

Note on implicit filters

Florimond Guéniat,¹ Lionel Mathelin,² and Yussuf Hussaini¹

¹*Department of mathematics, Florida State Univ., 32306-4510 Tallahassee, FL, US*

²*LIMSI-CNRS, 91403 Orsay cedex, France*

(Dated: 29 July 2015)

Keywords: Data assimilation, Model identification

I. INTRODUCTION

In this article, the numerical cost and the accuracy of Implicit Particle Filter (IPF) are quantified, using Burgers's equation and shallow water equations as scenarios. Identifying the uncertainties in a data assimilation problem is specifically addressed.

IPF aims at estimating or the state, or the initial conditions or or some parameters of a system, by coupling a model and observations. A noise component can be present in both the observations and the model. Such an aim is of importance in many domains, most notably meteorology^{ref}, where initial conditions compatible with observations leads to accurate predictions. If not the case, as expected from a chaotic system, the accuracy of the forecast can drop significantly. Seismography^{ref} is an other field of importance for data assimilation, *e.g.*, to locate epicenter. It has been also successfully applied to nuclear fusion¹ and agronomy².

When the model is linear, and the noise Gaussian, the most suited class of methods are the Kalman filters. Uncertainties are directly related to the parameters of the Gaussian, which are totally determined by the Kalman filter. The generalization of Kalman filters, ensemble Kalman filters,¹ can deal with non-linear model by using replica (particles) of the system, to represent the prior density of solutions. Nevertheless, it performs poorly when the errors and noises are not Gaussian,². 4D-Var and cost function-minimization techniques are very popular methods, giving usually accurate results, but they fail at quantifying the uncertainties in the results, with respect to the noise in the model and observations. The sensitivity analysis cannot be considered as accurate, as soon as the noise cannot be considered as small [cite Plessix, review](#). Moreover, results are actually biased if the posterior pdf is multimodal,³.

Indeed, IPF have been successfully used for data assimilation,⁴, system identification,⁵, and parameter estimation. It is particularly relevant for multimodal pdf ; it provides an unbiased estimate of the solutions,⁶. As for Markov-Chain Monte Carlo method, IPF relies on a set of particles for the quantification of the conditional pdf of the system. The main difference lies in the fact that only a few, weighted hence carefully chosen particles are considered. Identifying these particles is done through the maximisation of a tailor-constructed probability. It ensures that each particle is associated with relevant informations of the pdf. The moments of the estimated pdf, as in standard MCMC methods, might be used for quantifying the uncertainties.

In this article, the precision and accuracy of IPF is qualitatively shown, and a roadmap for tackling the computational costs issue is proposed.

II. 4D-VAR

A. Preliminaries

Let's consider the state vector \mathbf{x} of a system of interest, evolving on the $n_{\mathcal{D}}$ manifold \mathcal{D} . This field is modeled through the set of equations f :

$$f(\mathbf{x}, \dot{\mathbf{x}}, q, t) = 0. \quad (1)$$

The model might not be perfect. Flaws can come, from instance, from uncertainties in the model parameters q , or on some initial conditions. Consequently, a functional J is introduced, in order to fit the model by taking in account

¹ Cacuci, D. G., & Ionescu-Bujor, M. (2010). Best-Estimate model calibration and prediction through experimental data assimilation-I: Mathematical framework. Nuclear science and engineering, 165(1), 18-44.

² Guerif, M., & Duke, C. L. (2000). Adjustment procedures of a crop model to the site specific characteristics of soil and crop using remote sensing data assimilation. Agriculture, ecosystems & environment, 81(1), 57-69.

observations at hand:

$$J(\mathbf{x}, q) = \int_0^T j(\mathbf{x}, q, t) dt. \quad (2)$$

For instance, one may want to minimize the discrepancy between the model predictions and some observations y :

$$j(\mathbf{x}, q, t) = |h(\mathbf{x}, q, t) - y(t)|^2$$

In the data assimilation context, the initial conditions are actually derived from the array of parameters $q \in \mathbb{R}^{n_q}$, with a relation $g : \mathcal{D} \times \mathbb{R}^{n_q} \rightarrow \mathbb{R}$:

$$g(\mathbf{x}(0), q) = 0.$$

Minimizing J gives the best initial conditions, in the sense that Eq. (1) will produce fields as compatible as possible with the observations.

