

Data assimilation: the seamless integration of data into computational models

Jana de Wiljes

October 10, 2017

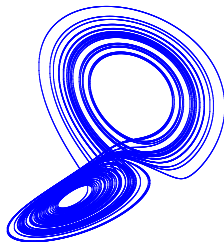
A toy atmospheric model

Lorenz equations:

$$\dot{x} = \sigma(y - x)$$

$$\dot{y} = x(\rho - z) - y$$

$$\dot{z} = xy - \beta z$$



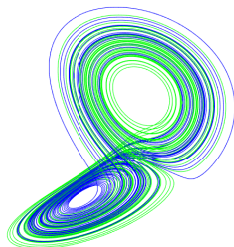
Uncertainty in initial conditions

Lorenz equations:

$$\dot{x} = \sigma(y - x)$$

$$\dot{y} = x(\rho - z) - y$$

$$\dot{z} = xy - \beta z$$



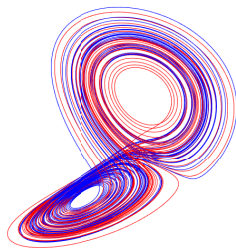
Uncertainty in parameters

Lorenz equations:

$$\dot{x} = \sigma(y - x)$$

$$\dot{y} = x(\rho - z) - y$$

$$\dot{z} = xy - \beta z$$



Numerical discretization and differentiation

Lorenz equations:

$$x_n = x_{n-1} + [\sigma(y_{n-1} - x_{n-1})]dt$$

$$y_n = y_{n-1} + [x_{n-1}(\rho - z_{n-1}) - x_{n-1}]dt$$

$$z_n = z_{n-1} + [x_{n-1}y_{n-1} - \beta z_{n-1}]dt$$

Model (deterministic)

Evolution equation

$$z_n = \Psi(z_{n-1}, \lambda)$$

where

$$z_0 \sim \mathcal{N}(m_0, C_0)$$

Model

Evolution equation

$$z_n = \Psi(z_{n-1}, \lambda) + \xi_{n-1}$$

where

$$z_0 \sim \mathcal{N}(m_0, C_0)$$

$$\xi_n \sim \mathcal{N}(0, B) \quad \text{i.i.d.} \quad \forall n$$

Parameter estimation

Augmented state space

$$\mathbf{z}_n = \Psi(\mathbf{z}_{n-1}, \lambda_{n-1}) + \xi_{n-1}$$

$$\lambda_n = \lambda_{n-1}$$

where

$$[\mathbf{z}_0, \lambda_0]^\top \sim \mathcal{N}(\mathbf{m}_0, \mathbf{C}_0)$$

$$\xi_n \sim \mathcal{N}(0, B) \quad \text{i.i.d.} \quad \forall n$$

L63 example

Augmented state space

$$z_n = \Psi(z_{n-1}, \lambda_{n-1}) + \xi_{n-1}$$

$$\lambda_n = \lambda_{n-1}$$

where

$$[z_0, \lambda_0]^\top \sim \mathcal{N}(m_0, C_0)$$

$$\xi_n \sim \mathcal{N}(0, B) \quad \text{i.i.d.} \quad \forall n$$

Observations

Partial and noisy data:

$$y_n = h(z_n) + \eta_n$$

where

$$\eta_n \sim \mathcal{N}(0, R) \quad \text{i.i.d.} \quad \forall n$$

Conditional probability

Definition (Conditional probability)

For a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ and events $A, B \in \mathcal{F}$ the conditional probability of B given A is

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A, B)}{\mathbb{P}(B)}.$$

Bayes theorem

Theorem (Bayes)

For a given probability space $(\Omega, \mathcal{F}, \mathbb{P})$ the following holds for two events A and B in \mathcal{F}

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(B|A)\mathbb{P}(A)}{\mathbb{P}(B)}.$$

