

## Retrieving of particulate matter from optical measurements: A semiparametric approach

B. Pelletier,<sup>1,2</sup> R. Santer,<sup>1,3</sup> and J. Vidot<sup>1,3</sup>

Received 5 October 2005; revised 3 March 2006; accepted 12 September 2006; published 24 March 2007.

[1] The fine particle abundance, i.e., particle matter (PM) concentration, is one of the indicators of air quality and is therefore subject to ground-based measurements. Complementary satellite aerosol remote sensing techniques provide one with maps of the aerosol optical thickness (AOT), which is sensitive to particle abundance. This paper investigates the problem of retrieving the PM concentration from the AOT, both on daily average values, on the basis of a large data set where data from the air quality networks are combined with ground-based measurements of the AOTs. It is found that a linear model fails at explaining the data well but that the performance may be significantly improved when such a linear relationship is conditioned on auxiliary parameters, mainly meteorological variables. The proposed model is expressed as an additive varying coefficient model (AVCM), which is defined as a linear model where the coefficients are additive functions of the auxiliary parameters. The model is represented using penalized smoothing splines, allowing for a proper control of the overall number of degrees of freedom via multiple smoothness parameters selection. The methodology is applied to data collected around Lille (France). The PM<sub>10</sub> concentrations are retrieved with an average uncertainty of less than 20%, leading to a correlation coefficient of 0.87 between fitted and expected PM<sub>10</sub>.

**Citation:** Pelletier, B., R. Santer, and J. Vidot (2007), Retrieving of particulate matter from optical measurements: A semiparametric approach, *J. Geophys. Res.*, 112, D06208, doi:10.1029/2005JD006737.

### 1. Introduction

[2] Air quality is of major concern and the relevant legislation is changing rapidly. The fine particle abundance is one of the indicators of air quality and is therefore subject to an official norm: PM<sub>10</sub> and PM<sub>2.5</sub> are the mass per unit volume corresponding to Particulate Matter (PM) of diameter respectively less than 10  $\mu\text{m}$  and 2.5  $\mu\text{m}$  [World Health Organization (WHO), 2000]. The size and composition of the ambient PM not only depend on the emission process, but also, and particularly for the finer fractions, on the atmospheric processes that the particles go through after emission. The particle mass is usually found in two size-related modes [Van Dingenen *et al.*, 2004; Putaud *et al.*, 2004]. The fine mode particles, up to around 1  $\mu\text{m}$ , generally originate from high-temperature processes and/or gas-to-particle formation processes in the atmosphere; these particles carry inorganic compounds (such as sulphates, nitrates, and elemental carbon) and organic compounds, including semivolatile components. Mechanical

processes such as erosion, corrosion and material abrasion give rise to coarser particles, usually larger than 1  $\mu\text{m}$  and called coarse mode particles. These particles carry, e.g., soil components and sea spray. Another fraction, the ultra fine particles (UFP), in size below 0.1  $\mu\text{m}$ , is better characterized by the number concentration (number of particles per  $\text{cm}^3$ ), because despite their large number they contribute only little to the particulate mass. Large and very small particles have a limited atmospheric residence time due to deposition or coagulation. Particles in the size range between approximately 0.1 and a few  $\mu\text{m}$  remain much longer in the atmosphere (typically several days to one week) and can consequently be transported over long distances (1000 or more kilometers). Particles are emitted directly from “primary” sources and are also formed in the atmosphere by reaction of precursor gases (“secondary sources”). The main precursor gases are SO<sub>2</sub>, NO<sub>x</sub>, VOC and NH<sub>3</sub>. Other common distinctions are natural/anthropogenic sources and combustion/noncombustion sources of aerosols [D’Almeida *et al.*, 1991].

[3] When inhaled, the larger particles contained in the PM<sub>10</sub> size fraction reach the upper part of the lung. The smaller particles of this size fraction (in particular PM<sub>2.5</sub> and PM<sub>10</sub>) penetrate more deeply into the lung and reach the alveolar region. A large body of scientific evidence has emerged that has strengthened the link between ambient PM exposure and health effects [Wilson and Sprengler, 1996]. New analyses have shown death being advanced by at least a few months on population average, at current PM con-

<sup>1</sup>Laboratoire Interdisciplinaire des Sciences de l’Environnement, Université du Littoral Côte d’Opale, Wimereux, France.

<sup>2</sup>Now at Institut de Mathématiques et de Modélisation de Montpellier, UMR CNRS 5149, Equipe de Probabilités et Statistique, Université Montpellier II, Montpellier, France.

<sup>3</sup>Association pour le Développement de la Recherche et de l’Innovation dans le Nord-Pas de Calais, Lille, France.

centrations in Europe, for causes such as cardiovascular and lung disease. Furthermore, there are robust associations between ambient PM and increases in lower respiratory symptoms and reduced lung function in children, and chronic obstructive pulmonary disease and reduced lung function in adults [Wyzga, 2002]. There is no evidence for a threshold below which ambient PM has no effect on health. In its recent review, the World Health Organization (WHO) has concluded that there is a causal relationship between PM exposure and health effects [WHO, 2000]. However, it has not been possible to establish a causal relationship between PM-related health effects and one single PM component. This is in spite of intensive research roughly over the last decade. Nevertheless, there is strong evidence to conclude that fine particles, usually measured as PM<sub>2.5</sub> in health effects studies, are more hazardous than larger ones [Schwartz *et al.*, 1996]. This does not imply that the coarse fraction of PM<sub>10</sub> is innocuous. PM characteristics found to contribute to toxicity include: metal content, presence of polycyclic aromatic hydrocarbons and other organic components, endotoxin content and small (less than 2.5  $\mu\text{m}$ ) and extremely small (less than 0.1  $\mu\text{m}$ ) size. Epidemiological studies suggest that a number of emission sources are associated with health effects, especially motor vehicles and coal combustion. Toxicological studies show that particles originating from internal combustion engines, coal burning, residual oil combustion and wood burning have strong inflammatory potential.

[4] European directives indicate for air quality the objectives to be reached at short and medium terms. For example, the European Community (EC) directive imposes an upper limit on PM<sub>10</sub> of 50  $\mu\text{g}/\text{m}^3$ , which cannot be overpassed more than 35 days a year [European Union, 1999]. This recommendation is effective since 1 January 2005. The norm will become even stricter in 2010 with a maximum of seven days. In order to respect this air quality regulation, air samples are analyzed in ground-based networks that are providing one with PM measurements.

[5] As an alternative to gravimetric measurements, one direction is to use the optical properties of the aerosols to estimate their abundance. The presence of more aerosols reduces the meteorological visibility. A radiometer at the top of the atmosphere (TOA) measures the solar direct irradiance  $E_s$  (in  $\text{W}/\text{m}^2/\mu\text{m}$ ) at different wavelengths. If the radiometer is installed at the ground level,  $E_s$  is reduced by the scattering and absorption processes through the atmospheric path, so that the solar irradiance  $E$  at ground level is expressed as:

$$E = E_s \cdot e^{-m\tau}. \quad (1)$$

[6] In this equation  $m$  is the air mass defined by  $m = 1/\cos(\theta_s)$ , where  $\theta_s$  denotes the solar zenith angle. The air mass describes the length of the atmospheric path and the key parameter is the optical thickness  $\tau$ . The contribution of the molecules is well known and when removed, we get the aerosol optical thickness (AOT)  $\tau_a$ .

[7] Earth Observations are interpreted at different levels. Level 2 are geophysical products and AOTs are provided since decades over ocean [Takayama and Takashima, 1986; Gordon and Wang, 1994] and more recently over land

[Kaufman *et al.*, 1997; Santer *et al.*, 1999]. Earth Observations offer a better spatial coverage and provide a cost-effective approach than ground-based measurements. The retrieving of particulate matter from optical measurements may be considered as a supplemental source of information, and particularly over water which is of concern is the understanding of particle transportation.

[8] In theory, the two quantities may be associated under the assumptions that (1) the aerosols are spherical, (2) the aerosols optical properties are identical in the atmospheric column, (3) the vertical distribution of the total abundance is known, (4) the chemical composition and density of the aerosols are known, and (5) the normalized particle size distribution is known. Assumption 1 simplifies the optical theory; assumption 3 allows the optical characterization of the atmospheric column into PM at the ground level, and assumption 4 allows one to determine the refractive index and to transform a number of particle into a mass.

[9] Let us now work out this association. Denote by  $n(r)$  the size distribution, and by  $n(r)dr$  (in  $\text{m}^{-3}$ ) the number of particles of radius between  $r$  and  $r + dr$ . The size distribution also depends on the altitude  $z$  (in m). At the ground level ( $z = 0$ ), we now introduce the size distribution  $n_0(r)$ , normalized to one particle:

$$n(r, 0) = N_0 n_0(r), \quad (2)$$

where  $N_0$  is the number of particles in one  $\text{m}^3$  (no unit).  $n_0(r)dr$  is in  $\text{m}^{-3}$ . Note by assumption 5,  $n_0(r)$  is known.

[10] By assumption 3, the vertical distribution is known. Suppose it is of the exponential form:

$$n(r, z) = n(r, 0) \exp\left(-\frac{z}{H_a}\right), \quad (3)$$

where  $H_a$  is the aerosols vertical scale height (m).

[11] For spherical particles, the Mie theory applies and we can compute the extinction coefficient  $\sigma_a$  (in  $\text{m}^{-1}$ ) corresponding to the normalized size distribution  $n_0(r)$  and to a particle refractive index. The vertical integration of  $\sigma_a$  corresponds to  $\tau_a$ , and we get:

$$\tau_a = N_0 H_a \sigma_a, \quad (4)$$

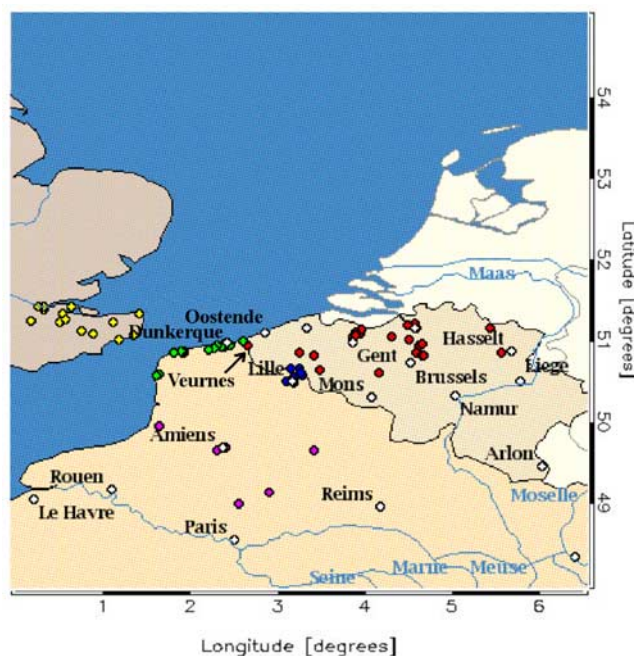
[12] For the normalized size distribution, the mass concentration  $m_o$ , of particles of radius less than  $X/2$  is expressed as:

$$m_o = \frac{4}{3} \pi d \int_0^{X/2} r^3 n_0(r) dr, \quad (5)$$

where  $d$  is the particle density ( $\mu\text{g} \cdot \text{m}^{-3}$ ), supposed to be constant with  $r$ . Finally, from equations (4) and (5), the association of the AOT and PM is obtain as the product of  $N_0$  and  $m_o$ :

$$PM_x = \frac{\tau_a m_o}{H_a \sigma_a}. \quad (6)$$

[13] Of course, through equation (6), we can convert  $\tau_a$  at any wavelength. This relation is linear for a given vertical



**Figure 1.** EXPER/PF zone with main cities (white dots) and PM monitoring stations of the air quality networks (colored dots).

distribution and a given aerosol model. In practice, one aerosol model may be selected from measurements of  $\tau_a$  at several wavelengths, but one a raw way. The vertical distribution can occasionally be measured at one place using a LIDAR or during aircraft campaigns, but certainly not on an operational basis.

