



1 **An Estimate of Inflation Factor and Analysis Sensitivity in**
2 **Ensemble Kalman Filter**

3

4 Guocan Wu^{1,2}

5

6

7 1 College of Global Change and Earth System Science, Beijing Normal University,

8 Beijing, China

9 2 Joint Center for Global Change Studies, Beijing, China

10



1 **Abstract**

2

3 The estimation accuracy of forecast error matrix is crucial to the assimilation
4 result. Ensemble Kalman filter (EnKF) is a widely used ensemble based assimilation
5 method, which initially estimate the forecast error matrix using a Monte Carlo
6 method with the short-term ensemble forecast states. However, this estimate needs to
7 be further improved using inflation technique.

8 In this study, the forecast error inflation factor is estimated based on cross
9 validation and the analysis sensitivity is also investigated. The improved EnKF
10 assimilation scheme is validated by assimilating spatially correlated observations to
11 the atmosphere-like Lorenz-96 model. The experiment results show that, the analysis
12 error is reduced and the analysis sensitivity to observations is improved.

13

14 **Key words:** data assimilation; ensemble Kalman filter; forecast error inflation;
15 analysis sensitivity; cross validation

16



1 1. Introduction

2

3 In the mathematical and physical research fields, data assimilation is a powerful
4 mechanism to estimate the true trajectory of a state variable, based on the effective
5 combination of the dynamic forecast system (numerical model) and the observations
6 (Miller et al. 1994). It can provide an analysis state, which is generally treated as the
7 weighted average of the model forecasts and observations, and is more close to the
8 true state than either of them. The weights are approximately proportional to the
9 inverse of the corresponding covariance matrices (Talagrand 1997). Therefore, the
10 performance of a data assimilation method significantly relies on whether the error
11 covariance matrices are estimated accurately. If this is the case, the analysis state can
12 be technically easily obtained by minimizing a cost function with many existing
13 optimization methods (Reichle 2008).

14 Ensemble Kalman filter (EnKF) is a very practical ensemble based assimilation
15 scheme, which estimates the forecast error covariance matrix using a Monte Carlo
16 method with the short-term ensemble forecast states (Burgers et al. 1998; Evensen
17 1994). Because of the limited ensemble size and large model error, the sampling
18 covariance matrix of the ensemble forecast states is usually an underestimate of the
19 true forecast error covariance matrix. It can lead that the filter over trusts the model
20 forecasts and excludes the observations, and can eventually result in the divergence of
21 the filter (Anderson; Anderson 1999; Constantinescu et al. 2007; Wu et al. 2014).
22 Therefore, using the inflation technique to enhance the estimate accuracy of the



1 forecast error covariance matrix becomes gradually important.

2 In early studies in forecast error inflation, researchers usually tune the inflation
3 factor by repeated assimilation experiments and select the estimated inflation factor
4 according to their experiences and prior knowledge (Anderson; Anderson 1999).
5 Hence such experimental determining is very empirical and subjective. In later
6 studies, the inflation factor can be estimated on-line based on the innovation statistic
7 (observation-minus-forecast (Dee 1995; Dee; Silva 1999)) with some different
8 conditions. The moment estimation can facilitate the calculation by solving an
9 equation of the innovation statistic and its realization (Li et al. 2009; Miyoshi 2011;
10 Wang; Bishop 2003). The maximum likelihood approach can obtain a better inflation
11 but has to calculate a high dimensional matrix determinant (Liang et al. 2012; Zheng
12 2009). The Bayesian approach assumes a prior distribution for the inflation factor but
13 limited to the spatially independent observational errors (Anderson 2007, 2009). This
14 study seeks to address the estimation of the inflation factor from the point view of
15 Cross Validation (CV).

16 In fact, the idea of Cross Validation (CV) is firstly involved in linear regression
17 (Allen 1974) and smoothing spline (Wahba; Wold. 1975). It is a common approach
18 that can be applied to estimate tuning parameters in generalized additive models,
19 nonparametric regression and kernel smoothing (Eubank 1999; Gentle et al. 2004;
20 Green; Silverman. 1994; Wand; Jones 1995). In cross validation, sample is cut into
21 several smaller data subsets, and some of them are used for modeling and analysis
22 while others are used for verification and validation. The widely used technique is to



1 remove only one data point and use the rest to estimate the value at this point so as to
2 test the estimation accuracy, which is also called Leave-Out-One Cross Validation
3 (Gu; Wahba 1991).

4 The basic motivation behind the Cross Validation is minimizing the prediction
5 error at the sampling points. The Generalized Cross Validation (GCV) is the modified
6 form of Cross Validation, which is more popular for choosing these turning
7 parameters (Craven; Wahba 1979). For instance, Gu and Wahba applied the Newton
8 method to optimize the Generalized Cross Validation score with multiple smoothing
9 parameters in a smoothing spline model (Gu; Wahba 1991). Wahba briefly reviewed
10 the properties of Generalized Cross Validation and carried out an experiment to
11 choose smoothing parameters in the context of variational data assimilation schemes
12 with Numerical Weather Prediction models (Wahba et al. 1995). Zheng and Basher
13 also used Generalized Cross Validation in thin-plate smoothing spline modeling of
14 spatial climate data and applied to south Pacific rainfalls (Zheng; Basher 1995). The
15 Generalized Cross Validation criterion also has been found to possess several
16 favorable properties, such as consistency of the relative loss (Gu 2002).

17 Intuitively, if the forecast error matrix is inflated properly, the assimilation
18 procedure can reassign the weights of the model forecasts and observations.
19 Therefore the analysis sensitivity is also investigated in this study. Generally speaking,
20 analysis sensitivity is how uncertainty in the output can be apportioned to different
21 source of uncertainty in the input (Saltelli et al. 2004; Saltelli et al. 2008). The
22 quantity can be introduced to the context of statistical data assimilation framework. It



1 indicates that the sensitivity of analysis to observations, which is complementary to
2 the sensitivity of analysis to model forecasts (Cardinali et al. 2004; Liu et al. 2009)..

3 This study focuses on the methodology part that can be potentially applied in
4 the geophysical research fields in the near future. This paper consists of 4 sections.
5 The conventional EnKF scheme is summarized and the improved EnKF with forecast
6 error inflation scheme is proposed in Section 2. The verification and validation are
7 conducted on an idealized model in Section 3. Discussion and conclusions are given
8 in Section 4.