Bayesian data assimilation ansatz

$$\mathbb{P}(\text{Model}|\text{Obs}) = \frac{\mathbb{P}(\text{Obs}|\text{Model})\mathbb{P}(\text{Model})}{\mathbb{P}(\text{Obs})}$$

Bayesian data assimilation ansatz

$$\mathbb{P}(\text{Model}|\text{Obs}) \propto \mathbb{P}(\text{Obs}|\text{Model})\mathbb{P}(\text{Model})$$

Bayesian data assimilation for densities

$$\begin{aligned}\pi(\mathbf{z}_{n+1}|\mathbf{y}_{1:n+1}) &= \pi(\mathbf{z}_{n+1}|\mathbf{y}_{1:n}, \mathbf{y}_{n+1}) \\ &= \frac{\pi(\mathbf{y}_{n+1}|\mathbf{y}_{1:n}, \mathbf{z}_{n+1})\pi(\mathbf{z}_{n+1}|\mathbf{y}_{1:n})}{\pi(\mathbf{y}_{n+1}|\mathbf{y}_{1:n})} \\ &= \frac{\pi(\mathbf{y}_{n+1}|\mathbf{z}_{n+1})\pi(\mathbf{z}_{n+1}|\mathbf{y}_{1:n})}{\pi(\mathbf{y}_{n+1}|\mathbf{y}_{1:n})}\end{aligned}$$

$$\implies \pi(\mathbf{z}_{n+1}|\mathbf{y}_{1:n+1}) \propto \pi(\mathbf{y}_{n+1}|\mathbf{z}_{n+1})\pi(\mathbf{z}_{n+1}|\mathbf{y}_{1:n}) \quad (1)$$

Special case

Linear model: Ψ is linear, e.g.,

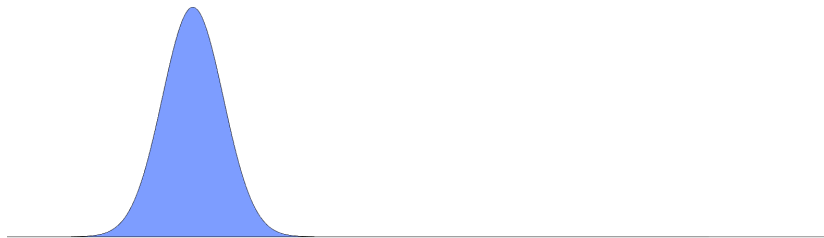
$$\mathbf{z}_n = A\mathbf{z}_{n-1} + \boldsymbol{\xi}_{n-1} \quad (2)$$

with $A \in \mathbb{R}^{N_z} \times \mathbb{R}^{N_z}$

Linear observation operator

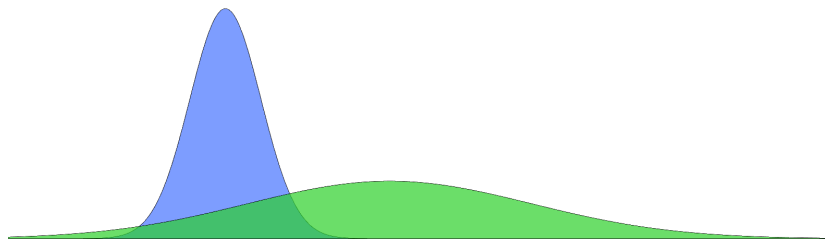
$$h = H \quad \text{with} \quad H \in \mathbb{R}^{N_y} \times \mathbb{R}^{N_y}$$

Linear Model



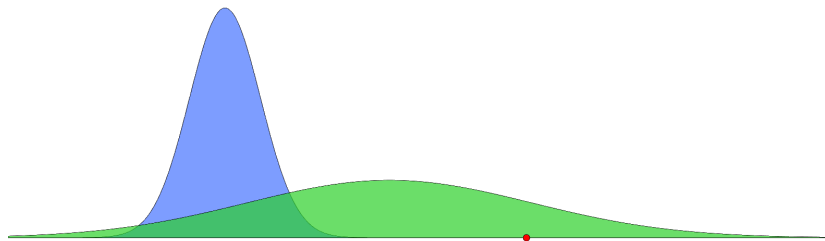
Initial distribution: $\mathbf{z}_0 \sim \mathcal{N}(\mathbf{m}_0, \mathbf{C}_0)$

