

Overview of Data Assimilation Methods

Massimo Bonavita

ECMWF

Massimo.Bonavita@ecmwf.int

Outline

- What is data assimilation and how does it work?
- Data assimilation ingredients: Observations and models
- Blending observations and model: the Bayes perspective
- A whirlwind introduction to DA methods in the geophysical sciences:
 - Particle Filters
 - Kalman Filters
 - Variational methods
 - Hybrid methods

Data Assimilation

NWP definition of DA: Process by which “optimal” initial conditions for numerical forecasts are estimated.

- The best analysis is the analysis that leads to the best forecast (different criteria for other applications, eg reanalysis)
- Optimal in a statistical sense: minimises error and/or maximises probability of the analysis being accurate
- Provides an estimate of initial uncertainties, typically through a Monte Carlo procedure (i.e., ensemble DA)
- Do it “quickly” – typically in less than 45 minutes on a large high performance computer (for global NWP; for limited area NWP available time is usually much less!) – and in a cycled manner

Data Assimilation

The goal of Data Assimilation is:

“Estimate the probability distribution function (pdf) of the Earth system state at the initial time”

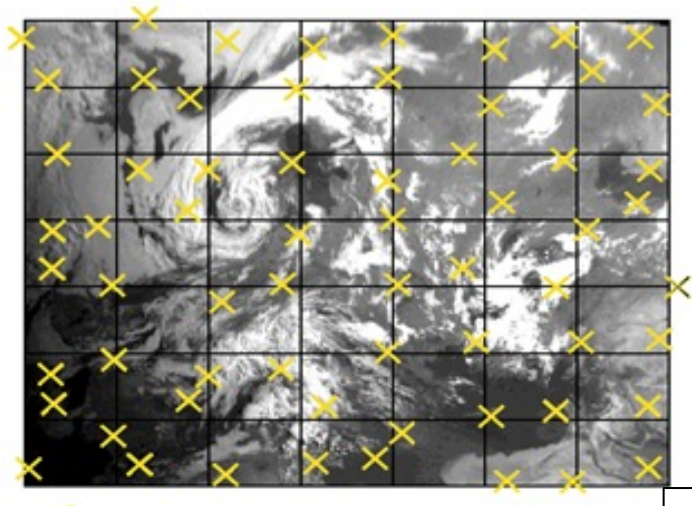
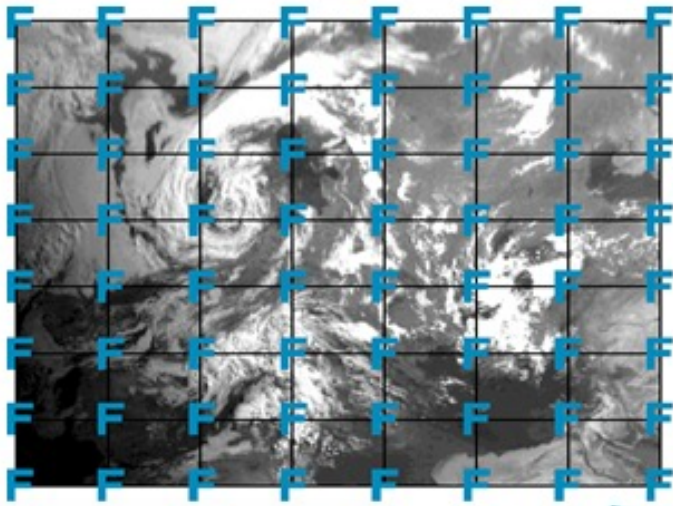
The initial state pdf can be explicitly **sampled** (Ensemble DA) and is usually summarised in terms of its central value (the “analysis”) and its uncertainty (the variance around the central estimate).

This representation of the initial pdf in terms of its first two moments (mean and covariance) is appropriate for \sim Gaussian (or at least unimodal) error distributions, it loses meaning for multimodal error distributions.

In large scale geophysical applications of DA we are forced to assume approx. Gaussianity to make the problem computationally tractable

Model Forecast (with errors)

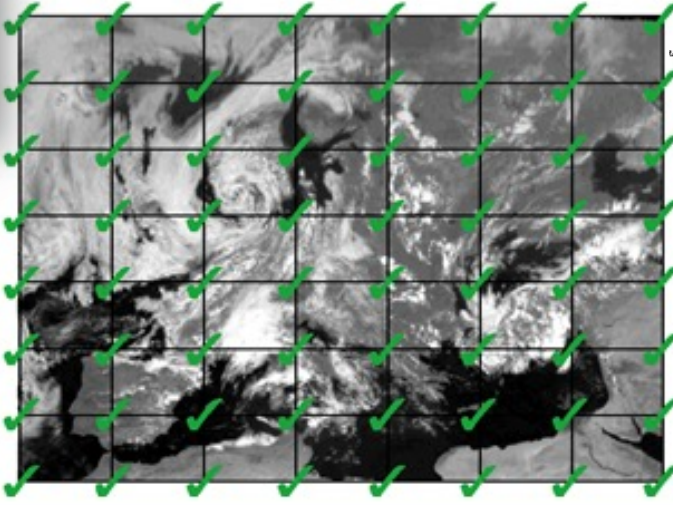
Observations (with errors)



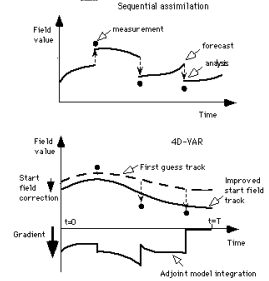
Computer (with a lot of CPUs/GPUs)



Clever people

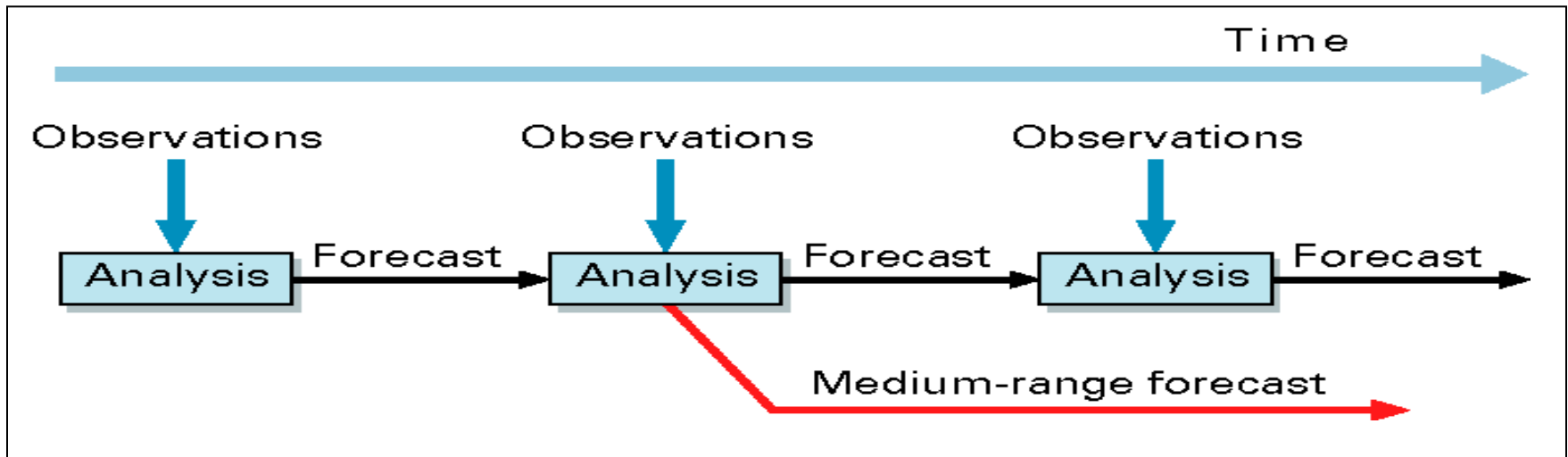


$$J(x) = (x - x_b)^T B^{-1} (x - x_b) + \sum_{i=1}^n (y_i - H_i[x_i])^T R_i^{-1} (y_i - H_i[x_i])$$



