Non-Gaussian Detection using Machine Learning with Data Assimilation Applications

Michael Goodliff¹, Steven James Fletcher², Anton Kliewer³, Andrew S. Jones³, and John M Forsythe³

¹University of Colorado ²Cooperative Institute for Reserach in the Atmosphere ³Colorado State University

November 21, 2022

Abstract

In most data assimilation and numerical weather prediction systems, the Gaussian assumption is prevalent for the behaviour of the random variables/errors that are involved. At the Cooperative Institute for Research in the Atmosphere (CIRA) theory has been developed for different forms of variational data assimilation schemes that enables the Gaussian assumption to be relaxed. For certain variable types, a lognormally distributed random variable can be combined in a mixed Gaussian-lognormal distribution to better capture the interactions of the errors of different distributions. However, assuming that a distribution can change in time, then developing techniques to know when to switch between different versions of the data assimilation schemes becomes very important. Given this ability to change the formulation of the data assimilation system enable us to select the more optimal scheme for the different distributed situations.

In this paper, we present results with a machine learning technique (the support vector machine) to switch between data assimilation methods based on the detection of a change in the Lorenz 1963 model's \$z\$ component's probability distribution. Given the machine learning technique's detection/prediction of a change in the distribution, then either a Gaussian or a mixed Gaussian-lognormal 3DVar based cost function is used to minimise the errors in this period of time. It is shown that switching from a Gaussian 3DVar to a lognormal 3DVar at lognormally-distributed parts of the attractor improves the data assimilation analysis error compared to using one distribution type for the entire attractor.

Non-Gaussian Detection using Machine Learning with Data Assimilation Applications

Michael R. Goodliff^{1,2}, Steven J. Fletcher³, Anton J. Kliewer³, Andrew S. Jones³, John M. Forsythe³

¹Cooperative Institute for Research in Environmental Sciences (CIRES) at the University of Colorado Boulder ²National Oceanic and Atmospheric Administration (NOAA) Physical Sciences Laboratory (PSL) ³Cooperative Institute for Research in the Atmosphere, Colorado State University, Fort Collins, USA

Key Points:

- ¹⁰ Data Assimilation
 - Machine Learning
 - Gaussianity

3

5 6

8

9

11

12

 $Corresponding \ author: \ Michael \ Goodliff, \verb"michael.goodliff@noaa.gov" \\$

13 Abstract

In most data assimilation and numerical weather prediction systems, the Gaussian 14 assumption is prevalent for the behaviour of the random variables/errors that are in-15 volved. At the Cooperative Institute for Research in the Atmosphere (CIRA) theory 16 has been developed for different forms of variational data assimilation schemes that en-17 ables the Gaussian assumption to be relaxed. For certain variable types, a lognormally 18 distributed random variable can be combined in a mixed Gaussian-lognormal distribu-19 tion to better capture the interactions of the errors of different distributions. However, 20 assuming that a distribution can change in time, then developing techniques to know 21 when to switch between different versions of the data assimilation schemes becomes 22 very important. Given this ability to change the formulation of the data assimila-23 tion system enable us to select the more optimal scheme for the different distributed 24 situations. 25

In this paper, we present results with a machine learning technique (the support 26 vector machine) to switch between data assimilation methods based on the detection 27 of a change in the Lorenz 1963 model's z component's probability distribution. Given 28 the machine learning technique's detection/prediction of a change in the distribution, 29 then either a Gaussian or a mixed Gaussian-lognormal 3DVar based cost function is 30 used to minimise the errors in this period of time. It is shown that switching from a 31 Gaussian 3DVar to a lognormal 3DVar at lognormally-distributed parts of the attractor 32 improves the data assimilation analysis error compared to using one distribution type 33 for the entire attractor. 34

35 1 Introduction

The assumption that variables, and their errors, are Gaussian distributed is com-36 monplace in areas such as numerical weather prediction and modelling. Research such 37 as that undertaken by Perron and Sura in Perron & Sura (2013) has shown that this 38 assumption is generally false for atmospheric variables, and that Gaussian variables 39 in the atmosphere are rare. The aforementioned statement was based on a sixty two 40 year long project from daily data taken from the National Centers for Environmental 41 Prediction and the National Center for Atmospheric Research (NCEP-NCAR), using 42 the Reanalysis I Project data set. Given this evidence, the need to be able to relax 43 the Gaussian assumption for the errors involved in the data assimilation schemes be-44 comes quite important if the analysis error is to be minimised, and thus a possible 45 improvement in the subsequent forecast. 46

Most of the current formulations of data assimilation, for example variational 47 methods such as 3DVar and 4DVar (which are based upon Bayes theorem Fletcher 48 (2017)), and ensemble methods such as the Ensemble Kalman Filter (ENKF) Evensen 49 & Van Leeuwen (1996), the (local) Ensemble Transform Kalman Filter ((L)ETKF) Ott 50 et al. (2004); Wang & Bishop (2003), which are based upon a control theory/weighted 51 least squares approach using ensemble members to approximate the analysis mean and 52 covariances, and the Maximum Likelihood Ensemble Kalman Filter (MLEF), Zupanski 53 (2005), which uses the Kalman filter equations combined with the 3DVar cost function, 54 all assume that the errors involved are Gaussian distributed. Other papers who look 55 into non-Gaussian data assimilation methods are local particle filters, van Leeuwen et 56 al. (2019), and Amezcua & Leeuwen (2014) who looked into Gaussian anamorphosis 57 on the EnKF. 58

