

Kalman Filtering: Part I

Instructor: Istvan Szunyogh

Class: February 23, 2007

Recommended Readings:

- Geir Evensen, 2006: *Data Assimilation: The Ensemble Kalman Filter*, Springer, 280 pages, is a nice handbook that also provides a good summary of the history. Cautionary notes
 - There have been many important developments since book has been completed (most likely in early 2005)
 - When considering the computational cost of the alternative computational algorithms, the book does not really consider that the algorithms are usually implemented on parallel computers (this is the case for an operational NWP model)
 - A little too much credit is claimed by the author—this limits the value of the book only as a source on history
- Brian Hunt, Eric Kostelich and Istvan Szunyogh, 2007: Efficient Data Assimilation for Spatiotemporal Chaos: a Local Ensemble Kalman Filter. *Physica D*. Available from the Weather-Chaos web page.

Mathematical Formulation, following Brian Hunt (the most elegant formulation I am aware of)

The Analysis Problem:

- Consider a system governed by the ordinary differential equation

$$\frac{dx}{dt} = F(t, \mathbf{x}), \quad (1)$$

where \mathbf{x} is an m -dimensional vector representing the state of the system at a given time.

- Suppose we are given a set of (noisy) observations of the system made at various times.
- We want to determine which trajectory $\{\mathbf{x}(t)\}$ of (1) “best” fits the observations. For any given t , this trajectory gives an estimate of the system state at time t .

Notation

- Let us assume that the observations are the result of measuring quantities that depend on the system state in a known way, with Gaussian measurement errors.
- An observation at time t_j is a triple $(\mathbf{y}_j^o, H_j, \mathbf{R}_j)$, where \mathbf{y}_j^o is a vector of observed values, and H_j and \mathbf{R}_j describe the relationship between \mathbf{y}_j^o and $\mathbf{x}(t_j)$:

$$\mathbf{y}_j^o = H_j(\mathbf{x}(t_j)) + \varepsilon_j,$$

where ε_j is a Gaussian random variable with mean $\mathbf{0}$ and covariance matrix \mathbf{R}_j .

- Here, a perfect model is assumed: the observations are based on a trajectory of (1), and our problem is simply to infer which trajectory produced the observations. In a real application, the observations come from a trajectory of the physical system for which (1) is only a model.

The maximum likelihood estimate for the trajectory that best fits the observations at times $t_1 < t_2 < \dots < t_n$.

- The likelihood of a trajectory $\mathbf{x}(t)$ is proportional to

$$\prod_{j=1}^n \exp(-[\mathbf{y}_j^o - H_j(\mathbf{x}(t_j))]^T \mathbf{R}_j^{-1} [\mathbf{y}_j^o - H_j(\mathbf{x}(t_j))]),$$

since the observational errors are normally distributed and are assumed to be independent at the different observation times. The most likely trajectory is the one that maximizes this expression.

- Equivalently, the most likely trajectory is the one that minimizes the “cost function”

$$J^o(\{\mathbf{x}(t)\}) = \sum_{j=1}^n [\mathbf{y}_j^o - H_j(\mathbf{x}(t_j))]^T \mathbf{R}_j^{-1} [\mathbf{y}_j^o - H_j(\mathbf{x}(t_j))]. \quad (2)$$

Thus, the “most likely” trajectory is also the one that best fits the observations in a least square sense.

Replacing the Trajectory with the State at a Particular Time

- (2) expresses the cost J^o as a function of the trajectory $\{\mathbf{x}(t)\}$. To minimize the cost, it is more convenient to write it as a function of the system state at a particular time t .
- Let $M_{t,t'}$ be the map that propagates a solution of (1) from time t to time t' . Then

$$J_t^o(\mathbf{x}) = \sum_{j=1}^n [\mathbf{y}_j^o - H_j(M_{t,t_j}(\mathbf{x}))]^T \mathbf{R}_j^{-1} [\mathbf{y}_j^o - H_j(M_{t,t_j}(\mathbf{x}))] \quad (3)$$

expresses the cost in terms of the system state \mathbf{x} at time t .

- To estimate the state at time t , we attempt to minimize J_t^o .

