



Potential and limitations of assimilating multiple collocated datasets with “optimal” error estimates

Statistical view

Annika Vogel^{1,2} and Richard Ménard¹

¹ Air Quality Research Division, Environment and Climate Change Canada

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



Canada 

THE ERROR ESTIMATION PROBLEM IN DATA ASSIMILATION

- 2 datasets

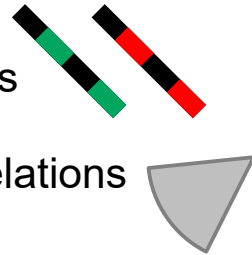
+1

- Given:
 - 1 innovation covariance



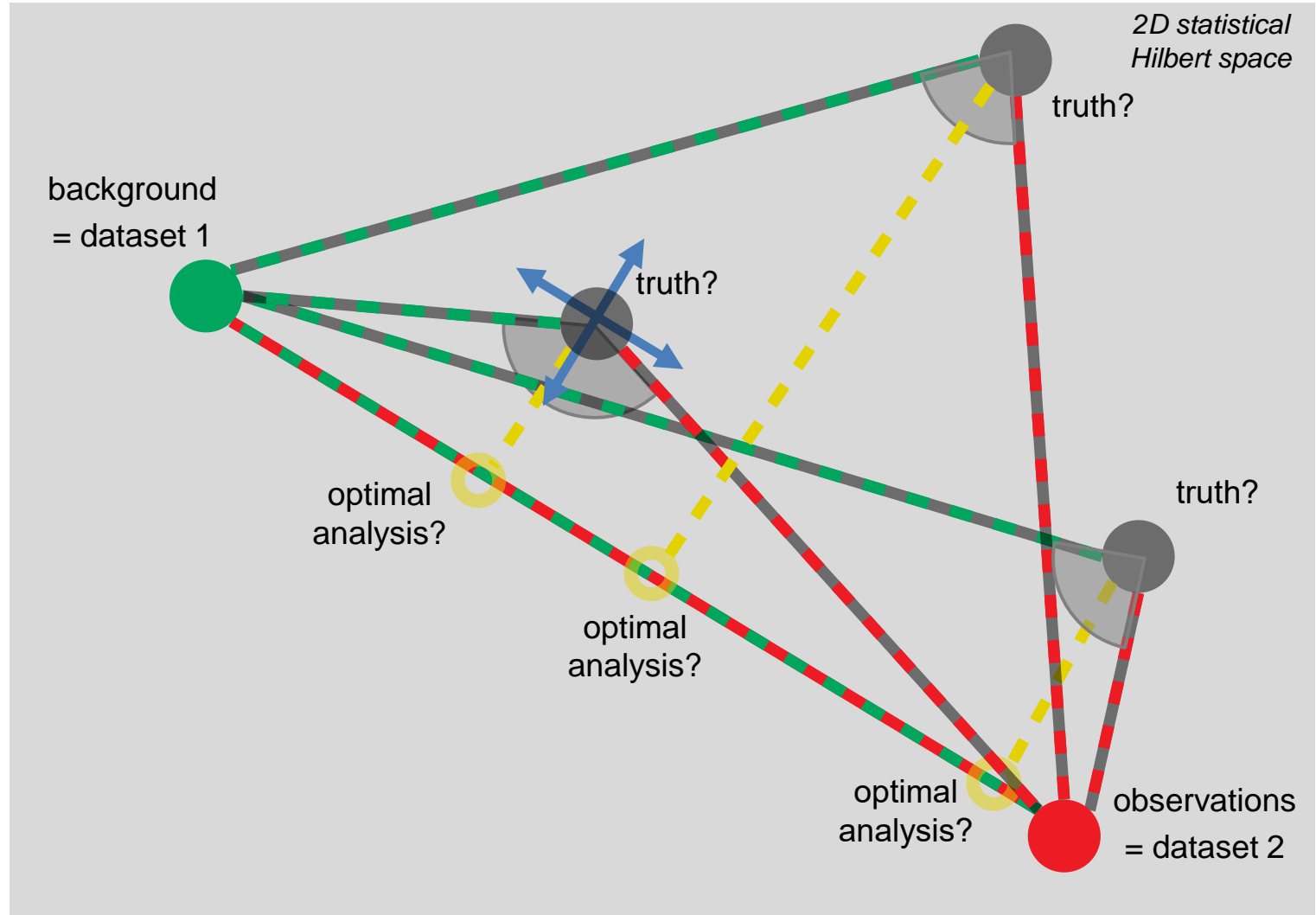
-3

- Unknown:
 - 2 error covariances
 - 1 error cross-correlations



-2





Underdetermined!



