

**Understanding data assimilation:**  
how observations and a model are  
weaved into the analysis via statistics

**Tomoko Matsuo**  
NCAR

Thanks: Doug Nychka, Jeff Anderson,  
Kevin Raeder, Alain Caya, Art Richmond, Gang Lu ...



NCAR

# What is data assimilation?!

## Combining Information

### prior knowledge of the state of system

empirical or physical models (e.g. physical laws)

*complete in space and time*  $x$

### observations

directly measured or retrieved quantities

*incomplete in space and time*  $y_o$

# Bayes Theorem

“Bayesian statistics provides a coherent probabilistic framework for most DA approaches” [e.g., Lorenc, 1986]

## prior knowledge

$$P(x) \sim N(x_f, \mathbf{P}_f)$$

$$x = x_f + \varepsilon_f$$

## observations

$$P(y_o | x) \sim N(H(x), \mathbf{R}) \quad y_o = H(x) + \varepsilon_o$$

Note: observations  $y$  conditioned upon the state  $x$

## posterior

$$P(x | y_o) \propto P(y_o | x)P(x)$$

$$P(x | y_o) \sim N(x_a, \mathbf{P}_a) \quad \mathbf{H} \text{ is linear}$$

$$\text{where } x_a = x_f + \mathbf{K}(y - \mathbf{H}x_f)$$

$$\mathbf{P}_a = (\mathbf{I} - \mathbf{K}\mathbf{H})\mathbf{P}_f$$

Assumption: Normal Distribution

# Importance of Covariance

[e.g. Rodgers, 2000]

