Atmos. Meas. Tech., 9, 741–751, 2016 www.atmos-meas-tech.net/9/741/2016/ doi:10.5194/amt-9-741-2016 © Author(s) 2016. CC Attribution 3.0 License.





# Notably improved inversion of differential mobility particle sizer data obtained under conditions of fluctuating particle number concentrations

Bjarke Mølgaard<sup>1</sup>, Jarno Vanhatalo<sup>2</sup>, Pasi P. Aalto<sup>1</sup>, Nønne L. Prisle<sup>1</sup>, and Kaarle Hämeri<sup>1</sup>

<sup>1</sup>Department of Physics, University of Helsinki, Helsinki, Finland <sup>2</sup>Department of Environmental Sciences, University of Helsinki, Helsinki, Finland

Correspondence to: Bjarke Mølgaard (bjarke.molgaard@helsinki.fi)

Received: 9 June 2015 – Published in Atmos. Meas. Tech. Discuss.: 7 October 2015 Revised: 20 January 2016 – Accepted: 2 February 2016 – Published: 29 February 2016

Abstract. The differential mobility particle sizer (DMPS) is designed for measurements of particle number size distributions. It performs a number of measurements while scanning over different particle sizes. A standard assumption in the data-processing (inversion) algorithm is that the size distribution remains the same throughout each scan. For a DMPS deployed in an urban area this assumption is likely to be violated most of the time, and the resulting size distribution data are unreliable. To improve the reliability, we developed a new algorithm using a statistical model in which the problematic assumption was replaced with more realistic smoothness assumptions, which were expressed through Gaussian process prior probabilities. We tested the model with data from a twin DMPS located at an urban background site in Helsinki and found that it provides size distribution data which are much more realistic. Furthermore, particle number concentrations extracted from the DMPS data were compared with data from a condensation particle counter. At 10 min resolution, the correlation for a period of 10 days was 0.984 with the new algorithm and 0.967 with the old one. Moreover, the time resolution was improved, and at 30 s resolution we obtained positive correlations for 89% of the scans. Thus, the quality of the inverted data was clearly improved.

## 1 Introduction

There is no direct way of measuring the size distribution of fine particles. To get information on the size distribution, mobility particle size spectrometers (Wiedensohler et al., 2012) select in turns particles of various electrical mobilities, and for each electrical mobility the number of particles in some volume is counted. To obtain the size distributions, the dependence of electrical mobilities on particle sizes is utilised. However, the electrical mobility depends also on the particle charge, so various combinations of particle size and charge give the same electrical mobility, and the inference of the actual particle size distribution is not trivial. The algorithms which have been developed for this purpose are generally known as inversion algorithms. The task can be split into two parts: determination of the transfer function, which gives the detection probabilities of particles in the sampled air, and the actual inversion. In this study we focused on the latter part. Further we restrict ourselves to considering differential mobility particle sizers (DMPSs) which differ from scanning mobility particle sizers (SMPSs) by changing the selected electrical mobility in discrete steps.

A DMPS typically performs a few tens of measurements while scanning over a wide range of electrical mobilities. Each measurement takes several seconds, and a typical waiting time between the measurements is around 10 s. Thus, each scan takes 5 to 10 min. The inversion must be based on some assumption about the time evolution of the aerosol during the scan. The simplest and most commonly used assumption is that the particle size distribution stays constant during each scan. In remote areas, where the particle size dis-



**Figure 1.** Size distributions from 2 March 2015 according to the old inversion algorithm, which assumes stationary particle size distributions during each scan.

tribution changes slowly most of the time, this assumption is reasonable most of the time. However, when a DMPS is deployed in a city with numerous nearby sources and a turbulent wind flow, this assumption is often far from the truth, and many of the derived particle size distributions are unreliable. Let us illustrate this with an example using inverted data from the twin DMPS (comprises two DMPSs) at the SMEAR III (Station for Measuring Ecosystem-Atmosphere Relations) station in Helsinki (Järvi et al., 2009) on 2 March 2015. The particle size distribution fluctuated substantially during daytime (Fig. 1), but at night-time there were periods without fluctuations. For these night-time periods, the particle size distribution was a rather smooth function of size, but in daytime unrealistically narrow peaks are present in the inverted data. In particular, the scan from 11:00 to 11:10 UTC +2 h is problematic, because of the strong, narrow peaks at 13 and 30 nm (Fig. 2). Single-charged particles of these two sizes are measured almost simultaneously in the two DMPSs, so the peaks at these sizes are most likely caused by a brief concentration peak. At 26 nm  $dN/d\log_{10}D_p$  appears to be low (Fig. 2; see also the light-blue spot close to the middle of Fig. 1), although the raw data indicate that the concentration was already elevated when the DMPS measured particles at this size. However, a fraction of the particles in the 30 nm bin will be detected in the measurement centred at 26 nm, and, assuming a stationary particle size distribution, the 30 nm particles were very abundant also during this measurement. Thus, assuming a stationary particle number size distribution, most of the particles detected in this measurement belonged to the 30 nm bin. Additionally, double-charged 39 nm particles contributed somewhat to the particle count. As a result, a small value was assigned to  $dN/d\log_{10}D_p$  at 26 nm.



**Figure 2.** Three particle number size distributions from 2 March 2015 according to the old inversion algorithm, which assumes stationary particle size distributions during each scan.

A few studies have addressed this issue of possible size distribution changes happening during a scan. Voutilainen and Kaipio (2001) presented an algorithm based on the Kalman filter. They let the particle size distribution change in discrete steps at each measurement based on the observation and a random walker. Subsequently, a smoother and a non-negativity constraint were applied. The algorithm was applied to synthetic data, and it reproduced a slowly varying size distribution well. Voutilainen and Kaipio (2002, 2005) parametrised the size distribution and replaced the random walker by estimations of the time evolution based on an aerosol model which took coagulation and condensation into account. The underlying assumption is that the DMPS continuously samples from the same aerosol, which changes in time due to these processes. Also this assumption is generally invalid in a city. Although the algorithm by Voutilainen and Kaipio (2005) was also shown to adjust to abrupt changes in the aerosol within a couple of minutes, it was not designed for use in urban locations.