[Vogel & Ménard, egusphere-2022-996]



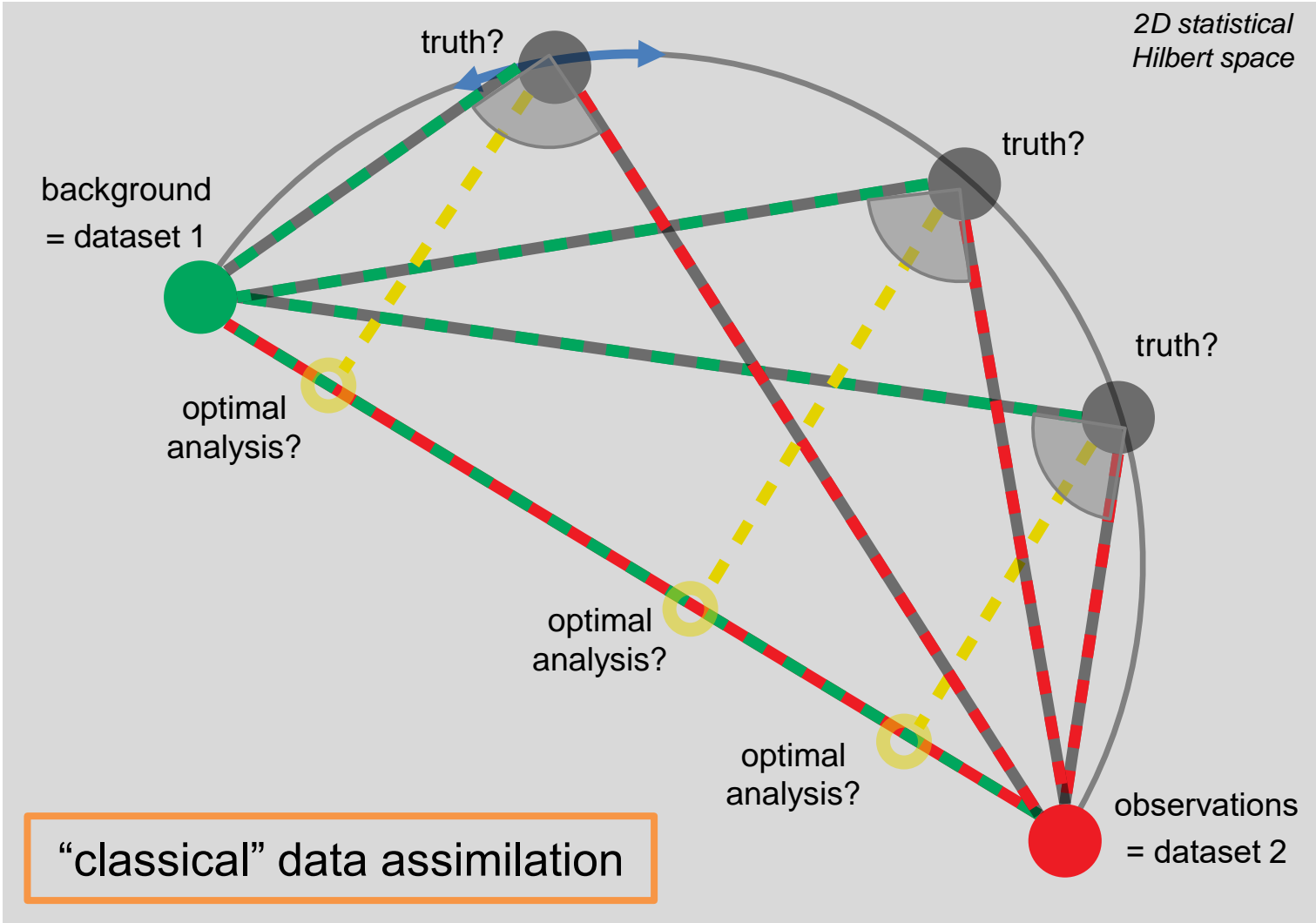
THE ERROR ESTIMATION PROBLEM IN DATA ASSIMILATION

- 2 datasets
- **Given:**
 - +1 – 1 innovation covariance 
- **Unknown:**
 - 3 – 2 error covariances 
 - 3 – 1 error cross-correlations 
- **Assumptions:**
 - +1 – uncorrelated errors 

-1

Still underdetermined!

– but analysis error only function of analysis state



[Vogel & Ménard, egusphere-2022-996]

ANALYSIS WITH CORRELATED ERRORS – KF EQUATIONS

- Standard KF analysis equations (for 2 datasets) assume uncorrelated errors between BG and OBS:

$$x_a = (\mathbb{1} - \mathbf{K}_1 \mathbf{H}_1) x_b + \mathbf{K}_1 y_1$$

$$\mathbf{K}_1 = \mathbf{B} \mathbf{H}_1^T (\mathbf{H}_1 \mathbf{B} \mathbf{H}_1^T + \mathbf{R}_1)^{-1}$$

$$\mathbf{A} = (\mathbb{1} - \mathbf{K}_1 \mathbf{H}_1) \mathbf{B} (\mathbb{1} - \mathbf{H}_1^T \mathbf{K}_1^T) + \mathbf{K}_1 \mathbf{R}_1 \mathbf{K}_1^T$$

x_b background, y_1 observation, x_a analysis state

\mathbf{B} background, \mathbf{R}_1 observation, \mathbf{A} analysis error covariance