Comments:

- In practice the observations do not have to be all collected at t_n . In a typical implementation, at t_n we assimilate all observations that were collected at times t in the window $t_n - \Delta t/2 < t < t_n + \Delta t/2$, where $\Delta t = t_j - t_{j-1}$, $j = 2, \dots, n$.
- For a nonlinear model, there is no guarantee that a unique minimum exists.
- Even if a minimum exist, evaluating J_t^o is apt to be computationally expensive, and minimizing it may be impractical.
- But, if both the model and the observation operators H_j are linear, the minimization is quite tractable, because J_t^o is then quadratic. Furthermore, one can compute the minimum by an iterative method, namely the Kalman Filter (Kalman 1960; Kalman and Bucy 1961).

Linear Scenario: the Kalman Filter

- In the linear scenario, we can write $M_{t,t'}(\mathbf{x}) = \mathbf{M}_{t,t'}\mathbf{x}$ and $H_j(\mathbf{x}) = \mathbf{H}_j\mathbf{x}$ where $\mathbf{M}_{t,t'}$ and \mathbf{H}_j are matrices.
- We now describe how to perform
 - a forecast step from time t_{n-1} to time t_n
 - followed by an analysis step at time t_n ,
 - in such a way that if we start with the most likely system state, given the observations up to time t_{n-1} , we end up with the most likely state given the observations up to time t_n .

The estimate of the state and the uncertainty at t_{n-1}

- Suppose the analysis at time t_{n-1} has produced a state estimate $\bar{\mathbf{x}}_{n-1}^a$ and an associated covariance matrix \mathbf{P}_{n-1}^a . In probabilistic terms, $\bar{\mathbf{x}}_{n-1}^a$ and \mathbf{P}_{n-1}^a represent the mean and covariance of a Gaussian probability distribution that represents the relative likelihood of the possible system states given the observations from time t_1 to t_{n-1} .
- Algebraically, what we assume is that for some constant c ,

$$\sum_{j=1}^{n-1} [\mathbf{y}_j^o - \mathbf{H}_j \mathbf{M}_{t_{n-1}, t_j} \mathbf{x}]^T \mathbf{R}_j^{-1} [\mathbf{y}_j^o - \mathbf{H}_j \mathbf{M}_{t_{n-1}, t_j} \mathbf{x}] = \quad (4)$$

$$= [\mathbf{x} - \bar{\mathbf{x}}_{n-1}^a]^T (\mathbf{P}_{n-1}^a)^{-1} [\mathbf{x} - \bar{\mathbf{x}}_{n-1}^a] + c.$$

In other words, the analysis at time t_{n-1} has “completed the square” to express the part of the quadratic cost function $J_{t_{n-1}}^o$ that depends on the observations up to that time as a single quadratic form plus a constant.

- The Kalman Filter determines $\bar{\mathbf{x}}_n^a$ and \mathbf{P}_n^a such that an analogous equation holds at time t_n .

The Kalman Filter I

- We propagate the analysis state estimate $\bar{\mathbf{x}}_{n-1}^a$ and its covariance matrix \mathbf{P}_{n-1}^a using the forecast model to produce a background state estimate $\bar{\mathbf{x}}_n^b$ and covariance \mathbf{P}_n^b for the next analysis:

$$\bar{\mathbf{x}}_n^b = \mathbf{M}_{t_{n-1}, t_n} \bar{\mathbf{x}}_{n-1}^a, \quad (5)$$

$$\mathbf{P}_n^b = \mathbf{M}_{t_{n-1}, t_n} \mathbf{P}_{n-1}^a \mathbf{M}_{t_{n-1}, t_n}^T. \quad (6)$$

- Under a linear model, a Gaussian distribution of states at one time propagates to a Gaussian distribution at any other time, and the equations above describe how the model propagates the mean and covariance of such a distribution.

