

# **NON-GAUSSIAN DATA ASSIMILATION**

**STEVEN FLETCHER**

**COOPERATIVE INSTITUTE FOR RESEARCH IN THE  
ATMOSPHERE,  
COLORADO STATE UNIVERSITY**

**Lecture for the Workshop on Applications of Remotely Sensed  
Observations in Data Assimilation,  
University of Maryland, August 3<sup>rd</sup> 2007**



# Overview of Lecture

- 1) **Descriptive Univariate Statistics**
- 2) **Multivariate Distribution Theory**
- 3) **Bayes Theorem and Gaussian Data Assimilation**
- 4) **Lognormal Distribution, Univariate and Multivariate theory**
- 5) **Lognormal Data Assimilation**
- 6) **Hybrid Data Assimilation**
- 7) **4D Lognormal Data Assimilation**

# **Who is Steven Fletcher?**

**B.Sc. (HONS) Mathematics and Statistics, 1998**

**M.Sc. Numerical Solutions to Differential Equations, 1999 (Prof. Nancy Nichols)**

**Ph.D. Numerical Weather Prediction, 2004  
(Prof. Nancy Nichols, Dr. Ian Roulstone)**

**All from the Department of Mathematics  
(Statistics), University of Reading, U. K.**

**Post-Doc, CIRA/CSU 2004 – 2006, (Dr. Milija  
Zupanski)**

**Research Scientist II 2006 - present**

## Univariate Descriptive Statistics

**We shall only be considering the continuous probability distributions for this lecture. The three descriptive statistics for cts. distributions are defined as**

**The cts. Mean, which is also the first moment of the distribution, it is also the minimum variance estimator and is given by**

$$\mu = E(X = x) = \int_a^b xf(x)dx$$
$$x \in [a, b]$$

**The next statistic is the median. This is the unbiased estimator. For cts. distributions this is given by**

$x_{med}$  is the values such that

$$\int_a^{x_{med}} f(z) dz = \frac{1}{2}$$

**The third and final descriptive statistic is the mode. This is the maximum likelihood estimator and is defined by**

$x_{mod}$  is the value of  $x$  such that

$$\left. \frac{df}{dx} \right|_{x=x_{mod}} = 0$$

## Univariate normal (Gaussian) distribution

The univariate normal distribution is regularly used to approximate other distributions through the central limit theorem, we do not discuss that here, but it has two definitions. The Standard Normal is defined by

$$f(z) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left\{-\frac{z^2}{2}\right\} \quad z \in (-\infty, \infty)$$

Whilst the biased version is given by

$$f(x) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left\{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}\right\} \quad x \in (-\infty, \infty)$$
$$\mu = E(x), \quad \sigma^2 = E(x^2) - (E(x))^2$$

## Multivariate Distribution Theory

Question: How do the three statistic operators transfer to multivariate theory? **NOT VERY WELL!!!**

MEAN: Intuition says

$$E(\mathbf{x}) = \int_a^b \int_c^d \dots \int_e^f \mathbf{x} f(\mathbf{x}) d\mathbf{x}$$

$$x_1 \in (a, b), x_2 \in (c, d), \dots, x_N \in (e, f)$$

**NOT DEFINED IN MATHEMATICS!!!!**

Therefore the expectation of a random vector is the vector of random expectations

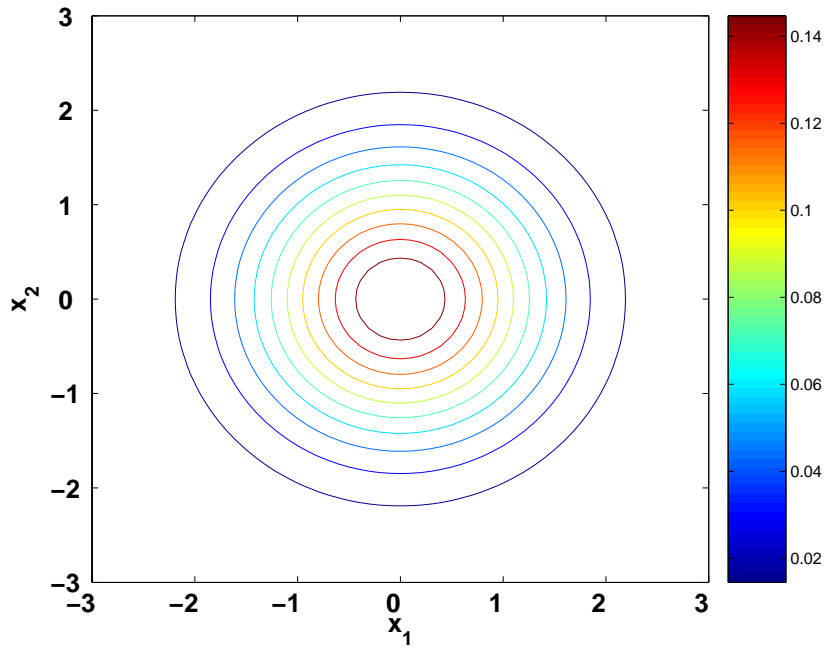
## MULTIVARIATE MEDIANS

This is the worst of the three statistics. Although it is *unbiased* it is also non-unique, even for the multivariate normal distribution. The multivariate definition is

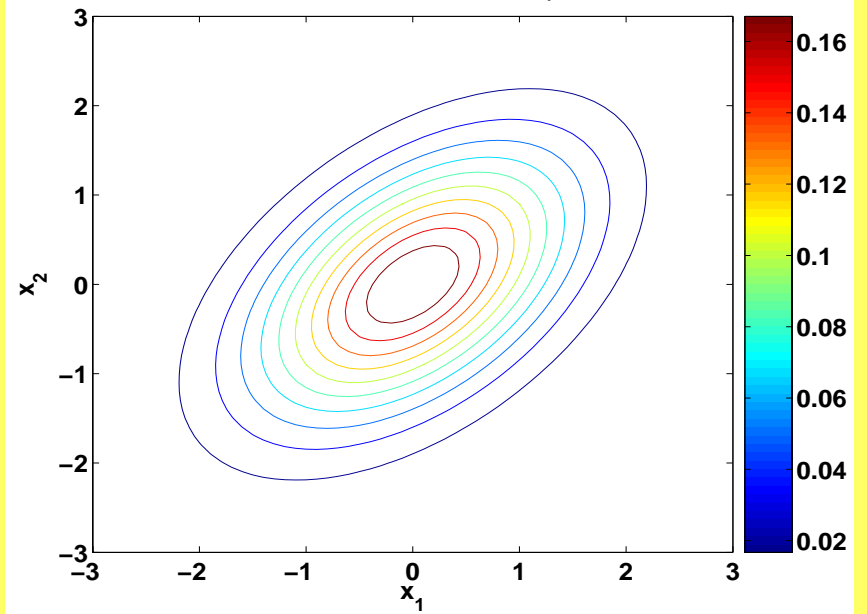
$$\mathbf{x}_{med} = \begin{pmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(N)} \end{pmatrix} = \int_a^{x^{(1)}} \int_c^{x^{(2)}} \cdots \int_e^{x^{(N)}} f(\mathbf{x}) d\mathbf{x}$$



**BIVARIATE UNIT NORMAL,  $\rho=0$**



**BIVARIATE UNIT NORMAL,  $\rho=0.5$**



**Is all hope lost? No**

**Our saviour is the mode!!!!**

**There is a simple multivariate extension of the definition of the mode from the univariate case to the multivariate.**

$\mathbf{x}_{\text{mod}}$  is the value of  $\mathbf{x}$  such that

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\mathbf{x}_{\text{mod}}} = \mathbf{0}$$

The Multivariate Normal Distribution is given by

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{N}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}$$
$$\mathbf{x} \in \mathfrak{R}^N$$

# BAYES THEOREM AND GAUSSIAN DATA ASSIMILATION

The basis of 3D VAR comes from Lorenc (1986) where he uses Bayes theorem

$$P(A|B) = P(B|A)P(A)$$

with the event A presenting the effect background knowledge and the event B the effects of the observations given the uncertainties with the background states. The right hand side of the equation above becomes the *ANALYSIS DISTRIBUTION* whose mean, mode or median we seek as the *ANALYSIS STATE OR TRUE SOLUTION*. For 3D VAR this is written as

$$P(\mathbf{x} = \mathbf{x}_a | \mathbf{y} = \mathbf{y}_t) = P(\mathbf{y} = \mathbf{y}_t | \mathbf{x} = \mathbf{x}_a)P(\mathbf{x} = \mathbf{x}_a)$$

$$P(\mathbf{x} - \mathbf{x}_b) \propto \exp\left\{-\frac{1}{2}(\mathbf{x} - \mathbf{x}_b)^T \mathbf{B}^{-1}(\mathbf{x} - \mathbf{x}_b)\right\}$$

$$P(\mathbf{y} = \mathbf{y}_t | \mathbf{x} = \mathbf{x}_a) \propto \exp\left\{-\frac{1}{2}(\mathbf{y} - \mathbf{h}(\mathbf{x}))^T \mathbf{R}^{-1}(\mathbf{y} - \mathbf{h}(\mathbf{x}))\right\}$$

**Where we have assumed that the background and observational errors are independent. Therefore multiplying the two distributions above and taking the negative logarithm we obtain the standard 3D VAR cost function**

$$J(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \mathbf{x}_b)^T \mathbf{B}^{-1}(\mathbf{x} - \mathbf{x}_b) + \frac{1}{2}(\mathbf{y} - \mathbf{h}(\mathbf{x}))^T \mathbf{R}^{-1}(\mathbf{y} - \mathbf{h}(\mathbf{x}))$$

# LOGNORMAL DISTRIBUTION: UNIVARIATE AND MULTIVARIATE THEORY

The lognormal distribution is for a class of positive definite random variables, i.e.

$$x \geq 0$$

The pdf is given by

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma x} \exp\left\{-\frac{1}{2} \frac{(\ln x - \mu)^2}{\sigma^2}\right\}$$
$$\mu = E(\ln x), \quad \sigma^2 = E((\ln x)^2) - (E(\ln x))^2$$

# PROPERTIES OF THE LOGNORMAL DISTRIBUTION UNIVARIATE

PROPERTY 1: IF  $X \sim \text{LN}(\mu, \sigma^2)$  THEN  $\ln X$  IS DISTRIBUTED  $N(\mu, \sigma^2)$

PROPERTY 2: THE MODE OF  $X$  IS  $\exp\{\mu - \sigma^2\}$

PROPERTY 3: THE MEDIAN OF  $X$  IS  $\exp\{\mu\}$

PROPERTY 4: THE MEAN OF  $X$  IS  $\exp\left\{\mu + \frac{\sigma^2}{2}\right\}$

PROPERTY 5: THE LOGNORMAL DISTRIBUTION CAN NOT BE UNIQUELY DETERMINED BY ITS MOMENTS!!!

## MULTIVARIATE LOGNORMAL DISTRIBUTION

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{N}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \left( \prod_{i=1}^N \left( \frac{1}{x_i} \right) \right) \times \exp \left\{ -\frac{1}{2} (\ln \mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\ln \mathbf{x} - \boldsymbol{\mu}) \right\}$$

$$x_{mea}(i) = E(x_i) = \exp \left\{ \mu_i + \frac{\sigma^2}{2} \right\}, i = 1, 2, \dots, N$$

$$\mathbf{x}_{med} = \exp\{\boldsymbol{\mu}\},$$

$$\mathbf{x}_{mod} = \exp\{\boldsymbol{\mu} - \boldsymbol{\Sigma}^T \mathbf{1}\}, \mathbf{1}^T = (1 \quad 1 \quad \dots \quad 1)$$



# PROPERTIES OF THE LOGNORMAL DISTRIBUTION

## MULTIVARIATE

PROPERTY 1: MEDIAN IS NON-UNIQUE

PROPERTY 2: MOMENTS DO NOT DETERMINE THE DISTRIBUTION UNIQUELY

PROPERTY 3: MEAN IS INDEPENDENT OF COVARIANCE AND IS UNDOUNDED WITH RESPECT TO THE VARIANCE

PROPERTY 4: MODE IS BOUNDED AND FINITE WITH RESPECT TO THE VARIANCE.

PROPERTY 5: MODE IS UNIQUE!!!

# LOGNORMAL DATA ASSIMILATION

## FLETCHER AND ZUPANSKI (2006a)

BY FOLLOWING LORENC (1986) AND COHN (1997) WE CAN DEFINE A COST FUNCTION FOR LOGNORMAL OBSERVATION ERRORS.

$$J(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \mathbf{x}_b)^T \mathbf{B}^{-1}(\mathbf{x} - \mathbf{x}_b) + \frac{1}{2}(\ln \mathbf{y} - \ln \mathbf{h}(\mathbf{x}))^T \mathbf{R}^{-1}(\ln \mathbf{y} - \ln \mathbf{h}(\mathbf{x})) + \sum_{i=1}^{N_o} (\ln y_i - \ln h_i(\mathbf{x})) \quad (1)$$

WHERE THE ERRORS ARE DEFINED BY

$$\boldsymbol{\varepsilon}_b = \mathbf{x} - \mathbf{x}_b \propto N(\mathbf{0}, \mathbf{B}), \quad \text{AND} \quad \boldsymbol{\varepsilon}_o = \frac{\mathbf{y}}{h(\mathbf{x})} \propto LN(\mathbf{0}, \mathbf{R}) \quad (2)$$

# MOTIVATION

**UNDERLYING ASSUMPTION IN VARIATIONAL, KALMAN FILTER AND ENSEMBLE DATA ASSIMILATION IS THAT THE VARIABLES, OBSERVATIONS AND HENCE ERRORS ARE NORMAL (GAUSSIAN) DISTRIBUTED, LORENC (1986), KALMAN (1960), EVENSEN (1994).**

**QUESTION 1: IS THIS TRUE FOR ALL SCALES IN THE ATMOSPHERE?**

**QUESTION 2: IS THIS TRUE FOR ALL TYPES OF OBSERVATIONS?**

## Question 1

**Synoptic scale:** Are all the synoptic variables Gaussian? **NO**

**Humidity:** To assimilate this variable we have to use the logarithm of the variable (Polavarapu *et al* 2005). This indicates that this variable is **LOGNORMAL!!**

**Wind component:** Combined with the moisture flux then this variable is showing sign that the wind components may be **LOGNORMAL!!** (Raymond 1997)

**Comment:** These variables are both positive definite!!

## Question 1

**Meso-scale:** Are all meso-scale variables Gaussian? **NO**

In the paper by Miles *et al.* (2000) there is a large summary of cloud variables which are not Gaussian, specifically **LOGNORMAL** and **GAMMA**. As early as the 1970s rain and cloud variables had been identified as **LOGNORMAL**, Mielke *et al.* (1977)

## Question 2

Are all observations Gaussian? **NO**

**Direct Observations:** The variables which we have already mentioned are not Gaussian and therefore a direct observation of them is also not.

**Retrievals:** By having to take the logarithm of the humidity to apply 1D Variational data assimilation (1D VAR) then the observation is not Gaussian. It is **LOGNORMAL!!**

**Radiances:** By having to *bias correct* the solution for the VAR scheme to work implies Gaussian assumption has been violated, (Derber and Wu, (1998), Harris and Kelly (2001)).