\mathbf{H}_1 observation operator, \mathbf{K}_1 Kalman Gain

- KF analysis equations with correlations (for 2 datasets)*: **What if the errors are not fully uncorrelated?**

$$x_a = (\mathbb{1} - \mathbf{K}_1 \mathbf{H}_1) x_b + \mathbf{K}_1 y_1$$

error cross-covariance: $\mathbf{X}_{b;1} := \overline{\epsilon_B \cdot \epsilon_1^T}$

$$\mathbf{K}_1 = (\mathbf{B} \mathbf{H}_1^T - \mathbf{X}_{b;1}) (\mathbf{H}_1 \mathbf{B} \mathbf{H}_1^T - \mathbf{H}_1 \mathbf{X}_{b;1} - \mathbf{X}_{b;1}^T \mathbf{H}_1^T + \mathbf{R}_1)^{-1}$$

$$\mathbf{A} = (\mathbb{1} - \mathbf{K}_1 \mathbf{H}_1) \mathbf{B} (\mathbb{1} - \mathbf{H}_1^T \mathbf{K}_1^T) + (\mathbb{1} - \mathbf{K}_1 \mathbf{H}_1) \mathbf{X}_{b;1} \mathbf{K}_1^T + \mathbf{K}_1 \mathbf{X}_{b;1}^T (\mathbb{1} - \mathbf{H}_1^T \mathbf{K}_1^T) + \mathbf{K}_1 \mathbf{R}_1 \mathbf{K}_1^T$$

– cross-covariance terms in Kalman gain and analysis error covariance matrix

* partly equivalent to serially correlated observation errors [Daley 1992, MWR,v120,pp.164]

ANALYSIS WITH CORRELATED ERRORS – KF EQUATIONS

- Standard KF analysis equations (for 2 datasets) assume uncorrelated errors between BG and OBS:

x_b background, y_1 observation, x_a analysis state

\mathbf{B} background, \mathbf{R}_1 observation, \mathbf{A} analysis error covariance

\mathbf{H}_1 observation operator, \mathbf{K}_1 Kalman Gain

innovation covariance: $\Gamma_1 := \mathbf{H}_1 \mathbf{B} \mathbf{H}_1^T + \mathbf{R}_1$

$$\mathbf{K}_1 = \mathbf{B} \mathbf{H}_1^T \Gamma_1^{-1}$$

$$\mathbf{A}^{-1} = \mathbf{B}^{-1} + \mathbf{H}_1^T \mathbf{R}_1^{-1} \mathbf{H}_1$$

- KF analysis equations with correlations (for 2 datasets)*: **What if the errors are not fully uncorrelated?**

substituted error covariances:

$$\tilde{\mathbf{B}} := \mathbf{B} - \mathbf{X}_{b;1} \mathbf{H}_1^{+T} \quad \tilde{\mathbf{R}}_1 := \mathbf{R}_1 - \mathbf{X}_{b;1}^T \mathbf{H}_1^T \quad \tilde{\mathbf{A}} := \mathbf{A} - \mathbf{X}_{b;1} \mathbf{H}_1^{+T}$$

+ pseudo-inverse

generalized
innovation
covariance:

$$\begin{aligned} \tilde{\Gamma}_1 &:= \mathbf{H}_1 \mathbf{B} \mathbf{H}_1^T - \mathbf{H}_1 \mathbf{X}_{b;1} - \mathbf{X}_{b;1}^T \mathbf{H}_1^T + \mathbf{R}_1 \\ &= \mathbf{H}_1 \tilde{\mathbf{B}} \mathbf{H}_1^T + \tilde{\mathbf{R}}_1 \end{aligned}$$

asymmetry
(of error cross-
covariance):

$$\mathbf{Y}_{b;1} := \mathbf{X}_{b;1} - \mathbf{H}_1^+ \mathbf{X}_{b;1}^T \mathbf{H}_1^T$$

$$\mathbf{K}_1 = \tilde{\mathbf{B}} \mathbf{H}_1^T \tilde{\Gamma}_1^{-1}$$

$$\tilde{\mathbf{A}}^{-1} = \tilde{\mathbf{B}}^{-1} + \mathbf{H}_1^T \tilde{\mathbf{R}}_1^{-T} \mathbf{H}_1 - \mathbf{H}_1^T \tilde{\mathbf{R}}_1^{-T} \mathbf{Y}_{b;1}^T \tilde{\mathbf{B}}^{-1}$$

- equivalent form of substituted Kalman gain
- additional “asymmetry-term” in substituted analysis error covariance

* partly equivalent to serially correlated observation errors [Daley 1992, MWR,v120,pp.164]

ANALYSIS WITH CORRELATED ERRORS - SENSITIVITIES

- Sensitivity of KF analysis variance to error correlation:

- scalar ($\mathbf{H} \rightarrow 1$)

$$x_{b;1} = \rho \sigma_b \sigma_1$$

$$\frac{\partial \sigma_a^2}{\partial \rho} = 2 \left(\overset{\text{background variance}}{\sigma_b \sigma_1} (\overset{\text{background variance}}{\sigma_b^2} - x_{b;1}) (\overset{\text{observation variance}}{\sigma_1^2} - x_{b;1}) (\overset{\text{cross-variance}}{\sigma_b^2} - 2 x_{b;1} + \overset{\text{cross-variance}}{\sigma_1^2}) \right)^{-2}$$

