

# Parameter Estimation in Size-Structured Aerosol Populations using Bayesian State Estimation

Matthew Ozon<sup>a</sup>, Aku Seppänen<sup>a</sup>, Jari Kaipio<sup>a,b</sup>, & Kari Lehtinen<sup>a,c</sup>

<sup>a</sup>Department of Applied Physics, University of Eastern Finland, Kuopio, Finland

<sup>b</sup>Department of Mathematics, Faculty of Science, University of Auckland, New Zealand

<sup>c</sup>Finnish Meteorological Institute, Kuopio, Finland



UNIVERSITY OF  
EASTERN FINLAND

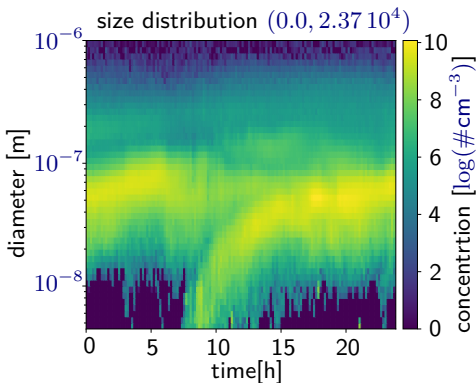


SCIENCE  
DEPARTMENT OF MATHEMATICS



FINNISH  
METEOROLOGICAL  
INSTITUTE

Conservation Principles, Data and Uncertainty in Atmosphere-Ocean  
Modelling 2019, Potsdam, April 3

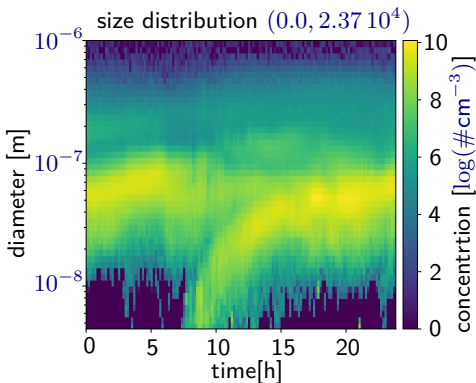


Goal: estimation of the characteristic parameters of aerosol systems along with their uncertainties, e.g.

- condensational growth rate
- nucleation rate
- linear loss rates

Why? Per se, because it is interesting, but also it is relevant, e.g. aerosols play a key role in cloud formation. Also, aerosols are known to disturb

- satellite data
- climate model (IPCC reports)
- radiative forcing



Evolution model of the form:

$$\frac{\partial u}{\partial t} = F(t, s, u; \theta)$$

If it is so important, methods should already exist, right? Well, it's not that simple! Indeed, some methods exist for estimating some parameters... but they are not satisfactory!

- manual recipes,
- no uncertainty estimation.

How to estimate the parameters? We have:

- data: time series of number concentration
- evolution model: GDE...

The perfect framework for an inverse problem.

How to describe the evolution of a population of aerosols whose characteristics depend on size? Size-structured Population Balance Equation (PBE).

$$\frac{\partial u}{\partial t} + \frac{\partial gu}{\partial s} = F(t, s, u; \theta) \quad (1)$$

This PBE is a scalar conservation law where  $u$  is the particle size density,  $g$  the growth rate and  $F$  is the term that describes the mechanisms that make the density evolve — it depends on some parameter  $\theta$ . We refer to this equation as the General Dynamic Equation for aerosols, or simply GDE.

**Note:** without going into details, each aerosol particle cannot be described only by a size. By nature, particles are complex objects of different shape, size and chemical composition, but we model them as spherical objects of equivalent volume without considering chemistry. Therefore, it cannot be a perfectly precise model (source of uncertainty, potentially modelled by SPDE).

## Condensational growth

Ambient vapor condenses onto the surface of particles, resulting in the growth of particles. Represented in the left hand side of equation (1) by  $g$  which is a speed.

## Nucleation (and sources)

New particles can be added to an aerosol system either by adding existing particle or by formation of particles from the ambient vapor i.e. nucleation. In either case:

$$F(t, s, u; J) = J \quad (1)$$

is a simple source term.