However, at the Cooperative Institute for Research in the Atmosphere (CIRA) at
 Colorado State University (CSU), there has been theory that has been developed, and
 tested, that allows for the Gaussian assumption for the distribution of the errors to be
 relaxed to a lognormal distribution. In Fletcher & Zupanski (2006a) the theory is pre sented for the case where there are lognormal observational errors in 3D. In Fletcher

& Zupanski (2006b) a mixed Gaussian-lognormal distribution is presented, and an as-64 sociated cost function that allows for the simultaneous minimisation of Gaussian and 65 lognormal errors is presented. The mixed approach was extended to the background 66 term in Fletcher & Zupanski (2007), and then tested with the Lorenz 1963 model 67 (Lorenz, 1963), Lorenz 63 hereafter, where it is shown here that the z component of 68 this model is not Gaussian distributed. The mixed distribution theory was extended to 69 a 4DVar type system in Fletcher (2010), and eventually shown for incremental 3DVar 70 and 4DVar in Fletcher & Jones (2014). In these papers it is shown that the lognormal 71 variant of 3DVar and 4DVar showed improvements in analysis accuracy over the tra-72 ditional Gaussian, and logarithmic transforms method for the z component, but that 73 there was also improvement in the analysis error for the x and y components, where 74 the errors associated with these components were assumed to be Gaussian distributed. 75

However, as shown in the first part of Goodliff et al. (2020), the trajectory of the 76 z component of the Lorenz 63 model changes distributions on different parts of the 77 underlying attractor, and as such if the data assimilation is to be optimised then these 78 changes need to be used to *switch* from the Gaussian to the mixed distribution-based 79 cost functions. In the second part of Goodliff et al. (2020) a support vector machine 80 and a neural network machine learning techniques were tested with the Lorenz 63 81 model to detect non-Gaussian behaviour. It was shown that these techniques were 82 very capable of detecting skewness, and differences in descriptive statistics, in order 83 to estimate, and predict, non-Gaussianity. 84

Recently, machine learning methods have become very popular in atmospheric 85 sciences, especially in areas such as numerical weather prediction and modelling (Scher 86 & Messori, 2018) to help find biases and correlations in data, and also to help reduce 87 analysis and forecast errors. Pasini and Pelino ((Pasini & Pelino, 2005) and Pasini 88 (2008)) used two Lorenz 63 attractors to analyse predictability. There have been many 89 other studies using machine learning methods to try and improve weather forecasting 90 and climate modelling and the reader is referred to Dueben & Bauer (2018); Rasp & 91 Lerch (2018); Scher (2018); Scher & Messori (2019); Weyn et al. (2019) for some of 92 these extra examples. 93

Given the progress made with machine learning techniques, and the need identi-94 fied above to be able to inform a data assimilation scheme to switch between different 95 versions of the cost function, this paper investigates a support vector machine, which is a supervised machine learning algorithm, to detect non-Gaussian probability den-97 sity functions in the Lorenz 63 model. This approach is applied to the z component 98 of the Lorenz attractor, where the skewness of said z variable is the target data, and 99 the x and y components of the attractor are our training data. We use this to then 100 apply a "switch" to our data assimilation method to change between a Gaussian fits 101 all cost function to a mixed Gaussian-lognormal based cost function, where the x and 102 y components are assumed to have Gaussian error throughout the experiment, and 103 that the z component is *switching* between a Gaussian and a lognormal distribution. 104 This switch changes the data assimilation methodology from the traditional Gaussian 105 3DVar to a mixed Gaussian-lognormal variant of 3DVar Fletcher & Zupanski (2007) 106 based on the distribution estimation given by the support vector machine. 107

The remainder of the paper is organised as follows: Section 2 will start with an overview of the Bayesian model for the variational data assimilation theories and the methods used. We will also discuss the machine learning methodology and how we use it with the data assimilation. Section 3 describes the Lorenz 63 model, and section 4 shows the experimentation using the mixed DA/ML scheme on the Lorenz 63 model to improve forecasts. The paper is concluded in section 5.

114 2 Methodology

117

123

¹¹⁵ In this section we shall present the different data assimilation and machine learn-¹¹⁶ing techniques that are used in the results presented later.