[14] In other words, the use of equation (6) for retrieving PM is not straightforward. Other approaches attempt to directly relate PM and AOT. Actually, there exist some indications that  $PM_{10}$  are proportional to the AOT at 550 nm, further denoted by  $\tau_a^{550}$ . *Chu et al.* [2003] linearly correlate daily averaged values of  $PM_{10}$  and  $\tau_a^{550}$  from AERONET with a reasonable success, on the basis of 29 points collected during a 3 months period. AERONET is an optical ground-based aerosol monitoring network [*Holben et al.*, 1998]. *Wang and Christopher* [2003] do the same intercomparison but between  $PM_{2.5}$  measurements and satellite-derived AOT. More precisely, mean hourly  $PM_{2.5}$ , measured with the Tapered-Element Oscillating Microbalance (TEOM, estimated error of  $\pm 1.5 \mu g \cdot m^{-3}$ ), and  $\tau_a^{550}$  derived from measurements with the Moderate Resolution Imaging Spectroradiometer (MODIS, level 3 product on  $10 \times 10 km^2$ ) are considered, and a linear correlation of  $R = 0.7$  is reported on 1095 points.

[15] Another approach consists in assimilating the AOT derived from a satellite in an aerodynamic model which outputs PM. *Sarigiannis et al.* [2003] use satellite-derived determination of  $PM_{10}$  concentration to determine the associated risk on public health using merging data techniques with meteorological data and pollutant dispersion. *Liu et al.* [2004] propose annual mean ground-level  $PM_{2.5}$  concentration maps using the Multiangle Imaging Spectroradiometer (MISR) AOT over the continuous United States.

[16] In this paper, the relationship between AOT and PM is investigated on the basis of a large set of data collected during four consecutive years in the region of Lille (France). This data is described in section 2, together with several preprocessing operations. In section 3, evidence that a direct correlation or a linear model is not sufficient for explaining the data is given. These results motivate the introduction of auxiliary parameters, which is investigated in section 4. A model is proposed in the form of an additive varying coefficient model (AVCM), that is fitted semiparametrically using penalized smoothing splines, and under the assumption of additivity of the components. The theoretical materials and details concerning the fitting procedure are postponed in Appendix A, at the end of this paper, while in section 4, the focus is on the methodology and the results. The application to a second site is also investigated in this section. Finally in section 5, a summary of the findings is given, as well as perspectives on future work.

## 2. Database

[17] The European project EXPER/PF (EXposition des populations de l'Eurorégion aux polluants atmosphériques: le cas des poussières fines) is a project for the development and the promotion of a cross-border (between France and Belgium) database on atmospheric particulate matter (more information may be found on the Web site <http://www.appanpc-asso.org/experpf/>). As part of this project, a large database of  $PM_{10}$  records from air quality networks has been built up. For the purpose of this work, we added optical ground measurements, and auxiliary meteorological parameters. The data collection processes are described in the next subsections, as well as several preprocessing steps applied on these data.

### 2.1. Particulate Matter Measurements

[18] The PM records are obtained from 3 air quality networks, namely AREMA and OPALAIR deployed in France, and VMM deployed in Belgium. The locations of the PM stations are given in Figure 1 which displays the project area comprised between  $50^\circ 30'N$  and  $53^\circ 30'N$  and  $1^\circ W$  and  $4^\circ E$ . Belgium stations (in red) are located in rural and urban areas, and French stations are close to the cities of Lille (in blue) and Dunkerque (in green).  $PM_{10}$  are measured using the TEOM instrument. What is important to say is that these measurements are performed on dry aerosols.

[19] The instruments are contained in stations, either mobile or fixed, which are equipped with a stainless sampling head situated on the roof. The manifold is constituted of a Teflon or glass cylinder protected from the rain by a metallic cover. Several short lines are connected on the manifold until each monitor. The air is aspirated in the principal line by a high-volume pump. The analyzed air is the one which is situated around the station. The pumping of air is a continuous process. Each monitor has its own pump and samples the necessary volume of air from the manifold. The particles are pumped through a special head, which selects those whose diameter is less than  $10 \mu m$ . The larger particles impact on the impaction plate and stop. The smaller ones are carried by the flow in the monitor. To determine particles of diameter less than  $2.5 \mu m$  a sharp cut cyclone is added to modify the flow pattern of the particles



according to their size. The PM monitors are calibrated using filters with a well-known mass, and the flow is regularly checked.

[20] The PM<sub>10</sub> are acquired every hour, and the EXPER/PF database gives access to measurements from all stations, collected during years 1999 to 2002.

## 2.2. Aerosol Optical Thickness Measurements

[21] The AOT ground measurements come from the AERONET network. The AERONET sites falling into the EXPER/PF area (see Figure 1) are located in the cities of Lille (50°36'N, 3°08'E), Oostende (51°13'N, 2°55'E) and Dunkerque (51°06'N, 2°61'E). AERONET data are provided through the Web site <http://aeronet.gsfc.nasa.gov>. For the study, level 1.5 AERONET measurements are used; they consist of AOTs at 440, 670, and 870 nm. These values are retrieved from direct Sun measurements (following equation (1)) at least every 15 min under clear sky conditions, and the data are automatically cloud screened [Smirnov *et al.*, 2000]. The accuracy of AOT measurements is between 0.01 and 0.02 in AOT measurement for an air mass equal to 1 [Dubovik *et al.*, 2000].

## 2.3. Data Preprocessing

[22] Let  $\tau_a = (\tau_a^{440}, \tau_a^{670}, \tau_a^{870})^T$  be the vector of aerosol optical thicknesses at 440, 670, and 870 nm. The dimension of the  $\tau_a$  vectors has been reduced from three to two by performing a Principal Component Analysis (PCA) [see, e.g., Fedorov *et al.*, 2003]. In fact, the PCA has been applied to the logarithm of  $\tau_a$ , because (1) the logarithm of the aerosol optical thickness presents an almost linear spectral dependence (in terms of the logarithm of the wavelength) and (2) PCA is a linear technique. The components of the first two eigenvectors  $\mathbf{u}^1$  and  $\mathbf{u}^2$  of the covariance matrix of  $\log \tau_a$  are given by Table 1. The projections of  $\log \tau_a$  on  $\mathbf{u}^1$  and  $\mathbf{u}^2$  will be denoted by  $\pi^1 \tau_a$  and  $\pi^2 \tau_a$ , respectively, i.e., we have

$$\pi^i \tau_a = u_i^1 \log \tau_a^{440} + u_i^2 \log \tau_a^{670} + u_i^3 \log \tau_a^{870}, \quad (7)$$

for  $i = 1, 2$ . From the coefficients in Table 1, it may be seen that  $\pi^1 \tau_a$  corresponds approximately to the mean level of the logarithm of the aerosol optical thickness, and that  $\pi^2 \tau_a$  is comparable to the difference of the logarithm of  $\tau_a^{440}$  and  $\tau_a^{870}$ . So  $\pi^2 \tau_a$  is almost proportional to the average slope of  $\log \tau_a$ , considered as a function of the wavelength.

[23] Concerning the PM<sub>10</sub> variable, a logarithm transformation has been applied. This is motivated by the fact that the distribution of PM<sub>10</sub> is highly skewed, and a log-transform allows to almost symmetrize the shape of its distribution. Variables transformations have a long history in the field of statistics [see, e.g., Atkinson, 1985; Carroll and Ruppert, 1988]. The logarithm of the PM<sub>10</sub> will be denoted by  $y$ .

## 3. Preliminary Results: Linear Approaches

### 3.1. Direct Correlation

[24] As mentioned in the Introduction, there exists some indications that PM and AOT may be related linearly. Because we have access to this large EXPER/PF database,

**Table 1.** Components of the First Two Eigenvectors  $\mathbf{u}^1$  and  $\mathbf{u}^2$  of the Covariance Matrix of  $\log \tau_a$

$\mathbf{u}^1$	$\mathbf{u}^2$
0.513	0.781
0.596	-0.028
0.617	-0.623

we first investigate this opportunity to directly correlate AOT at 550 nm and PM. This comparison has been done for the spatial average of the five PM stations located in the Lille area on a daily basis. Results are reported in Figure 2 for the AOT station located in Lille, and it may be seen that the plot is highly scattered.

[25] In fact, this was quite expected because most of the driven parameters (size distribution, vertical profile, ...) are ignored. The additional information that may be used is the Angstrom coefficient  $\alpha$  at two wavelengths  $\lambda$  and  $\lambda'$ , defined as follows:

$$\frac{\tau_a(\lambda)}{\tau_a(\lambda')} = \left( \frac{\lambda}{\lambda'} \right)^\alpha, \quad (8)$$

which thus may be obtained from two measurements of AOT at  $\lambda$  and  $\lambda'$ . The Angstrom coefficient may be compared to  $\pi^2 \tau_a$ . The Angstrom coefficient is an important parameter because it can be directly associated to the size distribution when expressed in terms of the Junge power law [Van de Hulst, 1957].

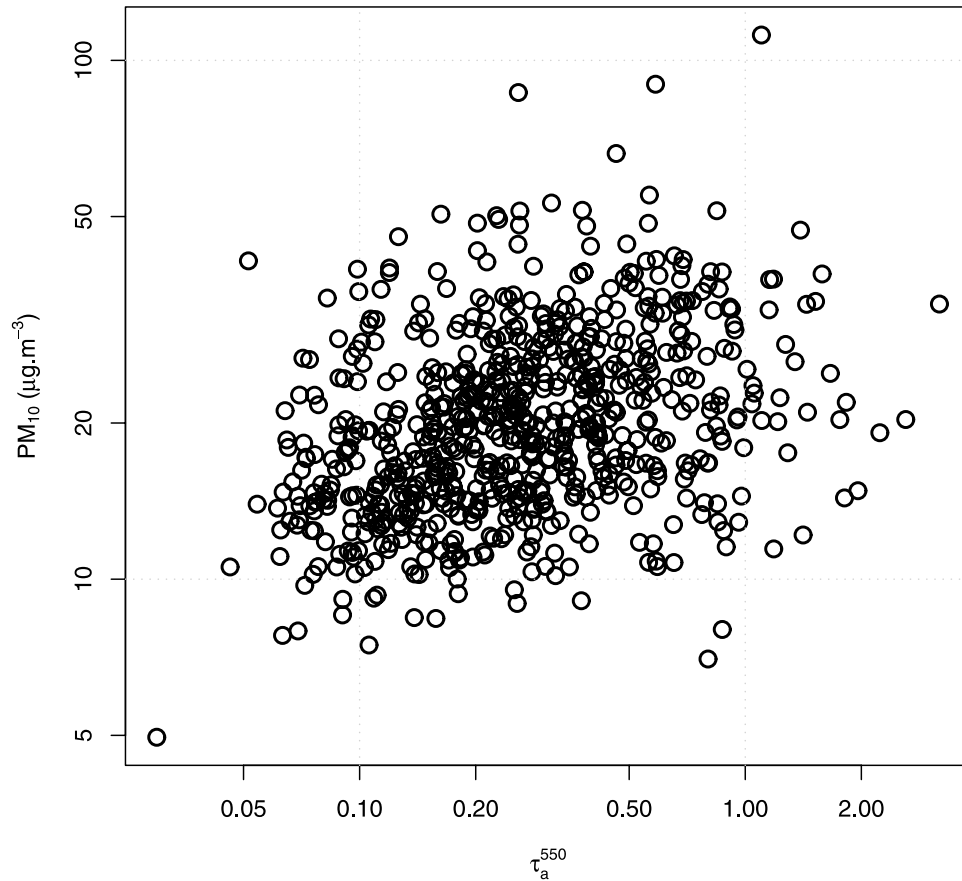
### 3.2. Linear Model

[26] The Angstrom coefficient in equation (8) corresponds to the slope of the line giving the logarithm of the AOT versus the logarithm of the wavelength. In the present case, the AOT is sampled at three wavelengths, and the empirical Angstrom coefficient is estimated by averaging the slopes between two samples. Thus the empirical Angstrom coefficient is none other than a linear combination of the logarithm of the AOTs. Under the assumption that equation (8) holds strictly, the triplet  $(\tau_a^{440}, \tau_a^{670}, \tau_a^{870})^T$  may be represented without any loss by a pair composed of (1) one AOT at a given wavelength,  $\tau_a^{440}$  say, and (2) the Angstrom coefficient. In practice however, the measurements may slightly deviate from the model in equation (8) for various reasons, including experimental noise, so that a better statistical representation (in terms of explained deviance) of the measured triplet is provided by the pair  $(\pi^1 \tau_a, \pi^2 \tau_a)$  of principal components. This pair also has the advantage over a pair of the form  $(\tau_a^{440}, \alpha)$  to have uncorrelated components. Nonetheless, the two pairs may be considered as almost equivalent in the sense that they allow for an approximate, yet accurate, reconstruction of the measured triplet.