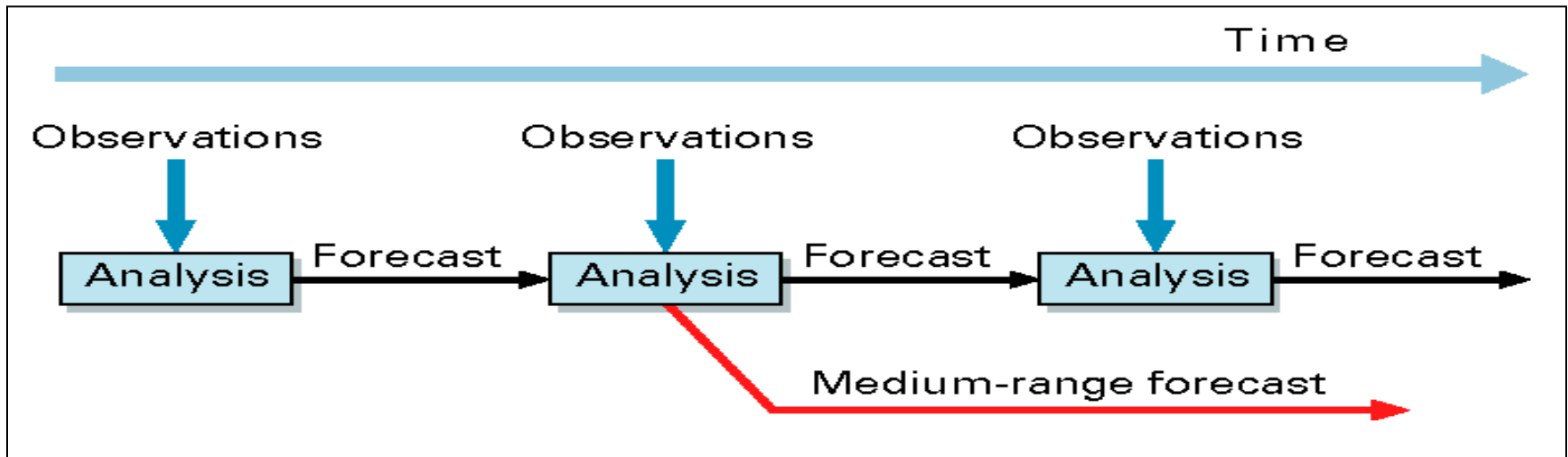
Analysis (with - smaller - errors)

The Data assimilation cycle



- An analysis is not produced by observations alone!
- The observations are used to correct errors in the short forecast from the previous analysis time (every 12 hours at ECMWF, 6-3 hourly in other global NWP systems; 3-1 hourly for higher resolution, limited area models).
- The short-range forecast carries information from past observations into the current analysis

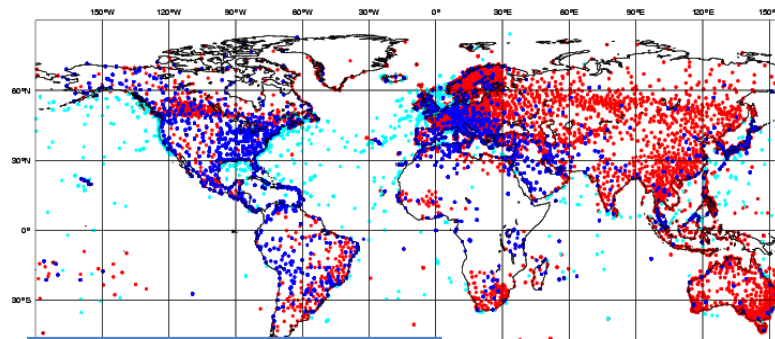
The Data assimilation cycle



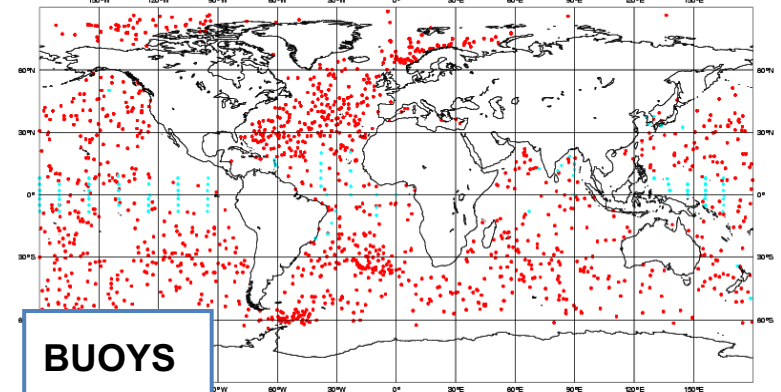
- At ECMWF, twice a day about 40,000,000 observations are used to correct the 250,000,000 variables that define the analysis state.
- This is done by a 4-dimensional adjustment in space and time based on the available observations (4D-Var); 4D-Var is computer intensive (approx. as much as the 10-day forecast)

