

DATA ASSIMILATION WITH MONTE CARLO MIXTURE KALMAN FILTER TOWARD SPACE WEATHER FORECASTING

Tomoyuki Higuchi

The Institute of Statistical Mathematics, 4-6-7 Minami-Azabu, Minato-ku, Tokyo 106-8569, Japan

ABSTRACT

Data Assimilation is a technique for a synthesis of information from a dynamic (numerical) model and observation data. It is an emerging area in earth sciences, particularly oceanography, stimulated by recent improvements in computational and modeling capabilities and the increase in the amount of available observations. Past studies for data assimilation employed a linear Gaussian state space model and applied Kalman filter. The Kalman filter based methods, however, do not allow for the strong nonlinear and/or non-Gaussian disturbance behaviors. We develop a new time dependent inversion method, which we call Monte Carlo mixture Kalman filter (MCMKF). A basic idea of MCMKF is as follows: (1) we prepare a finite number of competing state space models, each of which follows a different state space model, (2) we introduce a switching structure among these competing models. We address a capability of applying MCMKF to several topics in space physics.

1. DATA ASSIMILATION

Obtaining an accurate prediction of the upper ocean thermal structure is one of the important issue in modeling studies of the ocean circulation (and Earth's environment). A numerical ocean model-based experiment, i.e., simulation study has been conducted for this purpose. Dependent variables of the numerical ocean model can be considered as stochastic variables due to the uncertainty in the initial and boundary conditions and the imperfection of the numerical model. A natural idea to compensate for such insufficient information only via simulations is to combine observations with numerical models. Hence, a reasonable way of blending a numerical (physical) model and observation is now becoming a central issue in the earth science community.

Data Assimilation is a technique for a synthesis of information from a dynamic (numerical) model and observation data [1]. In statistical sense, data assimilation supposes two models: system model and observation model. The system and observation model correspond to large-scale numerical model-based simulations and large-scale satellite- and/or ground-based measurement systems, respectively. The data assimilation can be therefore formulated in the state

space model (SSM) [2, 3] as follows:

System model

$$\mathbf{x}_n = F_n \mathbf{x}_{n-1} + \mathbf{v}_n, \quad \mathbf{v}_n \sim N(\mathbf{0}, Q) \quad (1)$$

Observation model

$$\mathbf{y}_n = H_n \mathbf{x}_n + \mathbf{w}_n, \quad \mathbf{w}_n \sim N(\mathbf{0}, R), \quad (2)$$

where $N(\mathbf{0}, Q)$ and $N(\mathbf{0}, R)$ denote a normal distribution with a mean $\mathbf{0}$ and variance-covariance matrix of Q and R , respectively. The state of the ocean at time n is represented by a state vector \mathbf{x}_n of which dimension is n_x . \mathbf{y}_n is an observation (measurement) vector at time n . Its dimension is denoted by n_y . In data assimilation, n_x is huge such as $n_x \approx 10^5 \sim 10^6$. In contrast, n_y is much smaller than n_x ($n_x \gg n_y$), and then to get an optimal solution of \mathbf{x}_n becomes a time-dependent inversion problem. It should be noticed that n_y is extremely larger than those in the conventional state space models. For example, n_y exceeds a few thousands.

State Update —Blending a numerical (physical) model and observation—

$$\mathbf{x}_{n|n} = \mathbf{x}_{n|n-1} + K_n \mathbf{e}_{n|n-1}$$

The state of the ocean is represented by a state vector.

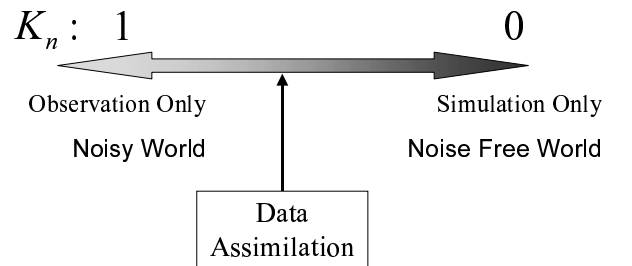


Fig. 1. Schematic Representation

Figure 1 shows a schematic representation of the data assimilation concept with a case of $n_x = n_y = 1$. At each time the state vector is updated according to this scheme (called the filtering step in the Kalman filter). $\mathbf{x}_{n|n-1}$ is a