Linear Model



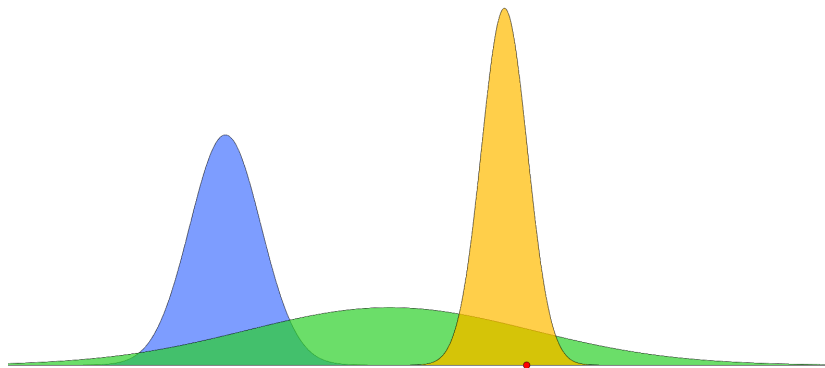
Prior distribution: $\mathcal{N}(\hat{m}_1, \hat{C}_1)$

Linear Model



Likelihood: $\mathcal{N}(\hat{z}_1, R)$

Linear Model



$$\text{Posterior: } \mathcal{N}(\mathbf{m}_1, \mathbf{C}_1) \propto \mathcal{N}(H\hat{\mathbf{z}}_1, R)\mathcal{N}(\hat{\mathbf{m}}_1, \hat{\mathbf{C}}_1)$$

Kalman filter

Two steps:

Kalman filter

Two steps:

- Forecast: $(m_n, C_n) \mapsto (\hat{m}_{n+1}, \hat{C}_{n+1})$

Kalman filter

Two steps:

- ▶ Forecast: $(m_n, C_n) \mapsto (\hat{m}_{n+1}, \hat{C}_{n+1})$
- ▶ Analysis: $(\hat{m}_{n+1}, \hat{C}_{n+1}) \mapsto (m_{n+1}, C_{n+1})$

Kalman filter

Two steps:

- ▶ Forecast: $(m_n, C_n) \mapsto (\hat{m}_{n+1}, \hat{C}_{n+1})$
- ▶ Analysis: $(\hat{m}_{n+1}, \hat{C}_{n+1}) \mapsto (m_{n+1}, C_{n+1})$

Forecast formulas

$$\hat{m}_{n+1} = A m_n$$

$$\hat{C}_{n+1} = A C_n A^T + B$$

Kalman filter

Two steps:

- Forecast: $(m_n, C_n) \mapsto (\hat{m}_{n+1}, \hat{C}_{n+1})$
- Analysis: $(\hat{m}_{n+1}, \hat{C}_{n+1}) \mapsto (m_{n+1}, C_{n+1})$

Forecast formulas

$$\hat{m}_{n+1} = A m_n$$

$$\hat{C}_{n+1} = A C_n A^\top + B$$

Analysis formulas

$$m_{n+1} = \hat{m}_{n+1} - K_{n+1}(H\hat{m}_{n+1} - y_{n+1})$$

$$C_{n+1} = \hat{C}_{n+1} - K_{n+1}H\hat{C}_{n+1}$$

Kalman gain

$$K_{n+1} = \hat{C}_{n+1}H^\top(R + H\hat{C}_{n+1}H^\top)^{-1}$$

Nonlinear Model

Problem: Kalman Filter is not applicable anymore

Ansatz: approximative Algorithms

1. **Extended Kalman Filter:** linearize model function

Nonlinear Model

Problem: Kalman Filter is not applicable anymore

Ansatz: approximative Algorithms

1. **Extended Kalman Filter:** linearize model function
2. **Monte Carlo Approximation:**