In this work, we developed a new inversion algorithm for processing DMPS data from locations with fluctuating particle number concentrations. The particle number size distribution was modelled as a function of time and particle size using a Gaussian process (GP) model (Rasmussen and Williams, 2006). Assumptions of smoothness in both dimensions were incorporated through a GP prior. We tested our new algorithm with data from the twin DMPS at the urban background station SMEAR III in Helsinki. For periods with considerable fluctuations, the time resolution and the reliability of the derived particle number size distribution were substantially improved. A demo version of our algorithm is provided in the Supplement.

#### B. Mølgaard et al.: Notably improved DMPS data inversion

#### 2 Methods

#### 2.1 Quantification of particle size distributions

Particle size distributions are usually described by the  $dN/dlog_{10}D_p$ , although it is not a distribution function in a mathematical sense. *N* is the product of the total particle number concentration and the cumulative distribution function of the particle diameters (Hinds, 1999). Thus,  $dN/dlog_{10}D_p$  is the product of the concentration and the probability density function on the logarithmic scale.

## 2.2 DMPS

The DMPS comprises a neutraliser (bi-polar charger), a differential mobility analyser (DMA), and a condensation particle counter (CPC). In the neutraliser, ionising radiation ensures that the particles in the sampled air reach the equilibrium charge distribution. This charge distribution is known and depends on the particle size. In the DMA, the voltage and airflow are adjusted to select particles with a certain electrical mobility  $Z = \frac{qC_c}{3\pi\eta D_p}$ , where q is the charge;  $D_p$  is the mobility (Stokes') diameter;  $C_c$  is the Cunningham slip correction factor, which depends on  $D_p$ ; and  $\eta$  is the dynamic viscosity of the air. So Z depends on two particle properties:  $D_p$  and q = ze, where z is an integer and e is the elementary charge. The particles selected by the DMA flow to the CPC which counts them. Usually, the flows are kept constant and the DMPS scans over a few tens of discrete DMA voltages in order to select particles of different electrical mobilities. At each of these voltages, a measurement is performed with the CPC.

#### 2.3 Transfer function

The transfer function  $T(D_p, U)$  is defined as the probability that a particle of diameter  $D_p$  will be detected in the CPC when the DMA voltage is U. This probability is the product of the following three probabilities:

- *P*<sub>DMA</sub>, the probability that the particle is selected by the DMA (Stolzenburg, 1988; Mamakos et al., 2007);
- *P*<sub>Pen</sub>, the probability that the particle penetrates all sampling lines without being deposited (Wiedensohler et al., 2012);
- $P_{CPC}$ , the detection probability for particles reaching the CPC.

The DMA is designed to select particles with electric mobilities in a narrow band. The electric mobility depends on  $D_p$  and the number z of charges. The probability of a certain number z of charges depends on the particle size, i.e. on  $D_p$ . Thus,  $P_{\text{DMA}}$  can be described as a function of  $D_p$  and the DMA voltage U. Diffusion is the main reason for deposition of particles in the sampling lines, and the particle diffusivity depends on the particle size, so  $P_{\text{Pen}}$  is also a function of  $D_p$ .  $P_{CPC}$  is close to unity (1) for most particles, but for the smallest particles it is lower. During a measurement *i* the DMA voltage is kept at a constant value  $U_i$ . We will define  $T_i(D_p) = T(D_p, U_i)$ .  $T_i$  has a few clear peaks and is zero for the rest of the interval. The largest peak is for the diameter which gives the selected electrical mobility Z when the number of charges z equals 1 or -1. The sign depends on the polarity of the DMA voltage. A second peak is observed for the diameter  $D_p$ , which gives the same electrical mobility for a particle with double charge. This peak is smaller than the first one, because particles are less likely to carry two charges than one charge. Triple-charged particles cause a third peak, which is smaller than the second peak, and subsequent peaks are even smaller.

#### 2.4 Inversion algorithm based on a GP model

A GP, or Gaussian random field, is a stochastic process that can be used to define probability distributions over functions, and it is a generalisation of the multivariate normal (Gaussian) distribution (O'Hagan, 1978; Rasmussen and Williams, 2006). It is defined by a mean and a covariance function, which determine the properties, such as the smoothness and variability of the process. GPs are widely used for interpolation and to model coloured (spatially correlated) noise in spatial statistics (Gelfand et al., 2010), and they have obtained increasing interest also in, e.g., statistics and machine learning due their good interpolation and smoothing properties as well as convenient marginalisation and conditioning properties (see Sect. 2.5.2) (Rasmussen and Williams, 2006; Vanhatalo et al., 2010).

Let us define a latent function  $f(t, u) = \log(dN/du)$ , where t is time and  $u = \log_{10}D_p$ . We use the Bayesian formalism and express our prior belief about its smoothness through a GP prior. Thus, the posterior, which is proportional to the product of the prior and the likelihood, is a probability distribution of f.

# 2.4.1 Likelihood function

Each measurement *i* provides a count  $y_i$  of particles, so the Poisson distribution  $\text{Poi}(y_i|\lambda_i) = \exp(-\lambda_i)\frac{\lambda_i^{y_i}}{y_i!}$  is suitable for each factor of the likelihood

$$p(\mathbf{y}|f) = \prod_{i} \exp(-\lambda_i) \frac{\lambda_i^{y_i}}{y_i!}.$$
 (1)

The rate parameter  $\lambda_i$  is obtained by integrating the product of  $dN/d\log_{10}D_p$  and the transfer function  $T_i$ . In terms of f and u defined above,

$$\lambda_i = Q \int_{-\infty}^{\infty} \int_{t_{i,\text{begin}}}^{t_{i,\text{end}}} \exp(f(t,u)) T_i(10^u) dt du, \qquad (2)$$

where Q is the sample flow rate and  $t_{i,\text{begin}}$  and  $t_{i,\text{end}}$  are the beginning and end times of measurement i. If f(t, u)

does not fluctuate much during each measurement (which has much shorter duration than a scan), the following approximation holds:

$$\lambda_i \approx V_i \int_{-\infty}^{\infty} \exp(f(t_i, u)) T_i(10^u) \mathrm{d}u, \qquad (3)$$

where  $V_i$  is the volume of sampled air and  $t_i$  is the middle of the time interval between  $t_{i,\text{begin}}$  and  $t_{i,\text{end}}$ . Because  $T_i$  has a few clear peaks and is zero for the rest of the interval, the integral can be well approximated by a sum. Let  $T_{i,j}$  equal the integral of  $T_i$  over the peak centred at size  $u_{i,j}$ , and let  $f_{i,j} = f(t_i, u_{i,j})$ . Then

$$\lambda_i \approx V_i \sum_j \exp(f_{i,j}) T_{i,j}.$$
(4)

The number of peaks to consider in this sum depends on the size of the selected particles. When the DMA selects particles with high electrical mobility (meaning very small particles),  $T_{i,2} \ll T_{i,1}$ , because these particles have very small probability of carrying two charges. On the other hand, for particles with diameters of a few hundred nanometres, multiple charges are common, and we considered particles with up to six charges (following the custom of the old inversion algorithm).