state estimation at time n only via simulations. $e_{n|n-1}$ is a prediction error between an actual observation and predictive value of the observation based on the result of simulations. K_n is a trade off parameter to control how the simulation model accommodates an actual observation. K_n is called the Kalman gain. When $K_n = 0$, an actual observation has no effect on a simulation process. In this case we totally rely on the simulation result. On the other hand, when $K_n = 1$, any discrepancy between the predictive and real values of observations is perfectly adjusted. In this case, it is difficult to identify a dynamics inherent to a simulation model from an estimation of the state vector, because a state vector is highly sensitive to the observation errors.

The SSM has given a platform in nonstationary time series and control studies for three decades after Kalman [3]. The Network Inversion Filter (NIF) [4] which is now recognized as one of the standard techniques in the time-dependent inversion to reveal the whole time history of fault slips event with an analysis of GPS network data, is also defined by the SSM. In the meanwhile, many phenomena in earth sciences have a nonlinearity and come along with non-Gaussian fluctuations. In particular, most of problems arising in space physics have to be discussed in this context, because a media in space, namely, plasma produces a wide variety of nonlinear and non-Gaussian phenomena. Prior to a realization of a space weather forecasting, the nonlinear non-Gaussian data assimilation method is required to evaluate how much the numerical model-based approach can describe the real world phenomena. The motivation of this work is to develop a new technique for non-linear non-Gaussian data assimilation method. Evensen [5] proposed the assimilation method for strongly nonlinear problem, which is usually called the ensemble Kalman filter (EnKF). The Monte Carlo method is used to estimate the variance of the prior probability distribution function (PDF). However, the EnKF also assumes that all the required PDFs are Gaussian. It should be noticed that the EnKF is different from the extended Kalman filter (EKF) which has been applied to weakly nonlinear problems [2].

2. MODEL AND BASIC IDEA

2.1. Conditional dynamic linear model (CDLM)

One of the generalization to develop the non-Gaussian non-linear data assimilation is to introduce the conditional dynamic linear model (CDLM) [6]. The CDLM can be defined as:

$$\mathbf{x}_n = F_n(I_n)\mathbf{x}_{n-1} + \mathbf{v}_n, \quad (3)$$

$$\mathbf{y}_n = H_n(I_n)\mathbf{x}_n + \mathbf{w}_n \quad (4)$$

where $\mathbf{v}_n \sim N(\mathbf{0}, Q_n(I_n))$ and $\mathbf{w}_n \sim N(\mathbf{0}, R_n(I_n))$. The indicator vector I_n is a discrete latent variable which takes

an integer value between $1 \sim M$. Usually a number of models treated in the Mixture Kalman filter is about $2 \sim 3$, but we consider a problem of dealing with a large number of models, $M \simeq 100$. Given I_n, F_n, H_n, Q_n , and R_n are known matrices of appropriate dimension. The CDLM is a direct generalization of the dynamic linear model (DLM) and retain a capability of dealing with outliers, sudden jumps, clutters, and other nonlinear features. The CDLM includes other types of generalization of DLM, e.g., Partial non-Gaussian state space model [7], Markov switching state space model, [8] and Dynamic linear models with switching.

2.2. Monte Carlo Mixture Kalman Filter (MCMKF)

We introduce a new filtering scheme and call it as Monte Carlo mixture Kalman filter (MCMKF) that allows us to choose the optimal model from many candidates or to average over many models [9].

The MCMKF algorithm requires a stochastic model which describes a time-dependent structure for I_n . In this study, I_n is assumed to follow a stationary Markov process, i.e.,

$$p(I_n | \mathbf{I}_{1:n-1}) = p(I_n | I_{n-1}) \quad (5)$$

where $\mathbf{I}_{i:j} = (I_i, I_{i+1}, \dots, I_j)$ and $p(\cdot)$ denotes probability density function. An evolution of I_n is realized by Markov switching model with transition probability given by

$$\pi_{ij} = \Pr(I_n = j | I_{n-1} = i) \quad (6)$$

where \Pr denotes realization probability. In the following, we present an algorithm that determines time evolution of I_n .