9
10

11 **2. Methodology**

12

13 **2.1. EnKF Algorithm**

14 For the sake of consistency, a nonlinear discrete-time dynamic forecast and linear
15 observation system can be expressed as (Ide et al. 1997),

$$16 \quad \mathbf{x}_i^t = M_{i-1}(\mathbf{x}_{i-1}^a) + \boldsymbol{\eta}_i, \quad (1)$$

$$17 \quad \mathbf{y}_i^o = \mathbf{H}_i \mathbf{x}_i^t + \boldsymbol{\epsilon}_i, \quad (2)$$

18 where i stands for the time index; $\mathbf{x}_i^t = \{x_{i,1}^t, x_{i,2}^t, \dots, x_{i,n}^t\}^T$ is the n -dimensional true
19 state vector at i -th time step; $\mathbf{x}_{i-1}^a = \{x_{i-1,1}^a, x_{i-1,2}^a, \dots, x_{i-1,n}^a\}^T$ is the n -dimensional
20 analysis state vector which is an estimate of \mathbf{x}_{i-1}^t , M_{i-1} is a nonlinear dynamic
21 forecast operator such as a numeric weather prediction model;
22 $\mathbf{y}_i^o = \{y_{i,1}^o, y_{i,2}^o, \dots, y_{i,p_i}^o\}^T$ is a p_i -dimensional observation vector; \mathbf{H}_i is the



1 observation operator matrix, $\boldsymbol{\eta}_i$ and $\boldsymbol{\varepsilon}_i$ are the forecast and observation error
 2 vectors, which are assumed to be time-uncorrelated, statistically independent of each
 3 other and have mean zero and covariance matrices \mathbf{P}_i and \mathbf{R}_i , respectively. The
 4 EnKF assimilation result is a series of analysis state \mathbf{x}_i^a that are sufficiently close to
 5 the corresponding true states \mathbf{x}_i^t , based on the information provided by M_i and
 6 \mathbf{y}_i^o .

7 Suppose the perturbed analysis state at previous time step $\mathbf{x}_{i-1}^{a(j)}$ has been
 8 estimated ($1 \leq j \leq m$ and m is the ensemble size), the detailed EnKF assimilation
 9 procedure is summarized as the following forecast step and analysis step (Burgers et
 10 al. 1998; Evensen 1994).

11 Step 1. Forecast Step.

12 The perturbed forecast states are generated by dynamic model forecast forward

$$13 \quad \mathbf{x}_i^{f(j)} = M_{i-1}(\mathbf{x}_{i-1}^{a(j)}). \quad (3)$$

14 The forecast state \mathbf{x}_i^f is defined to be the ensemble mean of $\mathbf{x}_i^{f(j)}$ and the forecast
 15 error covariance matrix is initially estimated as the sampling covariance matrix of
 16 perturbed forecast states

$$17 \quad \mathbf{P}_i = \frac{1}{m-1} \sum_{j=1}^m (\mathbf{x}_i^{f(j)} - \mathbf{x}_i^f)(\mathbf{x}_i^{f(j)} - \mathbf{x}_i^f)^T. \quad (4)$$

18 Step 2. Analysis Step.

19 The analysis state is estimated by minimizing the following cost function

$$20 \quad J(\mathbf{x}) = (\mathbf{x} - \mathbf{x}_i^f)^T \mathbf{P}_i^{-1} (\mathbf{x} - \mathbf{x}_i^f) + (\mathbf{y}_i^o - \mathbf{H}_i \mathbf{x})^T \mathbf{R}_i^{-1} (\mathbf{y}_i^o - \mathbf{H}_i \mathbf{x}), \quad (5)$$

21 which has the analytic form



1
$$\mathbf{x}_i^a = \mathbf{x}_i^f + \mathbf{P}_i \mathbf{H}_i^T (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} \mathbf{d}_i, \quad (6)$$

2 where

3
$$\mathbf{d}_i = \mathbf{y}_i^o - \mathbf{H}_i \mathbf{x}_i^f, \quad (7)$$

4 is the innovation statistic (observation-minus-forecast residual). In order to complete
 5 the ensemble forecast, the perturbed analysis state are calculated using perturbed
 6 observations (Burgers et al. 1998), that is

7
$$\mathbf{x}_i^{a(j)} = \mathbf{x}_i^{f(j)} + \mathbf{P}_i \mathbf{H}_i^T (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} (\mathbf{d}_i + \boldsymbol{\varepsilon}_i^{(j)}), \quad (8)$$

8 where $\boldsymbol{\varepsilon}_i^{(j)}$ is a normally distributed random variable with mean zero and covariance
 9 matrix \mathbf{R}_i . Here $(\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1}$ can be easily calculated using the
 10 Sherman-Morrison-Woodbury formula (Liang et al. 2012; Tippett et al. 2003).
 11 Finally, set $i = i + 1$ and return to Step 1 for the model forecast at next time step.

12

13 **2.2. Influence matrix and forecast error inflation**

14 It is recognized that the forecast error inflation scheme should be included in any
 15 ensemble based assimilation scheme, otherwise, the filter could diverge (Anderson;
 16 Anderson 1999; Constantinescu et al. 2007). The multiplicative inflation is one of the
 17 commonly used inflation techniques, which adjusts the initially estimated forecast
 18 error covariance matrix \mathbf{P}_i to $\lambda_i \mathbf{P}_i$ and then estimates the inflation factors λ_i
 19 properly.

20 In previous studies, there are many methods for estimating the inflation factor,
 21 such as the maximum likelihood approach (Liang et al. 2012; Zheng 2009), moment
 22 approach (Li et al. 2009; Miyoshi 2011; Wang; Bishop 2003) and Bayesian approach



1 (Anderson 2007, 2009). In this study, a new procedure for estimating the
 2 multiplicative inflation factors λ_i is proposed based on the following Generalized
 3 Cross Validation (GCV) function (Craven; Wahba 1979)

$$4 \quad GCV_i(\lambda) = \frac{\frac{1}{p_i} \mathbf{d}_i^T \mathbf{R}_i^{-1/2} (\mathbf{I}_{p_i} - \mathbf{A}_i(\lambda))^2 \mathbf{R}_i^{-1/2} \mathbf{d}_i}{\left[\frac{1}{p_i} \text{Tr}(\mathbf{I}_{p_i} - \mathbf{A}_i(\lambda)) \right]^2}, \quad (9)$$

5 where \mathbf{I}_{p_i} is the identity matrix with dimension $p_i \times p_i$, $\mathbf{R}_i^{-1/2}$ is the square root
 6 matrix of \mathbf{R}_i and

$$7 \quad \mathbf{A}_i(\lambda) = \mathbf{I}_{p_i} - \mathbf{R}_i^{1/2} (\mathbf{H}_i \lambda \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} \mathbf{R}_i^{1/2} \quad (10)$$

8 is the influence matrix (see Appendix for details).

