

Ensemble-based data assimilation method for count data

中野 慎也 モデリング研究系 教授

Variational data assimilation with count data

We denote the system state at time t_k as $\mathbf{x}_k \in \mathbb{R}^d$ and the observation obtained from \mathbf{x}_k as $\mathbf{y}_k \in \mathbb{R}^n$. The temporal evolution of the state given the numerical simulation model is written as

$$\mathbf{x}_k = \mathbf{f}_k(\mathbf{x}_{k-1}) \quad (1)$$

for $k = 1, \dots, K$.

We consider cases where each element of an observation vector \mathbf{y}_k follows a Poisson distribution and derive an ensemble-based variational method for such cases. The likelihood of \mathbf{x}_k is

$$P(\mathbf{y}_k | \mathbf{x}_k) = \prod_{j=1}^n \frac{h_{kj}(\mathbf{x}_k)^{y_{kj}}}{y_{kj}!} \exp[-h_{kj}(\mathbf{x}_k)], \quad (2)$$

where y_{kj} denotes the j -th element of the vector \mathbf{y}_k and the function h_{kj} provides a prediction of y_{kj} from \mathbf{x}_k .

We assume that all the elements of the state vector \mathbf{x}_k are non-negative and that the inequality $h_{kj}(\mathbf{x}_k) \geq 0$ is satisfied for each j . To ensure that each element $x_{0\ell}$ ($\ell = 1, \dots, d$) is positive, we represent $x_{0\ell}$ as follows

$$\xi_\ell = \log x_{0\ell}.$$

Combining ξ_ℓ for all ℓ into a vector ξ , the posterior distribution of ξ is

$$p(\xi | \mathbf{y}_{1:K}) \propto \frac{1}{\sqrt{(2\pi)^d |\mathbf{P}_\xi|}} \exp\left(-\frac{1}{2} [\xi - \bar{\xi}_b]^\top \mathbf{P}_\xi^{-1} [\xi - \bar{\xi}_b]\right) \times \prod_{k=1}^K \prod_{j=1}^n \frac{h_{kj}(\mathbf{x}_k)^{y_{kj}}}{y_{kj}!} \exp[-h_{kj}(\mathbf{x}_k)].$$

The log posterior density becomes

$$\log p(\xi | \mathbf{y}_{1:K}) = -\frac{1}{2} [\xi - \bar{\xi}_b]^\top \mathbf{P}_\xi^{-1} [\xi - \bar{\xi}_b] + \sum_{k=1}^K \sum_{j=1}^n [y_{kj} \log h_{kj}(\mathbf{x}_k) - \log y_{kj}! - h_{kj}(\mathbf{x}_k)] + \text{const.} \quad (3)$$

When we define the following composite function:

$$z_{kj}(\xi) = h_{kj} \circ \mathbf{f}_k \circ \mathbf{f}_{k-1} \circ \dots \circ \mathbf{f}_1 \circ \mathbf{e}(\xi) (= h_{kj}(\mathbf{x}_k)), \quad (4)$$

where \mathbf{e} denotes a d -dimensional function which satisfies

$$\mathbf{e}(\xi) = \begin{pmatrix} \exp(\xi_1) \\ \vdots \\ \exp(\xi_d) \end{pmatrix}.$$

To maximise $\log p(\xi | \mathbf{y}_{1:K})$, we minimise the following objective function:

$$J_P(\xi) = \frac{1}{2} [\xi - \bar{\xi}_b]^\top \mathbf{P}_\xi^{-1} [\xi - \bar{\xi}_b] + \sum_{k,j} [\log y_{kj}! + z_{kj}(\xi) - y_{kj} \log z_{kj}(\xi)]. \quad (5)$$

Approximation with an ensemble

We approximate the objective function J_P using a result of ensemble simulation. We generate N different values for ξ , $\{\xi_m^{(1)}, \dots, \xi_m^{(N)}\}$ to satisfy

$$\frac{1}{N} \sum_{i=1}^N \xi_m^{(i)} = \hat{\xi}_{m-1}. \quad (6)$$

where $\hat{\xi}_{m-1}$ is the $(m-1)$ -th estimate. We introduce a $d \times N$ matrix:

$$\Xi_m = \begin{pmatrix} \xi_m^{(1)} - \hat{\xi}_{m-1} & \dots & \xi_m^{(N)} - \hat{\xi}_{m-1} \end{pmatrix},$$

and represent ξ as follows:

$$\xi = \hat{\xi}_{m-1} + \Xi_m \mathbf{w}_m. \quad (7)$$

We obtain $z_{kj}(\xi_m^{(i)})$ for $i = 1, \dots, N$ for each time step k by running the ensemble simulation. Defining an N -dimensional row vector $\zeta_{kj,m}^\top$ as

$$\zeta_{kj,m}^\top = \left(z_{kj}(\xi_m^{(1)}) - z_{kj}(\hat{\xi}_{m-1}) \quad \dots \quad z_{kj}(\xi_m^{(N)}) - z_{kj}(\hat{\xi}_{m-1}) \right),$$

we can approximate $z_{kj}(\xi)$ as

$$z_{kj}(\xi) \approx z_{kj}(\bar{\xi}_b) + \zeta_{kj,m}^\top \mathbf{w}_m. \quad (8)$$

