The Ensemble Kalman filter



Part II: Practicalities

Alison Fowler (Based on notes by Ross Bannister)

Data-assimilation training course. 21st-24th June, University of Reading

Quick quiz!

Sampling error

- The ensemble Kalman Filter theory assumes that the ensemble is large enough to give an accurate estimate of the sample mean and covariance, \bar{x} and P.
- Even for a two variable model a large sample size is needed to accurately estimate the mean and covariance:



Example: True distribution $\mathbf{x} \sim N([1, 1], \mathbf{I}_2)$

There are numerous consequences to under sampling in the EnKF

1. There may be a **bias** in the **ensemble mean**.

2. The forecast ensemble spread, which defines \mathbf{P}^{f} , will be subject to sampling error (illustrated by the 2d example)

Recall

 $\mathbf{x}_{k}^{(i),a} = \mathbf{x}_{k}^{(i),f} + \mathbf{P}_{k}^{f} \mathbf{H}^{T} (\mathbf{H}\mathbf{P}_{k}^{f} \mathbf{H}^{T} + \mathbf{R})^{-1} (\mathbf{y}_{k} + \boldsymbol{\epsilon}_{y}^{(i)} - \mathbf{H}\mathbf{x}_{k}^{(i),f})$

• If the spread $(\mathbf{P}_k^{\mathbf{f}})$ is too large the analysis ensemble will over fit to the observations.

• If the spread (\mathbf{P}_k^f) is too small, the ensemble will under fit to the observations. If the ensemble repeatedly underestimates the forecast error and the information in the observations is ignored, then it is difficult to regain spread in the ensemble. This is called **'filter divergence'**.



1.

2.

*

A system without sampling error can correctly follow the observations within their error.

time

Illustration of sampling error leading to filter divergence

1. Sample estimate of prior spread is too small

2. Too much confidence in the prior means the analysis underfits the observations and the posterior spread is too small

3. Overconfident posterior leads to an overconfident prior at the next assimilation time, which is exasperated further by under sampling.

4. On each cycle the overconfidence is propagated and worsened until there is no spread in the prior and the analysis is no longer able to use the information in the observations.



3. The **correlation will be subject to sampling error**. Implying that observations can influence regions and variables that they shouldn't.





From Houtekamer & Mitchell (1998)



From Bannister, Migliorini & Dixon (2011)

4. $\mathbf{P}_k^{\mathrm{f}}$ is rank deficient

The analysis increments are in the sub-space spanned by the forecast ensemble

• The analysis increments are given by

$$\mathbf{x}_{k}^{(i),a} - \mathbf{x}_{k}^{(i),f} = \mathbf{P}_{k}^{f} \mathbf{H}^{T} (\mathbf{H} \mathbf{P}_{k}^{f} \mathbf{H}^{T} + \mathbf{R})^{-1} (\mathbf{y}_{k} + \boldsymbol{\epsilon}_{y}^{(i)} - \mathbf{H} \mathbf{x}_{k}^{(i),f})$$

- The analysis increments are therefore a linear combination of the forecast error perturbations.
- Therefore, even if the observations indicate otherwise, the analysis is restricted to space spanned by the ensemble which has at the most a dimension of *N*-1.

Possible solutions

- 1. Use more ensemble members (see Miyoshi et al. 2014)
- 2. Re-centre the ensemble around a deterministic analysis e.g., from 4DVar.
 - Addresses problem of bias in the ensemble mean
- 3. Ensemble inflation
 - Addresses problem of filter divergence
- 4. Localization
 - Addresses problem of spurious correlations
 - Splits problem into quasi-independent problems, increasing the rank of forecast perturbation matrix.
- 5. Combine ensemble with variational approaches (see tomorrow's lectures)
 - These are known as hybrid methods