9 The estimation procedure of inflation factors λ_i is implemented between the
 10 Step 1 and 2 in Section 2.1. Then the perturbed analysis states are modified to

$$11 \quad \mathbf{x}_i^{a(j)} = \mathbf{x}_i^{f(j)} + \lambda_i \mathbf{P}_i \mathbf{H}_i^T (\mathbf{H}_i \lambda_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} (\mathbf{d}_i + \boldsymbol{\varepsilon}_i^{(j)}). \quad (11)$$

12 The flowchart of the EnKF equipped with forecast error inflation based on GCV
 13 method is shown in Figure 1.

14

15 **2.3. Analysis sensitivity**

16 In EnKF, the analysis state (Eq. (6)) can be treated as a weighted average of the
 17 observation and forecast, that is,

$$18 \quad \mathbf{x}_i^a = \mathbf{K}_i \mathbf{y}_i^o + (\mathbf{I}_n - \mathbf{K}_i \mathbf{H}_i) \mathbf{x}_i^f \quad (12)$$

19 where $\mathbf{K}_i = \mathbf{P}_i \mathbf{H}_i^T (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1}$ is the Kalman gain matrix, and \mathbf{I}_n is the identity
 20 matrix with dimension $n \times n$. Then the normalized analysis vector can be expressed

21 as



1
$$\tilde{\mathbf{y}}_i^a = \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{K}_i \mathbf{R}_i^{1/2} \tilde{\mathbf{y}}_i^o + \mathbf{R}_i^{-1/2} (\mathbf{I}_{p_i} - \mathbf{H}_i \mathbf{K}_i) \mathbf{R}_i^{1/2} \tilde{\mathbf{y}}_i^f \quad (13)$$

2 where $\tilde{\mathbf{y}}_i^f = \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^f$ is the normalized projection of the forecast on the
 3 observation space. The sensitivities of the analysis to the observation and forecast are
 4 defined as

5
$$\mathbf{S}_i^o = \frac{\partial \tilde{\mathbf{y}}_i^a}{\partial \tilde{\mathbf{y}}_i^o} = \mathbf{R}_i^{1/2} \mathbf{K}_i^T \mathbf{H}_i^T \mathbf{R}_i^{-1/2}, \quad (14)$$

6
$$\mathbf{S}_i^f = \frac{\partial \tilde{\mathbf{y}}_i^a}{\partial \tilde{\mathbf{y}}_i^f} = \mathbf{R}_i^{1/2} (\mathbf{I}_{p_i} - \mathbf{K}_i^T \mathbf{H}_i^T) \mathbf{R}_i^{-1/2}, \quad (15)$$

7 respectively, which satisfy $\mathbf{S}_i^o + \mathbf{S}_i^f = \mathbf{I}_{p_i}$.

8 The elements of the matrix \mathbf{S}_i^o reflect the sensitivity of the normalized analysis
 9 state to the normalized observations. Its diagonal elements are the analysis
 10 self-sensitivities, and off-diagonal elements are cross-sensitivities. On the other hand,
 11 the elements of the matrix \mathbf{S}_i^f reflect the sensitivity of the normalized analysis state
 12 to the normalized forecast vector. The two quantities are complementary.

13 In fact, the sensitivity matrix \mathbf{S}_i^o is equal to the influence matrix \mathbf{A}_i (see
 14 Appendix B for detail proof), whose trace can be used to measure the “equivalent
 15 number of parameters” or “degrees of freedom for signal”. Similarly, it can be
 16 interpreted as a measurement of the amount of information extracted from the
 17 observations. The trace diagnostic can be used to analyze the sensitivities to
 18 observation or forecast vector (Cardinali et al. 2004). The Global Average Influence
 19 (GAI) at i -th time step is defined as the globally averaged observation influence, that
 20 is

21
$$GAI = \frac{\text{Tr}(\mathbf{S}_i^o)}{p_i} \quad (16)$$



1 where p_i is the total number of observations at the i -th time step.

2 In the conventional EnKF, the forecast error covariance matrix \mathbf{P}_i is initially
3 estimated using a Monte Carlo method with the short-term ensemble forecast states.
4 However, due to the limited ensemble size and large model error, the sampling
5 covariance matrix of perturbed forecast states is usually an underestimation of the
6 true forecast error covariance matrix. This will cause the assimilation systems over
7 trust the forecast state, and then exclude the observational information. That is why
8 the values of GAI are too small in conventional EnKF scheme. It will show that in
9 simulations, this problem will be alleviated to some extent through the inflation
10 adjustment on forecast error covariance matrix.

11

12 **2.4 Analysis RMSE**

13 In the following experiments, the “true” state \mathbf{x}_i^t is non-dimensional and can
14 be obtained by numerical solution of partial differential equations. In this case, the
15 distance of the analysis state to the true state can be defined as the analysis
16 root-mean-square error (RMSE), which is used to evaluate the accuracy of the
17 assimilation results. The RMSE at i -th time step is defined as

$$18 \quad \text{RMSE} = \sqrt{\frac{1}{n} \sum_{k=1}^n \left(x_{i,k}^a - x_{i,k}^t \right)^2}. \quad (17)$$

19 where $x_{i,k}^a$ and $x_{i,k}^t$ are the k -th component of the analysis state and true state at
20 i -th time step.

21

22 **3. Experimentations**



1

2 The proposed data assimilation scheme is validated using the Lorenz-96 model
3 (Lorenz 1996) with model error and a linear observation system as a test bed in this
4 section. The performances of the assimilation schemes described in Section 2 are
5 evaluated through the following experiments.

6

7 ***3.1. The dynamic forecast model and observation systems***

8 The Lorenz-96 model (Lorenz 1996) is a quadratic nonlinear dynamical system,
9 which has the properties relevant to realistic forecast problems and is governed by the
10 equation

$$11 \quad \frac{d\mathbf{X}_k}{dt} = (\mathbf{X}_{k+1} - \mathbf{X}_{k-2})\mathbf{X}_{k-1} - \mathbf{X}_k + F, \quad (18)$$

12 where $k=1,2,\dots,40$. The cyclic boundary conditions $\mathbf{X}_{-1} = \mathbf{X}_{K-1}$, $\mathbf{X}_0 = \mathbf{X}_K$,
13 $\mathbf{X}_{K+1} = \mathbf{X}_1$ is applied to make Eq. (18) to be well-defined for all values of k . The
14 Lorenz-96 model is “atmosphere-like”, since the three terms on the right-hand side of
15 Eq. (18) can be analogized to a nonlinear advection-like term, a damping term, and an
16 external forcing term respectively. It can be thought of as some atmospheric quantity
17 (e.g. zonal wind speed) distributed on a latitude circle. Therefore the Lorenz-96
18 model is widely used as a test bed to evaluate the performances of assimilation
19 schemes in many studies (Wu et al. 2013).

20 The time step is set as 0.05 non-dimensional unit when generate the numeric
21 solution, which is roughly equivalent to 6 hours in real time, assuming that the
22 characteristic time-scale of the dissipation in the atmosphere is 5 days (Lorenz 1996).