## Question 2

**Optical Depth:** In Stephens *et al.* 2002 the pdf of this variable is presented, clearly showing a **LOGNORMAL** structure.

**Infra-Red Flux Differences:** From the same paper.

**Cloud base height:** In Sengupta *et al.* 2004 these observations shows signs of a lognormal structure.

**Liquid water path:** Same paper, shows a sharp positive skewness associated with a **LOGNORMAL** distribution with large variance.

**Brightness Temperature:** Work currently being done with the GOES-R sounder (Grosso and Sengupta) shows, when scaled to the main region of observation, then this is **LOGNORMAL**.

# CURRENT TECHNIQUES

CURRENTLY THERE ARE TWO MECHANISM WHICH ARE USED TO OVERCOME THIS PROBLEM.

**TECHNIQUE 1:** TRANSFORMATION. IF  $X \sim \text{LN}(\mu, \sigma^2)$  THEN  $\ln X \sim \text{N}(\mu, \sigma^2)$ . **NOTE: THE MEAN AND VARIANCE ARE UNCHANGED.**

**TECHNIQUE 2:** FORCE THE GAUSSIAN ASSUMPTION AND BIAS CORRECT.

A THIRD TECHNIQUE WHICH IS EMPLOYED IN OTHER FIELDS IS TO USE A **GAUSSIAN SUM FILTER**.

A FOURTH APPROACH IS TO USE A **MAXIMUM ENTROPY CONDITION**.



# PROBLEMS ASSOCIATED WITH CURRENT TECHNIQUES

TRANSFORMATION:



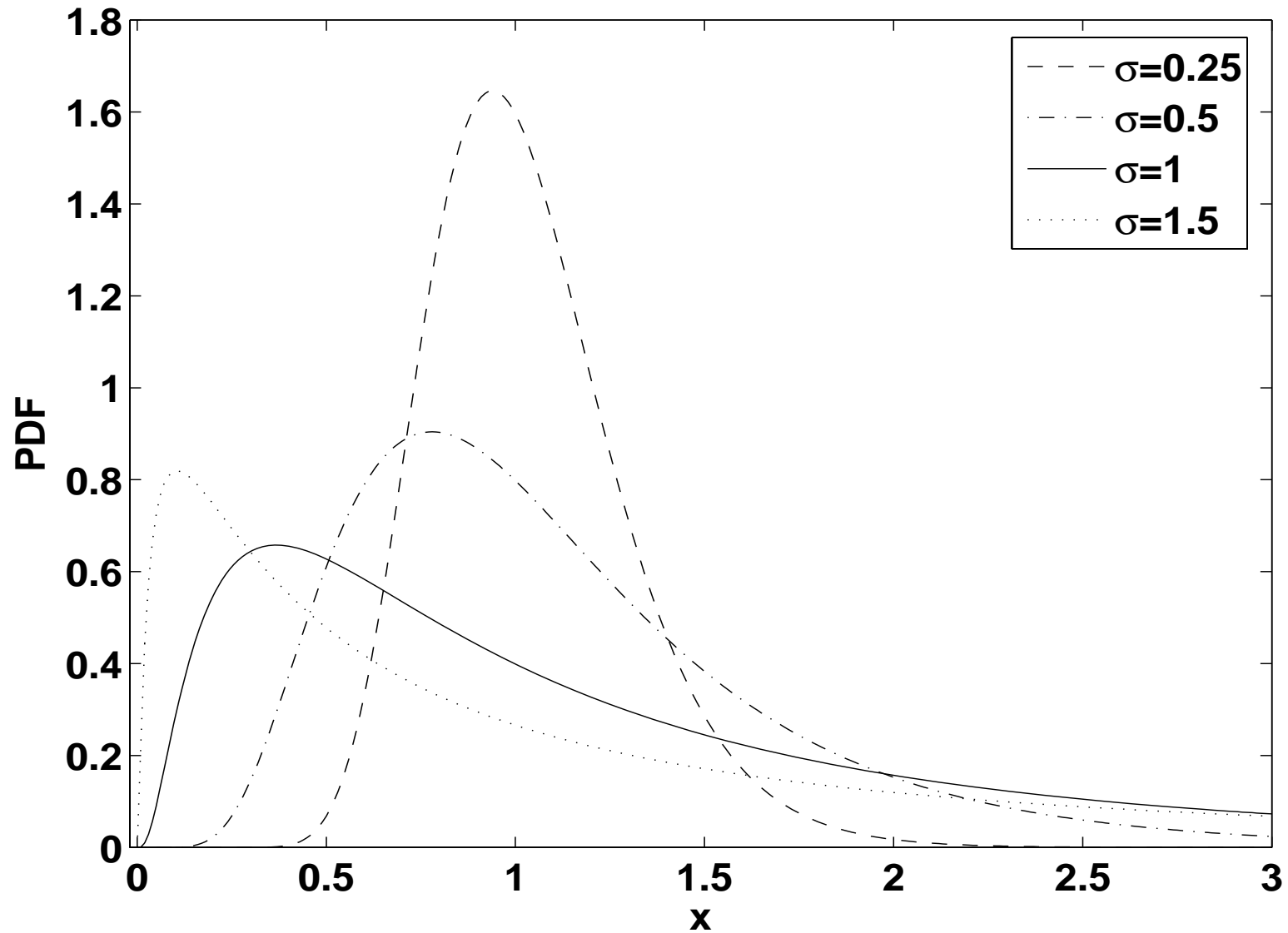
IMPACT 1: THE STATISTIC WHICH IS FOUND IS THE MEDIAN NOT THE MODE NOR THE MEAN.

IMPACT 2: WHEN TRANSFORMING ALL HIGHER ORDER MOMENT INFORMATION IS LOST.

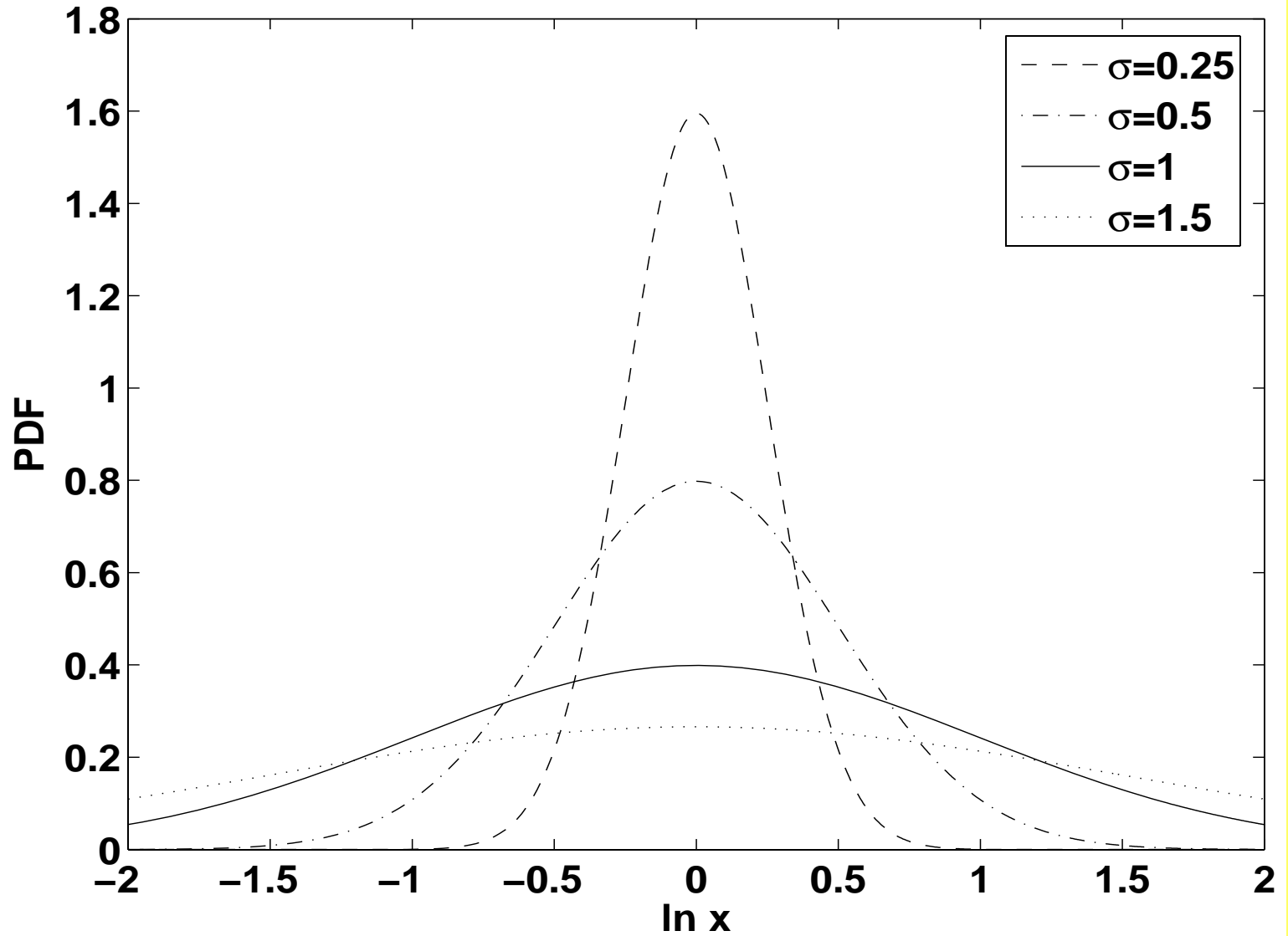
IMPACT 3: THE ANALYSIS STATE FOUND WILL OVER ESTIMATE THE TRUE 'MOST LIKELY STATE'



# PLOT OF LOGNORMAL DISTRIBUTIONS



# PLOT OF TRANSFORMED NORMAL DISTRIBUTIONS



# PROBLEMS ASSOCIATED WITH CURRENT TECHNIQUES

**FORCED GAUSSIAN:**



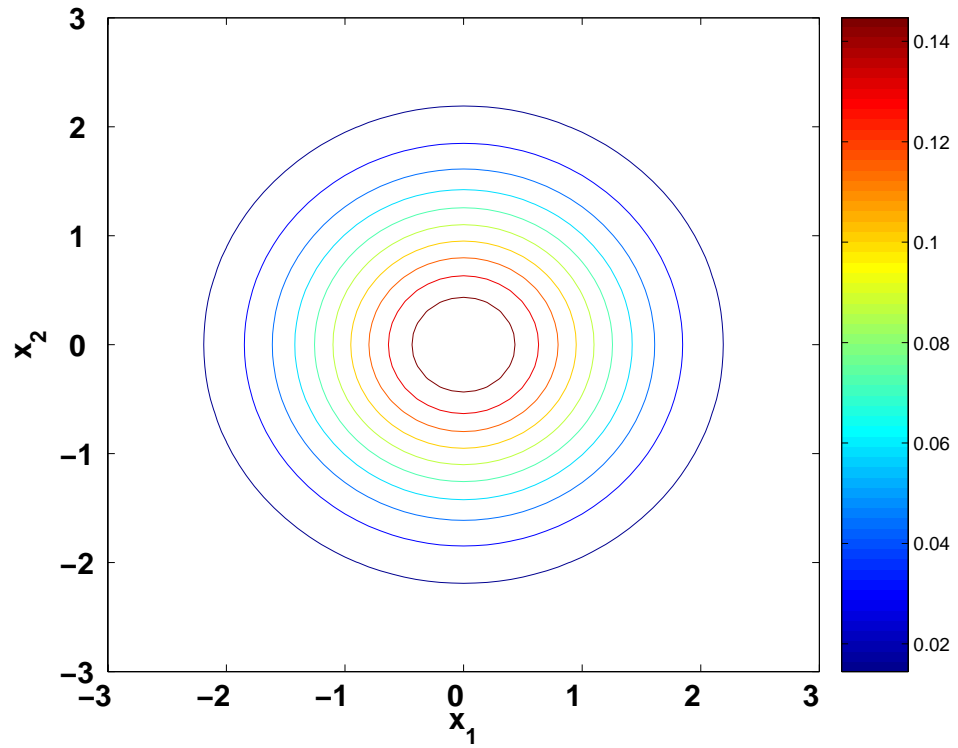
**IMPACT 1: WRONG PROBABILITIES ASSIGNED TO THE  
OUTLIERS.**

**IMPACT 2: PROBABILITIES ASSIGNED TO UNPHYSICAL  
VALUES.**

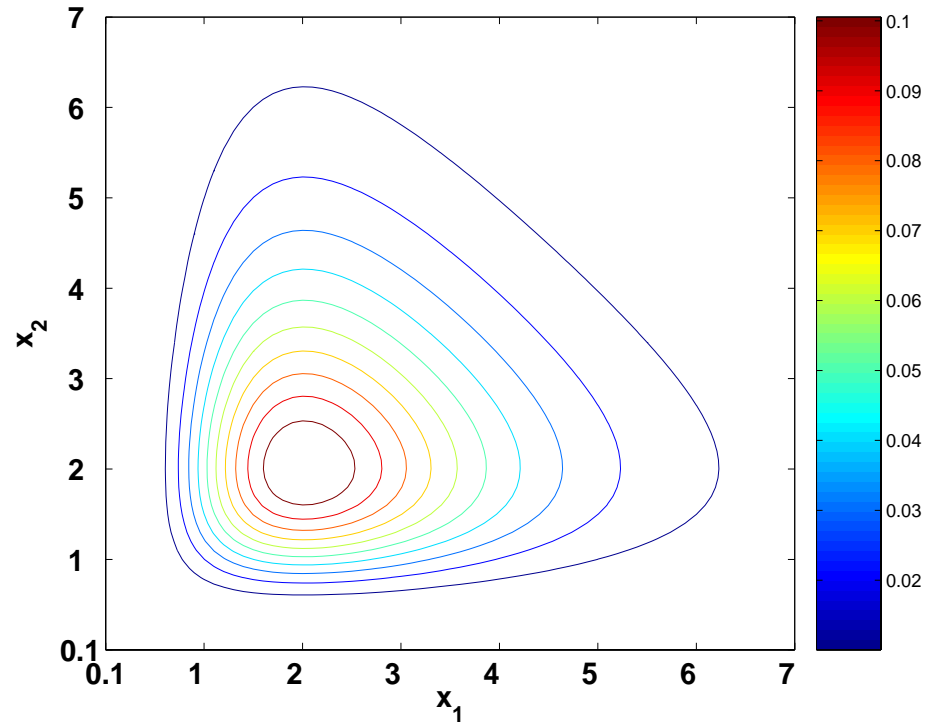
**IMPACT 3: WRONG STATISTICS USED TO  
APPROXIMATE THE VARIABLE'S DISTRIBUTION.**