Efficient and reliable methods for minimizing, such as the BFGS^{ref} algorithm, need the gradient of the functional with respect to the parameters. When $n_{\mathcal{D}}$ is large, approximating the gradient with finite differences become impractical. A powerful alternative to compute the gradient is to solve the adjoint equation^{ref}.

B. Computing the gradient of the cost functional

Identifying the minimum of the J relies on the gradient of the functional with respect to the parameters: $D_q J = \frac{D J}{D q}$. For that, one can introduce the Lagrangian \mathcal{L} , associated with the two lagrangian parameters λ and μ :

$$\mathcal{L}(\mathbf{x}, \dot{\mathbf{x}}, q, \lambda, \mu) = \int_0^T \left[j(\mathbf{x}, q, t) + \lambda^\dagger f(\mathbf{x}, \dot{\mathbf{x}}, q, t) \right] dt + \mu^\dagger g(\mathbf{x}(0), q) \quad (3)$$

where † is the transpose operator. Naturally, both \mathbf{x} and $\dot{\mathbf{x}}$ are considered as variables. As f and g are null by construction, λ and μ can be designed specifically to alleviate the computations.

The gradient $D_q \mathcal{L} = \frac{D \mathcal{L}}{D q}$ of \mathcal{L} is:

$$\begin{aligned} D_q \mathcal{L}(\mathbf{x}, \dot{\mathbf{x}}, q, \lambda, \mu) = & \int_0^T \left[\partial_{\mathbf{x}} j(\mathbf{x}, q, t) d_q \mathbf{x} + \partial_q j(\mathbf{x}, q, t) + \right. \\ & \left. \lambda^\dagger \partial_{\mathbf{x}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t) d_q \mathbf{x} + \lambda^\dagger \partial_{\dot{\mathbf{x}}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t) d_q \dot{\mathbf{x}} + \lambda^\dagger \partial_q f(\mathbf{x}, \dot{\mathbf{x}}, q, t) \right] dt + \\ & \mu^\dagger \partial_{\mathbf{x}(0)} g(\mathbf{x}(0), q) d_{\mathbf{x}(0)} + \mu^\dagger \partial_q g(\mathbf{x}(0), q). \end{aligned} \quad (4)$$

As a matter of facts, $D_q \mathcal{L} = D_q J$.

The term in $d_{\dot{\mathbf{x}}}$ cannot be easily estimated. An integration by parts gives:

$$\int_0^T \lambda^\dagger \partial_{\dot{\mathbf{x}}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t) d_q \dot{\mathbf{x}} dt = [\lambda^\dagger \partial_{\dot{\mathbf{x}}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t) d_q \mathbf{x}]_0^T - \int_0^T \left\{ \dot{\lambda}^\dagger \partial_{\dot{\mathbf{x}}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t) + \lambda^\dagger d_t \partial_{\dot{\mathbf{x}}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t) \right\} d_q \mathbf{x} dt$$

The term associated with $d_{\dot{\mathbf{x}}}$ can now be replaced in Eq. (4). Ordering terms in this equation leads to:

$$\begin{aligned} D_q \mathcal{L}(\mathbf{x}, \dot{\mathbf{x}}, q, \lambda, \mu) = & \int_0^T \left[\left(\partial_{\mathbf{x}} j(\mathbf{x}, q, t) + \lambda^\dagger \partial_{\mathbf{x}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t) - \left\{ \dot{\lambda}^\dagger \partial_{\dot{\mathbf{x}}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t) + \lambda^\dagger d_t \partial_{\dot{\mathbf{x}}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t) \right\} \right) d_q \mathbf{x} + \right. \\ & \left. \partial_q j(\mathbf{x}, q, t) + \lambda^\dagger \partial_q f(\mathbf{x}, \dot{\mathbf{x}}, q, t) \right] dt + \\ & \left\{ \mu^\dagger \partial_{\mathbf{x}} g(\mathbf{x}, q) - \lambda^\dagger \partial_{\dot{\mathbf{x}}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t) \right\} |_0 d_{\mathbf{x}(0)} + \mu^\dagger \partial_q g(\mathbf{x}(0), q) + \left\{ \lambda^\dagger \partial_{\dot{\mathbf{x}}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t) \right\} |_T d_q \mathbf{x}(T) \end{aligned} \quad (5)$$