2.1 Traditional 3DVar (3DVar-G)

Variational data assimilation methods estimate the most probable state of the system, which is the mode of the posterior probability distribution function. In traditional 3DVar, this comes from Gaussian statistics and is a combination of the background and the likelihood, with the background term written as:

$$p\left(\mathbf{x}\right) = \frac{1}{\left|\mathbf{B}_{c}\right|^{\frac{1}{2}} \left(2\pi\right)^{\frac{N}{2}}} \exp\left(-\frac{1}{2}\left(\mathbf{x} - \mathbf{x}^{b}\right)^{\top} \mathbf{B}_{c}^{-1} \left(\mathbf{x} - \mathbf{x}^{b}\right)\right),\tag{1}$$

where the background state, and initial state that is sought, are given by \mathbf{x}^{b} and \mathbf{x} , respectively, $\mathbf{B}_{c} \in \mathcal{R}^{N \times N}$ is the background error covariance matrix, and N is the total number of background variables. The likelihood for Gaussian errors is defined as

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{|\mathbf{R}|^{\frac{1}{2}} (2\pi)^{\frac{N_o}{2}}} \exp\left(-\frac{1}{2} \left(\mathbf{y} - \mathbf{h}(\mathbf{x})\right)^{\top} \mathbf{R}^{-1} \left(\mathbf{y} - \mathbf{h}(\mathbf{x})\right)\right), \qquad (2)$$

where \mathbf{y} is the observation, $\mathbf{h}(\mathbf{x})$ is the (non)-linear observation operator, $\mathbf{R} \in \mathcal{R}^{N_o \times N_o}$ is the observational error covariance matrix, and N_o is the total number of observations. The next step is to substitute (1) and (2) into Bayes' theorem

$$p(\mathbf{x}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{x}) p(\mathbf{x}),$$
 (3)

and seek the state that maximises the probability in (3). However, it is quite often easier to work with the equivalent problem that seeks the state that minimises the negative log-likelihood of (3), which for the Gaussian definitions presented above results in the following cost function,

$$J(x) = \frac{1}{2} \left(\mathbf{x} - \mathbf{x}^{b} \right)^{\top} \mathbf{B}_{c}^{-1} \left(\mathbf{x} - \mathbf{x}^{b} \right) + \frac{1}{2} \left(\mathbf{y} - \mathbf{h} \left(\mathbf{x} \right) \right)^{\top} \mathbf{R}^{-1} \left(\mathbf{y} - \mathbf{h} \left(\mathbf{x} \right) \right), \qquad (4)$$

that has to be minimised.

2.2 Mixed Gaussian-Lognormal 3DVar (3DVar-Mix)

The mixed Gaussian-lognormal 3DVar data assimilation scheme was first presented in Fletcher & Zupanski (2007) for both the background and likelihood components. This version of 3DVar uses a multivariate lognormal distribution based cost function for lognormal random variables that is derived through using a similar approach for Bayes theorem as presented above. Thus, for the lognormal approach the a priori probability density function is given by

$$p(\mathbf{x}) = \left(\prod_{i=1}^{N} \frac{1}{x_i}\right) \frac{1}{|\mathbf{B}_L|^{\frac{1}{2}} (2\pi)^{\frac{N}{2}}} \exp\left(-\frac{1}{2} \left(\ln \mathbf{x} - \ln \mathbf{x}^b\right)^\top \mathbf{B}_L^{-1} \left(\ln \mathbf{x} - \ln \mathbf{x}^b\right)\right),$$
(5)

where \mathbf{B}_L is the lognormal based background error covariance matrix, which is defined in terms of expectations of $\ln \mathbf{x}$ and not \mathbf{x} . The equivalent likelihood distribution for lognormal errors is given by

$$p\left(\mathbf{y}|\mathbf{x}\right) = \left(\prod_{i=1}^{N_o} \frac{\left(\mathbf{h}\left(\mathbf{x}\right)\right)_i}{\mathbf{y}_i}\right) \frac{1}{|\mathbf{R}|^{\frac{1}{2}} (2\pi)^{\frac{N_o}{2}}} \exp\left(-\frac{1}{2} \left(\ln \mathbf{y} - \ln \mathbf{h}\left(\mathbf{x}\right)\right)^\top \mathbf{R}_L^{-1} \left(\ln \mathbf{y} - \ln \mathbf{h}\left(\mathbf{x}\right)\right)\right).$$
(6)

This then results in the lognormal 3DVar cost function given by

$$J(x) = \frac{1}{2} \left(\ln \mathbf{x} - \ln \mathbf{x}^{b} \right)^{\top} \mathbf{B}_{L}^{-1} \left(\ln \mathbf{x} - \ln \mathbf{x}^{b} \right) + \left(\ln \mathbf{x} - \ln \mathbf{x}^{b} \right)^{\top} \mathbf{1}_{N} + \frac{1}{2} \left(\ln \mathbf{y} - \ln \mathbf{h}(\mathbf{x}) \right)^{\top} \mathbf{R}_{L}^{-1} \left(\ln \mathbf{y} - \ln \mathbf{h}(\mathbf{x}) \right) + \left(\ln \mathbf{y} - \ln \mathbf{h}(\mathbf{x}) \right)^{\top} \mathbf{1}_{N_{o}}.$$
(7)

Minimising this cost function gives us the solution to the lognormal 3DVar. For in depth information about lognormal 3DVar, refer to Fletcher & Zupanski (2007) and Fletcher (2010).