The Kalman Filter II

- Next, we want to rewrite the cost function $J_{t_n}^o$ given by (3) in terms of the background state estimate and the observations at time t_n . (This step is often formulated as applying Bayes' Rule to the corresponding probability density functions.) In (4), \mathbf{x} represents a system state at time t_{n-1} . In our expression for $J_{t_n}^o$, we want \mathbf{x} to represent a system state at time t_n

- Using (5) and (6) yields that part of the cost function at t_n that reflects the effect of observations collected up to t_n

$$\sum_{j=1}^{n-1} [\mathbf{y}_j^o - \mathbf{H}_j \mathbf{M}_{t_n, t_j} \mathbf{x}]^T \mathbf{R}_j^{-1} [\mathbf{y}_j^o - \mathbf{H}_j \mathbf{M}_{t_n, t_j} \mathbf{x}] = [\mathbf{x} - \bar{\mathbf{x}}_n^b]^T (\mathbf{P}_n^b)^{-1} [\mathbf{x} - \bar{\mathbf{x}}_n^b] + c.$$

- It follows that the total cost function at t_n is

$$J_{t_n}^o(\mathbf{x}) = [\mathbf{x} - \bar{\mathbf{x}}_n^b]^T (\mathbf{P}_n^b)^{-1} [\mathbf{x} - \bar{\mathbf{x}}_n^b] + [\mathbf{y}_n^o - \mathbf{H}_n \mathbf{x}]^T \mathbf{R}_n^{-1} [\mathbf{y}_n^o - \mathbf{H}_n \mathbf{x}] + c. \quad (7)$$

where the second term reflects the effects of observations collected at t_n

The Kalman Filter III

- To complete the data assimilation cycle, we determine the state estimate $\bar{\mathbf{x}}_n^a$ and its covariance \mathbf{P}_n^a so that

$$J_{t_n}^o(\mathbf{x}) = [\mathbf{x} - \bar{\mathbf{x}}_n^a]^T (\mathbf{P}_n^a)^{-1} [\mathbf{x} - \bar{\mathbf{x}}_n^a] + c'$$

for some constant c' .

- Equating the terms of degree 2 in \mathbf{x} , we get

$$\mathbf{P}_n^a = \left[(\mathbf{P}_n^b)^{-1} + \mathbf{H}_n^T \mathbf{R}_n^{-1} \mathbf{H}_n \right]^{-1}. \quad (8)$$

- Equating the terms of degree 1, we get

$$\bar{\mathbf{x}}_n^a = \mathbf{P}_n^a \left[(\mathbf{P}_n^b)^{-1} \bar{\mathbf{x}}_n^b + \mathbf{H}_n^T \mathbf{R}_n^{-1} \mathbf{y}_n^o \right]. \quad (9)$$

- The last equation in some sense (consider, for example, the case where \mathbf{H}_n is the identity matrix) expresses the analysis state estimate as a weighted average of the background state estimate and the observations, weighted according to the inverse covariance of each.

The Kalman Filter IV

Equations (8) and (9) can be written in many different but equivalent forms

- Using (8) to eliminate $(\mathbf{P}_n^b)^{-1}$ from (9) yields

$$\bar{\mathbf{x}}_n^a = \bar{\mathbf{x}}_n^b + \mathbf{P}_n^a \mathbf{H}_n^T \mathbf{R}_n^{-1} (\mathbf{y}_n^o - \mathbf{H}_n \bar{\mathbf{x}}_n^b) = \bar{\mathbf{x}}_n^b + \mathbf{K} (\mathbf{y}_n^o - \mathbf{H}_n \bar{\mathbf{x}}_n^b) \quad (10)$$

- The matrix $\mathbf{K} = \mathbf{P}_n^a \mathbf{H}_n^T \mathbf{R}_n^{-1}$ is called the “Kalman gain”. It multiplies the difference between the observations at time t_n and the values predicted by the background state estimate to yield the increment between the background and analysis state estimates.
- Rearranging (8) yields

$$\mathbf{P}_n^a = (\mathbf{I} + \mathbf{P}_n^b \mathbf{H}_n^T \mathbf{R}_n^{-1} \mathbf{H}_n)^{-1} \mathbf{P}_n^b = (\mathbf{I} - \mathbf{K} \mathbf{H}) \mathbf{P}_n^b. \quad (11)$$

This expression is better than the previous one from a practical point of view, since it does not require inverting \mathbf{P}_n^b .