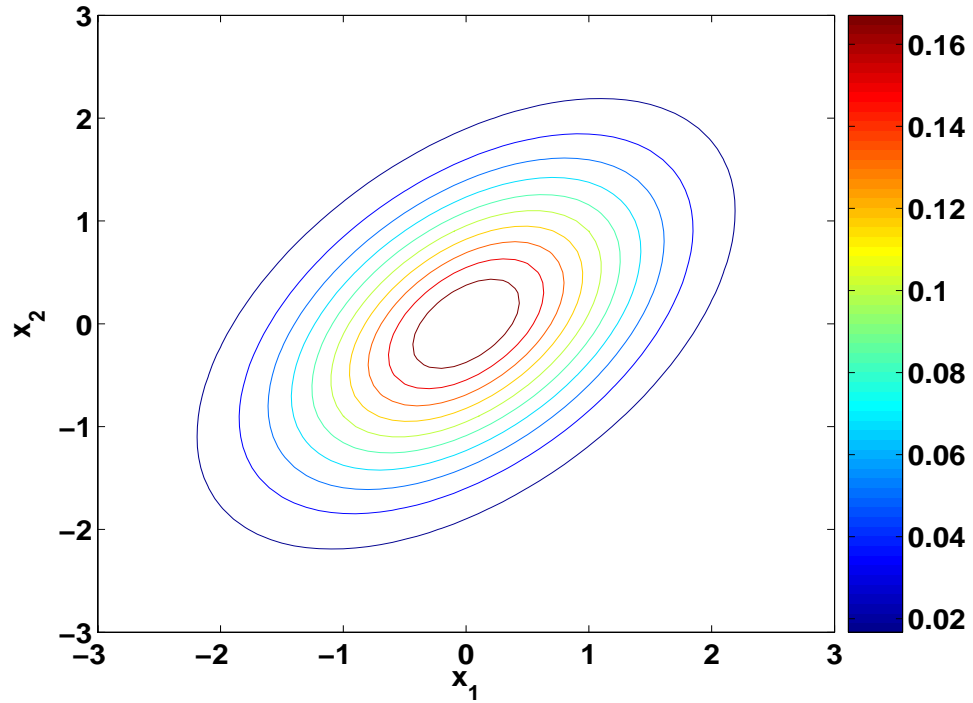
**BIVARIATE UNIT NORMAL,  $\rho=0$**



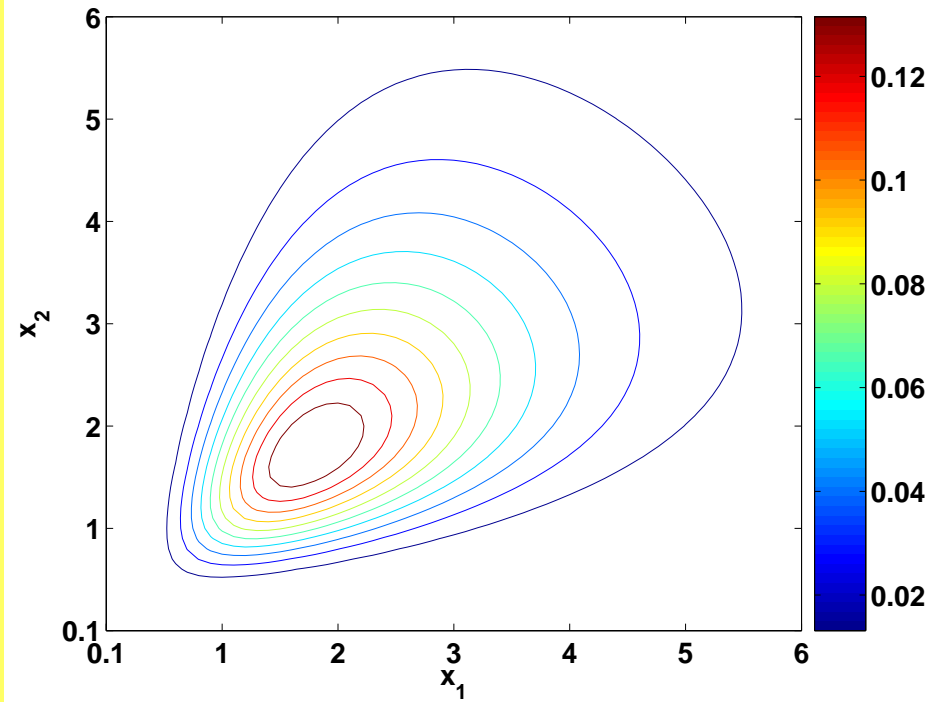
**BIVARIATE LOGNORMAL,  $\mu_1=\mu_2=1, \sigma_1=\sigma_2=0.5, \rho=0$**



**BIVARIATE UNIT NORMAL,  $\rho=0.5$**



**BIVARIATE LOGNORMAL,  $\mu_1=\mu_2=1, \sigma_1=\sigma_2=0.5, \rho=0.5$**



## EXAMPLE WITH THE LORENZ'63 MODEL

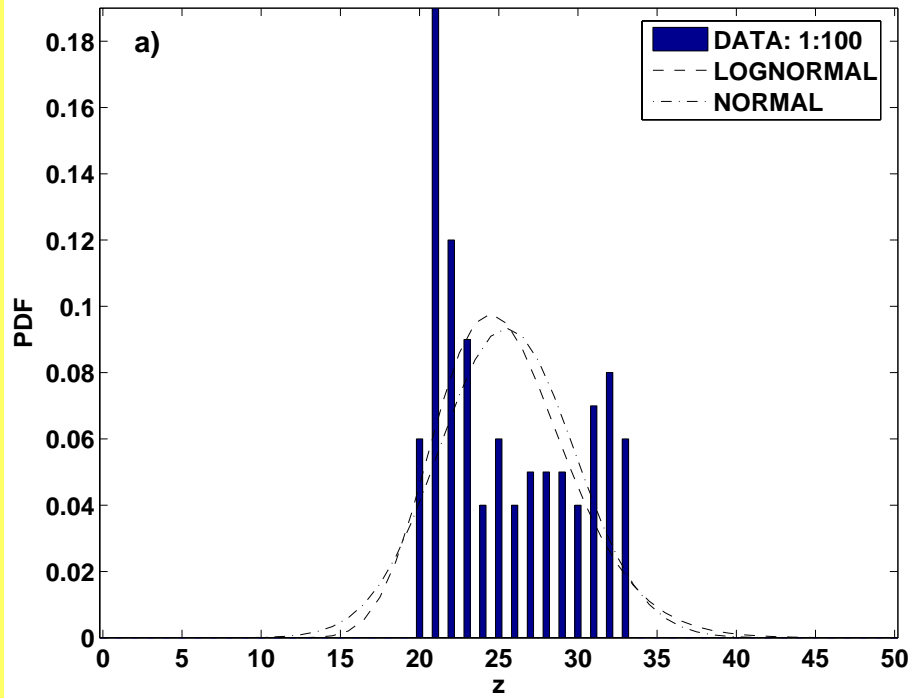
THE MODEL CONSIST OF THE COUPLED SYSTEM OF THREE NON-LINEAR PDES

$$\begin{aligned}\dot{x} &= -\sigma x + \sigma y \\ \dot{y} &= -xz + \rho x - y \\ \dot{z} &= xy - \beta z\end{aligned}$$

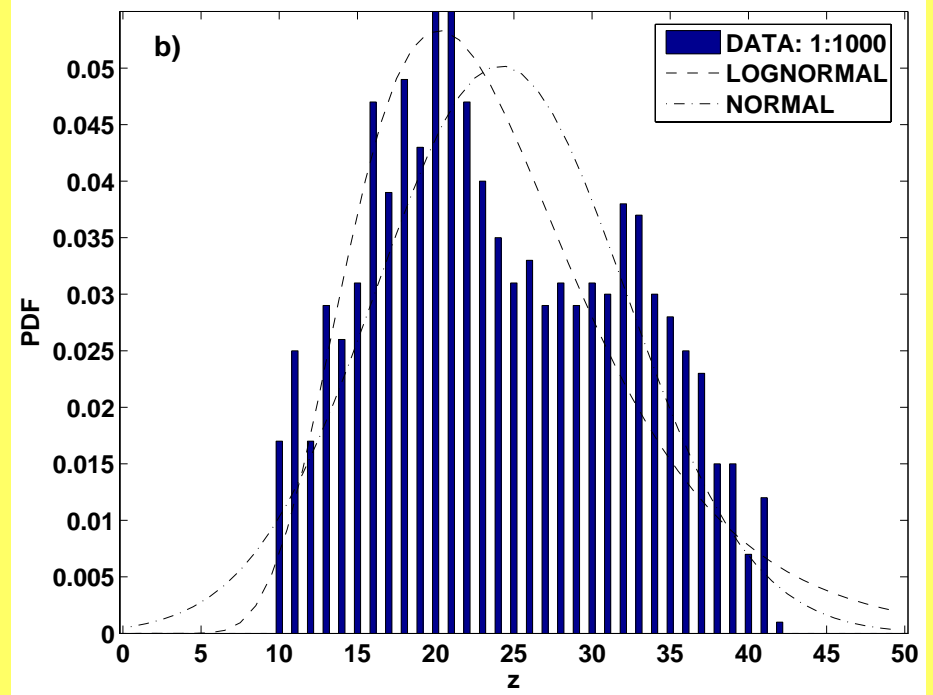
$$\beta = \frac{8}{3}, \quad \sigma = 10 \quad \text{AND} \quad \rho = 28$$

$$x_0 = -5.4458, \quad y_0 = -5.4841 \quad \text{AND} \quad z_0 = 22.5606$$

LOGNORMAL AND NORMAL PLOTS FOR z FIRST 100 TIME STEPS

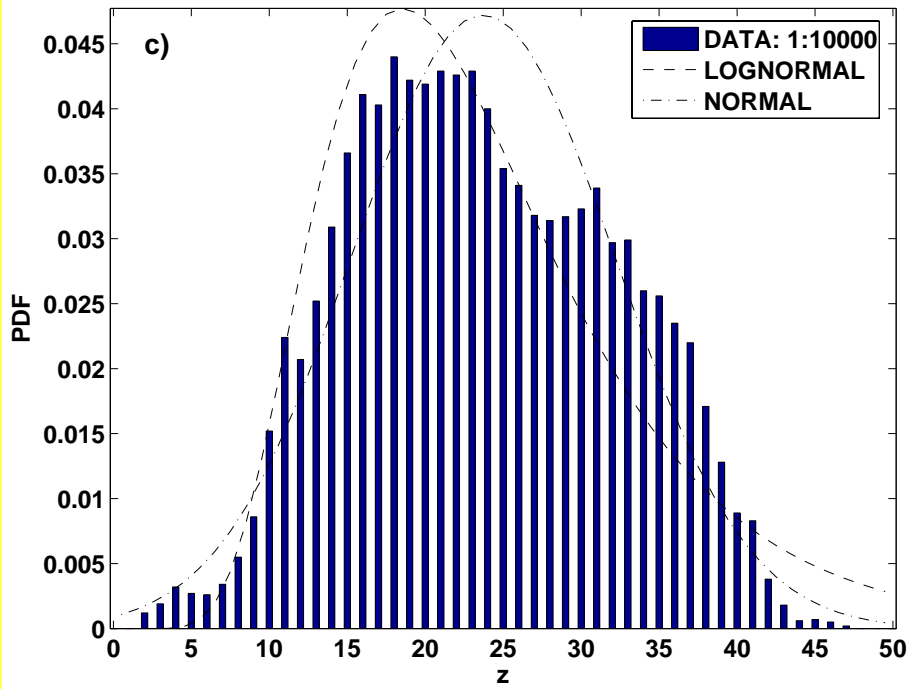


LOGNORMAL AND NORMAL PLOTS FOR z FIRST 1000 TIME STEPS

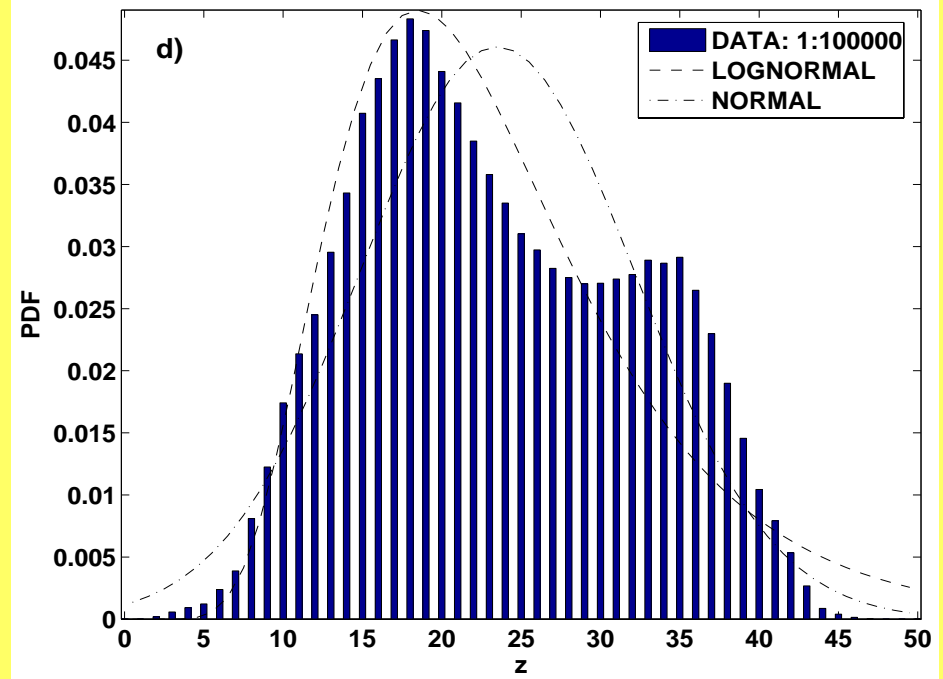




LOGNORMAL AND NORMAL PLOTS FOR z FIRST 10000 TIME STEPS



LOGNORMAL AND NORMAL PLOTS FOR z FIRST 100000 TIME STEPS



# LOGNORMAL DATA ASSIMILATION

We start by consider which statistic to use to best represent the underlying analysis pdf.

The three descriptive statistics are the **mode** 'most likely state', **median** 'unbiased state' and the **mean** 'minimum variance'.

Unlike with the Gaussian distribution and other symmetric distributions these three statistics are not identical so which one to use?

