

Aim of lecture

- Give an overview of the different approaches to solving the data assimilation problem by building on a common framework.
- Introduce the frequently used terminology
- Highlight the similarities and differences between the different approaches.
- More detailed descriptions of the algorithms will be presented in the rest of the week.

Recap aims of DA

The DA problem is to combine **prior knowledge** and **relevant observations** to give an **updated estimate of the state** of the atmosphere/oceans/land surface etc.



- To initialize a forecast the better the initial conditions the better the forecast.
- Create **reanalyses** to understand the recent past.
- Estimate parameters in the model to give a better understanding of the processes represented.

Linking observations to models

Horizontal Grid (Latitude-Longitude) Vertical Grid (Height or Pressure) Physical Processes in a Model T

 $\mathbf{x} \in \mathbb{R}^{n}$ is a vector of the model variables that we want to estimate. Referred to as the *STATE*.

> As **x** and **y** lie in different spaces it is necessary to define an OBSERVATION OPERATOR, $h: \mathbb{R}^n \to \mathbb{R}^p$, that maps **x** from state space to observation space.





The inverse problem

The observation operator, which maps from state to observation space, is often known as the *FORWARD MODEL*.

Reconstructing **x** from **y** is often referred to as the *INVERSE PROBLEM*.

If the observation operator were linear, we could write $\mathbf{y} = \mathbf{H}\mathbf{x}$, where $\mathbf{H} \in \mathbb{R}^{p \times n}$.

Usually **H** is non-square ($n \neq p$) and/or rank deficient. Therefore, we cannot invert **H** to find $\mathbf{x} = \mathbf{H}^{-1} \mathbf{y}$ directly.

Additionally, the observations are not perfect and contain errors.

Data assimilation provides a framework for solving the inverse problem by introducing prior information about \mathbf{x} .

Bayes' theorem

Most DA algorithms can be derived from Bayes' theorem:





y is a vector of the available observations

Bayes' theorem

Most DA algorithms can be derived from Bayes' theorem:

 $p(\mathbf{x}|\mathbf{y}) \propto p(\mathbf{x})p(\mathbf{y}|\mathbf{x})$

posterior prior likelihood

y is a vector of the available observations



Bayes' theorem: Scalar illustration $p(\mathbf{x}|\mathbf{y}) \propto p(\mathbf{x})p(\mathbf{y}|\mathbf{x})$



 The prior PDF, p(x), describes the probability of your state variables. Often, this knowledge comes from a previous forecast.

Bayes' theorem: Scalar illustration $p(\mathbf{x}|\mathbf{y}) \propto p(\mathbf{x})p(\mathbf{y}|\mathbf{x})$



 Likelihood PDF, p(y|x), describes the probability of observations given that they are measuring the state we are interested in. P(y|x)=L(x|y) so that we can think of it as a function of x.

Bayes' theorem: Scalar illustration $p(\mathbf{x}|\mathbf{y}) \propto p(\mathbf{x})p(\mathbf{y}|\mathbf{x})$



 The posterior PDF is given by multiplying the two together and normalising. Updating the prior with information from the observations has shifted the probability mass and reduced the range of probable values of x (i.e. the uncertainty in x is reduced!)

Bayes' theorem: 2 variable example

Let our state **x** be a vector of zonal and meridional winds, u and v at one location. Observe u only.



Compared to the Prior, the region of high probability in the Posterior for *u* is reduced and is shifted towards that of the likelihood.

Bayes' theorem: 2 variable example

How do you think the Posterior would change if the **Prior** was **correlated**?



When the Prior is correlated, **observations of one variable can be used to update the analysis of both variables**. This helps to **ensure the analysis is physically realistic**.

Bayes' theorem: 2 variable example

In this example the prior is non-Gaussian. Let us observe u^2 , how will this change the likelihood?



How to solve Bayes' theorem?

A naïve approach could be to discretize the whole of the state space, evaluate $p(x = x_i)$ and $p(y|x = x_i)$ and multiply (this is exactly what I did in the 2D examples)



Might miss the regions of high probability
 Waste effort evaluating the PDFs in regions of near-zero probability

We will discuss two different, more efficient, approaches to solving Bayes' theorem

- 1. Non-Parametric approaches: No assumptions are made about the form of Prior or likelihood and consequently the Posterior. These methods include MCMC and the Particle filter
- 2. Parametric approaches: The prior, likelihood and hence the Posterior are assumed to follow a given distribution. These methods include variational techniques and the EnKF.