1 The true state is derived by a fourth-order Runge-Kutta time integration scheme
2 (Butcher 2003). The forcing term is set as $F = 8$, so that the leading Lyapunov
3 exponent implies an error-doubling time of about 8 time steps, and the fractal
4 dimension of the attractor is 27.1 (Lorenz; Emanuel 1998). The initial value is chosen
5 to be $\mathbf{X}_k = F$ when $k \neq 20$ and $\mathbf{X}_{20} = 1.001F$.

6 In this study, the synthetic observations are assumed to be generated at all of the
7 40 model grids but every 4 time steps by adding random noises which are
8 multivariate-normally distributed with mean zero and covariance matrix \mathbf{R}_i to the
9 true states. The observation errors are assumed to be spatially correlated, which is the
10 most cases in remote sensing observations and radiances data. The variance of the
11 observation on each grid point is set $\sigma_o^2 = 1$ and the covariance of the observations
12 between the j -th and k -th grid points is

$$13 \quad \mathbf{R}_i(j, k) = \sigma_o^2 \times 0.5^{\min\{|j-k|, 40-|j-k|\}}. \quad (19)$$

14 The heat map of the observation error covariance matrix is shown in Figure 2.

15

16 **3.2. Assimilation schemes comparison**

17 Since model error is inevitable in practical dynamic forecast models, it is
18 reasonable for us to add model error to the Lorenz-96 model in the assimilation
19 process. Lorenz-96 model is a forced dissipative model with a parameter F that
20 controls the strength of the forcing (Eq. (18)). The model forecast changes very much
21 along with the change of F and is chaos with integer values of F larger than 3.
22 Therefore the forcing term is set as 7 to simulate the range of model error in the



1 assimilation schemes while retaining $F = 8$ when generating the “true” state. The
2 ensemble size is selected as 30.

3 To evaluate the analysis sensitivity, the GAI statistics (Eq. (16)) are calculated
4 and the results are plotted in Figure 3 over 2000 time steps, which is about equivalent
5 to 500 days in realistic problems. It clearly shows that, the percentage of the analysis
6 result relied on the observation is about 10% for the conventional EnKF, which is
7 increased to about 30% for the EnKF with forecast error inflation.

8 To evaluate the assimilation result, the analysis RMSE (Eq. (17)) and the
9 corresponding values of the GCV functions (Eq. (9)) are calculated and plotted in
10 Figures 4 and 5, respectively. It illustrates that, the analysis RMSE, as well as the
11 values of the GCV functions increase sharply if the forecast error inflation is adopted.
12 The variety of the analysis RMSE is very consistent with that of the value of the GCV
13 function for the EnKF with forecast error inflation scheme. The correlation
14 coefficient of the analysis RMSE and the value of the GCV function at the
15 assimilation time step is about 0.76, which indicating that, the GCV function seems to
16 be a good criterion to estimate the inflation factor.

17 The time-mean values of the GAI statistics, the GCV functions and the analysis
18 RMSE over 2000 time steps are listed in Table 1. These results illustrate that, the
19 forecast error inflation technique using the GCV function can indeed increase the
20 analysis sensitivity to the observations and reduce the analysis RMSE.

21

22



1 **4. Discussion and Conclusions**

2

3 As we all know that accurately estimating the error covariance matrix is one of
4 the most important steps in data assimilation, which has curial influence to the
5 assimilation results. In conventional EnKF assimilation scheme, the forecast error
6 covariance matrix is estimated as the sampling covariance matrix of the ensemble
7 forecast states. But due to the limited ensemble size and large model error, this initial
8 estimate is usually an underestimation, which can lead to that the filter over trusts the
9 forecasts and excludes the observations, and eventually the divergence of the filter.
10 Therefore the forecast error inflation with proper inflation factor is increasingly
11 important.

12 The multiplicative inflation is an effective inflation technique and the inflation
13 factor can be estimated under different assumptions. The moment approach can be
14 easily conducted based on the moment estimation of the innovation statistic. The
15 maximum likelihood approach can obtain a more accurate inflation factor than the
16 moment approach but with complicated calculations of high dimensional matrix
17 determinant. The Bayesian approach assumes a prior distribution for the inflation
18 factor but limited to the spatially independent observational errors. In this study, the
19 inflation factor is estimated from the point of view of cross validation and the
20 analysis sensitivity is detected.

21 The GCV function seems to be a good objective function that can well quantify
22 the goodness of fit of the error covariance matrix. In fact, cross validation, which can



1 evaluate and compare learning algorithms, is a widely used statistical method. The
2 most common form of cross validation is leave-out-one cross validation. For this
3 algorithm, all the data except for a single observation are used for training and the
4 comparison is made on that single observation. Generalized Cross Validation
5 estimate is a modified form of ordinary Cross Validation, which has the
6 rotation-invariant property relative to an orthogonal transform of the observations and
7 is a consistent estimate of the relative loss.

8 In this study, the idea of Cross Validation is adopted to the inflation factor
9 estimation in the improved EnKF assimilation scheme and validated on the
10 Lorenz-96 model. The values of the GCV function obviously decrease in the
11 proposed approach comparing with that in the conventional EnKF scheme. The
12 analysis RMSE in the proposed approach also is much smaller than that in the
13 conventional EnKF scheme. This suggested that the estimate inflation factor method
14 through minimizing the GCV function works very well.

15 The varieties of analysis sensitivity in the proposed approach and in the
16 conventional EnKF scheme are also investigated in this study. The influence matrix is
17 treated as the partial derivative of the normalized analysis state vector to the
18 normalized observational vector, which is also used in the GCV function.
19 Additionally, the time-mean Global Average Influence statistic is increased from
20 about 10% in the conventional EnKF scheme to about 30% in the proposed improved
21 EnKF assimilation scheme. This illustrated that the shortcoming of the assimilation
22 result excessively depending on the forecast and excluding the observation can be



1 mitigated in some extent. The relations of analysis state to forecast state and
2 observation are more reasonable.

3 It is notable that, the inflation factor is assumed to be constant in space in this
4 study, which may be not the case in the realistic assimilation problems. Therefore
5 persistently adjust all the state vectors using the same inflation factor could
6 systematically overinflate the ensemble variances in sparsely observed areas,
7 especially when the observations are unevenly distributed. In the case studies
8 conducted in Section 3, the observations are relatively evenly distributed and the
9 assimilation accuracy can be indeed improved by the forecast error inflation
10 technique. It mainly sheds light on the methodology and validate on Lorenz-96 model
11 to illustrate the feasibility in this study. In the near future, it will investigated that how
12 to modify the adaptive procedure to suit the system with unevenly distributed
13 observations and apply the proposed methodologies using more sophisticated
14 dynamic and observation systems.