Observations

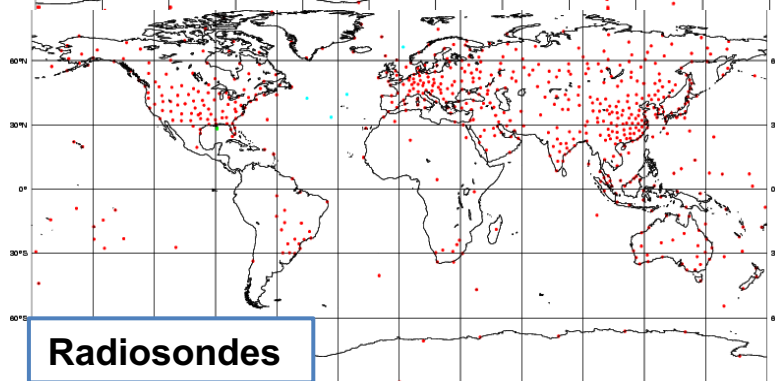
Distribution of in situ observations



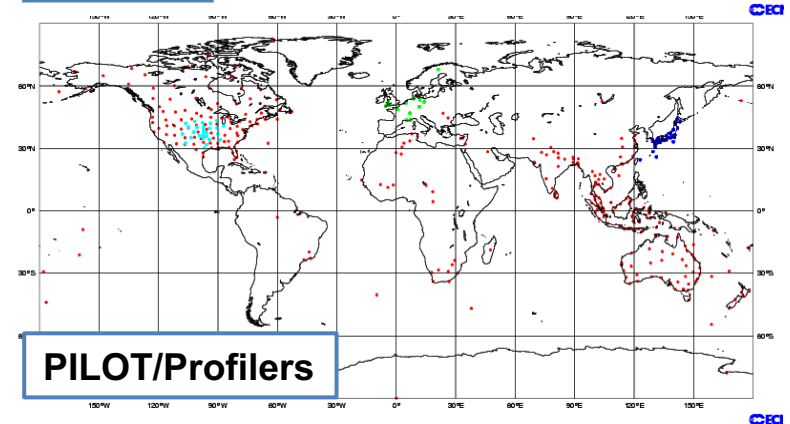
SYNOP/METAR/SHIP



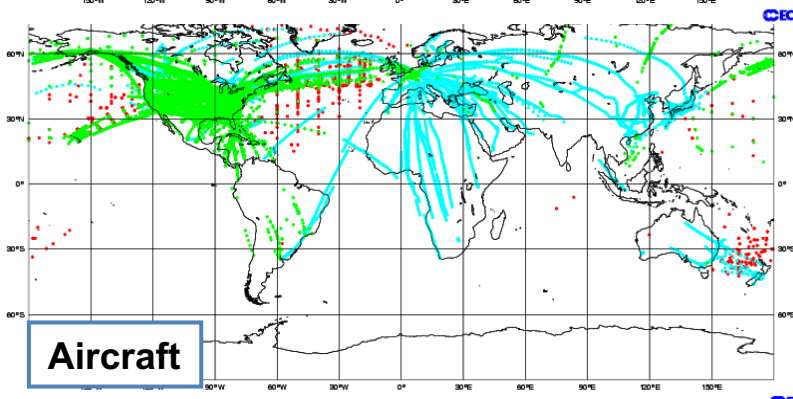
BUOYS



Radiosondes



PILOT/Profilers

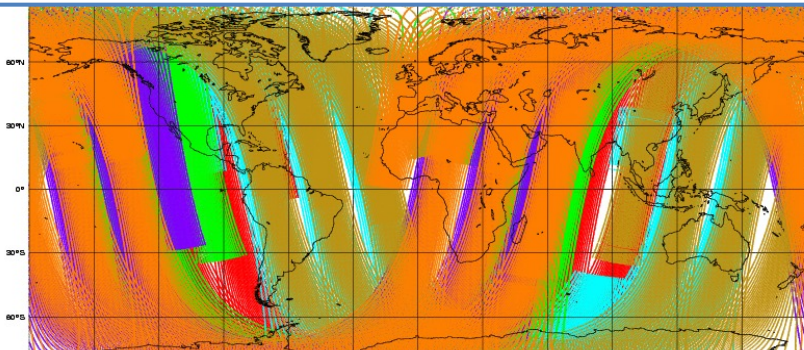


Aircraft

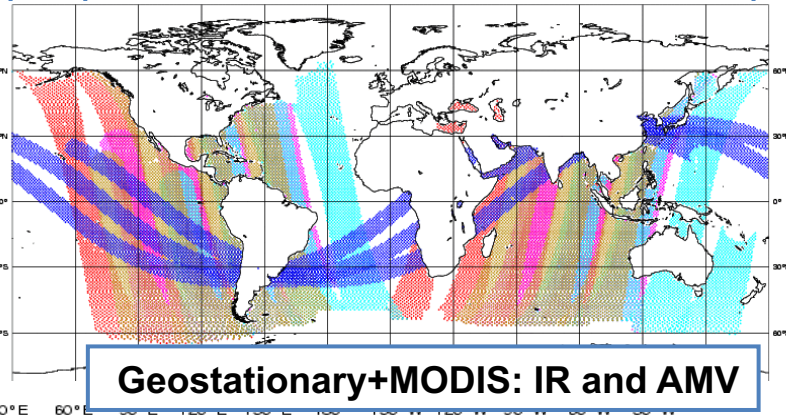
<https://www.ecmwf.int/en/forecasts/quality-our-forecasts/monitoring-observing-system>

Satellite data sources used by ECMWF's analysis

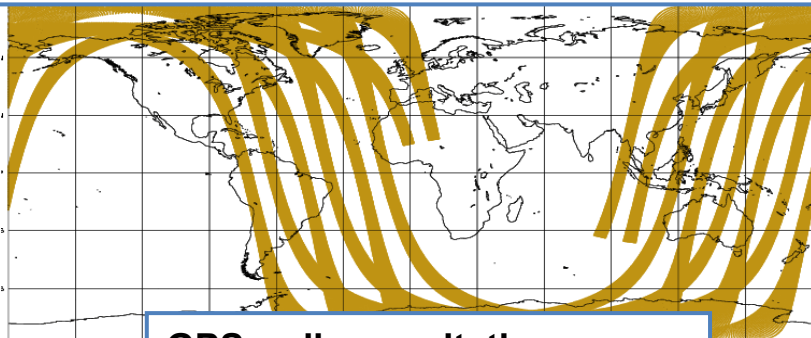
Sounders: NOAA AMSU-A/B, HIRS, AIRS, IASI, MHS



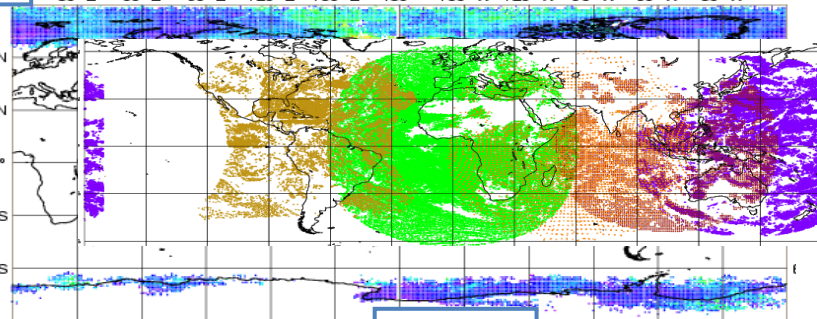
Imagers: SSMI, SSMIS, AMSR-E,



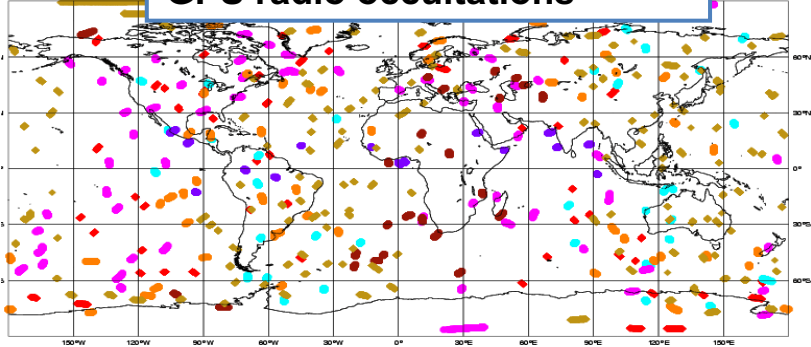
Scatterometer ocean low-level winds: ASCAT



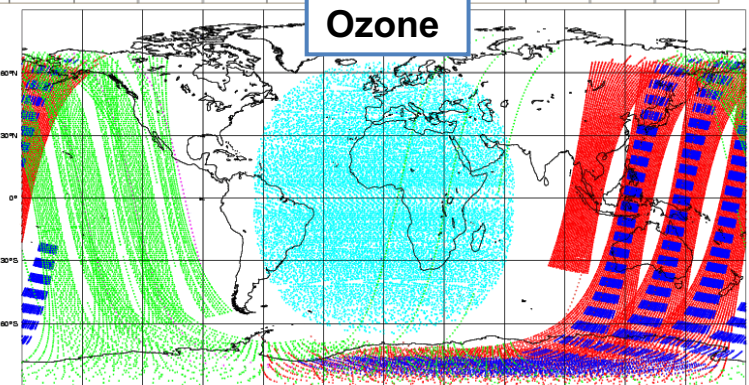
Geostationary+MODIS: IR and AMV



GPS radio occultations



Ozone



Observation errors

- Even the current global observing system is not able to observe the atmosphere completely (data voids): from a mathematical standpoint data assimilation is an under-determined problem
- Most satellite observations (e.g. radiances) are only indirectly related to the quantities of interest (i.e., grid point values of T, u, v, q, O_3, \dots)
- Majority of satellite observations have coarse vertical resolution
- Observations measure quantities not located at grid points