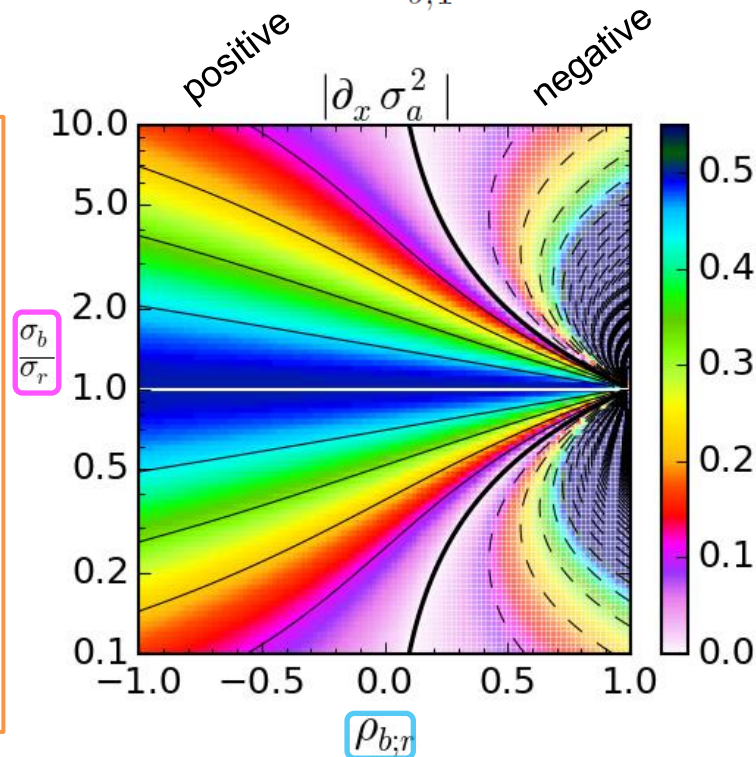
- Sensitivity of KF analysis variance to error cross-variance:

- define ratio of standard-deviations

$$\alpha := \frac{\sigma_b}{\sigma_1}$$

$$\frac{\partial \sigma_a^2}{\partial x_{b;1}} = 2 (\alpha - \rho) (\alpha^{-1} - \rho) (\alpha - 2\rho + \alpha^{-1})^{-2}$$

- Only a function of error correlation and ratio
- Symmetric around max. at $\alpha = 1$
- Decreasing with differences between BG and OBS
- Almost linear for small differences ($\alpha \in [\frac{1}{3}; 3]$)
- Negative sensitivities for high correlations
= when cross-variance > smaller variance
- Rapidly decreasing negative sensitivity
for high correlation and moderate $\alpha \in [\frac{1}{6}; 6]$



ASSIMILATION FROM THREE DATASETS

- 3 datasets

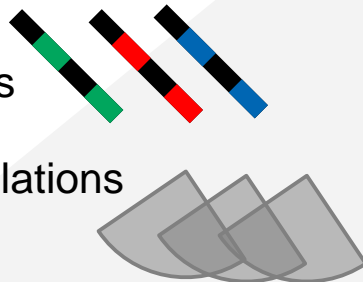
Given:

- 3 innovation covariances



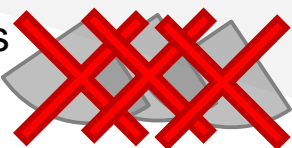
Unknown:

- 3 error covariances
- 3 error cross-correlations



Assumptions:

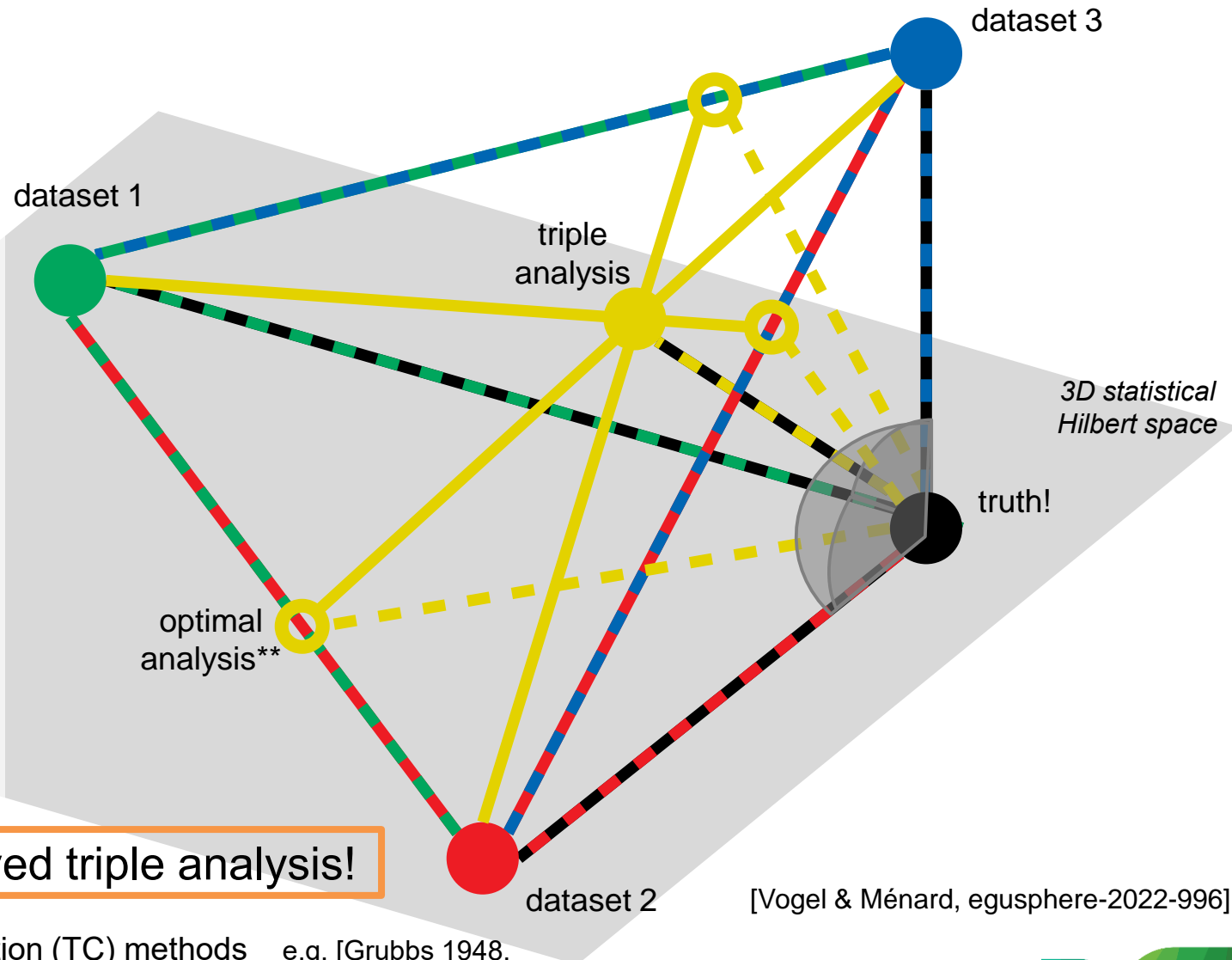
- uncorrelated errors



0

Fully determined!*

Improved triple analysis!



[Vogel & Ménard, egusphere-2022-996]

* equivalent to scalar 3 cornered hat (3CH) and triple collocation (TC) methods e.g. [Grubbs 1948, doi:10.1080/01621459.1948.10483261]; [Stoffelen 1998, doi:10.1029/97JC03180]; [Sjoberg et al. 2021, doi:10.1175/jtech-d-19-0217.1]

** equivalent to cross-validation approach [Ménard&Deshaies-Jacques 2018a&b, doi:10.3390/atmos9030086 & atmos9020070]



GENERALIZED ANALYSIS FROM MULTIPLE DATASETS – KF EQUATIONS

What if we have more than two datasets?

- Direct KF analysis equations (for 3 datasets)*:

$$x_a = \underbrace{(\mathbb{1} - \mathbf{W}_1 \mathbf{H}_1 - \mathbf{W}_2 \mathbf{H}_2)}_{:= \mathbf{W}_b} x_b + \mathbf{W}_1 y_1 + \mathbf{W}_2 y_2$$

$$\mathbf{A} = \mathbf{W}_b \mathbf{B} \mathbf{W}_b^T + \mathbf{W}_1 \mathbf{R}_1 \mathbf{W}_1^T + \mathbf{W}_2 \mathbf{R}_2 \mathbf{W}_2^T$$

W weights

define
Kalman gain-
like matrices:

$$\mathbf{K}_{d_1} := \mathbf{B} \mathbf{H}_1^T (\mathbf{H}_1 \mathbf{B} \mathbf{H}_1^T + \mathbf{R}_1)^{-1} = \mathbf{B} \mathbf{H}_1^T \mathbf{\Gamma}_1^{-1}$$

$$\mathbf{K}_{d_2} := \mathbf{B} \mathbf{H}_2^T (\mathbf{H}_2 \mathbf{B} \mathbf{H}_2^T + \mathbf{R}_2)^{-1} = \mathbf{B} \mathbf{H}_2^T \mathbf{\Gamma}_2^{-1}$$

e.g.: $\mathbf{W}_1 = (\mathbb{1} - \mathbf{W}_2 \mathbf{H}_2) \mathbf{B} \mathbf{H}_1^T (\mathbf{H}_1 \mathbf{B} \mathbf{H}_1^T + \mathbf{R}_1)^{-1}$

$$\mathbf{W}_1 = (\mathbb{1} - \mathbf{K}_{d_2} \mathbf{H}_2) \mathbf{K}_{d_1} (\mathbb{1} - \mathbf{H}_1 \mathbf{K}_{d_2} \mathbf{H}_2 \mathbf{K}_{d_1})^{-1}$$

- Sequential KF analysis equations (for 3 datasets)*:

$$x_{a_2} = \underbrace{(\mathbb{1} - \mathbf{K}_{a_2} \mathbf{H}_2) (\mathbb{1} - \mathbf{K}_{a_1} \mathbf{H}_1)}_{:= \mathbf{W}_b} x_b + \underbrace{(\mathbb{1} - \mathbf{K}_{a_2} \mathbf{H}_2) \mathbf{K}_{a_1}}_{:= \mathbf{W}_1} y_1 + \underbrace{\mathbf{K}_{a_2}}_{:= \mathbf{W}_2} y_2$$

$$\mathbf{W}_b = \mathbf{A}_2 \mathbf{B}^{-1}$$

$$\mathbf{A}_2^{-1} = \mathbf{B}^{-1} + \mathbf{H}_1^T \mathbf{R}_1^{-1} \mathbf{H}_1 + \mathbf{H}_2^T \mathbf{R}_2^{-1} \mathbf{H}_2$$

$$\mathbf{W}_1 = \mathbf{A}_2 \mathbf{H}_1^T \mathbf{R}_1^{-1}$$

$$\mathbf{W}_2 = \mathbf{A}_2 \mathbf{H}_2^T \mathbf{R}_2^{-1}$$

- summation-like generalization of analysis state and error covariance
- weight(/gain) reduced by weight of additional dataset, equivalent analysis-based form
- equivalence of direct and sequential form