Proper choices for λ and μ allow to simplify the expression of $D_q \mathcal{L}$. The choice of $\lambda(T) = 0$ actually nullifies the term $\{\lambda^\dagger \partial_{\dot{\mathbf{x}}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t)\} |_T d_q \mathbf{x}(T)$. Then, λ can be chosen as the solution of the so-called adjoint equation:

$$\left(\partial_{\mathbf{x}} j(\mathbf{x}, q, t) + \lambda^\dagger \partial_{\mathbf{x}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t) - \left\{ \dot{\lambda}^\dagger \partial_{\dot{\mathbf{x}}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t) + \lambda^\dagger d_t \partial_{\dot{\mathbf{x}}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t) \right\} \right) = 0$$

integrated in backward time. Finally, the Lagrange parameter μ is set so it nullifies the component associated with $dq_{\mathbf{x}(0)}$:

$$\{\mu^\dagger \partial_{\mathbf{x}} g(\mathbf{x}, q) - \lambda^\dagger \partial_{\dot{\mathbf{x}}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t)\} |_0 = 0.$$

Then, computing $D_q \mathcal{L}$, hence $D_q J$, is achieved by the integration of:

$$D_q J(\mathbf{x}, q) = \int_0^T \left[\partial_q j(\mathbf{x}, q, t) + \lambda^\dagger \partial_q f(\mathbf{x}, \dot{\mathbf{x}}, q, t) \right] dt + \mu^\dagger \partial_q g(\mathbf{x}(0), q). \quad (6)$$

III. SYSTEM

A. Burgers' equation

The Burgers' equation **stochastic part ?** is:

$$\partial_t \mathbf{x} + \frac{1}{2} \partial_z \mathbf{x}^2 - \nu \partial_{zz}^2 \mathbf{x} = 0.$$

\mathbf{x} evolves on an $n_{\mathcal{D}}$ -dimensional manifold \mathcal{D} . We consider that it has been spatially discretized. The field \mathbf{x} can hence be modeled by $\hat{\mathbf{x}}$ through a set of ordinary differential equations $f(\hat{\mathbf{x}}, \dot{\hat{\mathbf{x}}}, q, t) = 0$.

$$\dot{\hat{\mathbf{x}}} + L(\hat{\mathbf{x}}, q, t) + NL(\hat{\mathbf{x}}, q, t) = 0 \equiv f(\hat{\mathbf{x}}, \dot{\hat{\mathbf{x}}}, q, t),$$

where $q \in \mathbb{R}^q$ are the parameters of the model. The solution is discretized with n_x points linearly spaced: $\hat{\mathbf{x}} \in \mathcal{R}^{n_x}$, following an implicit scheme, in order to be solved by Newton iterations. The $\hat{}$ symbol is dropped in the following for the sake of readability.

$$\frac{\mathbf{x}_i^{n+1} - \mathbf{x}_i^n}{\Delta t} + \mathbf{x}_i^{n+1} \frac{\mathbf{x}_i^{n+1} - \mathbf{x}_{i-1}^{n+1}}{\Delta z} - \nu \frac{\mathbf{x}_{i+1}^{n+1} - 2\mathbf{x}_i^{n+1} + \mathbf{x}_{i-1}^{n+1}}{\Delta z^2} = 0 \quad (7)$$

More information can be found in Sec. A

B. Shallow water equations

We consider the SWE equations (in its Ito form ?):

$$\begin{aligned} \partial_t u &= -u \partial_x u - v \partial_y u - f v - g \partial_x h - b u \\ \partial_t v &= -u \partial_x v - v \partial_y v + f u - g \partial_y h - b v \\ \partial_t h &= -\partial_x ((H + h) u) - \partial_y ((H + h) v), \end{aligned} \quad (8)$$

Ito ? 0 where u and v are the surface velocity, h (resp. H) is the deviation of height (resp. mean height) from the bottom. f is the Coriolis coefficient and ν is the viscous drag coefficient. g is the acceleration due to gravity. This equation is spatially discretized.