However, in the results that will be presented later in this paper, the mixed Gaussian-lognormal approach is utilised which comes from the mixed Gaussian-lognormal probability density function derived in Fletcher & Zupanski (2006b) where the multi-variate distribution that is used for the a priori distribution is given by

$$p\left(\mathbf{x}\right) = \left(\prod_{i=p+1}^{N} \frac{1}{x_{i}}\right) \frac{1}{|\mathbf{B}_{mx}|^{\frac{1}{2}} (2\pi)^{\frac{N}{2}}} \times \exp\left(-\frac{1}{2} \left(\begin{array}{c} \mathbf{x}_{p} - \mathbf{x}_{bp} \\ \ln \mathbf{x}_{q} - \ln \mathbf{x}_{bq} \end{array}\right)^{\top} \mathbf{B}_{mx}^{-1} \left(\begin{array}{c} \mathbf{x}_{p} - \mathbf{x}_{bp} \\ \ln \mathbf{x}_{q} - \ln \mathbf{x}_{bq} \end{array}\right)\right),$$
(8)

where p is the number of Gaussian random variables, q is the number of lognormal random variables, such that N = p+q. The mixed distribution error covariance matrix here is defined as

$$\mathbf{B}_{mx} \equiv \begin{pmatrix} \left(\boldsymbol{\varepsilon}_{bp}^{G}\right) \left(\boldsymbol{\varepsilon}_{bp}^{G}\right)^{\mathsf{T}} & \left(\boldsymbol{\varepsilon}_{bp}^{G}\right)^{\mathsf{T}} & \left(\boldsymbol{\varepsilon}_{bq}^{L}\right)^{\mathsf{T}} \\ \left(\boldsymbol{\varepsilon}_{bq}^{L}\right) \left(\boldsymbol{\varepsilon}_{bp}^{G}\right)^{\mathsf{T}} & \left(\boldsymbol{\varepsilon}_{bq}^{L}\right) \left(\boldsymbol{\varepsilon}_{bq}^{L}\right)^{\mathsf{T}} \end{pmatrix}.$$
$$\boldsymbol{\varepsilon}_{bp}^{G} \equiv \mathbf{x}_{p} - \mathbf{x}_{bp}, \quad \boldsymbol{\varepsilon}_{bq}^{L} \equiv \ln \mathbf{x}_{q} - \ln \mathbf{x}_{bq}, \tag{9}$$

where

and the superscripts G and L denote the Gaussian and lognormal components. The mixed Gaussian-lognormal distribution that would be used for the likelihood of Gaussian and lognormal errors is given by

$$p\left(\mathbf{y}|\mathbf{x}\right) \equiv \left(\prod_{i=p+1}^{N} \frac{\mathbf{h}_{i}\left(\mathbf{x}\right)}{\mathbf{y}_{i}}\right) \frac{1}{|\mathbf{R}_{\mathbf{mx}}|^{\frac{1}{2}} (2\pi)^{\frac{N_{o}}{2}}} \times \exp\left(-\frac{1}{2} \left(\begin{array}{c} \mathbf{y}_{p} - \mathbf{h}_{p}\left(\mathbf{x}\right)\\ \ln \mathbf{y}_{q} - \ln \mathbf{h}_{q}\left(\mathbf{x}\right) \end{array}\right)^{\top} \mathbf{R}_{mx}^{-1} \left(\begin{array}{c} \mathbf{y}_{p} - \mathbf{h}_{p}\left(\mathbf{x}\right)\\ \ln \mathbf{y}_{q} - \ln \mathbf{h}_{q}\left(\mathbf{x}\right) \end{array}\right)\right),$$
(10)

where the observation covariance matrix is assumed to be diagonal and the associated variances in these entries are calculated as per their distribution that they are associated with. Thus the associated mixed Gaussian-lognormal cost function is given by

$$J(\mathbf{x}) = \frac{1}{2} \begin{pmatrix} \mathbf{x}_{p} - \mathbf{x}_{bp} \\ \ln \mathbf{x}_{q} - \ln \mathbf{x}_{bq} \end{pmatrix}^{\top} \mathbf{B}_{mx}^{-1} \begin{pmatrix} \mathbf{x}_{p} - \mathbf{x}_{bp} \\ \ln \mathbf{x}_{q} - \ln \mathbf{x}_{bq} \end{pmatrix} + \begin{pmatrix} \mathbf{x}_{p} - \mathbf{x}_{bp} \\ \ln \mathbf{x}_{q} - \ln \mathbf{x}_{bq} \end{pmatrix}^{\top} \begin{pmatrix} \mathbf{0}_{p} \\ \mathbf{1}_{q} \end{pmatrix} + \frac{1}{2} \begin{pmatrix} \mathbf{y}_{p} - \mathbf{h}_{p}(\mathbf{x}) \\ \ln \mathbf{y}_{q} - \ln \mathbf{h}_{q}(\mathbf{x}) \end{pmatrix}^{\top} \mathbf{R}_{mx}^{-1} \begin{pmatrix} \mathbf{y}_{p} - \mathbf{h}_{p}(\mathbf{x}) \\ \ln \mathbf{y}_{q} - \ln \mathbf{h}_{q}(\mathbf{x}) \end{pmatrix} + \begin{pmatrix} \mathbf{y}_{p} - \mathbf{h}_{p}(\mathbf{x}) \\ \ln \mathbf{y}_{q} - \ln \mathbf{h}_{q}(\mathbf{x}) \end{pmatrix}^{\top} \begin{pmatrix} \mathbf{0}_{p} \\ \mathbf{1}_{q} \end{pmatrix}.$$

$$(11)$$

It should be noted that it could well be the case that the number of state variables that are Gaussian or lognormal may not be the same as those of the observational errors. Another important feature to note here is that the mode of the mixed distribution is a function of the error covariance matrices, and it is shown in Fletcher (2017) that the Gaussian components become a function of the covariances with the lognormal components, which then enables a relationship between the Gaussian and lognormal components, which is not present in the mode of a Gaussian fits all approach.