#### 2.4.2 Prior

The particle number size distribution is assumed to be a smooth function of the particle size, and it is assumed to vary smoothly over time. These properties are modelled by giving a GP prior for the latent function

$$f(t,u) \sim \operatorname{GP}(\mu(t,u), k(u,u')k(t,t')), \tag{5}$$

where  $\mu(t, u)$  is the mean function and k(u, u')k(t, t')is the covariance function such that k(u, u') = Cov(f(t, u), f(t, u')) and k(t, t') = Cov(f(t, u), f(t', u)).

We assume that the mean function is constant,  $\mu(t, u) = \mu$ , so that it represents the average of f and give it a Gaussian prior  $\mu \sim N(0, \sigma_{\mu}^2)$ . This implies that the prior can be written as  $f(t, u) \sim \text{GP}(0, \sigma_{\mu}^2 + k(u, u')k(t, t'))$ . The covariance function along the particle size follows the Matérn covariance function with 5/2 degrees of freedom (Rasmussen and Williams, 2006):

$$k(u,u') = \sigma^2 \left( 1 - \frac{\sqrt{5}|u-u'|}{l_u} + \frac{5|u-u'|^2}{3l_u^2} \right) e^{-\sqrt{5}|u-u'|/l_u}, \quad (6)$$

where  $\sigma^2$  governs the magnitude of process variation and  $l_u$  governs the autocorrelation length of the GP along the particle size dimension. The covariance function along the time domain is exponential:

$$k(t,t') = e^{-|t-t'|/l_t},$$
(7)

where  $l_t$  is the autocorrelation length of the GP along the time dimension. The Matérn and exponential covariance functions lead to a stationary process in particle size and time dimension. The exponential covariance function corresponds to a continuous-time autoregressive model of order one and is mean-square continuous, but not mean-square differentiable (see, e.g., Rasmussen and Williams, 2006). The Matérn covariance function with 5/2 degrees of freedom is twice mean-square differentiable, for which reason our construction leads to a process that is smoother along the particle size than time dimension.

The prior variance of mean  $\sigma_{\mu}^2 = 10$ , leading to relatively flat (vague) prior distribution. The covariance function parameters,  $\theta = \{\sigma, l_t, l_u\}$ , are given weakly informative half Student *t* priors (Gelman, 2006) so that  $\sigma, l_t, l_u \sim$  Student  $t_+(\nu = 4, s^2)$ , which is the Student *t* distribution scaled and restricted to positive values. The scale parameters  $s^2$  are 3, 0.01 days, and 0.25.

#### 2.5 Implementation

#### 2.5.1 Data from the SMEAR III station in Helsinki

We used data from the urban background station SMEAR III (Järvi et al., 2009). In some wind directions, traffic emissions affect the sampled aerosol substantially (see the map in Fig. 3). The particle size distributions were measured with a twin DMPS (two DMPSs running in parallel). Each DMPS used a Hauke-type DMA and a butanol CPC from TSI. In each scan, DMPS-1 performed 15 measurements in the size range 3-40 nm using a short DMA (10.9 cm) and CPC model 3025, and DMPS-2 performed 30 measurements in the range 15-820 nm using a long DMA (28 cm) and CPC model 3010. Following the custom used with the old inversion algorithm, we discarded the first three DMPS-2 measurements due to high uncertainty in the transfer function and thereby reduced the size range to 23-820 nm. The measurements varied in duration from 5.6 to 70.2 s in DMPS-1 and from 4.8 to 9.2 s in DMPS-2. The longest measurements were for the smallest particle sizes. Between consecutive measurements, there was a lag time of about 12 s, which was needed for the voltage change and for flushing the sampling tube between the DMA and the CPC. The time stamps of individual measurements were not recorded until recently, so we had to reconstruct them. The uncertainty of the reconstructed time stamps are estimated to be 2s at the end of a scan and less than that at the beginning of a scan.

We got the transfer function from the old inversion algorithm used at University of Helsinki. For each measurement *i*, we integrated the transfer function  $T_i$  for each peak separately to get the values  $T_{i,j}$  in Eq. (4). However, as mentioned in the description of the likelihood, we ignored some of its minor peaks, and therefore the number of terms included in the sum in Eq. (4) varied from one to six. We ignored peaks for which  $T_{i,j} < \alpha T_{i,1}$ , where  $\alpha$  was chosen



Figure 3. Neighbourhood of the SMEAR III station.

as  $10^{-4}$  for DMPS-1, and  $10^{-3}$  for DMPS-2. Furthermore, we used the common assumption that the concentration of particles larger than 1 µm is negligible (Wiedensohler et al., 2012). With these choices, we got about 130 training inputs  $\mathbf{x}_{i,j} = (t_i, u_{i,j})$  for each scan as illustrated in Fig. 4.

In our pre-processing of the data we also had to reconstruct the particle counts in all measurements by multiplying the saved concentrations, sample flows, and durations of measurements. We rounded the results of this multiplication to get integer counts. This reconstruction may be affected by rounding errors, which, however, are of secondary importance.

