



# An automated baseline correction protocol for infrared spectra of atmospheric aerosols collected on polytetrafluoroethylene (Teflon) filters

Adele Kuzmiakova<sup>1</sup>, Ann M. Dillner<sup>2</sup>, and Satoshi Takahama<sup>1</sup>

<sup>1</sup>ENAC/IIIE Swiss Federal Institute of Technology Lausanne (EPFL), Lausanne, Switzerland

<sup>2</sup>University of California, Davis, California, USA

Correspondence to: Satoshi Takahama (satoshi.takahama@epfl.ch)

Received: 8 December 2015 – Published in Atmos. Meas. Tech. Discuss.: 29 January 2016

Revised: 23 April 2016 – Accepted: 9 May 2016 – Published: 21 June 2016

**Abstract.** A growing body of research on statistical applications for characterization of atmospheric aerosol Fourier transform infrared (FT-IR) samples collected on polytetrafluoroethylene (PTFE) filters (e.g., Russell et al., 2011; Ruthenburg et al., 2014) and a rising interest in analyzing FT-IR samples collected by air quality monitoring networks call for an automated PTFE baseline correction solution. The existing polynomial technique (Takahama et al., 2013) is not scalable to a project with a large number of aerosol samples because it contains many parameters and requires expert intervention. Therefore, the question of how to develop an automated method for baseline correcting hundreds to thousands of ambient aerosol spectra given the variability in both environmental mixture composition and PTFE baselines remains. This study approaches the question by detailing the statistical protocol, which allows for the precise definition of analyte and background subregions, applies nonparametric smoothing splines to reproduce sample-specific PTFE variations, and integrates performance metrics from atmospheric aerosol and blank samples alike in the smoothing parameter selection. Referencing 794 atmospheric aerosol samples from seven Interagency Monitoring of PROtected Visual Environment (IMPROVE) sites collected during 2011, we start by identifying key FT-IR signal characteristics, such as non-negative absorbance or analyte segment transformation, to capture sample-specific transitions between background and analyte. While referring to qualitative properties of PTFE background, the goal of smoothing splines interpolation is to learn the baseline structure in the background region to predict the baseline structure in the analyte re-

gion. We then validate the model by comparing smoothing splines baseline-corrected spectra with uncorrected and polynomial baseline (PB)-corrected equivalents via three statistical applications: (1) clustering analysis, (2) functional group quantification, and (3) thermal optical reflectance (TOR) organic carbon (OC) and elemental carbon (EC) predictions. The discrepancy rate for a four-cluster solution is 10 %. For all functional groups but carboxylic COH the discrepancy is  $\leq 10\%$ . Performance metrics obtained from TOR OC and EC predictions ( $R^2 \geq 0.94$ , bias  $\leq 0.01 \mu\text{g m}^{-3}$ , and error  $\leq 0.04 \mu\text{g m}^{-3}$ ) are on a par with those obtained from uncorrected and PB-corrected spectra. The proposed protocol leads to visually and analytically similar estimates as those generated by the polynomial method. More importantly, the automated solution allows us and future users to evaluate its analytical reproducibility while minimizing reducible user bias. We anticipate the protocol will enable FT-IR researchers and data analysts to quickly and reliably analyze a large amount of data and connect them to a variety of available statistical learning methods to be applied to analyte absorbances isolated in atmospheric aerosol samples.

## 1 Introduction

Measurement and quantification of atmospheric aerosol composition and abundance provide a basis from which we can monitor regional air quality, predict potential impacts on health and climate, and deduce formation mechanisms to reduce uncertainties in climate models for simulating al-

