

Preface

Data assimilation is the science of combining information from prior knowledge in the form of a numerical dynamical model with new knowledge in the form of observations to obtain a best description of the system at hand. It is used to predict the future of the system, infer best parameter values, and to evaluate and compare different models. This ‘best’ description needs to contain information about the uncertainty, and the most general form is in terms of a probability distribution over the space of all possible model states. The basic mathematical formulation of the data assimilation problem is based on Bayes theorem, which states that this best probability distribution, called the posterior, is a point wise multiplication of the probability distribution of our prior knowledge from the numerical model with the probability distribution of the observations given each possible state of the model. The method is applied in almost all branches of science, although often under different names. Indeed, inverse modelling can be seen as a specific branch of data assimilation (or the other way around), as long as the model is dynamic in nature. Data assimilation can also often be formulated in terms of filtering and smoothing problems for stochastic processes.

An application of great practical relevance can be found in numerical weather forecasting, where atmospheric models and observations are combined every 6 to 12 hours to provide the best starting point for future forecasts. Other important applications can be found in all branches of the geosciences, such as oceanography, atmospheric pollution, marine biogeochemistry, ozone, seasonal forecasting, climate forecasting, sea-ice, glaciers and ice caps, ecology, land surface, etc. It is also used in oil reservoir modelling and seismology, and is typically referred to as history matching in those fields. Industrial applications are also widespread; think about all processes that need automatic control. Medical applications are growing too; data assimilation constitutes, for instance, an emerging field in the neurosciences.

Weather forecasting has been the driving force behind many recent theoretical and practical advances in data assimilation algorithms. The reason for this is twofold: the dimension of the system is huge, typically a billion nowadays, and the turnaround time is very short, the actual data assimilation can only take up to one hour, the rest of the 6- or 12- hour cycle is used to collect the billions of

observations, perform quality control, which means throwing away close to 95% of the data, and preparing the observation-model interface. Currently used data-assimilation methods, which fulfil these operational constraints, can be divided into two categories: sequential methods and variational methods. This division is somewhat arbitrary, as will become clear shortly. The sequential methods are based on the Kalman filter. In the Kalman filter, the assumption is made that the probability distributions involved are all Gaussian. (In fact, the Kalman filter can also be derived assuming a linear update of the system, but that description falls outside the Bayesian framework.) The advantage of this approach is that only the first two moments of the distributions are needed. However, the size of the system in numerical weather prediction is too large to use the Kalman filter directly, simply because the second moments, the covariance matrix, need a billion squared entries. We have no supercomputer that can store that amount of numbers at present. Also, propagation of the covariance matrix under the model equations is prohibitively expensive. Perhaps, a bigger problem is that the Kalman filter is only justified for linear models. This limitation motivated the development of the ensemble Kalman filters starting in the 1990s in which the probability distribution is represented by a finite set of model states that are propagated with the full nonlinear model equations in between observations. Only at observation times, the ensemble of model states is assumed to represent a Gaussian and the Kalman filter update is implemented directly on the ensemble of model states. The Gaussian approximation can be justified from maximum entropy considerations given that only the ensemble mean and covariance matrix can be estimated from the available data. Furthermore, quite sophisticated methods have been developed to ensure efficient implementation, e.g. to avoid having to compute or store the full covariance matrix at any point in the algorithm. The finite ensemble size, typically 10–100 members can be afforded, leads to rank deficient matrices, and methods like localisation and inflation are used to counter this problem. These are to a large extent ad hoc, and this is a very active area of research.

The variational methods search for the mode of the posterior distribution. This can be the marginal posterior distribution at the time of the observations, leading to 3DVar, or the joint-in-time posterior distribution over a time window, in which case the method is called 4DVar. Again very sophisticated numerical techniques have been developed to solve this optimisation problem in billion-dimensional spaces. Unfortunately, these methods rely on linearisations and Gaussian assumptions on the prior and observation errors. Furthermore, the methods do not provide an uncertainty estimate, or only at very large cost.

Although numerical weather forecasting is quite successful, the introduction of convection resolving models and more complex observation networks leads to new challenges and the present-day methods will struggle. In particular, there is a strong need to move away from Gaussian data assimilation methods towards non-parametric methods, which are applicable to strongly nonlinear problems. (In numerical weather prediction, the so-called hybrids are becoming popular, combining ensemble Kalman filter and variational ideas, but this does not necessarily make the methods applicable to more nonlinear problems.) Fully non-parametric methods

for probability distributions of arbitrary shape do exist and are based on sequential Monte Carlo methods, in which ensembles of model states, called samples, are generated to represent the posterior distribution. While these methods are extremely useful for small dimensional systems, they quickly suffer from the so-called curse of dimensionality, in which it is very unlikely for these states to end up in the high-probability areas of the posterior distribution.

Two solutions have been suggested to solve the curse-of-dimensionality problem. The first one is based on exploring the proposal density freedom in Monte Carlo methods. Instead of drawing samples from the prior distribution, one can draw samples from a proposed distribution and either accept them with a certain probability related to this proposal or change their weights relative to the other samples. This proposal density is chosen such that the samples will be from the high-probability area of the posterior distribution by construction.

Another option is to try to reduce the size of the problem by the so-called localisation. This reduces the influence of an observation to its direct neighbourhood, so that the actual data assimilation problem for each observation is of much smaller dimension. It is a standard method in ensemble Kalman filtering for high-dimensional systems, and it is key to their success.

This volume of *Frontiers in Applied Dynamical Systems* focuses on these two potential solutions to the nonlinear data assimilation problem for high-dimensional systems. Both contributions start from particle filters. A particle filter is a sequential Monte Carlo method in which the samples are called particles. It is a fully non-parametric method and applicable to strongly nonlinear systems. Particle filters have already found widespread applications ranging from speech recognition to robotics to, recently, the geosciences. The contribution of van Leeuwen focuses on the potential of proposal densities for efficiently implementing particle filters. It discusses the issues with present-day particle filters and explores new ideas for proposal densities to resolve them. A particle filter that works well in systems of any dimension is proposed and implemented for a high-dimensional example.

The contribution by Cheng and Reich discusses a unified framework for ensemble transform particle filters. This allows one to bridge successful ensemble Kalman filters with fully nonlinear particle filters and allows for a proper introduction of localisation in particle filters, which has been lacking up to now.

While both approaches introduce tuneable parameters into particle filters (such as the localisation radius), the proposed methods are capable of capturing strongly non-Gaussian behaviour in high dimensions. Both approaches are quite general and can be explored further in many different directions, making them both potential candidates for solving the full problem (for instance by combining them). We hope that they will form the basis of many new exciting ideas that push this field forward. As mentioned above, the number of application areas is huge, so high impact is guaranteed.

Reading, United Kingdom
Potsdam, Berlin, Germany
Potsdam, Berlin, Germany

Peter Jan van Leeuwen
Yuan Cheng
Sebastian Reich

Nonlinear Data Assimilation

Van Leeuwen, P.J.; Cheng, Y.; Reich, S.

2015, XII, 118 p. 19 illus., 15 illus. in color., Softcover

ISBN: 978-3-319-18346-6