[27] One approach for improving the correlation is to use the two pieces of information, and to express  $y$  as a linear model in  $\pi^1 \tau_a$  and  $\pi^2 \tau_a$ , i.e., as:

$$y = a_0 + a_1 \pi^1 \tau_a + a_2 \pi^2 \tau_a, \quad (9)$$

where  $a_0$ ,  $a_1$ , and  $a_2$  are scalar parameters. A model of this form has been adjusted on the data set (composed of



**Figure 2.** Direct relation ship between the aerosol optical thickness at 550 nm and the PM for the Lille area.

724 measurements points; see details in the next section) via a standard mean square fitting procedure. Performance statistics are reported in Table 2, which gives the empirical mean and standard deviation of the errors, together with the squared correlation coefficient between fitted and expected  $\log PM_{10}$ . Both are daily averaged values. As shown in this table, the standard deviation of the residuals is equal to about 0.35. Note that since the mean of the residuals is null, the standard deviation of the residuals is equal to the root mean squared error in natural logarithm of the model, which in turn corresponds, up to a second-order term, to the root mean squared relative error of the model. Thus  $PM_{10}$  concentrations are retrieved by this linear model with an average uncertainty of 35%. The squared correlation coefficient between fitted and expected  $\log PM_{10}$  is equal to 0.269, so only 27% of the deviance is explained by the model. As expected, the introduction of an additional variable slightly improves the correlation but at a poor level (compare Figure 2 with Figure 3).

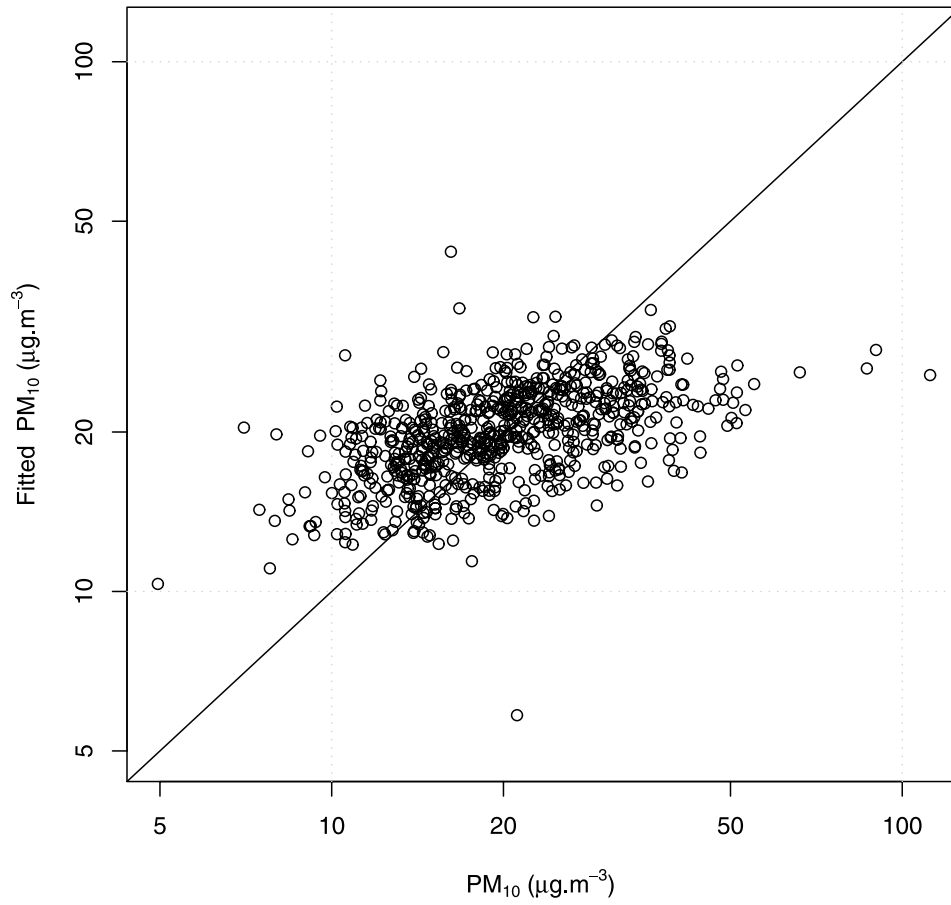
[28] So if there are some indications in the literature that a linear relationship may hold, it is found here that linearity can no longer be assumed, since models of this form fail at explaining accurately the recorded data. There may be several possible reasons for this, including the incorrectness of the linearity assumption itself. On the other hand, the general knowledge we have about aerosols conducts to associate aerosol characteristics to simple generic param-

eters. For instance, the wind direction and speed is a good indicator of the nature of the aerosols. In the region of the study, there exist two dominant wind directions: when the wind blows from S-W, maritime aerosols are expected, while N-E wind directions correspond to continental aerosol. Another interesting parameter is the date, which allows one to trace seasonal phenomena, related to the height of the mixing layer, the air relative humidity, ... Consequently, the linearity assumption is not necessarily inappropriate, but may need to be conditioned on outcomes of several auxiliary parameters. This is the general idea developed in the next section, where auxiliary parameters, namely several meteorological variables and the Julian date, are involved in the retrieval of PM from AOTs. These auxiliary parameters are considered as modifying variables, and the proposed

**Table 2.** Performance of the Models: Mean and Standard Deviation  $\sigma(\varepsilon)$  of the Residuals and Squared Correlation Coefficient  $R^2$  Between Fitted and Expected  $PM_{10}$ <sup>a</sup>

Model	Mean	$\sigma(\varepsilon)$	$R^2$
Direct	0.000	0.348	0.269
AVCM 1	0.000	0.199	0.758
AVCM 2	0.000	0.211	0.727

<sup>a</sup> Direct corresponds to the direct association between AOT at 550 nm as reported in Figure 2. The complete AVCM is denoted by AVCM 1, and the AVCM with dropped terms is denoted by AVCM 2.



**Figure 3.** Fitted versus expected  $PM_{10}$  around Lille, with a  $45^\circ$  line added, as obtained with the linear model of equation (10).

model is expressed in the form of a varying coefficient model.

#### 4. Retrieving With an Additive Varying Coefficient Model (AVCM)

##### 4.1. Auxiliary Meteorological Parameters

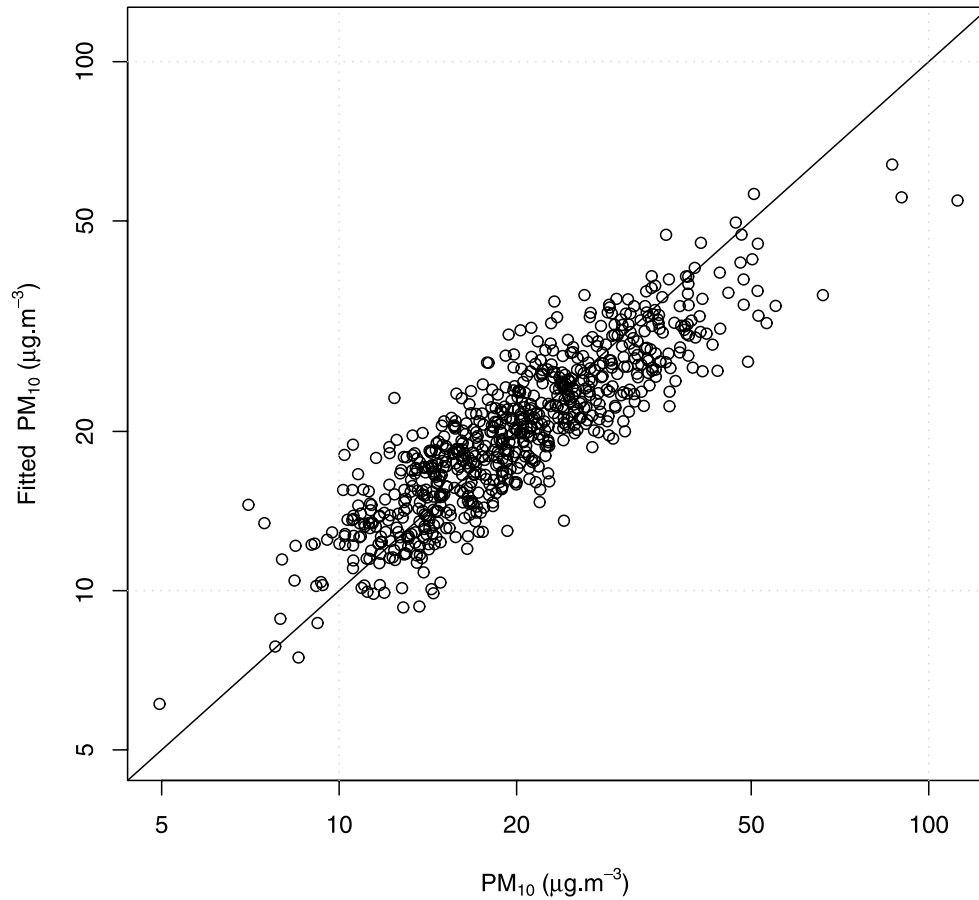
[29] Ancillary meteorological parameters are provided by the National Center for Environmental Prediction (NCEP) through the Distributed Active Archive Center (DAAC on <http://daac.gsfc.nasa.gov>). The different ancillary data are zonal and meridional wind speed components ( $m.s^{-1}$ ), surface pressure (millibars), relative humidity (%) and precipitable water ( $kg.m^{-2}$ ). The data is available every 6 hours on a global scale with a resolution of  $1^\circ$ . The zonal wind speed is defined as the west-to-east component of the wind vector (counted positive eastward), and the meridional wind speed is defined as the south-to-north wind speed (counted positive northward). The meteorological values at AERONET sites are obtained by linear interpolation on the  $1^\circ$  by  $1^\circ$  grid of the NCEP data.

[30] The data described above and collected over years 1999 to 2002 have been merged into a data set used to fit a semiparametric AVCM. In this data set, each sample is composed of the daily average of (1) the vector  $\tau_a = (\tau_a^{440}, \tau_a^{670}, \tau_a^{870})^T$  of aerosol optical thicknesses at 440, 670, and 870 nm, (2) the wind vector, (3) the pressure, (4) the relative humidity, (5) the precipitable water, and (6) the  $PM_{10}$ ,

spatially averaged over the five measurement stations. Note that the local variation of PMs is smoothed by spatial averaging. The time range of this data set corresponds to 1366 consecutive days, with missing data in the time series of aerosol optical thicknesses, due to bad sky conditions. After removing these missing data, we obtained a data set composed of 724 simultaneous occurrences of aerosol optical thicknesses,  $PM_{10}$ , and meteorological variables, on a daily scale. The wind vector, the pressure, the relative humidity, and the precipitable water, will be denoted by  $x_w$ ,  $x_p$ ,  $x_h$ , and  $x_{wp}$ , respectively. Since the measurements have been collected over time, the data is expected to be subject to some time-dependent effects, and an additional variable, further denoted by  $x_d$ , that counts the Julian day of the measurement is introduced. With these notations, the data set used to fit the statistical model is composed of 724 occurrences of  $\pi^1 \tau_a$ ,  $\pi^2 \tau_a$ ,  $x_w$ ,  $x_p$ ,  $x_h$ ,  $x_{wp}$ ,  $x_d$ , and  $y$ . Another empirical model has been developed by Liu *et al.* [2005] to retrieve  $PM_{2.5}$  from AOT and auxiliary data. The main differences of this model are the introduction of the planetary boundary layer as a parameter of the model and the use of regression technique to fit the data. Overall, this model explained 48% of the variability in daily  $PM_{2.5}$ .

##### 4.2. Model Characteristics

[31] To explain the data described above, we consider a varying coefficient model [Hastie and Tibshirani, 1993] where the logarithm of  $PM_{10}$  is related linearly to the



**Figure 4.** Fitted versus expected  $\text{PM}_{10}$  around Lille as obtained with AVCM1.

projections of the logarithm of the aerosol optical thicknesses, and where the coefficients of the linear relationship are allowed to vary over time and with the meteorological variables, i.e., we have:

$$y = f_0 + f_1 \pi^1 \tau_a + f_2 \pi^2 \tau_a, \quad (10)$$

where  $f_0$ ,  $f_1$ , and  $f_2$  are functions of the variables,  $\mathbf{x}_w$ ,  $x_p$ ,  $x_h$ ,  $x_{wv}$ , and  $x_d$ . However, this model still is very general, since the unspecified functions  $f_0$ ,  $f_1$ , and  $f_2$  depend on a large number of variables; 6 variables, indeed. Consequently it may be difficult to fit, and may not differ much from a multivariate regression model, which would diminish the interpretability of the modifying functions  $f_0$ ,  $f_1$ , and  $f_2$ .