The Nonlinear Scenario: The Extended Kalman Filter

- Many approaches to data assimilation for nonlinear problems are based on the Kalman Filter, or at least on minimizing a cost function similar to (7).
- At a minimum, a nonlinear model forces a change in the forecast equations (5) and (6), while nonlinear observation operators H_n force a change in the analysis equations (10) and (11)
- The Extended Kalman Filter (see, for example, Jazwinski 1970) computes $\bar{\mathbf{x}}_n^b = M_{t_{n-1}, t_n}(\bar{\mathbf{x}}_{n-1}^a)$ using the nonlinear model, but computes \mathbf{P}_n^b using the linearization $\mathbf{M}_{t_{n-1}, t_n}$ of M_{t_{n-1}, t_n} around $\bar{\mathbf{x}}_{n-1}^a$. The analysis then uses the linearization \mathbf{H}_n of H_n around $\bar{\mathbf{x}}_n^b$.

Difficulties with the Implementation of the Extended Kalman Filter

- It is not easy to linearize the dynamics for a complex, high-dimensional model, such as a global weather prediction model.
- The number of model variables m is several million, and as a result the $m \times m$ matrix inverse required by the analysis cannot be performed in a reasonable amount of time.
- The use of the linear evolution equations can lead to an unbounded linear instability (see chapter 4.2.3 in Evensen 2006).

Practical Implementations at the NWP Centers

- Approaches used in operational weather forecasting generally eliminate for pragmatic reasons the time iteration of the Kalman Filter.
- NCEP/NWS: data assimilation is done every 6 hours with a “3D-VAR” method, in which the background covariance \mathbf{P}_n^b is replaced by a constant matrix \mathbf{B} . The 3D-VAR cost function also includes a nonlinear observation operator H_n , and is minimized numerically to produce the analysis state estimate \mathbf{x}_n^a .
- The “4D-VAR” method (e.g., Le Dimet and Talagrand 1986; Talagrand and Courtier 1987) used by the European Centre for Medium-Range Weather Forecasts uses a cost function that includes a constant-covariance background term as in 3D-VAR together with a sum like (2) accounting for the observations collected over a 12 hour time span.

Ensemble Kalman Filtering:

The key idea of ensemble Kalman filtering (e.g., Evensen 1994; Evensen 2006) is to choose at time t_{n-1} an ensemble of initial conditions whose spread around $\bar{\mathbf{x}}_{n-1}^a$ characterizes the analysis covariance \mathbf{P}_{n-1}^a , propagate each ensemble member using the nonlinear model, and compute \mathbf{P}_n^b based on the resulting ensemble at time t_n . Thus like the Extended Kalman Filter, the (approximate) uncertainty in the state estimate is propagated from one analysis to the next.

Potential Advantages of an Ensemble Kalman Filter:

- 3D-VAR does not propagate the uncertainty at all
- 4D-VAR propagates uncertainty only for a limited time.
- 4D-VAR and the Extended Kalman Filter require linearization of the model. (The EnKF does not require such linearization.)

Potential Disadvantage of an Ensemble Kalman Filter:

- Assuming that computational resources restrict the number of ensemble members k to be much smaller than the number of model variables m , this can be a severe limitation. On the other hand, if this limitation can be overcome, then the analysis can be performed in a much lower-dimensional space (k versus m).

Notation I

- We start with an ensemble $\{\mathbf{x}_{n-1}^{a(i)} : i = 1, 2, \dots, k\}$ of m -dimensional model state vectors at time t_{n-1} .
- We assume the ensemble to be chosen so that its average represents the analysis state estimate.
- We evolve each ensemble member according to the non-linear model to obtain a background ensemble $\{\mathbf{x}_n^{b(i)} : i = 1, 2, \dots, k\}$ at time t_n :

$$\mathbf{x}_n^{b(i)} = M_{t_{n-1}, t_n}(\mathbf{x}_{n-1}^{a(i)}).$$

- For the rest of this lecture, we will discuss what to do at the analysis time t_n , and so we drop the subscript n . Thus, for example, H and \mathbf{R} will represent respectively the observation operator and the observation error covariance matrix at the analysis time.
- Let ℓ be the number of scalar observations used in the analysis.

