

# Introduction to Data Assimilation, Subgrid-scale Parameterization and Predictability

**Christian Franzke**

Meteorological Institute

Center for Earth System Research and Sustainability

University of Hamburg

Email: [christian.franzke@uni-hamburg.de](mailto:christian.franzke@uni-hamburg.de)

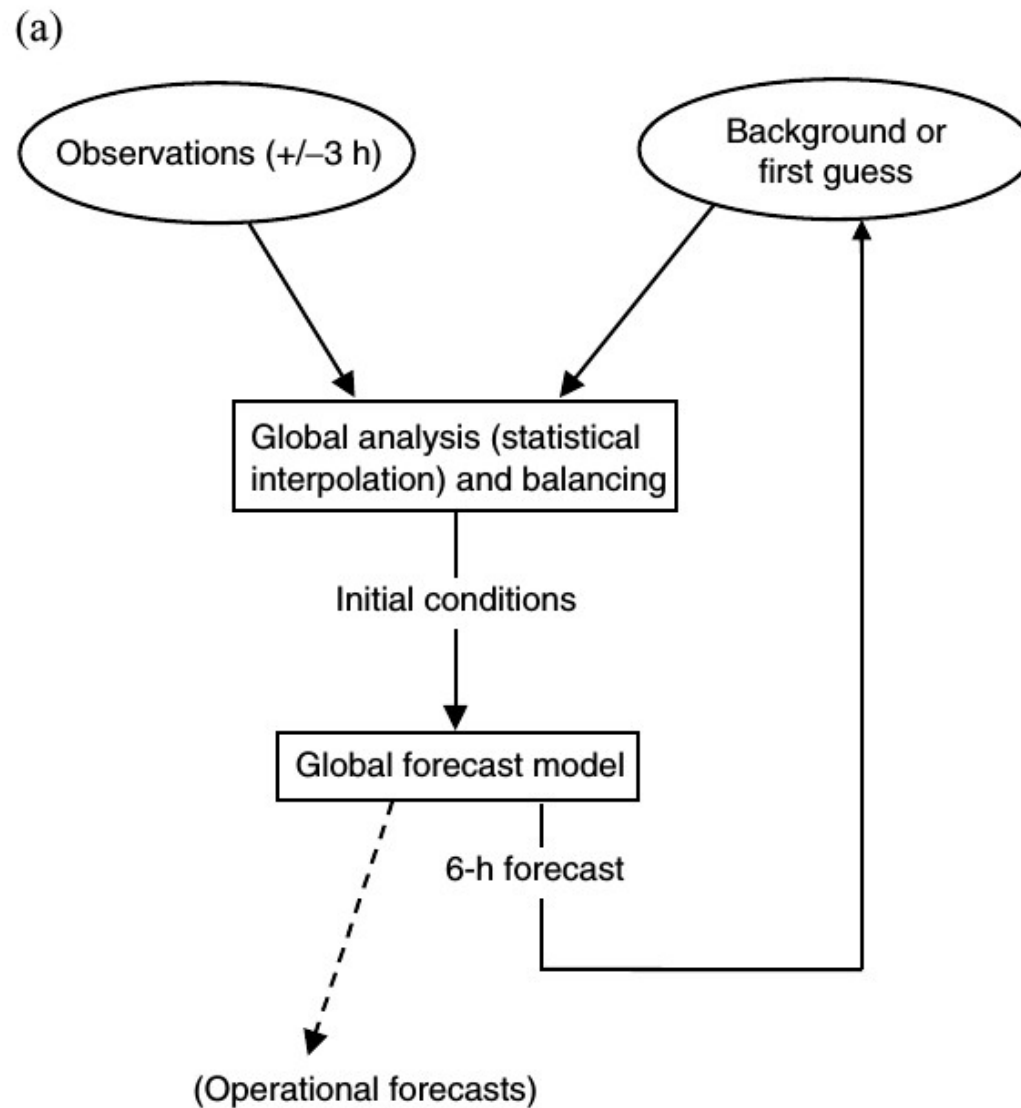
# Outline

---

- Multivariate Data Assimilation
- Ensemble Kalman Filter

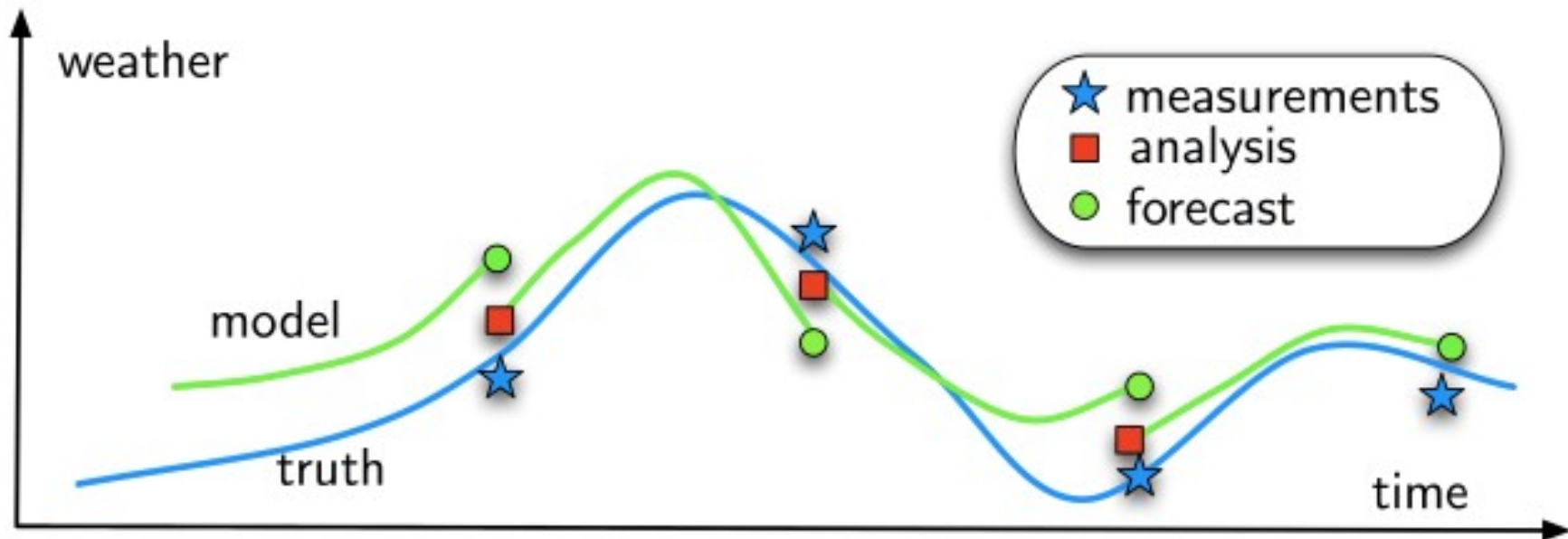
# Forecast Cycle

---



# Data Assimilation

---



# Optimal Interpolation

---

$$\mathbf{x}_t - \mathbf{x}_b = \mathbf{W}[\mathbf{y}_o - H(\mathbf{x}_b)] - \boldsymbol{\varepsilon}_a = \mathbf{W}\mathbf{d} - \boldsymbol{\varepsilon}_a$$

$$\boldsymbol{\varepsilon}_a = \mathbf{x}_a - \mathbf{x}_t$$

- Truth, analysis and background are vectors of length  $n$
- Weights  $\mathbf{W}$  are  $(n \times p)$  matrix
- $\mathbf{y}_o$ : Observation
- $\mathbf{x}_t$ : Truth
- Observations are vector of length  $p$
- Observational operator  $H$  can be nonlinear

# Optimal Interpolation

---

Observational increments vector:

$$\mathbf{d} = \mathbf{y}_o - H(\mathbf{x}_b)$$

We define background and analysis error as:

$$\boldsymbol{\varepsilon}_b(x, y) = \mathbf{x}_b(x, y) - \mathbf{x}_t(x, y)$$

$$\boldsymbol{\varepsilon}_a(x, y) = \mathbf{x}_a(x, y) - \mathbf{x}_t(x, y)$$

# Optimal Interpolation

---

Observational errors (at irregularly spaced points)

$$\boldsymbol{\varepsilon}_{oi} = \mathbf{y}_o(\mathbf{r}_i) - \mathbf{y}_t(\mathbf{r}_i) = \mathbf{y}_o(\mathbf{r}_i) - H[\mathbf{x}_t(\mathbf{r}_i)]$$