IV. PARTICLE METHODS AND IMPORTANCE SAMPLING

A. Preliminaries

Consider a stochastic dynamical system, in the Ito formalism, evolving on an $n_{\mathcal{D}}$ -dimensional manifold \mathcal{D} :

$$d\mathbf{x}(t) = \mathbf{f}(\mathbf{x}(t)) dt + \mathbf{G} dW, \quad \mathbf{x} \in \mathcal{D}, \quad (9)$$

with x the state of the system and $\mathfrak{f} : \mathcal{D} \rightarrow \mathcal{D}$ the flow operator. dW is a an $n_{\mathcal{D}}$ -dimensional Brownian motion (Wiener process) and \mathbf{G} a diffusion $n_{\mathcal{D}} \times n_{\mathcal{D}}$ matrix. With a proper time scheme³ associated with a time step δt , Eq. (9) becomes:

$$\mathbf{x}(t_0 + (n+1)\delta t) = \mathfrak{F}(\mathbf{x}(t_0 + n\delta t)) + \mathbf{G}\sqrt{\delta t}E, \quad (10)$$

where E is an $n_{\mathcal{D}}$ -dimensional realization of the Wiener process. This model is to be understand as space and time discretization of a set of partial differential equations.

Let $h : \mathcal{D} \rightarrow \mathbb{R}^{n_p}$ be a sensor function, so an observation on the system is:

$$y(t) = h(x, t) + \sqrt{\mathbf{S}}Y, \quad (11)$$

Y is an n_p -dimensional, standard-normal random variable, and $\sqrt{\mathbf{S}}$ a covariance matrix. Let also consider, without any loss of generality, that $\mathbf{G} = g\mathbf{Id}$ and $\sqrt{\mathbf{S}} = \sqrt{s}\mathbf{Id}$. In the following, the time will be discretized accordingly to δt . For any variable \star , $\star(t_0 + n\delta t)$ will be shortened as \star_n .

B. Probabilities derived from the Ito formulation

The fact that noises in Eqs. (10) and (11) are Gaussian does not mean that the probability density function (pdf) associated with the state of the system is Gaussian, as it will be shown in the following.

The probability that the $n_{\mathcal{D}}$ -dimensional random variable X^{n+1} coincide with a solution of the system described in Eq. (9), at time $n+1$ and with the state \mathbf{x}_n as initial condition, is given by:

$$P_{int}^{\mathbf{x}_n, n+1}(X^{n+1}) = P(X^{n+1} | \mathbf{x}_n).$$

This probability only depends on the system equations, *i.e.* Eqs. (9) and (10). When dealing with a Wiener process, the conditional probability is:

$$P_{int}^{\mathbf{x}_n, n+1}(X^{n+1}) \propto \exp\left(-\|X^{n+1} - \mathfrak{F}(\mathbf{x}_n)\|_2^2 / 2\delta t g^2\right). \quad (12)$$

Note that if the noise is not Gaussian (*i.e.* not described by a Wiener process), it only affects the right side of Eq. (12).

On an other hand, the probability that the observation of the variable X^{n+1} is compatible with the observation of the system is given by:

$$P_{obs}^{n+1}(X^{n+1}) = P(h(X^{n+1}) | y_{n+1}).$$

Again, this probability only depends on the sensor function, see Eq. (11):

$$P_{obs}^{n+1}(X^{n+1}) \propto \exp\left(-\|h(X^{n+1}) - y_{n+1}\|_2^2 / 2s\right). \quad (13)$$

From the Bayes' theorem, by combining Eq. (12) and Eq. (13), the probability $P_{SI}^{\mathbf{x}_n, n+1}$ that X^{n+1} is a possible outcome from the previous state \mathbf{x}_n , given the observation y_{n+1} , is:

$$P_{SI}^{\mathbf{x}_n, n+1}(X^{n+1}) \propto P_{int}^{\mathbf{x}_n, n+1}(X^{n+1}) P_{obs}^{n+1}(X^{n+1}) \quad (14)$$

From this writing, as a product, the pdf of $P_{SI}^{\mathbf{x}_n, n+1}$ is not normal.