$$\pi(\mathbf{z}_n | \mathbf{y}_{1:n}) = \frac{1}{M} \sum_{i=1}^M \delta(\mathbf{z} - \mathbf{z}_n^i)$$

where

$$\mathbf{z}_n^i \sim \pi(\mathbf{z}_n | \mathbf{y}_n)$$

Nonlinear Model

Problem: Kalman Filter is not applicable anymore

Ansatz: approximative Algorithms

1. **Extended Kalman Filter:** linearize model function
2. **Monte Carlo Approximation:**

$$\pi(\mathbf{z}_n | \mathbf{y}_{1:n}) = \frac{1}{M} \sum_{i=1}^M \delta(\mathbf{z} - \mathbf{z}_n^i)$$

where

$$\mathbf{z}_n^i \sim \pi(\mathbf{z}_n | \mathbf{y}_n)$$

This ansatz leads to a variety of filters e.g.,

Nonlinear Model

Problem: Kalman Filter is not applicable anymore

Ansatz: approximative Algorithms

1. **Extended Kalman Filter:** linearize model function
2. **Monte Carlo Approximation:**

$$\pi(\mathbf{z}_n | \mathbf{y}_{1:n}) = \frac{1}{M} \sum_{i=1}^M \delta(\mathbf{z} - \mathbf{z}_n^i)$$

where

$$\mathbf{z}_n^i \sim \pi(\mathbf{z}_n | \mathbf{y}_n)$$

This ansatz leads to a variety of filters e.g.,

- ▶ Particle filters

Nonlinear Model

Problem: Kalman Filter is not applicable anymore

Ansatz: approximative Algorithms

1. **Extended Kalman Filter:** linearize model function
2. **Monte Carlo Approximation:**

$$\pi(\mathbf{z}_n | \mathbf{y}_{1:n}) = \frac{1}{M} \sum_{i=1}^M \delta(\mathbf{z} - \mathbf{z}_n^i)$$

where

$$\mathbf{z}_n^i \sim \pi(\mathbf{z}_n | \mathbf{y}_n)$$

This ansatz leads to a variety of filters e.g.,

- ▶ Particle filters (**curse of dimensionality**)

Nonlinear Model

Problem: Kalman Filter is not applicable anymore

Ansatz: approximative Algorithms

1. **Extended Kalman Filter:** linearize model function
2. **Monte Carlo Approximation:**

$$\pi(\mathbf{z}_n | \mathbf{y}_{1:n}) = \frac{1}{M} \sum_{i=1}^M \delta(\mathbf{z} - \mathbf{z}_n^i)$$

where

$$\mathbf{z}_n^i \sim \pi(\mathbf{z}_n | \mathbf{y}_n)$$

This ansatz leads to a variety of filters e.g.,

- ▶ Particle filters (**curse of dimensionality**)
- ▶ Ensemble Kalman filter

Nonlinear Model

Problem: Kalman Filter is not applicable anymore

Ansatz: approximative Algorithms

1. **Extended Kalman Filter:** linearize model function
2. **Monte Carlo Approximation:**

$$\pi(\mathbf{z}_n | \mathbf{y}_{1:n}) = \frac{1}{M} \sum_{i=1}^M \delta(\mathbf{z} - \mathbf{z}_n^i)$$

where

$$\mathbf{z}_n^i \sim \pi(\mathbf{z}_n | \mathbf{y}_n)$$

This ansatz leads to a variety of filters e.g.,

- ▶ Particle filters (**curse of dimensionality**)
- ▶ Ensemble Kalman filter(**underlying Gaussian assumption**)