ternative scenarios relevant to climate change adaptation or policy decision-making (Drouet et al., 2015; Monks et al., 2009; Isaksen et al., 2009; Goldstein and Galbally, 2007; Kanakidou et al., 2005). Atmospheric aerosols, or particulate matter (PM), occur as complex mixtures of inorganic salts, crustal elements, sea spray, organic compounds, black carbon, and water (Seinfeld and Pandis, 2006), and a combination of analytical techniques are required to resolve their physical and chemical characteristics (Kulkarni et al., 2011). A useful and relatively inexpensive strategy is to collect atmospheric aerosol particles onto a substrate for offline analysis in the laboratory. Amongst different substrates, polytetrafluoroethylene (PTFE) filters have been extensively used in both measurement campaigns (Maria et al., 2002, 2003; Takahama et al., 2011; Frossard et al., 2014; Russell, 2003) and routine monitoring networks, such as the IMPROVE network in pristine and rural areas or the Chemical Speciation Network/Speciation Trends Network in urban and suburban areas in the United States (Dillner and Takahama, 2015a). Advantages of PTFE substrates include their stability, hydrophobicity, and negligible carbon gas adsorption (Turpin et al., 1994; Gilardoni et al., 2007; Ruthenburg et al., 2014). As such, they are amenable to gravimetric mass, elemental analysis, and detailed chemical speciation analysis (e.g., Surratt et al., 2007).

Carbonaceous particulate matter (PM) composition collected on PTFE filters is characterized by Fourier transform infrared (FT-IR) spectroscopy. Organic functional groups in PM absorb mid-infrared (IR) radiation in specific segments of the spectrum. The amount of light absorbed is proportional to the moles of the functional group (Beer–Lambert law). The absorption at the characteristic frequency of a particular type of bond is measured directly through PTFE filters (Maria et al., 2003; Griffiths and De Haseth, 2007). The high-frequency region ( $> 1500\text{ cm}^{-1}$ ) contains stretching and bending modes of important functional groups, such as alkane (consisting of saturated aliphatic C–CH bonds found in hydrocarbon chains), carboxylic acid (COH and C=O found in carboxylic acids and diacids), carbonyl (C=O found in ketones, aldehydes, and esters), hydroxyl (COH found in straight chain alcohols), and amine (C–NH<sub>2</sub> found in primary amines) (Russell et al., 2011). The fingerprint region ( $< 1500\text{ cm}^{-1}$ ) contains absorption bands organonitrate (CONO<sub>2</sub>) and organosulfate compounds (COSO<sub>3</sub>) (Day et al., 2010; Hawkins and Russell, 2010) but is outside the scope of our study.

A growing number of papers in recent years have been published to introduce and apply different statistical applications for atmospheric aerosol characterization from the infrared spectra. One of the applications includes unsupervised clustering of discrete spectra categories to quantify source contributions, such as fossil fuels, biomass vegetation, or biomass burning, to the total organic PM mass. Spectral clustering has been used in several atmospheric aerosol measurement campaigns and data analysis studies (Russell et al.,

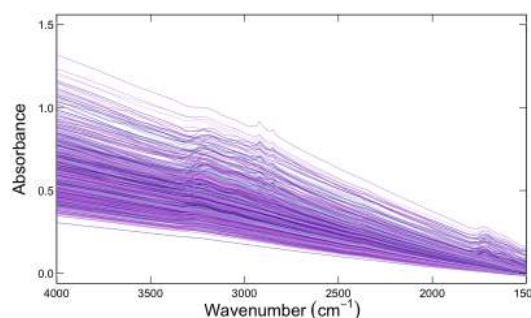
2009; Liu et al., 2009; Takahama et al., 2011; Ruthenburg et al., 2014). Cluster analysis of spectra have compared favorably with source class interpretation from factor analysis (Russell et al., 2009; Liu et al., 2009; Takahama et al., 2011; Russell et al., 2011), which attribute variations in the spectra matrix to varying contributions from an underlying set of components, and multiple linear regression with pre-determined factor sources (Takahama et al., 2011). Another approach, which has a long record in use for quantification of functional group composition and source apportionment in atmospheric aerosol samples (e.g., Russell et al., 2011), is fitting individual Gaussian line shapes to quantify alcohol COH, carboxylic COH, alkane CH, carbonyl CO, and amine NH functional groups (Takahama et al., 2013). Finally, functional groups (Coury and Dillner, 2008; Ruthenburg et al., 2014) and organic and elemental carbon content equivalent to that of thermal optical reflectance (TOR) (Dillner and Takahama, 2015a, b) have been estimated from partial least squares (PLS) calibration applied to infrared spectra. However, all applications but PLS regression require baseline-corrected infrared spectra (without PTFE interferences) to apply the Beer–Lambert law-type analysis and account for variations in analyte (aerosol) absorbance only. Aside from these statistical applications, removing PTFE interferences is a necessary step for visual inspection and comparison of similarity of aerosol composition in FT-IR spectra.

