

PARTIAL FILTERING IN STRONGLY NON-LINEAR ECOSYSTEM MODELLING

S. Losa, J. Schröter and M. Wenzel
Svetlana.Losa@awi.de



Abstract

We investigate advantages and specific difficulties in implementation of particle filters when assimilating data into strongly non-linear and non-Gaussian systems. With respect to the last circumstance, a Sequential Importance Resampling (SIR) filter seems to be of a great advantage since the filter updates probability of the particles (according to their agreement with the assimilated data) and thus, allows one to use the full forecast and data errors statistics. Still the most challenging thing in such a filtering is sampling the particles which approximate the continuous probability density function (pdf) evolved according to a stochastic dynamical model. We consider different sampling strategies which depend on our prior knowledge of the system and determine the filter performance cost and quality. All the experiments have been carried out with an ocean biogeochemical model which indeed is an example of a non-Gaussian system.

Data assimilation for stochastic systems

(Prediction)

$$\frac{dx}{dt} = M(x, y) + F(t) + \epsilon$$

Dynamical model
Predicting, one is uncertain in:
initial condition $x(0)$, model parameters p , external forcing F

$$\text{Forecast - a PDF - } p(x(t), p) \text{ defined on a } X_t \times P \text{ space, } x(t) \in X_t, p \in P$$

Global definition of PDF over the $X_t \times P$ space, here $X_t \cup P$ is a space of model trajectories:

$$p^*(x(t), p) = C \int_{X_0} \int_{P_0} p^*(x(t, x_0), p) \delta(x_0) \delta(p) dx_0 dp$$

$$p(x, p) = C \int_{X_0} \int_{P_0} p^*(x(t, x_0), p) \delta(x_0) \delta(p) dx_0 dp$$

M is a number of time step Δt
Improving the forecast with data assimilation globally – problem of Mem(X|dim(P) dimension)

(Analysis)

Improving the forecast with data assimilation globally – problem of Mem(X|dim(P) dimension)
Possible simplification:

1. Neglecting the system errors ϵ reduces the dimension of the problem by M .
2. Sequential approach

$$p^*(x(t), p | d_t) = C \delta(d_t - x(t)) p^*(x(t), p)$$

Evolving the forecast error statistics p^* 's still a problem of dim(X|dim(P) dimension)
Monte-Carlo methods is to avoid the dimension problem.

(The error of Monte Carlo approximation is $O(1/\sqrt{N})$, here N is a number of samples)

Ensemble Initialization

Spread of the initial ensemble reflects uncertainties in knowledge of a prior system and parameter pdf

$$\delta(x_0) \propto \prod_{k=1}^K \delta(x_0 - x_k(0)),$$

$$\delta(p) \propto \prod_{k=1}^K \delta(p - p_k),$$

δ - Dirac functions

Practically:

An ensemble of K members is generated from an exponential distribution

$$y = \text{JACOBI}(\frac{1}{\lambda})$$

mean of the distribution is assumed to be a first guess.

The idea – approximating the continuous probability density function (pdf) with an ensemble of δ -functions – particles, which evolves according to a stochastic dynamical model.

One of the particle filters. The Sequential Importance Resampling filter (SIR) has been first introduced by Rubin (1988), implemented for dynamical systems by Gordon et al. (1993); for sequential parameter estimation in stochastic systems, introduced by Kivman (2003); applied for basin scale problem by van Leeuwen (2003, 2004);

References

- Brasseur, P., Bahrel, P., Bertino, L., Birol, F., Brankart, J.-M., Ferry, N., Losa, S., Remy, E., Schroeter, J., Skachko, S., Testu, C.-E., Tranchant, B., van Leeuwen, P.-J., Verron, J., 2005. Data assimilation in operational ocean forecasting systems: the MERCATOR and MERSEA developments. *Quarterly Journal of the Royal Meteorological Society*, 131(613), 3561-3582.
- Kivman, G. A., 2003. Sequential parameter estimation for stochastic systems. *Nonlinear Process. Geophys.* 10, 253-256.
- van Leeuwen, P.J., 2003. A truly variance minimizing filter for large-scale applications. *Mon. Weather Rev.* 131, 2071-2084.
- Rubin D.B., 1988. Using the SIR algorithm to simulate posterior distribution in Bayesian Statistics 3 (Eds J.M. Bernardo, M.H. Degroot, D.V. Lindley and, A.F.M. Smith). Oxford Univ. Press., 395-402.