The MCMKF algorithm consists of two steps. First, temporal variation of the probability distribution of indicator variable I_n is determined. Second, temporal variation of the probability distribution of the state vector \mathbf{x}_n is estimated following the history of I_n .

Let $\mathbf{y}_{i:j}$ and $\mathbf{I}_{i:j}$ be a set of data vectors and indicator variable from time t_i to time t_j , respectively, i.e., $\mathbf{y}_{i:j} = (\mathbf{y}_i, \mathbf{y}_{i+1}, \dots, \mathbf{y}_j)$ and $\mathbf{I}_{i:j} = (I_i, I_{i+1}, \dots, I_j)$. In MCMKF, two conditional joint distributions of $\mathbf{I}_{1:n}$: (i) predictive distribution $p(\mathbf{I}_{1:n} | \mathbf{y}_{1:n-1})$ and (ii) filter distribution $p(\mathbf{I}_{1:n} | \mathbf{y}_{1:n})$, are approximated by many ‘‘particles’’ that can be considered as independent realizations from each distribution. Let $\mathbf{I}_{1:i|k}^{(j)} = (I_{1|k}^{(j)}, I_{2|k}^{(j)}, \dots, I_{i|k}^{(j)})$ be the j th realization of the conditional distribution $p(\mathbf{I}_{1:i} | \mathbf{y}_{1:k})$. Each distribution is approximated by N_p ($N_p \gg 1$) realizations as follows:

$$\left\{ \mathbf{I}_{1:n|n-1}^{(m)} \right\}_{m=1}^{N_p} \sim p(\mathbf{I}_{1:n} | \mathbf{y}_{1:n-1}) \quad (7)$$

$$\left\{ \mathbf{I}_{1:n|n}^{(m)} \right\}_{m=1}^{N_p} \sim p(\mathbf{I}_{1:n} | \mathbf{y}_{1:n}) \quad (8)$$

where

$$\Pr(\mathbf{I}_{1:n} = \mathbf{I}_{1:n|n-1}^{(j)} | \mathbf{y}_{1:n-1}) = \frac{1}{N_p}, \quad (9)$$

$$\Pr(\mathbf{I}_{1:n} = \mathbf{I}_{1:n|n}^{(j)} | \mathbf{y}_{1:n}) = \frac{1}{N_p}. \quad (10)$$

In this study, we refer to $\{\mathbf{I}_{1:n|n-1}^{(m)}\}_{m=1}^{N_p}$ and $\{\mathbf{I}_{1:n|n}^{(m)}\}_{m=1}^{N_p}$ as ‘‘approximated predictive distribution’’ and ‘‘approximated filter distribution’’, respectively.

3. RECURSIVE CALCULATION

3.1. Indicator variable estimation

In this subsection, we show that an approximated predictive distribution at time n is obtained from an approximated filter distribution at time $n-1$. We assume that $\{\mathbf{I}_{1:n-1|n-1}^{(m)}\}_{m=1}^{N_p}$ and $\mathbf{y}_{1:n-1}$ are given. Then the probability $\Pr(\mathbf{I}_{1:n} = \mathbf{I}_{1:n|n-1}^{(j)} | \mathbf{y}_{1:n-1})$ is given as

$$\begin{aligned} \Pr(\mathbf{I}_{1:n} = \mathbf{I}_{1:n|n-1}^{(j)} | \mathbf{y}_{1:n-1}) \\ = \Pr(I_n = I_{n|n-1}^{(j)} | I_{n-1} = I_{n-1|n-1}^{(j)}) \frac{1}{N_p}. \end{aligned} \quad (11)$$

(11) indicates that $\{\mathbf{I}_{1:n|n-1}^{(m)}\}_{m=1}^{N_p}$ is obtained by sampling a realization $I_{n|n-1}^{(j)}$ with probability or weight $\Pr(I_n = I_{n|n-1}^{(j)} | I_{n-1} = I_{n-1|n-1}^{(j)})$, and setting $\mathbf{I}_{1:n|n-1}^{(j)} = (I_{1:n-1|n-1}^{(j)}, I_{n|n-1}^{(j)})$. Note that $\Pr(I_n = I_{n|n-1}^{(j)} | I_{n-1} = I_{n-1|n-1}^{(j)})$ is given by the Markovian transition probability defined by (6).