[32] So we restrict the above varying coefficient model by assuming that the functions  $f_0$ ,  $f_1$ , and  $f_2$  are additive, i.e., for  $i = 0, 1, 2$ , we let:

$$f_i(\mathbf{x}_w, x_p, x_h, x_{wv}, x_d) = f_i^w(\mathbf{x}_w) + f_i^p(x_p) + f_i^h(x_h) + f_i^{wv}(x_{wv}) + f_i^d(x_d), \quad (11)$$

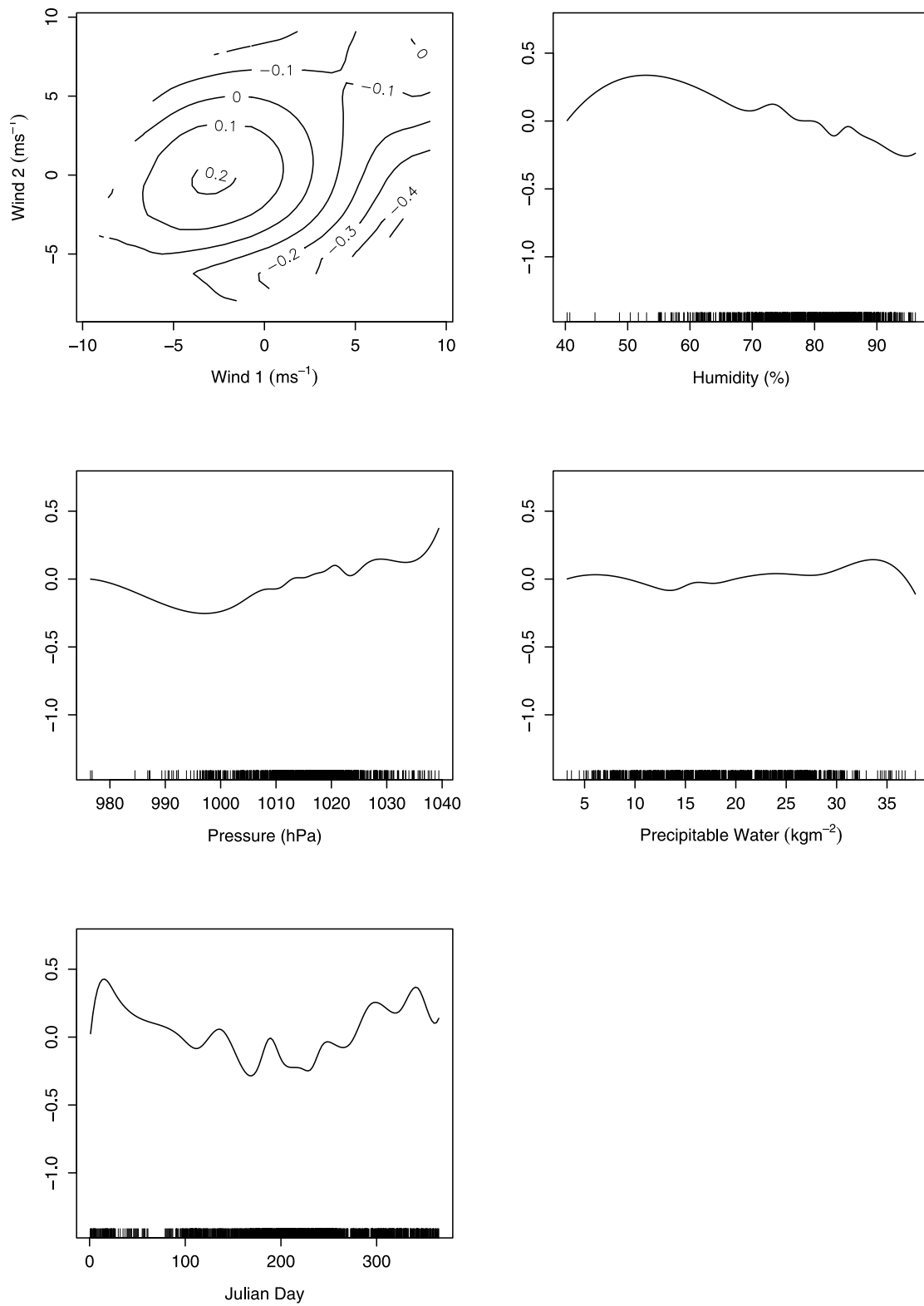
where,  $f_i^w$ ,  $f_i^p$ ,  $f_i^h$ ,  $f_i^{wv}$ , and  $f_i^d$  are unknown functions, either univariate (case of the pressure, relative humidity, water vapor, and day components) or bivariate (case of the wind component). The additivity assumption has proved useful in a variety of multivariate modeling situations, leading to

so-called additive models [Friedman and Stuetzle, 1981; Hastie and Tibshirani, 1990; Ruppert et al., 2003].

[33] The functions  $f_i^w$ ,  $f_i^p$ ,  $f_i^h$ ,  $f_i^{wv}$ , and  $f_i^d$  constitute the free parameters of the model and have to be estimated from the data. It is generally desirable that their shapes remain largely unspecified, to provide enough flexibility, while controlling the resulting number of degrees of freedom, to avoid overfitting. This may be achieved via penalized smoothing splines. Basically, a spline function is composed of polynomial pieces that are connected at some points, called knots. The theoretical materials on fitting this model are postponed to Appendix A. Here, the focus is on the results, i.e., on the shapes of these functions, more than on the process by which they are constructed.

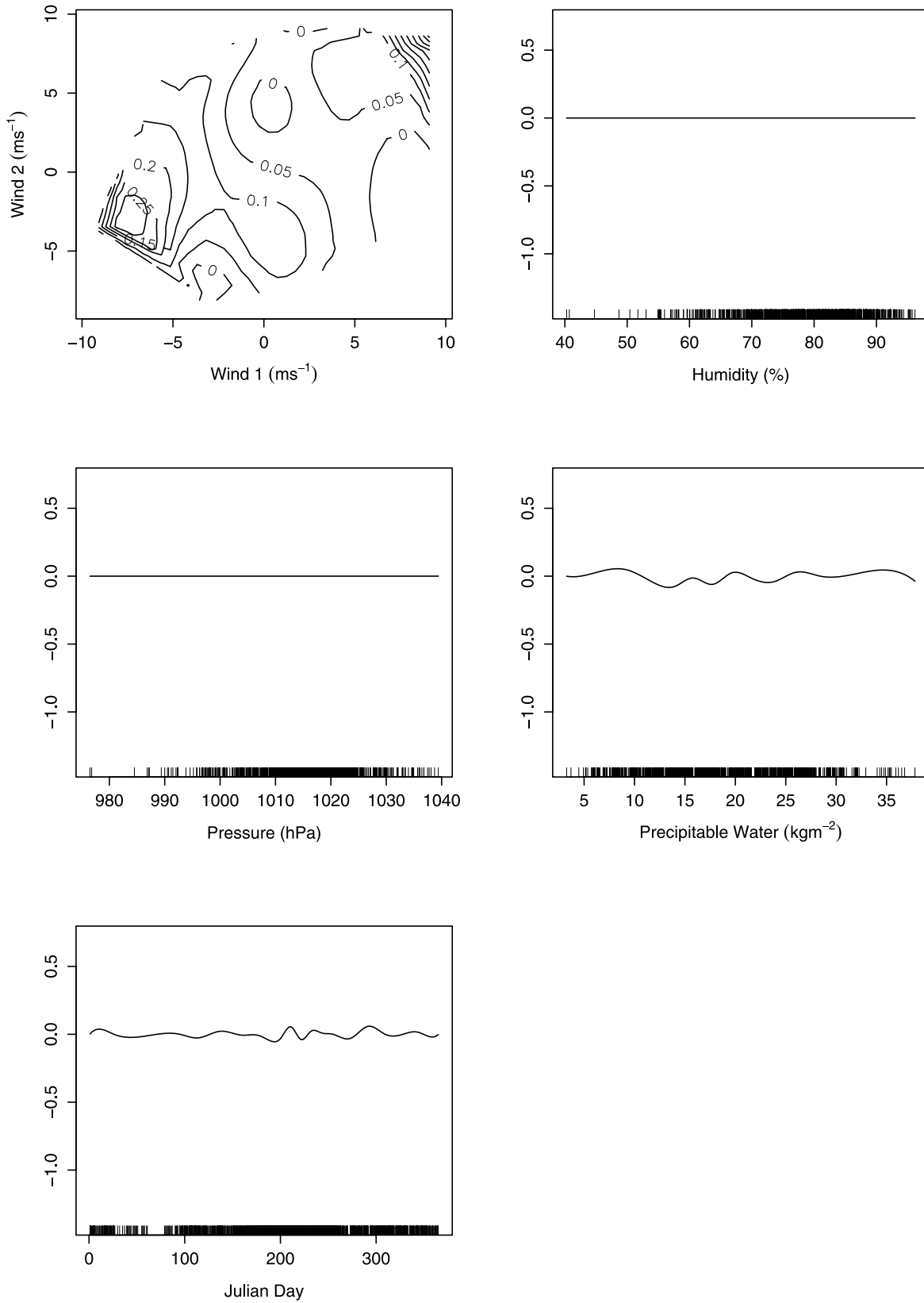
#### 4.3. Results

[34] One AVCM, further denoted by AVCM1, has been adjusted on the data described above. Performance statistics are reported in Table 2, where the improvement brought by AVCM1 over the linear model may be noticed. The mean of the residuals is negligible, and the standard deviation of the residuals is slightly improved to about 0.20, but the  $R^2$  is increased to 0.758. Hence above 75% of the deviance is explained by the model, and this corresponds to a correlation coefficient of 0.87. Fitted versus expected  $\text{PM}_{10}$  are plotted on Figure 4 using a decimal logarithmic scale, for convenience.