Non-parametric approaches

NERC/NCEO/DARC Training course on data assimilation and its interface with machine learning 2025

MCMC

Markov Chain Monte Carlo (MCMC) are a class of algorithms that allow you to sample from the Posterior distribution.

The basic algorithm is:

- 1. Initialise, i = 0, select \mathbf{x}_i .
- 2. Draw a new sample from a proposal distribution $\tilde{\mathbf{x}} \sim q[\tilde{\mathbf{x}}; \mathbf{x}_i]$.
- 3. Evaluate $\tilde{\mathbf{y}} = h(\tilde{\mathbf{x}})$ and $P(\tilde{\mathbf{x}})$
- 4. Decide whether to keep $\tilde{\mathbf{x}}$ as part of the sample. Compute $\rho = \frac{P(\tilde{\mathbf{x}})P(\mathbf{y}|\tilde{\mathbf{x}})q[\tilde{\mathbf{x}};\mathbf{x}_i]}{P(\mathbf{x}_i)P(\mathbf{y}|\mathbf{x}_i)q[\mathbf{x}_i;\tilde{\mathbf{x}}]}$ and $\nu \sim U[0,1]$. If $\rho > \nu$ accept proposed sample $\mathbf{x}_{i+1} = \tilde{\mathbf{x}}$, else it is rejected.
- 5. Set i=i+1 and return to 2.





The accept/reject procedure ensures that samples that provide a better fit to the observations are immediately accepted, those that provide a similar fit are considered, and those that lead to simulated observations that are very different from the measurements are rejected.

Regions with relatively high probability are therefore preferentially sampled, whereas regions with low probability are avoided, and a sample of the posterior distribution is produced using far fewer iterations than direct computation of the PDF.

The theoretical underpinnings of the MCMC algorithm can be found in Mosegaard and Tarantola(2002) and Tarantola (2005).

The Particle filter

Like the MCMC the Particle filter (PF, also known as sequential MC) aims to produce a sample representation of the posterior distribution.

The basic idea is:

- 1. Draw a sample from the prior distribution $\mathbf{x}_i \sim P(\mathbf{x})$, for i = 1, ...N.
- 2. For each sample evaluate $\mathbf{y}_i = h(\mathbf{x}_i)$.
- 3. For each sample compute $w_i = P(\mathbf{y}|\mathbf{x}_i)$.
- 4. Normalise the weights $w_i = w_i / \sum_i w_i$

The posterior is then given by

$$P(\mathbf{x}|\mathbf{y}) \approx \sum_{i} w_i \delta(\mathbf{x} = \mathbf{x}_i)$$





Application of Bayes' to large problems

Non-parametric methods that make no assumption about the nature of the prior and the likelihood are inefficient for large-scale problems because the number of samples needed to represent the whole distribution grows exponentially with the size of the state to be estimated.

CURSE OF DIMENSIONALITY

The efficiency of the MCMC method can be quantified by the acceptance/reject ratio (0.33 in the example).

The efficiency of the particle filter can be measured by the effective sample size:

$$ess = \frac{1}{\sum_{i=1}^{N} w_i^2}$$

If the weights are very uneven the *ess* can approach 1 and the weighted sample will be overwhelmed by sampling noise (4.7 with N=50 in the example).

To increase the *ess*, PFs are being developed to increase the chance that the sample is the region of high likelihood, see van Leeuwen et al. 2019.

Parametric approaches - Gaussian assumption

In NWP we are interested in applying Bayes' theorem to approximately 10⁸ dimensions.

In many cases, it is appropriate to assume that the prior and likelihood are Gaussian

$$f(\mathbf{x}; \boldsymbol{\mu} \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$$

Mean vector

Gaussian assumption

Mean of the prior distribution $p(\mathbf{x}) = \frac{1}{(2\pi)^{n/2} |\mathbf{B}|^{1/2}} exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}^{\mathbf{b}})^{\mathsf{T}} \mathbf{B}^{-1}(\mathbf{x} - \mathbf{x}^{\mathbf{b}})\right)$ od:

The prior:

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{(2\pi)^{n/2}} e^{xp} \left(-\frac{1}{2} (\mathbf{y} - h(\mathbf{x}))^T \mathbf{R}^{-1} (\mathbf{y} - h(\mathbf{x})) \right)$$