Next we show that an approximated filter distribution at time n is obtained from an approximated predictive distribution at time n . Given the observation \mathbf{y}_n , the probability $\Pr(\mathbf{I}_{1:n} = \mathbf{I}_{1:n|n-1}^{(j)} | \mathbf{y}_{1:n-1})$ is updated as follows:

$$\Pr(\mathbf{I}_{1:n} = \mathbf{I}_{1:n|n-1}^{(j)} | \mathbf{y}_{1:n}) = \frac{w_n^{(j)}}{\sum_{j=1}^{N_p} w_n^{(j)}} \quad (12)$$

where

$$w_n^{(j)} = p(\mathbf{y}_n | \mathbf{I}_{1:n} = \mathbf{I}_{1:n|n-1}^{(j)}, \mathbf{y}_{1:n-1}). \quad (13)$$

Equation (12) means that the filter distribution $p(\mathbf{I}_{1:n} | \mathbf{y}_{1:n})$ is approximated by giving weight proportional to $w_n^{(j)}$ to the j th particle of approximated predictive distribution. For the next prediction step, it is necessary to represent the approximated filter distribution with equally weighted particles $\{\mathbf{I}_{1:n|n}^{(m)}\}_{m=1}^{N_p}$. This is achieved by generating N_p particles $\{\mathbf{I}_{1:n|n}^{(m)}\}_{m=1}^{N_p}$ by resampling $\{\mathbf{I}_{1:n|n-1}^{(m)}\}_{m=1}^{N_p}$ with probability proportional to $\{w_n^{(m)}\}_{m=1}^{N_p}$.

3.2. Kalman filter

Since $\mathbf{I}_{1:n-1|n-1}^{(j)}$ is assumed to be given, the CDLM (3) and (4) reduces to a linear Gaussian state space model and thus $\mathbf{x}_{n-1|n-1}^{(j)}$ and $V_{n-1|n-1}^{(j)}$ are calculated by Kalman filter. $\mathbf{x}_{n|n-1}^{(j)}$ and $V_{n|n-1}^{(j)}$ are also calculated by Kalman filter using $\mathbf{x}_{n-1|n-1}^{(j)}$, $V_{n-1|n-1}^{(j)}$ and $I_{n|n-1}^{(j)}$. Then we have

$$\begin{aligned} p(\mathbf{x}_{n-1} | \mathbf{I}_{1:n-1} = \mathbf{I}_{1:n-1|n-1}^{(j)}, \mathbf{y}_{1:n-1}) \\ \sim N(\mathbf{x}_{n-1|n-1}^{(j)}, V_{n-1|n-1}^{(j)}) \end{aligned} \quad (14)$$

$$\begin{aligned} p(\mathbf{x}_n | \mathbf{I}_{1:n} = \mathbf{I}_{1:n|n-1}^{(j)}, \mathbf{y}_{1:n-1}) \\ \sim N(\mathbf{x}_{n|n-1}^{(j)}, V_{n|n-1}^{(j)}). \end{aligned} \quad (15)$$

Similarly the predictive distribution of data also becomes a Gaussian:

$$\begin{aligned} p(\mathbf{y}_n | \mathbf{I}_{1:n} = \mathbf{I}_{1:n|n-1}^{(j)}, \mathbf{y}_{1:n-1}) \\ \sim N(\mathbf{y}_{n|n-1}^{(j)}, W_{n|n-1}^{(j)}) \end{aligned} \quad (16)$$

where

$$\mathbf{y}_{n|n-1}^{(j)} = H_n(I_{n|n-1}^{(j)}) \mathbf{x}_{n|n-1}^{(j)} \quad (17)$$

$$\begin{aligned} W_{n|n-1}^{(j)} = H_n(I_{n|n-1}^{(j)}) V_{n|n-1}^{(j)} H_n^T(I_{n|n-1}^{(j)}) \\ + R_n(I_{n|n-1}^{(j)}). \end{aligned} \quad (18)$$

The left hand side of (16) is the weight $w_n^{(j)}$ defined in (13). Thus $w_n^{(j)}$ follows the Gaussian distribution with mean $\mathbf{y}_{n|n-1}^{(j)}$ and covariance matrix $W_{n|n-1}^{(j)}$.