# HYBRID DISTRIBUTION

## FLETCHER AND ZUPANSKI (2006b)

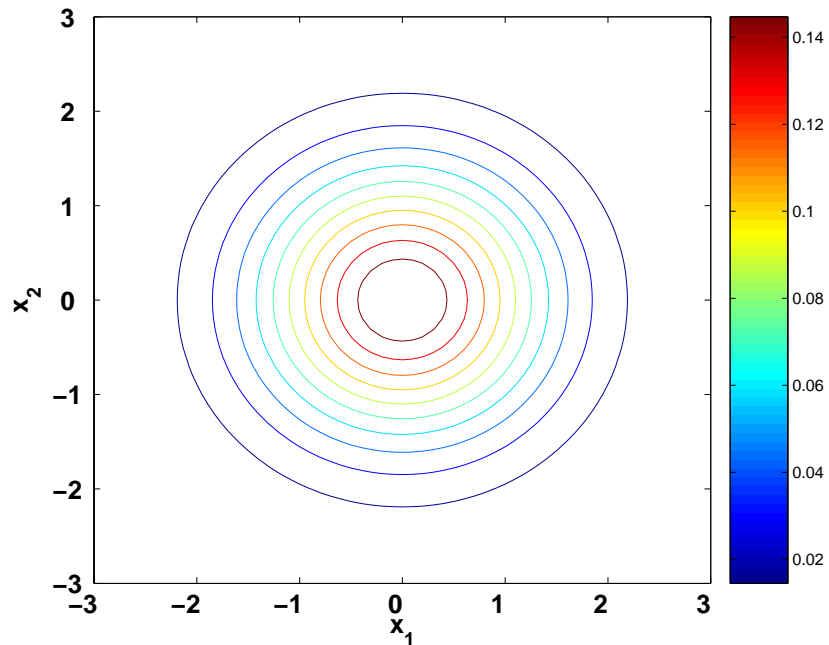
**CAN DEFINE A HYBRID NORMAL-LOGNORMAL  
MULTIVARIATE PROBABILITY DENSITY FUNCTION OF  
THE FORM**

$$f_{p,q}(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{N}{2}} |\mathbf{R}|^{\frac{1}{2}}} \left( \prod_{i=p+1}^N \frac{1}{x_i} \right) \exp \left\{ -\frac{1}{2} (\hat{\mathbf{x}} - \boldsymbol{\mu})^T \mathbf{R}^{-1} (\hat{\mathbf{x}} - \boldsymbol{\mu}) \right\} \quad (3)$$

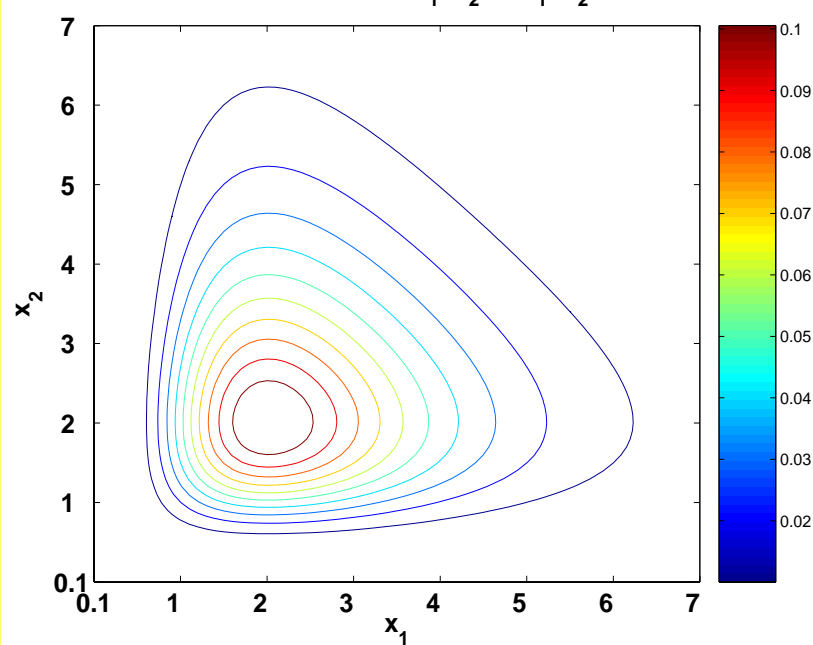
**WHERE**

$$\hat{\mathbf{x}} = \begin{pmatrix} \mathbf{x}_p \\ \ln \mathbf{x}_q \end{pmatrix}$$

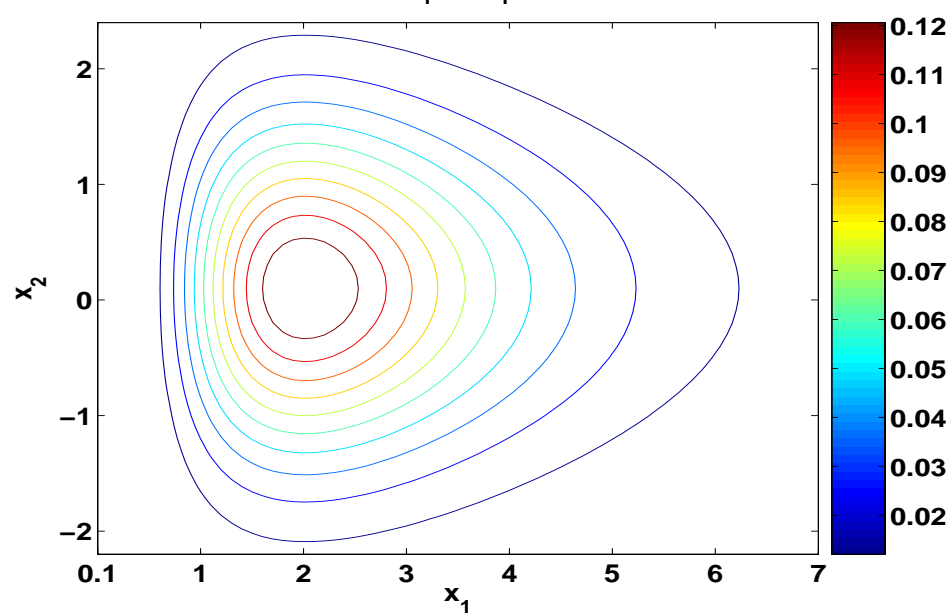
**BIVARIATE UNIT NORMAL,  $\rho=0$**



**BIVARIATE LOGNORMAL,  $\mu_1=\mu_2=1, \sigma_1=\sigma_2=0.5, \rho=0$**



**HYBRID DISTRIBUTION,  $\mu_1=1, \sigma_1=0.5$ , UNIT NORMAL,  $\rho=0$**



# HYBRID ASSIMILATION

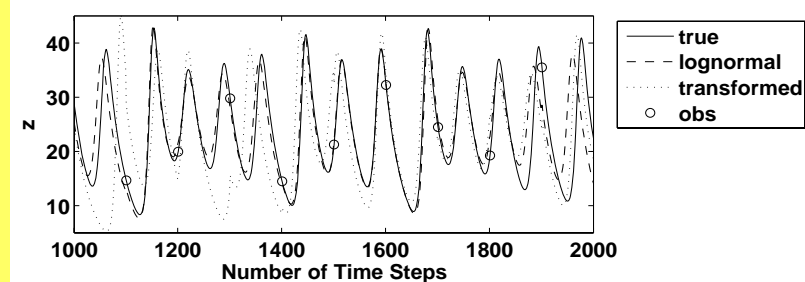
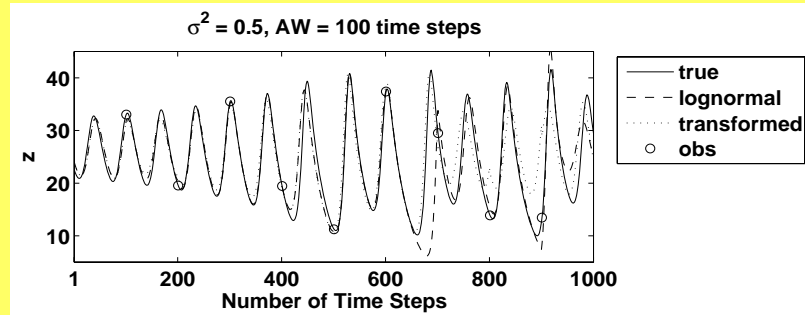
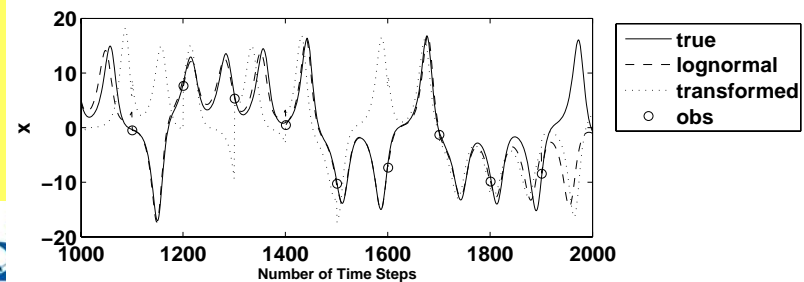
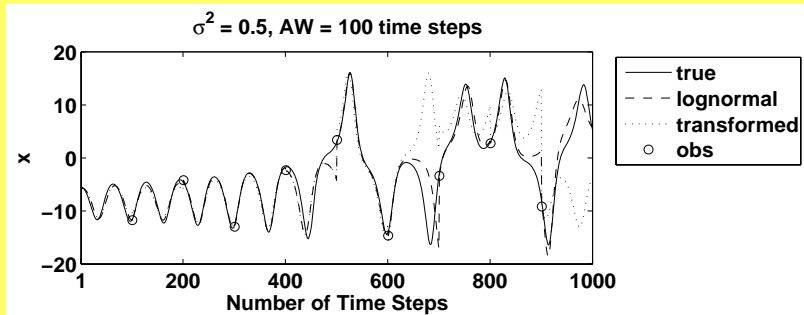
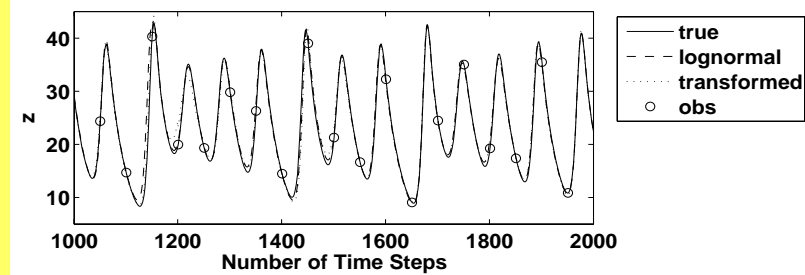
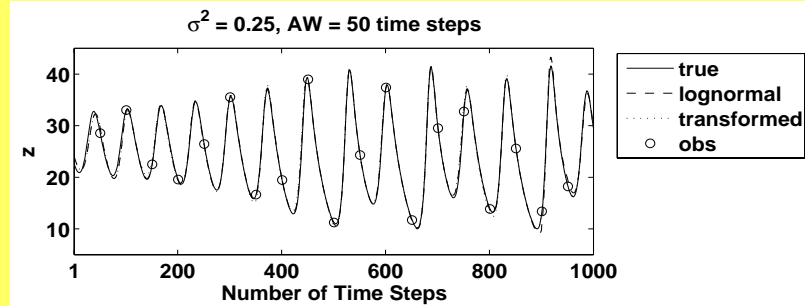
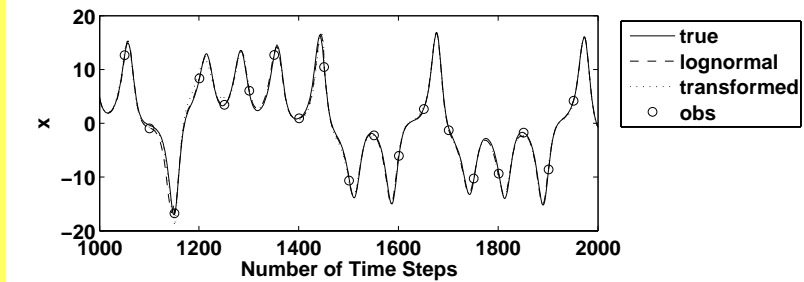
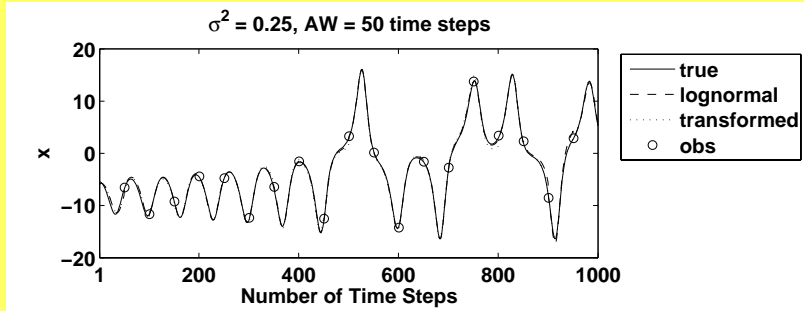
## FLETCHER AND ZUPANSKI 2006b

FROM THE DISTRIBUTION DEFINED IN (3) IT IS POSSIBLE TO DEFINE A COST FUNCTION FOLLOWING THE DERIVATION IN LORENC (1986). THIS IS DEFINED WITH A HYBRID BACKGROUND AS

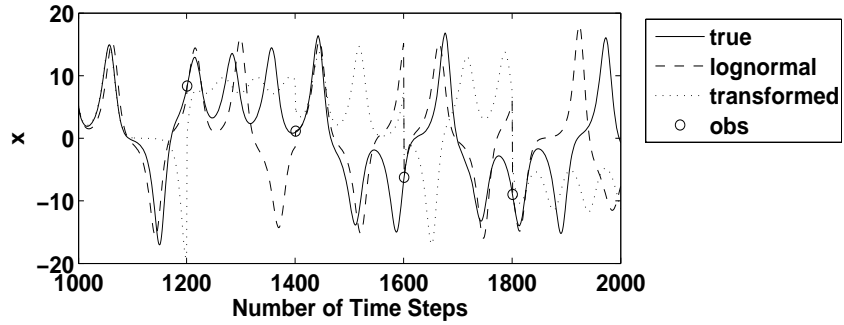
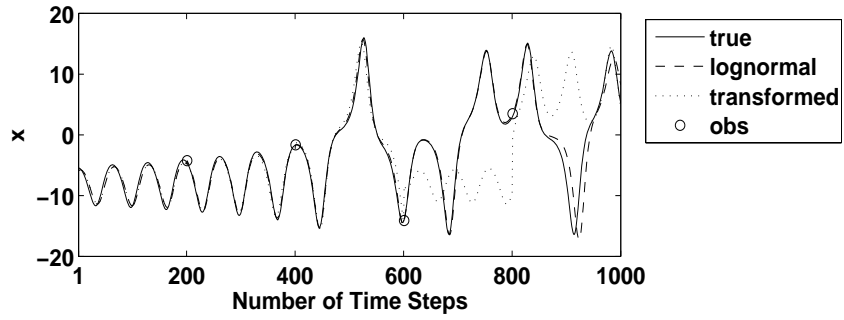
$$J(\mathbf{x}) = \frac{1}{2} \hat{\boldsymbol{\varepsilon}}_b^T \mathbf{B}^{-1} \hat{\boldsymbol{\varepsilon}}_b + \frac{1}{2} \hat{\boldsymbol{\varepsilon}}_o^T \mathbf{R}^{-1} \hat{\boldsymbol{\varepsilon}}_o + \sum_{i=p_1+1}^N \hat{\boldsymbol{\varepsilon}}_{bi} + \sum_{j=p_2+1}^{N_o} \hat{\boldsymbol{\varepsilon}}_{oj} \quad (4)$$