## Linear losses

Particles can be removed from the system by various mechanism such as sedimentation or wall deposition. It is well described by a linear damping term:

$$F(t, s, u; \lambda) = \lambda(t, v)u(t, v) \quad (2)$$

## Coagulation

Particle of a given size may be created or removed by coagulation — when two particles collide and stick, they form a new particle.

$$F(t, v, u; \beta) = \underbrace{\int_{v_0}^{v-v_0} \beta_v(s, v-s) u(t, s) u(t, v-s) ds}_{\text{coagulation source}} - \underbrace{u(t, v) \int_{v_0}^{\infty} \beta_v(v, s) u(t, s) ds}_{\text{coagulation sink}}$$

where  $\beta_v$  is a collision frequency factor.

The continuous form of the GDE is not really suitable for our purpose (parameter estimation from time series of number concentration), hence we define:

$$\forall i \in \llbracket 0, T \rrbracket, k \in \llbracket 1, K \rrbracket, \quad N_i^k = \int_{\Omega_i} u(s, k\Delta_t) ds \quad (1)$$

the number of particles in the size range  $\Omega_i$  per unit of volume. Considering a logarithmic scale, and using the Euler time discretization scheme, we obtain the following time-and-size discrete evolution equations:

$$N_1^{k+1} = N_1^k + \Delta_t^k \left( J^k - \left( \frac{g_1^k}{\Delta_1} + \lambda_1 \right) N_1^k - C_1^{\text{sink}}(N^k) N_1^k \right) + \varepsilon_1^k \quad (2)$$

$$\begin{aligned} N_i^{k+1} = N_i^k + \Delta_t^k \left( \frac{g_{i-1}^k}{\Delta_{i-1}} N_{i-1}^k - \left( \frac{g_i^k}{\Delta_i} + \lambda_i \right) N_i^k \right. \\ \left. + C_i^{\text{source}}(N^k) - C_i^{\text{sink}}(N^k) N_i^k \right) + \varepsilon_i^k \end{aligned} \quad (3)$$

**Note** that both discretization steps — time and size — add errors. The overall errors/uncertainties are encompassed in the terms  $\varepsilon_i^k$ .

**Note** that I leave out the parameter evolution for now.

## DMA

This device acts as a selector of near monodisperse size distribution around a given size  $d_i$ ; its size discrimination power determines the sets  $\{d_i\}_{i \in \llbracket 1, N \rrbracket}$  and  $\{\Delta_i\}_{i \in \llbracket 1, N \rrbracket}$ . For each channel, we denote the time invariant kernel  $\psi_i$ , which models the efficiency of the device. The number concentration at the outlet of the DMA is approximated by:

$$z_i^k = \frac{1}{\Delta_t} \int_{t_0 + (k-1)\Delta_t}^{t_0 + k\Delta_t} \int_{\omega_i} \psi_i(s) u(s, t) ds dt + \iota_i^k = \varphi_i^k + \iota_i^k \quad (4)$$

where  $\omega_i$  is the support of  $\psi_i$  — where it is not null — and  $\iota_i^k$  accounts for the model uncertainties. The number of particle is then obtained by multiplying by the volume.



## CPC

Let the number concentration of particles at the inlet of a CPC is  $z_i^k$  (coming from the  $i^{\text{th}}$  channel of the DMA), the output  $y_i^k$  is modeled by:

$$y_i^k = \frac{\tilde{y}_i^k}{V}, \quad \text{with} \quad \tilde{y}_i^k \sim \text{Poisson}(V z_i^k) \quad (4)$$

where  $V$  is the volume of sample used in the CPC for counting. In most cases, the number of particle in the CPC is large enough ( $V z_i^k > 20$ ), thus the Poisson distribution can satisfactorily be approximated by:

$$y_i^k \sim \mathcal{N}(z_i^k, \frac{z_i^k}{V}). \quad (5)$$

Note that from this model it is clear that the quality of the measurement is directly link to the volume of the sample.