The problem of background removal is ubiquitous in nearly all spectroscopies (e.g., FT-IR, nuclear magnetic resonance (NMR), and Raman spectroscopies) and their respective applications that quantify chemical quantities based on the shape and distribution of spectral peaks (Schulze et al., 2005; Rinnan et al., 2009; Bacsik et al., 2004). A general formulation of the problem is to partition an observed spectroscopic signal into two components: one that varies smoothly (baseline) and one that is zero except in specific, localized regions (analyte). However, background correction represents an ill-posed problem; we do not know the exact proportions of the baseline and analyte in the observed signal. As a result, a realistic approach is to implement a baseline model representation capable of capturing underlying physical phenomena causing the baseline specific to the spectroscopy type. While many such investigations have been made in FT-IR biospectroscopy (Baker et al., 2014; Trevisan et al., 2012; Felten et al., 2015), single-compound, gas-phase FT-IR (Shao and Griffiths, 2007; Griffiths et al., 2009; Zhao et al., 2015), NMR (Golotvin and Williams, 2000; Xi and Rocke, 2008), and Raman spectroscopies (Weakley et al., 2012; Liland et al., 2010; Rowlands and Elliott, 2011a), the background removal in FT-IR atmospheric aerosol samples remains a far less-studied topic. Therefore, we evaluate existing classes of background correction methods to identify the most promising one based on ambient aerosol spectral characteristics.

Existing techniques include frequency decomposition (via Fourier transform, wavelets, or digital filters) to separate the

baseline component from the analyte absorption in the frequency domain (Shao et al., 2003). While the frequency decomposition techniques have been shown to successfully correct biological (Trevisan et al., 2012) or single-compound FT-IR spectra (Shao and Griffiths, 2007; Griffiths et al., 2009), they do not apply in the PM context, where spectral features are not well separated due to broad analyte absorption regions in condensed-phase aerosol samples. Another existing class, numerical differentiation (e.g., first or second differentiation, or Savitzky–Golay derivation) (Schulze et al., 2005; Rinnan et al., 2009), leads to noise amplification and requires additional smoothing that is sensitive to the signal-to-noise ratio for a specific set of samples. Furthermore, as a result of negative values from the derivative transformation, transformed spectra are difficult to visually interpret for spectroscopists. The interpolation approach (Liland et al., 2010; Ruckstuhl et al., 2001; Mazet et al., 2005; de Rooi and Eilers, 2012; Schirm and Watzig, 1998; Peng et al., 2010; Takahama et al., 2013) uses sample-specific PTFE signals on background regions where analyte absorption is not expected, and interpolates through analyte regions to identify their relative contributions at each wavelength.

Widely used interpolation methods for aerosol analysis on PTFE filters (Maria et al., 2003; Gilardoni et al., 2007; Takahama et al., 2013) are not scalable to projects with a large number (hundreds of thousands) of aerosol samples. The most modern implementation of these methods (Takahama et al., 2013) addresses the challenges described above by prescribing a set of four default background regions and a polynomial model for the variation. Background regions can be adjusted for each sample to improve accuracy; yet, this leads to additional costs in labor and variability across users. For example, users with extensive FT-IR baseline correction experience may feel comfortable using visual inspection to identify the background and analyte regions in Fig. 1. Others, on the other hand, may prefer to look at past examples or conduct a brief literature search on the presence and locations of absorbing functional groups. Alternatively, for a fixed background region, a non-negativity constraint can be imposed to alleviate issues of unrealistic spectral features that can arise from incorrect specification of the background. However, this adaptation to handle negative analyte absorbance can lead to positive bias in certain regions of ambient sample spectra, or overall in blank sample spectra for which the mean absorbance should be zero (as an average of positive and negative values). Finally, as the predefined polynomial forms are unable to account for all PTFE interferences, the method requires subtraction of blank filter spectra to remove some PTFE features a priori. Therefore, the baseline correction is performed on the residual spectra rather than the original. However, because blank filters themselves exhibit variability, there is no perfect PTFE reference and the subtraction procedure may impart additional bias. Additionally, collecting blank PTFE filters increases FT-IR analysis costs and time.