15

16 **Appendix A**

17 From Eq. (2), the normalized observation equation can be defined as

$$18 \quad \tilde{\mathbf{y}}_i^o = \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^t + \tilde{\boldsymbol{\varepsilon}}_i, \quad (\text{A1})$$

19 where $\tilde{\mathbf{y}}_i^o = \mathbf{R}_i^{-1/2} \mathbf{y}_i^o$ is the normalized observation vector and $\tilde{\boldsymbol{\varepsilon}}_i \sim N(\mathbf{0}, \mathbf{I})$, \mathbf{I}_{p_i} is
20 the identity matrix with dimension $p_i \times p_i$. Similarly, the normalized analysis vector
21 is $\tilde{\mathbf{y}}_i^a = \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^a$ and the influence matrix \mathbf{A}_i relates the normalized observation
22 vector to the normalized analysis vector, ignoring the normalized forecast state in the



1 observation space (Gu 2002), i.e.,

$$2 \quad \tilde{\mathbf{y}}_i^a - \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^f = \mathbf{A}_i \left(\tilde{\mathbf{y}}_i^o - \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^f \right). \quad (\text{A2})$$

3 Since the analysis state \mathbf{x}_i^a is given by Eq. (5), it can be easily checked that the
 4 influence matrix \mathbf{A}_i is given by

$$5 \quad \mathbf{A}_i = \mathbf{I}_{p_i} - \mathbf{R}_i^{1/2} \left(\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i \right)^{-1} \mathbf{R}_i^{1/2}. \quad (\text{A3})$$

6 If the initial forecast error covariance matrix is inflated as described in Section 2.2,
 7 the influence matrix is treated as the following function of λ

$$8 \quad \mathbf{A}_i(\lambda) = \mathbf{I}_{p_i} - \mathbf{R}_i^{1/2} \left(\mathbf{H}_i \lambda \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i \right)^{-1} \mathbf{R}_i^{1/2}, \quad (\text{A4})$$

9 The principle of cross validation aims at minimizing the estimated error at the
 10 observation grid point. Lacking an independent validation data set, the alternative
 11 strategy commonly used is to minimize the squared distance between the normalized
 12 observation value and the analysis value while not using the observation on the same
 13 grid point, that is the following objective function

$$14 \quad V_i(\lambda) = \frac{1}{p_i} \sum_{k=1}^{p_i} \left(\tilde{\mathbf{y}}_{i,k}^o - \left(\mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^{a[k]} \right)_k \right)^2, \quad (\text{A5})$$

15 where $\mathbf{x}_i^{a[k]}$ is the minima of the following “delete-one” objective function

$$16 \quad \left(\mathbf{x} - \mathbf{x}_i^f \right)^T \left(\lambda \mathbf{P}_i \right)^{-1} \left(\mathbf{x} - \mathbf{x}_i^f \right) + \left(\mathbf{y}_i^o - \mathbf{H}_i \mathbf{x} \right)_{-k}^T \mathbf{R}_{i,-k}^{-1/2} \left(\mathbf{y}_i^o - \mathbf{H}_i \mathbf{x} \right)_{-k}. \quad (\text{A6})$$

17 The subscript $-k$ means a vector (matrix) with its k -th element (k -th row and column)
 18 deleted. Instead of minimizing Eq. (A6) p_i times, the objective function (Eq. (A5))
 19 has another more simple expression (Gu 2002)

$$20 \quad V_i(\lambda) = \frac{1}{p_i} \sum_{k=1}^{p_i} \frac{\left(\tilde{\mathbf{y}}_{i,k}^o - \left(\mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^a \right)_k \right)^2}{\left(1 - a_{k,k} \right)^2}, \quad (\text{A7})$$

21 where $a_{k,k}$ is the element at the site pair (k, k) of the influence matrix $\mathbf{A}_i(\lambda)$. Then,



1 substituting $a_{k,k}$ by the average $\frac{1}{p_i} \sum_{k=1}^{p_i} a_{k,k} = \frac{1}{p_i} \text{Tr}(\mathbf{A}_i(\lambda))$ and ignoring the
 2 constant to get the following generalized cross validation (GCV) statistic (Gu 2002)

$$3 \quad GCV_i(\lambda) = \frac{\frac{1}{p_i} \mathbf{d}_i^T \mathbf{R}_i^{-1/2} (\mathbf{I}_{p_i} - \mathbf{A}_i(\lambda))^2 \mathbf{R}_i^{-1/2} \mathbf{d}_i}{\left[\frac{1}{p_i} \text{Tr}(\mathbf{I}_{p_i} - \mathbf{A}_i(\lambda)) \right]^2}. \quad (\text{A8})$$

4

5 **Appendix B**

6 The sensitivities of the analysis to the observation is defined as

$$7 \quad \mathbf{S}_i^o = \frac{\partial \tilde{\mathbf{y}}_i^a}{\partial \tilde{\mathbf{y}}_i^o} = \mathbf{R}_i^{1/2} \mathbf{K}_i^T \mathbf{H}_i^T \mathbf{R}_i^{-1/2}, \quad (\text{B1})$$

8 Substitute the Kalman gain matrix $\mathbf{K}_i = \mathbf{P}_i \mathbf{H}_i^T (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1}$ into \mathbf{S}_i^o , then

$$\begin{aligned} \mathbf{S}_i^o &= \mathbf{R}_i^{1/2} \mathbf{K}_i^T \mathbf{H}_i^T \mathbf{R}_i^{-1/2} \\ &= \mathbf{R}_i^{1/2} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} \mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T \mathbf{R}_i^{-1/2} \\ &= \mathbf{R}_i^{1/2} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i - \mathbf{R}_i) \mathbf{R}_i^{-1/2} \\ 9 \quad &= \mathbf{R}_i^{1/2} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i) \mathbf{R}_i^{-1/2} - \mathbf{R}_i^{1/2} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} \mathbf{R}_i \mathbf{R}_i^{-1/2} \quad (\text{B2}) \\ &= \mathbf{I}_{p_i} - \mathbf{R}_i^{1/2} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} \mathbf{R}_i^{1/2} \\ &= \mathbf{A}_i \end{aligned}$$

10 Therefore, the sensitivity matrix \mathbf{S}_i^o is equal to the influence matrix \mathbf{A}_i .

11

12 **Acknowledgements.** This work is supported by the National Basic Research
 13 Program of China (Grant Nos. 2015CB953703 and 2010CB950703) and the National
 14 Natural Science Foundation of China (Grant No. 41405098).

