

Mixed Lognormal-Gaussian Incremental Variational Data Assimilation

Steven J. Fletcher, Anton J. Kliever, Andrew S. Jones and
John M. Forsythe

Cooperative Institute for Research in the Atmosphere (CIRA)

Colorado State University (CSU)

Steven.fletcher@colostate.edu



2nd June 2015

Sensitivity Analysis and Data Assimilation Workshop



Outline of Talk

- 1) Mixed Lognormal-Gaussian Full Field VAR
- 2) Do we linearise the Bayesian problem or find the Bayesian problem for the increment?
- 3) Lognormal Incremental VAR - Geometric Tangent Linear Theory
- 4) Mixed Multiplicative-Additive Lognormal-Gaussian VAR
- 5) Examples with Lorenz 63 model
- 6) Comparison to a Gaussian incremental system
- 7) Conclusions

Lognormal Full Field Theory

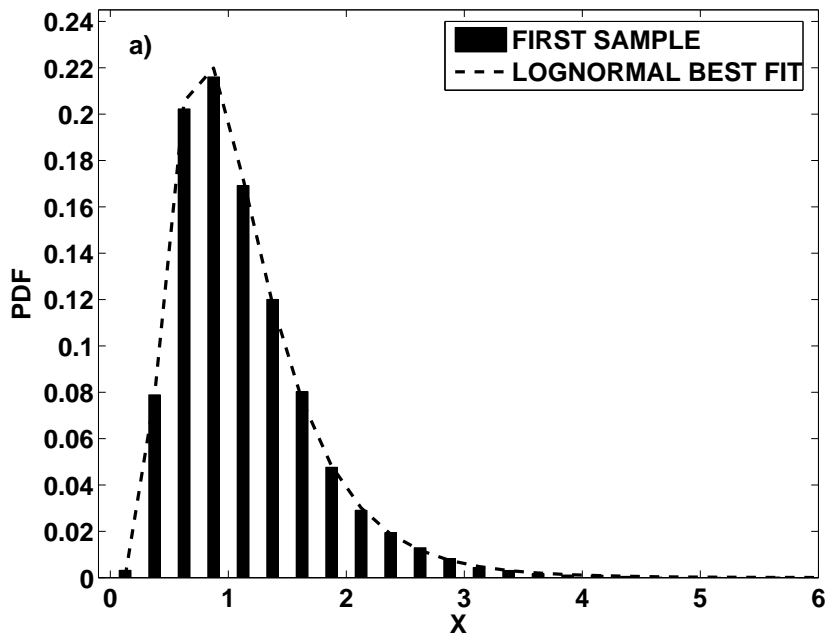
The starting point for lognormal based data assimilation is the definition of the errors. This was first proposed in Cohn (1997) where, due to the geometric behaviour of the lognormal distribution the errors could not be defined as the difference between two variables.

NOTE: There is no known distribution of the difference between two lognormally distributed independent random variables.

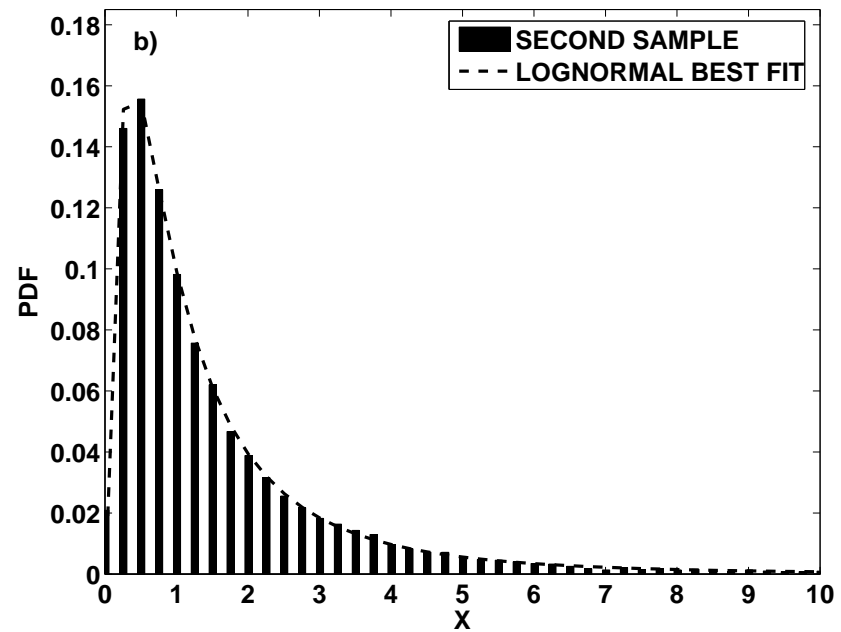
However, it is known that the distribution of the differences is **NEITHER** a Gaussian distribution nor a lognormal distribution.

Therefore we use the property that the ratio (and product) of two lognormal independent random variables is also a lognormally distributed random variable.

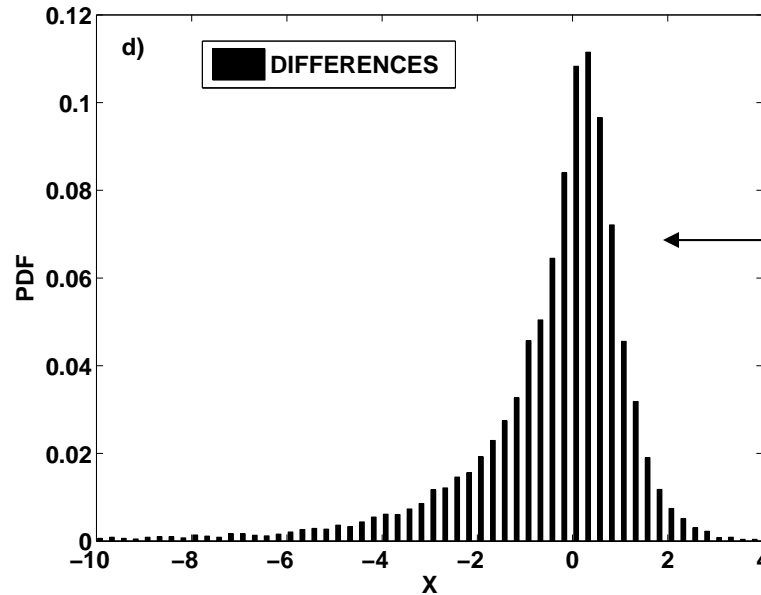
RANDOM LOGNORMAL SAMPLE WITH $\mu=0, \sigma=0.5$, SAMPLE=20,000