In order to compare observations (\mathbf{y}) and model (\mathbf{x}) we need to perform spatial and temporal interpolations of the model fields and (for satellite observations) transform model fields into radiances: we call this set of operations the **observation operator** (\mathcal{H}):

$$\mathbf{y} = \mathcal{H}(\mathbf{x})$$

Note: we project model fields into observed quantities, not observations into model variables (this second operation is called a “retrieval”)

Observation errors

- Observations are affected by different types of errors
- Denoting \mathbf{y}^* as the true observations of the model state ($\mathbf{y}^* = \mathcal{H}(\mathbf{x}^*)$):

$$\mathbf{y} - \mathbf{y}^* = \varepsilon_o = \varepsilon_G + \varepsilon_M + \varepsilon_R + \varepsilon_H$$

ε_G = **Gross errors** (incorrect coding of observation, duplicates, incorrect location, wrong cloud clearing, etc.).

ε_M = **Measurement errors** (instrument noise)

ε_R = **Representativity errors** (e.g., in situ observations compared to grid point model value)

ε_H = **Observation operator** (Forward model) errors (e.g., errors in the radiative transfer model, interpolation errors, etc.)

Observation errors

$$\mathbf{y} - \mathbf{y}^* = \varepsilon_o = \varepsilon_G + \varepsilon_M + \varepsilon_R + \varepsilon_H$$

- ε_G (gross errors) are dealt with by **Observation Quality Control** techniques (to be discussed later this week); Some of these checks are applied before ingesting the observations (Climatological checks, consistency checks, first guess checks), others are part of the analysis algorithm itself (buddy checks, variational quality control)

- Observations are assumed to be un-biased:

$$\langle \varepsilon_o \rangle = 0$$

- Biases are dealt with specific **Bias Correction** techniques (to be discussed this week), which can be part of the analysis algorithm itself (e.g., Variational Bias Correction)

- The covariance matrix of the observation errors is denoted as \mathbf{R} :

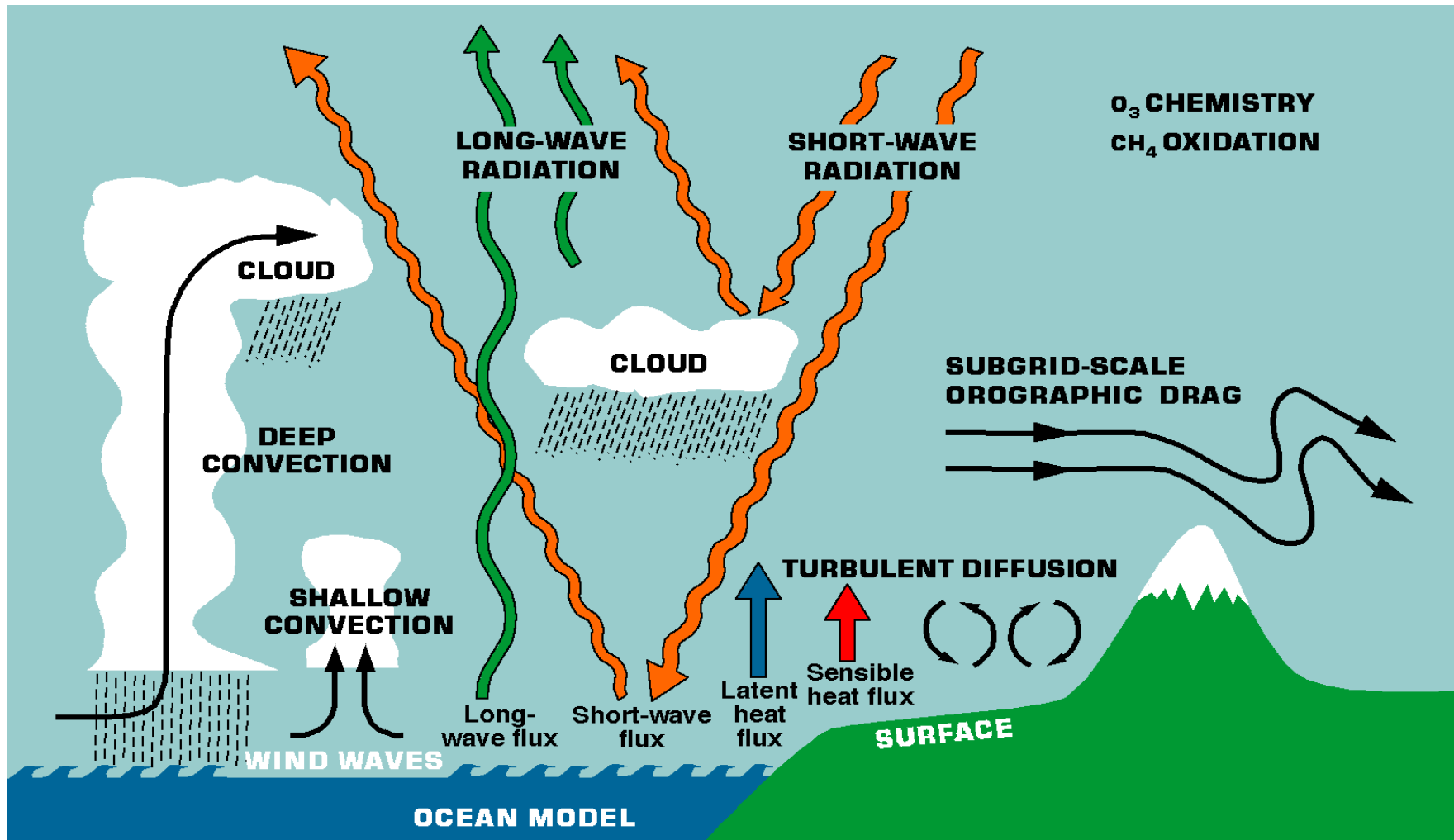
$$\langle \varepsilon_o \varepsilon_o^T \rangle = \mathbf{R}$$

The forecast model

The forecast model is a very important part of the data assimilation system

- The short-range forecast connecting successive analysis updates carries information from past observations to the current analysis time (this is called the “**background**”): the better the model the more accurate the background state
- A good model starting from accurate previous analysis will produce an accurate background ➡ the analysis will make only **small corrections** to the background
- In fact when the analysis makes **large corrections** to the background state is usually a sign that something interesting is happening... (e.g., rapid development not present in the forecast; suspect observations)
- In modern data assimilation methods (4D-Var, EnKF, PF) the analysed state is constructed so as to **respect the physical and dynamical balances of the model** ➡ the model is an integral part of the analysis algorithm