# Optimal Interpolation

---

We define error covariance matrices:

$$\mathbf{P}_a = \mathbf{A} = E \left\{ \boldsymbol{\varepsilon}_a \boldsymbol{\varepsilon}_a^T \right\}$$

$$\mathbf{P}_b = \mathbf{B} = E \left\{ \boldsymbol{\varepsilon}_b \boldsymbol{\varepsilon}_b^T \right\}$$

$$\mathbf{P}_o = \mathbf{R} = E \left\{ \boldsymbol{\varepsilon}_o \boldsymbol{\varepsilon}_o^T \right\}$$



# Optimal Interpolation

---

The nonlinear observation operator  $H$  that transforms model variables into observed variables can be linearized as:

$$H(\mathbf{x} + \delta\mathbf{x}) = H(\mathbf{x}) + \mathbf{H}\delta\mathbf{x}$$

where  $\mathbf{H}$  is a  $p \times n$  matrix, denoting the linear observation operator with elements

$$h_{i,j} = \partial H_i / \partial x_j$$

# Optimal Interpolation

---

We assume that the background is a good approximation of the truth, so that the analysis and the observations are equal to the background values plus small increments

$$\begin{aligned}\mathbf{d} &= \mathbf{y}_o - H(\mathbf{x}_b) = \mathbf{y}_o - H(\mathbf{x}_t + (\mathbf{x}_b - \mathbf{x}_t)) \\ &= \mathbf{y}_o - H(\mathbf{x}_t) - \mathbf{H}(\mathbf{x}_b - \mathbf{x}_t) = \boldsymbol{\varepsilon}_o - \mathbf{H}\boldsymbol{\varepsilon}_b\end{aligned}$$

The  $\mathbf{H}$  matrix transforms vectors from model to observation space.

Its transpose or adjoint  $\mathbf{H}^T$  transforms vectors from observation to model space.

# Optimal Interpolation

---

Best unbiased estimation: Least Squares

$$\mathbf{W} = E(\mathbf{xy}^T) [E(\mathbf{yy}^T)]^{-1}$$



# How to estimate $W$ ?

---

Multiple Regression

# How to estimate W?

---

## Multiple Regression

Assume we have two time series of vectors

$$\mathbf{x}(t) = \begin{bmatrix} x_1(t) \\ x_2(t) \\ \vdots \\ x_n(t) \end{bmatrix} \quad \mathbf{y}(t) = \begin{bmatrix} y_1(t) \\ y_2(t) \\ \vdots \\ y_p(t) \end{bmatrix}$$

# How to estimate $\mathbf{W}$ ?

---


## Multiple Regression

We now derive the best linear unbiased estimation of  $\mathbf{x}$  in terms of  $\mathbf{y}$ , the optimal value of  $\mathbf{W}$