WHERE

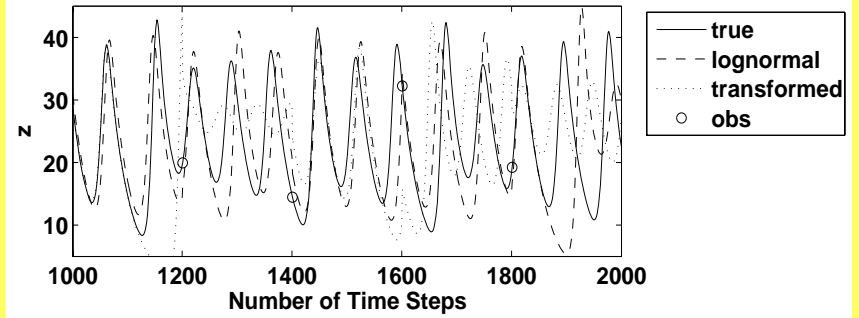
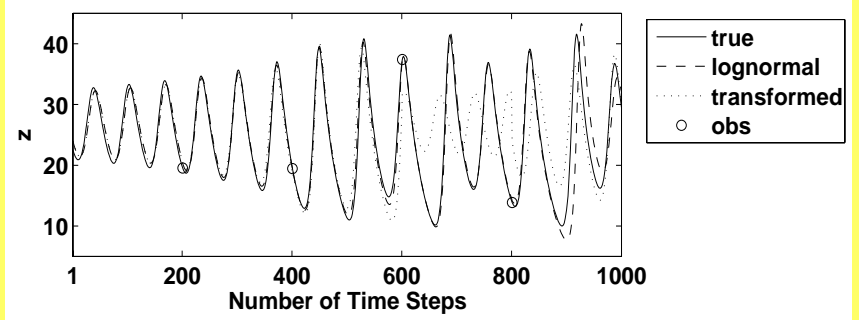
$$\hat{\boldsymbol{\varepsilon}}_b = \begin{pmatrix} \mathbf{x}_{p_1} - \mathbf{x}_{b,p_1} \\ \ln \mathbf{x}_{q_1} - \ln \mathbf{x}_{b,q_1} \end{pmatrix} \quad \hat{\boldsymbol{\varepsilon}}_o = \begin{pmatrix} \mathbf{y}_{p_2} - \mathbf{h}_{p_2}(\mathbf{x}) \\ \ln \mathbf{y}_{q_2} - \ln \mathbf{h}_{q_2}(\mathbf{x}) \end{pmatrix} \quad (5)$$



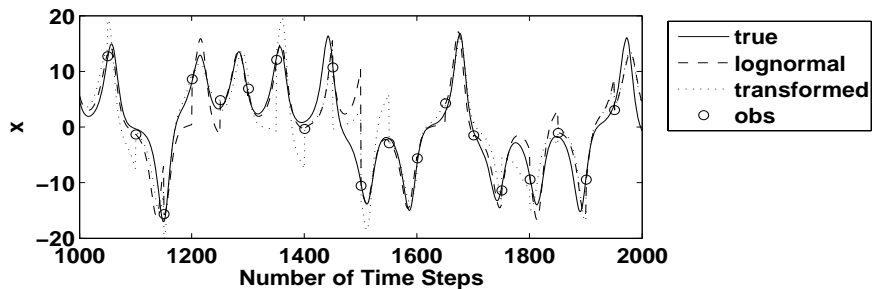
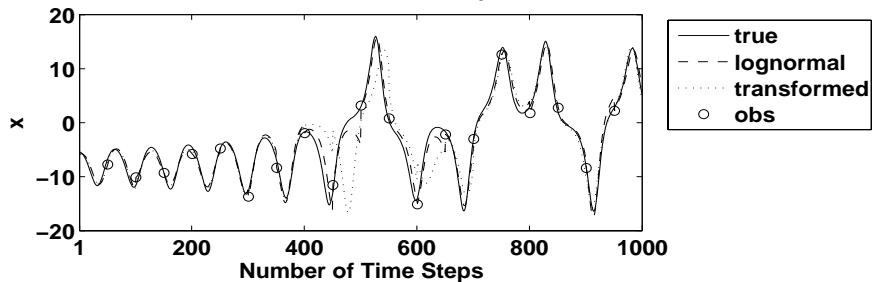
$\sigma^2 = 0.25$ , AW = 200 time steps



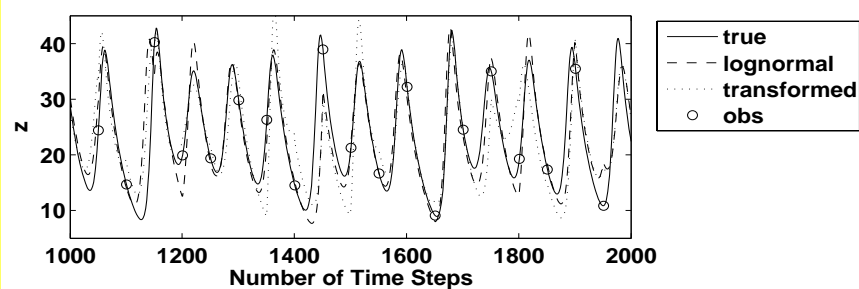
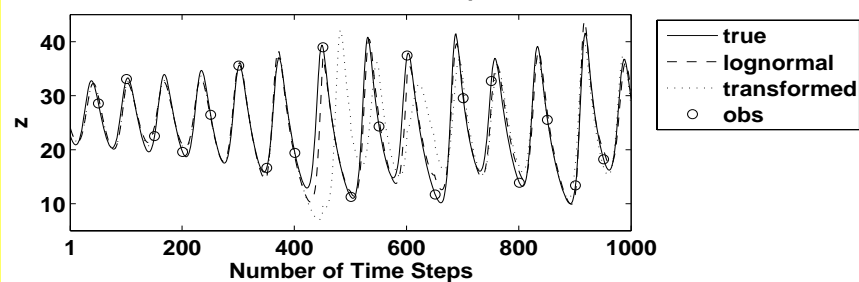
$\sigma^2 = 0.25$ , AW = 200 time steps



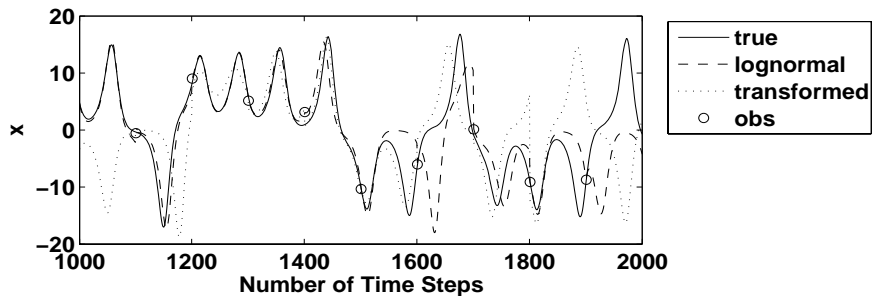
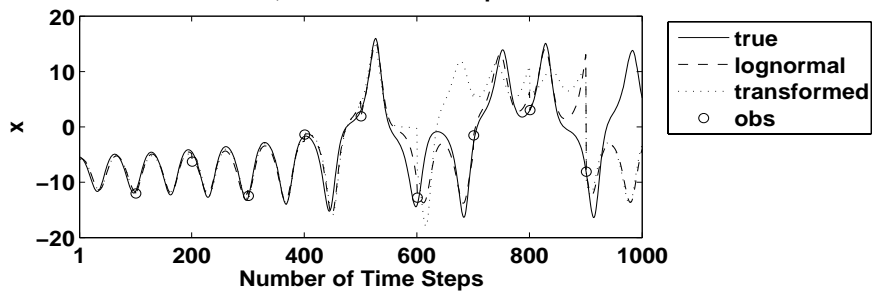
$\sigma^2 = 1, AW = 50$  time steps



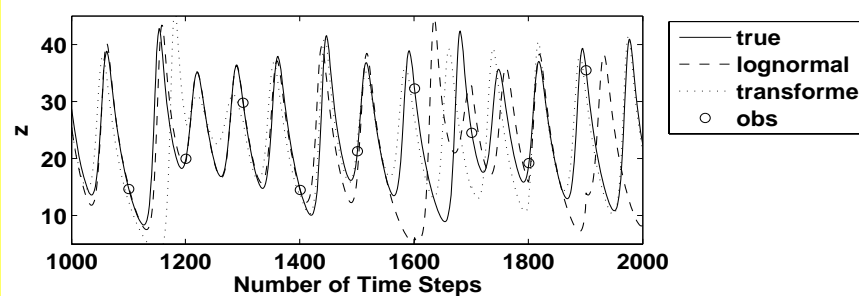
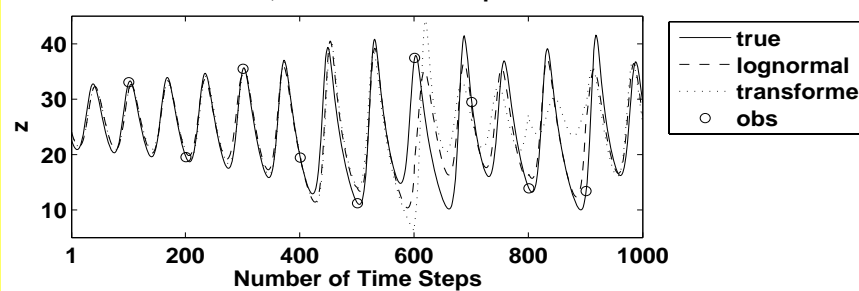
$\sigma^2 = 1, AW = 50$  time steps



$\sigma^2 = 1, AW = 100$  time steps

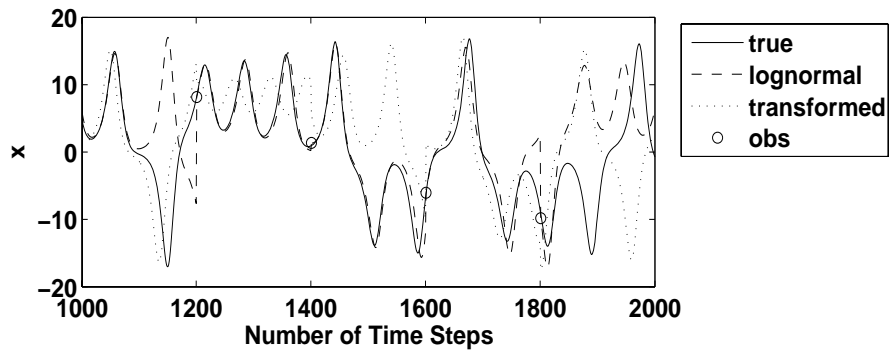
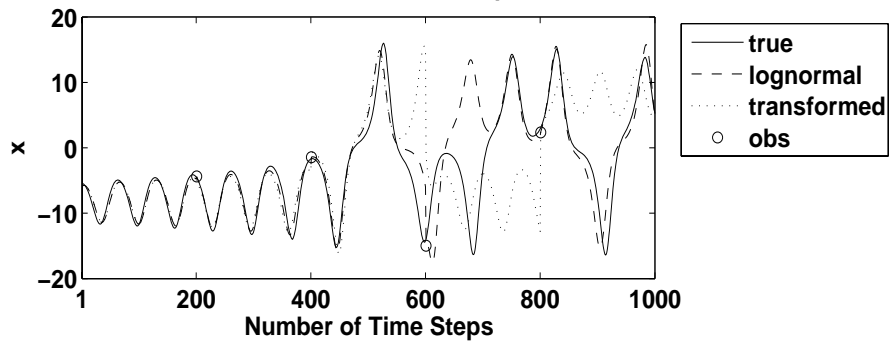


$\sigma^2 = 1, AW = 100$  time steps

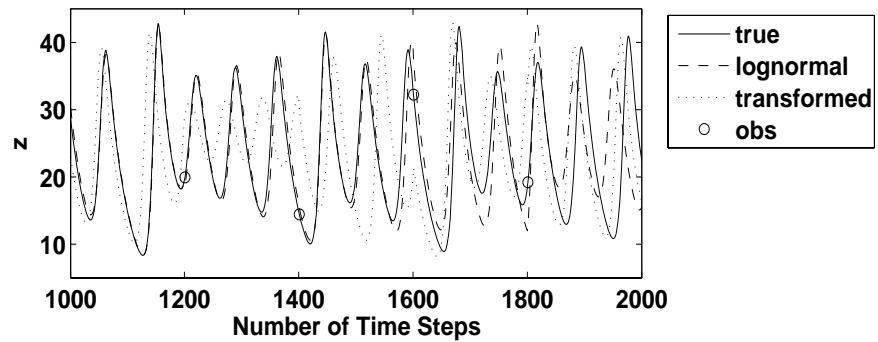
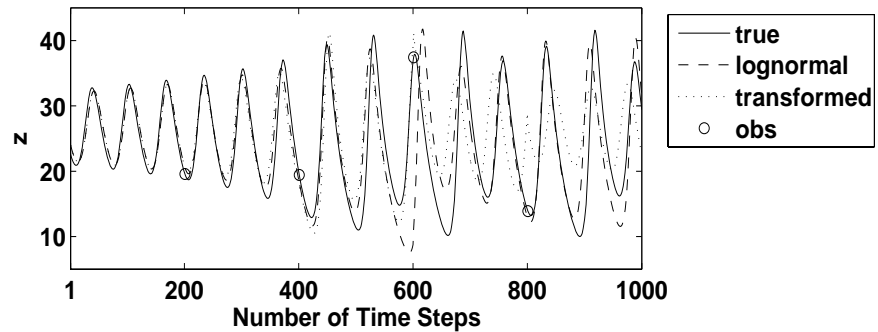




$\sigma^2 = 1$ , AW = 200 time steps



$\sigma^2 = 1$ , AW = 200 time steps



## **APPLICATION TO SATELLITE DATA** **ASSIMILATION**

**The cost functions which we have presented enable us to avoid the bias correction which is applied in both direct radiance data assimilation and in retrievals.**

**Given this lognormal framework we can quantify the bias corrections which are currently applied.**

## TRANSFORMATION (RETRIEVALS)

By transforming to the logarithm of the state variable and the observations which are a mixture then the associated cost function is the standard Gaussian cost function,

$$J(\tilde{\mathbf{x}}) = \frac{1}{2}(\tilde{\mathbf{x}} - \tilde{\mathbf{x}}_b)^T \mathbf{B}^{-1}(\tilde{\mathbf{x}} - \tilde{\mathbf{x}}_b) + \frac{1}{2}(\tilde{\mathbf{y}} - \tilde{\mathbf{h}}(\tilde{\mathbf{x}}))^T \mathbf{R}^{-1}(\tilde{\mathbf{y}} - \tilde{\mathbf{h}}(\tilde{\mathbf{x}}))$$

We should have the following cost function

$$J(\mathbf{x}) = \frac{1}{2} \hat{\boldsymbol{\varepsilon}}_b^T \mathbf{B}^{-1} \hat{\boldsymbol{\varepsilon}}_b + \frac{1}{2} \hat{\boldsymbol{\varepsilon}}_o^T \mathbf{R}^{-1} \hat{\boldsymbol{\varepsilon}}_o + \sum_{i=p_1+1}^N \hat{\boldsymbol{\varepsilon}}_{bi} + \sum_{j=p_2+1}^{N_o} \hat{\boldsymbol{\varepsilon}}_{oj}$$

## SOLUTIONS

The non-linear solutions to the two cost functions on the previous slide are derived in Fletcher and Zupanski (2007) and Lorenc (1986) but are stated here as

$$\tilde{\mathbf{x}} = \tilde{\mathbf{x}}_b + \mathbf{B}\tilde{\mathbf{H}}^T\tilde{\mathbf{W}}_o^T\mathbf{R}^{-1}(\tilde{\mathbf{y}} - \tilde{\mathbf{h}}(\tilde{\mathbf{x}})) \quad (6)$$

$$\hat{\mathbf{x}} = \hat{\mathbf{x}}_b - \mathbf{B} \left[ \begin{pmatrix} \mathbf{0}_{p_1} \\ \mathbf{I}_{q_1} \end{pmatrix} - \mathbf{W}_b^{-1} \mathbf{H}^T \mathbf{W}_o^T \mathbf{R}^{-1} \left( \hat{\boldsymbol{\varepsilon}}_o + \mathbf{R} \begin{pmatrix} \mathbf{0}_{p_2} \\ \mathbf{I}_{q_2} \end{pmatrix} \right) \right] \quad (7)$$

$$\mathbf{W}_b = \frac{\partial \hat{\boldsymbol{\varepsilon}}_b}{\partial \mathbf{x}}, \quad \mathbf{W}_o = \frac{\partial \hat{\boldsymbol{\varepsilon}}_o}{\partial \mathbf{h}}, \quad \tilde{\mathbf{W}}_o = \frac{\partial \tilde{\boldsymbol{\varepsilon}}_o}{\partial \tilde{\mathbf{h}}}, \quad \mathbf{H} = \frac{\partial \mathbf{h}}{\partial \mathbf{x}} \quad \text{and} \quad \tilde{\mathbf{H}} = \frac{\partial \tilde{\mathbf{h}}}{\partial \tilde{\mathbf{x}}}$$

## WRONG DISTRIBUTION

By assuming the wrong distribution we have to see how much of an effect this has. If we assume a Normal background and a set of lognormal observations then if we apply a Taylor series expansion of  $\varepsilon_0$  then we obtain the following cost function

$$J(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \mathbf{x}_b)^T B^{-1}(\mathbf{x} - \mathbf{x}_b) + \frac{1}{2}(\mathbf{y} - h(\mathbf{x}))^T W_o^T \bar{R}^{-1} W_o (\mathbf{y} - h(\mathbf{x})) \quad (8)$$

Where the solution to (8) is given by

$$\mathbf{x} = \mathbf{x}_b + B H^T W_o^T \bar{R}^{-1} W_o (\mathbf{y} - h(\mathbf{x}))$$

## 4D LOGNORMAL DATA ASSIMILATION

Unlike with the three dimensional version of variational data assimilation there is not probability model set out for 4D VAR. Until Fletcher 2007, where we are able to define 4D VAR as a Bayesian Network.