The forecast model is a very important part of the data assimilation system



Physical processes in the ECMWF model

Model errors

- Despite their increasing complexity and sophistication models are not perfect (yet)!
- Sources of model error include: missing physical processes, errors in parametrizations of physical processes, discretisation errors (from continuous PDEs to discrete formulation), error in the forcing fields, etc.,
- We define **model error** as (* denotes true state, i is the time index):

$$\mathbf{x}_i^* = \mathcal{M}(\mathbf{x}_{i-1}^*) + \boldsymbol{\eta}_i$$

- Model error can in general have non zero mean:

$$\langle \boldsymbol{\eta}_i \rangle \neq 0$$

- The covariance matrix of the model errors is denoted as \mathbf{Q} :

$$\langle \boldsymbol{\eta}_i \boldsymbol{\eta}_i^T \rangle = \mathbf{Q}_i$$

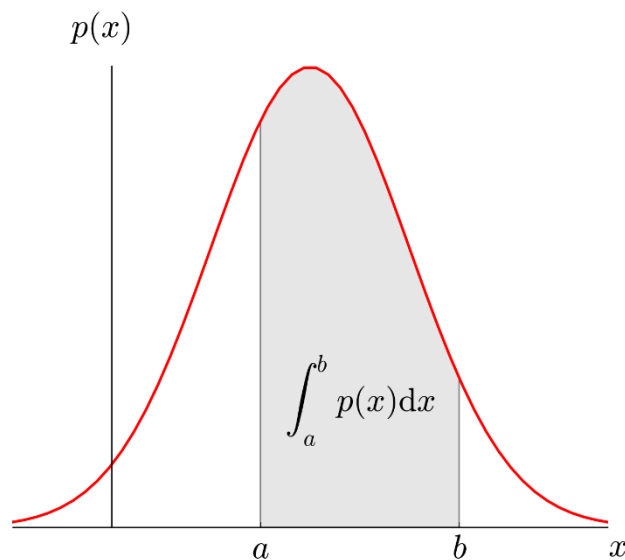
- The treatment of **model error in DA** will be discussed in more detail in a lecture later this week

Blending observations and model information: the Bayes perspective

The Bayes perspective

- Both observations and models are affected by random errors*
- This means that they should be described as **random variables**
- All we can/need to know about random variables are their **probability distribution functions**:

$$\Pr[a \leq X \leq b] = \int_a^b p(x)dx$$



* Assume here that systematic errors have been separately dealt with

The Bayes perspective

- **Bayes law** descends directly by the definition of conditional probabilities:

$$p(A, B) = p(A|B)p(B) = p(B|A)p(A)$$

⇒

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)}$$

Where:

$p(A, B)$ = probability of events A and B both happening (**joint** prob. distribution)

$p(A|B)$ = probability of event A given that B is true (**conditional** prob. distribution)

$p(A), p(B)$ = probability of event A (B) happening (**marginal** prob. distribution)

The Bayes perspective

- An illustration (http://en.wikipedia.org/wiki/Base_rate_fallacy):

The police have been issued with breathalysers which never fail to detect a drunk person but have a 5% rate of false positives. Prior campaigns have shown that, on average, one in one thousand drivers drives drunk. If the police stop a driver at random, and he/she results positive to the breathalyser, what is the probability that he/she is actually drunk?

Event A: being a drunk driver. Probability of being a drunk driver, before being tested: $p(A) = 0.001$

Event B: testing positive to the breathalyser. The probability of testing positive is 1 for the drunken subset of the drivers (0.001) and 0.05 for the sober subset of the drivers (0.999):

$$p(B) = (1 * 0.001) + (0.999 * 0.05) = 0.05095$$

Probability of testing positive to the breathalyser when drunk: $p(B|A) = 1$

Probability of being drunk after testing positive to the breathalyser, $p(A|B)$:

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)} = \frac{1 * 0.001}{0.05095} = 0.0198$$

In words: out of 1000 people stopped by the police, about 51 will result positive, but the probability that anyone of them is actually drunk is less than 2%. (This shows how Bayesian thinking can be useful even beyond data assimilation!)

The Bayes perspective

- Another illustration: the Monty Hall problem (https://en.wikipedia.org/wiki/Monty_Hall_problem):

This is a brain teaser inspired by the American quiz show *Let's Make a Deal* and named after its original host, Monty Hall. The version of the problem that appeared in the Parade mag1990 read:

"Suppose you're on a game show, and you're given the choice of three doors. Behind one door is a car, behind the others, goats. You pick a door, say #1, and the host, who knows what's behind the doors, opens another door, say #3, which has a goat. He says to you, "Do you want to pick door #2?" Is it to your advantage to switch your choice of doors?"

Let us use our Bayesian tools to see what the savvy game show participants should do! Let us indicate with $p(1)$, $p(2)$, $p(3)$ the probability that the car is behind door 1,2,3. Initially $p(1)=p(2)=p(3)=1/3$. To check whether it is a good idea to switch door, we are interested in is the probability of the car being behind door=2 after the host has chosen door=3 and we (the guest) have chosen door=1, in symbols:

$$p(2|H = 3, G = 1)$$

The Bayes perspective

- Another illustration: the Monty Hall problem (continued)

$$p(2|H = 3, G = 1) = \frac{p(2, H=3, G=1)}{p(H=3, G=1)} \quad (\text{from the definition of conditional prob.})$$

$$\frac{p(2, H=3, G=1)}{p(H=3, G=1)} = \frac{p(H=3 | 2, G=1) p(2, G=1)}{p(H=3, G=1)} \quad (\text{again from def. of conditional prob.})$$

Now note that a) the probability that the host chooses door=3, given the car is in 2 and we have chosen door=1 is 1; and b) the probability that the car is behind door=2 is independent of our choice of door=1:

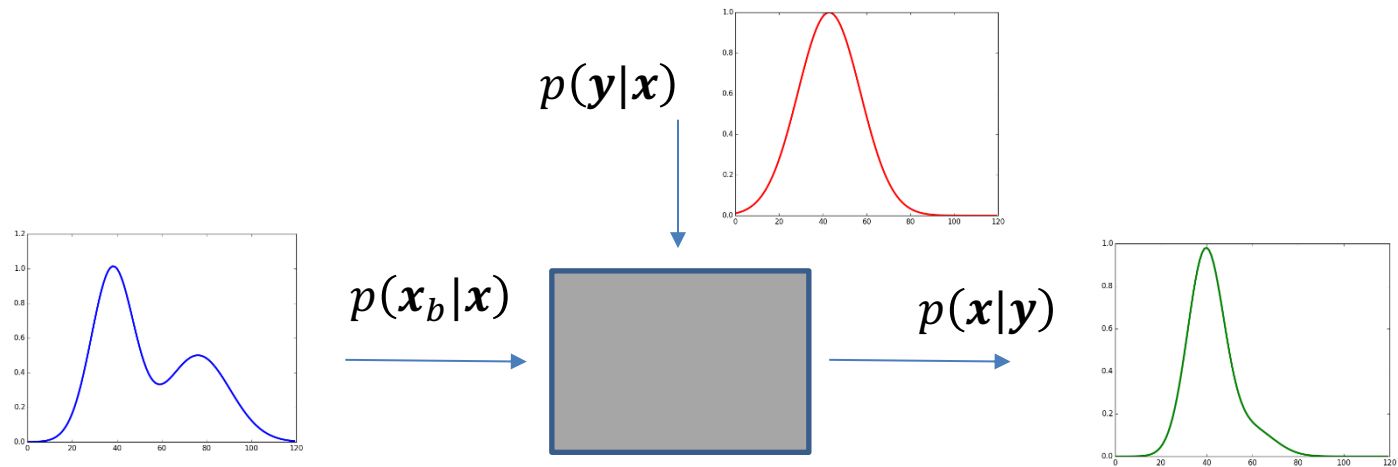
$$\frac{p(H=3 | 2, G=1) p(2, G=1)}{p(H=3, G=1)} = \frac{p(2) p(G=1)}{p(H=3 | G=1) p(G=1)} = \frac{p(2)}{p(H=3 | G=1)} = \frac{1/3}{1/2} = \frac{2}{3}$$

It does make sense to switch our choice to door=2!