134

2.3 Mixing Machine Learning into Data Assimilation (3DVar-ML)

The machine learning method used in this study to classify our data is the Support Vector Machine (Nello & Shawe-Taylor, 2000). This method separates classified training data with a hyperplane. The support vector machine is a method of supervised learning where we supply a training set and a target set to train our model. In this experiment, we use the radial basic function (RBF) kernel, and train the machine learning algorithm for 50,000 time steps of the Lorenz 63 model.

In this experiment, we predict the probability density function of the z component of the Lorenz 63 model based on the values of x and y components of the model. Using x and y as the training data and the z-score of the target data, we have shown (Goodliff et al., 2020) that we can be highly precise in our predictions of the distribution of z. The z-score (skewness statistic $\sqrt{\beta_1}$) is calculated and estimated by methods shown in D'agostino et al. (1990) from the standardised skewness:

$$\sqrt{\beta_1} = \frac{E(X-\mu)^3}{\sigma^3} \tag{12}$$

where μ and σ are the mean and standard distribution, respectively. Here, a negative z-score represents a left (negative) skewed distribution, a positive z-score represents a right (positive) skewed distribution, and a z-score of zero refers to a symmetric distribution.

The window length in this study refers to the data around the observation point (example, a window length of 11 will be the data 5 points either side of the observation + the observation point). The z-score affects observation generation (see below) and which version of 3DVar the machine learning algorithm will choose at each observation point.

Through using the support vector machine to detect the probability of the trajectory, this enables us to utilise this as a switch to decide which data assimilation method is best at the current point in time. The optimal data assimilation method is used with the machine learning prediction as:

$$Method = \begin{cases} 3DVar-G, & \text{if } z\text{-score} < 1, \\ 3DVar-Mix, & \text{otherwise.} \end{cases}$$
(13)

150 3 Lorenz 63

As mentioned in the introduction, for the study that we shall present in the next section, we will be using the Lorenz 1963 model (Lorenz, 1963). This model is a good choice due to it's simplicity for a dynamic model which also exhibits chaotic behaviour. The model is very sensitive to the initial conditions from which it starts, and as such can give very different answers even by being out by a few decimal places from the true state, Fletcher (2017). These model equations are as given by

$$\frac{dx}{dt} = -\sigma \left(x - y \right), \tag{14}$$

$$\frac{dy}{dt} = \rho x - y - zx, \tag{15}$$

$$\frac{dz}{dt} = xy - \beta z, \tag{16}$$

where x = x(t), y = y(t), z = z(t) are the state variables (where t is time) and $\sigma = 10$, $\rho = 28$ and $\beta = 8/3$ are parameters. We start the machine learning training, the true run, and the data assimilation from different initial states on the attractor.

¹⁶⁰ 4 Experimentation

The experimentation starts with running the support vector machine algorithm on the Lorenz 63 model. We train the support vector machine on variables x and y, and the skewness of z as the output (target) data. The training of the machine learning method is performed over 50,000 time steps to obtain a somewhat robust fit for the system. This approach is based on the method in Goodliff et al. (2020).

To generate the observations, we use the machine learning fit to determine the probability function at observation time. If the observation is on a positively skewed area of the attractor, that is to say that the z-score ≥ 1 , the observation is generated using a lognormal distribution function, else, our observations are generated from a Gaussian distribution. Thus the observations for the three components of the Lorenz 63 model are of the form:

$$obs_x = x_t + G_x(0, \sigma_{xx}), \qquad (17)$$

$$obs_y = x_t + G_y(0, \sigma_{yy}), \qquad (18)$$

$$obs_z = x_t * exp\left(G_z\left(0, \sigma_{zz}\right)\right), \tag{19}$$

where $obs_{x,y,z}$ are the observations, x_t is the truth, and $G_{x,y,z}(0, \sigma_{xx,yy,zz})$ is a Gaussian based random number generated with a standard deviation σ . The square of these standard deviations, the variance, will form the diagonal entries of the observational error covariance matrices, where we are assuming that the observations are uncorrelated, and as such the **R** matrices will only be diagonal. In this study, **R** = 1

We then run our three data assimilation schemes: 3DVar-G, 3DVar-Mix, and 3DVar-ML. Each method is run over 5000 time steps, with 50 runs. Running the system for this long negates any biases generated by randomness (Goodliff et al., 2015). This is done over a mixture of observation periods to test different linearities, here we use (4, 8, 12, 16, 20, 24, 28), and with different window lengths (9, 13, 17, 21, 25, 29) for the machine learning skewness detection (Goodliff et al., 2020). The background error covariance matrices, **B**, is fully flow dependent.