$$\mathbf{x}_a(t) = \mathbf{W}\mathbf{y}(t)$$

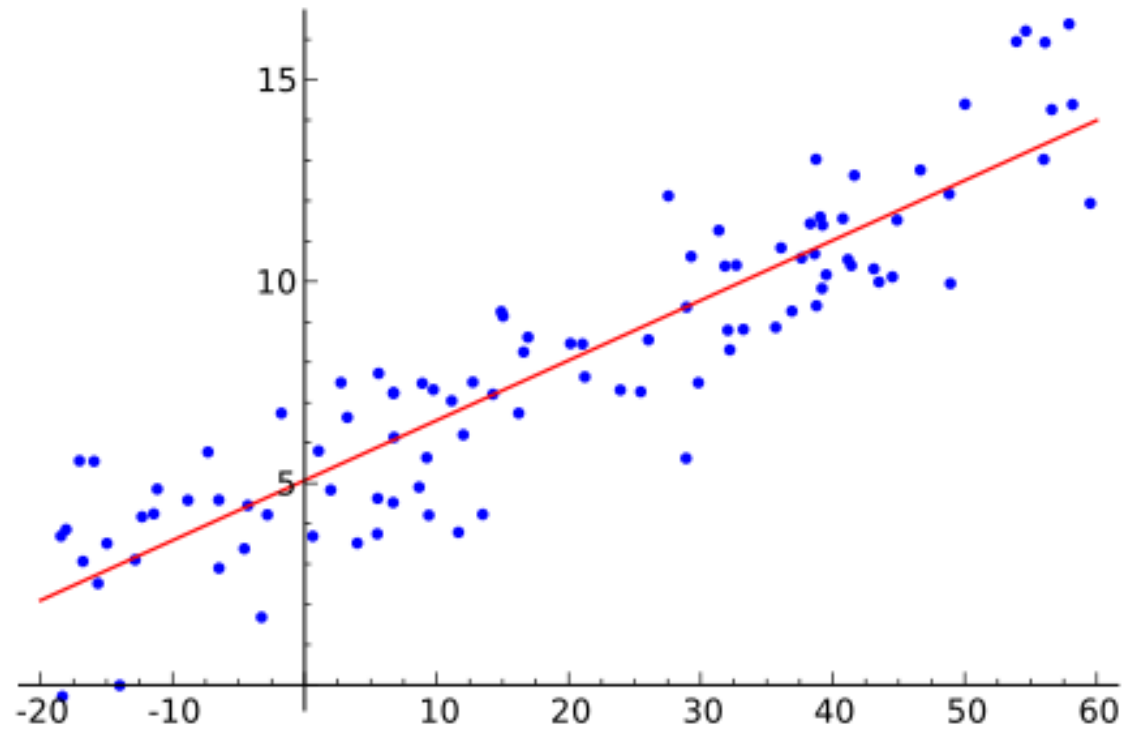
which approximates the true relationship

$$\mathbf{x}(t) = \mathbf{W}\mathbf{y}(t) - \varepsilon(t)$$

linear regression  (analysis) error

# How to estimate $W$ ?

---



# How to estimate W?

---

W minimizes mean squared error  $E(\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon})$

$$x_i(t) = \sum_{k=1}^p w_{ik} y_k(t) - \varepsilon_i(t)$$

$$\rightarrow \sum_{i=1}^n \varepsilon_i^2(t) = \sum_{i=1}^n \left[ \sum_{k=1}^p w_{ik} y_k(t) - x_i(t) \right]^2$$



# How to estimate W?

---

Derivative with respect to weight matrix

$$\begin{aligned}\frac{\partial \sum_{i=1}^n \varepsilon_i^2(t)}{\partial w_{ij}} &= 2 \left[ \sum_{k=1}^p w_{ik} y_k(t) - x_i(t) \right] [y_j(t)] \\ &= 2 \left[ \sum_{k=1}^p w_{ik} y_k(t) y_i(t) - x_i(t) y_j(t) \right]\end{aligned}$$

in Matrix form

$$\frac{\partial \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon}}{\partial w_{ij}} = 2 \left\{ [\mathbf{W} \mathbf{y}(t) \mathbf{y}^T(t)]_{ij} - [\mathbf{x}(t) \mathbf{y}^T(t)]_{ij} \right\}$$

# How to estimate $\mathbf{W}$ ?

---

We get the normal equations

$$\mathbf{W} E(\mathbf{y}\mathbf{y}^T) - E(\mathbf{x}\mathbf{y}^T) = 0$$

or

$$\mathbf{W} = E(\mathbf{x}\mathbf{y}^T) [E(\mathbf{y}\mathbf{y}^T)]^{-1}$$

which gives the best unbiased estimation

$$\mathbf{x}_a(t) = \mathbf{W}\mathbf{y}(t)$$

# Optimal Interpolation

---

Best unbiased estimation: Least Squares

$$\mathbf{W} = E(\mathbf{xy}^T) [E(\mathbf{yy}^T)]^{-1}$$

Estimate of weight matrix:

$$\begin{aligned}\mathbf{W} &= E\{(\mathbf{x}_t - \mathbf{x}_b)[\mathbf{y}_o - H(\mathbf{x}_b)]^T\} (E\{[\mathbf{y}_o - H(\mathbf{x}_b)][\mathbf{y}_o - H(\mathbf{x}_b)]^T\})^{-1} \\ &= E[(-\boldsymbol{\varepsilon}_b)(\boldsymbol{\varepsilon}_o - \mathbf{H}\boldsymbol{\varepsilon}_b)^T] \{E[(\boldsymbol{\varepsilon}_o - \mathbf{H}\boldsymbol{\varepsilon}_b)(\boldsymbol{\varepsilon}_o - \mathbf{H}\boldsymbol{\varepsilon}_b)^T]\}^{-1}\end{aligned}$$

$$\mathbf{W} = \mathbf{BH}^T (\mathbf{R} + \mathbf{HBH}^T)^{-1}$$

# Optimal Interpolation

---

The analysis covariance matrix is:

$$\begin{aligned}\mathbf{P}_a &= E\{\boldsymbol{\varepsilon}_a \boldsymbol{\varepsilon}_a^T\} = E\{\boldsymbol{\varepsilon}_b \boldsymbol{\varepsilon}_b^T + \boldsymbol{\varepsilon}_b (\boldsymbol{\varepsilon}_o - \mathbf{H} \boldsymbol{\varepsilon}_b)^T \mathbf{W}^T \\ &\quad + \mathbf{W} (\boldsymbol{\varepsilon}_o - \mathbf{H} \boldsymbol{\varepsilon}_b) \boldsymbol{\varepsilon}_b^T + \mathbf{W} (\boldsymbol{\varepsilon}_o - \mathbf{H} \boldsymbol{\varepsilon}_b) (\boldsymbol{\varepsilon}_o - \mathbf{H} \boldsymbol{\varepsilon}_b)^T \mathbf{W}^T\} \\ &= \mathbf{B} - \mathbf{B} \mathbf{H}^T \mathbf{W}^T - \mathbf{W} \mathbf{H} \mathbf{B} + \mathbf{W} \mathbf{R} \mathbf{W}^T + \mathbf{W} \mathbf{H} \mathbf{B} \mathbf{H}^T \mathbf{W}^T\end{aligned}$$

# Optimal Interpolation

---

The analysis covariance matrix is:

$$\begin{aligned}\mathbf{P}_a &= E\{\boldsymbol{\varepsilon}_a \boldsymbol{\varepsilon}_a^T\} = E\{\boldsymbol{\varepsilon}_b \boldsymbol{\varepsilon}_b^T + \boldsymbol{\varepsilon}_b (\boldsymbol{\varepsilon}_o - \mathbf{H} \boldsymbol{\varepsilon}_b)^T \mathbf{W}^T \\ &\quad + \mathbf{W} (\boldsymbol{\varepsilon}_o - \mathbf{H} \boldsymbol{\varepsilon}_b) \boldsymbol{\varepsilon}_b^T + \mathbf{W} (\boldsymbol{\varepsilon}_o - \mathbf{H} \boldsymbol{\varepsilon}_b) (\boldsymbol{\varepsilon}_o - \mathbf{H} \boldsymbol{\varepsilon}_b)^T \mathbf{W}^T\} \\ &= \mathbf{B} - \mathbf{B} \mathbf{H}^T \mathbf{W}^T - \mathbf{W} \mathbf{H} \mathbf{B} + \mathbf{W} \mathbf{R} \mathbf{W}^T + \mathbf{W} \mathbf{H} \mathbf{B} \mathbf{H}^T \mathbf{W}^T\end{aligned}$$

This can be written as:

$$\mathbf{P}_a = (\mathbf{I} - \mathbf{W} \mathbf{H}) \mathbf{B}$$

# Summary of OI

---

Basic equations of OI:

$$\mathbf{x}_a = \mathbf{x}_b + \mathbf{W}[\mathbf{y}_o - H(\mathbf{x}_b)] = \mathbf{x}_b + \mathbf{W}\mathbf{d}$$

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

$$\mathbf{P}_a = (\mathbf{I}_n - \mathbf{W}\mathbf{H})\mathbf{B}$$

# Summary of OI

---

## Interpretation

$$\mathbf{x}_a = \mathbf{x}_b + \mathbf{W}[\mathbf{y}_o - H(\mathbf{x}_b)] = \mathbf{x}_b + \mathbf{W}\mathbf{d}$$

The analysis is obtained by adding to the first guess the product of the optimal weight matrix and the innovation (difference between observation and first guess).