**Figure 5.** Intercept components of the AVCM1 fitted on Lille data.

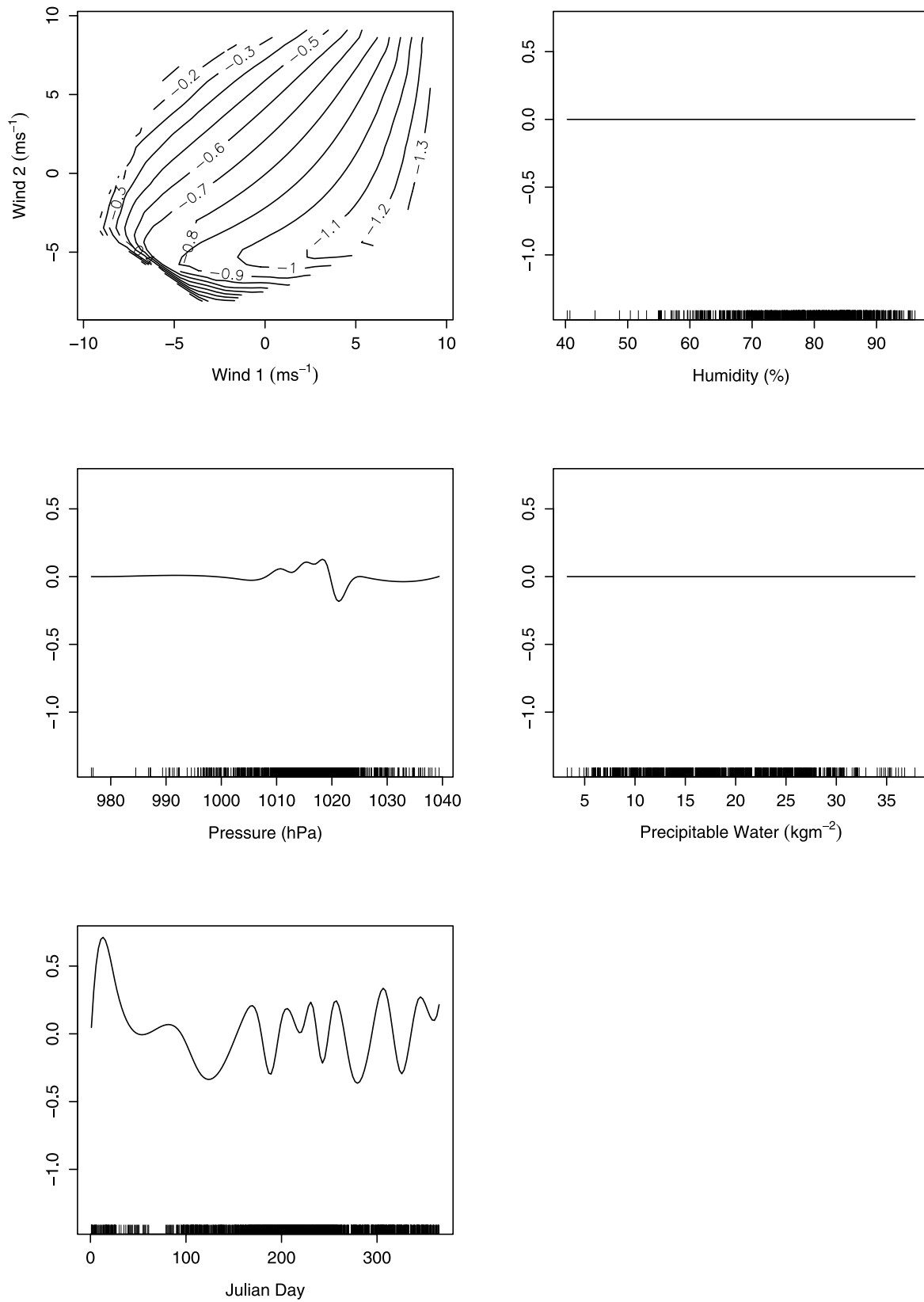




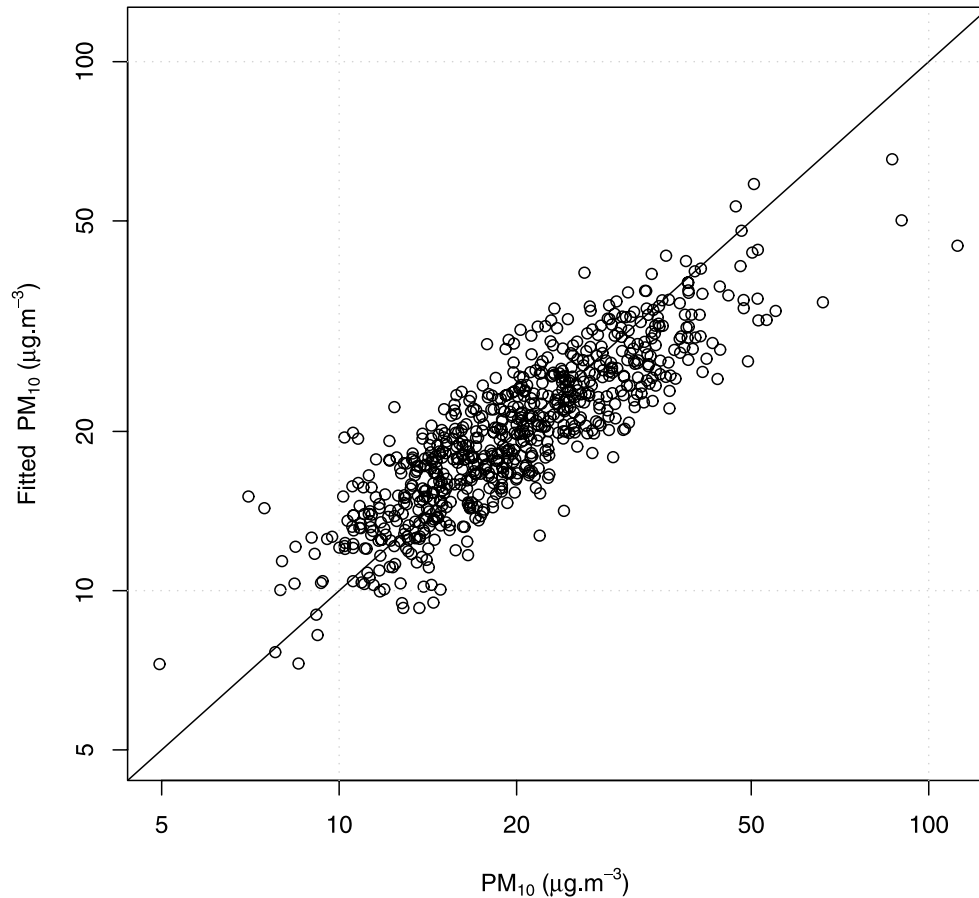
**Figure 6.**  $\pi^1\tau_a$ -component of the AVCM1 fitted on Lille data.

[35] The terms of the three components (i.e., the functions  $f_0, f_1$ , and  $f_2$ .) of AVCM1 are displayed in Figure 5 (intercept component), Figure 6 (component on  $\pi^1\tau_a$ ), and Figure 7 (component on  $\pi^2\tau_a$ ). Since each term in a component is

designed to sum up to zero over the observations, these curves are to be interpreted as variations around the mean level. In the three components, the influence of the wind appears to be directional. In the intercept component, it may



**Figure 7.**  $\pi^2\tau_a$  -component of the AVCM1 fitted on Lille data.



**Figure 8.** Fitted versus expected PM for the AVCM2 with dropped terms and fitted on Lille data.

be seen that the  $PM_{10}$  decreases with increasing relative humidity, and increases with increasing pressure and water vapor. Time-dependent effects are revealed by the terms corresponding to the day number.

[36] The curves in Figures 6 and 7 suggest that some variables have few to almost no influence on the coefficients related to  $\pi^1\tau_a$  and  $\pi^2\tau_a$ , and that their associated terms may safely be dropped. This is the case for all but the wind term in the component on  $\pi^1\tau_a$  (see Figure 6), and for all but the wind and the Day Number terms in the component on  $\pi^2\tau_a$  (see Figure 7). So we designed a similar AVCM, further denoted by AVCM2, with these terms removed, and we obtained almost identical results. The performance statistics take comparable values; see Table 2 and the scatterplot of fitted versus expected  $PM_{10}$  in Figure 8. The components of this second AVCM are displayed in Figures 9–11, and present shapes being very similar to the ones in Figures 5–7.

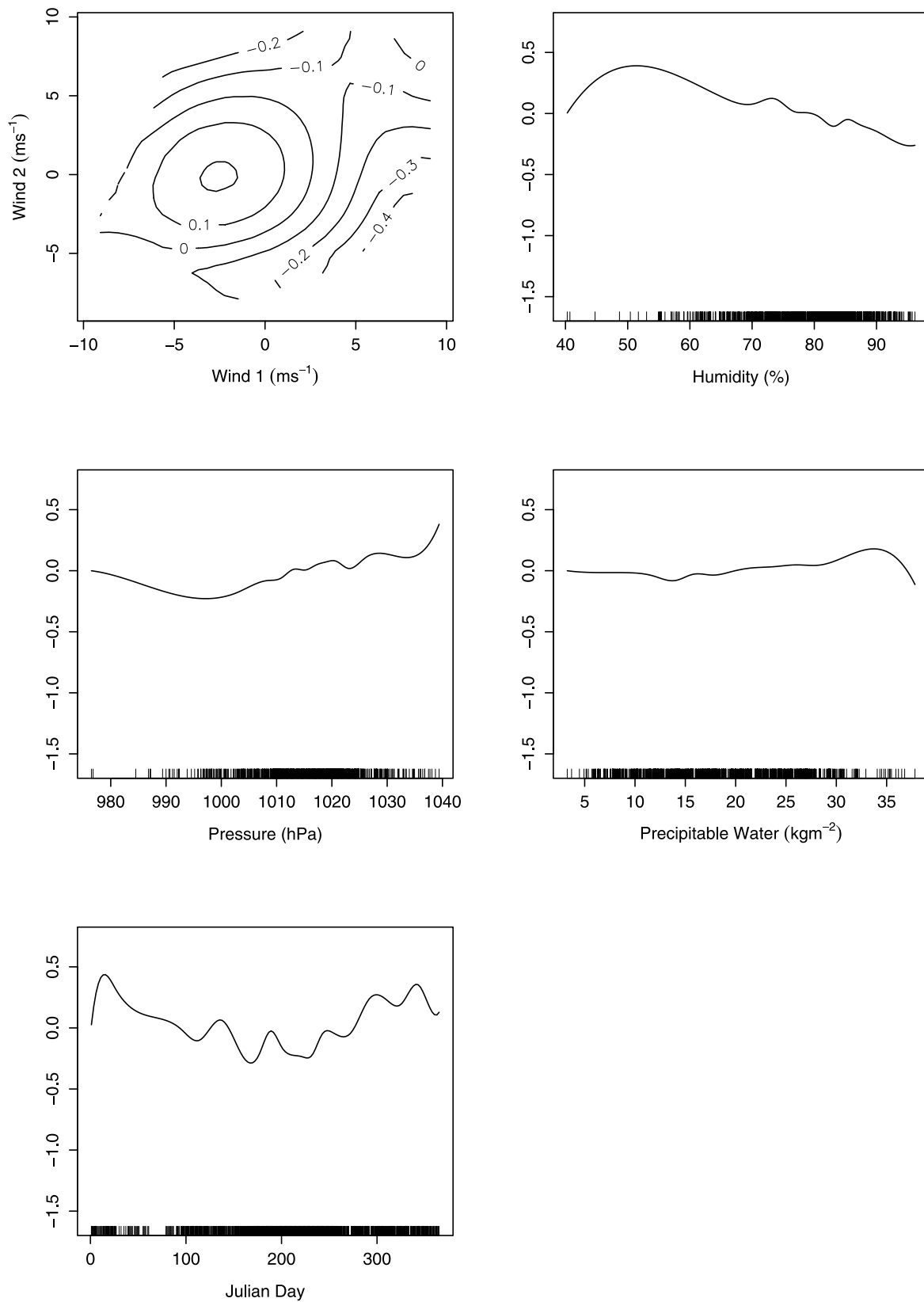
[37] Another illustration of the results is given in Figure 12 through the correlation between AOT and  $PM_{10}$  (in a log scale) at different times of the year. We now have a relationship between  $PM_{10}$  and the AOTs which depends on auxiliary parameters. For the AMCV2, it is the date and the wind speed vector. In order to investigate the role of one specific parameter, we set all the others to constant values. The two plots are for  $\alpha = -1$  and  $\alpha = 0$ . Note that only two significant parameters are derived from the three AOTs via the PCA. AOT at 440 nm is in x axis, while the relation is

plotted for these two values of  $\alpha$ . The date is the parameter we selected to illustrate the relation  $PM_{10}$  versus AOT at 440 nm. This relation is almost linear. The wind speed is set to zero. The variability induced by date includes different physical processes related to the following:

[38] 1. The vertical distribution of the aerosols impacts on the conversion of the AOT (representative of the atmospheric column) into PM (measured at ground level). The altitude of the mixing layer is lower in winter; so for a given AOT and a given aerosol model, the PM value is expected to be larger in winter than in summer, which is in accordance with the results in Figure 12.

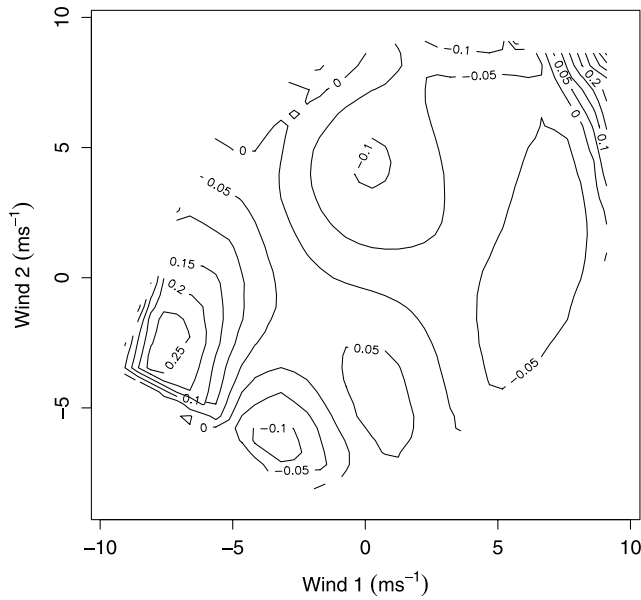
[39] 2. The air is drier in winter compared to summer. That impact on the particle size but this effect is partially described in  $\alpha$ .

