# Basics of Data Assimilation



▷ Chris Snyder (NCAR)

Why Data Assimilation?

We need initial conditions for NWP

We have limited observations

Simplest: interpolate between observations

▷ objective analysis

# Function Fitting -



▷ field to be estimated

## Function Fitting.



 $\triangleright$  observations at discrete points; no obs error

## Function Fitting



▷ observations at discrete points; no obs error

## Function Fitting.



▷ cubic-spline interpolation

## Function Fitting.



▷ interpolation gives good estimate of scales resolved by obs

## Function Fitting



▷ imperfect observations (std. dev. 0.2)

## Function Fitting.



▷ small scales lost in observational noise

# Why Data Assimilation? (II) \_\_\_\_

Real observing networks

- ▷ observations are imperfect
- obs may be *indirect:* measured variables differ from forecast-model variables
- $\triangleright~$  inhomogeneous in space and time  $\rightarrow$  gaps

Use other information to fill gaps

- ▷ e.g., previous forecast (*background* or *first guess*)
- ▷ How to combine with observations?

Observations (and forecasts!) have errors

 $\triangleright$  How to account for this?

Indirect observations

▷ How to spread information from observed variables to model variables?

# A Basic Example \_\_\_\_

#### Two instruments, both measure a quantity $\boldsymbol{x}$

 $\triangleright$  observations have form

$$\begin{array}{rcl} y_1 &=& x + \epsilon_1, \\ y_2 &=& x + \epsilon_2 \end{array}$$

 $\triangleright~$  observation errors  $\epsilon_1,~\epsilon_2$  are random but we know something about their statistics

Given observations  $y_1$  and  $y_2$ , estimate x

▷ since estimate is imperfect, want estimate of its accuracy as well.

# The Minimum Variance Solution \_

Acknowledge probabilistic aspects of problem

- $\triangleright$  obs errors are random
- $\triangleright$  can't know true value of x; can only evaluate expected errors

Seek unbiased  $x_a$  with minimum expected squared error

$$\triangleright \quad x_a - x \text{ is error}$$

▷ require

$$E(x_a - x) = 0$$

⊳ minimize

$$E\left((x_a - x)^2\right) = \operatorname{Var}\left(x_a - x\right)$$

## The Minimum Variance Solution (cont.) \_

Spse  $x_a$  is linear combination of observations,

$$x_a = \alpha_1 y_1 + \alpha_2 y_2 = (\alpha_1 + \alpha_2) x + \alpha_1 \epsilon_1 + \alpha_2 \epsilon_2$$

Then,

$$E(x_a - x) = 0 \quad \Rightarrow \quad \alpha_1 + \alpha_2 = 1 \quad \text{if } E(\epsilon_i) = 0$$

Next,

$$E((x_{a} - x)^{2}) = E\left[\{\alpha_{1}(y_{1} - x) + \alpha_{2}(y_{2} - x)\}^{2}\right]$$
  
=  $E\left[\alpha_{1}^{2}\epsilon_{1}^{2} + \alpha_{2}^{2}\epsilon_{2}^{2}\right]$  if  $E(\epsilon_{1}\epsilon_{2}) = 0$   
=  $\alpha_{1}^{2}\sigma_{1}^{2} + (1 - \alpha_{1})^{2}\sigma_{2}^{2}$ ,

where  $\sigma_i^2 = \text{Var}\left(\epsilon_i\right)$ 

# The Minimum Variance Solution (cont.) \_

Minimizing w.r.t.  $\alpha {\rm 's}$  then gives

$$\alpha_1 = \sigma_2^2 / (\sigma_1^2 + \sigma_2^2), \quad \alpha_2 = \sigma_1^2 / (\sigma_1^2 + \sigma_2^2)$$

Properties

- > solution depends only on observation-error variances
- ▷ observations with large errors receive small weight
- ▷ expected squared error

$$E((x_a - x)^2) = \alpha_1^2 \sigma_1^2 + \alpha_2^2 \sigma_2^2 = \sigma_1^2 \sigma_2^2 / (\sigma_1^2 + \sigma_2^2)$$

 $\triangleright$  error of  $x_a$  is (on average) smaller than error of either observation

Best linear unbiased estimator (BLUE)

# The Bayesian Solution.

#### True state is unknown

- $\triangleright$   $\;$  observations, models both have random errors
- ▷ wish to calculate  $p(x|y_1, y_2)$ , probability density of x given  $y_1$  and  $y_2$ . Also known as the *posterior* density.

#### Bayes rule

- $\triangleright \quad p(x|y_1, y_2) \propto p(y_2|x, y_1)p(x|y_1)$
- $\triangleright$  if obs errors independent,  $p(y_2|x,y_1) = p(y_2|x)$

## Bayes Illustrated

 $\triangleright p(x|y_1)$  (blue)



# Bayes Illustrated (cont.)

 $\triangleright p(x|y_1)$  (blue)

$$\triangleright \quad p(y|x) \text{ for } y = 0.75 \text{ (red)}$$



## Bayes Illustrated (cont.)

- $\triangleright p(x|y_1)$  (blue)
- $\triangleright \quad p(y|x) \text{ for } y = 0.75 \text{ (red)}$
- $\triangleright \quad p(x|y) \propto p(y|x)p(x) \text{ (black)}$



## Bayes Solution for Gaussian Errors

Suppose densities are Gaussian

 $\triangleright~$  e.g., of the form  $p(x)=c\exp(-\frac{1}{2}(x-\bar{x})^2/\sigma^2)$ 