Generalization is straightforward. Having several observations $\{y_{n+i}\}_{i \in \{1, \dots, r\}}$ available at hand, between time t_n and t_{n+r} , is common. Consequently, Eq. (14) becomes:

$$P_{SI}^{\mathbf{x}_n, n+r}(\mathbf{X}_{n+1}^{n+r}) \propto \prod_{i=1}^r P_{int}^{X^{n+i-1}, n+i}(X^{n+i}) \prod_{i \in \{1, \dots, r\}} P_{obs}^{n+i}(X^{n+i}), \quad (15)$$

where $\mathbf{X}_{n+1}^{n+r} = \{X^{n+1}, \dots, X^{n+r}\}$ is now a $r \times n_{\mathcal{D}}$ -dimensional random variable, coinciding with a potential trajectory of the state between $n+1$ and $n+r$. and with $X^n = \mathbf{x}_n$.

³ Kloeden, P. E., Platen, E. (1999). Numerical solution of stochastic differential equations. Berlin: Springer, Klauder Petersen scheme

The probability P_{DA}^n that the system is in a given state X^n at the given time n can be derived from Eq. (15):

$$P_{DA}^n(\mathbf{X}_{n+1}^{n+r}) \propto \prod_{i=1}^r P_{int}^{X^{n+i-1}, n+i}(X^{n+i}) \prod_{i \in \{1, \dots, r\}} P_{obs}^{n_i}(X^{n+i}), \quad (16)$$

with $\mathbf{X}_{n+1}^{n+r} = \{X^n, \dots, X^{n+r}\}$.

Data assimilation, as system identification, aims at identifying the state of the considered system, at any desired time. Consequently, $\chi \equiv \text{argmax}(P_{DA}^n)$ (resp. $\chi \equiv \text{argmax} P_{SI}^{\mathbf{x}_n, n+r}$) corresponds to the best solution of the data assimilation (resp. system identification) problem,⁴⁻⁶. The form of Eqs. (12) and (13) implies that minimizing $F_{SI} = -\log P_{SI}^{\mathbf{x}_n, n+r}$ and $F_{DA} = -\log P_{DA}^n$ is numerically preferable.

In the following, the notation P (resp. F) will refer indistinctly to $P_{SI}^{\mathbf{x}_n, n+r}$ or P_{DA}^n (resp. F_{SI} or F_{DA}), when no confusion is possible. Similarly, X will refer to the considered random variable, and its dimension will be short-named $n_{\mathcal{D}}$.

Identifying the maximum can be, alternatively, efficiently done by 4D-VAR, [refs](#), when the distribution associated with P_{DA} is known to be mono-modal⁴.

For state identification, as for data assimilation, the initial conditions \mathbf{x}_n in Eq. (15) are often unknown. This problem is usually tackled by considering m particles, *i.e.* replicas of the system⁵. For each particle i , an initial condition \mathbf{x}_n^i is guessed. Maxima are computed from Eq. (15), for all the initial guesses.

By the law of large numbers, it allows to estimate the moments of the pdf, that is, for system identification, the mean state \mathbf{x}_n , by averaging the results.

Estimating the moments of the a given pdf might need a consequent number of particles [refs](#), and both the storage and the identification of samples from that pdf might be near to impossible. Actually, when strictly following the law of large numbers, all the particles has the same importance. An option is, instead of drawing random particles, to identify particles more representative of the probability density function (pdf) associated with P . Hence, less particles are needed for an accurate estimation of the pdf properties.

C. Importance sampling

Intuitively, a particle carries more information on the probability, as soon as it is associated with a high probability in Eq. (14) (resp. (15) or (16)). Moreover, estimations of the pdf properties rely on drawing samples from the pdf. A sample has to be a possible outcome of the physical system described by Eq. (9) for being meaningful.

Importance sampling,^{4,5}, [ref is](#) allows to tackle both these issues. The key is to consider an easy-to-sample pdf p . Some samples X_i are drawn from it. These samples are mapped with the the actually hard-to-sample pdf P . It also provides a quantification of the importance of each sample with respect to the pdf, namely the weight of each particle, so the estimation of the moments are not biased.