prior

$$P(x) \sim N(x_f, \mathbf{P}_f)$$

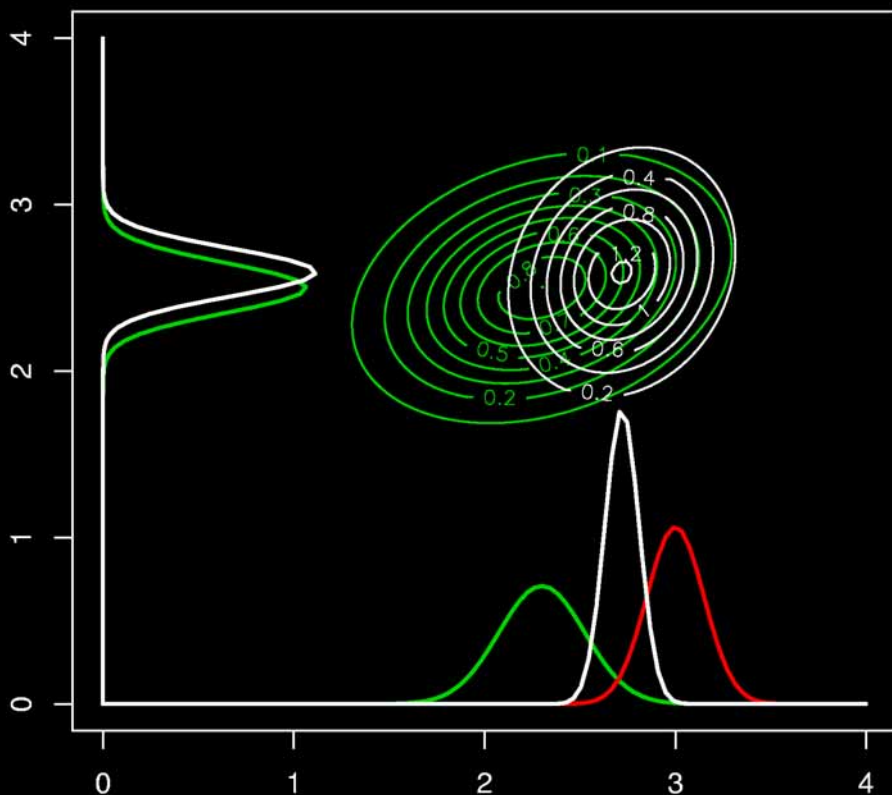
posterior

$$P(x | y_o)$$

observations

$$P(y_o | x) \sim N(\mathbf{H}x, \mathbf{R})$$

→  
update



$$x_f = \begin{pmatrix} 2.3 & 2.5 \end{pmatrix}$$

$$\mathbf{P}_f = \begin{pmatrix} 0.225 & 0.05 \\ 0.05 & 0.15 \end{pmatrix}$$

$$\mathbf{H} = \begin{pmatrix} 1 & 0 \end{pmatrix}$$

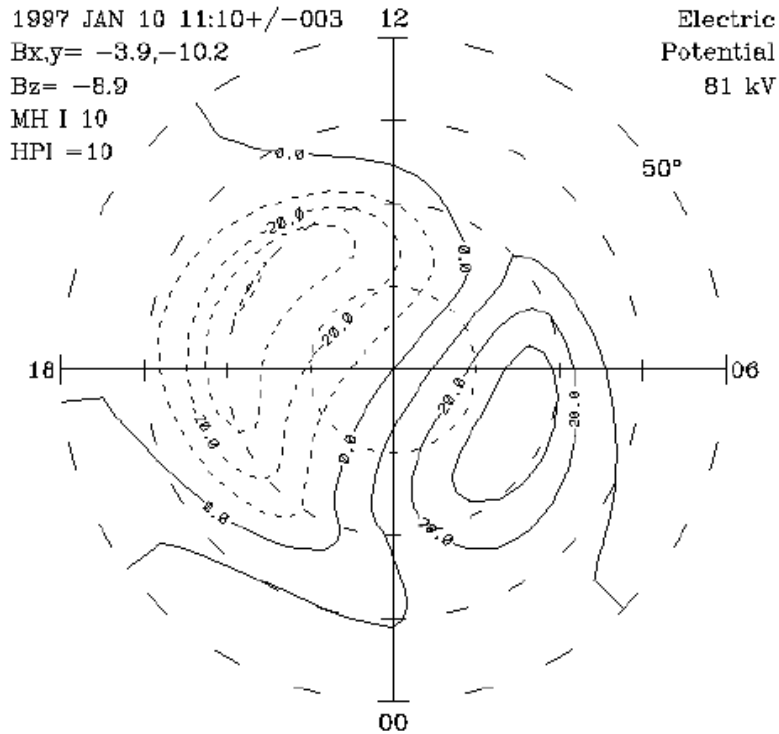
$x_1$  : observed

$x_2$  : unobserved

# Assimilative Mapping of Ionospheric Electrodynamics [Richmond and Kamide, 1988]

$$x_a = x_b + \mathbf{K}(y - \mathbf{H}x_b)$$

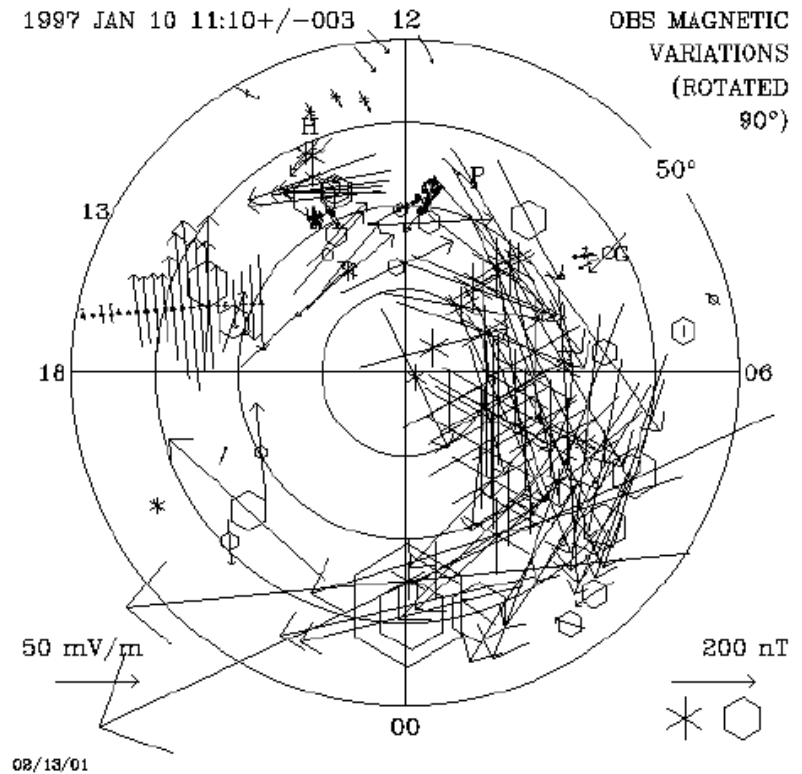
## prior knowledge



[Foster et al., 1986]