RANDOM LOGNORMAL SAMPLE WITH $\mu=0, \sigma=1$, SAMPLE=20,000



DIFFERENCE BETWEEN THE TWO LOGNORMAL SAMPLES



**ASSUMED
GAUSSIAN
APPROACH**

**NOTE: DIFFERENCE
IS NOT
A GAUSSIAN**

Fletcher (2010)

The lognormal distributed observational errors from Cohn (1997) are defined as

$$\varepsilon_{o,i} = \frac{y_i}{h_i(x)} \quad (1)$$

where $i = 1, \dots, N_o$, y_i are the observations, and $h_i(x)$ is the non-linear observation operator.

The definition for lognormally distributed background errors comes from Fletcher and Zupanski (2007) and they are defined by

$$\varepsilon_{b,j} = \frac{x_j^t}{x_{b,j}} \quad (2)$$

These are combined with the standard Gaussian error definition to form the mixed lognormal-Gaussian error definition (Fletcher and Zupanski, 2006b, 2007)

FULL FIELD MIXED GAUSSIAN-LOGNORMAL 3D AND 4D VAR

The resulting 3DVAR and 4DVAR cost functions are

$$J(x) = \frac{1}{2} \varepsilon_b^T B^{-1} \varepsilon_b + \varepsilon_b^T \begin{pmatrix} \mathbf{0}_{bp} \\ \mathbf{1}_{bq} \end{pmatrix} + \frac{1}{2} \varepsilon_o^T R^{-1} \varepsilon_o + \varepsilon_o^T \begin{pmatrix} \mathbf{0}_{op} \\ \mathbf{1}_{oq} \end{pmatrix} \quad (3)$$

where

$$\varepsilon_b = \begin{pmatrix} x_p^t - x_{bp} \\ \ln x_q^t - \ln x_{bq} \end{pmatrix} \text{ and } \varepsilon_o = \begin{pmatrix} y_{op} - h_{op}(x) \\ \ln y_{oq} - \ln h_{oq}(x) \end{pmatrix}$$

Fletcher and
Zupanski
(2007)

and

$$J(x_0) = \frac{1}{2} \varepsilon_{0b}^T B^{-1} \varepsilon_{0b} + \varepsilon_{0b}^T \begin{pmatrix} \mathbf{0}_{bp} \\ \mathbf{1}_{bp} \end{pmatrix} + \frac{1}{2} \sum_{i=1}^{N_o} \varepsilon_{0i}^T R_i^{-1} \varepsilon_{0i} + \sum_{i=1}^{N_o} \varepsilon_{0i}^T \begin{pmatrix} \mathbf{0}_{opi} \\ \mathbf{1}_{oqi} \end{pmatrix} \quad (4)$$

where

$$\varepsilon_{0b} = \begin{pmatrix} x_p^t(t_0) - x_{bp}(t_0) \\ \ln x_q^t(t_0) - \ln x_{bq}(t_0) \end{pmatrix} \text{ and } \varepsilon_{0i} = \begin{pmatrix} y_{pi} - h_{pi}(M_i(x(t_0))) \\ \ln y_{qi} - \ln h_{qi}(M_i(x(t_0))) \end{pmatrix}$$

Fletcher (2010)

Do we linearise the Bayesian problem or find the Bayesian problem for the increment?

In Song *et al.* (2012) the first version of a lognormal distribution incremental VAR is presented. The starting point in Song *et al.* (2012) is to define the relationship between the true state and the background as

$$x^t = x_b \circ e^{\Delta x} \quad (5)$$

NOTE: Eqn (5) is equivalent to a log-linearisation of the full field cost function. However, this raises the point that for both sides to be lognormal, the increment in (5) has to be Gaussian distributed.

Therefore, the associated cost function for the increment should be from a Gaussian distribution. Which is the median approach in Song *et al.* (2012) and shows positive results against an assumed Gaussian approach.

The problem here is that this increment is used to linearize the background component, which means that the Bayesian problem is not consistent for the Gaussian increment to make (5) lognormal.

Lognormal Incremental VAR

The starting point is to realize that unlike in the Gaussian framework we **can not** apply an **additive increment** if we want to maintain a lognormal framework, we have to define the relationship between the true state and the background as