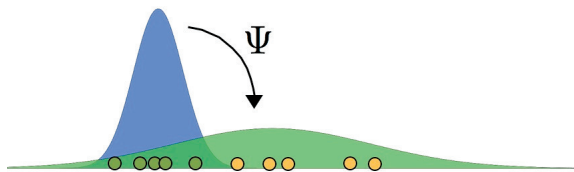
Ensemble Kalman Filter



$$\mathcal{N}(m_0, C_0) \text{ with } m_0 \approx \frac{1}{M} \sum_{i=1}^M z_0^i$$

$$C_0 \approx \frac{1}{M} \sum_{i=1}^M (z_0^i - m_0)(z_0^i - m_0)^\top$$

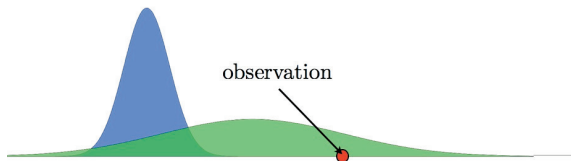
Ensemble Kalman Filter



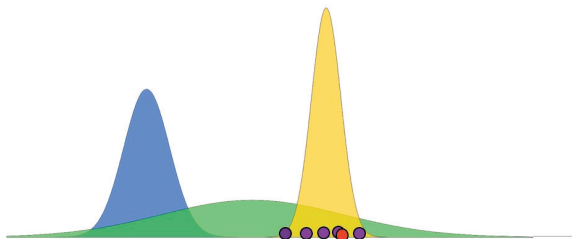
$$\mathcal{N}(\hat{\mathbf{m}}_1, \hat{\mathbf{C}}_1) \text{ with } \hat{\mathbf{m}}_1 \approx \frac{1}{M} \sum_{i=1}^M \hat{\mathbf{z}}_1^i = \frac{1}{M} \sum_{i=1}^M \Psi(\mathbf{z}_0^i)$$

$$\hat{\mathbf{C}}_1 \approx \frac{1}{M} \sum_{i=1}^M (\hat{\mathbf{z}}_1^i - \hat{\mathbf{m}}_1)(\hat{\mathbf{z}}_1^i - \hat{\mathbf{m}}_1)^\top$$

Ensemble Kalman Filter



Ensemble Kalman Filter



$$\mathcal{N}(m_1, C_1) \text{ with } m_1 \approx \frac{1}{M} \sum_{i=1}^M z_1^i$$

$$C_1 \approx \frac{1}{M} \sum_{i=1}^M (z_1^i - m_1)(z_1^i - m_1)^\top$$

Ensemble Kalman Filter

Approximation: $\pi(z_n | y_{1:n})$

Ensemble Kalman Filter

Approximation: $\pi(\mathbf{z}_n | \mathbf{y}_{1:n})$

Ansatz:: propagate samples $\hat{\mathbf{z}}_{n+1}^i$ with Kalman formula

$$\mathbf{z}_{n+1}^i = \hat{\mathbf{z}}_{n+1}^i - \mathbf{K}_{n+1}(\mathbf{H}\hat{\mathbf{z}}_{n+1}^i - \tilde{\mathbf{y}}_{n+1}^i)$$

Ensemble Kalman Filter

Approximation: $\pi(\mathbf{z}_n | \mathbf{y}_{1:n})$

Ansatz:: propagate samples $\hat{\mathbf{z}}_{n+1}^i$ with Kalman formula

$$\mathbf{z}_{n+1}^i = \hat{\mathbf{z}}_{n+1}^i - \mathbf{K}_{n+1}(\mathbf{H}\hat{\mathbf{z}}_{n+1}^i - \tilde{\mathbf{y}}_{n+1}^i)$$

Need:: perturbed observations

$$\tilde{\mathbf{y}}_{n+1}^i = \mathbf{y}_{n+1} + \boldsymbol{\epsilon}_{n+1}^i$$

with $\boldsymbol{\epsilon}_{n+1}^i \sim \mathcal{N}(0, R)$ i.i.d. to get the correct mean and covariance
in the linear case for $M \rightarrow \infty$

Ensemble Kalman Filter

Works well in practice: e.g., EnKF is used for operational NWP for \mathbf{z}_n^i with dimension 10^9 only using $M = 100$