**Figure 1.** 794 FT-IR atmospheric aerosol spectra collected on PTFE filter. Each spectrum is color-differentiated.

However, a separation of atmospheric aerosol absorbance bands from the PTFE baseline via interpolation can be very complex, and therefore difficult to quantify precisely and reason about. We break down the issue into two separate problems. The first problem is determining sample-specific bounds for analyte and background subregions in the atmospheric aerosol spectra. Atmospheric PM mixtures are thought to comprise  $10^4$ – $10^5$  atmospheric organic species (Hamilton et al., 2004; Goldstein and Galbally, 2007; Kroll et al., 2011), leading to broad, overlapping IR absorption bands (features on the order of  $10$ – $10^2$   $\text{cm}^{-1}$ ) of different functional groups that absorb within similar wavenumber regions (Coury and Dillner, 2008). In Fig. 1 we show 794 atmospheric PM samples collected on PTFE filters, each differentiated by color. The overlapping absorbance bands can be seen as smoothly varying features in regions at  $\sim 3700$ – $2200$  and  $\sim 1820$ – $1500$   $\text{cm}^{-1}$ , superimposed on a sloping baseline. The range is only indicative; the wavenumber specificity is further limited by a variability in ambient PM mixture composition. As composition varies as a function of PM source and date, several of these functional groups may be absent in the sample at hand. Due to the absence of structurally distinguishable features to indicate the onset of analyte contributions, it is challenging to pinpoint the exact locations of analyte absorption. The second problem is reproducing the structure of the PTFE baseline. PTFE scattering represents the largest source of variation of the FT-IR signal when particles are collected (McClenny et al., 1985). The extent of variation in slope and shape of baseline can vary substantially among individual samples (Fig. 1). Baseline variations due to PTFE fiber stretching are unique to each sample and do not follow a prescribed or universal pattern, rendering standardized baseline preprocessing methods, for example pre-scan subtraction, standard normal variate, and multiplicative scattering correction (Rinnan et al., 2009), insufficient. Due to a lack of structural specificity of the underlying PTFE signal, we need a sample-adaptive model.

Naturally, this problem raises the question of how to develop an automated method for baseline correcting hundreds or thousands of ambient aerosol FT-IR spectra given the

**Table 1.** Notation for variables

Category	Symbol	Description
Smoothing splines model formalization (Sect. 2.1)	$w$	weight
	$y$	observed absorbance
	$\hat{y}$	fitted absorbance
	$x$	wavenumber
	$j$	an index to denote the number of wavenumbers
Spectral signal decomposition (Sect. 2.1)	$\lambda$	smoothing penalty
	$B$	background component in observed absorbances
	$A$	analyte component in observed absorbances
	$\mathcal{W}_A$	set of wavenumbers with absorbances
	$\mathcal{W}_B$	set of wavenumbers without absorbances
Smoothing splines parameter selection (Sects. 2.1, 2.3.2, 3.1)	$W_1-W_4$	specific wavenumbers to denote boundaries between analyte and background components
	EDF <sub>T</sub>	desired (target) EDF parameter defined by a user prior to applying the model
	EDF <sub>A</sub>	actual EDF parameter computed by the algorithm to match the user-specified EDF <sub>T</sub>
	EDF*	optimal EDF parameter selected from a range of EDF <sub>T</sub>