## SMPS

The full measurement device can be summarized by the model:

$$y_i^k = \frac{\tilde{y}_i^k}{V}, \quad \text{with} \quad \tilde{y}_i^k \sim \mathcal{Poisson}(V(\varphi_i^k + \iota_i^k)). \quad (4)$$

Assuming that the device operates under normal conditions — that is it is actually counting something — and that the DMA model has no flaw, then the model becomes:

$$y_i^k = \varphi_i^k + \frac{1}{V}\tilde{\iota}_i^k, \quad \text{with} \quad \tilde{\iota}_i^k \sim \mathcal{N}(0, V\varphi_i^k). \quad (5)$$

Note that the  $N_i^k$ 's used in the GDE correspond to the case:

$$\psi_i(s) = \begin{cases} \frac{1}{|\omega_i|} & \text{if } s \in \omega_i \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

For the sake of clarity, we consider the above case, hence, the measurement model is:

$$y_i^k = N_i^k + \frac{1}{V}\tilde{\iota}_i^k \quad \text{with} \quad \tilde{\iota}_i^k \sim \mathcal{N}(0, VN_i^k) \quad (7)$$

## Kalman Filter (KF)

Estimation of the expected state and its uncertainty

$$X^{k|k} = \mathbb{E}[X^k | \mathcal{Y}_k] \quad \text{and} \quad \Gamma^{k|k} = \text{Cov}[X^k | \mathcal{Y}_k], \quad (8)$$

**Initialization**

Set  $X^{0|0} \in \mathbb{R}^N$  and  $\Gamma^{0|0} \in \mathbb{R}^{N \times N}$  (Prior knowledge)

**Recursion**

While  $k \leq K$ , do

*Prediction*

$$X^{k|k-1} = F^{k-1} (X^{k-1|k-1}) \quad (\text{state expectation})$$

$$\Gamma^{k|k-1} = \partial F^{k-1} \Gamma^{k-1|k-1} (\partial F^{k-1})^T + \Gamma_w^{k-1} \quad (\text{state covariance})$$

*Calculation of Kalman's gain*

$$K^k = \Gamma^{k|k-1} (H^k)^T (H^k \Gamma^{k|k-1} (H^k)^T + \Gamma_v^k)^{-1}$$

*Filtering*

$$\Gamma^{k|k} = (I - K^k H^k) \Gamma^{k|k-1}$$

$$X^{k|k} = X^{k|k-1} + K^k (y^k - H^k X^{k|k-1})$$

*Update iterator*

$$k \leftarrow k + 1$$

end(while)

## Fixed Interval Kalman Smoother (FIKS)

$$X^{k|K} = \mathbb{E}[X^k | \mathcal{Y}_K] \quad \text{and} \quad \Gamma^{k|K} = \text{Cov}[X^k | \mathcal{Y}_K], \quad (8)$$

**Initialization**

Run KF and store all variables

Set  $X_{smo}^K = X^{K|K}$  and  $\Gamma_{smo}^K = \Gamma^{K|K}$

$k \leftarrow K - 1$

**Recursion**

While  $k \geq 1$ , do

    Compute smoothing gain

$$K_{smo}^k = \Gamma^{k|k} (\partial F^{k+1})^T (\Gamma^{k+1|k})^{-1}$$

    Smoothing

$$X_{smo}^k = X^{k|k} + K_{smo}^k (X_{smo}^{k+1} - X^{k+1|k})$$

$$\Gamma_{smo}^k = \Gamma^{k|k} + K_{smo}^k (\Gamma_{smo}^{k+1} - \Gamma^{k+1|k}) (K_{smo}^k)^T$$

    Update iterator

$$k \leftarrow k - 1$$

end(while)

We have almost all the elements to run the algorithm:

- Evolution model of the number concentration
- Measurement model
- Data: simulation or measurement
- Initial guesses: may depend on the user

however, we still miss one part:

- Evolution model of the parameters

We'll have to create some based on what we know of the system.

**Time invariant** Some parameter are time invariant, so their time evolution model is a random walk:

$$k \geq 1, \quad p^{k+1} = p^k + \eta^k, \quad \text{with} \quad \eta^k \sim \mathcal{N}(0, \Gamma_\eta) \quad (8)$$

where  $\Gamma_\eta$  is the covariance of the model uncertainty.