3.3. Fixed lag smoother

By using the prediction and the filtering algorithm recursively, we finally obtain N_p particles $\{\mathbf{I}_{1:N_e|N_e}^{(m)}\}_{m=1}^{N_p}$ that approximate $p(\mathbf{I}_{1:N_e} | \mathbf{y}_{1:N_e})$, the posterior distribution of $\mathbf{I}_{1:N_e}$ conditioned on all of available data. Here, N_e is the number of observations. $p(\mathbf{I}_{1:N_e} | \mathbf{y}_{1:N_e})$ is called smoother distribution of $\mathbf{I}_{1:N_e}$. A sequence of each particle, $\mathbf{I}_{1:N_e|N_e}^{(j)} = [I_{1|N_e}^{(j)}, I_{2|N_e}^{(j)}, \dots, I_{N_e|N_e}^{(j)}]$, is called the trajectory.

In the particle filter, the repetition of resampling gradually decreases the number of different realizations as time passes. Therefore the shape of the distribution of the state deteriorates as a time passes. Kitagawa [10] showed that this difficulty can be eliminated by employing fixed L-lag smoother rather than fixed interval smoother. Bergman *et al.* [7] presented the other way to avoid this difficulty.

Thus following Kitagawa (1996), we modify the MCMKF filtering algorithm as follows. For fixed L , generate N_p particles $\{\mathbf{I}_{n-L:n|n}^{(m)}\}_{m=1}^{N_p}$ by the resampling of $\{\mathbf{I}_{n-L:n|n-1}^{(m)}\}_{m=1}^{N_p}$ with probability proportional to $\{w_n^{(m)}\}_{m=1}^{N_p}$ defined in (13).

It is recommended to take L not so large (say, 10 or 20 at the largest 50) [10, 11]. We adopt $L = 20$ in our application. This filtering algorithm is conceptually similar to the storing state vector algorithm in the Monte Carlo filter proposed by Kitagawa [10]. He applied the Monte Carlo approximation directly to the distribution of the state, whereas we apply the approximation to the distribution of the indicator variable.

3.4. State Vector Estimation

We present here an algorithm to estimate the state using all the N_p trajectories $\{\mathbf{I}_{1:N_e|N_e}^{(m)}\}_{m=1}^{N_p}$. In this case, $F_n(I_n^{(j)})$, $Q_n(I_n^{(j)})$, $H_n(I_n^{(j)})$ and $R_n(I_n^{(j)})$ ($j = 1, \dots, N_p$) in (3) and (4) reduce to sets of known matrices which have different time evolutions corresponding to trajectories. Thus the CDLM defined by (3) and (4) reduces to the conventional linear Gaussian state space model to which Kalman filter is applicable for state estimation. $\{\mathbf{x}_{n+1|n}^{(j)}, V_{n+1|n}^{(j)}\}_{j=1}^{N_p}$, $\{\mathbf{x}_{n|n}^{(j)}, V_{n|n}^{(j)}\}_{j=1}^{N_p}$ and $\{\mathbf{x}_{n|N_e}^{(j)}, V_{n|N_e}^{(j)}\}_{j=1}^{N_p}$ are recursively obtained by Kalman filter. Given $\{\mathbf{x}_{n|N_e}^{(j)}, V_{n|N_e}^{(j)}\}_{j=1}^{N_p}$, distribution of the final estimate for \mathbf{x}_n , $p(\mathbf{x}_n|\mathbf{y}_{1:N_e})$, is written as

$$p(\mathbf{x}_n|\mathbf{y}_{1:N_e}) = \frac{1}{N_p} \sum_{j=1}^{N_p} N(\mathbf{x}_{n|N_e}^{(j)}, V_{n|N_e}^{(j)}). \quad (19)$$