Throughout the development of the mixed distribution approach it became ap-184 parent that there was a sensitivity to the definition of the background error covariance 185 matrix that impacted the ability for the mixed based approach to minimize. The rea-186 son for this problem is due to the property that the mode of the mixed distribution 187 is a function (sum) of the covariances. To over come this problem a flow dependent 188 approach was applied in Fletcher & Zupanski (2007) and all subsequent publications 189 associated with the mixed distribution based data assimilation schemes. This flow 190 dependency is achieved through using the averages and covariance averages from the 191 differences between the previous background trajectory and the current trajectory 192 through the time to the next cycle analysis time. This has been shown through the 193 non-Gaussian development to help stablise the mixed approach. To highlight the im-194 pact of not updating the background error covariance it can be seen in the results in 195



Figure 1. Plot comparing 3DVar-G (green), 3DVar-Mix (yellow) and 3DVar-ML (blue) with an observation period of 4 time steps, with skewness window lengths of 9 (left) and 29 (right) points. X-axis shows number of runs (each run has 5000 observations) and the y-axis is RMSE, with a rounded cumulative RMSE for each method in the legend.

Kliewer et al. (2016) that when the dynamics become more Gaussian rather lognormal the Gaussian retrieval has a smaller root mean square error than the mixed approach, but when the dynamics appear more lognormal then the mixed approach was optimal. This is an indicator that flow dependency helps improve the performance of the lognormal approach. However, because this was a retrieval system and not a model, there was no way to time evolve the solution from the previous retrieval time and hence a climatological background error covariance matrix was used.

In figure 1, we compare the three data assimilation methods with an observation 203 period of 4 time steps, and skewness window lengths of 9 points (left) and 29 points 204 (right). It can be seen that the 3DVar-ML outperforms both the 3DVar-G and 3DVar-205 Mix in both scenarios. On the left plot, we see 3DVar-G also is more accurate (in terms 206 of combined RMSE for x, y and z) than the 3DVar-Mix. On the right plot, we see the 207 opposite. In this case, 3DVar-Mix outperforms 3DVar-G. Comparing both plots, the 208 shorter skewness window length is more accurate than having a longer window length. This could be due to the skewness being accurate for the current observation, but as 210 the skewness window length increases, more information from different parts of the 211 attractor will be added to the distribution calculations. 212

By increasing the observation period to 28 time steps, it can be seen in figure 2 how the methods work in a more nonlinear setting. On the left plot, with a skewness window length of 9 points, 3DVar-ML outperforms both other methods, this result is also the case in the right plot where the skewness window length is 29 points.

By comparing figures 1 and 2, the common result is that 3DVar-ML outperforms both 3DVar-G and 3DVar-Mix. It is also seen that as we increase the observation period and skewness window length, the RMSE increases. The observation period correlation to increase RMSE values is due to the greater nonlinearity of the problem. As the data assimilation problem becomes more nonlinear, finding the minimum of the cost function becomes a more challenging problem (Goodliff et al., 2015).

In figure 3 we compare the RMSE of each method at different window lengths. As the observation period increases, the RMSE increases. This is expected in data assimilation due to nonlinear problems being harder to solve for the data assimilation methods. Again, from the results shown in figure 3, the common result for all window



Figure 2. Plot comparing 3DVar-G (green), 3DVar-Mix (yellow) and 3DVar-ML (blue) with an observation period of 28 time steps, with skewness window lengths of 9 (left) and 29 (right) points. X-axis shows number of runs (each run has 5000 observations) and the y-axis is RMSE, with a rounded cumulative RMSE for each method in the legend.



Figure 3. RMSE (y-axis) of all methods with different skewness window lengths, as a function of different observation periods (x-axis), over 50 runs.



Figure 4. Skewness Window Length by Observation Period. This graph shows the improvement in RMSE from 3DVar-G to 3DVar-ML

lengths is that 3DVar-ML outperforms 3DVar-G and 3DVar-Mix at all observation
 periods.

Figure 4 shows the percent improvement in RMSE comparing 3Dvar-G and 3DVar-ML over the all ranges set above, the improvement is higher in a more linear setting, with higher skewness window radii. It can be seen that a larger skewness window length improves the RMSE more than a shorter skewness window length. This could be due to the larger windows having more data, so that it is better able to describe the probability skewness.

²³⁵ 5 Conclusion

In this paper, we have used a machine learning technique to improve the pre-236 dictability of 3DVar when there is a change in the underlying distribution for the 237 background error distribution from Gaussian to lognormal and back to Gaussian again. 238 This improvement was achieved through using a support vector machine to detect and 239 predict non-Gaussian distributions on the z component of the Lorenz 63 model. This 240 model was used due to its simplicity, while being a chaotic system, as it is often used 241 to simulate the behaviour of the atmosphere. To determine the improvement through 242 using the support vector machine approach three data assimilation methods were com-243 pared: a Gaussian fits all 3DVar, referred to as 3DVar-G, a mixed Gaussian-lognormal 244 variant, which was referred to as 3DVar-Mix, and finally a version which used a sup-245 port vector machine to switch between Gaussian and lognormal variants for the z246 component, where this formulation was referred to as 3DVar-ML. 247

The support vector machine approach showed promising results when used in conjunction with 3DVar. It has been shown before that certain areas of the Lorenz 63 attractor do better with a lognormal variant of 3DVar, Fletcher & Zupanski (2007), due to those areas being lognormally distributed. Here, we have shown that assimilating certain areas of the attractor, depending on their probability density function (either Gaussian or lognormal distributions), can show improvements with respect to the analysis root mean square error.