The Gaussian 4D VAR is defined through a calculus of variation problem with initial conditions found through the adjoint. We can do this as well for the lognormal through the inner product

$$g_1(\mathbf{x}_0) = \iiint_A \sum_{i=1}^{N_o} \frac{1}{2} \left\langle \ln y_i - \ln h_i(\mathbf{M}_i(\mathbf{x}_0)), \mathbf{R}_i^{-1} (\ln y_i - \ln h_i(\mathbf{M}_i(\mathbf{x}_0))) \right\rangle$$

As with the Gaussian case we know that the first variation of the functional defined on the previous is equivalent to

$$\begin{aligned}\delta g_1(\mathbf{x}_0) &= \sum_{i=1}^{N_0} \left\langle \ln y_i - \ln h_i(\mathbf{M}_i(\mathbf{x}_0)), \mathbf{W}_{o,i} \mathbf{H}_i \mathbf{M}_i \mathbf{R}_i^{-1} \delta \mathbf{x}_0 \right\rangle \\ &= \left\langle \nabla g_1(\mathbf{x}_0), \delta \mathbf{x}_0 \right\rangle\end{aligned}$$

Through using the properties of inner products we get that the gradient is

$$\nabla g_1(\mathbf{x}_0) = \sum_{i=1}^{N_0} \left( \mathbf{W}_{o,i} \mathbf{H}_i \mathbf{M}_i \mathbf{R}_i^{-1} \right)^T (\ln y_i - \ln h_i(\mathbf{M}_i(\mathbf{x}_0)))$$

What is wrong with the set up on the previous slide?

The solution is a median and not the mode and hence is independent of the variance.

We need to define the functional as

$$g_2(\mathbf{x}_0) = \iiint_A \sum_{i=1}^{N_o} \frac{1}{2} \langle \ln \mathbf{y}_i - \ln \mathbf{h}_i(\mathbf{M}_i(\mathbf{x}_0)) + \mathbf{R}^T \mathbf{1}, \mathbf{R}_i^{-1} (\ln \mathbf{y}_i - \ln \mathbf{h}_i(\mathbf{M}_i(\mathbf{x}_0))) \rangle$$

Which then has a gradient of

$$\nabla g_2(\mathbf{x}_0) = \sum_{i=1}^{N_o} \left( \mathbf{W}_{o,i} \mathbf{H}_i \mathbf{M}_i \mathbf{R}_i^{-1} \right)^T \left( \ln \mathbf{y}_i - \ln \mathbf{h}_i(\mathbf{M}_i(\mathbf{x}_0)) + \mathbf{R}_i^T \mathbf{1} \right)$$



**We are only able to define these different functionals because of our previous work with the 3D VAR lognormal data assimilation method and the properties of the three distribution estimators. However, if we wanted to extend this to variables that are not normal or lognormally distributed we need a probability model.**

**In Fletcher 2007 we present a model through using Bayesian networks theory. The main point about this approach is that it allows us to simplify the expression for Bayes Theorem extended to multiple events which is given by**

$$P(x_0, x_1, x_2, \dots, x_{N_0} | y_1, y_2, y_3, \dots, y_{N_0}) = \\ P(x_0)P(x_1|x_0)P(y_1|x_1, x_0)P(x_2|y_1, x_1, x_0) \\ \dots P(y_{N_0}|x_{N_0}, y_{N_0-1}, x_{N_0}, \dots, y_1, x_1, x_0)$$

**Bayesian networks allow us to remove terms that are not conditioned on other random variables. From the diagram on the board for the perfect model we have that all the model states are only dependent on the initial conditions. We also have that the observations are only dependent on the model state at the time which allows further simplification to give us**

$$P(\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N_0} | \mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3, \dots, \mathbf{y}_{N_0}) =$$

$$P(\mathbf{x}_0) \prod_{i=1}^{N_0} P(\mathbf{y}_i | \mathbf{x}_0)$$

For the multivariate Gaussian case we have

$$P(\mathbf{x}_0) \propto \exp\left\{-\frac{1}{2}(\mathbf{x}_0 - \mathbf{x}_{b,0})^T \mathbf{B}_0^{-1}(\mathbf{x}_0 - \mathbf{x}_{b,0})\right\}$$

$$P(\mathbf{y}_i | \mathbf{x}_0) \propto \exp\left\{-\frac{1}{2}(\mathbf{y}_i - \mathbf{h}_i(\mathbf{M}_i(\mathbf{x}_0)))^T \mathbf{R}_i^{-1}(\mathbf{y}_i - \mathbf{h}_i(\mathbf{M}_i(\mathbf{x}_0)))\right\}$$

For the multivariate lognormal case we have

$$P(\mathbf{x}_\theta) \propto \left( \prod_{j=1}^N \frac{\mathbf{x}_{\theta,i}}{\mathbf{x}_{b,\theta,j}} \right) \times$$
$$\exp \left\{ -\frac{1}{2} (\ln \mathbf{x}_\theta - \ln \mathbf{x}_{b,\theta})^T \mathbf{B}_0^{-1} (\ln \mathbf{x}_\theta - \ln \mathbf{x}_{b,\theta}) \right\}$$
$$P(\mathbf{y}_i | \mathbf{x}_\theta) \propto \left( \prod_{k=1}^{N_{o,i}} \frac{h_{i,k}(\mathbf{M}(\mathbf{x}_\theta))}{y_{i,k}} \right) \times$$
$$\exp \left\{ -\frac{1}{2} (\mathbf{y}_i - \mathbf{h}_i(\mathbf{M}_i(\mathbf{x}_\theta)))^T \mathbf{R}_i^{-1} (\mathbf{y}_i - \mathbf{h}_i(\mathbf{M}_i(\mathbf{x}_\theta))) \right\}$$

## REFERENCES

- COHN, S.E., 1997:** AN INTRODUCTION TO ESTIMATOR ERROR THEORY. *J. Met. Soc. Japan*, **75**, 420-436.
- DERBER, J.C. AND WU, W.-S., 1998:** THE USE OF TOVS CLOUD-CLEARED RADIANCES IN THE NCEP SSI ANALYSIS SYSTEM. *Mon. Wea. Rev.*, **126**, 2287-2299.
- EVENSON, G., 1994:** SEQUENTIAL DATA ASSIMILATION WITH A NONLINEAR QUASI-GEOSTROPHIC MODEL USING MONTE-CARLO METHODS TO FORECAST ERROR STATISTICS. *J. Geophys. Res.* **99** (C5), 10,143-10,162.
- FLETCHER, S.J. AND ZUPANSKI, M. 2006a:** A DATA ASSIMILATION METHOD FOR LOGNORMALLY DISTRIBUTED OBSERVATIONAL ERRORS. In Print: *Q. J. R. Meteor. Soc.*
- FLETCHER, S.J. AND ZUPANSKI, M. 2006b:** A HYBRID MULTIVARIATE NORMAL AND LOGNORMAL DISTRIBUTION FOR DATA ASSIMILATION. *Atmos. Sci. Letters*. **7**, 43-46.
- FLETCHER, S.J. AND ZUPANSKI, M. 2007:** A CAVEAT CONCERNING THE USE OF NORMAL FRAMEWORKS FOR LOGNORMAL VARIABLES AND OBSERVATIONS IN VARIATIONAL DATA ASSIMILATION. Submitted *Mon. Wea. Rev.*
- HARRIS, B.A. AND KELLY, G. 2001:** A SATELLITE RADIANCE-BIAS CORRECTION SCHEME FOR DATA ASSIMILATION. *Q. J. Roy. Meteor. Soc.* **127**, 1453-1486.
- HEYDE, C.C., 1963:** ON A PROPERTY OF THE LOGNORMAL DISTRIBUTION. *J. Roy. Stats. Soc. Ser. B.*, **25**, 392-393.
- KALMAN, R.E., 1960:** A NEW APPROACH TO LINEAR FILTERING AND REDICTION PROBLEMS. *Trans. ASME. J. Basic. Eng.*, **82**, 35-45.
- LORENC, A.C., 1986:** ANALYSIS METHODS FOR NUMERICAL WEATHER PREDICTION. *Q. J. Roy. Meteor. Soc.* **112**, 1177-1194.

## REFERENCES

- MIELKE Jr., P.W., WILLIAMS, J.S. AND WU, S.-U., 1997:** COVARIANCE ANALYSIS TECHNIQUES BASED UPON BIVARIATE LOG-NORMAL DISTRIBUTION WITH WEATHER MODIFICATION APPLICATION. *J. Appl. Meteorol.*, **16**, 183-187.
- MILES N.L., VERLINDE, J. AND CLOTHIAUX, E.E., 2000:** CLOUD DROPLET SIZE DISTRIBUTION IN LOW-LEVEL STRATIFORM CLOUDS. *J. Atmos. Sci.*, **57**, 295-311.
- POLAVARAPU, S., REN, S., ROCHON, Y., SANKEY, D., EK, N., KOSHYK, J. AND TARASICK, D., 2005:** DATA ASSIMILATION WITH THE CANADIAN MIDDLE ATMOSPHERE MODEL. *Atmosphere-Ocean* **43(1)**, 77-100.
- RAYMOND, W.H., 1997:** A THEORETICAL EVALUATION OF THE RELEVANCE OF LOGNORMAL DISTRIBUTIONS FOR THE MOISTURE FLUX AND WIND COMPONENTS. *Mon. Wea. Rev.*, **125**, 3018-3023
- STEPHENS, G.L. AND COAUTHORS, 2002:** THE CLOUDSAT MISSION AND THE A-TRAIN. *Bull. Amer. Meteor. Soc.* **83**, 1771-1190.

# **MINIMISATION AND PRECONDITIONING**

## Outline of lecture

**1) One dimension Newton-Rhapson**

**2) Unconstrained optimisation**

**3) Newton's method for non-linear minimisation**

**4) Wolfe Condition**

**5) Quasi-Newton methods**

**6) Conjugate Gradients**

**7) Preconditioning**

**8) Control Variable Transforms**



## Newton's method for one-dimensional problems

We wish to solve the problem  $f(x)=0$  but we can not analytically find the solutions so we approximate the gradient with a tangent line to the function given a good guess to the true solution. We obtain the iterative formula from a Taylor series expansion of the function in the vicinity of the true root.

$$f(x) = f(x_n) + (x - x_n)f'(x_n) + \frac{(x - x_n)^2}{2} f''(\xi)$$

$$x \leq \xi \leq x_n. \text{ Let } x = \alpha \Rightarrow f(\alpha) = 0$$

$$\therefore \alpha = x_n - \frac{f(x_n)}{f'(x_n)}$$

## Newton's method in multiple dimensions

Like with one dimension case the multi-dimensional version of Newton's method come from a Taylor series which then gives

$$\mathbf{x}_{n+1} = \mathbf{x}_n - F^{-1}(\mathbf{x}_n) f(\mathbf{x}_n)$$
$$F^{-1}(\mathbf{x}_n) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_N} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_N} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_N}{\partial x_1} & \frac{\partial f_N}{\partial x_2} & \dots & \frac{\partial f_N}{\partial x_N} \end{pmatrix}$$

## UNCONSTRAINED OPTIMISATION

In this class of problems we are seeking the minimum/maximum of a continuous function of several variables. To constrain the solution we have that a point  $\alpha$  is called a **STRICT LOCAL MINIMUM** of the function  $f$  if

$$f(x) > f(\alpha) \quad \forall x \text{ close to } \alpha \text{ and } x \neq \alpha$$

Generally an initial guess of  $\alpha$  will be well known and we also assume that the function is twice differentiable with respect to the variables.

**We now reformulate Newton's method for the function to use the famous calculus result that at the minima then the gradient of the function is zero, i.e.**

$$\frac{\partial f(\mathbf{a})}{\partial x_i} = 0, \quad i = 1, 2, \dots, N$$

**Thus we are having to solve the non-linear system of equations given by**

$$\frac{\partial f(\mathbf{x})}{\partial x_i} = 0, \quad i = 1, 2, \dots, N$$

**Which in vector notation this problem looks like**

$$\nabla f(\mathbf{x}) = \left( \frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \dots \quad \frac{\partial f}{\partial x_N} \right)^T = \mathbf{0}$$

**To formulate the minimisation scheme we consider a Taylor series expansion of the gradient of the function.**

$$\begin{aligned}0 &= \nabla f(\alpha) \approx \nabla f(x_n) + (x_{n+1} - x_n)H(x_n) \\ -\nabla f(x_n) &= (x_{n+1} - x_n)H(x_n) \\ (x_{n+1} - x_n) &= -H^{-1}(x_n)\nabla f(x_n) \\ x_{n+1} &= x_n - H^{-1}(x_n)\nabla f(x_n)\end{aligned}$$

**Where**

$$H(x)_{i,j} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j}, \quad 1 \leq i, j \leq N$$

**Is the Hessian matrix. The problem for the weather prediction community is that the size of this matrix is  $10^7$  by  $10^7$  and therefore we can not analytically find G. Therefore we need other methods. These are called the descent methods.**

At the point  $\mathbf{x}_n$  pick a direction  $\mathbf{d}_n$  such that  $f(\mathbf{x}_n)$  will decrease as  $\mathbf{x}$  moves away from  $\mathbf{x}_n$  in the direction,  $\mathbf{d}_n$ .