Focus of this lecture

Ensemble inflation

Ways to inflate

- Additive inflation (Mitchell and Houtekamer, 2000; Corazza et al., 2003)
 - At each model time step add a random perturbation using similar ideas to representing model error given in the last lecture

$$\mathbf{x}_k^{(i)} = M_{t_{k-1} \to t_k}(\mathbf{x}_{k-1}^{(i)}) + \mathbf{\eta}_k^{(i)}$$
, where $\mathbf{\eta} \sim N(\mathbf{0}, \mathbf{Q})$

• Multiplicative inflation (Anderson and Anderson, 1999)

$$\mathbf{P}_{\text{inflated}}^{\text{f}} = (1+\rho)^2 \mathbf{P}^{\text{f}}, \qquad \rho > 0$$

- Relaxation to prior ensemble (Zhang *et al.,* 2004; Whitaker and Hamill, 2012)
 - Only accept part of the spread reduction proposed by

$$\mathbf{P}_k^{\mathrm{a}} = (\mathbf{I} - \mathbf{K}_k \mathbf{H}) \mathbf{P}_k^{\mathrm{f}}$$

• Two methods: relaxation to prior perturbation and relaxation to prior spread



Tuning the inflation factor – validation of ensemble spread

Method 1: rank histograms (Hamill, T., 2001):

- For the ensemble to be reliable it is assumed that it is sampling the same distribution as the truth.
- A rank histogram is constructed by considering a point in space that is well observed.
 - The values of the ensemble members at that point are ranked from highest to lowest creating N-1 bins.

Mem.

• Then each observation is binned to give a frequency diagram.

Interpretation:

- Concave shape- the ensemble is under spread
- Convex shaped- the ensemble is overspread
- Flat- the ensemble is correctly spread
- Asymmetric- the ensemble is biased



Rank histograms for surface precipitation rate rate. From Migliorini et al. (2011).

0.15

Tuning the inflation factor – validation of ensemble spread

Method 2: Covariance matching

• This checks that the sum of the spread in the background ensemble and observation error variance match with the variance of the innovations (e.g. Houtekamer et al. 2005).

 $E\left[\left(y - \mathbf{H}\mathbf{x}^{\mathrm{f}}\right)\left(y - \mathbf{H}\mathbf{x}^{\mathrm{f}}\right)^{\mathrm{T}}\right] \approx \mathbf{R} + \mathbf{H}\mathbf{P}^{\mathrm{f}} \mathbf{H}^{\mathrm{T}}$



The idea of covariance matching has led to various schemes for adaptive covariance inflation e.g. Kotsuki et al. 2017

FIG. 5. Comparison of error amplitudes that have been averaged over a 10-day experimental period. Shown are the predicted innovation std dev (solid) that should match the observed innovation std dev (dashed–dotted). The predicted std dev is computed from the rms observational error (dotted) and the rms ensemble spread (dashed).

Localisation



The aim of localisation is to restrict the influence of observations to just a physically realistic region.

Two ways of doing this are:

- **P**^f-localisation
 - This modifies the forecast error covariance matrix to reduce long-range correlations.
- R-localisation
 - This restricts observations which are allowed to influence each grid point.



P^f- localization (Houtekamer and Mitchell, 2001)

- In practice cannot act on $P^{\rm f}$ directly

$$\mathbf{K} = \boldsymbol{\rho} \circ (\mathbf{P}^{f} \ \mathbf{H}^{T}) (\boldsymbol{\rho} \circ (\mathbf{H}\mathbf{P}^{f} \ \mathbf{H}^{T}) + \mathbf{R})^{-1}$$

- Need to choose function ρ and length scales, this may be state-dependent
- Not clear how to define distance between observations which have no clearly defined location in space, e.g. satellite observations
- Not clear how to deal with multivariate covariances
- Can affect the balance e.g. to conserve geostrophic balance length scale O(1000)km must be used in the horizontal

R-localization (Hunt et al. 2007)

 Localisation is performed by gradually increasing the observation error variances using the positive exponential function:

$$f_{Rloc} = exp\left[\frac{d(i,j)^2}{2L^2}\right]$$

where

- d(i,j,) is the distance between observation *i* and model grid point *j*.
- *L* is the localisation lengthscale.
- This is the method used by the LETKF
- The optimal lengthscale for R-localisation is found to be shorter than for Pf-localisation (Greybush, 2011).
- Using the optimal lengthscales, R-localisation and Pflocalisation have comparable performance in terms of analysis rmse and balance.