15

16



1 References

2

- 3 Allen, D. M., 1974: The relationship between variable selection and data augmentation and a method
4 for prediction. *Technometrics*, **16**, 125-127.
- 5 Anderson, J. L., 2007: An adaptive covariance inflation error correction algorithm for ensemble filters.
6 *Tellus*, **59A**, 210-224.
- 7 Anderson, J. L., 2009: Spatially and temporally varying adaptive covariance inflation for ensemble
8 filters. *Tellus*, **61A**, 72-83.
- 9 Anderson, J. L., and S. L. Anderson, 1999: A Monte Carlo implementation of the nonlinear filtering
10 problem to produce ensemble assimilations and forecasts. *Monthly Weather Review*, **127**, 2741-2758.
- 11 Burgers, G., P. J. Leeuwen, and G. Evensen, 1998: Analysis scheme in the ensemble kalman filter.
12 *Monthly Weather Review*, **126**, 1719-1724.
- 13 Butcher, J. C., 2003: *Numerical methods for ordinary differential equations*. JohnWiley & Sons, 425
14 pp.
- 15 Cardinali, C., S. Pezzulli, and E. Andersson, 2004: Influence - matrix diagnostic of a data assimilation
16 system. *Quarterly Journal of the Royal Meteorological Society*, **130**, 2767-2786.
- 17 Constantinescu, E. M., A. Sandu, T. Chai, and G. R. Carmichael, 2007: Ensemble-based chemical data
18 assimilation I: general approach. *Quarterly Journal of the Royal Meteorological Society*, **133**,
19 1229-1243.
- 20 Craven, P., and G. Wahba, 1979: Smoothing noisy data with spline functions. *Numerische Mathematik*,
21 **31**, 377-403.
- 22 Dee, D. P., 1995: On-line estimation of error covariance parameters for atmospheric data assimilation.
23 *Monthly Weather Review*, **123**, 1128-1145.
- 24 Dee, D. P., and A. M. Silva, 1999: Maximum-likelihood estimation of forecast and observation error
25 covariance parameters part I: methodology. *Monthly Weather Review*, **127**, 1822-1834.
- 26 Eubank, R. L., 1999: *Nonparametric regression and spline smoothing*. Marcel Dekker, Inc., 338 pp.
- 27 Evensen, G., 1994: Sequential data assimilation with a nonlinear quasi-geostrophic model using
28 Monte Carlo methods to forecast error statistics. *Journal of Geophysical Research*, **99**, 10143-10162.
- 29 Gentle, J. E., W. Hardle, and Y. Mori, 2004: *Handbook of computational statistics: concepts and
30 methods*. Springer, 1070 pp.
- 31 Green, P. J., and B. W. Silverman., 1994: *Nonparametric Regression and Generalized Additive Models*.
32 Vol. 182, Chapman and Hall,.
- 33 Gu, C., 2002: *Smoothing Spline ANOVA Models*. Springer-Verlag, 289 pp.
- 34 Gu, C., and G. Wahba, 1991: Minimizing GCV/GML scores with multiple smoothing parameters via the
35 Newton method. *SIAM Journal on Scientific and Statistical Computation*, **12**, 383-398.
- 36 Ide, K., P. Courtier, M. Ghil, and A. C. Lorenc, 1997: Unified notation for data assimilation operational
37 sequential and variational. *Journal of the Meteorological Society of Japan*, **75**, 181-189.
- 38 Li, H., E. Kalnay, and T. Miyoshi, 2009: Simultaneous estimation of covariance inflation and
39 observation errors within an ensemble Kalman filter. *Quarterly Journal of the Royal Meteorological
40 Society*, **135**, 523-533.
- 41 Liang, X., X. Zheng, S. Zhang, G. Wu, Y. Dai, and Y. Li, 2012: Maximum Likelihood Estimation of Inflation
42 Factors on Error Covariance Matrices for Ensemble Kalman Filter Assimilation. *Quarterly Journal of the*



- 1 *Royal Meteorological Society*, **138**, 263-273.
- 2 Liu, J., E. Kalnay, T. Miyoshi, and C. Cardinali, 2009: Analysis sensitivity calculation in an ensemble
- 3 Kalman filter. *Quarterly Journal of the Royal Meteorological Society*, **135**, 1842-1851.
- 4 Lorenz, E. N., 1996: Predictability a problem partly solved.
- 5 Lorenz, E. N., and K. A. Emanuel, 1998: Optimal sites for supplementary weather observations
- 6 simulation with a small model. *Journal of the Atmospheric Sciences*, **55**, 399-414.
- 7 Miller, R. N., M. Ghil, and F. Gauthiez, 1994: Advanced data assimilation in strongly nonlinear
- 8 dynamical systems. *Journal of the Atmospheric Sciences*, **51**, 1037-1056.
- 9 Miyoshi, T., 2011: The Gaussian approach to adaptive covariance inflation and its implementation
- 10 with the local ensemble transform Kalman filter. *Monthly Weather Review*, **139**, 1519-1534.
- 11 Reichle, R. H., 2008: Data assimilation methods in the Earth sciences. *Advances in Water Resources*,
- 12 **31**, 1411-1418.
- 13 Saltelli, A., S. Tarantola, F. Campolongo, and M. Ratto, 2004: *Sensitivity Analysis in Practice: A Guide to*
- 14 *Assessing Scientific Models*. JohnWiley & Sons, 219 pp.
- 15 Saltelli, A., and Coauthors, 2008: *Global Sensitivity Analysis: The Primer*. John Wiley & Sons, 292 pp.
- 16 Talagrand, O., 1997: Assimilation of Observations, an Introduction. *Journal of the Meteorological*
- 17 *Society of Japan*, **75**, 191-209.
- 18 Tippett, M. K., J. L. Anderson, C. H. Bishop, T. M. Hamill, and J. S. Whitaker, 2003: Notes and
- 19 correspondence ensemble square root filter. *Monthly Weather Review*, **131**, 1485-1490.
- 20 Wahba, G., and S. Wold., 1975: A completely automatic french curve. *Communications in Statistics*, **4**,
- 21 1-17.
- 22 Wahba, G., R. J. Donald, F. Gao, and J. Gong, 1995: Adaptive Tuning of Numerical Weather Prediction
- 23 Models: Randomized GCV in Three- and Four-Dimensional Data Assimilation. *Monthly Weather*
- 24 *Review*, **123**, 3358-3369.
- 25 Wand, M. P., and M. C. Jones, 1995: *Kernel Smoothing*. Chapman and Hall, 212 pp.
- 26 Wang, X., and C. H. Bishop, 2003: A comparison of breeding and ensemble transform kalman filter
- 27 ensemble forecast schemes. *Journal of the Atmospheric Sciences*, **60**, 1140-1158.
- 28 Wu, G., X. Zheng, L. Wang, S. Zhang, X. Liang, and Y. Li, 2013: A New Structure for Error Covariance
- 29 Matrices and Their Adaptive Estimation in EnKF Assimilation. *Quarterly Journal of the Royal*
- 30 *Meteorological Society*, **139**, 795-804.
- 31 Wu, G., X. Yi, X. Zheng, L. Wang, X. Liang, S. Zhang, and X. Zhang, 2014: Improving the Ensemble
- 32 Transform Kalman Filter Using a Second-order Taylor Approximation of the Nonlinear Observation
- 33 Operator. *Nonlinear Processes in Geophysics*, **21**, 955-970.
- 34 Zheng, X., 2009: An adaptive estimation of forecast error statistic for Kalman filtering data
- 35 assimilation. *Advances in Atmospheric Sciences*, **26**, 154-160.
- 36 Zheng, X., and R. Basher, 1995: Thin-plate smoothing spline modeling of spatial climate data and its
- 37 application to mapping south Pacific rainfall. *Monthly Weather Review*, **123**, 3086-3102.
- 38
- 39