Vector of observations

observation error covariance matrix

Observation operator, mapping from state to observation space

Applying Bayes' theorem the posterior is:

$$p(\mathbf{x}|\mathbf{y}) \propto exp\left(-\frac{1}{2}(\mathbf{x}-\mathbf{x}^{\mathrm{b}})^{T}\mathbf{B}^{-1}(\mathbf{x}-\mathbf{x}^{\mathrm{b}}) - \frac{1}{2}(\mathbf{y}-h(\mathbf{x}))^{T}\mathbf{R}^{-1}(\mathbf{y}-h(\mathbf{x}))\right)$$

If *h* is linear then the posterior is also Gaussian and can be parameterized according to its mean (xa, the analysis) and its (analysis error) covariance matrix.

Bayes' theorem scalar Gaussian illustration



The **mean** of p(y|x) and p(x) are the **observed** and **background/forecast** model values respectively.

The **standard deviations** of p(y|x) and p(x) are the **uncertainties** of the observed and background values respectively.

Bayes' theorem scalar Gaussian illustration



The uncertainty of the posterior is smaller than either the likelihood or the prior.

The mean = mode of p(x|y) gives the '**analysis**', which is also the minimum variance estimate!

Analytical solution – the Kalman equations

The analysis (the maximum a posteriori state) can be derived analytically as

 $\mathbf{x}^{a} = \mathbf{x}^{b} + \mathbf{K}(\mathbf{y} - h(\mathbf{x}^{b}))$ where $\mathbf{K} = \mathbf{B}\mathbf{H}^{T}(\mathbf{H}\mathbf{B}\mathbf{H}^{T} + \mathbf{R})^{-1}$

and \mathbf{H} is the linearized observation operator.

K (known as the Kalman gain) prescribes the weight given to the observations versus the prior. We see that the K increases as the prior uncertainty (B) increases and the observation uncertainty (R) decreases.

The analysis error covariance can also be derived analytically

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

We see as \mathbf{K} increases \mathbf{P}^{a} decreases.

In practice, we can not evaluate these expressions directly.

Variational and Ensemble Kalman Filter techniques

NERC/NCEO/DARC Training course on data assimilation and its interface with machine learning 2025

Variational (Var) DA

Var algorithms aim to find Maximum a-posterioiri (MAP) state/parameter set, which is also the mean of the Posterior distribution (assuming Gaussianity) and hence the same as the minimum variance estimate (Lorenc 1986).

In Var the analysis (the MAP estimate) is found by minimising the following cost function:

$$J(\mathbf{x}) = -const. \ln P(\mathbf{x}|\mathbf{y})$$

$$= -\ln(P(\mathbf{x})P(\mathbf{y}|\mathbf{x}))$$

$$= -\ln P(\mathbf{x}) - \ln P(\mathbf{y}|\mathbf{x})$$

$$P(\mathbf{x}|\mathbf{y}) \sim N(\mathbf{x}^{b}, \mathbf{B})$$

$$P(\mathbf{y} \mid \mathbf{x}) \sim N(\mathbf{y}, \mathbf{R})$$

$$J(\mathbf{x}) = (\mathbf{x} - \mathbf{x}^b)^{\mathrm{T}} \mathbf{B}^{-1} (\mathbf{x} - \mathbf{x}^b) + (\mathbf{y} - h(\mathbf{x}))^{\mathrm{T}} \mathbf{R}^{-1} (\mathbf{y} - h(\mathbf{x}))$$

If the observation operator is linear then the cost function is quadratic

Variational DA

1 variable Example:



The minimum of the cost function, x^a , is known as the analysis.

The assimilation window (4DVar)

The general form of variational DA is called 4DVar. This allows for observations over a window to be assimilated.



The 4DVar cost function is:

Can think of $h(m(\mathbf{x}_0, t_0, t_i))$ as a generalized ob operator:

$$J(\mathbf{x}_0) = (\mathbf{x}_0 - \mathbf{x}_0^b)^{\mathrm{T}} \mathbf{B}^{-1} (\mathbf{x}_0 - \mathbf{x}_0^b) + (\hat{\mathbf{y}} - \hat{h}(\mathbf{x}_0))^{\mathrm{T}} \hat{\mathbf{R}}^{-1} (\hat{\mathbf{y}} - \hat{h}(\mathbf{x}_0))$$