[40] These two effects are also used by *Eck et al.* [2005] in an attempt to correlate  $PM_{10}$  and AOTs. Note that the date component presents some wiggleness (with intraseasonal variability) that is difficult to interpret. Indeed the date variability is likely to be mostly traced in the meteorological parameters. So under the assumption of validity of the AVCM, the date component represents a time-dependent modifying effect on the log linear PM-AOT relationship that cannot be explained from the meteorological parameter under the shape constraint, namely that of additivity. Thus the date component may be thought of as a residual component that traces over time a variability due to complex mechanisms, including the physical process described



**Figure 9.** Intercept component of the AVCM2 with dropped terms fitted on Lille data.





**Figure 10.** Wind term of the component related to  $\pi^1\tau_a$  of the AVCM2 with dropped terms and fitted on Lille data.

above, that are not characterized by meteorological variables only. By essence, the physical interpretation of the AVCM is difficult. If it is not the case, then the only two auxiliary parameters we isolated may directly explain the  $PM_{10}$  and AOT correlation through simple physical law.

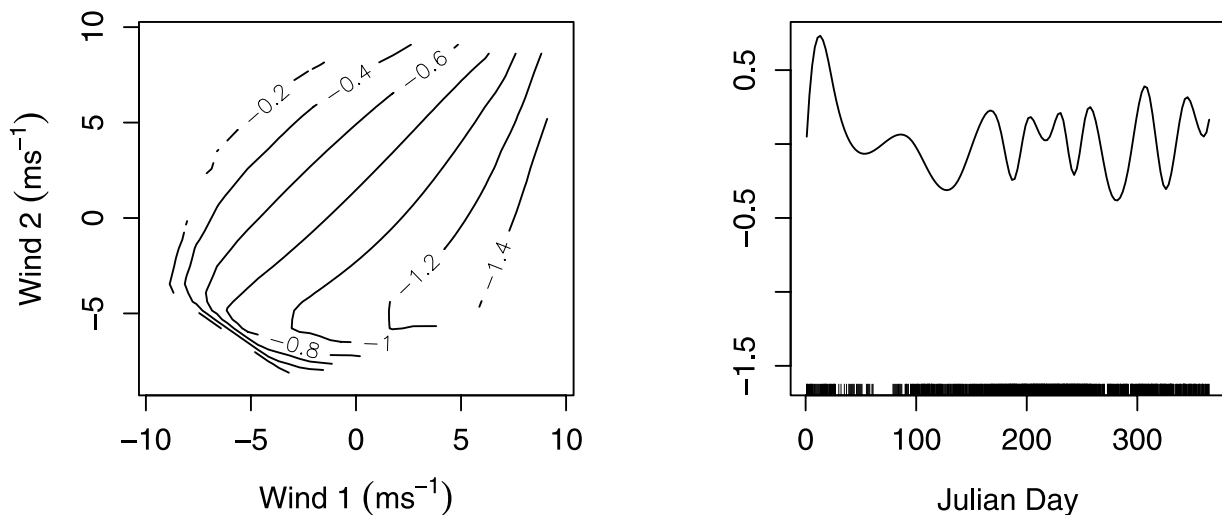
#### 4.4. Application to a Second Site

[41] We applied the two AVCM to a set of data collected on a second site during year 2003. The  $PM_{10}$  measurements were recorded in the city of Veurnes. For the aerosol optical thicknesses, as there is no station in the city of Veurnes, we computed the averages of the optical measurements from the two closest available stations, namely those located in the cities of Oostende and of Dunkerque. They are both about 20 km from the city of Veurnes. This second data set

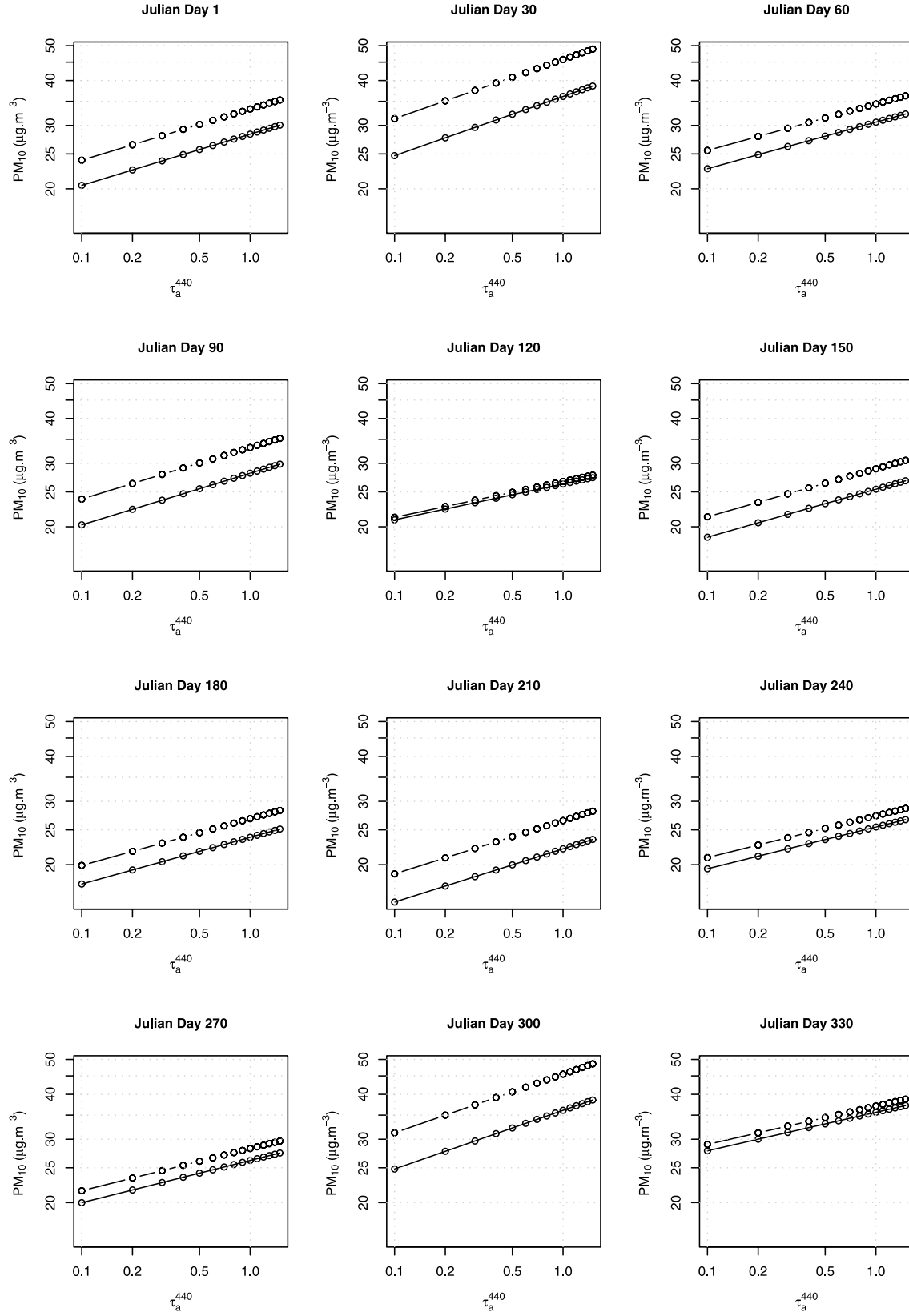
is composed of 82 samples. Retrieved versus expected  $PM_{10}$  are shown in Figure 13 for the AVCM2 model; but results with the AVCM1 are very similar. As expected, the correlation between retrieved and expected values is somewhat worsened. Retrieved values are weakly negatively biased, the average of the error (in natural logarithm) being on the order of minus up to two percents. The standard deviation of the error is on the order of 35%, and the  $R^2$  between retrieved and expected values is equal to 30%. These statistics are summarized in Table 3.

[42] To assess the validity of the AVCM approach, one AVCM of the type defined above has been fitted to this data. Because of the lower number of data points, some components are removed and the model is expressed as in equations (10) and (11); simplified with  $f_1$  and  $f_2$  as scalars. The empirical mean and standard deviation of the residuals are equal to 0 and 0.2249, respectively, and the squared correlation between fitted and expected PM is equal to 0.6992. These results are also plotted in Figure 14.

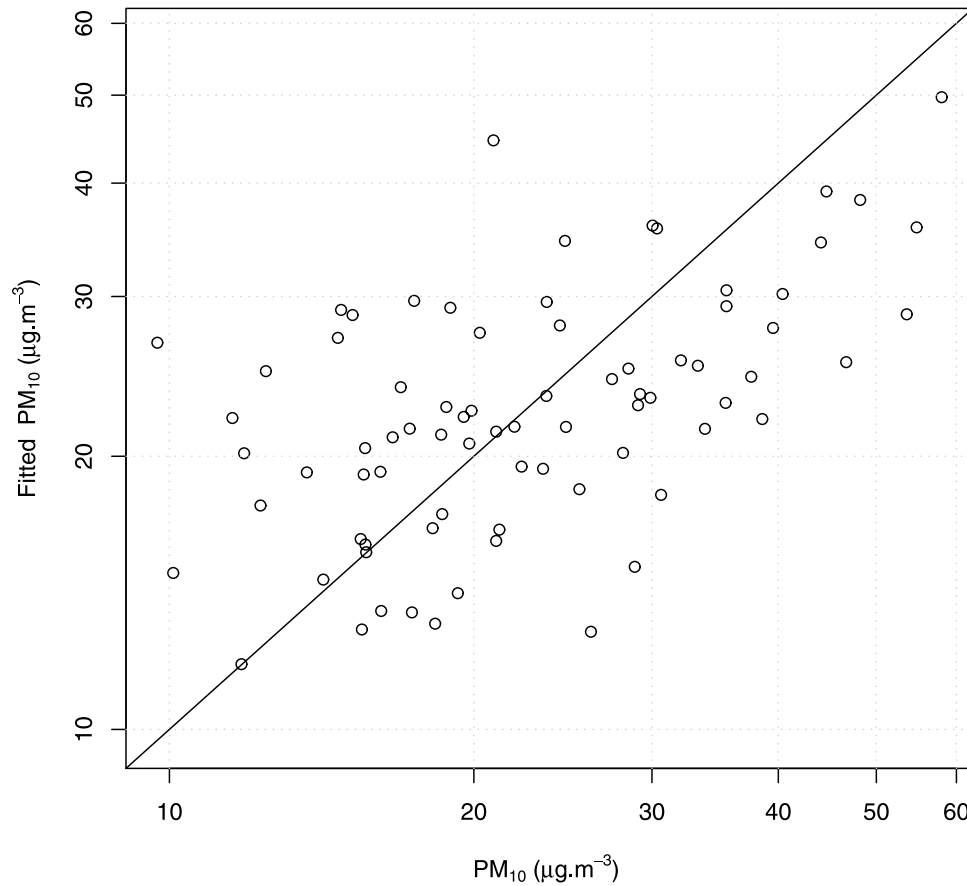
[43] These results reveal the importance of the localization, since the performance is decreased by a factor 2 when the AVCM adjusted on Lille data is applied on Veurnes data. However, these results suggest that the AVCM approach is reproducible in other locations provided that a sufficiently large data set has been set up, since the AVCM adjusted on Veurnes data explains about 70% of the deviance (75% for AVCM1 on Lille data). Thus there is a slight degradation of about 5% in terms of the explained deviance, which may be due to several factors. First, it is well known that PM measurements have to be corrected, and that the corrective factor varies from one country to one another, as it is the case between Belgium and France. Therefore, if the AOT are quite universal, the method developed on the “French” PM may not be directly applied to the “Belgium” PM, and since logarithms of PM are considered in this work, this would result in systematic biases. Second, and contrary to the case of Lille data, the CIMEL radiometer is not located in the same place as the PM measurements stations. In fact, Veurnes is located halfway between Dunkerque and Oostende, but it



**Figure 11.** Wind and Julian Day terms related to the  $\pi^2\tau_a$  component for the AVCM2 with dropped terms and fitted on Lille data.



**Figure 12.** PM10 versus AOT following the AVCM2 for different Julian days. Solid line is for  $\alpha = -1$ , and dashed line is for  $\alpha = 0$ .



**Figure 13.** Application to the second site. Predicted versus expected PM for the AVCM2 with dropped terms.

is a small town with no industrial activities compared to the two others.