**Second order** If a parameter is known to evolve smoothly with time, it can be modelled as a second order stochastic process such as:

$$G^k = \begin{bmatrix} p^k \\ p^{k-1} \end{bmatrix} = \begin{bmatrix} 2r_p & -r_p^2 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} p^{k-1} \\ p^{k-2} \end{bmatrix} + \begin{bmatrix} \eta^k \\ 0 \end{bmatrix} = B(r_p)G^{k-1} + \begin{bmatrix} \eta^k \\ 0 \end{bmatrix} \quad (8)$$

where  $r_p$  is the smoothness lever and  $\eta^k \sim \mathcal{N}(0, \sigma_\eta^2)$  with  $\sigma_\eta$  controlling the amplitude of the process. The latter is given by its covariance matrix defined by:

$$\Gamma_G^k = \text{cov}(G^k) = B(r_p)\Gamma_G^{k-1}B(r_p)^T + \begin{bmatrix} \sigma_\eta^2 & 0 \\ 0 & 0 \end{bmatrix}. \quad (9)$$

We are interested in the asymptotic behavior, i. e.  $\Gamma_G^k = \Gamma_G^{k-1} = \Gamma_G^\infty = \begin{bmatrix} \sigma_p^2 & c \\ c & \sigma_p^2 \end{bmatrix}$ , and how the variance of  $\eta$  controls the variance of  $p$ ,  $\sigma_p^2$ . By expanding the previous relation, we find that:

$$\sigma_\eta^2 = \sigma_p^2 \left( 1 - r_p^2 \left( 4 + r_p^2 \left( 1 - \frac{8}{1 + r_p^2} \right) \right) \right), \quad c = \frac{2r_p}{1 + r_p^2} \sigma_p^2. \quad (10)$$

**Size correlation** Some parameters are distributed, yet, the size dependence may be unknown or only approximately known. For most parameter, it is safe to assume that size dependence is continuous, and even rather smooth.

Let  $p^k \in \mathbb{R}^N$  follow a random walk described by eq. (8), the covariance  $\Gamma_\eta$  convey the size dependence information, and it can be constructed as:

$$\Gamma_\eta = \tilde{D}^{\frac{1}{2}} \bar{D}^{-\frac{1}{2}} \bar{\Gamma} \bar{D}^{-\frac{1}{2}} \tilde{D}^{\frac{1}{2}} \quad (8)$$

where  $\bar{\Gamma}$  is the Toeplitz matrix build with the sequence  $(\bar{\sigma}_i)_{i \in \llbracket 1, N \rrbracket}$  which determines how the size dependence evolves with the size difference. We choose the sequence

$$\forall i \in \llbracket 1, N \rrbracket, \quad \bar{\sigma}_i = e^{\frac{1-i}{\delta}} \quad (9)$$

with  $\delta$  so that only the first  $\delta$  neighboring sizes significantly contribute to the evolution of one variable. The decay is exponential in size index, so it is actually linearly/polynomially decaying in the diameter space. The diagonal matrix  $\bar{D} = \bar{\sigma}_1^2 I_N$  normalize the covariance  $\bar{\Gamma}$  and the diagonal matrix  $\tilde{D} = \text{diag}([\sigma_{\eta,1}^2 \sigma_{\eta,2}^2 \dots \sigma_{\eta,N}^2])$  scales the variance of each size.



What if we know the possible range of a parameter?

Lower bound

$$p = a + \frac{1}{\alpha} \log(1 + e^{\alpha\zeta}), \quad (10)$$

Range

$$p = a + \frac{b - a}{1 + \frac{1}{\alpha} e^{-\alpha\zeta}} \quad (11)$$

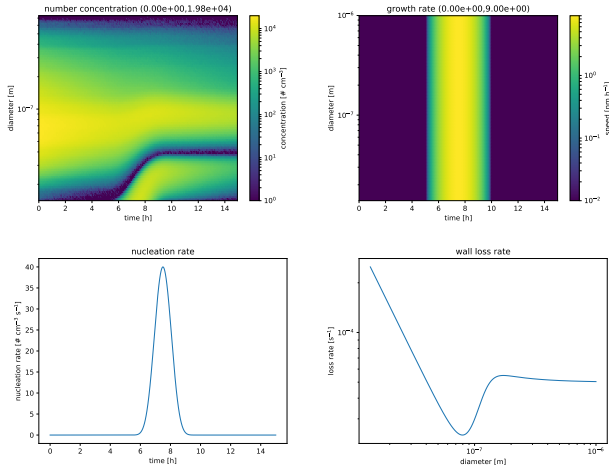
We use simulated data in order to evaluate the performance of the method.