Bayes rule involves multiplication of densities

- products of Gaussians amount to adding exponents
- ▷ for two observations, posterior density proportional to  $\exp[-\frac{1}{2}J(x)] = \exp[-\frac{1}{2}(x-y_1)^2/\sigma_1^2 - \frac{1}{2}(x-y_2)^2/\sigma_2^2]$

Yields same answer as minimum-variance approach

- ▷ left as exercise
- $\triangleright \quad \text{hint: write } \exp[-\frac{1}{2}J(x)] = \exp[-\frac{1}{2}(x-x_a)^2/\sigma_a^2]$

# Background Forecasts \_

Often, have prior information in addition to observations

- $\triangleright$  climatology
- ▷ forecast from earlier time (first guess, background)

In NWP, background forecasts often as accurate as observations

▷ propagates information from previous obs forward in time

Need to estimate error statistics for DA

## Many variables \_\_\_\_

Notation

- x = continuous state of system *projected* onto discrete basis, e.g. grid-point values or Fourier coefficients
- $\triangleright$  y = all observations concatenated into single vector

# Observation Operator and Observation Errors \_\_\_\_

### Assimilation requires relation of observations to **x**

 $\triangleright$  observation model

 $\mathbf{y} = H(\mathbf{x}) + \epsilon$ 

- $\triangleright$  *H* is observation (or *forward*) operator mapping state onto obs
- e.g., interpolation for in-situ obs of state variables or radiative-transfer integrals for remotely sensed radiances
- $\triangleright$  *H* may be nonlinear; errors may not be additive

#### Three sources of observation error

- ▷ measurement error: noise in instrument or uncertainty in location
- $\triangleright~$  errors of representativeness: obs influenced by scales not represented in discrete basis of  ${\bf x}$
- ▷ observation operator: relation of obs to **x** may be incorrectly specified
- $\triangleright \quad \epsilon \text{ is sum of all three effects}$

# Multivariate Gaussians.

Probability density function

$$\triangleright \quad p(\mathbf{x}) = c \exp\left((\mathbf{x} - \overline{\mathbf{x}})^T \mathbf{P}^{-1} (\mathbf{x} - \overline{\mathbf{x}})\right)$$

- $\triangleright$  completely specifed by mean  $\overline{\mathbf{x}}$  and covaraince matrix  $\boldsymbol{P}$
- $\triangleright$  **P** describes statistical relations between elements of **x**

### Standard assumptions

- $\triangleright$  Gaussian forecast errors:  $\mathbf{x} \sim N(\mathbf{x}_b, \mathbf{B})$
- $\triangleright~$  Gaussian obs errors, linear obs operator:  $\mathbf{y}=\mathbf{H}\mathbf{x}+\epsilon,$  with  $\epsilon\sim N(\mathbf{0},\mathbf{R})$

## Bayes rule for Gaussians -

- Products of Gaussian yield Gaussians, so posterior (or analysis) pdf is also Gaussian. Thus, need formulas for its mean and covariance
- ▷ analysis equations:

 $\mathbf{x}_a = (\mathbf{I} - \mathbf{K}\mathbf{H})\mathbf{x}_b + \mathbf{K}\mathbf{y}_o \quad ; \quad \mathbf{A} = (\mathbf{I} - \mathbf{K}\mathbf{H})\mathbf{B},$ 

▷ Kalman gain

$$\mathbf{K} = \mathbf{B}\mathbf{H}^{T}(\mathbf{H}\mathbf{B}\mathbf{H}^{T} + \mathbf{R})^{-1}$$
 [c.f.,  $\alpha_{2} = \sigma_{1}^{2}/(\sigma_{1}^{2} + \sigma_{2}^{2})$ ]

▷ equivalently, compute  $\mathbf{x}_a$  as minimizer of  $J(\mathbf{x}) = (\mathbf{x} - \mathbf{x}_b)^T \mathbf{B}^{-1} (\mathbf{x} - \mathbf{x}_b) + (\mathbf{y}_o - \mathbf{H}\mathbf{x})^T \mathbf{R}^{-1} (\mathbf{y}_o - \mathbf{H}\mathbf{x})^T$ 

## Importance of Covariances.

- ▷ 2D:  $\mathbf{x} = (x_1, x_2)$
- ▷ forecast/prior



## Importance of Covariances

- $\triangleright$  observation,  $y = x_1^t + \text{Gaussian noise} = 1.4$
- $\triangleright$  observation likelihood independent of  $x_2$



## Importance of Covariances

- $\triangleright$  analysis/posterior,  $[x_1, x_2|y]$
- $\triangleright$  Cov( $x_1, x_2$ ) provides information on unobserved variable



Practical Considerations

Dimension of state vector may exceed  $10^8$ 

Implement Bayes rule in n dimensions

- $\triangleright$  pdfs are functions of n variables
- $\triangleright~$  even discretizing each variable with 10 "points" requires  $10^n$  total degrees of freedom

Implement Gaussian update

- $\triangleright$  densities specified by their mean and covariances
- $\triangleright$  requires *n* floating point numbers for mean, n(n-1)/2 for covariance
- ▷ moreover, background error covariances poorly known

Necessity for approximation and simplification

# Additional Matters \_\_\_\_

Quality control

- $\triangleright$   $\;$  identify observations with gross errors
- ▷ crucial in operational schemes

### Computational issues

- ▷ tractable covariance models
- $\triangleright$  minimization

#### Propagation of information in time

- ▷ role of dynamical systems
- $\triangleright$  sources and evolution of forecast error

Nonlinearity and non-Gaussianity