## 5. Discussion and Conclusions

[44] In this work, a model is proposed to derive PM from AOT by varying the coefficients of a standard linear model in function of several auxiliary parameters. The variables in the linear model are two linear combinations of the logarithms of the AOTs at 440, 670, and 870 nm, which together are almost equivalent, in terms of provided information, to the AOTs at two distinct wavelengths, as discussed in section 3.2. Physically, the AOTs at two wavelengths are needed because the color of the aerosols is directly associated with the size distribution. The auxiliary parameters are composed of the Julian date and of the meteorological parameters.

[45] This approach is quite successful, with 70% and 75% of the deviance being explained by the models adjusted on Veurnes and Lille data, respectively. Nonetheless, further investigations are required to better assess the influence of the meteorological parameters on the relationship, since the meteorological grid of  $1^\circ$  by  $1^\circ$  used in this study is maybe too coarse for this local correlation, and better spatial resolutions are certainly needed.

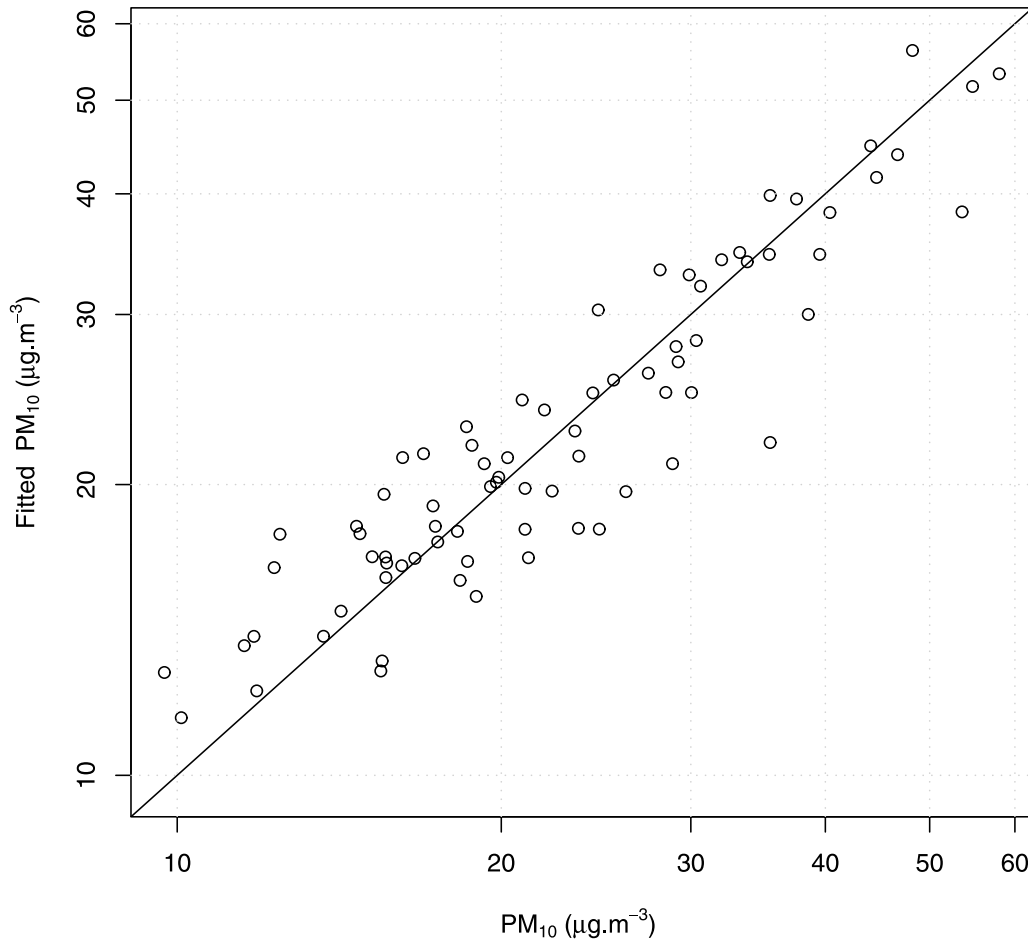
[46] As discussed in the previous section, the localization of the models is strong. In fact, this is typical of modeling situations dealing with longitudinal data with between-

subject differences (herein, one subject corresponds to one measurement site). In this case, it might be interesting to develop a unique model that explains the two data sets simultaneously. This might be achieved by incorporating random effects in the AVCM, which would lead to a mixed effects varying coefficient model. A mixed effect model is composed of fixed effects that represent a general phenomenon, common to all subjects, while random effects account for statistical differences between subjects. This point of view might be relevant when PM retrieval is sought on a finite number of locations. Alternatively, when the goal is to provide one with a continuous map of ground PM, it might be more appropriate to adopt the point of view of spatial processes, i.e., to consider the observations as realizations of spatial random fields, and to involve spatial coordinates in the retrieving model.

[47] This point of view is particularly relevant in the perspective of applying a similar methodology to satellite-derived AOTs. Presently, one limitation is the need of the

**Table 3.** Comparison Statistics for the Application of the Two AVCMs to the Second Site: Mean and Standard Deviation  $\sigma(\varepsilon)$  of the Residuals and Squared Correlation Coefficient  $R^2$  Between Predicted and Expected  $PM_{10}$

Model	Mean	$\sigma(\varepsilon)$	$R^2$
AVCM 1	−0.021	0.35	0.30
AVCM 2	−0.012	0.36	0.30



**Figure 14.** Fitted versus expected PM for the AVCM adjusted on Veurnes data. The errors have a null mean and a standard deviation of 0.2249, and the squared correlation coefficient is equal to 0.6992.

AOTs at two wavelengths. Such products are available over water, but some progress still need to be done over land.

## Appendix A: Theoretical Materials

[48] Appendix A presents the theoretical materials behind smoothing splines and varying coefficient models. The exposition is intentionally constrained to the setting of the paper, and is organized as follows. The first paragraph introduces penalized smoothing splines and their use for scatterplot smoothing. Next, varying coefficient models with additive components represented with smoothing splines are presented.

### A1. Penalized Smoothing Splines

[49] Splines functions are very attractive for nonparametric modeling and regression analysis [see, e.g., *Eilers and Marx*, 1996; *Wood*, 2003], and the references therein for materials on the subject.

[50] Basically, a spline function is composed of polynomial pieces that are connected at some points, called knots, in a certain manner, typically related to the order of differentiability of the function. A B-spline of order 3 of a real variable  $x$  is constructed on the basis of 5 knots, say  $\kappa_0$ ,  $\kappa_1$ ,  $\kappa_2$ ,  $\kappa_3$ , and  $\kappa_4$ , and consists of three cubic parts joined at the knots  $\kappa_1$ ,  $\kappa_2$ , and  $\kappa_3$ . At these joining knots, the values

of the cubic parts match, as well as their derivatives up to order 2. Outside the interval  $(\kappa_0; \kappa_4)$ , the B-spline equals zero. Now given  $M + 1$  knots  $\kappa_0, \dots, \kappa_M$ , a third-degree B-splines basis may be constructed this way; its dimension is  $M + 3$ . Now consider a smoothing problem where, given a data set of  $N$  pairs  $(x_i, y_i)$ , a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  that “fits the data well” is to be estimated. Let  $\varphi_1, \dots, \varphi_M$  be a B-spline basis of dimension  $M$ . A model  $f$  based on spline functions is defined as being a linear combination of these basis functions, i.e., it is expressed as

$$f(x) = \sum_{k=1}^M \alpha_k \varphi_k(x). \quad (\text{A1})$$

[51] The least squares estimates of the coefficients (further denoted by  $\hat{\alpha}_k$ ) are obtained by minimizing the objective function

$$\frac{1}{2} \sum_{i=1}^N (y_i - f(x_i))^2 = \frac{1}{2} \sum_{i=1}^N \left( y_i - \sum_{k=1}^M \alpha_k \varphi_k(x_i) \right)^2, \quad (\text{A2})$$

with respect to the  $\alpha_k$ , which may be solved via some elementary matrix calculus as follows.

[52] Let  $\mathbf{y} = (y_1, \dots, y_N)^T$ . Similarly, arrange the outputs  $\tilde{y}_i = f(x_i)$  of the model for the  $N$  samples in the vector  $\tilde{\mathbf{y}} =$



$(\tilde{y}_1, \dots, \tilde{y}_N)^T$ . Let  $\mathbf{C}$  be the design matrix with  $N$  rows and  $M$  columns, where the entry on the  $i$ th row and  $j$ th column is defined to be  $C_{ij} = \varphi_j(x_i)$ . Let  $\alpha = (\alpha_1, \dots, \alpha_M)^T$ . Then we have  $\tilde{\mathbf{y}} = \mathbf{C}\alpha$  and the objective function takes the following expression:

$$\frac{1}{2} \sum_{i=1}^N (y_i - f(x_i))^2 = \frac{1}{2} \|\mathbf{y} - \mathbf{C}\alpha\|^2. \quad (\text{A3})$$

[53] Then the parameter vector  $\hat{\alpha}$  minimizing the objective function is given by

$$\hat{\alpha} = (\mathbf{C}^T \mathbf{C})^{-1} \mathbf{C}^T \mathbf{y}. \quad (\text{A4})$$

[54] When the number of knots is large, i.e., when the dimension of the basis is large relatively to the number of data points, this type of model is very flexible and may lead to overfitting. The idea behind penalized smoothing splines is to penalize the roughness of the model, by minimizing the following objective function:

$$\frac{1}{2} \|\mathbf{y} - \mathbf{C}\alpha\|^2 + \lambda R(\alpha), \quad (\text{A5})$$

where  $R(\alpha)$  is a roughness penalty term, and where  $\lambda$  is a smoothness parameter that controls the trade-off between goodness-of-fit and parsimony. In the context of B-splines, a quadratic penalty based on the integrated squared second derivative of the model is a popular choice. Such a penalty is defined by  $\frac{1}{2} \alpha^T \mathbf{D} \alpha$ , where  $\mathbf{D}$  is the matrix with entries

$$D_{ij} = \int \varphi_i''(x) \varphi_j''(x) dx. \quad (\text{A6})$$

[55] Then for a fixed smoothness parameter, the solution to the minimization problem is given by

$$\hat{\alpha} = (\mathbf{C}^T \mathbf{C} + \lambda \mathbf{D})^{-1} \mathbf{C}^T \mathbf{y}. \quad (\text{A7})$$

[56] So for a given smoothness parameter, the parameter estimation problem is solved using elementary matrix algebra. However, it is in no way obvious to estimate a suitable smoothing parameter, as attested by the important statistical literature on the subject of model selection statistics, the most common of whose being cross validation or generalized cross validation. Actually, large values of the smoothness parameter lead to smooth models, thus having few degrees of freedom. Conversely, small values of the smoothness parameters lead to rough models with a large number of degrees of freedom. Smoothness parameter estimation for the AVCM is discussed in the next paragraph.

[57] So far, only modeling of univariate functions has been exposed. As mentioned above, representing multivariate functions with B-splines is fairly straightforward by considering tensor products of the marginal basis. For instance, let  $x$  and  $y$  two real variables, and let  $\varphi_1, \dots, \varphi_M$  and  $\psi_1, \dots, \psi_K$  be two univariate splines basis for respectively  $x$  and  $y$ . Then the bivariate tensor product spline basis is composed of the bivariate functions  $\varphi_i \otimes \psi_j$  for  $i = 1, \dots, M$  and  $j = 1, \dots, K$ . They are defined by  $(\varphi_i \otimes \psi_j)(x, y) = \varphi_i(x) \psi_j(y)$ . This is also

in the form of a linear model. Denoting by  $\mathbf{C}_x$  and  $\mathbf{C}_y$  the design matrices of the marginal basis for the  $x$  and  $y$  variables, respectively, the design matrix relative to the tensor product spline basis is simply given by the tensor product of  $\mathbf{C}_x$  and  $\mathbf{C}_y$ , i.e., we have  $\mathbf{C} = \mathbf{C}_x \otimes \mathbf{C}_y$ . Concerning penalty matrices for multidimensional splines, there exist different solutions. For bivariate terms, the solution adopted in this work is to have a bipenalty, i.e., one penalty along each direction. In this case, and denoting by  $\mathbf{D}_x$  and  $\mathbf{D}_y$  the marginal penalty matrices, the objective function to be minimized is now given by:

$$\frac{1}{2} \|\mathbf{y} - \mathbf{C}\alpha\|^2 + \lambda_x \frac{1}{2} \alpha^T \mathbf{D}_x \alpha + \lambda_y \frac{1}{2} \alpha^T \mathbf{D}_y \alpha, \quad (\text{A8})$$

where with a bipenalty, two smoothing parameters are involved to control the trade-off, as explained above.