# Summary of OI

---

## Interpretation

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

The optimal weight matrix is given by the background error covariance in the observation space ( $\mathbf{B}\mathbf{H}^T$ ) multiplied by the inverse of the total error covariance.



# Summary of OI

---

## Interpretation

$$\mathbf{P}_a = (\mathbf{I}_n - \mathbf{WH})\mathbf{B}$$

The error covariance of the analysis is given by the error covariance of the background, reduced by a matrix equal to the identity matrix minus the optimal weight matrix.

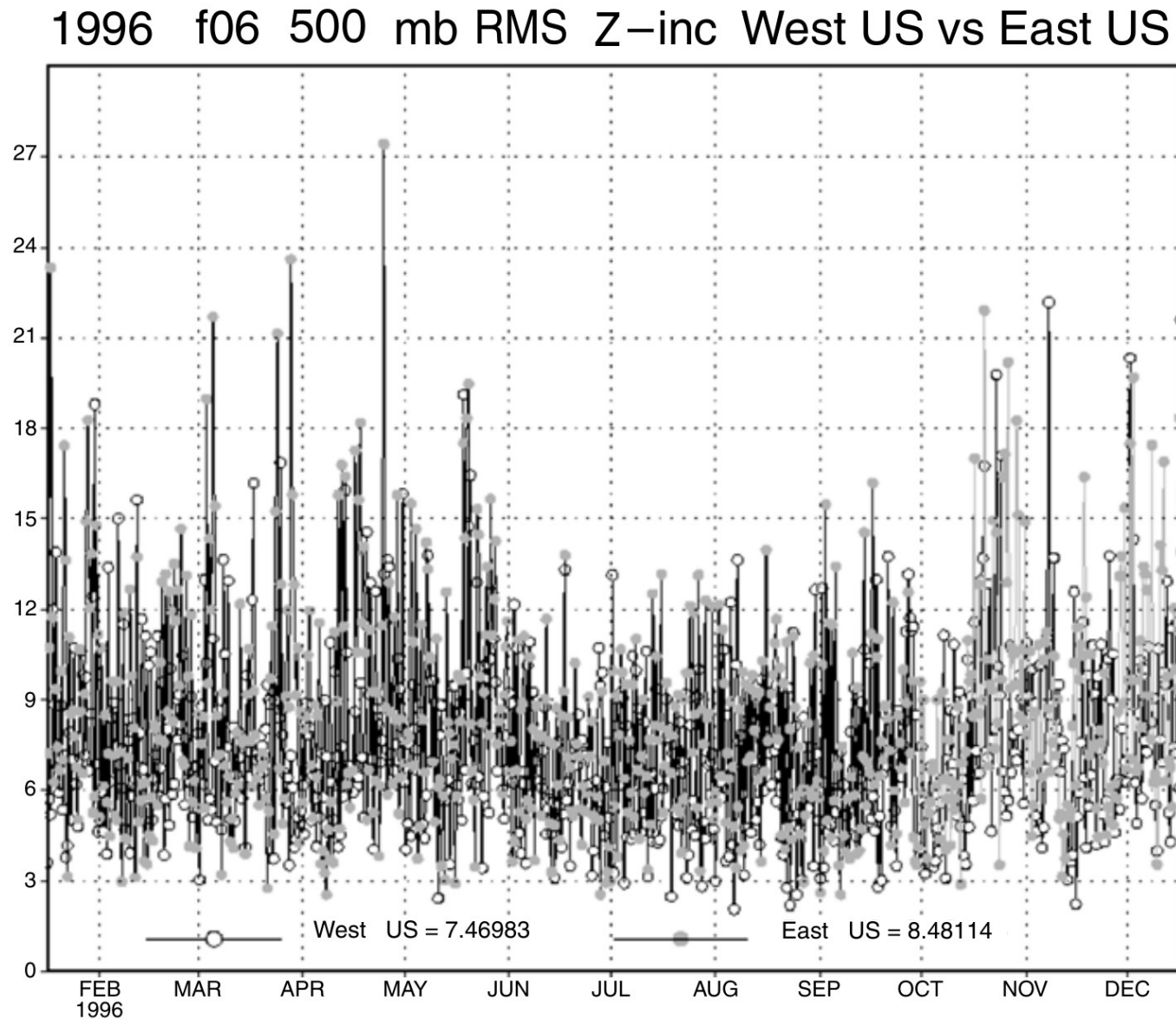
# 3D-VAR ~ OI

---

- OI: Find optimal weight matrix **W** that minimizes the analysis error covariance
- 3D-VAR: Find optimal analysis field that minimizes a cost function

$$2J(\mathbf{x}) = (\mathbf{x} - \mathbf{x}_b)^T \mathbf{B}^{-1} (\mathbf{x} - \mathbf{x}_b) + [\mathbf{y}_o - H(\mathbf{x})]^T \mathbf{R}^{-1} [\mathbf{y}_o - H(\mathbf{x})]$$

# DA with evolving error covariance



# 3D-VAR ~ 4D-VAR

---

## 3D-VAR

$$2J(\mathbf{x}) = (\mathbf{x} - \mathbf{x}_b)^T \mathbf{B}^{-1}(\mathbf{x} - \mathbf{x}_b) + [\mathbf{y}_o - H(\mathbf{x})]^T \mathbf{R}^{-1}[\mathbf{y}_o - H(\mathbf{x})]$$