For real world applications applying the support vector machine machine learn-255 ing method to choose different data assimilation types could be a way to relax the 256 Gaussian assumption for the background and observational error distributions. These 257 results could then imply that the most optimal assimilation method could be changing 258 dynamically in time, to be consistent with the more physical behavior of the errors. 259 By having this flexibility, we hypothesise that it may improve the forecast for non-260 Gaussian variables, such as used in water vapour mixing ratio retrievals Kliewer et 261 al. (2016), as well as in operational numerical weather prediction in the prediction of 262 humidity and possible certain hydrometeors. Outside of the discipline of atmosphere 263 sciences, areas that use the Gaussian assumption for data assimilation in non-Gaussian 264 systems such as space weather (example: solar winds, Lang et al. (2017)) and ocean 265 dynamics (example, ocean-biogeochemistry assimilation Goodliff et al. (2019)) could 266 also benefit through changing the underlying cost function in their data assimilation 267 systems. Implementation for this method into a high-dimensional geophysical appli-268 cation would be computationally low cost (except for training, which is usually done 269 once, and offline). The switch would act in real time, giving the predicted optimal ver-270 sion of 4DVar for each variable. In future work, we shall apply non-Gaussian detection 271 to augment data assimilation in numerical weather prediction models to determine the 272 sensitivity of this training data as well as to quantify the improvement in the forecast. 273

²⁷⁴ 6 Acknowledgements

This work is supported by the National Science Foundation grant AGS-1738206 at CIRA/CSU.

277 **References**

- Amezcua, J., & Leeuwen, P. J. V. (2014). Gaussian anamorphosis in the analysis
 step of the enkf: a joint state-variable/observation approach. *Tellus A: Dynamic Meteorology and Oceanography*, 66(1), 23493. Retrieved from https://doi.org/
 10.3402/tellusa.v66.23493 doi: 10.3402/tellusa.v66.23493
- D'agostino, R. B., Belanger, A., & D'agostino Jr., R. B. (1990). A suggestion for using powerful and informative tests of normality. *The American Statistician*, 44(4), 316-321. Retrieved from https://www.tandfonline.com/doi/abs/10.1080/ 00031305.1990.10475751 doi: 10.1080/00031305.1990.10475751
- Dueben, P. D., & Bauer, P. (2018, OCT 1). Challenges and design choices for global
 weather and climate models based on machine learning. *Geoscientific Model De-* velopment, 11(10), 3999-4009. doi: {10.5194/gmd-11-3999-2018}
- Evensen, G., & Van Leeuwen, P. (1996). Assimilation of Geosat altimeter data for
 the Agulhas current using the ensemble Kalman filter with a quasigeostrophic
 model. Mon. Wea. Rev., 124, 85-96.
- Fletcher, S. J. (2010, MAY). Mixed Gaussian-Lognormal Four-Dimensional Data Assimilation. *Tellus Series A-Dynamic Meteorology and Oceanography*, 62(3), 266-287. doi: {10.1111/j.1600-0870.2010.00439.x}
- Fletcher, S. J. (2017). Data assimilation for the geosciences: From theory to applications. Elsevier.
- Fletcher, S. J., & Jones, A. S. (2014). Multiplicative and additive incremental variational data assimilation for mixed lognormal-gaussian errors. *Monthly Weather Review*, 142(7), 2521-2544. Retrieved from https://doi.org/10.1175/MWR-D-13
 -00136.1 doi: 10.1175/MWR-D-13-00136.1
- Fletcher, S. J., & Zupanski, M. (2006a). A data assimilation method for lognormally distributed observational errors. Q. J. Roy. Meteor. Soc., 132, 2505– 2519.
- ³⁰⁴ Fletcher, S. J., & Zupanski, M. (2006b). A hybrid normal and lognormal distribu-