*Observations

uncertainty,

background

characterised by **B**

observation

4DVar

The minimum of the cost function (the analysis) can be found iteratively by searching in the direction of the gradient

$$\mathbf{x}_0^{k+1} = \mathbf{x}_0^k + \alpha \nabla J(\mathbf{x}_0^k)$$



The gradient of the cost function is given by

$$\nabla J(\mathbf{x}_0^k) = \mathbf{B}^{-1}(\mathbf{x}_0 - \mathbf{x}_0^k) + \sum_{i=1}^p \mathbf{M}_{t_0 \to t_i}^{\mathrm{T}} \mathbf{H}^{\mathrm{T}} \mathbf{R}_i^{-1} (\mathbf{y}_i - h(M_{t_0 \to t_i}(\mathbf{x}_0))),$$

where **M** is the tangent linear of the forecast model, $\mathbf{M} \in \mathbb{R}^{n \times n} = \frac{\partial \mathbf{x}_i}{\partial \mathbf{x}_0}$ evaluated at \mathbf{x}_0^k . \mathbf{M}^T is the model adjoint.

An alternative approach

A major disadvantage of variational techniques is that **B** is generally not updated on each assimilation cycle. In variational DA **B** is only designed to represent a climatological estimate of the error covariances. However, the errors can be highly flow-dependent.

This has motivated the development of the Ensemble Kalman Filter (EnKF).

The EnKF also based on linear/Gaussian theory so is a retractable method for large systems.

Unlike Var the EnKF represents ${f B}$ using an ensemble (and calls it ${f P}^{
m f}$).

The ensemble Kalman filter (EnKF)

The EnKF (Evensen 1994) merges KF theory with Monte Carlo estimation methods



The ensemble Kalman filter (EnKF)

The EnKF makes use of an ensemble approximation of the forecast error covariance matrix (previously the B matrix) to allow us to compute the Kalman equations directly.

$$\mathbf{P}^{f} = \frac{1}{N-1} \sum_{i=1}^{N} \left(\mathbf{x}_{k}^{(i),f} - \bar{\mathbf{x}}_{k}^{f} \right) \left(\mathbf{x}_{k}^{(i),f} - \bar{\mathbf{x}}_{k}^{f} \right)^{\mathrm{T}} = \frac{1}{N-1} \mathbf{X}^{\prime f} \underbrace{\mathbf{X}^{\prime f}}_{N-1} \mathbf{X}^{\prime f} \underbrace{\mathbf{X}^{\prime f}}_{N} \underbrace{\mathbf{X}^{\prime f}}_$$

The update step for the mean: $\bar{\mathbf{x}}^a = \bar{\mathbf{x}}^f + \mathbf{K}(\mathbf{y} - h(\bar{\mathbf{x}}^f)), \qquad \mathbf{K} = \mathbf{P}^f \mathbf{H}^T (\mathbf{H} \mathbf{P}^f \mathbf{H}^T + \mathbf{R})^{-1}$

EnKF – the update step for the covariance

The analysis error covariance matrix can be estimated in two ways

• Stochastic methods: Samples from the posterior distribution are created

 $\mathbf{x}_{k}^{(i),a} = \mathbf{x}_{k}^{(i),f} + \mathbf{K}_{k}(\mathbf{y}_{k} + \mathbf{\epsilon}_{y}^{(i)} - h(\mathbf{x}_{k}^{(i),f})), \text{ where } \mathbf{\epsilon}_{y} \sim N(\mathbf{0}, \mathbf{R}). \text{ Then the sample}$

covariance matrix is computed as:

$$\mathbf{P}^{a} = \frac{1}{N-1} \sum_{i=1}^{N} \left(\mathbf{x}_{k}^{(i),a} - \bar{\mathbf{x}}_{k}^{a} \right) \left(\mathbf{x}_{k}^{(i),a} - \bar{\mathbf{x}}_{k}^{a} \right)^{\mathrm{T}}$$

• **Deterministic methods**: An analysis perturbation matrix is generated by applying a transformation to the forecast perturbation matrix $\mathbf{X}_{k}^{'a} = \mathbf{X}_{k}^{'f}\mathbf{T}_{k}$. The transform **T** is derived to ensure $\mathbf{P}^{a} = (\mathbf{I} - \mathbf{K}\mathbf{H})\mathbf{P}^{f}$ holds.

$$\mathbf{P}^a = \frac{1}{N-1} \mathbf{X}^{\prime a} (\mathbf{X}^{\prime a})^T$$

EnKF and Sample errors

Because the ensemble size is limited, $N \ll n$. This means that the ensemble estimate of \mathbf{P}^{f} will be low rank and affected by sample errors.