This is another example of how the probability of the new piece of information (the Host's choice of door=3) has modified the a-priori probability of where the car might be.

The Bayes perspective

- At an abstract level, we can think of the analysis process as updating our prior knowledge about the state, represented by a background forecast and its pdf, with new observations, represented by their values and their respective pdfs:



$$p(x|y) = \frac{p(y|x)p(x)}{p(y)} = \frac{p(y|x)p(x_b|x)}{p(y)} \propto p(y|x)p(x_b|x)$$

- $p(x_b|x)$ = **prior pdf** (encapsulate our knowledge about the state before new observations)
- $p(y|x)$ = **observations likelihood** (pdf of the observations conditioned on the state)
- $p(x|y)$ = **posterior pdf** (updated pdf of the state after the analysis)
- $p(y)$ = **marginal pdf of the observations** (does not depend on x : normalising constant in Bayes' law)

Particle Filters

$$p(\mathbf{x}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{x})p(\mathbf{x}_b|\mathbf{x}) \quad (1)$$

- In principle an analysis update requires being able to compute the product pdf of the random variables \mathbf{y} , \mathbf{x}_b . This is usually not possible to do explicitly unless we choose very specific functional forms for the pdfs
- We thus need to make approximations
- One idea is to use Monte Carlo methods to sample and propagate the pdfs in (1):
Particle Filters
- In Particle Filters, pdfs are sampled by a collection of “particles” (i.e., model states) with assigned weights:

$$p(\mathbf{x}) \sim \sum_{i=1, N} w_i \delta(\mathbf{x} - \mathbf{x}_i) \quad (2)$$

Particle Filters

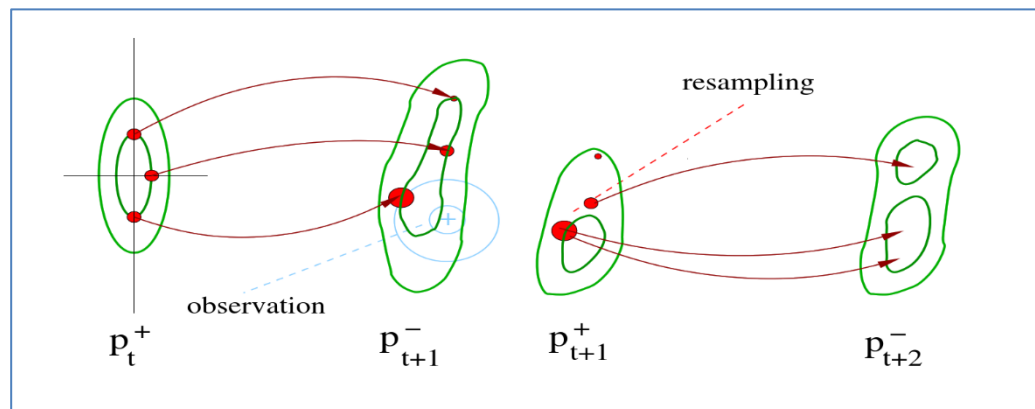
- The pdf is propagated in time by integrating the different particles with the model:

$$p(\mathbf{x}_b|\mathbf{x}) \sim \sum_{i=1, N} w_i \delta(\mathbf{x} - M(\mathbf{x}_i)) \quad (3)$$

- In the analysis update the weights of the particles are updated according to the observations' likelihood:

$$w_i^a \propto w_i p(\mathbf{y}|\mathbf{x}_i)$$

- The ensemble of particles is usually resampled, i.e. high-weight particles are duplicated and low-weight particles discarded
- The Particle Filter described here is one of the most basic implementation (Bootstrap Particle Filter, Gordon et al., 1993)



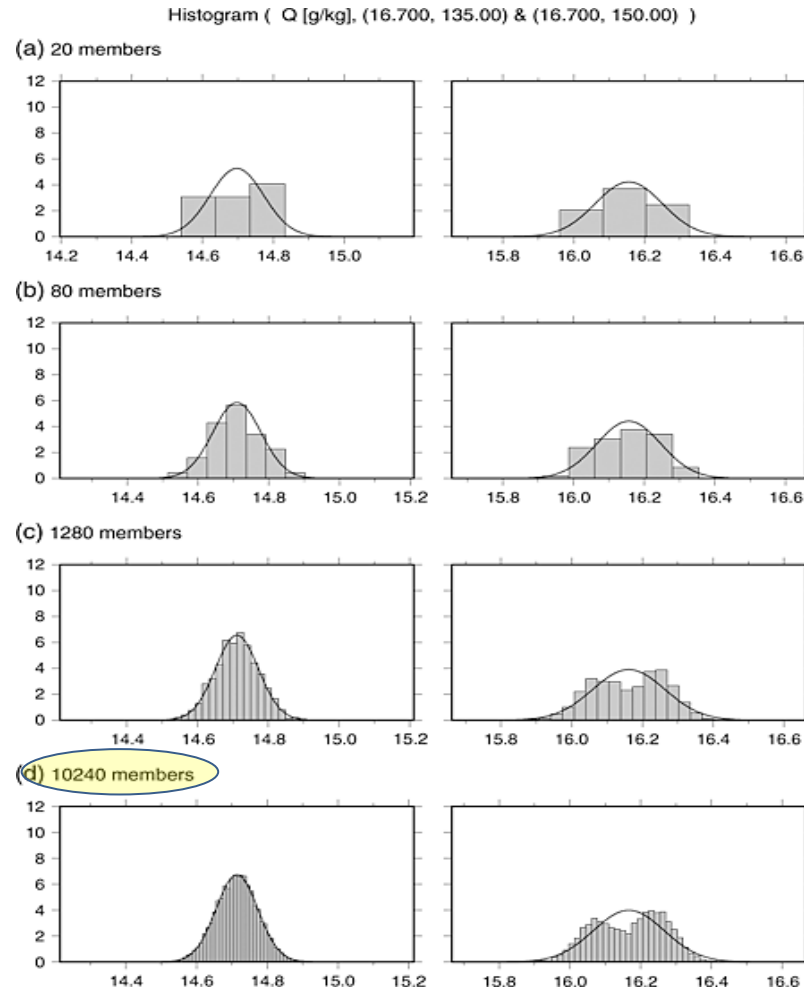
From M. Bocquet

Particle Filters

- Particle Filters work well for very small state space sizes **and** observation sizes ($N \sim 10$ to 100)
- For larger state space and/or observation sizes the required number of particles increases **exponentially** (Snyder et al., 2015)
- A large body of contemporary research is devoted to reduce the computational demands of particle filters for high dimensional systems
- One of the main themes of PF research is how to prevent the particles from diverging from the true state and becoming too unlikely, i.e. uninformative about the true state
- One of the ideas is to also use observations (and not only the model) to “guide” the particles’ evolution from $t=t_{n-1}$ to $t=t_n$; many variants possible (Ades and van Leeuwen, 2014)
- Another over-arching idea is to introduce some form of localisation in the PF (similar to what is done for the EnKF): see Farchi and Bocquet, 2018, for a review

Particle Filters

- Regardless of the assimilation algorithm the number of particles (ensemble members) needed to reliably resolve non-Gaussian pdfs is very high:



Histograms of a 6 h ensemble forecast for specific humidity (g kg^{-1}) for a intermediate AGCM. Miyoshi et al., 2014