$$x^t = x_b \circ \Delta x \quad (6)$$

This then allows the background component of the cost function to be

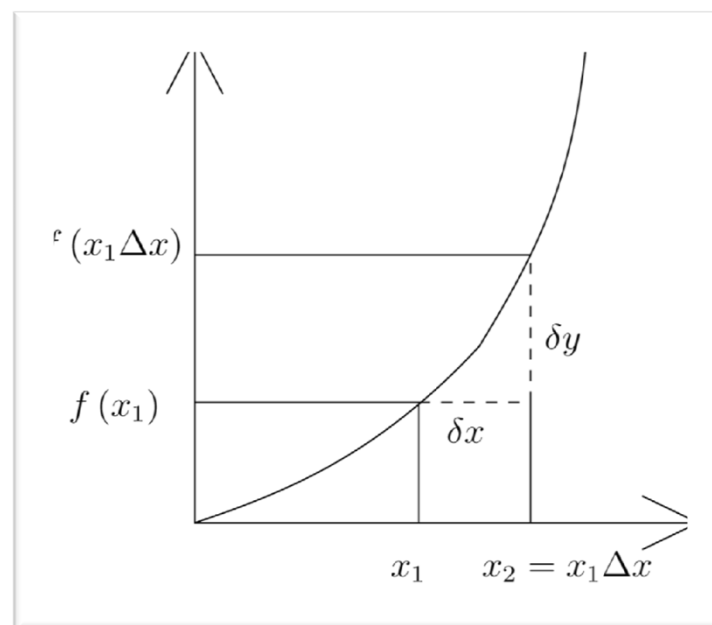
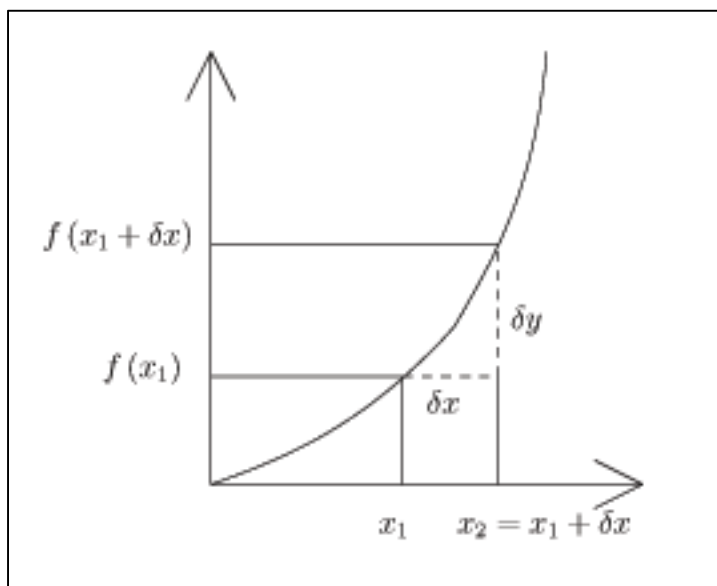
$$J_b(\Delta x) = \frac{1}{2} \ln \Delta x^T B^{-1} \ln \Delta x + \ln \Delta x^T \mathbf{1} \quad (7)$$

Therefore, (7) is consistent with a Bayesian problem for a lognormally distributed increment, **but is not a linearisation.**

Geometric Tangent Linear Modeling

We have to address how we are to linearise $h(x_b \circ \Delta x)$. Again as with the background we require the observational component to be as consistent with a lognormal distribution as possible.

Therefore, we can not simply use the standard tangent linear approximation. To overcome this problem, in Fletcher and Jones (2014), we define the Geometric Tangent Linear approximation as



Geometric Tangent Linear Modeling

This then allows us to be able to use standard derivative results for a multiplicative increment. Which means that for lognormal 3DVAR the observation operator can be approximated by

$$h(x_b \circ \Delta \mathbf{x}) \approx h(x_b) + \frac{\partial h}{\partial \mathbf{x}} x_b (\Delta \mathbf{x} - \mathbf{1}) \quad (7)$$

And for 4DVAR by

$$h_i(M_i(x_b(t_0) \circ \Delta \mathbf{x}(t_0))) \approx h(M_i(x_b(t_0))) + \frac{\partial h_i}{\partial \mathbf{x}} \frac{\partial M_i}{\partial \mathbf{x}} x_b(t_0) (\Delta \mathbf{x}(t_0) - \mathbf{1}) \quad (8)$$

Fletcher and Jones (2014).

Mixed Multiplicative-Additive Lognormal-Gaussian Incremental VAR

As we do not live in a just Gaussian or lognormal world, we have to combine the two approaches, as we have done for the full field (Fletcher, 2010). We therefore define our incremental vector as

$$\Delta \mathbf{x}_{\text{mx}} = \begin{pmatrix} \delta \mathbf{x}_{bp}(t_0) \\ \Delta \mathbf{x}_{bq}(t_0) \end{pmatrix} \quad (9)$$

This then gives the following 4DVAR cost function (Fletcher and Jones, 2014).

$$\begin{aligned}
J(\Delta x_{mx}) &= \frac{1}{2} \Delta x_{mx}^T \mathbf{B} \Delta x_{mx} + \Delta x_{mx}^T \begin{pmatrix} \mathbf{0}_{bp} \\ \mathbf{1}_{bq} \end{pmatrix} \\
&+ \frac{1}{2} \sum_{i=1}^{N_o} \left(\begin{array}{c} y_{opi} - h_{opi}(\mathbf{M}_i(\mathbf{x}_b(t_0))) - H_{opi} \bar{\mathbf{M}}_i \Delta x_{mx} \\ \ln y_{oqi} - \ln h_{oqi}(\mathbf{M}_i(\mathbf{x}_b(t_0))) - W_{oi}^{-1} H_{oqi} \bar{\mathbf{M}}_i \Delta x_{mx} \end{array} \right)^T \mathbf{R}_i^{-1} \\
&\quad \times \left(\begin{array}{c} y_{opi} - h_{opi}(\mathbf{M}_i(\mathbf{x}_b(t_0))) - H_{opi} \bar{\mathbf{M}}_i \Delta x_{mx} \\ \ln y_{oqi} - \ln h_{oqi}(\mathbf{M}_i(\mathbf{x}_b(t_0))) - W_{oi}^{-1} H_{oqi} \bar{\mathbf{M}}_i \Delta x_{mx} \end{array} \right) \\
&+ \left(\begin{array}{c} y_{opi} - h_{opi}(\mathbf{M}_i(\mathbf{x}_b(t_0))) - H_{opi} \bar{\mathbf{M}}_i \Delta x_{mx} \\ \ln y_{oqi} - \ln h_{oqi}(\mathbf{M}_i(\mathbf{x}_b(t_0))) - W_{oi}^{-1} H_{oqi} \bar{\mathbf{M}}_i \Delta x_{mx} \end{array} \right)^T \begin{pmatrix} \mathbf{0}_{opi} \\ \mathbf{1}_{oqi} \end{pmatrix}
\end{aligned}$$