The weights of X_i are:

$$w(X_i) = \frac{P(X_i)}{p(X_i)}. \quad (17)$$

For being efficient, it is necessary that the maxima of p and P coincide. If not, some particles will be associated with low weights, and are hence of a low importance. Because of potential storage limitations and computational costs, such particles are to be avoided, and the objective is to identify particles with comparable weights.

In order to construct a p density, let consider ξ , an intermediate random variable. ξ is drawn from a $n_{\mathcal{D}}$ -dimensional Gaussian pdf. Such a variable is, of course, far from being a possible outcome of the physical system. It will be used as an intermediate for identifying a high probability sample X , solution of the algebraic equation,⁵:

$$F(X) - \Phi_F = \xi^t \xi / 2, \quad (18)$$

where $\Phi_F = F(\chi)$ is the given minimum of F , identified for instance with a 4D-VAR method.

The sample X has then an high probability to remains close to χ , as ξ is drawn for a 0 mean pdf.

⁴ if note, the estimation might be biased. Weir 2013

⁵ *e.g.*, in Markov Chain Monte Carlo simulations

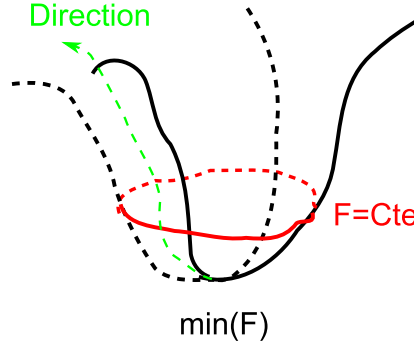


FIG. 1. Illustration of the methodology for drawing samples. The red curve corresponds to iso-values of $F(X) = F_{min} + 1/2\xi^T \xi$. The green curve illustrates the random direction associated with ξ . The sample is found at the intersection of these two curves.

In order to identify the weight from Eq. (17), one has to compute the probability $p_X(X)$ of drawing X from Eq. (18). By the substitution rule, p_X is related to the probability $P_\xi(\xi)$ of drawing the sample ξ by:

$$p_X(X) \left| \det \frac{\partial X}{\partial \xi} \right| = P_\xi(\xi).$$

Consequently, with Eq. (17), the weight associated to the sample X is:

$$w(X) = \frac{P(X)}{P_\xi(\xi) / \left| \det \frac{\partial X}{\partial \xi} \right|}. \quad (19)$$

D. Random map

Solving Eq. (18) is a necessary step. For that, Morzfeld *et al.*⁵ proposed to use the following ansatz:

$$X = X(\lambda) = \chi + \lambda L^t \xi / \sqrt{\xi^t \xi}, \quad (20)$$

where λ is a real parameter, and L^t is the Choleski decomposition of H^{-1} , H being the Hessian of F , evaluated in χ . It minimizes the dispersion of the jacobian $\left| \det \frac{\partial X}{\partial \xi} \right|$ when considering different particles, and hence the dispersion of weights,⁵.

A BFGS-algorithm is then used for identifying the λ minimizing Eq. (18).

Computing the weight needs the evaluation of the jacobian of the maps between the reference sample X and the random variable ξ . Following⁵, let be $\rho = \sqrt{\xi^t \xi}$.

Differentiating Eq. (20) leads to:

$$\begin{aligned} \frac{\partial X}{\partial \xi} &= \frac{\partial}{\partial \xi} \lambda L^t \xi / \sqrt{\xi^t \xi} \\ &= L^t \left(\frac{\partial \lambda}{\partial \xi} \xi / \sqrt{\xi^t \xi} + \lambda \frac{\partial}{\partial \xi} \xi / \sqrt{\xi^t \xi} \right) \end{aligned} \quad (21)$$

V. RESULTS

A. Burgers' equation

To illustrate the methodology discussed above, we now consider the Burgers' equation presented in Sec. III A. The objective is the identification of the initial conditions, knowing only some observations. The accuracy and the computational costs are estimated. Three main parameters are allowed to vary: the noises in the model g , the noises in the observations s , and finally the dimensions (spatial ? temp?).