4. MCMKF IN SPACE PHYSICS

Many phenomena in space physics tend to be discussed in terms of a complex system in which a nonlinear non-Gaussian fluctuations (disturbances) play an important role. The nonlinear non-Gaussian data assimilation method needs to be developed in an attempt to realize a quantitative prediction of the space environment, i.e., space weather forecasting. In order to cope with this request, we propose the CDLM and apply the MCMKF for its state estimation. Here we give a brief summary of how the MCMKF is applied to problems in space physics.

We prepare several simulation models (simulation schemes) for understanding the specific scientific problem. Each simulation model at time n is specified by an indicator variable I_n . $Q(I_n)$ can be possibly obtained by EnKF. Then we can formulate a system model for each simulation model. The observation model is naturally introduced when the state and observation vectors, \mathbf{x}_n and \mathbf{y}_n , are defined. The problem left to us is to give $Q(I_n)$. One of the easiest way to define $Q(I_n)$ is that we utilize information on statistical fluctuations in each component of the observation vector. We denote each component of \mathbf{y}_n by $y_n(j)$ ($j = 1, \dots, n_y$). To obtain R begins with applying a lowpass filter to each observed time series $y_n(j)$. The smoothed time series is specified by $\tilde{y}_n(j)$. Then a residual series for each component is determined by $e_n(j) = y_n(j) - \tilde{y}_n(j)$. Then $Q(I_n)$

is defined by a sample variance-covariance matrix of $e_n(j)$: $Q_{jj'}(I_n) = \langle e_n(j) \cdot e_n(j') \rangle$, where $Q_{jj'}$ denotes the jj' element of Q .

In conclusion, we address that the MCMKF is designed to deal with the CDLM and then can be applicable to a wide variety of the nonlinear non-Gaussian state space models. The MCMKF allows us to integrate various type of time series models and to generate a flexible time series model automatically.

5. REFERENCES

- [1] I. Fukumori, "Data assimilation by models," in *Satellite altimetry and earth sciences: A Handbook of Techniques and Applications*, L.L. Fu and A. Cazenave, Eds., pp. 237–265. Academic Press, 2001.
- [2] B.D.O. Anderson and J.B. Moore, *Optimal Filtering*, Prentice-Hall, New Jersey, 1979.
- [3] G. Kitagawa and W. Gersch, *Smoothness Priors Analysis of Time Series*, Springer-Verlag, New York, 1996.
- [4] P. Segall and M. Matthews, "Time dependent inversion of geodetic data," *Journal of Geophysical Research*, vol. 102, pp. 22391–22409, 1997.
- [5] G. Evensen, "Sequential data assimilation with a nonlinear quasi-geostrophic model using Monte Carlo methods to forecast error statistics," *Journal of Geophysical Research*, vol. 99, pp. 10143–10162, 1994.
- [6] R. Chen and J.S. Liu, "Mixture Kalman filters," *J. R. Statist. Soc. B.*, vol. 62, pp. 493–508, 2000.
- [7] N. Bergman, A. Doucet, and N. Gordon, "Optimal estimation and Cramér-Rao bounds for partial non-Gaussian state space models," *Annals of Institute of Statistical Mathematics*, vol. 53, pp. 97–122, 2001.
- [8] C.-J. Kim and Ch.R. Nelson, *State space models with regime switching*, MIT press, Cambridge, 1999.
- [9] T. Higuchi and J. Fukuda, "Monte Carlo mixture Kalman filter and its application to space-time inversion," in *13th IFAC Symposium on System Identification*, 2003, pp. 1299–1304.
- [10] G. Kitagawa, "Monte Carlo filter and smoother for non-Gaussian nonlinear state space model," *Journal of Computational and Graphical Statistics*, vol. 5, pp. 1–25, 1996.
- [11] T. Higuchi and G. Kitagawa, "Knowledge discovery and self-organizing state space model," *IEICE Transactions on Information and Systems*, vol. E83-D, pp. 36–43, 2000.