Example with the Lorenz 63 model

The Lorenz model is given by the following non-linear system of three ordinary differential equation

$$\begin{aligned}\dot{x} &= \alpha(y - x) \\ \dot{y} &= \rho x - y - xz \\ \dot{z} &= xy - \beta z\end{aligned}$$

The system is linearized and then descretized using the modified Euler scheme. The adjoint of this scheme is calculated analytically.

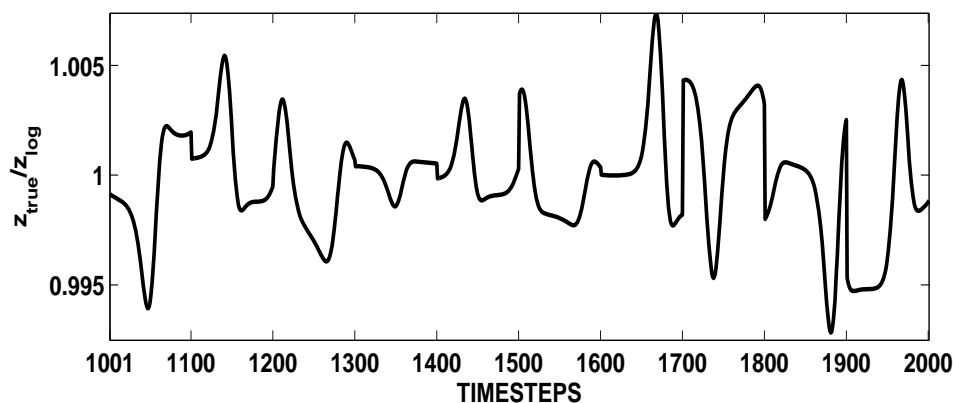
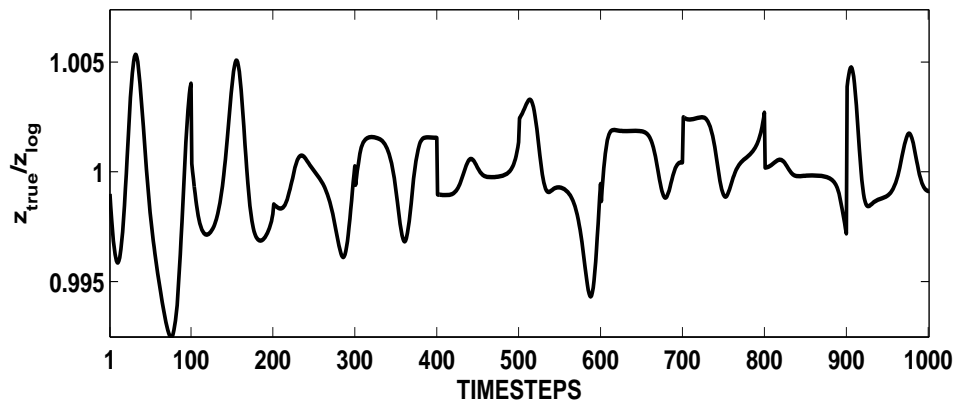
The observations are calculated by adding, or multiplying random perturbations from a Gaussian distribution, for the lognormal the background state is multiplied by the exponential of the scaled Gaussian increment

The initial conditions for the true solution, the initial background and the true increment at the initial time are

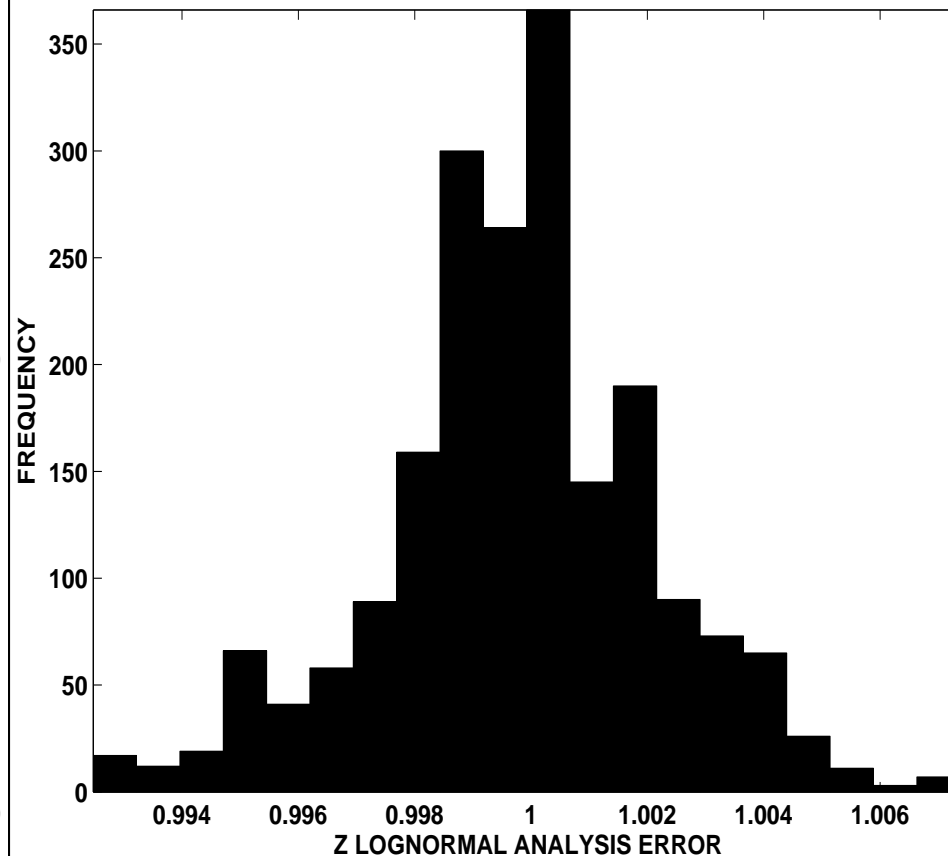
STATE	TRUE	BACKGROUND	INCREMENT
x	-5.4458	-5.9	0.4542
y	-5.4841	-5.0	-.4841
z	22.5606	24.0	0.94