We processed data from 26 February to 7 March 2015 in batches of eight scans. After fitting the model to the data, for the post-processing we defined a grid with 5 s time resolution and 59 points covering diameters from 3 to 1000 nm. For this grid, we calculated expected values and variance of f(E[f]) and Var[f]. In our posterior approximation (Sect. 2.5.2)  $dN/dlog_{10}D_p = dN/du =$  $\exp(f)$  is log-normally distributed, so  $E[dN/dlog_{10}D_p] =$  $\exp(E[f] + 0.5 \operatorname{Var}[f])$ . By numerical integration over the particle size we obtained expected values for the particle number concentration. To estimate the uncertainties of these concentrations, for each time point we first drew a sample of 200 size distributions from the posterior and calculated particle number concentrations based on these, and then we calculated 80% posterior intervals for the particle number concentration. Consecutive batches were overlapping each other, having two scans in common. The post-processed results from the individual batches were merged. For the 10 min in the middle of the overlap, the merged results were calculated as weighted averages with the weights gradually changing from one batch to the next one.

We did all calculations on a normal desktop computer. For each batch the model fitting took about 2 min, and another 2 min were spent on the post-processing. The sampling of



**Figure 4.** Training inputs for the scan on 2 March between 11:00 and 11:10 UTC +2 h. For clarity we put  $D_p$  instead of u on the vertical axis (y axis).

size distributions was the most time-consuming part of the post-processing.

Independent measurements of particle number concentrations were obtained with a CPC (TSI 3787 water CPC), which detected particles larger than 5 nm. The time resolution fluctuated a bit and was approximately 5 s. These data were only used for evaluating the results obtained with our inversion algorithm.

# 2.5.2 Inference

Given the model description and data, we approximate the posterior distribution (Gelman et al., 2013) of f(t, u) as follows. Let  $\mathbf{y} = \{y_1, ..., y_n\}^T$  denote the *n* counts of particles at times of measurement  $\mathbf{t} = \{t_1, ..., t_n\}^T$ ; let  $\mathbf{f} = \{\mathbf{f}_{1,.}^T, ..., \mathbf{f}_{n,.}^T\}^T$  denote all the latent variables needed to define the likelihood; and let  $\mathbf{u} = \{\mathbf{u}_{1,.}^T, ..., \mathbf{u}_{n,.}^T\}$  denote the corresponding log particle sizes. Here,  $\mathbf{f}_{i,.} = \{f_{i,j}\}_{j:T_{i,j} > \alpha T_{i,1}}$  and  $\mathbf{u}_{i,.} = \{u_{i,j}\}_{j:T_{i,j} > \alpha T_{i,1}}$  denote all the latent variables and sizes corresponding to time  $t_i$  for which  $T_{i,j}$  is greater than the threshold  $\alpha T_{i,1}$ . Due to the marginalisation property of GP, the prior for the latent variables is  $\mathbf{f} \mid t, \mathbf{u}, \mathbf{\theta} \sim N(\mathbf{0}, \mathbf{K})$ , where  $K_{l,m} = \text{Cov}(f_l, f_m)$ . The conditional posterior of the latent variables, given the hyperparameters, is then

$$p(\boldsymbol{f}|\boldsymbol{y}, \boldsymbol{t}, \boldsymbol{u}, \boldsymbol{\theta}) \propto N(\boldsymbol{f}|\boldsymbol{0}, \mathbf{K}) \prod_{i=1}^{n} p(y_i|\boldsymbol{f}),$$
(8)

where  $p(y_i|f) = \text{Poi}(y_i|V_i\sum_j \exp(f_{i,j})T_{i,j})$ . Motivated by the Laplace approximation in other GP applications (Rasmussen and Williams, 2006; Vanhatalo et al., 2010, 2013), we approximate the conditional posterior with a secondorder Taylor expansion of  $\log p(f|y, t, u, \theta)$  around the mode  $\hat{f} = \arg \max_f p(f|y, t, u, \theta)$ , which gives a Gaussian

#### B. Mølgaard et al.: Notably improved DMPS data inversion

approximation:

$$p(\boldsymbol{f}|\boldsymbol{y}, \boldsymbol{t}, \boldsymbol{u}, \boldsymbol{\theta}) \approx q(\boldsymbol{f}|\boldsymbol{y}, \boldsymbol{t}, \boldsymbol{u}, \boldsymbol{\theta}) = N(\boldsymbol{f}|\hat{\boldsymbol{f}}, \boldsymbol{\Sigma}), \tag{9}$$

where  $\Sigma^{-1} = -\nabla \nabla \log(p(f|y, t, u, \theta))|_{f=\hat{f}}$  is the Hessian of the negative log-conditional posterior at the mode. The mode  $\hat{f}$  is found by a modification of a Newton algorithm. The aim is to maximise  $\Psi(f) = \log p(y|f) + \log p(f|t, u, \theta)$ , for which the basic Newton iteration is

$$\boldsymbol{f}^{\text{new}} = \boldsymbol{f}^{\text{old}} - (\nabla \nabla \Psi)^{-1} \nabla \Psi \tag{10}$$

$$= f^{\text{old}} + (\mathbf{K}^{-1} + \mathbf{W})^{-1} (\nabla \log p(\mathbf{y}|f^{\text{old}}) - \mathbf{K}^{-1}f^{\text{old}}), \quad (11)$$

where  $\mathbf{W} = -\nabla \nabla \log p(\mathbf{y} | \mathbf{f}^{\text{old}})$ . We initialised the optimisation with  $\mathbf{f} = \mathbf{0}$ . Direct calculation of the inverse of  $\boldsymbol{\Sigma} = \mathbf{K}^{-1} + \mathbf{W}$  might be numerically unstable, so we used the form

$$\boldsymbol{\Sigma}^{-1} = \mathbf{L} \left( \mathbf{I} + \mathbf{L}^{\mathrm{T}} \mathbf{W} \mathbf{L} \right)^{-1} \mathbf{L}^{\mathrm{T}},$$
(12)

where  $\mathbf{L}\mathbf{L}^{\mathrm{T}} = \mathbf{K}$  is the Cholesky decomposition of the covariance matrix. Moreover, the likelihood is not a logconcave function of f, for which reason  $\Sigma$  may not be positive definite at early iteration steps far from mode  $\hat{f}$ . In fact, in our experience this is the usual case. For this reason we check whether  $\mathbf{I}+\mathbf{L}^{\mathrm{T}}\mathbf{W}\mathbf{L}$  is positive definite and, if not, make the Newton iteration

$$f^{\text{new}} = f^{\text{old}} + (\mathbf{K}^{-1} + \widetilde{\mathbf{W}})^{-1} (\nabla \log p(\mathbf{y} | f^{\text{old}}) - \mathbf{K}^{-1} f^{\text{old}}),$$
(13)

where  $\widetilde{W}_{l,m} = \max(W_{l,m}, 0)$  if l = m and  $\widetilde{W}_{l,m} = 0$  otherwise. Here, the implementation uses the numerically more stable form  $(\mathbf{K}^{-1} + \widetilde{\mathbf{W}})^{-1} = \mathbf{K} - \mathbf{K}\widetilde{\mathbf{W}}^{1/2} (\mathbf{I} + \widetilde{\mathbf{W}}^{1/2} \mathbf{K}\widetilde{\mathbf{W}}^{1/2})^{-1} \widetilde{\mathbf{W}}^{1/2} \mathbf{K}$  obtained using the Sherman–Morrison–Woodbury lemma (Rasmussen and Williams, 2006).