Hybrid methods

Hybrid methods combine the best parts of the EnKF (flow-dependent P^f) with the best parts of variational methods (full rank B).

The earliest hybrid method was proposed by Hamil and Snyder (2004), in which the representation of the error covariance of the prior information is a weigthed combination of the flow-dependent estimate from the EnKF, **P**^f, and the full rank estimate used in variational methods **P**^s

$$\mathsf{P}_{hybrid} = eta \mathsf{P}^s + (1 - eta) \mathsf{P}^f$$

where β is a tunable parameter.

Note localisation and inflation of the ensemble are still necessary.

Summary

- Ensemble data assimilation relies on a sample estimate of the mean and covariance of forecast distribution. This allows it to provide a flow-dependent estimate of the forecast uncertainty.
- If the ensemble size is much smaller than the size of the state then sampling error becomes an issue
 - Rank deficiency
 - Analysis increments lie in the subspace of the ensemble
 - Filter divergence
 - Spurious correlations
- To make ensemble DA practical need
 - Ensemble inflation
 - Localisation
 - ...Hybrid methods

Further reading

- Anderson JL, Anderson SL. 1999. A Monte Carlo implementation of the nonlinear filtering problem to produce ensemble assimilations and forecasts. Mon. Weather Rev. 127: 2741–2758.
- Bannister et al. 2011: Ensemble prediction for nowcasting with a convection-permitting model II: forecast error statistics, Tellus 63A, 497-512.
- Corazza et al.. 2003. Use of the breeding technique to estimate the structure of the analysis 'errors of the day'. Nonlinear Processes Geophys. 10: 233–243.
- Greybush, 2011: Balance and ensemble Kalman filter localisation techniques. Mon. Wea. Rev., 139, 511-522.
- Hamill, 2001: Interpretation of Rank Histograms for Verifying Ensemble Forecasts. Mon. Wea. Rev., 129,550-560.
- Hamill, 2006: Ensemble-based atmospheric data assimilation. In *Predictability of Weather and Climate*, Palmer T, Hagedorn R (eds). Cambridge University Press: Cambridge; 124–156.
- Houtekamer and Zhang (2016) Review of the Ensemble Kalman Filter for Atmospheric Data Assimilation . Mon. Weather Rev., 144, 4489–4532
- Houtekamer and Mitchell 1998: Data assimilation using an ensemble Kalman Filter technique, Mon. Wea. Rev. 126, 796-811.
- Houtekamer and Mitchell 2001: A sequential ensemble Kalman filter for atmospheric data assimilation. Mon. Wea. Rev., 129, 123–137
- Hunt et al. 2007: Efficient data assimilation for spatiotemporal chaos: A local ensemble transform Kalman filter. Physica D, 230, 112–126,
- Kotsuki et al. 2017: Adaptive covariance relaxation methods for ensemble data assimilation: experiments in the real atmosphere. Q.J.R.Meteorol.Soc.143: 2001 2015.
- Migliorini et al. 2011: Ensemble prediction for nowcasting with a convection-permitting model I: description of the system and the impact of radarderived surface precipitation rates, Tellus 63A, 468-496.
- Mitchell and Houtekamer 2000. An adaptive ensemble Kalman filter. Mon. Weather Rev. 128: 416–433.
- Miyoshi et al. 2014: The 10,240-member ensemble Kalman filtering with an intermediate AGCM. Geophys. Res. Lett., 41, 5264-5271.
- Whitaker and Hamill 2012. Evaluating methods to account for system errors in ensemble data assimilation. Mon. Weather Rev. 140: 3078–3089.
- Zhang et al. 2004. Impacts of initial estimate and observation availability on convective-scale data assimilation with an ensemble Kalman filter. Mon. Weather Rev. 132: 1238–1253.