variability in environmental mixture composition and PTFE baselines. This study approaches the question by detailing the statistical protocol, which allows for the precise definition of analyte and background subregions, applies nonparametric smoothing splines to model sample-specific PTFE variations, and integrates performance metrics from PM and blank samples alike in the smoothing parameter selection. Referencing an extensive set of atmospheric aerosol samples, in Sect. 2 we start by identifying key FT-IR signal characteristics (such as non-negative absorbance or analyte segment transformation), which reduce signal variations to fundamental features to capture sample-specific transitions between background and analyte. To reproduce sample-specific variations in PTFE background and analyte structures, we develop a nonparametric, adaptive model: interpolation based on smoothing splines regulated by the roughness parameter. While referring to qualitative properties of the baseline (such as smoothness), the goal is to learn the baseline structure in the background region to predict the baseline structure in the analyte region. In Sect. 3 we evaluate the model both at the physical and application layers. We establish the initial model feasibility by using near-zero blank absorbance and non-negative analyte absorbance as our physical criteria. Further, by comparing smoothing splines baseline (SSB)-corrected spectra with polynomial baseline (PB)-corrected spectra via three different applications, (1) visual and clustering analysis, (2) functional group quantification, and (3) organic and elemental carbon prediction, we are able to discern which variations in quantities obtained from SSB-corrected spectra are due to inherent variations already present and which are added due to the new baseline approximation. We

close with a summary of the baseline correction procedure extendible to the fingerprint region or spectra acquired on other substrates in Sect. 4.

## 2 Methods

Section 2.1 introduces smoothing splines in the context of FT-IR signal. Sections 2.2 and 2.3 detail the modeling protocol, including formalizing bounds for analyte and background regions and selecting smoothing parameters. Sections 2.4 and 2.5 describe the data set and applications we used for smoothing splines model evaluation.

### 2.1 Smoothing splines model description

For the sake of clarity, Table 1 summarizes notation for commonly used variables pertaining to specific categories in implementing the smoothing splines model. The proposed interpolation method uses smoothing splines, a popular nonparametric regression technique, which has been applied in different steps in spectral signal analysis: data exploration, model building, testing parametric models, and diagnosis (Rouh et al., 1993; Rowlands and Elliott, 2011b; Pouillet et al., 2007; Persson et al., 1992; Katajamaa and Oresic, 2007; Fourmond et al., 2009). Their expression is obtained by minimizing the following two-part objective function:

$$\underset{\hat{y}}{\text{minimize}} \sum_{j=1}^n w_j (y_j - \hat{y}_j)^2 + \lambda \int_a^b (\hat{y}''(x))^2 dx, \quad (1)$$

where  $w$  is weight at wavenumber  $j$ ,  $y$  and  $\hat{y}$  are observed and fitted absorbances at wavenumber  $j$ , and  $\lambda$  is a smoothing penalty. Minimizing this criterion over the entire spectrum leads to a unique solution, which is a natural cubic spline with knots at the unique values of the wavenumbers  $x$  for  $j = 1, 2, \dots, N$  (Hastie et al., 2009). The explicit solution in form of the natural spline eliminates the knot selection problem without leading to over-parameterization due to the smoothing penalty constraint. The advantage of smoothing splines is their capacity to operate both locally, through  $w$  representations for each wavenumber  $j$ , and globally, through a single  $\lambda$  representation over the entire wavenumber domain.

The first, least squares term,  $\sum_{j=1}^n w_j (y_j - \hat{y}_j)^2$ , represents the similarity measure consisting of the squared distance between observed absorbance values and interpolating function values. The advantage of locally moderated weights lies in allowing us to choose whether absorbance at a particular wavenumber  $j$  should be included in determining the fitted baseline. We define weights as follows. Let us decompose the original spectral signal into a two-component mixture:

$$y_j = \begin{cases} B_j + A_j & \text{if } j \in \mathcal{W}_A \text{ (analyte region)} \\ B_j & \text{if } j \in \mathcal{W}_B \text{ (background region)}. \end{cases} \quad (2)$$

Here  $B_j$  denotes the background component comprising baseline, noise, and, if present, any remaining local, high-frequency interference (Takahama et al., 2013).  $A_j$  denotes the analyte component,  $\mathcal{W}_A$  denotes the set of wavenumbers with analyte absorbance, and  $\mathcal{W}_B$  denotes the set of wavenumbers without analyte absorbance. We then wish to select observations that represent solely the background component and exclude those that contain the analyte contribution:

$$w_j = \begin{cases} 0 & \text{if } j \in \mathcal{W}_A \\ 1 & \text{if } j \in \mathcal{W}_B. \end{cases} \quad (3)$$