The hyperparameters,  $\theta$ , are set to their approximate maximum a posterior (MAP) estimate  $\hat{\theta} = \arg \max_{\theta} q(y|t, u, \theta) p(\theta)$ , where  $q(y|t, u, \theta)$  is the approximate marginal likelihood of the hyperparameters,

$$q(\mathbf{y}|\mathbf{t}, \mathbf{u}, \boldsymbol{\theta}) \approx p(\mathbf{y}|\mathbf{t}, \mathbf{u}, \boldsymbol{\theta}) = \int p(\mathbf{y}|f) p(f|t, \mathbf{u}, \boldsymbol{\theta}) \mathrm{d}f. \quad (14)$$

The integral on the right-hand side is not analytically tractable, for which reason we use the Laplace approximation a second time. We form a second-order Taylor expansion of  $\Psi(f)$  around  $\hat{f}$  so that  $\Psi(f) \approx \Psi(\hat{f}) - \frac{1}{2}(f - \hat{f})^{\mathrm{T}} \Sigma^{-1}(f - \hat{f})$ . Now the marginal likelihood can be approximated with a

Gaussian integral over f multiplied by a constant

$$q(\mathbf{y}|\mathbf{t}, \mathbf{u}, \boldsymbol{\theta}) = \exp(\Psi(\mathbf{f}))$$
$$\int \exp\left(-\frac{1}{2}(\mathbf{f} - \hat{\mathbf{f}})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1}(\mathbf{f} - \hat{\mathbf{f}})\right) \mathrm{d}\mathbf{f}. \quad (15)$$

The logarithm of the marginal likelihood is then (see Appendix A of Vanhatalo et al., 2010)

$$\log q(\mathbf{y}|\mathbf{t}, \mathbf{u}, \boldsymbol{\theta}) = -\frac{1}{2}\hat{f}^{\mathrm{T}}\mathbf{K}^{-1}\hat{f} + \log p(\mathbf{y}|\hat{f}) - \frac{1}{2}\log(|\mathbf{K}||\boldsymbol{\Sigma}|)$$
(16)

$$= -\frac{1}{2}\hat{\boldsymbol{f}}^{\mathrm{T}}\mathbf{K}^{-1}\hat{\boldsymbol{f}} + \log p(\boldsymbol{y}|\hat{\boldsymbol{f}}) - \frac{1}{2}\log\left(|\mathbf{I} + \mathbf{L}^{\mathrm{T}}\mathbf{W}\mathbf{L}|\right).$$
(17)

The MAP estimate of the hyperparameters can now be searched for by maximising  $\log q(\mathbf{y}|\mathbf{t}, \mathbf{u}, \boldsymbol{\theta}) + \log p(\boldsymbol{\theta})$ .

It is possible to analytically solve the gradients of  $\log q(y|t, u, \theta)$  with respect to  $\theta$  (see Rasmussen and Williams, 2006), which allows the use of gradient-based optimisation. We used the scaled conjugate gradient method available in the Matlab toolbox GPstuff (Vanhatalo et al., 2013) and optimised the hyperparameters on a log scale.

After finding  $\hat{\theta}$  and constructing the Gaussian approximation for the conditional posterior  $p(f|y, t, u, \hat{\theta})$ , we can use these approximations to calculate the (approximate) posterior predictive distribution of f(t, u) at any  $\{t, u\}$ . Due to the marginalisation properties of a GP, the posterior predictive mean and variance of f(t, u) can be calculated exactly if we know the posterior mean and variance of f (Vanhatalo, 2010). Because we cannot solve these quantities exactly, we approximate the posterior predictive mean as (Vanhatalo, 2010)

$$E[f(t,u)|\mathbf{y}, t, u, \hat{\boldsymbol{\theta}}] = \boldsymbol{k}^{\mathrm{T}} \mathbf{K}^{-1} E[f|\mathbf{y}, \hat{\boldsymbol{\theta}}] \approx \boldsymbol{k}^{\mathrm{T}} \mathbf{K}^{-1} \hat{f}$$
$$= \boldsymbol{k}^{\mathrm{T}} \nabla \log p(\mathbf{y}|\hat{f}), \qquad (18)$$

where k is a vector with elements  $k_l = \text{Cov}(f(t, u), f_l)$  and the last equality comes from the fact that

$$\nabla \left( \log p(\mathbf{y}|f) + \log p(f|t, \boldsymbol{u}, \hat{\boldsymbol{\theta}}) \right) |_{f=\hat{f}}$$
  
=  $\nabla \log p(\mathbf{y}|\hat{f}) - \mathbf{K}^{-1}\hat{f} = 0.$  (19)