The sample errors will affect

- The estimate of the **error variances** so that the wrong weight is given to the observations. When the the error variances are underestimated this can lead to filter divergence!!!!
- The **error correlations** such that observations can update variables which they do not contain information about.



) N = 128



From Houtekamer & Mitchell (1998)

EnKF and Sample errors

To mitigate these two problems, we use a combination of

- variance inflation (e.g. Mitchell and Houtekamer, 2000; Anderson and Anderson, 1999, Whitaker and Hamill, 2012) and
- **covariance localization (**Houtekamer and Mitchell, 2001, Hunt et al. 2007).





Summary



There are many different data assimilation algorithms proposed. The ones presented all aim to solve Bayes' theorem to find the probability distribution of the state consistent with the uncertainty in the observations and model.

Many methods used operationally rely on the Gaussian assumption to make the DA problem tractable for large dimensions and efficient to run as part of a cycled forecasting system.

These methods include the variational and Ensemble Kalman Filter techniques.

Summary



However, there are many reasons why the Gaussian assumption may not hold:

- Non-linear model
- Non-linear observation operator
- Non-Gaussian errors e.g. if the variable is bounded

Assuming Gaussianity in these cases can result in:

- High probabilities assigned to unphysical states that may result in numerical instabilities
- Low probabilities assigned to important regimes
- Extreme events not represented

Even within Var and EnKF, there is room for the relaxation of the Gaussian assumption e.g. inner/outer loops, Gaussian anamorphosis.

Can reformulate the problem to not make Gaussian assumptions e.g. the particle filter.

A note on parameter estimation

Each of the DA methods can be modified to estimate parameters instead of (or as well as) estimating the initial state. For example,

- Model parameters describing physical processes
- Parameters to describe the bias correction of the model or observations

Need to consider

- How the prior uncertainty is represented. If you updating the parameters using variational or EnKF methods, does a Gaussian error make sense?
- How can the observations be related to these parameters? Often, you do not observe parameters directly, so instead rely on their errors being correlated with state variables that are observed.

DA software

On Wednesday and Thursday we will dedicate more time to learning about Variational and EnKF algorithms with computer practicals.

Our code is available on <u>https://github.com/darc-reading/darc-</u> <u>training-2025/</u>

Tomorrow Yumeng and Chris will give a lecture DA software packages

Schematic comparison of techniques



MCMC all samples 7 accepted 6 > 5 > 4 3 2 2 2 8 u sample from prior weighted sample from posterior × likelihood likelihood Particle Filter > 5 > 5 6 А 8 2 6 8 2 4 u u sample from prior 0 sample from posterior likelihood likelihood 6 EnKF > 5 > 56 2 6 8 4

u

u

Can combine the best bits of (hybridise) the different algorithms

NERC/NCEO/DARC Training cour and its interface with machine learning 2025

References

- Mosegaard and Tarantola (2002): Probabilistic Approach to Inverse Problems. International Geophysics
- Tarantola (2005): Inverse problem theory and methods for model parameter estimation. SIAM
- van Leeuwen et al. (2019): Particle filters for high-dimensional geoscience applications: A review. *Q J R Meteorol Soc.*
- Lorenc (1986): Analysis methods for numerical weather prediction. Q.J.R. Meteorol. Soc
- Evensen (1994): Sequential data assimilation with a nonlinear quasi-geostrophic model using Monte Carlo methods to forecast error statistics, *J. Geophys. Res.*
- Mitchell and Houtekamer, (2000): An Adaptive Ensemble Kalman Filter. Mon. Wea. Rev.
- Anderson and Anderson, (1999): A Monte Carlo Implementation of the Nonlinear Filtering Problem to Produce Ensemble Assimilations and Forecasts. *Mon. Wea. Rev.*
- Whitaker and Hamill, (2012) Evaluating Methods to Account for System Errors in Ensemble Data Assimilation. *Mon. Wea. Rev.*
- Houtekamer and Mitchell, (2001): A Sequential Ensemble Kalman Filter for Atmospheric Data Assimilation. Mon. Wea. Rev.
- Hunt et al. (2007): Efficient data assimilation for spatiotemporal chaos: A local ensemble transform Kalman filter, Physica D: Nonlinear Phenomena