Other conceptually analogous variants for determining weights exist. Some researchers define weights as posterior probabilities from mixture models (de Rooi and Eilers, 2012). Some researchers use curve fitting with asymmetric weights (Liland et al., 2011; Felten et al., 2015; Peng et al., 2010; Mazet et al., 2005). While differing in the requirement of a priori knowledge on the assignment of observations to different components, all frameworks, including ours, propose that greater weight is given to those observations representing the background only, and smaller or no weight is given to those containing contribution from analyte peaks. Therefore, the aim of the least squares term is to extract the structural information from the neighboring background regions to infer the baseline structure in the analyte region.

The second term of the objective criterion,  $\lambda \int_a^b (\hat{y}''(x))^2 dx$ , is a regularization term. It constrains  $\hat{y}$  to vary smoothly on a global level. Overall, the objective

function trades off fit to the spectral data with the smoothness via the tuning parameter,  $\lambda$ . For smaller values of  $\lambda$  more weight is given to fitting the squared error term of the criterion. When  $\lambda = 0$  the unique minimizer is a natural cubic spline, which will interpolate the original response,  $y_j$ . Conversely, for greater values of  $\lambda$  more weight is given to keeping the curvature small. When  $\lambda \rightarrow \infty$ , the unique minimizer is a second-degree polynomial. A spectrum of  $\lambda$  values ranging from 0 to  $\infty$  will generate a family of models, from interpolation to the parametric polynomial model.

When faced with a problem of how much smoothing should be applied to fit the spectral data on hand, effective degrees of freedom (EDF) represents a more physically interpretable metric to parameterize the regularization of the smoothing spline than  $\lambda$  (Cantoni and Hastie, 2002). Consider writing the  $n$  vector of fitted values,  $\hat{y}$ , as

$$\hat{y} = \mathbf{N}(\mathbf{N}^T \mathbf{N} + \lambda \boldsymbol{\Omega}_N)^{-1} \mathbf{N}^T \mathbf{y} = \mathbf{S}_\lambda \mathbf{y}. \quad (4)$$

Here  $\mathbf{N}$  denotes an  $n \times n$  design matrix of the cubic spline basis functions evaluated at the observed values  $\mathbf{x}_j$  and  $\boldsymbol{\Omega}_N$  is  $\int_l^m N_l''(x) N_m''(x) dx$ . A linear operator referred to as a smoother matrix,  $\mathbf{S}_\lambda$  differentially shrinks influence of  $\mathbf{y}$  toward their alignment with the corresponding basis functions. Consequently, the EDF of a smoothing spline is defined as the sum of eigenvalues of  $\mathbf{S}_\lambda$ :

$$\text{EDF}_\lambda = \sum_{j=1}^n \{\mathbf{S}_\lambda\}_{jj}. \quad (5)$$

EDF is bounded between 2 and  $n$ . If  $\lambda = 0$ ,  $\mathbf{S}_\lambda$  becomes the  $n \times n$  identity matrix, and  $\text{EDF}_\lambda = n$ . Conversely, if  $\lambda = \infty$ ,  $\mathbf{S}_\lambda$  becomes the projection matrix from linear regression on  $\mathbf{x}$ , and  $\text{EDF}_\lambda = 2$ . The advantage of reformulating the smoothing parameter in EDF over  $\lambda$  is that its span is bounded and defined with respect to the number of wavenumbers in the region we want to baseline-correct.

If a desired (target) EDF is defined by a user, smoothing splines models are usually fitted via the backfitting algorithm to search for the actual EDF closest to the target. At convergence, the solution can be formulated as

$$\text{EDF}_A = \text{argmin}_\lambda \left( \text{EDF}_T - \sum_{j=1}^n \{\mathbf{S}_\lambda\}_{jj} \right)^2, \quad (6)$$

where  $\text{EDF}_A$  represents the actual EDF determined from  $\sum_{j=1}^n \{\mathbf{S}_\lambda\}_{jj}$  (Eq. 5) which minimizes the departure from the target EDF,  $\text{EDF}_T$ . The backfitting procedure is implemented in the `smooth.spline` function of the R statistical package (R Core Team, 2014), which we used to develop our baseline correction model. Thus, the user-defined  $\text{EDF}_T$  will form a basis for model parameter solutions from which the optimal parameter,  $\text{EDF}^*$ , will be chosen (Sect. 2.3).