Similarly, the posterior predictive variance is approximated as

$$\operatorname{Var}[f(t,u)|\mathbf{y}, t, u, \hat{\theta}] = \operatorname{Var}[f(t,u)] - \mathbf{k}^{\mathrm{T}} \left( \mathbf{K}^{-1} - \mathbf{K}^{-1} \operatorname{Cov}[f|\mathbf{y}, t, u, \hat{\theta}] \mathbf{K}^{-1} \right) \mathbf{k} \quad (20)$$

$$\approx \operatorname{Var}[f(t,u)] - \boldsymbol{k}^{\mathrm{T}} \left( \mathbf{K}^{-1} - \mathbf{K}^{-1} \left( \mathbf{K}^{-1} + \mathbf{W} \right) \mathbf{K}^{-1} \right) \boldsymbol{k}$$
(21)

$$= \operatorname{Var}[f(t, u)] - \boldsymbol{k}^{\mathrm{T}} \left( \mathbf{K} + \mathbf{W}^{-1} \right)^{-1} \boldsymbol{k}$$
(22)

$$= \operatorname{Var}[f(t, u)] - \boldsymbol{k}^{\mathrm{T}} \left( \mathbf{W} - \mathbf{W} \mathbf{L} (\mathbf{I} + \mathbf{L}^{\mathrm{T}} \mathbf{W} \mathbf{L})^{-1} \mathbf{L}^{\mathrm{T}} \mathbf{W} \right) \boldsymbol{k}, \qquad (23)$$

where the first equality is given in Vanhatalo (2010) and the last two are based on the Sherman–Morrison–Woodbury



**Figure 5.** The size distribution obtained with the old inversion algorithm and expected size distributions (new inversion) for a period with little fluctuation on 26 February.

lemma and Eq. (12). Given the approximate posterior mean and variance for f(t, u) at any  $\{t, u\}$ , it is natural to approximate the posterior distribution  $p(f(t, u)|\mathbf{y}, t, u, \hat{\theta})$  with a Gaussian distribution with the above mean and variance. The above-described Laplace approximation has been shown to produce accurate estimates for the marginal likelihood  $p(\mathbf{y}|t, u, \theta)$  and conditional posterior  $p(f|\mathbf{y}, t, u, \theta)$  in several models with similar structure (Tierney and Kadane, 1986; Rue et al., 2009; Vanhatalo et al., 2010, 2013).

#### 3 Results and discussion

We will evaluate the results from our algorithm both by looking at some illustrative examples and by comparing resulting particle number concentrations for the whole period with CPC data.

As a first test of the algorithm, let us consider periods without fluctuations, meaning periods for which the old algorithm performed well. As expected, for these periods our results agree well with the results from the old algorithm. The only clear difference is that we obtain smoother size distributions with our new algorithm as in the example in Fig. 5. Especially, in the region below 20 nm the old algorithm gave an uneven result due to low count statistics (ranging from 0 to 7 particles in each measurement). The smoother size distributions seem more plausible and were, indeed, obtained using a proper description of the count statistics in the likelihood function as described in Sect. 2.4.1 (the smoothness was also affected by the prior).

In our next example (2 March 10:30–12:00 UTC +2 h) the evolution of the size distribution was as shown in Fig. 6. Clearly, the total particle number concentration fluctuated a lot, and some changes in the size are also seen. At any given time,  $dN/dlog_{10}D_p$  changes smoothly as a function of size. The variances in Fig. 7 reflect the distance in time and



**Figure 6.** Upper panel: expected size distributions on 2 March between 10:30 and 12:00 UTC +2 h. Lower panel: size distributions for the same period obtained with the old inversion algorithm. The ticks on the time axis (*x* axis) denote the beginnings of scans.



**Figure 7.** Posterior variance of f on 2 March between 10:30 and 12:00 UTC +2 h. The ticks on the time axis (x axis) denote the beginnings of scans.

particle size to the nearest measurements: the further these measurements are, the greater the variance is. In Fig. 4 we showed the training inputs for one scan, and in Fig. 7 low-variance areas appear around the training inputs of nine subsequent scans.

The time evolution of the particle number concentrations obtained with the DMPS agrees well with the CPC measurements (Fig. 8), although the CPC generally shows lower



Figure 8. Particle number concentrations on 2 March between 10:30 and 12:00 UTC +2 h.

concentration. The reason for this difference is that the CPC measurements are not corrected for particle losses in the sampling lines. The peaks generally occur at the same time, although a few of the peaks in the CPC data are not reflected in concentrations obtained from the DMPS. For instance, the peak at 10:49 UTC +2 h is only seen in the CPC data. At this time, DMPS-1 measured 3 nm particles, and DMPS-2 measured particles with diameters of more than 500 nm. No clear signs of an elevated concentration are seen in these measurements, and thus no inversion algorithm will be able to reproduce this peak. In general, with the set-up of our twin DMPS, peaks occurring during the last 3-5 min of a scan are often not observed. On the other hand, when both DMPSs measure in the range from 10 to 50 nm, fluctuations in the concentration are most likely observed and our algorithm is able to extract these fluctuating concentrations well. For instance, the peak occurring around 11:03 UTC +2 h was well observed by the twin DMPS. With the old inversion algorithm, this peak was clearly a problem (as described in the Introduction), but with our new algorithm we obtained good agreement with the CPC data. Some expected size distributions are plotted together with the result of the old inversion algorithm in Fig. 9. Obviously, our new algorithm provided size distributions which are much more realistic, but we do not know how close these estimates are to the actual size distributions, because we have no size distribution data from other instruments. In general, the size information during peaks is limited because of the low number of DMPS measurements available.

Even less size information is available for the brief concentration peak occurring between 11:31:25 and 11:32:00 UTC +2 h. According to the CPC measurements, the top of the peak occurred at 11:31:50 UTC +2 h, which was during the waiting time in both DMPSs, so the peak is not as high according to the DMPS measurements. The few DMPS measurements which were affected by this concentra-



**Figure 9.** Expected size distribution before, during, and after the peak on 2 March at 11:03 UTC +2 h, and the result of the old inversion algorithm for the scan between 11:00 and 11:10 UTC +2 h.

tion peak were all for particles in the size range 19 to 23 nm, and these suggested higher  $dN/dlog_{10}D_p$  at 19 nm than at 23 nm, although this difference could be due to temporal fluctuation. Our algorithm gives the size distribution seen in Fig. 10 at 11:31:45 UTC +2 h; as expected,  $dN/dlog_{10}D_p$ shows a decrease between 19 and 23 nm. The sampled size distributions give examples of what the size distribution may have looked like, and they have maxima between 13 and 20 nm. This seems reasonable given the available size information and the general low values of  $dN/dlog_{10}D_p$ at the smallest diameters (Fig. 6). The accumulation mode seen at diameters around 150 nm in Fig. 10 is due to the smoothing in time. Figure 6 shows such a mode for about half an hour around this time. Considering this accumulation mode, our algorithm suggests that its concentration