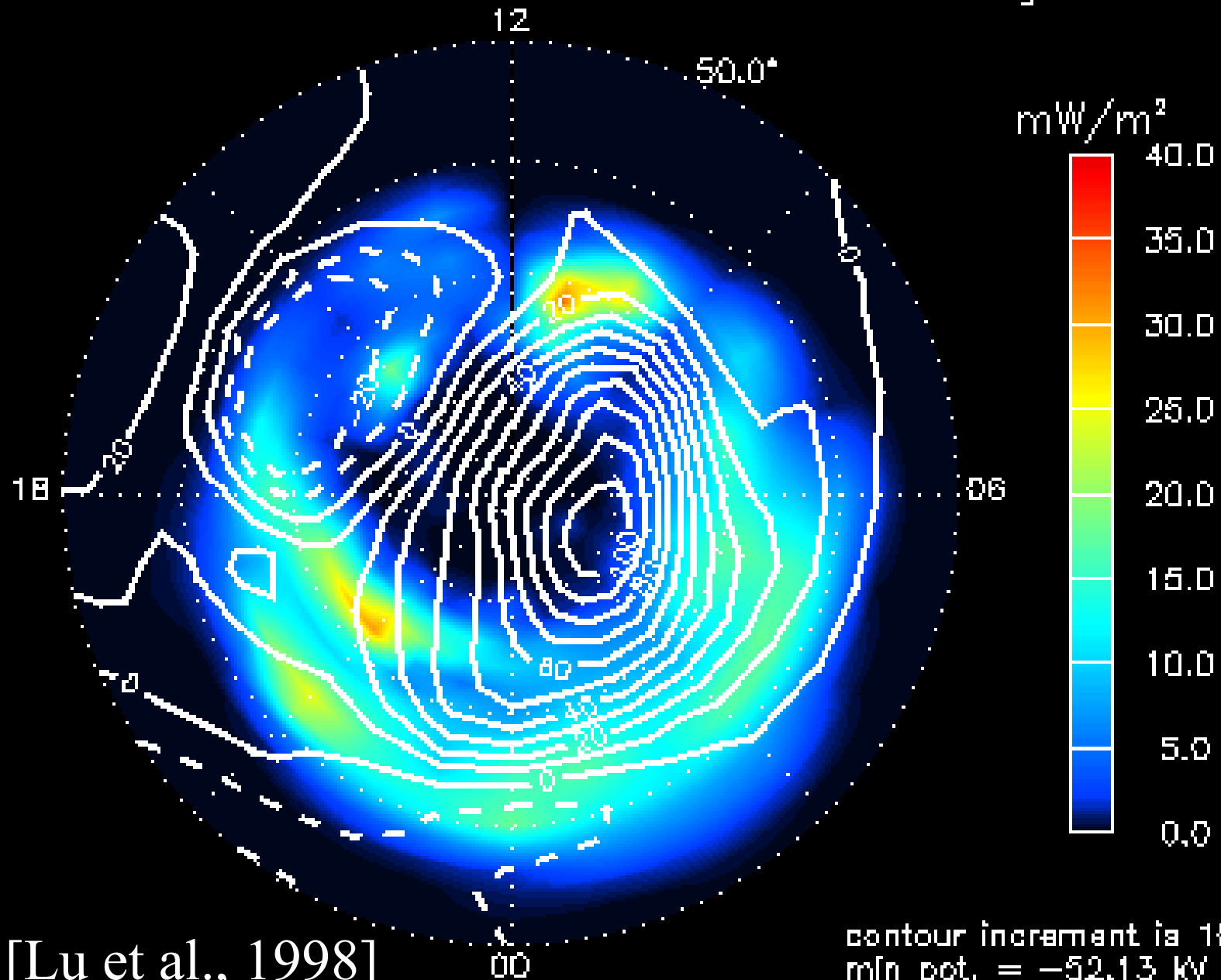
. OF 10,000 PPLS, SP+ 7.1412

## observations



Energy Flux (NH)  
with contours of Electric Potential

97:01:10 11:10 UTC  
data averaged over  $\pm 3$  mins