We start by letting

$$\mathbf{x}_{n+1} = \mathbf{x}_n + s\mathbf{d}_n$$

Usually  $s$  is chosen as the smallest positive relative minimum of  $\phi(s)$  therefore with each iteration

$$f(\mathbf{x}_{n+1}) < f(\mathbf{x}_n)$$

How do we choose  $d_n$  ?

1) The method of steepest descent uses

$$d_n = -\nabla f(x_n)$$

However this is not known for fast convergence.

2) Quasi - Newton methods are such that approximations are made to the Hessian matrix.

However we need a condition of the terms. These are the Wolfe conditions

## Wolfe Conditions

Let  $f : \mathfrak{R}^N \rightarrow \mathfrak{R}$  be a smooth objective function, and let  $\mathbf{d}_n$  be a given search direction.

A step length  $s_n$  is said to satisfy the Wolfe conditions if

i)  $f(\mathbf{x}_n + s_n \mathbf{d}_n) \leq f(\mathbf{x}_n) + c_1 s_n \mathbf{d}_n^T \nabla f(\mathbf{x}_n)$  (Armijo Condition)

ii)  $\mathbf{d}_n^T \nabla f(\mathbf{x}_n + s_n \mathbf{d}_n) \geq c_2 \mathbf{d}_n^T \nabla f(\mathbf{x}_n)$  (Curvature Condition)

There are problems with the second condition and so is made stronger by the following adjustment

$$\left| \mathbf{d}_n^T \nabla f(\mathbf{x}_n + s_n \mathbf{d}_n) \right| \leq c_2 \left| \mathbf{d}_n^T \nabla f(\mathbf{x}_n) \right|$$

i) combined with the new condition above is referred to as the strong Wolfe conditions .



## Quasi-Newton Methods

step i) :  $\Delta \mathbf{x}_n = -s_n H_n^{-1} \nabla f(\mathbf{x}_n)$  where  $s_n$  satisfies the  
Strong Wolfe Conditions

step ii) :  $\mathbf{x}_{n+1} = \mathbf{x}_n + \Delta \mathbf{x}_n$

step iii) : Calculate the new gradient  $\nabla f(\mathbf{x}_{n+1})$  at the new point.

This then enables us to update the Hessian by the approximation

$$\mathbf{y}_n = \nabla f(\mathbf{x}_{n+1}) - \nabla f(\mathbf{x}_n)$$

However, this is a very basic approximation. Therefore this has  
been much research into different approximations to the Hessian

matrix. One very well used approach is the

Broyden - Fletcher - Goldfarb - Shanno (BFGS) method.

## Broyden-Fletcher-Goldfarb-Shanno (BFGS) Method

The limited memory version of this algorithm is used in many data assimilation methods. This approach is used in NCEP's Meso-scale 4D VAR system, (ETA VAR), it is also used in the Cooperative Institute for Research in the Atmosphere (CIRA)/Colorado State University's 4D VAR system, RAMDAS – Regional Atmospheric Modeling Data Assimilation System, Zupanski *et al.* (2005). This is also used in the United Kingdom's Met Office 4D VAR operational weather prediction systems from synoptic to cloud resolving data assimilation systems.

The Hessian matrix is approximated by the addition of two more matrices. The algorithm starts from an initial guess for the true state  $\mathbf{x}_0$  and the Hessian matrix  $\mathbf{B}_0$ .

Then the algorithm is given by the following 4 steps.

step 1) Obtain  $\mathbf{s}_n$  by solving  $\mathbf{B}_n \mathbf{s}_n = -\nabla f(\mathbf{x}_n)$ .

step 2) Perform a line search to find the optimal  $\alpha_n$  in the direction found in the first step then update the state

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \alpha_n \mathbf{s}_n$$

Step 3)  $\mathbf{y}_n = \nabla f(\mathbf{x}_{n+1}) - \nabla f(\mathbf{x}_n)$

Step 4)  $\mathbf{B}_{n+1} = \mathbf{B}_n + \frac{\begin{pmatrix} \mathbf{y}_n \mathbf{y}_n^T \end{pmatrix}}{\begin{pmatrix} \mathbf{y}_n^T \mathbf{s}_n \end{pmatrix}} - \frac{\begin{pmatrix} \mathbf{B}_n \mathbf{s}_n \mathbf{s}_n^T \mathbf{B}_n^T \end{pmatrix}}{\begin{pmatrix} \mathbf{s}_n^T \mathbf{B}_n \mathbf{s}_n \end{pmatrix}}$

## Non-linear Conjugate Gradients methods

These approaches are used instead of Quasi-Newton methods as they avoid the Hessian approximations but are not as fast to converge as the Newton methods. These approaches are seeking solutions to the problem

$$f(\mathbf{x}) = |\mathbf{Ax} - \mathbf{b}|^2$$

Where the minimum occurs when the gradient is zero i.e.

$$\nabla f(\mathbf{x}) = 2\mathbf{A}^T(\mathbf{Ax} - \mathbf{b}) = 0$$

The non-linear CG methods seek a solution to the gradient problem same as the Newton methods.

So given a cts function to minimise the gradient indicates the direction of a maximum increment. Therefore we start in the opposite (steepest descent) direction.

$$1) \Delta \mathbf{x}_n = -\nabla f(\mathbf{x}_n)$$

2) compute the scalar  $\beta_n$  from the following two choices

$$\beta_n^{\text{FR}} = \frac{\Delta \mathbf{x}_n^T \Delta \mathbf{x}_n}{\Delta \mathbf{x}_{n-1}^T \Delta \mathbf{x}_{n-1}}$$

$$\beta_n^{\text{PR}} = \frac{\Delta \mathbf{x}_n^T (\Delta \mathbf{x}_n - \Delta \mathbf{x}_{n-1})}{\Delta \mathbf{x}_{n-1}^T \Delta \mathbf{x}_{n-1}}$$

$$3) \Delta \mathbf{x}_n = \Delta \mathbf{x}_n + \beta_n \Delta \mathbf{x}_{n-1}$$

4) Optimise  $\alpha_n \min f(\mathbf{x}_n + \alpha_n \Delta \mathbf{x}_n)$

$$5) \mathbf{x}_{n+1} = \mathbf{x}_n + \alpha_n \Delta \mathbf{x}_n$$

# Minimization in ensemble DA

**Direct solution of Extended Kalman Filter:**

$$\mathbf{x}_a = \mathbf{x}_b + \mathbf{P}_f \mathbf{K}^T (\mathbf{K} \mathbf{P}_f \mathbf{K}^T + \mathbf{R})^{-1} [\mathbf{y} - \mathbf{K}(\mathbf{x}_b)]$$

$$\mathbf{P}_a = \mathbf{P}_f - \mathbf{P}_f \mathbf{K}^T (\mathbf{K} \mathbf{P}_f \mathbf{K}^T + \mathbf{R})^{-1} \mathbf{K} \mathbf{P}_f$$

**Equivalent to solving a *quadratic* minimization problem (Lorenz 1986):**

$$J = \frac{1}{2} (\mathbf{x} - \mathbf{x}_b)^T \mathbf{P}_f^{-1} (\mathbf{x} - \mathbf{x}_b) + \frac{1}{2} [\mathbf{y} - \mathbf{K}(\mathbf{x})]^T \mathbf{R}^{-1} [\mathbf{y} - \mathbf{K}(\mathbf{x})]$$

**Subject to**

$$\mathbf{K}(\mathbf{x}) = \mathbf{K}(\mathbf{x}_b) + \mathbf{K}(\mathbf{x} - \mathbf{x}_b)$$

$$\mathbf{K} = \left( \frac{\partial \mathbf{K}}{\partial \mathbf{x}} \right)$$

Ensemble DA is a minimization process. Due to perfect Hessian preconditioning, the minimum solution is obtained in a single minimization iteration (for near-quadratic problem).

# Minimization in variational DA

Variational DA is a minimization algorithm

$$J = \frac{1}{2}(\mathbf{x} - \mathbf{x}_b)^T \mathbf{P}_f^{-1}(\mathbf{x} - \mathbf{x}_b) + \frac{1}{2}[\mathbf{y} - K(\mathbf{x})]^T \mathbf{R}^{-1}[\mathbf{y} - K(\mathbf{x})]$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{P}^{-1} \mathbf{g}_k$$

$\mathbf{P}$  - preconditioning matrix

- Gradient based minimization (CG, quasi-Newton, truncated Newton)
  - Adjoint
  - Multiple tangent linear models
  - Finite differences
- Difficult preconditioning due to large dimensions of the problem
- Preconditioning impacts not only the minimization convergence (cost), but also the accuracy of the solution (quality)

# Hessian preconditioning

*Hessian and inverse Hessian* (linear contribution)

$$\mathbf{H} = \frac{\partial^2 J}{\partial \mathbf{x}^2} = \mathbf{P}_f^{-1} + \mathbf{K}^T \mathbf{R}^{-1} \mathbf{K} = \mathbf{P}_f^{-\frac{T}{2}} (\mathbf{I} + \mathbf{A}) \mathbf{P}_f^{-\frac{1}{2}} \quad \mathbf{H} = \mathbf{E} \mathbf{E}^T$$

$$\mathbf{H}^{-1} = \left( \frac{\partial^2 J}{\partial \mathbf{x}^2} \right)^{-1} = \mathbf{P}_f^{1/2} (\mathbf{I} + \mathbf{A})^{-1} \mathbf{P}_f^{T/2}$$

$$\mathbf{A} = \mathbf{P}_f^{T/2} \mathbf{K}^T \mathbf{R}^{-1} \mathbf{K} \mathbf{P}_f^{1/2}$$

*Change of variable (preconditioning):*

$$\mathbf{x} - \mathbf{x}_B = \mathbf{E}^{-T} \boldsymbol{\zeta} = \mathbf{P}_f^{1/2} (\mathbf{I} + \mathbf{A})^{-T/2} \boldsymbol{\zeta}$$

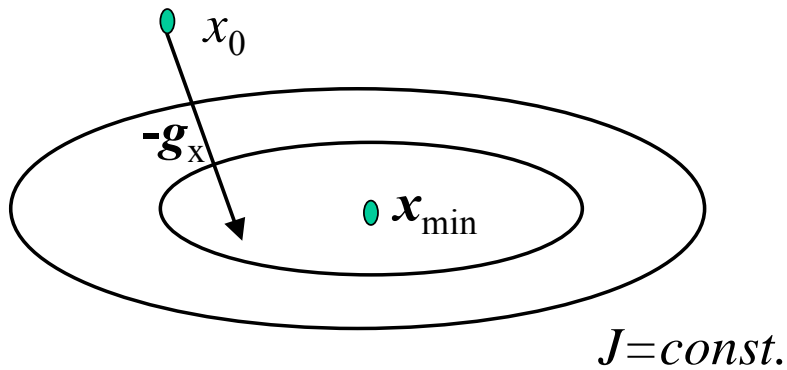
$$\mathbf{H}_\zeta = \mathbf{E}^{-1} \mathbf{H} \mathbf{E}^{-T} = \mathbf{E}^{-1} \mathbf{E} \mathbf{E}^T \mathbf{E}^{-T} = \mathbf{I}$$