1. accuracy

- noise in model
- noise in observation
- dimensions

2. Computational costs

- noise in model
- noise in observation
- dimensions

BIBLIOGRAPHY

¹G. Evensen, *Data assimilation: the ensemble Kalman filter*. Springer Science & Business Media, 2009.

²R. Miller, E. Carter, and S. Blue, “Data assimilation into nonlinear stochastic models,” *Tellus A*, vol. 51, no. 2, pp. 167–194, 1999.

³F. Rabier and P. Courtier, “Four-dimensional assimilation in the presence of baroclinic instability,” *Quarterly Journal of the Royal Meteorological Society*, vol. 118, no. 506, pp. 649–672, 1992.

⁴A. Chorin, M. Morzfeld, and X. Tu, “Implicit particle filters for data assimilation,” *Communications in Applied Mathematics and Computational Science*, vol. 5, no. 2, pp. 221–240, 2010.

⁵M. Morzfeld, X. Tu, E. Atkins, and A. Chorin, “A random map implementation of implicit filters,” *Journal of Computational Physics*, vol. 231, no. 4, pp. 2049–2066, 2012.

⁶E. Atkins, M. Morzfeld, and A. Chorin, “Implicit particle methods and their connection with variational data assimilation,” *Monthly Weather Review*, vol. 141, no. 6, 2013.

⁷E. Lorenz, “Deterministic nonperiodic flow,” *J. Atmos. Sci.*, vol. 20, no. 2, pp. 130–141, 1963.

⁸T. El Moselhy and Y. Marzouk, “Bayesian inference with optimal maps,” *Journal of Computational Physics*, vol. 231, no. 23, pp. 7815–7850, 2012.

Appendix A: Burgers’ equation and adjoint equations

We recall here the expression from Eq. (7):

$$\frac{\mathbf{x}_i^{n+1} - \mathbf{x}_i^n}{\Delta t} + \mathbf{x}_i^{n+1} \frac{\mathbf{x}_i^{n+1} - \mathbf{x}_{i-1}^{n+1}}{\Delta z} - \nu \frac{\mathbf{x}_{i+1}^{n+1} - 2\mathbf{x}_i^{n+1} + \mathbf{x}_{i-1}^{n+1}}{\Delta z^2} = 0$$

The i th component of f is then:

$$f_i(\mathbf{x}, \dot{\mathbf{x}}, q, t) = \dot{\mathbf{x}}_i + \mathbf{x}_i \frac{\mathbf{x}_i - \mathbf{x}_{i-1}}{\Delta z} - \nu \frac{\mathbf{x}_{i+1} - 2\mathbf{x}_i + \mathbf{x}_{i-1}}{\Delta z^2}.$$

Consequently, the partial derivatives with respect to \mathbf{x} are:

$$\partial_{\mathbf{x}_j} f_i(\mathbf{x}, \dot{\mathbf{x}}, q, t) = \frac{1}{\Delta z} (2\mathbf{x}_i \delta_{i,j} - \delta_{i-1,j}) + \frac{1}{\Delta z^2} (\delta_{i+1,j} - 2\delta_{i,j} + \delta_{i-1,j}).$$

The partial derivative with respect to $\dot{\mathbf{x}}$ is:

$$\partial_{\dot{\mathbf{x}}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t) = Id,$$

and then

$$dt \partial_{\dot{\mathbf{x}}} f(\mathbf{x}, \dot{\mathbf{x}}, q, t) = 0.$$

The cost function is associated with:

$$j(\mathbf{x}, q, t) = \|\mathbf{x} - \mathbf{y}(t)\|^2.$$

Hence, the partial derivatives with respect to q and \mathbf{x} are:

$$\partial_q j(\mathbf{x}, q, t) = 0,$$

$$\partial_{\mathbf{x}} j(\mathbf{x}, q, t) = 2 * (\mathbf{x} - y(t)).$$

The parameters q are actually the initial conditions. Hence, n_q is equal to $n_{\mathcal{D}}$. As a matter of fact, g is:

$$g(\mathbf{x}(0), q) = \mathbf{x}(0) - q = 0.$$

Then, the partial derivatives of g follow:

$$\partial_q g(\mathbf{x}(0), q) = -Id,$$

and

$$\partial_u g(\mathbf{x}, q)|_0 = Id.$$