Data and weighting

The relative weights might be calculated under the assumption of Gaussian

$$w_i = \frac{\exp(-\frac{1}{2\sigma^2}(x_i - \bar{x})^2)}{\sum_{j=1}^N \exp(-\frac{1}{2\sigma^2}(x_j - \bar{x})^2)}$$

or Lorenz data errors (van Leeuwen, 2004)

$$w_i = C(1 + (x_i - \bar{x})^2 / \sigma^2)^{-\beta}$$

where σ^2 is the variance of the observation.

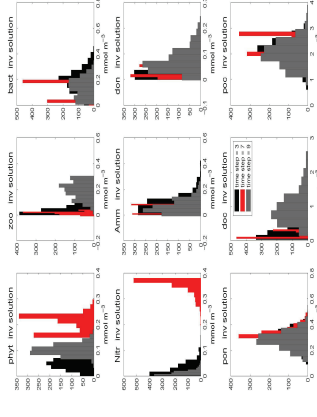
Resampling, practically:

IMSL library, DRNGDIA (Pseudo-CM) generates pseudorandom numbers from a general discrete distribution.

Advantages

When doing particle filtering, one updates probability of the particles according to their agreement with the observed data.

Particle filters make use of the full forecast and data errors statistics; they are truly variance minimizing methods (van Leeuwen 2003).

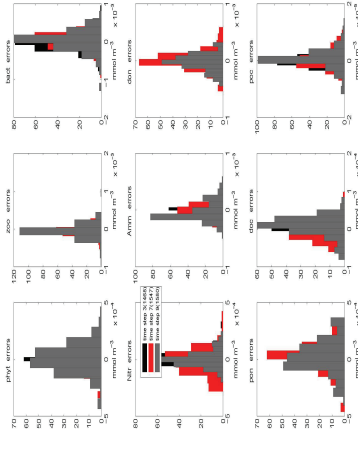


Problems

The SIR filter is known to suffer from degeneration of the ensemble (van Leeuwen, 2003) if either the system noise does not provide sufficient spreading of states which are resampled several times or the ensemble badly approximates the true prior distribution (the distance between the best member and the true state is too big). This problem is even more pronounced in the case of simultaneous state-parameter estimation where regenerating the number of samples in the parameter space is needed.

To avoid this the following procedure is used →

Figure below depicts estimated probability density of model errors, which, as one can see, have appeared to be non Gaussian. The more correct the model errors are accounted for, the better biological model parameters estimates and therefore the model data forecast are (Rabier et al. 2005).



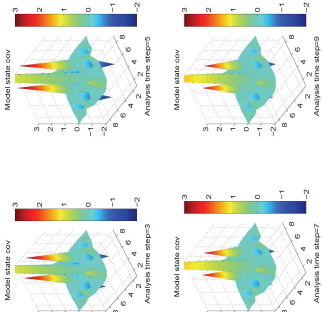
Example

We consider a 9-compartment model describing dynamics of phytoplankton (phy), zooplankton (zoo), bacteria (bact) and cycle of particular and organic matter (poc and dom), similar to what was used in Losa et al (2003). The model was constrained by data of the Bermuda Atlantic Time series Study, particularly, by measurements of nitrate, chlorophyll, dissolved organic nitrogen and carbon concentrations for the period December 1988 to January 1994. All the data were averaged over the ocean upper mixed layer (UML). The UML thickness were estimated by means of an analysis of BATS temperature profiles for the same period. The UML depth is determined as the depth at which the temperature is 50.5 °C less than that at the surface.

← Figure to the left depicts probability density obtained at different analysis time steps (days 1468, 1547, 1580) for 9 components of the ecosystem mode (ensemble size is equal to 1000).

Figure to the right → shows covariances between the model components calculated after 1468th day of the model integration.

← The pdfs of model states, parameters and variance of model noise obtained after 1547th day of model integration are then used to generate a new ensemble of 200 members to initialize the model (next data assimilation step)



Model noise generation and jittering model parameters

Levels E of the model noise ϵ might be considered as additional parameters to be optimized $E \in P$.

If, at an analysis step, parameter values are resampled (r many times), a new parameter ensemble can be redrawn (West, 1993) from a smoothed approximation of the posterior probability density

$$p^*(p(t_n)) = K^{-1} \sum_{k=1}^K \delta(p(t_n) - p_k^*(t_n))$$

Practically, either from

a uniform distribution within the interval $p \pm \sigma_p$, ... one has to specify β - nearest smaller value, $p +$ - nearest higher value)

or a normal distribution with a variance ... one has to specify;