Notation II

- For the background state estimate and its covariance we use the sample mean and covariance of the background ensemble:

$$\bar{\mathbf{x}}^b = k^{-1} \sum_{i=1}^k \mathbf{x}^{b(i)},$$

$$\mathbf{P}^b = (k-1)^{-1} \sum_{i=1}^k (\mathbf{x}^{b(i)} - \bar{\mathbf{x}}^b)(\mathbf{x}^{b(i)} - \bar{\mathbf{x}}^b)^T = (k-1)^{-1} \mathbf{X}^b (\mathbf{X}^b)^T, \quad (12)$$

where \mathbf{X}^b is the $m \times k$ matrix whose i th column is $\mathbf{x}^{b(i)} - \bar{\mathbf{x}}^b$.

Notation III

- The analysis must determine not only a state estimate $\bar{\mathbf{x}}^a$ and covariance \mathbf{P}^a , but also an ensemble $\{\mathbf{x}^{a(i)} : i = 1, 2, \dots, k\}$ with the appropriate sample mean and covariance:

$$\bar{\mathbf{x}}^a = k^{-1} \sum_{i=1}^k \mathbf{x}^{a(i)},$$

$$\mathbf{P}^a = (k-1)^{-1} \sum_{i=1}^k (\mathbf{x}^{a(i)} - \bar{\mathbf{x}}^a)(\mathbf{x}^{a(i)} - \bar{\mathbf{x}}^a)^T = (k-1)^{-1} \mathbf{X}^a (\mathbf{X}^a)^T, \quad (13)$$

where \mathbf{X}^a is the $m \times k$ matrix whose i th column is $\mathbf{x}^{a(i)} - \bar{\mathbf{x}}^a$.

- We will describe how to determine $\bar{\mathbf{x}}^a$ and \mathbf{P}^a for a (possibly) nonlinear observation operator H in a way that agrees with the Kalman Filter equations (10) and (11) in the case that H is linear.

Localization I

- If the ensemble has k members, then the background covariance matrix \mathbf{P}^b given by (12) describes nonzero uncertainty only in the k -dimensional subspace spanned by the ensemble, and a global analysis will allow adjustments to the system state only in this subspace.
- If the system is high-dimensionally unstable, then forecast errors will grow in directions not accounted for by the ensemble, and these errors will not be corrected by the analysis.
- In a sufficiently small local region, the system may behave like a low-dimensionally unstable system driven by the dynamics in neighboring region; such behavior was observed for a global weather model
- By allowing the local analyses to choose different linear combinations of the ensemble members in different regions, the global analysis is not confined to the k -dimensional ensemble space and instead explores a much higher dimensional space (Fukumori 2002; Ott et al. 2004)

Localization II

- Also, limited sample size provided by an ensemble will produce spurious correlations between distant locations in the background covariance matrix \mathbf{P}^b (Houtekamer and Mitchell 1998; Hamill et al. 2001). Unless they are suppressed, these spurious correlations will cause observations from one location to affect, in an essentially random manner, the analysis an arbitrarily large distance away.
- Localization also allows the analysis to be done more efficiently as a parallel computation (Keppene 2000; Ott et al. 2004).
- Localization is generally done either explicitly, considering only the observations from a region surrounding the location of the analysis (Kalnay and Toth 1994; Houtekamer and Mitchell 1998; Keppene 2000; Anderson 2001; Ott et al. 2004), or implicitly, by multiplying the entries in \mathbf{P}^b by a distance-dependent function that decays to zero beyond a certain distance (Hamill et al. 2000; Houtekamer and Mitchell 2001; Whitaker and Hamill 2002).