## 4D-VAR

$$J[\mathbf{x}(t_0)] = \frac{1}{2}[\mathbf{x}(t_0) - \mathbf{x}^b(t_0)]^T \mathbf{B}_0^{-1}[\mathbf{x}(t_0) - \mathbf{x}^b(t_0)] \\ + \frac{1}{2} \sum_{i=0}^N [H(\mathbf{x}_i) - \mathbf{y}_i^o]^T \mathbf{R}_i^{-1}[H(\mathbf{x}_i) - \mathbf{y}_i^o]$$

# 4D-VAR

---

$$J[\mathbf{x}(t_0)] = \frac{1}{2}[\mathbf{x}(t_0) - \mathbf{x}^b(t_0)]^T \mathbf{B}_0^{-1} [\mathbf{x}(t_0) - \mathbf{x}^b(t_0)] \\ + \frac{1}{2} \sum_{i=0}^N [H(\mathbf{x}_i) - \mathbf{y}_i^o]^T \mathbf{R}_i^{-1} [H(\mathbf{x}_i) - \mathbf{y}_i^o]$$

The cost function includes a term measuring the distance to the background at the beginning of the interval and a summation over time of the cost function for each observation increment

# 4D-VAR

---

4D-Var seeks an initial condition such that the forecast best fits the observations within the assimilation interval.

4D-Var assumes **perfect model**  
→ Big disadvantage

# Kalman Filter ~ OI

---

In OI the forecast error covariance is computed once and for all.

Assumption: Forecast errors are statistically stationary

Kalman filter propagates error covariance forward  
→ State-dependent error covariance

# Extended Kalman Filter

---

Forecast step:

Nonlinear model

$$\mathbf{x}^f(t_i) = M_{i-1}[\mathbf{x}^a(t_{i-1})]$$

$$\mathbf{P}^f(t_i) = \mathbf{L}_{i-1} \mathbf{P}^a(t_{i-1}) \mathbf{L}_{i-1}^T + \mathbf{Q}(t_{i-1})$$