[Lu et al., 1998]

# Inhomogeneous / anisotropic covariance

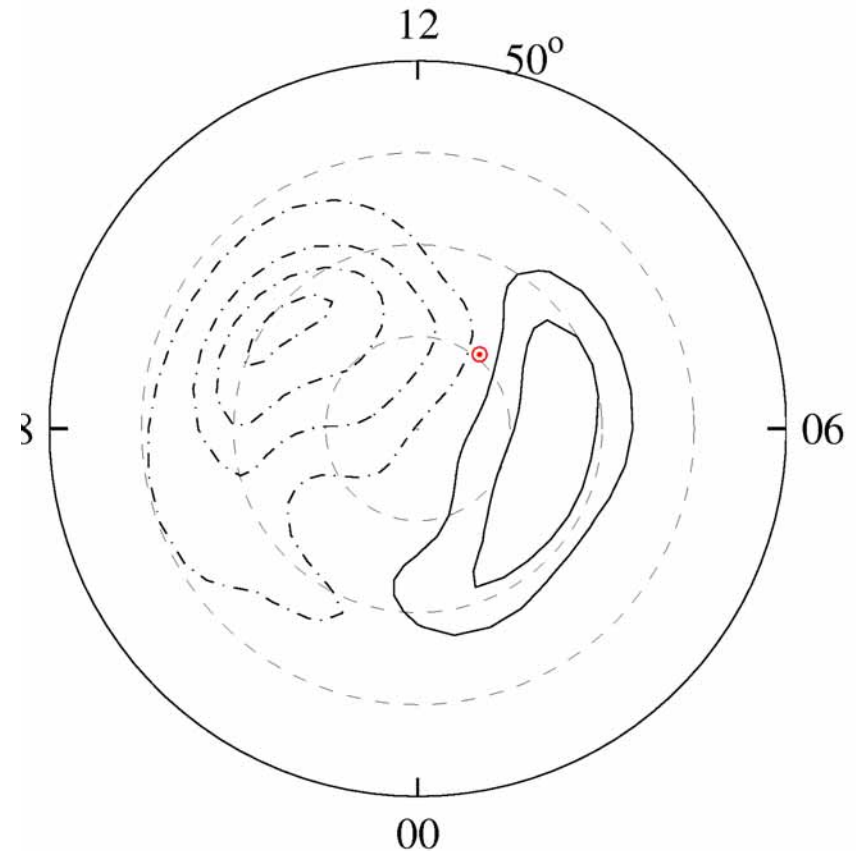
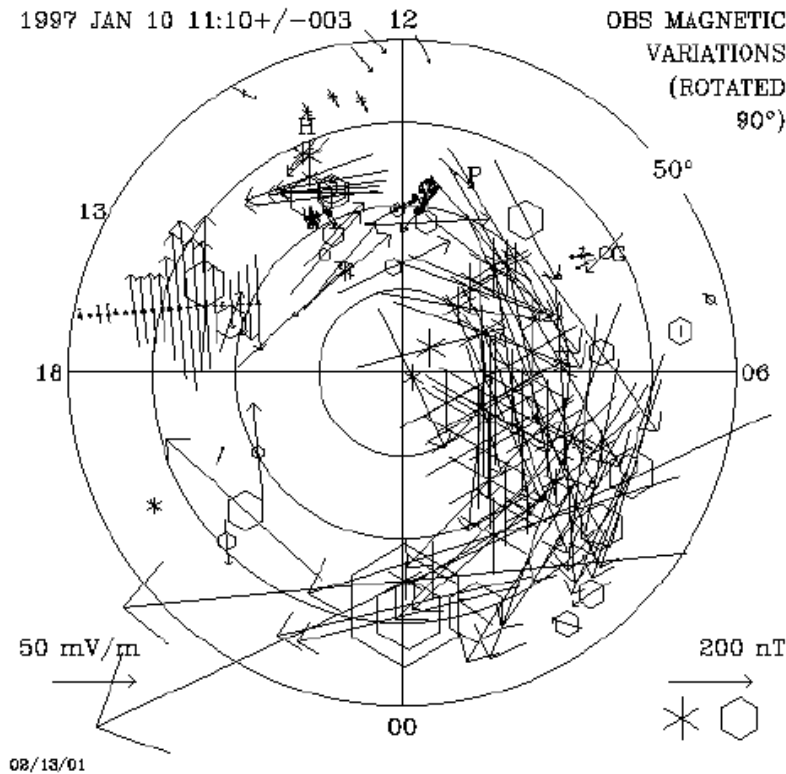
Adaptive Covariance Estimation Using Maximum likelihood Method