Results for 20 cycles of 100ts with few accurate obs

Z ANALYSIS ERROR FOR $N_o=5$, $CYC=20$, $\sigma_{ox,y}=0.1$, $\sigma_{oz}=0.0025$



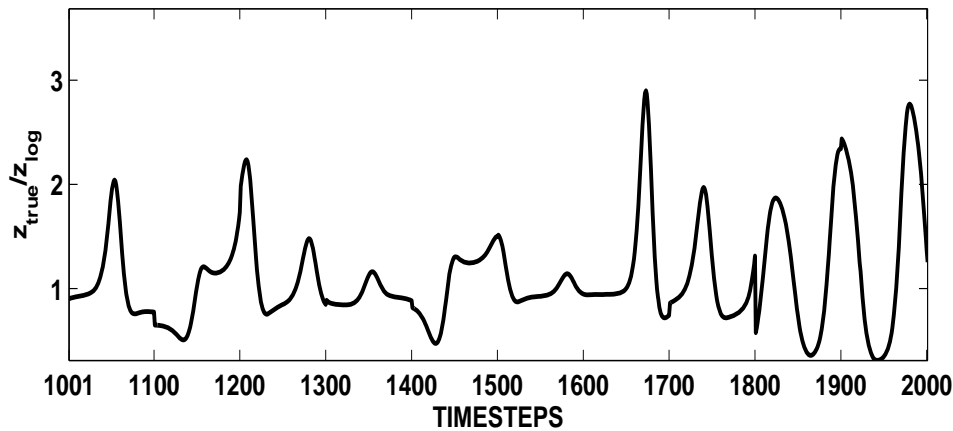
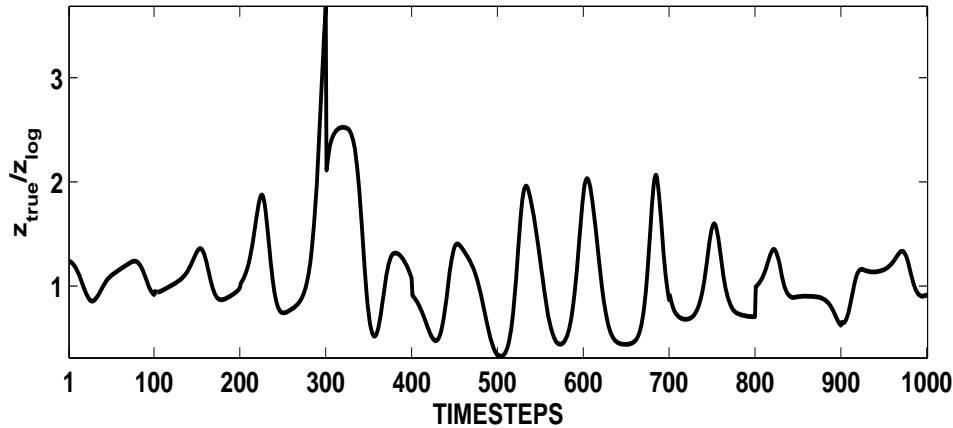
Z ANALYSIS ERROR DISTRIBUTION FOR $N_o=5$, $CYC=20$, $\sigma_{ox,y}=0.1$, $\sigma_{oz}=0.0025$



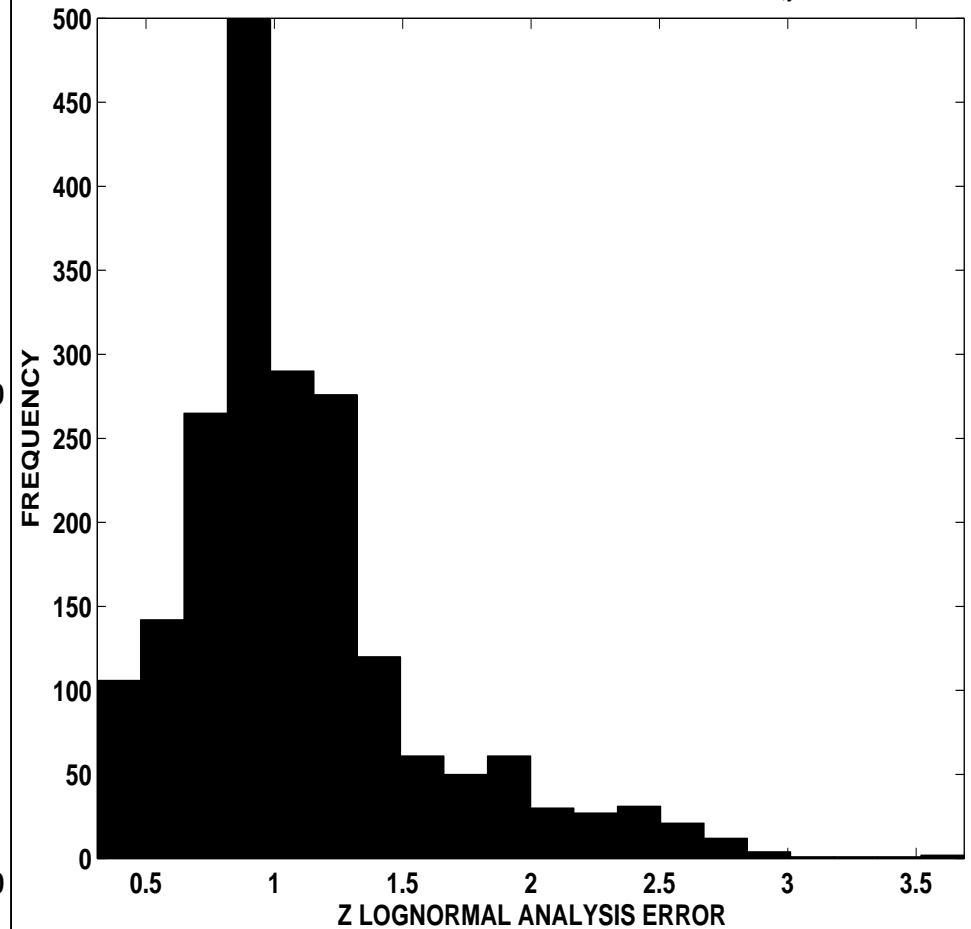
Fletcher and Jones (2014)