Ideal Hessian  
preconditioning

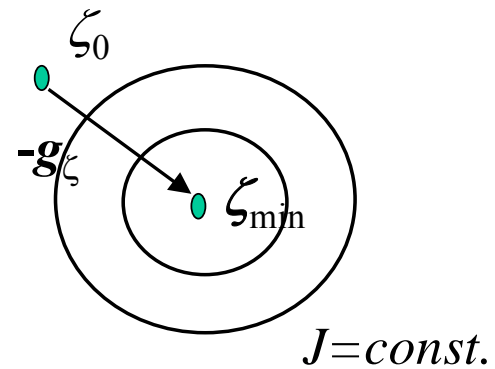


# Ideal Hessian Preconditioning

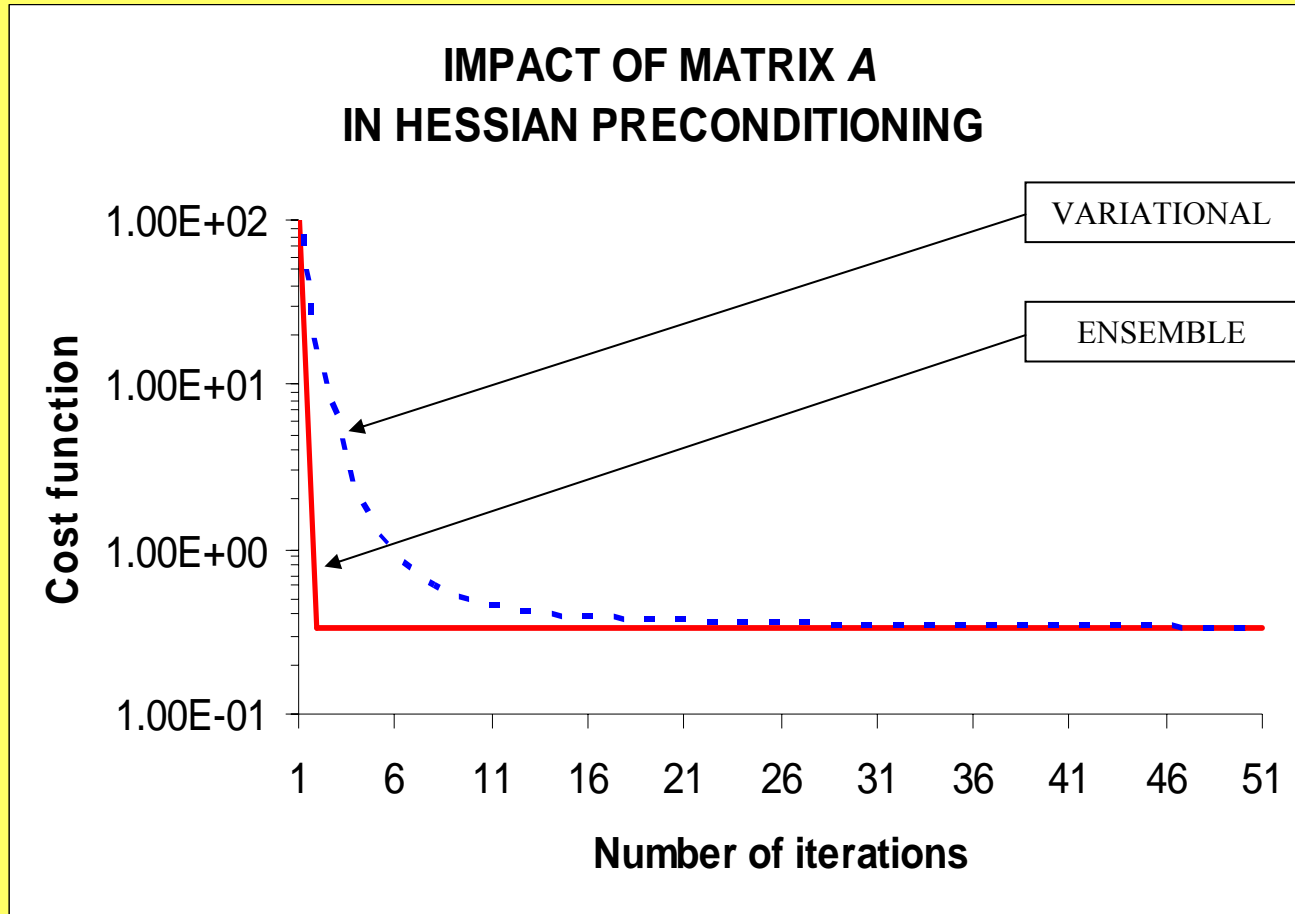
*Physical space ( $x$ )*



*Preconditioning space ( $\zeta$ )*



# Hessian Preconditioning



$$P_{VAR}^{-1} = P_f$$

$$P_{ENS}^{-1} = P_f^{1/2} (I + A)^{-1} P_f^{T/2}$$

# Impact of preconditioning in NCEP's Eta 4DVAR system

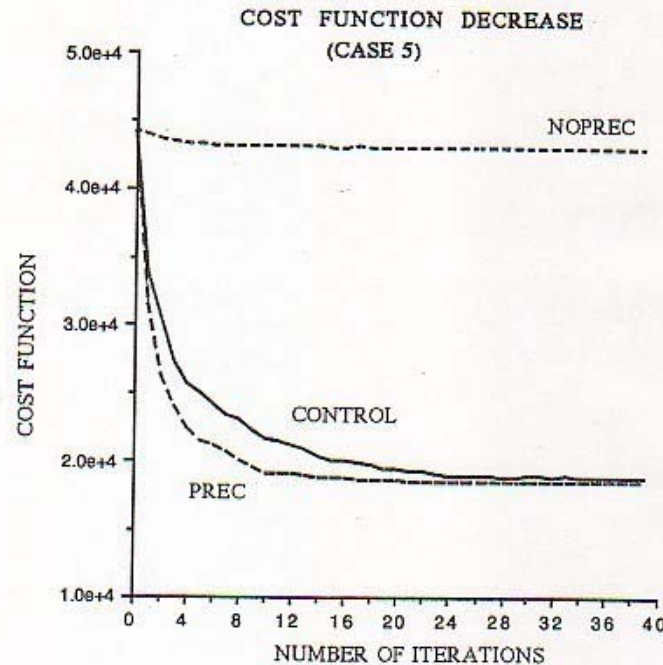


FIG. 2. The cost function decrease during the first 40 iterations of minimization algorithm, for the assimilation period 0000–1200 UTC 24 November 1994 (case 5). The three experiments shown include CONTROL, with background variance used as a preconditioner (full line); NOPREC, with no preconditioning used (short-dashed line); and PREC, with the proposed preconditioning method used (long-dashed line). *(Zupanski 1996, Mon. Wea. Rev.)*

Incorrect solution without preconditioning !

# PRECONDITIONING IN 4DVAR (NCEP ETA, CSU RAMS)

- Controls the magnitude of control variable adjustment

$$d = P^{-1} f(g)$$

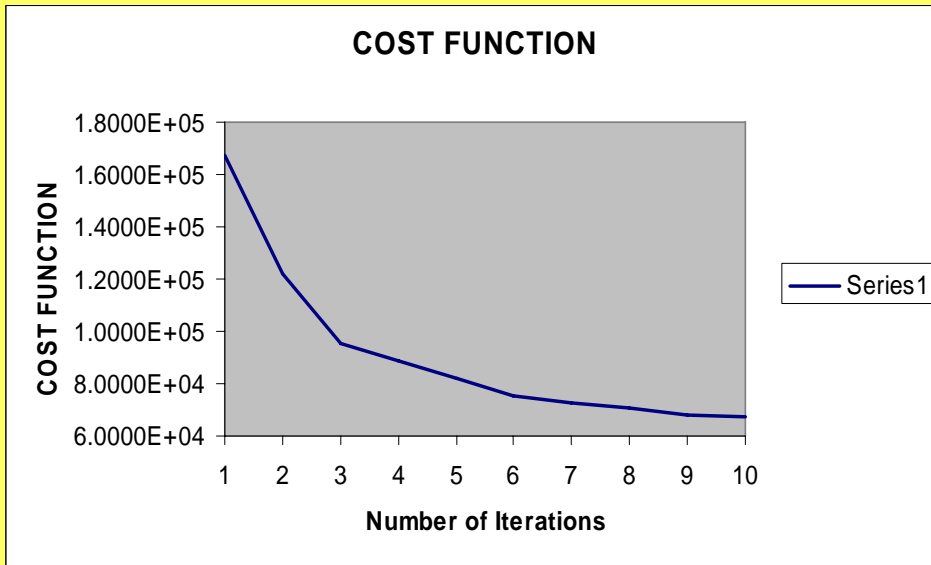
**Synoptic situation dependent (D=diagonal empirical matrix)**

$$\frac{\partial^2 J}{\partial x^2} = B^{1/2} (I + D)^{-1} B^{T/2}$$

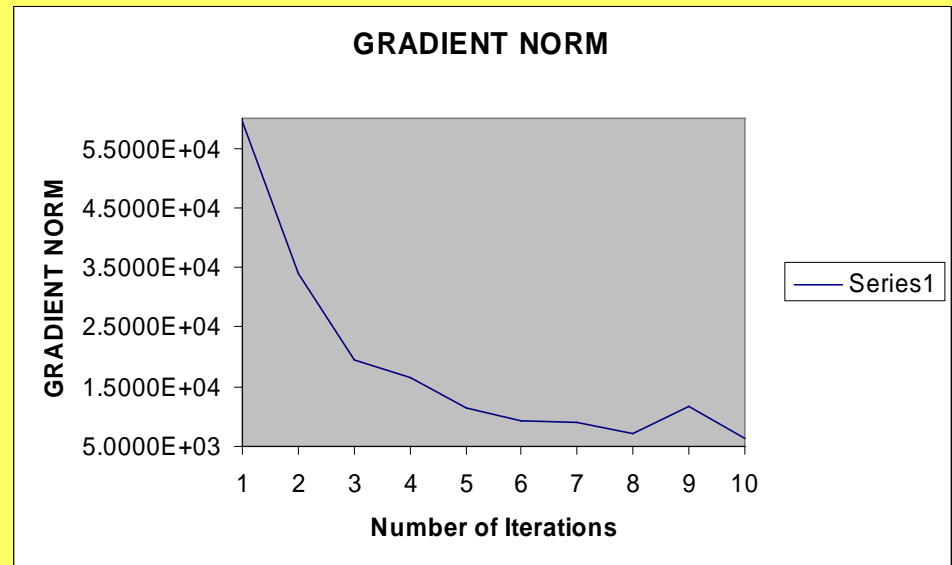
# 4DVAR RESULTS WITH *RAMS* MODEL

- **6-HOUR ASSIMILATION PERIOD:** 03/08/2002 12 Z – 18 Z
- **WRF observations** at 15Z and 18Z
- **Control Variables:** *INITIAL CONDITIONS + MODEL ERROR*
  - - *perturbation Exner function*
  - - *potential temperature*
  - - *horizontal winds*
  - - *total mixing ratio*
- **4DVAR analysis:** **END** of the assimilation interval
- **CSU RAMS** non-hydrostatic NWP model, Level 2 Microphysics
- **Adjoint** with explicit microphysics
- 15 km horizontal resolution, 31 vertical level
- 120 X 80 grid points (1800 km X 1200 km)
- **Calculations** performed on Linux PC Cluster (8 CPU)
- **10 MINIMIZATION ITERATIONS**
- **ONE** minimization iteration cost ~ 10 RAMS model integrations (6-h forecast)

# MINIMIZATION ITERATIONS



Smooth Convergence



# COST FUNCTION

## *Cost Function*

$$J(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \mathbf{x}_B)^T \mathbf{P}_f^{-1}(\mathbf{x} - \mathbf{x}_B) + \frac{1}{2} \sum_{n=1}^N [K(M(\mathbf{x})) - \mathbf{y}]_n^T \mathbf{R}_n^{-1} [K(M(\mathbf{x})) - \mathbf{y}]_n$$

- $\mathbf{x}$  – *control variable*  $\sim O(10^8)$
- $n$  – time index ( $n=1, \dots, N$ )
- $K$  – (non-linear) observation operator
- $M$  – (non-linear) NWP model
- $\mathbf{y}$  – observation vector
- $\mathbf{R}$  – observation error covariance
- $\mathbf{x}_B$  – first-guess (background) model state
- $\mathbf{P}_f$  – forecast error covariance

# HESSIAN

*Hessian matrix* (linear contribution)

$$\mathbf{H} = \frac{\partial^2 J}{\partial \mathbf{x}^2} = \mathbf{P}_f^{-1} + \sum_{n=1}^N \mathbf{M}^T \mathbf{K}_n^T \mathbf{R}_n^{-1} \mathbf{K}_n \mathbf{M} = \mathbf{P}_f^{-\frac{T}{2}} (\mathbf{I} + \mathbf{A}) \mathbf{P}_f^{-\frac{1}{2}}$$

$$\mathbf{A} = \sum_{n=1}^N \mathbf{P}_f^{T/2} \mathbf{M}^T \mathbf{K}_n^T \mathbf{R}_n^{-1} \mathbf{K}_n \mathbf{M} \mathbf{P}_f^{1/2}$$

$$\mathbf{H} = \mathbf{E} \mathbf{E}^T$$

*Change of variable (preconditioning):*

$$\mathbf{x} - \mathbf{x}_B = \mathbf{E}^{-T} \boldsymbol{\zeta} = \mathbf{P}_f^{1/2} (\mathbf{I} + \mathbf{A})^{-T/2} \boldsymbol{\zeta}$$

$$\mathbf{H}_\zeta = \mathbf{E}^{-1} \mathbf{H} \mathbf{E}^{-T} = \mathbf{E}^{-1} \mathbf{E} \mathbf{E}^T \mathbf{E}^{-T} = \mathbf{I}$$



# PRECONDITIONING *ISSUES* IN APPLICATIONS TO REALISTIC ATMOSPHERIC AND OCEANIC DATA ASSIMILATION

- (1) *Large dimensions of the Hessian*
- (2) *Unavailable matrix representation of the Hessian*
- (3) *Inversion of  $10^8 \times 10^8$  matrix not computationally feasible*
- (4) *Computationally expensive calculation of the cost-function and gradient*
- (5) ***Preconditioning is necessary: Hessian Condition Number  $\sim O(10^{15})$***

# PRECONDITIONING *SOLUTIONS* IN APPLICATIONS TO REALISTIC ATMOSPHERIC AND OCEANIC DATA ASSIMILATION

(1) *(Square-root) forecast error covariance matrix:*

$$\mathbf{E}^{-T} \cong \mathbf{P}_f^{1/2}$$

- *assumption:*

$$\mathbf{A} = \mathbf{I}$$

*Hessian Condition Number*  $\sim O(10^2)$ -  $O(10^3)$

(2) *Empirical preconditioning:*

$$\mathbf{E}^{-T} \cong \mathbf{P}_f^{1/2} (\mathbf{I} + \mathbf{D})^{-\frac{T}{2}}$$

- *assumption:*

$$\mathbf{A} = \mathbf{D}$$

*D* – *empirical diagonal matrix*

*Hessian Condition Number*  $\sim O(10^1)$

# EMPIRICAL PRECONDITIONING

(Zupanski - Tellus 1993, MWR 1996)

*Iterative minimization:*  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$  ;  $\mathbf{H}\mathbf{d}_k = -\mathbf{g}_k$

*USE TAYLOR EXPANSION OF THE COST FUNCTION:*

$$J_{k+1} = J_k + \alpha \mathbf{g}^T \mathbf{d} + \frac{1}{2} \alpha \mathbf{d}^T \mathbf{H} \mathbf{d}$$

$$\Delta J = J_k - J_{k+1} = \alpha \left(1 - \frac{\alpha}{2}\right) (-\mathbf{g}^T \mathbf{d})$$

# EMPIRICAL PRECONDITIONING – cont.1

*NOTE: Two (or more) components of the cost-function and the inner product*

$$J = J^B + J^{obs} \Rightarrow \Delta J = \Delta J^B + \Delta J^{obs}$$

$$(-\mathbf{g}^T \mathbf{d}) = (-\mathbf{g}^T \mathbf{d})^B + (-\mathbf{g}^T \mathbf{d})^{obs}$$

$$\mathbf{H} = \frac{\partial^2 J^B}{\partial \mathbf{x}^2} + \frac{\partial^2 J^{obs}}{\partial \mathbf{x}^2} = \mathbf{P}_f^{-1} + \mathbf{P}_f^{-\frac{T}{2}} \mathbf{A} \mathbf{P}_f^{-\frac{1}{2}}$$

## *REMARKS:*

- *Taylor expansion of the cost function can be applied to the total cost function, or to its components*
- *In this application, the background has well-defined covariance ( $\mathbf{P}_f$ )*
- *Apply preconditioning approximation only to the observational component of the cost function*

$$\mathbf{H}^{obs} = \frac{\partial^2 J^{obs}}{\partial \mathbf{x}^2} = \mathbf{P}_f^{-\frac{T}{2}} \mathbf{A} \mathbf{P}_f^{-\frac{1}{2}}$$

# EMPIRICAL PRECONDITIONING – cont.2

*Assumption:*

*MOST IMPORTANT VARIABILITY OF THE HESSIAN MATRIX IS*

- (1) In vertical coordinate direction*
- (2) For different physical variables*

$$\Delta J_L = \alpha \left(1 - \frac{\alpha}{2}\right) (-\mathbf{g}^T \mathbf{d})_L$$

*INTRODUCE UNKNOWN DIAGONAL MATRIX  $\mathbf{D}$ :*

- (1) Assume the Hessian preconditioning in the form*

$$(\mathbf{E}\mathbf{E}^T)^{-1} = \mathbf{P}_f^{1/2} \mathbf{D}^{-1} \mathbf{P}_f^{T/2}$$

- (2) Matrix  $\mathbf{D}$  elements change ONLY in vertical, and for physical variables*
- (3) NOTE: Only observational component of the cost function considered, the prior has known Hessian*

# EMPIRICAL PRECONDITIONING – cont.3

*Descent direction (unknown D):*

$$\mathbf{d} = -\mathbf{P}_f^{1/2} \mathbf{D}^{-1} \mathbf{P}_f^{T/2} \mathbf{g}$$

*Gradient norm:*

$$(-\mathbf{g}^T \mathbf{d})_L = (\mathbf{g}^T \mathbf{P}_f^{\frac{1}{2}} \mathbf{D}^{-1} \mathbf{P}_f^{\frac{T}{2}} \mathbf{g})_L \cong D_L^{-1} (\mathbf{g}^T \mathbf{P}_f \mathbf{g})_L$$

*Cost function:*

$$\Delta J_L = \beta J_L$$

$$\beta < 1$$

*Assumed decrease of the cost function  
(most often  $\beta \sim 0.5$ )*

*In common applications,  $\mathbf{P}_f$  defined over horizontal grid only*

$$(\mathbf{g}^T \mathbf{P}_f^{\frac{1}{2}} \mathbf{D}^{-1} \mathbf{P}_f^{\frac{T}{2}} \mathbf{g})_L = \mathbf{g}_L^T (\mathbf{P}_f^{\frac{1}{2}})_L \mathbf{D}_L^{-1} (\mathbf{P}_f^{\frac{T}{2}})_L \mathbf{g}_L = D_L^{-1} (\mathbf{g}^T \mathbf{P}_f \mathbf{g})_L$$