* equivalent to multi-model KF e.g. [Logutov&Robinson 2005, doi:10.1256/qj.05.99]; [Narayan et al. 2012, doi:10.1016/j.jcp.2012.06.002]

GENERALIZED ANALYSIS FROM MULTIPLE DATASETS – VAR EQUATIONS

What if we have more than two datasets?

- Direct 3D-var equations (for 3 datasets)*:

$$J_2(x) = \frac{1}{2} \left[(x - x_b)^T \mathbf{B}^{-1} (x - x_b) + (\mathbf{H}_1 x - y_1)^T \mathbf{R}_1^{-1} (\mathbf{H}_1 x - y_1) + \underline{(\mathbf{H}_2 x - y_2)^T \mathbf{R}_2^{-1} (\mathbf{H}_2 x - y_2)} \right]$$

$$\nabla_x J_2(x) = \mathbf{B}^{-1} (x - x_b) + \mathbf{H}_1^T \mathbf{R}_1^{-1} (\mathbf{H}_1 x - y_1) + \underline{\mathbf{H}_2^T \mathbf{R}_2^{-1} (\mathbf{H}_2 x - y_2)}$$

- Sequential 3D-var equations (for 3 datasets) assuming vanishing gradient in 1st assimilation:

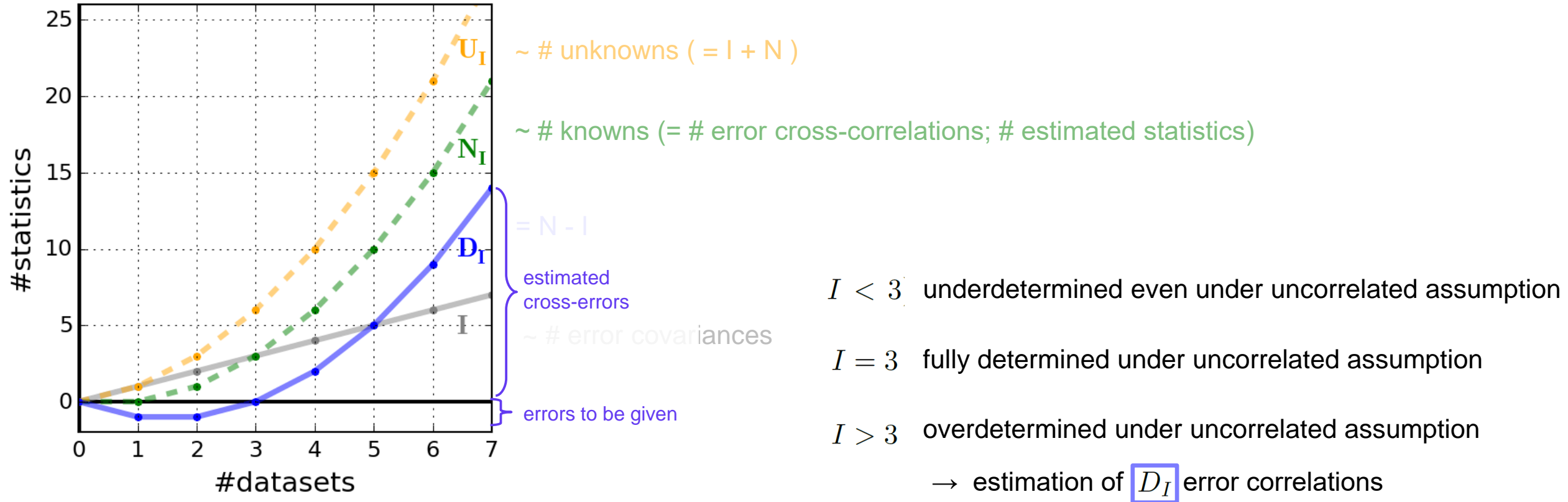
$$J_2(x) = \frac{1}{2} \left[\underbrace{(x - x_b)^T \mathbf{B}^{-1} (x - x_b) + (\mathbf{H}_1 x - y_1)^T \mathbf{R}_1^{-1} (\mathbf{H}_1 x - y_1) + \underline{(\mathbf{H}_2 x - y_2)^T \mathbf{R}_2^{-1} (\mathbf{H}_2 x - y_2)}}_{\rightarrow \text{dependent on optimization variable } x} + \underbrace{\left(\mathbf{B}^{-1} x_b + \mathbf{H}_1^T \mathbf{R}_1^{-1} y_1 \right)^T \left(\mathbf{B}^{-1} + \mathbf{H}_1^T \mathbf{R}_1^{-1} \mathbf{H}_1 \right)^{-1} \left(\mathbf{B}^{-1} x_b + \mathbf{H}_1^T \mathbf{R}_1^{-1} y_1 \right) - \left(x_b^T \mathbf{B}^{-1} x_b + y_1^T \mathbf{R}_1^{-1} y_1 \right)}_{\rightarrow \text{independent of } x} \right]$$

$$\nabla_x J_2(x) = \mathbf{B}^{-1} (x - x_b) + \mathbf{H}_1^T \mathbf{R}_1^{-1} (\mathbf{H}_1 x - y_1) + \underline{\mathbf{H}_2^T \mathbf{R}_2^{-1} (\mathbf{H}_2 x - y_2)}$$