Summarizing in Table 2, we argue that smoothing splines offers a more adaptive and realistic basis for modeling PTFE

**Table 2.** Comparison of key background modeling characteristics pertaining to the proposed and current models

Characteristics	Proposed method	Current method
Functional form	Smoothing splines	Polynomial
Type	Nonparametric	Parametric
Representations	Global (EDF <sub>T</sub> ) and local ( $w_j$ )	Global ( $n$ th degree of a polynomial)
Requires pre-scans?	No	Yes
Requires user's input?	No	For every scan

**Table 3.** The relationship between FT-IR spectrum features and smoothing splines parameters to model those features

Segment	Region type	Spectrum characteristics		Model parameters	
		Wavenumber range (cm <sup>-1</sup> )	Type of modeled baseline	Weights	EDF
1	Background upper	[4000, $W_1$ ]	Fitted	$w_j = 1$	EDF*
	Analyte	[ $W_1$ , $W_2$ ]	Predicted	$w_j = 0$	
	Background lower	[ $W_2$ , 1820]	Fitted	$w_j = 1$	
2	Background upper	[2000, $W_3$ ]	Fitted	$w_j = 1$	EDF*
	Analyte	[ $W_3$ , $W_4$ ]	Predicted	$w_j = 0$	
	Background lower	[ $W_4$ , $W_4 - 1(\Delta\tilde{\nu})$ ] <sup>a</sup>	Fitted	$w_j = 1$	

<sup>a</sup> The lower background region consists of a single wavenumber adjacent to  $W_4$  (Sect. 2.2.1).

variations than the current method by combining local and global representations. We apply smoothing splines to specific segments where each analyte region is sandwiched by neighboring background regions containing a smoothly varying baseline. As a result, each segment then contains an accurate basis for baseline prediction in the analyte region using an optimal smoothing parameter, EDF\*.

## 2.2 FT-IR baseline correction protocol

Using the smoothing splines theory described above, we formalize the baseline correction protocol in Table 3. The weights  $w_j$  from Eq. (3), i.e.,  $w_j = 0$  in the analyte region and  $W_A$  and  $w_j = 1$  in the background region  $W_B$ , are determined by sample-specific bounds for analyte and background regions,  $W_1$  to  $W_4$ . Fig. 2 illustrates a road map for our protocol. In Step 1, we divide a raw spectrum into two segments. Segment 1 includes the domain from 4000 to 1820 cm<sup>-1</sup>, to capture the maximum extent of the background regions surrounding the first analyte region. Segment 2 includes the domain from 2000 to 1500 cm<sup>-1</sup> and captures a sufficient extent of background regions surrounding the second analyte region. We set  $W_2$  to 2220 cm<sup>-1</sup>, which universally marks the start of the carbon dioxide (CO<sub>2</sub>) absorbance band (Pavia et al., 2008).

In Step 2, we perform a geometric transformation, which will be used to determine and verify some of the bounds for analyte and background regions:  $W_1$  in Segment 1 and  $W_3$  to  $W_4$  in Segment 2. As a linear operation, this geometric transformation preserves the actual absorbance magnitudes. Let  $\mathbf{a}$  denote an vector of raw absorbances corresponding to a seg-

ment selected in Step 1 illustrated in Fig. 2. First we rotate  $\mathbf{a}_j$  about a point  $a_1$  such that  $a_1 = a_1^R = a_N^R$ , where  $a_j^R$  denotes the rotated vector element and  $\mathbf{R}$  denotes the corresponding rotation matrix:

$$\mathbf{a}^R = \mathbf{R}\mathbf{a}, \text{ where } \mathbf{R} = \begin{bmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{bmatrix}$$