The Gaussian approximation

- Not making assumptions on the shape of the prior and the likelihood pdf makes the Bayesian problem difficult (i.e., analytically and computationally intractable)
- Usual choice is to assume a **Gaussian distribution** for the both the observations' likelihood and the prior pdf of the background forecast
- Why Gaussian?
 1. Mathematically tractable problem;
 2. Full distribution characteristics defined by only its first two moments (mean and covariance);
 3. Supported by the Central Limit Theorem;
 4. Least committed distribution for given first and second moments (i.e., we are making the least amount of hypotheses on the shape of the pdf for a given assumed variance)

The Gaussian approximation

- Usual choice is to assume a **Gaussian distribution** for the both the observations' likelihood and the prior pdf

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{(2\pi)^{N/2}|\mathbf{R}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{H}(\mathbf{x}))^T (\mathbf{R})^{-1}(\mathbf{x}_b - \mathbf{H}(\mathbf{x}))\right)$$

$$p(\mathbf{x}_b|\mathbf{x}) = \frac{1}{(2\pi)^{N/2}|\mathbf{P}_B|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x}_b - \mathbf{x})^T (\mathbf{P}_B)^{-1}(\mathbf{x}_b - \mathbf{x})\right)$$

$$p(\mathbf{x}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{x})p(\mathbf{x}_b|\mathbf{x}) \propto \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{H}(\mathbf{x}))^T (\mathbf{R})^{-1}(\mathbf{y} - \mathbf{H}(\mathbf{x})) - \frac{1}{2}(\mathbf{x}_b - \mathbf{x})^T (\mathbf{P}_B)^{-1}(\mathbf{x}_b - \mathbf{x})\right)$$

- where $\langle \varepsilon_o \varepsilon_o^T \rangle = \mathbf{R}$ and $\langle \varepsilon_b \varepsilon_b^T \rangle = \mathbf{P}_B$ are the covariances of the errors of the observations and of the prior (background forecast)
- Under this assumption **the posterior (analysis) distribution $p(\mathbf{x}|\mathbf{y})$** can also be expressed as a **Gaussian** distribution (to be shown during the week)

Kalman Filter methods

- Once we know (at least in principle!) the form of the posterior distribution $p(\mathbf{x}|\mathbf{y})$ we have a choice:

- 1) Either we can solve for the **mean** and the **covariance** of the posterior distribution:

$$\mathbf{x}_a = \int \mathbf{x} p(\mathbf{x}|\mathbf{y}) d\mathbf{x}$$

$$\mathbf{P}_a = \int (\mathbf{x} - \mathbf{x}_a)(\mathbf{x} - \mathbf{x}_a)^T p(\mathbf{x}|\mathbf{y}) d\mathbf{x}$$

Methods based on this approach include **Optimum Interpolation** (O.I.), **Kalman Filter**, **Ensemble Kalman Filter** (EnKF). These will all be discussed this week. The analysis found through this approach is referred to as the **minimum variance** solution or the best linear unbiased estimate (**BLUE**).

Note: Kalman Filter based methods can be derived without making any assumptions about the Gaussianity of the errors. However only if all error distributions are Gaussian will the KF provide the correct posterior distribution (i.e. Bayes posterior pdf).

Variational methods

- 2) Alternatively we might choose to estimate the **mode** of the posterior distribution $p(\mathbf{x}|\mathbf{y})$, i.e. find the analysis \mathbf{x}_a as the state that corresponds to the maximum of the posterior distribution (\Rightarrow the most probable state):

$$\mathbf{x}_a = \arg \max_{\mathbf{x}} (p(\mathbf{x}|\mathbf{y}))$$

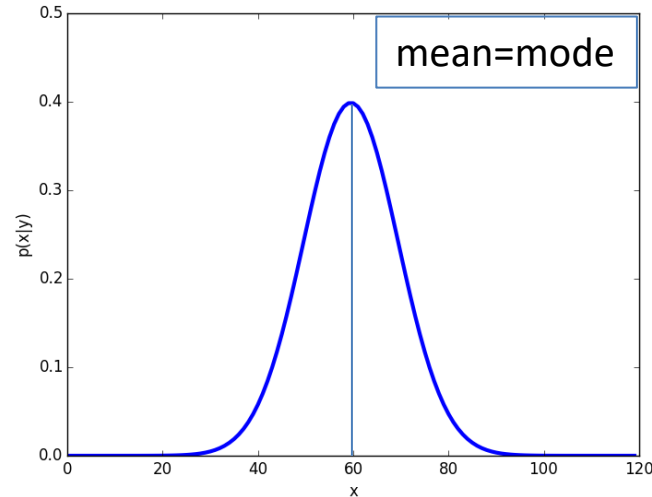
This way of attacking the problem leads to the **variational approach** (3D-Var, 4D-Var). They will be covered extensively in this week's lectures. The solution found in this way is called the **maximum a-posteriori probability** state (**MAP**).

In the variational framework the linear and Gaussian assumptions can be relaxed, i.e. the full nonlinear analysis problem can be decomposed into a series of linear Gaussian problems (incremental 4D-Var, to be discussed later this week).

However there is no guarantee of convergence!

Kalman Filter vs Variational methods

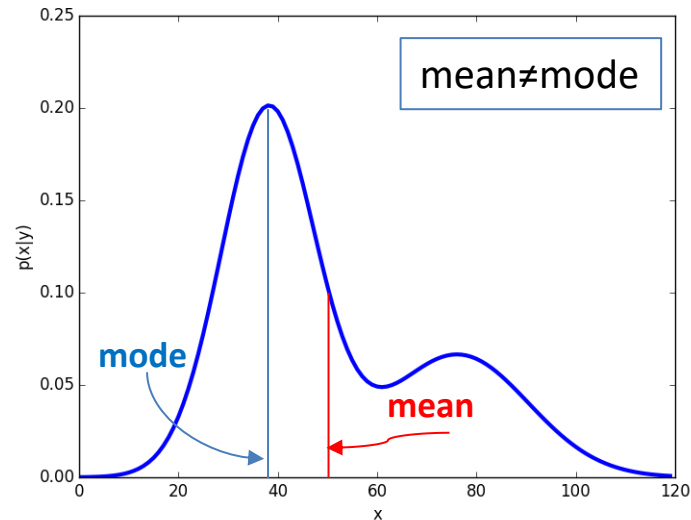
- For a Gaussian pdf the mean and the mode coincide:



- Thus if all the system statistics are Gaussian the minimum variance and maximum a-posteriori solutions coincide

Kalman Filter vs Variational methods

- For non-Gaussian pdfs the mean and the mode of the distribution generally differ:



- In non-Gaussian assimilation problems the minimum variance and maximum likelihood solutions will differ
- Which solution is better is problem dependent
- The **more non-Gaussian** the problem the more one needs information about the **whole posterior pdf**, not only its first two moments!