- tion for data assimilation. Atmos. Sci. Lett., 7, 43–46.
- Fletcher, S. J., & Zupanski, M. (2007). Implications and impacts of transforming lognormal variables into normal variables in VAR. *METEOROLOGISCHE ZEITSCHRIFT*, 16(6), 755-765. (7th International Workshop on Adjoint Applications in Dynamic Meteorology, Univ Obergurgl, Tyrol, AUSTRIA, OCT 08-13, 2006) doi: {10.1127/0941-2948/2007/0243}
- Goodliff, M. R., Amezcua, J., & Van Leeuwen, P. J. (2015). Comparing hybrid data assimilation methods on the Lorenz 1963 model with increasing non-linearity. *Tellus Series A-Dynamic Meteorology and Oceanography*, 67. doi: {10.3402/tellusa .v67.26928}
- Goodliff, M. R., Bruening, T., Schwichtenberg, F., Li, X., Lindental, A., & Nerger,
 L. (2019). Temperature assimilation into a coastal ocean-biogeochemical model:
 assessment of weakly and strongly coupled data assimilation. Ocean Dynamics,
 69, 1217–1237. Retrieved from https://doi.org/10.1007/s10236-019-01299-7
 doi: 10.1007/s10236-019-01299-7
- Goodliff, M. R., Fletcher, S. J., Kliewer, A. J., Forsythe, J. M., & Jones, A. S.
- (2020). Detection of non-gaussian behavior using machine learning techniques: A
 case study on the lorenz 63 model. Journal of Geophysical Research: Atmospheres,
 125(2), e2019JD031551. Retrieved from https://agupubs.onlinelibrary.wiley
 .com/doi/abs/10.1029/2019JD031551 (e2019JD031551 10.1029/2019JD031551)
 doi: 10.1029/2019JD031551
- Kliewer, A. J., Fletcher, S. J., Jones, A. S., & Forsythe, J. M. (2016, JAN). Comparison of Gaussian, logarithmic transform and mixed Gaussian-log-normal distribution based 1DVAR microwave temperature-water-vapour mixing ratio retrievals.
 Quarterly Journal of the Royal Meteorological Society, 142(694, A), 274-286. doi: {10.1002/qj.2651}
- Lang, M., Browne, P., van Leeuwen, P. J., & Owens, M. (2017). Data assimilation in the solar wind: Challenges and first results. Space Weather, 15(11), 1490-1510. Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/10
 .1002/2017SW001681 doi: 10.1002/2017SW001681
- Lorenz, E. N. (1963). Deterministic Nonperiodic Flow. Journal of The Atmospheric Sciences, 20(2), 130-141. doi: {10.1175/1520-0469(1963)020(0130:DNF)2.0
 .CO;2}
- Nello, C., & Shawe-Taylor, J. (2000). An introduction to support vector machines
 and other kernel-based learning method. Cambridge University Press, Cambridge,
 U.K.,.
- Ott, E., Hunt, B. R., Szunyogh, I., Zimin, A. V., Kostelich, E. J., Corazza, M., ...
- Yorke, J. A. (2004). A local ensemble transform Kalman filter for atmospheric data assimilation. *Tellus*, 56A, 415–428.
- Pasini, A. (2008, MAY-JUN). External forcings and predictability in Lorenz model:
 An analysis via neural network modelling. Nuovo Cimento Della Societa Italiana
 di Fisica C-Colloquia on Physics, 31(3), 357-370. doi: {10.1393/ncc/i2009-10312
 -1}
- Pasini, A., & Pelino, V. (2005). Can we estimate atmospheric predictability by
 performance of neural network forecasting? The toy case studies of unforced and
 forced Lorenz models., 69-74. (IEEE International Conference on Computational
 Intelligence for Measurement Systems and Applications, Messina, ITALY, JUL
 20-22, 2005)
- Perron, M., & Sura, P. (2013). Climatology of non-gaussian atmospheric statis tics. Journal of Climate, 26(3), 1063-1083. Retrieved from https://doi.org/10
 .1175/JCLI-D-11-00504.1 doi: 10.1175/JCLI-D-11-00504.1
- Rasp, S., & Lerch, S. (2018, NOV). Neural Networks for Postprocessing Ensemble
 Weather Forecasts. Monthly Weather Review, 146(11), 3885-3900. doi: {10.1175/
- 358 MWR-D-18-0187.1}

- Scher, S. (2018, NOV 28). Toward Data-Driven Weather and Climate Forecasting:
 Approximating a Simple General Circulation Model With Deep Learning. Geophysical Research Letters, 45(22), 12616-12622. doi: {10.1029/2018GL080704}
- Scher, S., & Messori, G. (2018). Predicting weather forecast uncertainty with machine learning. Quarterly Journal of the Royal Meteorological Society, 144(717), 2830-2841. Retrieved from https://rmets.onlinelibrary.wiley.com/doi/abs/ 10.1002/qj.3410 doi: 10.1002/qj.3410
- Scher, S., & Messori, G. (2019, JUL 10). Weather and climate forecasting with
 neural networks: using general circulation models (GCMs) with different complex ity as a study ground. *Geoscientific Model Development*, 12(7), 2797-2809. doi:
 {10.5194/gmd-12-2797-2019}
- van Leeuwen, P. J., Künsch, H. R., Nerger, L., Potthast, R., & Reich, S. (2019).
 Particle filters for high-dimensional geoscience applications: A review. *Quarterly Journal of the Royal Meteorological Society*, 145(723), 2335-2365. Retrieved from https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.3551 doi: https://doi.org/10.1002/qj.3551
- Wang, X., & Bishop, C. H. (2003). A comparison of breeding and ensemble transform Kalman filter ensemble forecast schemes. J. Atmos. Sci., 60, 1140–1158.
- Weyn, J. A., Durran, D. R., & Caruana, R. (2019). Can machines learn to predict
- weather? using deep learning to predict gridded 500-hpa geopotential height from historical weather data. Journal of Advances in Modeling Earth Systems, 11(8),
- 2680-2693. Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/ abs/10.1029/2019MS001705 doi: 10.1029/2019MS001705
- Zupanski, M. (2005). Maximum Likelihood Ensemble Filter. Part I: Theoretical Aspects. Mon. Wea. Rev, 133, 1710-1726.