Results with same window lengths and same number of obs but less accurate

Z ANALYSIS ERROR FOR $N_0=5$, $CYC=20$, $\sigma_{ox,y}=5$, $\sigma_{oz}=0.125$

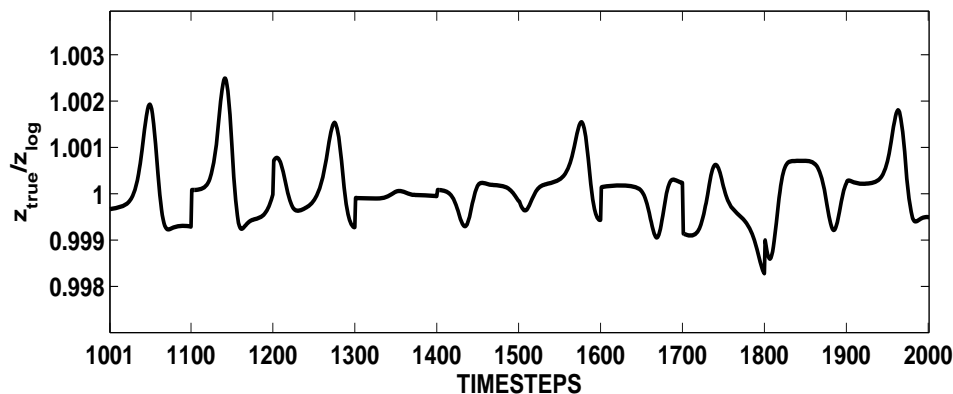
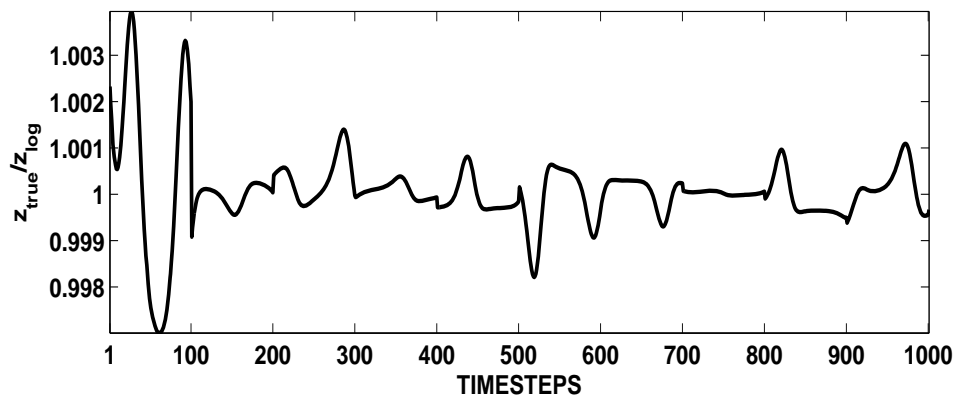


Z ANALYSIS ERROR DISTRIBUTION FOR $N_0=5$, $CYC=20$, $\sigma_{ox,y}=5$, $\sigma_{oz}=0.125$