[Dee 1995; Dee and da Saliva 1999]

$$x_a = x_b + \mathbf{K}(y - \mathbf{H}x_b)$$

$$\mathbf{K} = \mathbf{P}_b(\alpha) \mathbf{H}^T \left[ \mathbf{H} \mathbf{P}_b(\alpha) \mathbf{H}^T + \mathbf{R} \right]^{-1}$$

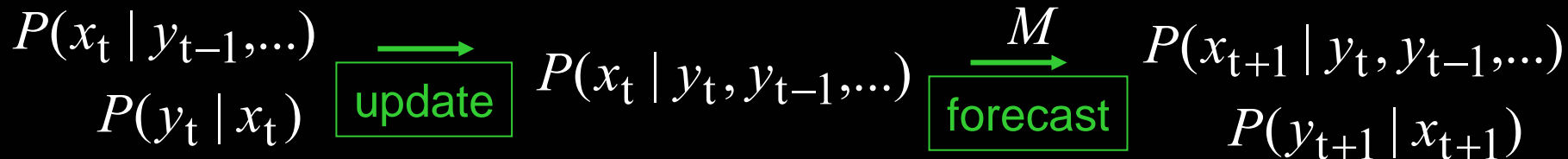
## observations



[Matsuo et al., 2002; 2005]