1 Table 1. The time-mean values of the GAI statistics, the GCV functions and the
2 analysis RMSE over 2000 time steps.

3

EnKF schemes	Conventional EnKF	EnKF with forecast inflation
GAI	10.78%	29.21%
GCV	31.14	3.29
RMSE	4.01	1.10

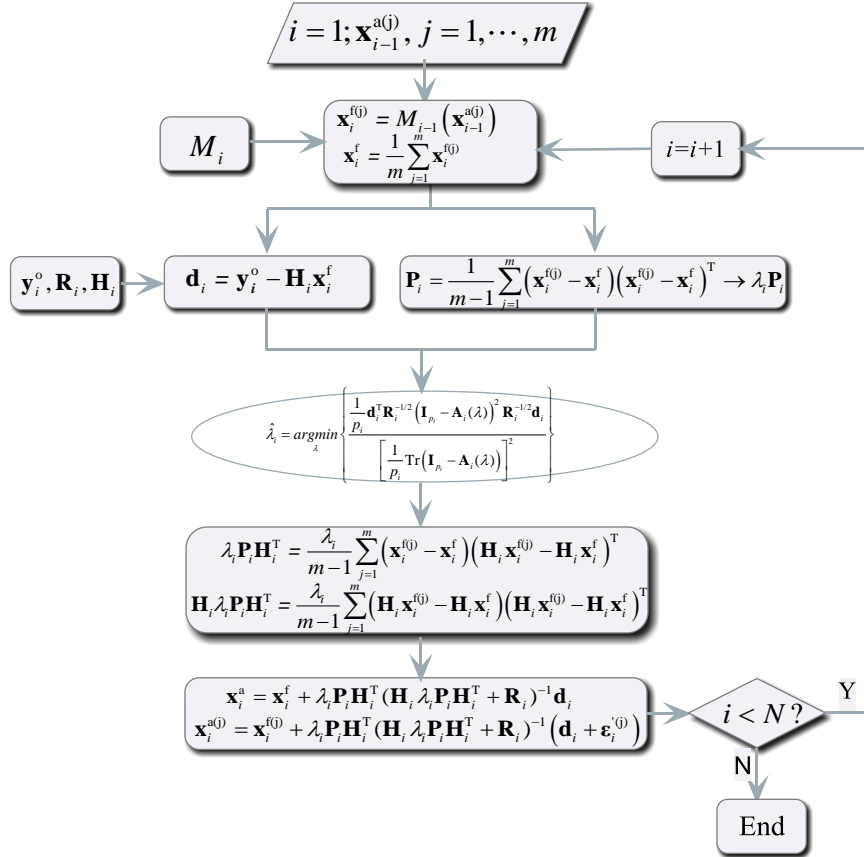
4

5

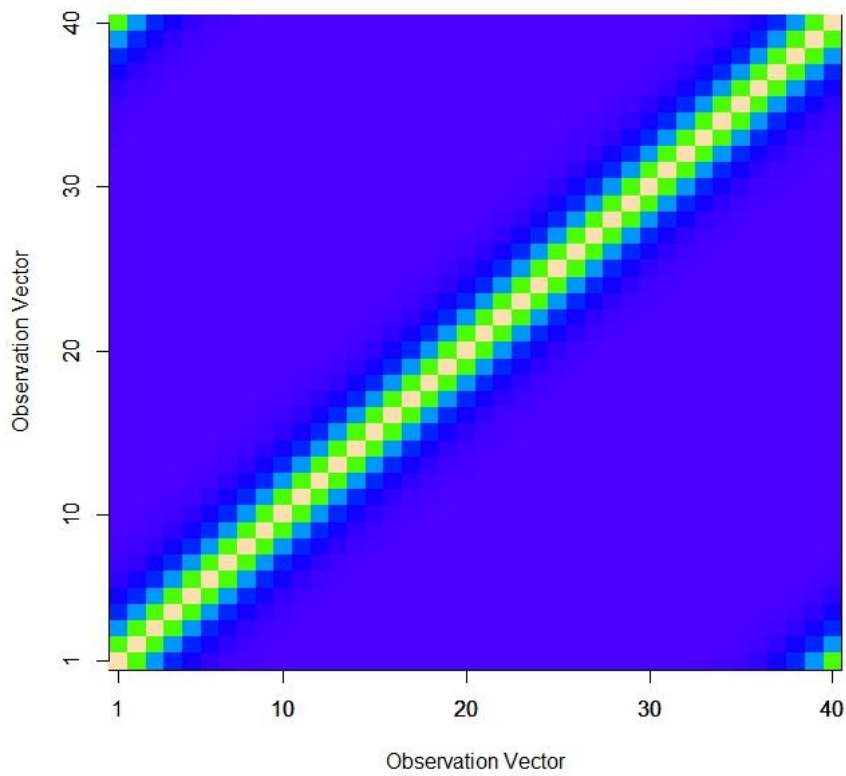
6



- 1 **Figure captions**
- 2 Figure 1. Flowchart of the proposed assimilation scheme.
- 3 Figure 2. The heat map of the observation error covariance matrix.
- 4 Figure 3. The GAI statistics of the conventional EnKF scheme and the improved
- 5 EnKF with forecast error inflation scheme.
- 6 Figure 4. The analysis RMSE of the conventional EnKF scheme and the improved
- 7 EnKF with forecast error inflation scheme.
- 8 Figure 5. The values of the GCV functions of the conventional EnKF scheme and the
- 9 improved EnKF with forecast error inflation scheme.
- 10