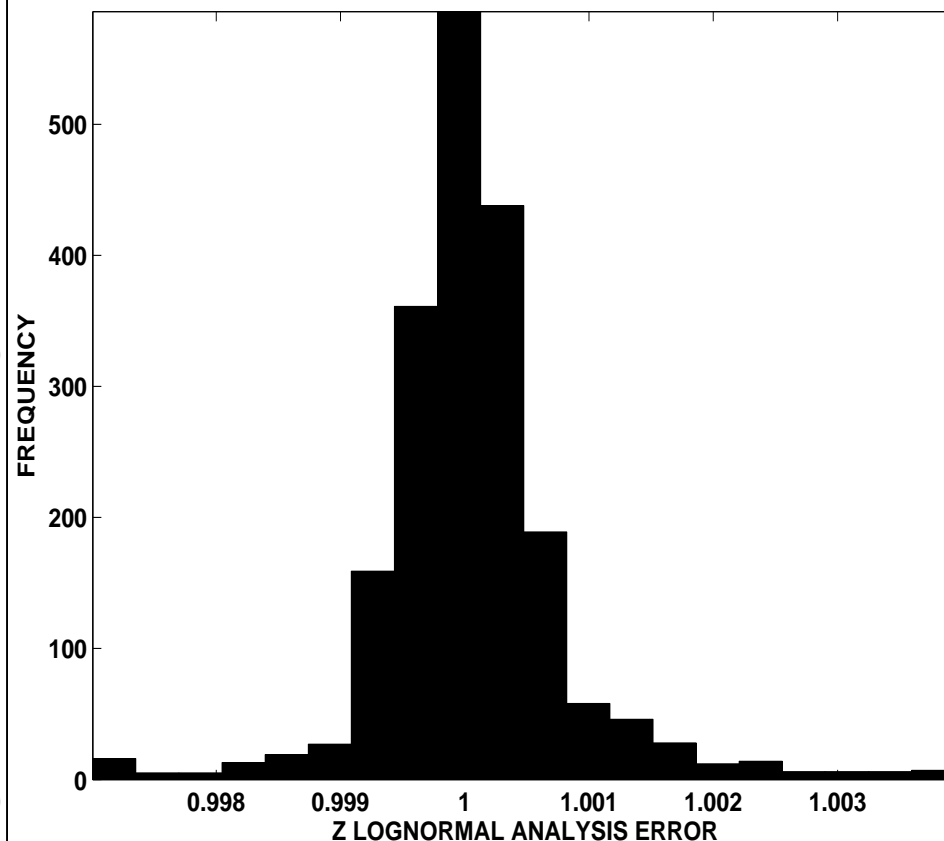


Results for same number of assimilation windows but with accurate observations every other time step

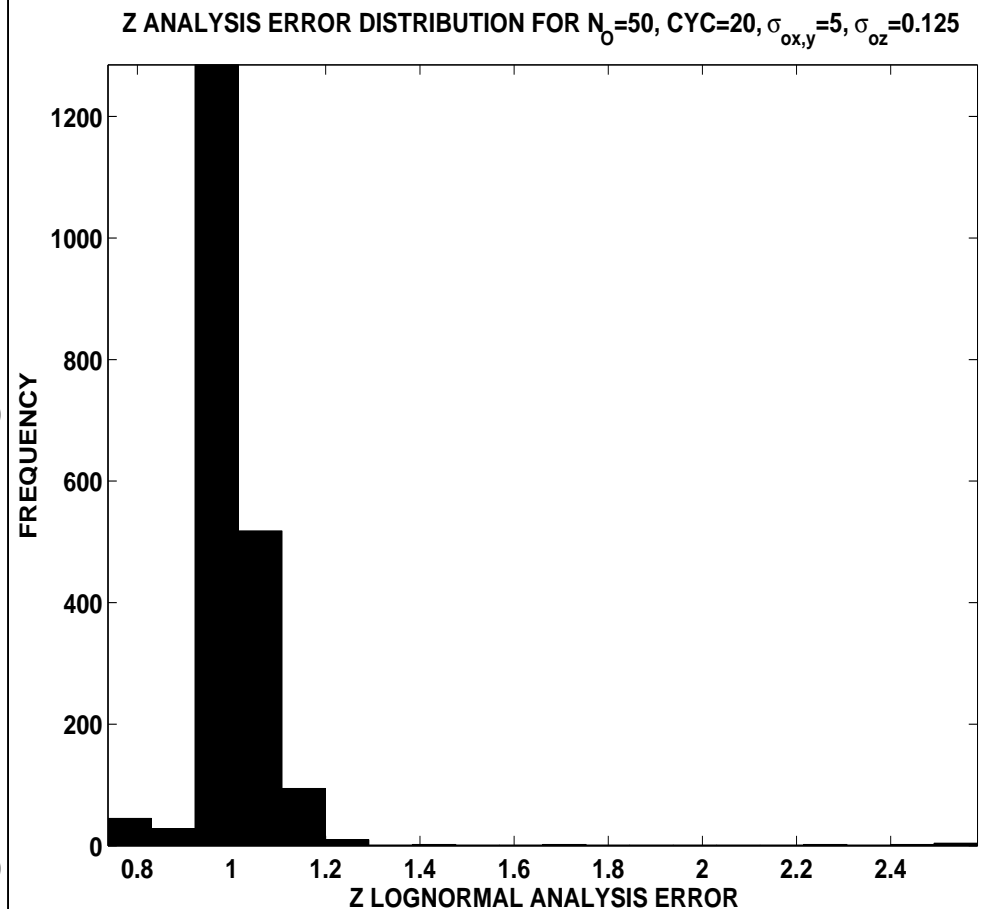
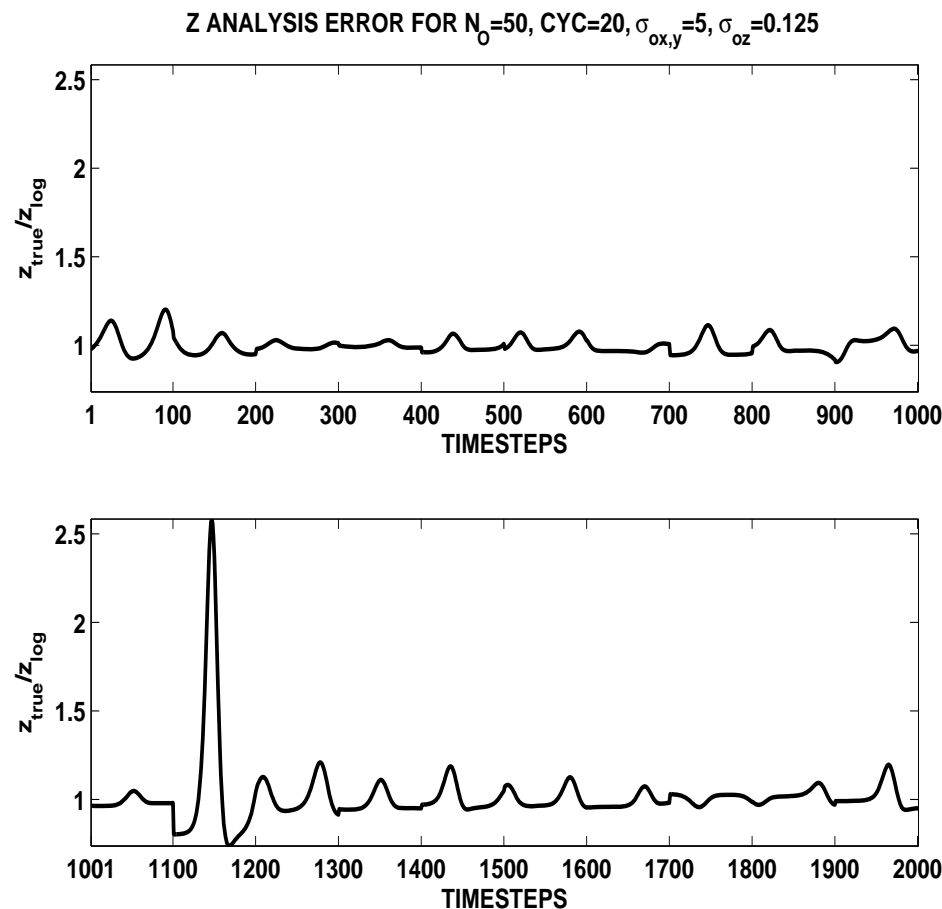
Z ANALYSIS ERROR FOR $N_o=50$, $CYC=20$, $\sigma_{ox,y}=0.1$, $\sigma_{oz}=0.0025$