# Use of Dynamics

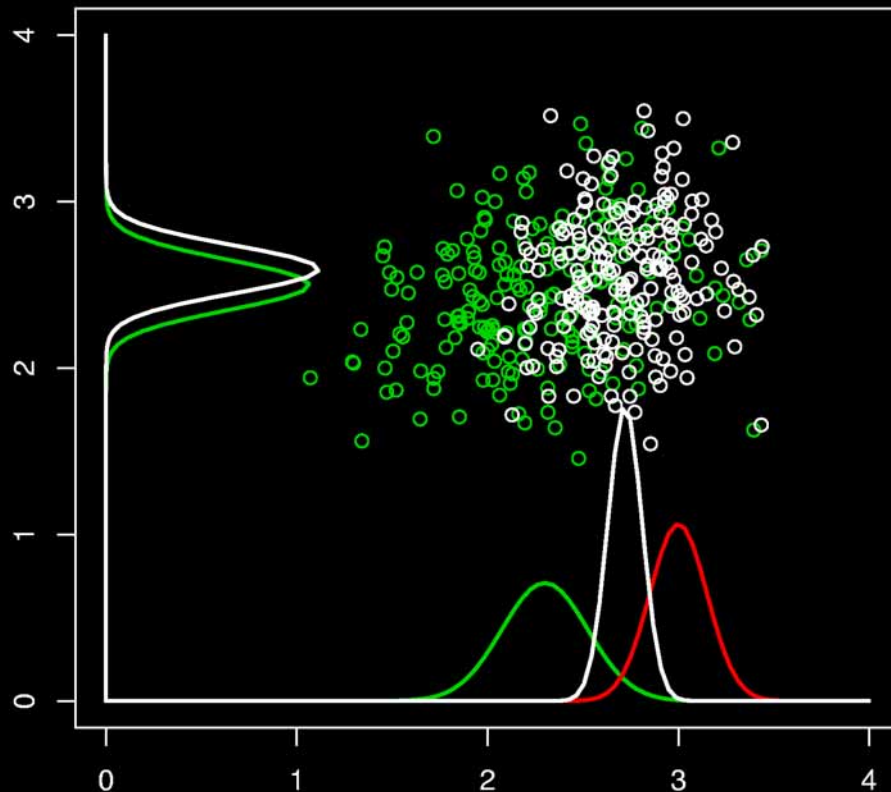
$$x_{t+1} = M(x_t)$$



## Ensemble Kalman Filter

Let's work with samples!

Challenge posed by  
the size of the covariance  
matrix ( $10^{12} - 10^{16}$ )