Figure 10. Expected size distribution with 95 % posterior intervals and five size distributions sampled from the posterior for 2 March at 11:31:45 UTC +2 h.

fluctuations are simultaneous with the fluctuations at smaller sizes (see Fig. 6). This is caused by the smoothing in size. Fast fluctuations are not observed when the DMPS measures in the accumulation mode, so we do not expect them to occur in the accumulation mode at other times either. The accumulation mode particles have a long lifetime in the atmosphere and may originate from distant sources, while the particles smaller than 25 nm most likely originate from nearby traffic emissions (Hussein et al., 2014). However, we used a stationary covariance function with two length scales  $l_u$  and  $l_t$ . Because of differences in lifetime and origin, it would make more sense to use a non-stationary covariance function with a long timescale at diameters of a few hundreds of nanometres and a short timescale for smaller particles. In practice, however, implementing such a covariance function is not straightforward. With the current covariance function,  $l_t$  will be a compromise between the actual timescales at different sizes. Thus, we expect too much smoothing in the time dimension at small diameters and, as noted above, too little at larger diameters.

To evaluate the performance for the processed 10-day period, we compared mean particle number concentrations obtained from the DMPS and the CPC data at 10 min and 30 s resolution. At 10 min resolution the correlation between the means from the two instruments was 0.984. For comparison, when processing the DMPS data with the old inversion algorithm, the obtained correlation (0.967) was twice as far from 1 (i.e. from perfect correlation). At the higher time resolution we calculated correlations for each scan separately, and Fig. 11 shows a histogram of these correlations. Clearly, for most scans there is a good correlation, and our algorithm extracts information of the time evolution, which was lost with the old algorithm. However, for 11% of the scans, the correlation is negative, reflecting a disagreement between



**Figure 11.** Histogram of correlations between 30 s mean particle number concentrations obtained from the DMPS and the CPC. The correlations are calculated for each scan separately.



Figure 12. Particle number concentrations on 7 March between 10:30 and 11:00 UTC +2 h.

time evolutions obtained from the DMPS with our algorithm and from the CPC. We have investigated all eight cases for which the correlation is smaller than -0.75. In two cases the concentration was almost constant according to both instruments, and it seems that small fluctuations caused the negative correlation by chance. In the remaining cases, concentration changes not observed by the DMPS seem to be at least part of the reason.

Let us illustrate this with an example (Fig. 12). For the time interval 10:40-10:50 UTC +2 h, the correlation was -0.87. The CPC measurements show that the concentration was higher after 10:45 UTC +2 h than before, and at 10:50 UTC +2 h it started decreasing. According to our analysis of the DMPS data, most of the particles had diameters between 7 and 70 nm during this period, but particles in this range were not measured by the DMPS after

10:45 UTC +2 h, so the slightly elevated concentration was not observed. This elevation ended at 10:50 UTC, and the counts during the following scan suggest a lower concentration, so our algorithm suggests a smooth decrease of the concentration in the time interval 10:45-10:50 UTC +2 h. Also during the first 5 min of the scan (10:40-10:45 UTC+2 h)we observe decreasing concentration according to the DMPS and increasing concentration according to the CPC. However, it seems that our model fits that data well also during this period, meaning that the counts y agree well with the rate parameters  $\lambda$  (data not shown). If we specifically consider the first 2 min of each of the scans in Fig. 12, the CPC showed on average slightly lower concentration during 10:40-10:42 UTC +2 h than during the other 2 min periods. However, comparing these three periods, one finds that the DMPS counts were clearly highest during 10:40-10:42 UTC +2 h. This gives some indication that the relatively high concentration in the beginning of this scan is supported by the DMPS data. The negative correlation during this scan seems to originate from the measurements rather than from any problem in our inversion algorithm. Our investigation of other scans with negative correlations did not suggest problems with the inversion either. So despite these occasional negative correlations, it seems that our model extracts the information of the concentration time evolution well from the available DMPS measurements.

The results above are based on a few simplifying assumptions. We assumed that the particle concentration only changes a little during each measurement. This is not necessarily always the case, but the approximation in Eq. (3) is correct at least at some point during the measurement. For DMPS-2 most of the measurements are short ( $\sim 5$  s), and any error arising from this approximation can be considered as a minor error in the timing. For DMPS-1, each measurement at the smallest sizes last around 1 min, but strong fluctuations are rare at these sizes. In principle, we could have split the time intervals into smaller pieces and summed up their contribution, but the minor improvement would not have justified the extra computational cost. We ignored particles larger than 1 µm, but with the chosen data this seems to be a minor issue. We also ignored the uncertainties of sample and sheath flows, so our uncertainties are somewhat underestimated. The sample flows affect the likelihood directly through Eq. (2), and all flows affect the transfer function. Other small inaccuracies in the transfer function may arise from inaccurate determination of diffusional losses and differing charging probabilities of non-spherical particles, such as agglomerates from diesel exhaust (Maricq, 2008).

In summary, our algorithm extracts well the time evolution of the particle number concentration from the available DMPS data, and in the absence of fluctuations the obtained size distributions fit well with results from the old algorithm. During fluctuations, only little information about the particle sizes is available, and the uncertainties of the size distributions are considerable. Due to a lack of independent size distribution data, a quantitative evaluation of the size distributions obtained for periods with fluctuation was impossible, but there is no doubt that these size distributions are much closer to the truth than the ones obtained with the old algorithm.

In principle, this method should work for the SMPS as well, but we expect the implementation to be more difficult. The continuous scan needs to be divided into a number of counting intervals. If the counting intervals are long, the peaks of the transfer function will be much wider. On the other hand, if the counting intervals are short, the number of training inputs in our model will be high, and our algorithm will be much slower.

#### 4 Conclusions

We have developed a new algorithm (provided in the Supplement) based on a Gaussian process model for processing DMPS data, and we tested it with data from a twin DMPS in an urban background location. Our algorithm derives  $dN/dlog_{10}D_p$  as a function of  $D_p$  and t based on DMPS measurements and smoothness assumptions. Because these assumptions are more realistic than the assumption of a stationary aerosol, the derived size distributions are also much more realistic. We compared particle number concentrations with independent CPC measurements and found a good agreement.