Yet: mathematical foundation largely missing

Recent study: accuracy results for EnKF for idealized setting: $H = Id$ and observational error small

Particle Filter

Problem: sampling from $\pi(\mathbf{z}_n | \mathbf{y}_{1:n})$ to approximate posterior via

$$\pi(\mathbf{z}_n | \mathbf{y}_{1:n}) = \frac{1}{M} \sum_{i=1}^M \delta(\mathbf{z} - \mathbf{z}_n^i)$$

Particle Filter

Problem: sampling from $\pi(\mathbf{z}_n | \mathbf{y}_{1:n})$ to approximate posterior via

$$\pi(\mathbf{z}_n | \mathbf{y}_{1:n}) = \frac{1}{M} \sum_{i=1}^M \delta(\mathbf{z} - \mathbf{z}_n^i)$$

Idea: sampling from $\pi(\mathbf{z}_n | \mathbf{y}_{1:n-1})$ instead i.e.,

$$\pi(\mathbf{z}_n | \mathbf{y}_{1:n}) = \sum_{i=1}^M w_n^i \delta(\mathbf{z} - \hat{\mathbf{z}}_n^i)$$

Particle Filter

Problem: sampling from $\pi(\mathbf{z}_n | \mathbf{y}_{1:n})$ to approximate posterior via

$$\pi(\mathbf{z}_n | \mathbf{y}_{1:n}) = \frac{1}{M} \sum_{i=1}^M \delta(\mathbf{z} - \mathbf{z}_n^i)$$

Idea: sampling from $\pi(\mathbf{z}_n | \mathbf{y}_{1:n-1})$ instead i.e.,

$$\pi(\mathbf{z}_n | \mathbf{y}_{1:n}) = \sum_{i=1}^M w_n^i \delta(\mathbf{z} - \hat{\mathbf{z}}_n^i)$$

Bayes:

$$\pi(\mathbf{z}_{n+1} | \mathbf{y}_{1:n}) \propto \pi(\mathbf{y}_n | \mathbf{z}_n) \pi(\mathbf{z}_n | \mathbf{y}_{1:n-1}) \quad (3)$$

Particle Filter

Problem: sampling from $\pi(\mathbf{z}_n | \mathbf{y}_{1:n})$ to approximate posterior via

$$\pi(\mathbf{z}_n | \mathbf{y}_{1:n}) = \frac{1}{M} \sum_{i=1}^M \delta(\mathbf{z} - \mathbf{z}_n^i)$$

Idea: sampling from $\pi(\mathbf{z}_n | \mathbf{y}_{1:n-1})$ instead i.e.,

$$\pi(\mathbf{z}_n | \mathbf{y}_{1:n}) = \sum_{i=1}^M w_n^i \delta(\mathbf{z} - \hat{\mathbf{z}}_n^i)$$

Bayes:

$$\pi(\mathbf{z}_{n+1} | \mathbf{y}_{1:n}) \propto \pi(\mathbf{y}_n | \mathbf{z}_n) \pi(\mathbf{z}_n | \mathbf{y}_{1:n-1}) \quad (3)$$

Weighting: unnormalized weights

$$\tilde{w}_n^i = \pi(\mathbf{y}_n | \mathbf{z}_n^i) w_n^i \text{ with } w_0^i = \frac{1}{M}$$

and normalized weights

$$w_n^i = \frac{\tilde{w}_n^i}{\sum_{j=1}^M \tilde{w}_n^j}$$

Resampling

Problem: weights w_n^i become very small

Ansatz: resampling

Input: w_n^i

For ($k = 1 : M$)

1. Draw a number $u \in [0, 1]$ from the uniform distribution $U[0, 1]$
2. Compute $i^* \in \{1, \dots, M\}$ which satisfies

$$i^* = \arg \min_{i \geq 1} \sum_{j=1}^i w_j \geq u \quad (4)$$

3. Set $\xi_{i^*} = \xi_i^* + 1$

Return ξ_i