Tangent linear

Analysis step:

$$\mathbf{x}^a(t_i) = \mathbf{x}^f(t_i) + \mathbf{K}_i \mathbf{d}_i$$

$$\mathbf{P}^a(t_i) = (\mathbf{I} - \mathbf{K}_i \mathbf{H}_i) \mathbf{P}^f(t_i)$$

$$\mathbf{d}_i = \mathbf{y}_i^o - H[\mathbf{x}^f(t_i)]$$



# Kalman Filter

---

1) Define ensemble of observations

$$y_j^o = y^o + \varepsilon_j$$

# Kalman Filter

---

1) Define ensemble of observations

$$y_j^o = y^o + \varepsilon_j$$

2) Define ensemble covariance matrix

$$R_e = \overline{\varepsilon \varepsilon^T}$$

# Kalman Filter

---

1) Define ensemble of observations

$$y_j^o = y_j^f + \varepsilon_j$$

2) Define ensemble covariance matrix

$$R_e = \overline{\varepsilon \varepsilon^T}$$

3) Analysis step

$$y_j^a = y_j^f + P_e^f H^T (H P_e^f H^T + R_e)^{-1} (y_j^o - H y_j^f)$$

# Kalman Filter

---

3) Analysis step

$$y_j^a = y_j^f + P_e^f H^T (H P_e^f H^T + R_e)^{-1} (y_j^o - H y_j^f)$$

4) Analysis and ensemble mean are identical

$$\bar{y}_j^a = \bar{y}_j^f + P_e^f H^T (H P_e^f H^T + R_e)^{-1} (\bar{y}_j^o - H \bar{y}_j^f)$$

5) Kalman gain (Optimal weight) matrix

$$K_e = P_e^f H^T (H P_e^f H^T + R_e)^{-1}$$

# Kalman Filter

---

5) Kalman gain (Optimal weight) matrix

$$K_e = P_e^f H^T (H P_e^f H^T + R_e)^{-1}$$

6) Error covariance

$$P_e^a = (I - K_e H) P_e^f$$

# Kalman Filter

---

5) Kalman gain (Optimal weight) matrix

$$K_e = P_e^f H^T (H P_e^f H^T + R_e)^{-1}$$

6) Error covariance

$$P_e^a = (I - K_e H) P_e^f$$

7) Model error covariance  $Q$

# Kalman Filter

---

5) Kalman gain (Optimal weight) matrix

$$K_e = P_e^f H^T (H P_e^f H^T + R_e)^{-1}$$

6) Error covariance

$$P_e^a = (I - K_e H) P_e^f$$

7) Model error covariance  $Q$

8) Ensemble error covariance evolves as

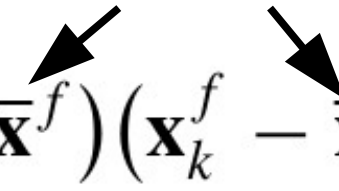
$$P_e^{k+1} = F P_e^k F^T + Q + \text{n.l.} \quad (F: \text{tangent linear operator})$$

# Ensemble Kalman Filter

---

- An ensemble of  $K$  data assimilation cycles is carried out simultaneously
- All cycles assimilate the same observations
- Different perturbations are added to observations  
→ ensemble forecasts
- This ensemble is used to estimate forecast error covariance

Ensemble mean

$$\mathbf{P}^f \approx \frac{1}{K-1} \sum_{k=1}^K (\mathbf{x}_k^f - \bar{\mathbf{x}}^f) (\mathbf{x}_k^f - \bar{\mathbf{x}}^f)^T$$




# Ensemble Kalman Filter

---

- Evolve each ensemble member forward using the nonlinear model perturbed by noise

$$y^a = M(y^f) + \varepsilon$$

- Compute ensemble mean and covariance

$$\bar{y}^f = 1/N \sum y^f$$

$$P^f = y^f (y^f)^T = 1/(N-1) \sum (y^f - \bar{y}^f)(y^f - \bar{y}^f)^T$$

- Update analysis

$$y^a = y^f + P^f K (y^o - H y^f)$$

$$K = P^f H^T (H P^f H^T + R)^{-1}$$

Where  $R$  is the covariance from the observations

# Excercise

---

Write a Ensemble Kalman Filter for Lorenz 96