The higher accuracy of the particle number size distributions can benefit studies of aerosols in urban locations and other places with fluctuating size distributions. The higher time resolution is useful, for instance, when attempting to pinpoint sources, given that other data, such as wind observations, exist at a good time resolution. Particle number size distributions at a high time resolution can be obtained with other instruments as well, but this algorithm offers an improvement both for existing and future DMPS data without any need to purchase new hardware.

# The Supplement related to this article is available online at doi:10.5194/amt-9-741-2016-supplement.

Acknowledgements. J. Vanhatalo and N. L. Prisle were funded by the Academy of Finland (grants 266349 and 257411, respectively).

Edited by: S. Malinowski

#### References

Gelfand, A. E., Diggle, P. J., Fuentes, M., and Guttorp, P.: Handbook of Spatial Statistics, CRC Press, Boca Raton, FL, USA, 620 pp., 2010.

#### B. Mølgaard et al.: Notably improved DMPS data inversion

- Gelman, A.: Prior distributions for variance parameters in hierarchical models, Bayesian Analysis, 1, 515–533, doi:10.1214/06-BA117A, 2006.
- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., and Rubin, D. B.: Bayesian Data Analysis, 3rd edn., Chapman and Hall/CRC, Boca Raton, FL, USA, 675 pp., 2013.
- Hinds, W. C.: Aerosol Technology: Properties, Behavior, and Measurement of Airborne Particles, John Wiley and Sons, Inc., Hoboken, NJ, USA, 504 pp., 1999.
- Hussein, T., Mølgaard, B., Hannuniemi, H., Martikainen, J., Järvi, L., Wegner, T., Ripamonti, G., Weber, S., Vesala, T., and Hämeri, K.: Fingerprints of the urban particle number size distribution in Helsinki, Finland: local versus regional characteristics, Boreal Environ. Res., 19, 1–20, 2014.
- Järvi, L., Hannuniemi, H., Hussein, T., Junninen, H., Aalto, P. P., Hillamo, R., Mäkelä, T., Keronen, P., Siivola, E., Vesala, T., and Kulmala, M.: The urban measurement station SMEAR III: Continuous monitoring of air pollution and surface-atmosphere interactions in Helsinki, Finland, Boreal Environ. Res., 14 (Supplement A), 86–109, 2009.
- Mamakos, A., Ntziachristos, L., and Samaras, Z.: Diffusion broadening of DMA transfer functions. Numerical validation of Stolzenburg model, J. Aerosol Sci., 38, 747–763, doi:10.1016/j.jaerosci.2007.05.004, 2007.
- Maricq, M. M.: Bipolar diffusion charging of soot aggregates, Aerosol Sci. Tech., 42, 247–254, doi:10.1080/02786820801958775, 2008.
- O'Hagan, A.: Curve fitting and optimal design for prediction, J. Roy. Stat. Soc. B Met., 40, 1–42, 1978.
- Rasmussen, C. E. and Williams, C. K. I.: Gaussian Processes for Machine Learning, The MIT Press, Cambridge, MA, USA, available at: www.gaussianprocess.org/gpml (last access: 1 October 2015), 2006.
- Rue, H., Martino, S., and Chopin, N.: Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations, J. Roy. Stat. Soc. B, 71, 1–35, doi:10.1111/j.1467-9868.2008.00700.x, 2009.
- Stolzenburg, M. R.: An ultrafine aerosol size distribution measuring system, PhD thesis, Department of Mechanical Engineering, University of Minnesota, USA, 1988.

- Tierney, L. and Kadane, J. B.: Accurate approximations for posterior moments and marginal densities, J. Am. Stat. Assoc., 81, 82–86, 1986.
- Vanhatalo, J.: Speeding Up the Inference in Gaussian Process Models, PhD thesis, School of Science and Technology, Aalto University, Finland, 126 pp., 2010.
- Vanhatalo, J., Pietiläinen, V., and Vehtari, A.: Approximate inference for disease mapping with sparse Gaussian processes, Stat. Med., 29, 1580–1607, doi:10.1002/sim.3895, 2010.
- Vanhatalo, J., Riihimäki, J. P., Hartikainen, J., Jylänki, P., Tolvanen, V., and Vehtari, A.: GPstuff: Bayesian Modeling with Gaussian Processes, J. Mach. Learn. Res., 14, 1175–1179, 2013.
- Voutilainen, A. and Kaipio, J. P.: Estimation of non-stationary aerosol size distributions using the state-space approach, J. Aerosol Sci., 32, 631–648, doi:10.1016/S0021-8502(00)00110-5, 2001.
- Voutilainen, A. and Kaipio, J. P.: Estimation of time-varying aerosol size distributions – exploitation of modal aerosol dynamical models, J. Aerosol Sci., 33, 1181–1200, doi:10.1016/S0021-8502(02)00062-9, 2002.
- Voutilainen, A. and Kaipio, R.: Sequential Monte Carlo estimation of aerosol size distributions, Comput. Stat. Data An., 48, 887– 908, doi:10.1016/j.csda.2004.03.011, 2005.
- Wiedensohler, A., Birmili, W., Nowak, A., Sonntag, A., Weinhold, K., Merkel, M., Wehner, B., Tuch, T., Pfeifer, S., Fiebig, M., Fjäraa, A. M., Asmi, E., Sellegri, K., Depuy, R., Venzac, H., Villani, P., Laj, P., Aalto, P., Ogren, J. A., Swietlicki, E., Williams, P., Roldin, P., Quincey, P., Hüglin, C., Fierz-Schmidhauser, R., Gysel, M., Weingartner, E., Riccobono, F., Santos, S., Grüning, C., Faloon, K., Beddows, D., Harrison, R., Monahan, C., Jennings, S. G., O'Dowd, C. D., Marinoni, A., Horn, H.-G., Keck, L., Jiang, J., Scheckman, J., McMurry, P. H., Deng, Z., Zhao, C. S., Moerman, M., Henzing, B., de Leeuw, G., Löschau, G., and Bastian, S.: Mobility particle size spectrometers: harmonization of technical standards and data structure to facilitate high quality long-term observations of atmospheric particle number size distributions, Atmos. Meas. Tech., 5, 657–685, doi:10.5194/amt-5-657-2012, 2012.