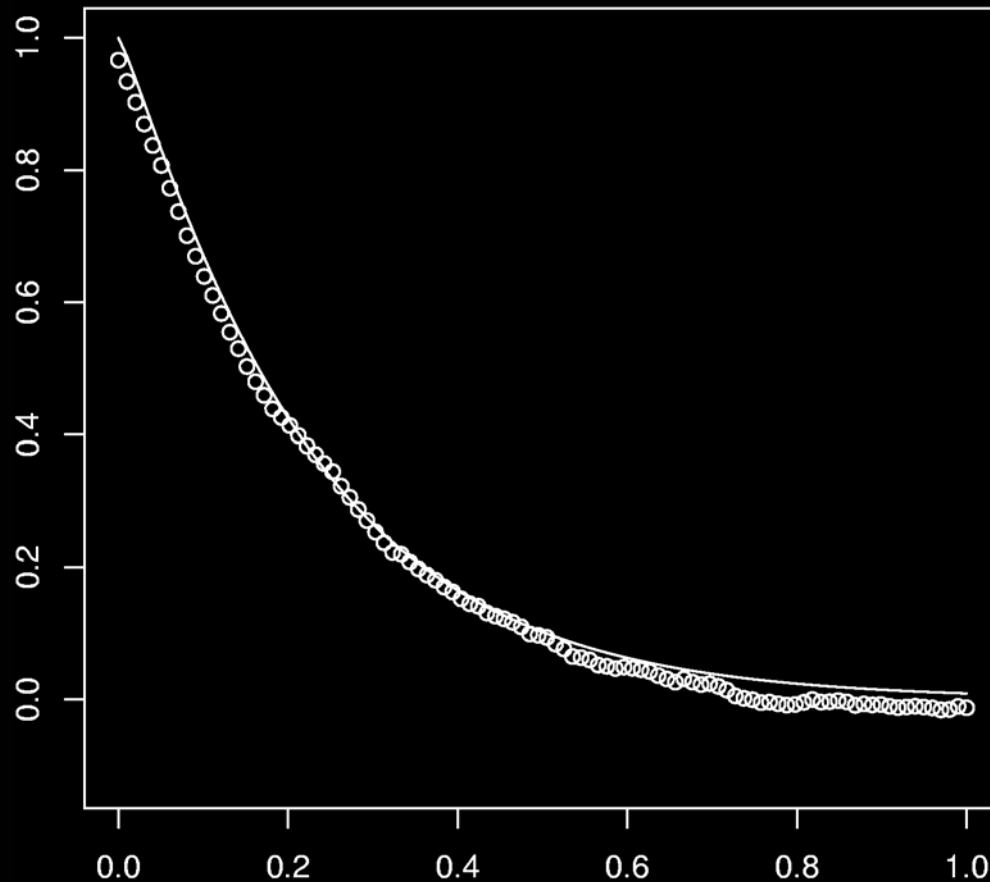




# Issues with sampling error

[e.g., Furrer and Bengtsson, 2005 ]

spurious correlations in the area of large lag distance.



# EnKF v.s. 3D-Var comparisons [Caya et al., 2005]

$$x_a = x_f + \underline{\mathbf{K}(y - \mathbf{H}x_f)}$$

EnKF

T (K) 850 (hPa) 146827 days 0 sec



With Localization

T (K) 850 (hPa) 146827 days 0 sec



3D-Var

T (K) 850 (hPa) 2003-01-01 00:00:00



**Issue with sampling error:** covariance localization (tapering) is necessary to remove spurious correlations in the area far from observation location.

# Summary

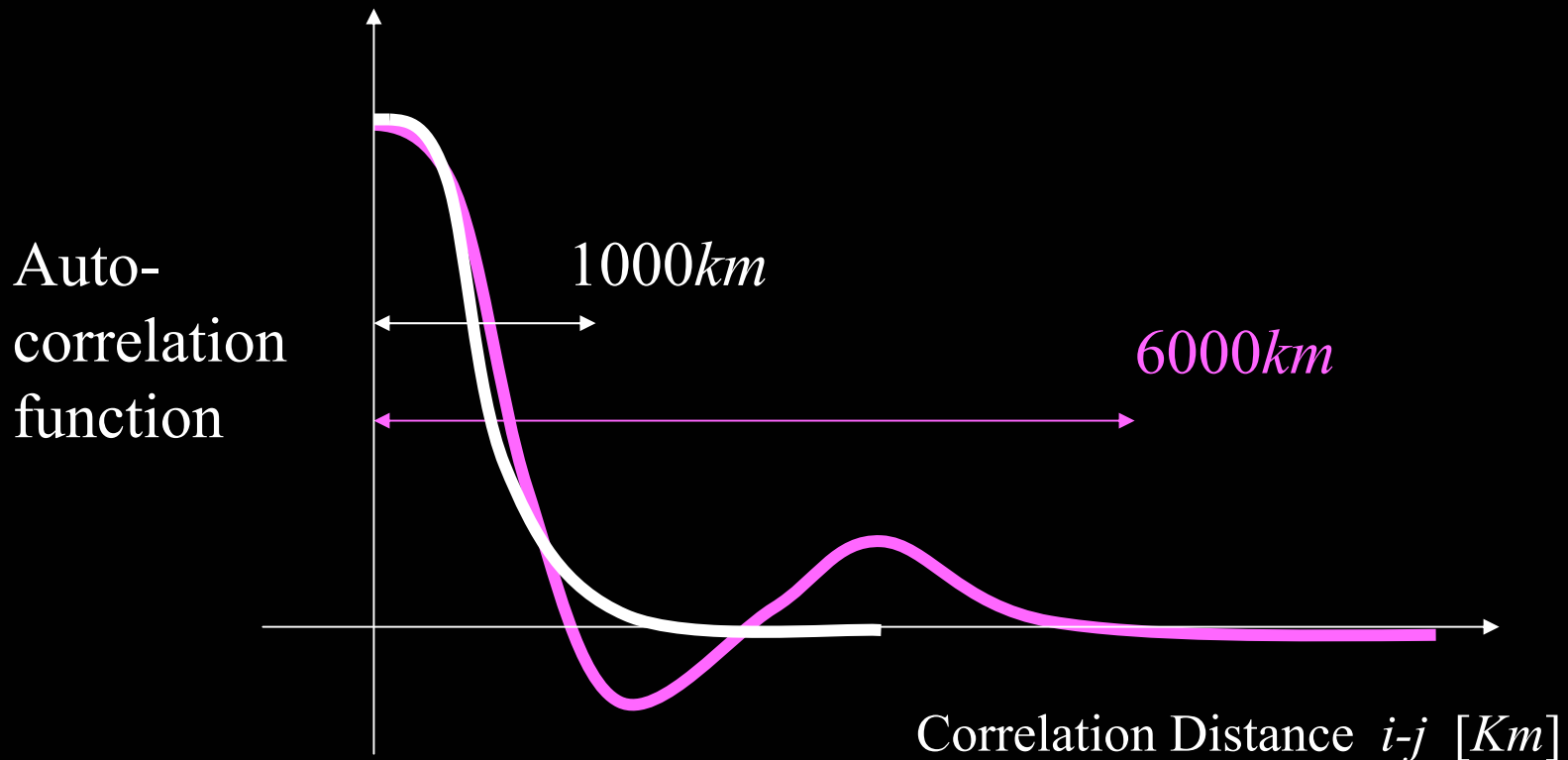
- Bayesian statistics as an overarching framework.
- By confronting a model with observations via first/second moment statistics, data assimilation
  - *improves the state estimation.*
  - *provides a means to evaluate the quality of the model and the value of observations.*
- Inhomogeneous and anisotropic covariance.
- Ensemble Kalman Filter does not require linearization of forward operator (H) and model (M), and has an advantage in capturing flow-dependent covariance structure.
  - *See <http://www.image.ucar.edu/DAReS/DART>*

