, telling us how valid our assumption of zero error covariance is. In the above, the relatively small spread above 10 km suggests that assumption is valid while the increased spread near the surface suggests that there are non-zero error covariances.

## The two-cornered hat relations

The two-cornered hat – or "triple co-location method" – error variance relation for $X_{1,i}$ is found by taking Eqs. (9)-(3a)-4·(7a):

$$\text{Var}\,[X'_{1,i} + X'_{2,i}] - \text{Var}\,[X'_{1,i} - X'_{2,i}] - 4\text{Var}\,[X'_{1,i}]$$
$$= -4\text{Var}\,[\varepsilon_{X1,i}] + 4\text{Cov}\,[\varepsilon_{X1,i}\varepsilon_{X2,i}] - 4E\,[\varepsilon_{X1,i}T'_i] + 4E\,[\varepsilon_{X2,i}T'_i]$$
$$\implies 4\text{Var}\,[\varepsilon_{X1,i}]$$
$$= 4\text{Var}\,[X'^2_{1,i} - X'_{1,i}X'_{2,i}] + 4\text{Cov}\,[\varepsilon_{X1,i}\varepsilon_{X2,i}] - 4E\,[(\varepsilon_{X1,i} - \varepsilon_{X2,i})T'_i]$$
$$\implies \text{Var}\,[\varepsilon_{X1,i}]$$
$$= \frac{1}{2}\left(\text{Var}\,[X'_{1,i} - X'_{2,i}] + \text{Var}\,[X'_{1,i}] - \text{Var}\,[X'_{2,i}]\right)$$
$$+ \text{Cov}\,[\varepsilon_{X1,i}, \varepsilon_{X2,i}] - E\,[(\varepsilon_{X1,i} - \varepsilon_{X2,i})T'_i]$$

Note that unlike the 3CH method and its generalization to $N$ data sets, the above relies on knowledge or assumptions about the reference data set $T_i$.

## Summary

For data sets that can be cast following Eqs. (1a-c), we may derive the 3CH relations for error variance. By assuming that there is no error covariance between the three data sets, we may apply this method to observations.

Here we show that:

- The 3CH method can be generalized for error variance estimation using $N$ different data sets.
- The single estimate of error variance using the $N$ data sets is equal to the mean of the $\frac{(N-1)(N-2)}{2}$ estimates of error variances for each individual data set using the 3CH method. These estimates can be used to evaluate our assumption of no error covariance between any given triplet of data sets.
- The 3CH method contains as a subset three estimates of error variance using only two sets of data – sometimes called the "triple co-location method." Application of this method requires making assumptions about the reference data set that will lead to decreased accuracy.

## References

- Anthes and Rieckh (2018), DOI: 10.5194/amt-11-4239-2018
- Rieckh and Anthes (2018), DOI: 10.5194/amt-11-4309-2018