- summation-like generalization of cost function and gradient
- constant additional term in sequential cost function → equivalence

* equivalent to multi-model KF e.g. [Logutov&Robinson 2005, doi:10.1256/qj.05.99]; [Narayan et al. 2012, doi:10.1016/j.jcp.2012.06.002]

GENERALIZED ERROR ESTIMATION PROBLEM



Consequences:

- Relative number of assumed statistics reduces with increasing number of datasets.
- Absolute number of assumptions increases with number of datasets (system is never closed).

Estimation of all error covariances and some cross-correlations

[Vogel & Ménard, egusphere-2022-996]



GENERALIZED ANALYSIS WITH CORRELATED ERRORS – KF EQUATIONS

- Direct KF analysis equations with correlation (for 3 datasets):

What if we have more than two correlated datasets?

substituted error covariances:

$$\begin{aligned} \tilde{\mathbf{B}}_1 &:= \mathbf{B} - \mathbf{X}_{b;1} \mathbf{H}_1^+ T & \tilde{\mathbf{B}}_2 &:= \mathbf{B} - \mathbf{X}_{b;2} \mathbf{H}_2^+ T & \tilde{\Gamma}_{1;2} &:= \mathbf{H}_1 \mathbf{B} \mathbf{H}_2^T - \mathbf{H}_1 \mathbf{X}_{b;2} - \mathbf{X}_{b;1}^T \mathbf{H}_2^T + \mathbf{X}_{1;2} \\ \tilde{\Gamma}_1 &:= \mathbf{H}_1 \mathbf{B} \mathbf{H}_1^T - \mathbf{H}_1 \mathbf{X}_{b;1} - \mathbf{X}_{b;1}^T \mathbf{H}_1^T + \mathbf{R}_1 & \tilde{\Gamma}_2 &:= \mathbf{H}_2 \mathbf{B} \mathbf{H}_2^T - \mathbf{H}_2 \mathbf{X}_{b;2} - \mathbf{X}_{b;2}^T \mathbf{H}_2^T + \mathbf{R}_2 \end{aligned}$$

$$\mathbf{W}_1 = \left(\tilde{\mathbf{B}}_1 \mathbf{H}_1^T - \tilde{\mathbf{B}}_2 \mathbf{H}_2^T \tilde{\Gamma}_2^{-1} \tilde{\Gamma}_{1;2}^T \right) \left(\tilde{\Gamma}_1 - \tilde{\Gamma}_{1;2} \tilde{\Gamma}_2^{-1} \tilde{\Gamma}_{1;2}^T \right)^{-1} \quad \text{compare: 2 datasets} \quad \mathbf{K}_1 = \tilde{\mathbf{B}} \mathbf{H}_1^T \tilde{\Gamma}_1^{-1}$$

- Direct KF analysis equations with correlation (for I datasets):

$$\mathbf{W}_1 = \left[\mathbf{B} \mathbf{H}_1^T - \mathbf{X}_{b;1} - \sum_{j=2}^{I-1} \mathbf{W}_j \left(\mathbf{H}_j \mathbf{B} \mathbf{H}_1^T - \mathbf{X}_{j;b} \mathbf{H}_1^T - \mathbf{H}_j \mathbf{X}_{b;1} + \mathbf{X}_{j;1} \right) \right] \left(\mathbf{H}_1 \mathbf{B} \mathbf{H}_1^T - \mathbf{H}_1 \mathbf{X}_{b;1} - \mathbf{X}_{1;b} \mathbf{H}_1^T + \mathbf{C}_1 \right)^{-1} \quad \dots???$$

- solution consistent with uncorrelated form and correlated form for 2 datasets
- significant increase in complexity
- use sequential form for assimilating multiple correlated datasets

GENERALIZED ANALYSIS WITH CORRELATED ERRORS - SENSITIVITIES

- Sensitivity of generalized analysis eq. :