Ensemble-based variational method for count data

The objective function J_P in Eq. (5) can thus be approximated as

$$J_P(\xi) \approx \frac{1}{2} \left(\hat{\xi}_{m-1} - \bar{\xi}_b + \Xi_m \mathbf{w}_m \right)^\top \mathbf{P}_\xi^{-1} \left(\hat{\xi}_{m-1} - \bar{\xi}_b + \Xi_m \mathbf{w}_m \right) + \sum_{k,j} [\log y_{kj}! + z_{kj}(\bar{\xi}_b) + \zeta_{kj,m}^\top \mathbf{w}_m - y_{kj} \log (z_{kj}(\bar{\xi}_b) + \zeta_{kj,m}^\top \mathbf{w}_m)].$$

We minimise an objective function with a regularisation term:

$$\hat{J}_P(\xi) = \frac{\sigma_m^2}{2} \mathbf{w}_m^\top \mathbf{w}_m + \frac{1}{2} \left(\hat{\xi}_{m-1} - \bar{\xi}_b + \Xi_m \mathbf{w}_m \right)^\top \mathbf{P}_\xi^{-1} \left(\hat{\xi}_{m-1} - \bar{\xi}_b + \Xi_m \mathbf{w}_m \right) + \sum_{k,j} [\log y_{kj}! + z_{kj}(\bar{\xi}_b) + \zeta_{kj,m}^\top \mathbf{w}_m - y_{kj} \log (z_{kj}(\bar{\xi}_b) + \zeta_{kj,m}^\top \mathbf{w}_m)]. \quad (9)$$

At the m -th iteration, the optimal \mathbf{w}_m is obtained by minimising \hat{J}_P as follows:

$$\hat{\mathbf{w}}_m = - \left(\sigma_m^2 \mathbf{I}_N + \Xi_m^\top \mathbf{P}_\xi^{-1} \Xi_m + \sum_{k,j} \left[\frac{y_{kj}}{z_{kj}(\bar{\xi}_b)^2} \zeta_{kj,m} \zeta_{kj,m}^\top \right] \right)^{-1} \times \left(\Xi_m^\top \mathbf{P}_\xi^{-1} [\hat{\xi}_{m-1} - \bar{\xi}_b] + \sum_{k,j} \left[\zeta_{kj,m} - \frac{y_{kj}}{z_{kj}(\bar{\xi}_b)} \zeta_{kj,m} \right] \right). \quad (10)$$

The m -th estimate of ξ is obtained as

$$\hat{\xi}_m = \hat{\xi}_{m-1} + \Xi_m \hat{\mathbf{w}}_{m+1}, \quad (11)$$

and therefore the estimate of \mathbf{x}_0 at the m -th iteration becomes

$$\hat{\mathbf{x}}_{0,m} = \mathbf{e}(\hat{\xi}_m). \quad (12)$$

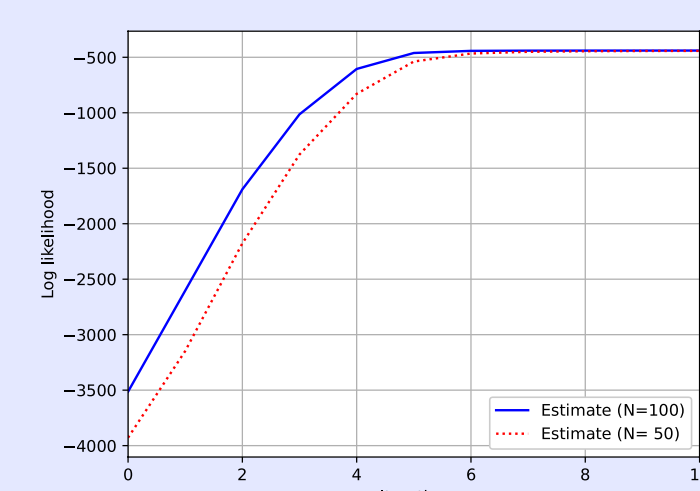
Experiment with a one-dimensional fluid model

We consider the one-dimensional advection and diffusion equation:

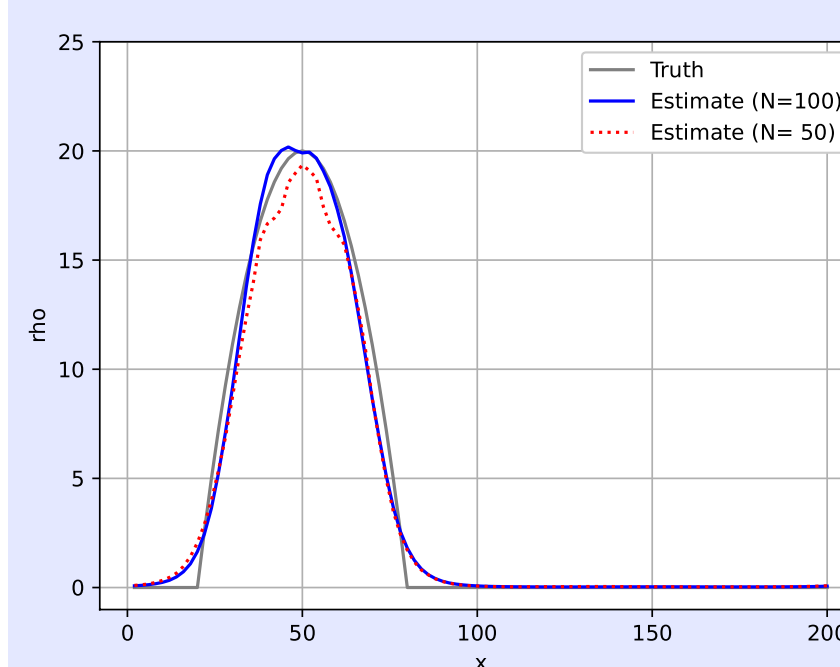
$$\frac{\partial \rho}{\partial t} = -u \frac{\partial \rho}{\partial x} + v \frac{\partial^2 \rho}{\partial x^2}, \quad (13)$$

where ρ is the fluid density, t indicates the time, x indicates the location of one-dimensional space, u is the velocity, and v is the diffusion coefficient.

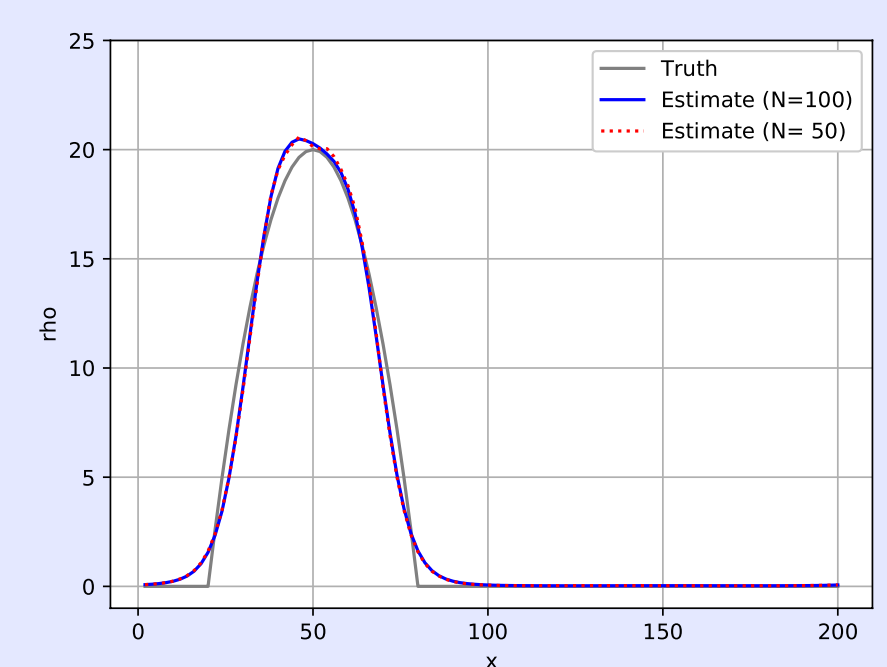
- ▶ The range of x is $0 \leq x < 200$ and a periodic boundary condition $\rho(t, x = 200) = \rho(t, x = 0)$ is imposed for any t . We also assume $u = 2$ and $v = 2$.
- ▶ The resolution is $\Delta x = 2$ in space and $\Delta t = 0.2$ in time.
- ▶ The Lax–Wendroff scheme is used for calculating the temporal evolution.
- ▶ The observations were obtained with a spatial resolution of $2\Delta x$ and a time interval of $20\Delta t$ until $t = 100\Delta t = 20$. Each of the data follows the Poisson distribution with a mean of $\rho(t, x)$.



The log likelihood of the estimate at each iteration. The solid line shows the result with the ensemble size $N = 100$ and the dotted line shows the result with $N = 50$.



The estimated after 6 iterations.



The estimated after 10 iterations.

The solid black line indicates the result with $N = 100$ and the dotted black line indicates the result with $N = 50$. The true initial state is shown with the solid gray line.

See Nakano (2003), Proc. Inst. Statist. Math., v. 70, p. 235 (in Japanese) for detail.