Choosing the analysis ensemble

Once $\bar{\mathbf{x}}^a$ are \mathbf{P}^a specified, there are still many possible choices of an analysis ensemble (or equivalently, a matrix \mathbf{X}^a that satisfies (13) and the sum of whose columns is zero):

- “Perturbed Observations” method, often called Ensemble Kalman Filter (EnKF, Burgers et al.1998; Houtekamer and Mitchell 1998)
- Square Root Filters
 - Ensemble Adjustment Kalman Filter (EnAKF, Anderson 2001)
 - Ensemble Transform Kalman Filter (ETKF, Bishop et al., 2001)
 - Ensemble Square Root Filter (EnSQR, Whitaker and Hamill 2002)
 - Local Ensemble Kalman Filter (LEKF, Ott et al., 2004) and Local Ensemble Transform Kalman Filter (LETKF, Hunt et al. 2007)

The “perturbed observation” methods

- The perturbed observation techniques apply (10) to both the ensemble mean

$$\bar{\mathbf{x}}^a = \bar{\mathbf{x}}^b + \mathbf{K}(\mathbf{y}^\circ - \mathbf{H}\bar{\mathbf{x}}^b), \quad (14)$$

and the the ensemble perturbations

$$\mathbf{x}^{a(i)} = \mathbf{x}^{b(i)} + \mathbf{K}(\mathbf{y}^{\circ(i)} - \mathbf{H}\mathbf{x}^{b(i)}), \quad (15)$$

where $\mathbf{y}^{\circ(i)}$ denotes the set of randomly generated observational errors for the i -th ensemble member.

Is it necessary to perturb the observations?

- Yes. The first ensemble Kalman filter by Evensen (1994) was incorrectly formulated in the sense that it used \mathbf{y}^o and \mathbf{K} in both (14) and (15), which led to a collapse of the ensemble
- Unless the observations are artificially perturbed so that each ensemble member is updated using different random realization of the perturbed observations, this approach generates an analysis ensemble whose sample covariance is smaller than \mathbf{P}^a (Burgers et al. 1998; Houtekamer and Mitchell 1998). When $\mathbf{y}^{o(i)} = 0$, (15) becomes

$$\mathbf{x}^{a(i)} = (\mathbf{I} - \mathbf{KH})\mathbf{x}^{b(i)}, \quad (16)$$

- Based on this definition of the ensemble perturbations, the analysis error covariance matrix is

$$\mathbf{P}_n^a = (\mathbf{I} - \mathbf{KH})\mathbf{P}_n^b(\mathbf{I} - \mathbf{KH})^T. \quad (17)$$

instead of the correct

$$\mathbf{P}_n^a = (\mathbf{I} - \mathbf{KH})\mathbf{P}_n^b. \quad (18)$$

Ensemble square-root filters

As it was observed by many, the main weakness of the "perturbed observations" approach is that it estimates a known object, \mathbf{R} , with a statistical sample. Thus, it may require a large ensemble to recover \mathbf{R} from the randomly generated observational noise. Ensemble square-root filters use deterministic algorithms to generate an analysis ensemble with the desired sample mean and covariance. As such, their analyses coincide exactly with the standard Kalman Filter in the linear scenario of the previous section.

The Ensemble Square Root Kalman Filter

- The basic idea of Whitaker and Hamill (2002) is to use the unperturbed observations \mathbf{y}° in both (14) and (15), and modify \mathbf{K} in (15) such that the resulting set of perturbations satisfy (13).

- The update equation for the perturbations is

$$\mathbf{x}^{a(i)} = \mathbf{x}^{b(i)} + \tilde{\mathbf{K}}(\mathbf{y}^\circ - \mathbf{H}\mathbf{x}^{b(i)}), \quad (19)$$

where $\tilde{\mathbf{K}} = \alpha\mathbf{K}$, $0 \leq \alpha \leq 1$ and the update equations (14) and (19) are applied serially (one by one) to the observations.

- The factor α is determined by the Potter (1964) formula

$$\alpha = \left(1 + \sqrt{\mathbf{R}/(\mathbf{H}\mathbf{P}^b\mathbf{H}^T + \mathbf{R})}\right)^{-1} \quad (20)$$

- Note that $\mathbf{H}\mathbf{P}^b\mathbf{H}^T$ and \mathbf{R} are scalars representing the background and observational uncertainties at the observational locations.

General Comments on the serial algorithms

- The assimilation of an observation can change several state vector components simultaneously, thus, in principle, any state vector component can change until there are observations left to assimilate. This may make an efficient implementation on a parallel computer difficult.
- However, implementation of a data thinning algorithm or a quality control algorithm may be easier for a serial algorithm.