- scalar ($\mathbf{H} \rightarrow 1$)
- background and $N=I-1$ observations with common variance
- no cross-correlation among obs., common cross-correlation to background
- define common ratio of standard-deviations

$$x_{b;r} = \rho \sigma_b \sigma_r \quad \forall r$$

$$\alpha := \frac{\sigma_b}{\sigma_r}$$

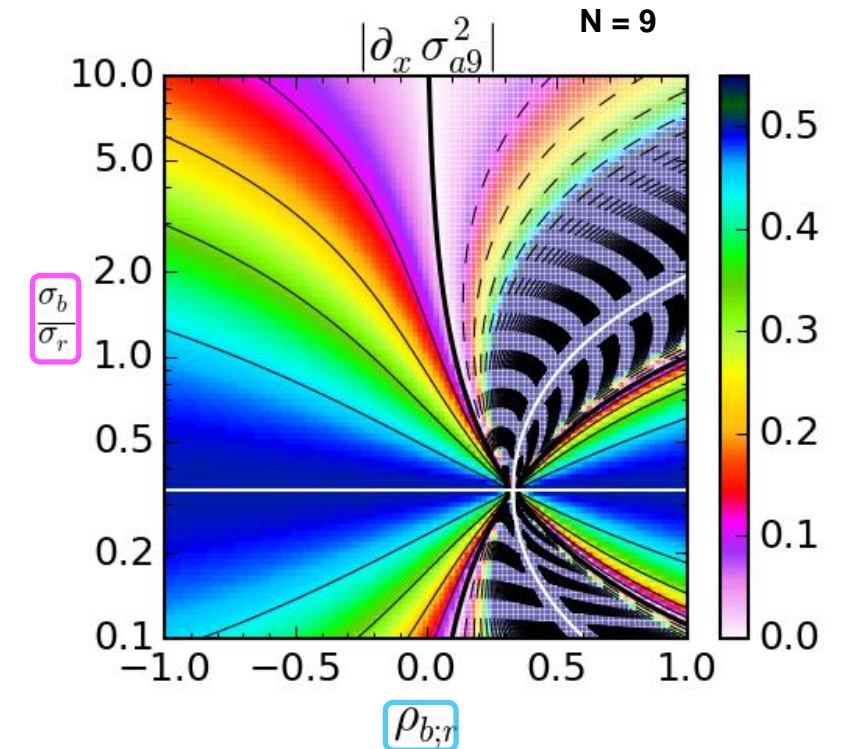
- analysis error variance and cross-variance from obs.-only

$$\sigma_{ar}^2 = \left[\sum_{n=1}^N (\sigma_n^2)^{-1} \right]^{-1} = \frac{1}{N} \sigma_r^2, \quad x_{b;ar} = x_{b;1} + \sum_{n=2}^N K_n (x_{b;n} - x_{b;1}) \xrightarrow{\text{scalar cross-variance}} 0 = x_{b;r}$$

- generalized sensitivity of final analysis error variance

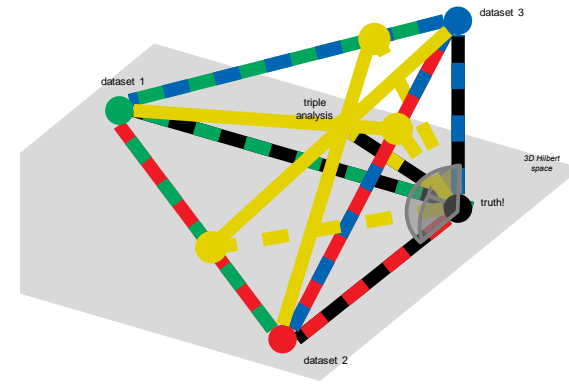
$$\frac{\partial \sigma_a^2}{\partial x_{b;r}} = 2(\alpha - \rho) \left(\frac{1}{N\alpha} - \rho \right) \left(\alpha - 2\rho + \frac{1}{N\alpha} \right)^{-2}$$

- symmetry line decreases (were $\sigma_{ar}^2 = \sigma_b^2$)
- common line of vanishing sensitivity (were $\alpha = \rho$)
- significant negative sensitivities already for decreasing correlations



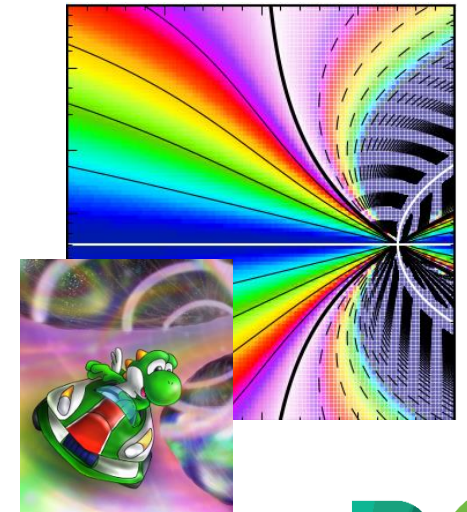
SUMMARY / CONCLUSIONS

- Error estimation:
 - Multiple collocated datasets enable estimation of optimal error statistics, incl. cross-correlations
 - “Datasets” may include multiple forecasts, observations, ...
 - Number of estimated error statistics increases with number of datasets
 - Some assumptions and conditions remain



Manuscript: Vogel & Ménard, “**How far can the statistical error estimation problem be closed by collocated data?**”
under review @ NPG: [egusphere-2022-996](https://www.nature.com/egusphere/egusphere2022996)

- Assimilation:
 - Improved analysis of multiple datasets
 - Direct assimilation of uncorrelated datasets equivalent to sequential form
 - Assimilation of correlated datasets becomes expensive
 - Critical sensitivities for high correlations and large differences in errors



“We cannot have too many datasets! (If we have appropriate assimilation algorithms...)”

Contact: annika.vogel@ec.gc.ca

<https://www.pinterest.ca/pin/469711436108621063/>