Z ANALYSIS ERROR DISTRIBUTION FOR $N_o=50$, $CYC=20$, $\sigma_{ox,y}=0.1$, $\sigma_{oz}=0.0025$

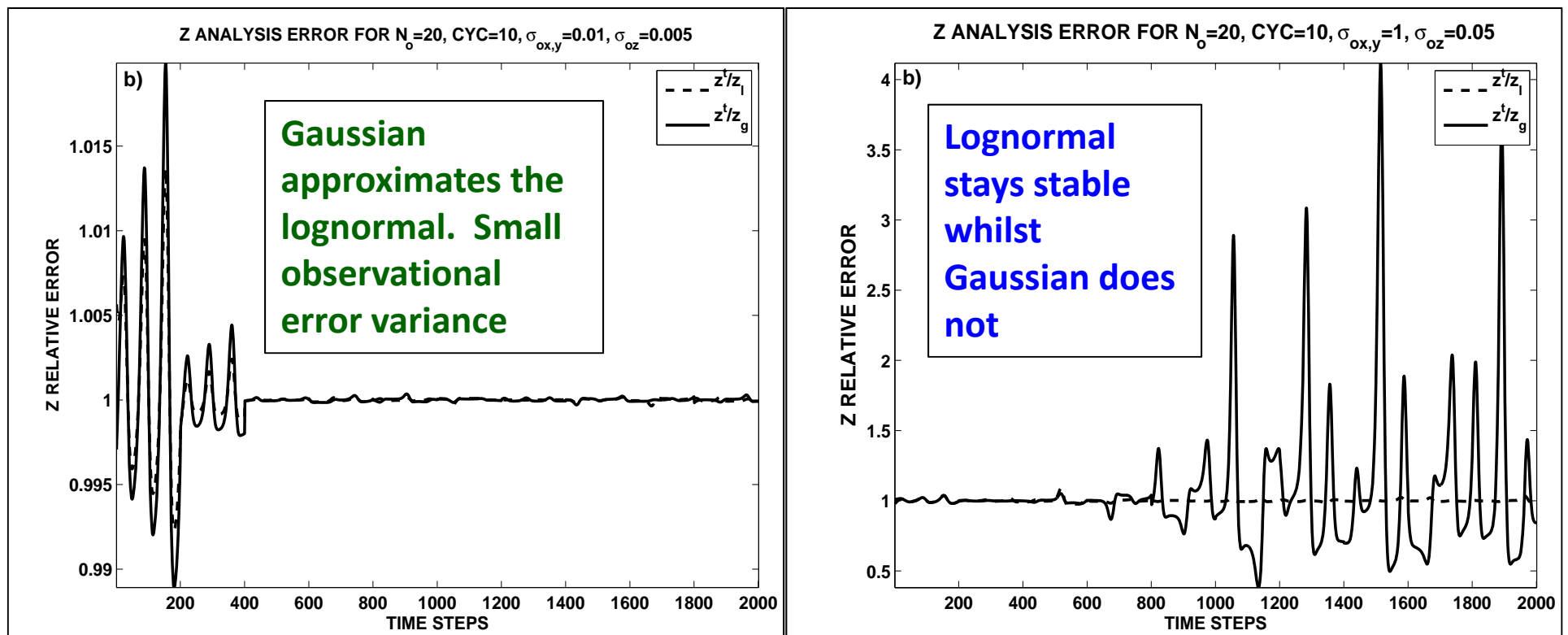


Results for same number of assimilation windows but with accurate observations every other time step



Comparison to a full Gaussian incremental system

These results are from Fletcher and Jones (2014) paper where here we are presenting results from two of the experiments, the first to highlight where the two systems are similar and the case where the lognormal converges but the Gaussian does not.



Conclusions

- It is possible to define an incremental version of the lognormal full field 3D and 4DVAR.
- This is possible through relating a multiplicative increment to an additive one from the gradient definition.
- Have been able to test in a Lorenz 63 model with just inner loops and no QC of observations and monitoring of the tangent linear assumption
- Tested with different number of observations with different variances and with different assimilation window lengths
- Next is to apply this theory in the WRF-GSI system which does not need the adjoint of the WRF model
- Have applied the full field mixed distribution formulation in a temperature-mixing ratio microwave retrieval system and have shown positive results compared to a Gaussian only version (Kliwer *et al.*, 2015).