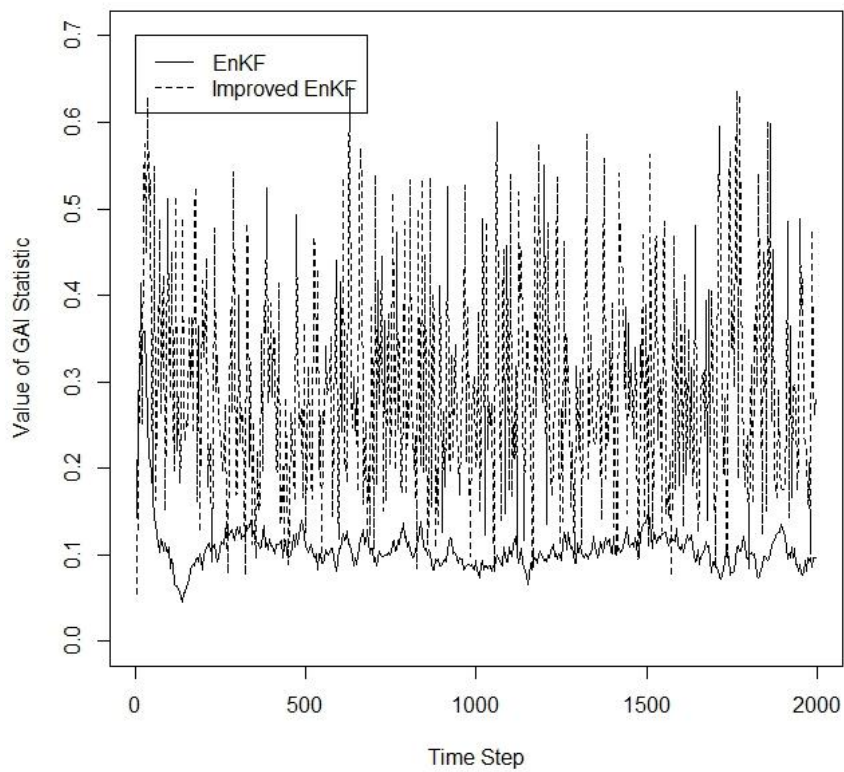
1
 2 Figure.1. Flowchart of the proposed assimilation scheme.
 3



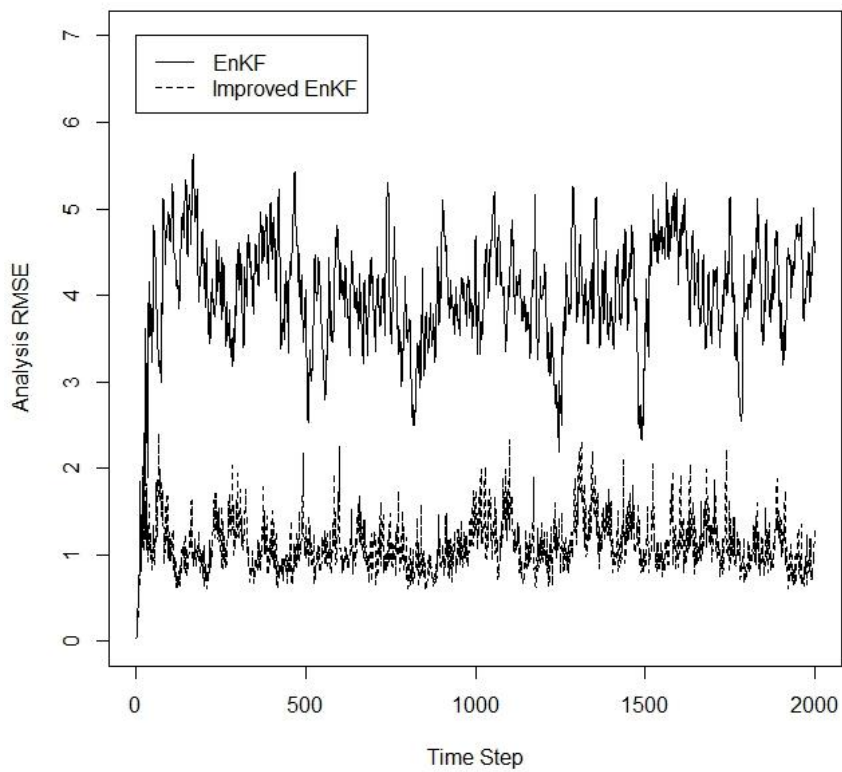
1

2 Figure 2. The heat map of the observation error covariance matrix.

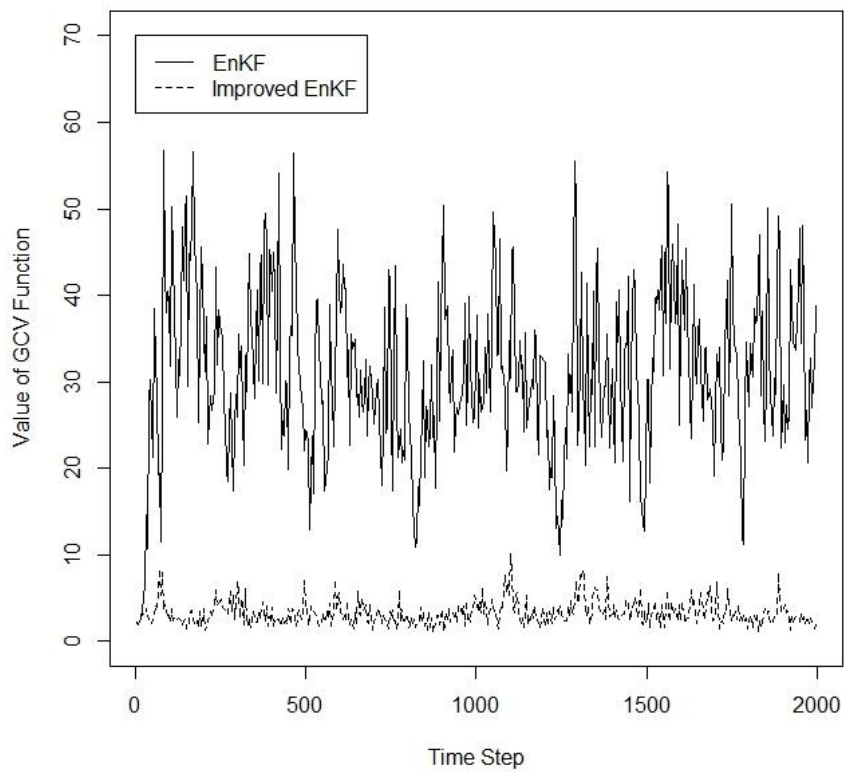
3



1
2 Figure 3. The GAI statistics of the conventional EnKF scheme and the improved
3 EnKF with forecast error inflation scheme.
4



1
2 Figure 4. The analysis RMSE of the conventional EnKF scheme and the improved
3 EnKF with forecast error inflation scheme.
4



1

2 Figure 5. The values of the GCV functions of the conventional EnKF scheme and the

3 improved EnKF with forecast error inflation scheme.

4