- dense discretization of the size space
- non-approximated measurement model (Gaussian kernel and Poisson noise)

The dense discretization of the size space may lead to spurious oscillation (or diverge) if the following condition on the time step is not met:

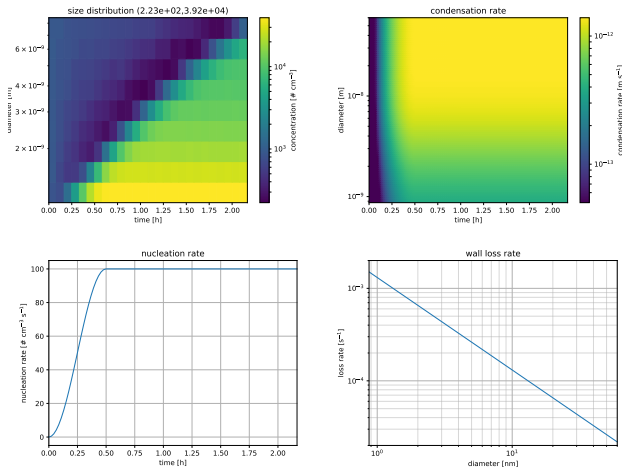
$$0 < \Delta_t^k < \frac{1}{\max_i \left\{ \frac{g_i^k}{\Delta_i} + \lambda_i + C_i^{\text{sink}}(N^k) \right\}}. \quad (12)$$

## Nucleation event



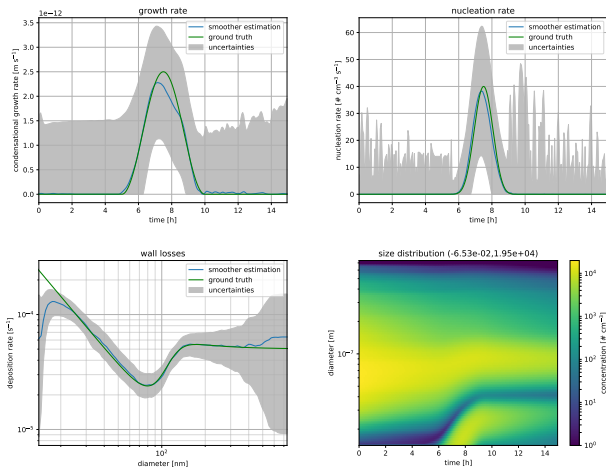
**Figure:** Simulation of a nucleation event. a) Number concentration contour plot, b) growth rate, c) nucleation rate at 14.1nm, and d) wall loss rate.

## CLOUD simulation



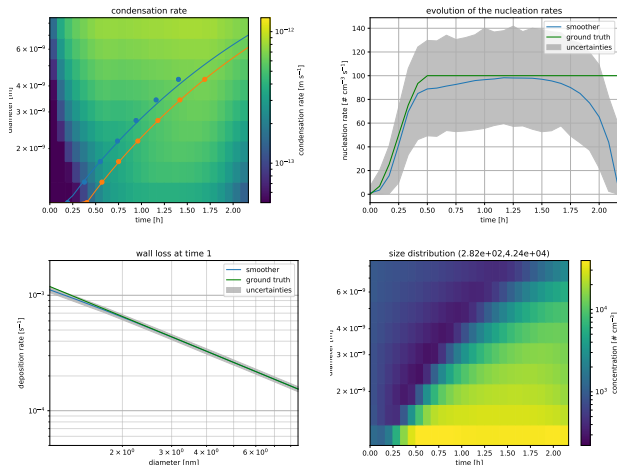
**Figure:** Simulation of a steady state. 1) Number concentration contour plot, 2) growth rate, 3) nucleation rate, and 4) wall loss rate.

# Nucleation event



**Figure:** Estimation of the parameter of the GDE for aerosols from a simulated nucleation event data: 1) growth rate, 2) nucleation rate, 3) loss rate, and 4) number concentration.

## CLOUD simulation



**Figure:** Estimation of the parameter of the GDE for aerosols from a simulated transition to steady state data: 1) growth rate, 2) nucleation rate, 3) loss rate, and 4) number concentration.

Message to take back home:

- Aerosols can disturb everything... at least the models
- The Fixed Interval Kalman Smoother is a suitable tool for the estimation of the GDE parameters (distributed or not) along with meaningful uncertainties
- The requirements for applying the method are “weak”: 1) the model must be well approximated by their Jacobian, and 2) the errors can be approximated as gaussian
- Need surrogate evolution models for the parameters of interest (unless someone comes up with a physically relevant evolution model).

# Thanks!