# Challenges and Future

- Observation is still sparse... (blessing?!)
- Dissipative system and strongly forced system in comparison with meteorological and oceanic systems.  
(forcing prediction is key to forecasting)
- Observing system design analysis or adaptive observation [e.g., Bishop et al., 2001]. (feedback to the design of observational campaigns )

# Why is data assimilation in a data sparse region challenging?

Large Correlation Distance  
Inhomogeneous & Anisotropic



# Adaptive covariance estimation using maximum likelihood method

## OI analysis: optimal estimation of $\alpha$

$$\alpha_a = \mathbf{K}^{OI} \mathbf{y}', \quad \text{where } \mathbf{y}' = \mathbf{y}_o - \mathbf{H}(\mathbf{x}_b)$$
$$\mathbf{K}^{OI} = [(\mathbf{EOF})^T \mathbf{R}^{-1} \mathbf{EOF} + \mathbf{P}_b^{-1}]^{-1} (\mathbf{EOF})^T \mathbf{R}^{-1}$$

Background error covariance:

$$\mathbf{P}_b = \langle \alpha \cdot \alpha^T \rangle$$
$$\approx \text{diag}(\mathbf{P}_b)_{\nu\nu} \approx \zeta_1 \nu^{-\zeta_2} \quad \nu = 1, \dots, 11$$

Observational error covariance:

$$\mathbf{R} \approx \text{diag}(\mathbf{R}) \approx f(\zeta_3, \zeta_4)$$

## Maximum-likelihood method: optimal estimation of $\zeta$

Innovation covariance:

[Dee 1995; Dee and da Silva, 1999]

$$\langle \mathbf{y}' \cdot \mathbf{y}'^T \rangle = \mathbf{R} + \mathbf{EOF} \mathbf{P}_b (\mathbf{EOF})^T \approx \mathbf{S}(\zeta) \quad \{\zeta_k \mid k = 1 \rightarrow 4\}$$

Cost function:

$$J(\zeta) = \log \det \mathbf{S}(\zeta) + \mathbf{y}'^T \mathbf{S}^{-1}(\zeta) \mathbf{y}'$$