Hybrid DA methods

- Both Variational and Kalman Filter based analysis methods require estimates of the background state and its error covariances ($p(\mathbf{x}_b|\mathbf{x}) \sim \mathcal{N}(\mathbf{x}_b, \mathbf{P}_B)$)
- The background state is usually provided by an integration of the forecast model started from the previous analysis:

$$\mathbf{x}_b^t = \mathcal{M}(\mathbf{x}_a^{t-1})$$

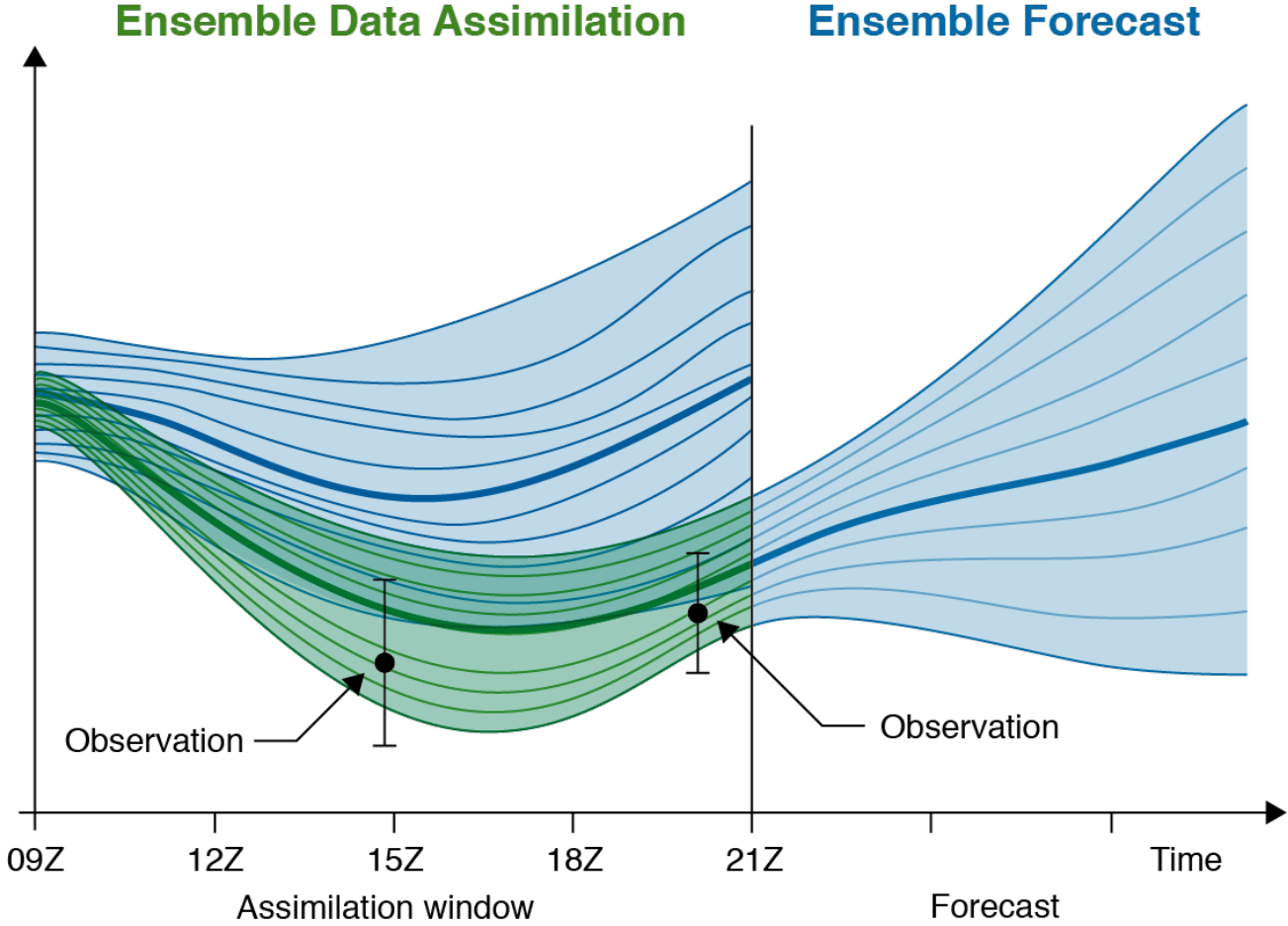
- The background (and analysis) **error statistics** are usually sampled with **Monte Carlo methods**: an ensemble of states is used to estimate the errors statistics
- Each ensemble member is advanced using a perturbed version of the model:

$$\mathbf{x}_{b,i}^t = \mathcal{M}(\mathbf{x}_{a,i}^{t-1}) + \boldsymbol{\eta}_i \quad i = 1, 2, \dots, N_{ens}$$

- Each ensemble member is usually updated using perturbed observations (though there are methods that can avoid this):

$$\mathbf{y}_i = \mathbf{y} + \boldsymbol{\varepsilon}_{o,i} \quad i = 1, 2, \dots, N_{ens}$$

Hybrid DA methods



Hybrid DA methods

- Data assimilation systems that have an ensemble data assimilation component used for the estimation of analysis and background errors are called **hybrid data assimilation systems** (Note: this definition is not universal!)
- We will discuss the various options for the ensemble data assimilation component in two dedicated lectures
- All major global NWP Centres run some form of hybrid data assimilation for atmospheric DA: a variational analysis cycle to estimate the mean/mode of the analysis pdf coupled with an ensemble data assimilation system to give a flow-dependent estimate of the second moments (covariances) of the error distributions.
- The ensemble DA component not only serves the purpose of estimating the background errors used in the analysis update, but it also provides a Monte Carlo sampling of the analysis pdf from which **ensemble forecasts** can be run

Summary

- Data assimilation in NWP aims to optimally blend information from observations and model to produce an accurate and physically consistent estimate of the initial state of the atmosphere and of the other components of the Earth System
- Both observations and models are affected by systematic and random errors: these need to be evaluated and taken into account in order to produce a statistically optimal analysis
- The Bayesian approach provides a unified theoretical framework for data assimilation
- Particle Filters provide a Monte Carlo implementation of the Bayes' Law in data assimilation. Asymptotically correct for $N_{ens} \rightarrow \infty$, but not (yet) applicable to high dimensional systems
- A Gaussian assumption on the error statistics is usually made to make the problem tractable in realistic geophysical DA
- Kalman Filter type methods and Variational methods can both be derived from Bayes' Law under these assumptions: they lead to the same solution for linear, Gaussian problems
- Hybrid data assimilation methods currently used in global NWP combine a variational analysis system with an ensemble data assimilation component for error estimation

Bibliography

- Ades M., Van Leeuwen P., 2014: The equivalent weights particle filter in a high dimensional system. *Q. J. R. Meteorol. Soc.*, 141: 484–503.
- Bocquet M., Pires C. A., Wu L., 2010: Beyond Gaussian Statistical Modeling in Geophysical Data Assimilation. *Mon. Wea. Rev.*, 138, 2997–3023, doi: 10.1175/2010MWR3164.1.
- Farchi, A. and M. Bocquet, 2018: Comparison of local particle filters and new implementations. *Nonlin. Processes Geophys.*, 25, 765–807. doi:10.5194/npg-25-765-2018
- Gordon, N.J., D.J. Salmond and A.F.M. Smith, 1993: A novel approach to nonlinear/non-Gaussian Bayesian state estimation. In: *IEE Proceedings on Radar and Signal Processing*. Vol. 140. pp. 107-113
- Miyoshi, T., K. Kondo, and T. Imamura, 2014: The 10,240-member ensemble Kalman filtering with an intermediate AGCM, *Geophys. Res. Lett.*, 41, 5264–5271, doi:10.1002/2014GL060863.
- Snyder, C., T. Bengtsson, and M. Morzfeld, 2015: Performance Bounds for Particle Filters Using the Optimal Proposal. *Mon. Wea. Rev.*, 143, 4750–4761, doi: 10.1175/MWR-D-15-0144.1.
- Wikle, C. K., and M. Berliner, 2007: A Bayesian tutorial for data assimilation. *Physica D*, 230, 1-16, doi:10.1016/j.physd.2006.09.017