## A2. AVCM

[58] A varying coefficient model is a linear model, the parameters of which are functions of so-called modifying variables [Hastie and Tibshirani, 1993], and in this paper, the components of the introduced varying coefficient models are additive; see Friedman and Stuetzle [1981], Hastie and Tibshirani [1990], and Ruppert et al. [2003] for materials on additive models. More specifically, the AVCM defined above by equation (11) has three components: one for the intercept, one related to  $\pi^1 \tau_a$ , and one related to  $\pi^2 \tau_a$ . In turn, each component is the sum of five terms: one bivariate (case of the wind term), and four univariate (case of the pressure, the relative humidity, the precipitable water, and the day number), i.e., for  $i = 0, 1, 2$ , we have, see equation (12) a total of 15 terms, that are represented using B-splines basis as defined above. So to each term there corresponds a design matrix, and either one or two penalty matrices with associated smoothness parameters, depending on the case (univariate or bivariate). Let  $\mathbf{C}_i^w$ ,  $\mathbf{C}_i^p$ ,  $\mathbf{C}_i^h$ ,  $\mathbf{C}_i^{pw}$ , and  $\mathbf{C}_i^d$  be the design matrices of the terms of the  $i$ th component, where components 0, 1, and 2 correspond to the intercept,  $\pi^1 \tau_a$ , and  $\pi^2 \tau_a$ , respectively. These matrices all have  $N = 724$  rows. Let  $\mathbf{C}_i = \text{diag}(\mathbf{C}_i^w, \mathbf{C}_i^p, \mathbf{C}_i^h, \mathbf{C}_i^{pw}, \mathbf{C}_i^d)$ , for  $i = 0, 1, 2$ . For  $i = 1$  and  $i = 2$ , define the matrix  $\hat{\mathbf{C}}_i$  where the  $k$ th entry is the  $k$ th entry of  $\mathbf{C}_i$  times the  $k$ th sample of  $\pi^i \tau_a$  in the data set. Now let  $\mathbf{C} = [\mathbf{C}_0 \ \hat{\mathbf{C}}_1 \ \hat{\mathbf{C}}_2]$ ; this is the design matrix of the AVCM. Denote by  $\mathbf{D}_i^{w,x}$ ,  $\mathbf{D}_i^{w,y}$ ,  $\mathbf{D}_i^p$ ,  $\mathbf{D}_i^h$ ,  $\mathbf{D}_i^{pw}$ ,  $\mathbf{D}_i^w$ , and  $\mathbf{D}_i^d$ , the penalty matrices associated with the terms of the  $i$ th component, where the first two correspond to the penalty matrices associated with the marginal bases of the wind vector components. Denote also by  $\alpha_i^w$ ,  $\alpha_i^p$ ,  $\alpha_i^h$ ,  $\alpha_i^{pw}$ , and  $\alpha_i^d$  the parameters of the terms of the  $i$ th component, and by  $\alpha$  a concatenation of these vectors. Then the objective function of the AVCM to be minimized is given by

$$\frac{1}{2} \|\mathbf{y} - \mathbf{C}\alpha\|^2 + \frac{1}{2} \sum_{i=0}^2 \left( \lambda_i^{w,x} (\alpha_i^w)^T \mathbf{D}_i^{w,x} \alpha_i^w + \lambda_i^{w,y} (\alpha_i^w)^T \mathbf{D}_i^{w,y} \alpha_i^w + \lambda_i^p (\alpha_i^p)^T \mathbf{D}_i^p \alpha_i^p + \lambda_i^h (\alpha_i^h)^T \mathbf{D}_i^h \alpha_i^h + \lambda_i^{pw} (\alpha_i^{pw})^T \mathbf{D}_i^{pw} \alpha_i^{pw} + \lambda_i^d (\alpha_i^d)^T \mathbf{D}_i^d \alpha_i^d \right), \quad (\text{A9})$$

where all the  $\lambda_i$  s are smoothing parameters. They are in a number of 18. Given the smoothness parameters, the estimation of the parameter vector  $\alpha$  is obtained by similar

linear operations as defined above. The main point is to estimate the smoothing parameters. This has been achieved by minimizing the Generalized Cross Validation (GCV) score according to the stable and efficient method of Wood [2004; see also Wood, 2000]. The procedure is also available as a function in the mgcv package for the R statistical software. Both are available at <http://www.r-project.org>.

[59] **Acknowledgments.** We first thank the European Commission which supports the EXPER/PF and ATTMA projects to allow us to build the required database. D. Ramon, from the HYGEOS Company, greatly helped us in providing the database extraction. We also thank the PIs of the local AERONET stations: P. Goloub for Lille, P. Goloub and J. F. Léon for Dunkerque, and K. Ruddick for Oostende.

## References

- Atkinson, A. (1985), *Plots, Transformations, and Regression*, Clarendon, Oxford, U. K.
- Carroll, R., and D. Ruppert (1988), *Transformation and Weighting in Regression*, CRC Press, Boca Raton, Fla.
- Chu, D. A., Y. J. Kaufman, G. Zibordi, J. D. Chern, J. Mao, C. Li, and B. N. Holben (2003), Global monitoring of air pollution over land from the Earth Observing System-Terra Moderate Resolution Imaging Spectroradiometer (MODIS), *J. Geophys. Res.*, 108(D21), 4661, doi:10.1029/2002JD003179.
- D'Almeida, G. A., P. Koepke, and E. P. Shettle (1991), *Atmospheric Aerosols: Global Climatology and Radiative Characteristics*, 561 pp., A. Deepak, Hampton, Va.
- Dubovik, O., A. Smirnov, B. N. Holben, M. D. King, Y. J. Kaufman, T. F. Eck, and I. Slutsker (2000), Accuracy assessments of aerosol optical properties retrieved from Aerosol Robotic Network (AERONET) Sun and sky radiance measurements, *J. Geophys. Res.*, 105(D8), 9791–9806.
- Eck, T. F., et al. (2005), Columnar aerosol optical properties at AERONET sites in central eastern Asia and aerosol transport to the tropical mid-Pacific, *J. Geophys. Res.*, 110, D06202, doi:10.1029/2004JD005274.
- Eilers, P. H. C., and B. D. Marx (1996), Flexible smoothing with B-splines and penalties, *Stat. Sci.*, 11(2), 89–121.
- European Union (1999), Council Directive 1999/30/EC of 22 April 1999 relating to limit values for sulphur dioxide, nitrogen dioxide and oxides of nitrogen, particulate matter and lead in ambient air (399L0030), *Off. J.*, L 163, 41–60.
- Fedorov, V., A. Herzberg, and S. Leonov (2003), Component-wise dimension reduction, *J. Stat. Plann. Inference*, 114, 81–93.
- Friedman, J., and W. Stuetzle (1981), Projection pursuit regression, *J. Am. Stat. Assoc.*, 76, 817–823.
- Gordon, H. R., and M. Wang (1994), Retrieval of water-leaving radiances and aerosol optical thickness over the oceans with SeaWiFS: A preliminary algorithm, *Appl. Opt.*, 33, 443–452.
- Hastie, T., and R. Tibshirani (1990), *Generalized Additive Models*, CRC Press, Boca Raton, Fla.
- Hastie, T., and R. Tibshirani (1993), Varying-coefficient models, *J. R. Stat. Soc., Ser. B.*, 55(4), 757–796.
- Holben, B., et al. (1998), AERONET—A federated instrument network and data archive for aerosol characterization, *Remote Sens. Environ.*, 66, 1–16.
- Kaufman, Y. J., D. Tanré, L. A. Remer, E. F. Vermote, A. Chu, and B. N. Holben (1997), Operational remote sensing of tropospheric aerosol over land from EOS moderate resolution imaging spectroradiometer, *J. Geophys. Res.*, 102(D14), 17,051–17,068.
- Liu, Y., R. J. Park, D. J. Jacob, Q. Li, V. Kilaru, and J. A. Sarnat (2004), Mapping annual mean ground-level PM<sub>2.5</sub> concentrations using Multi-angle Imaging Spectroradiometer aerosol optical thickness over the continuous United States, *J. Geophys. Res.*, 109, D22206, doi:10.1029/2004JD005025.
- Liu, Y., J. A. Sarnat, V. Kilaru, D. J. Jacob, and P. Koutrakis (2005), Estimating ground-level PM<sub>2.5</sub> in the eastern United States using satellite remote sensing, *Environ. Sci. Technol.*, 39, 3269–3278.
- Putaud, J.-P., et al. (2004), A European aerosol phenomenology—2: Chemical characteristics of particulate matter at kerbside, urban, rural and background sites in Europe, *Atmos. Environ.*, 38, 2579–2595.
- Ruppert, D., M. Wand, and R. Carroll (2003), *Semiparametric Regression*, Cambridge Series in Statistical and Probabilistic Mathematics, Cambridge Univ. Press, New York.
- Santer, R., V. Carrère, P. Dubuisson, and J. C. Roger (1999), Atmospheric corrections over land for MERIS, *Int. J. Remote Sens.*, 20, 1819–1840.
- Sarigiannis, D., N. I. Sifakis, N. Soulakellis, M. Tombrou, and K. Schäfer (2003), Satellite-derived determination of PM<sub>10</sub> concentration and of the associated risk on public health, in *Remote Sensing of Clouds and the Atmosphere VIII*, edited by K. P. Schäfer et al., *Proc. SPIE Int. Soc. Opt. Eng.*, 5235, 408–416.
- Schwartz, J., and D. W. Dockery, and L. M. Neas (1996), Is daily mortality associated specifically with fine particles?, *J. Air Waste Manage. Assoc.*, 46, 927–939.
- Smirnov, A., Y. Holben, T. Eck, O. Dubovik, and I. Slutsker (2000), Cloud-screening and quality control algorithms for the aeronet database, *Remote Sens. Environ.*, 73, 337–349.
- Takayama, Y., and T. Takashima (1986), Aerosol optical thickness of yellow sand over the Yellow Sea derived from NOAA satellite data, *Atmos. Environ.*, 20, 631–638.
- Van de Hulst, H. C. (1957), *Light Scattering by Small Particles*, John Wiley, Hoboken, N. J.
- Van Dingenen, R., et al. (2004), A European aerosol phenomenology—1: Physical characteristics of particulate matter at kerbside, urban, rural and background sites in Europe, *Atmos. Environ.*, 38, 2561–2577.
- Wang, J., and A. Christopher (2003), Intercomparison between satellite-derived aerosol optical thickness and PM<sub>2.5</sub> mass: Implications for air quality studies, *Geophys. Res. Lett.*, 30(21), 2095, doi:10.1029/2003GL018174.
- World Health Organization (2000), Air Quality Guidelines for Europe, 2nd ed., 288 pp., WHO Regional Office for Europe, Copenhagen, *WHO Regional Publications, European Series, No. 91*, Copenhagen. (Available at <http://www.euro.who.int/document/c71922.pdf>)
- Wilson, R., and J. D. Sprengler (1996), *Particles in our Air: Concentrations and Health Effects*, Harvard Univ. Press, Cambridge, Mass.
- Wood, S. N. (2000), Modelling and smoothing parameter estimation with multiple quadratic penalties, *J. R. Stat. Soc., Ser. B*, 65(2), 413–428.
- Wood, S. N. (2003), Thin plate regression splines, *J. R. Stat. Soc., Ser. B*, 65(1), 95–114.
- Wood, S. N. (2004), Stable and efficient multiple smoothing parameter estimation for generalized additive models, *J. Am. Stat. Assoc.*, 99(467), 673–686.
- Wynga, E. R. (2002), Air pollution and health: Are particulates the answer?, paper presented at Conference on PM<sub>2.5</sub> and Electric Power Generation: Recent Findings and Implications, Natl. Energy and Technol. Lab., Pittsburgh, Pa., 9–10 Apr.

B. Pelletier, Institut de Mathématiques et de Modélisation de Montpellier, UMR CNRS 5149, Equipe de Probabilités et Statistique, Université Montpellier II, CC 051, Place Eugène Bataillon, F-34095 Montpellier Cedex 5, France.

R. Santer and J. Vidot, LISE, Université du Littoral Côte d'Opale, 32, avenue FOCH, F-62930 Wimereux, France. ([santer@univ-littoral.fr](mailto:santer@univ-littoral.fr))