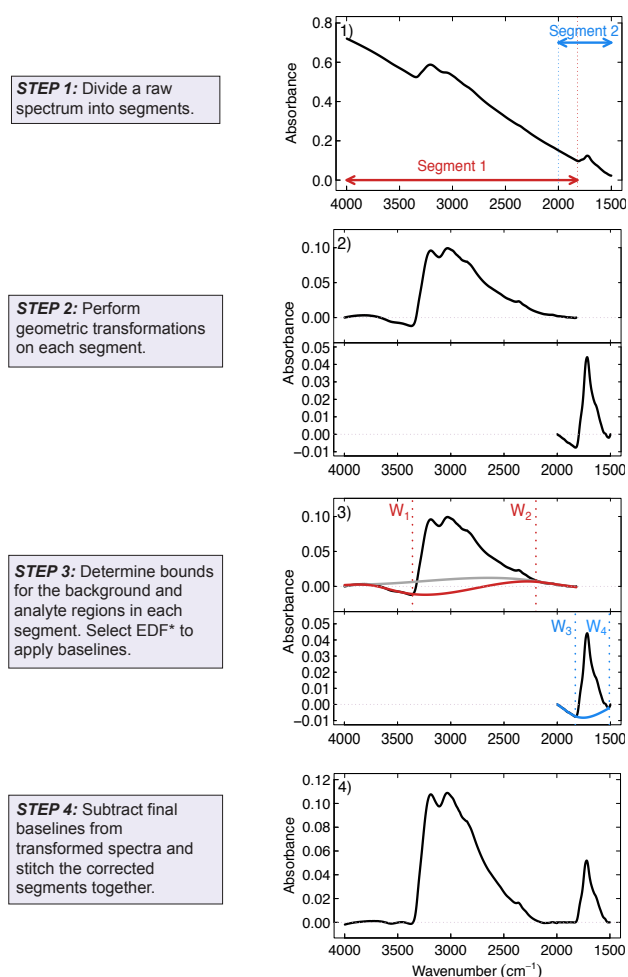
$$\text{and } \theta = \arctan\left(\frac{\nu_N - \nu_1}{a_N - a_1}\right).$$

Second, we translate  $\mathbf{a}_j^R$  such that  $a_1^* = a_N^* = 0$ , where  $\mathbf{a}_j^*$  denotes the resulting translated vector:

$$\mathbf{a}^* = \mathbf{a}^R - a_1^R.$$

Projecting raw absorbances on the local platform axis ( $a_1 = a_N = 0$ ) offers a valuable means of numerically representing a raw spectrum, without appealing to underlying PTFE structural specification. The geometric transformation is a key component in our protocol. First, it allows us to analytically separate background from the analyte in  $W_4$  by determining a local minimum. Second, it provides visually recognizable verification valuable for further method developments, if need be (e.g., precise  $W_1$ ,  $W_3$ , and  $W_4$  are difficult to recognize in raw data in Fig. 1). For instance, the concept is extendible to application developments for baseline correction in the fingerprint region (Day et al., 2010), which is outside the scope of our current study.

In Step 3, we determine specific bounds,  $W_1$  to  $W_4$ , for analyte and background regions,  $W_A$  and  $W_B$ . The benefits of determining sample-specific  $W_1$  and  $W_4$  are twofold. First,



**Figure 2.** (1) Uncorrected spectrum partitioned into two segments: Segment 1, 4000–1820  $\text{cm}^{-1}$  and Segment 2: 2000–1500  $\text{cm}^{-1}$ . (2) Transformed segments with zero first and last absorbance values. (3) Upper panel: initial baseline (gray), final baseline estimated iteratively via a non-negativity constraint (red). Red vertical lines delineate background and analyte regions:  $W_1 = 3360 \text{ cm}^{-1}$  and  $W_2 = 2220 \text{ cm}^{-1}$ . Lower panel: final baseline (blue). Blue vertical lines delineate background and analyte regions:  $W_3 = 1820 \text{ cm}^{-1}$  and  $W_4 = 1520 \text{ cm}^{-1}$ . (4) Resultant corrected spectrum.

certain analytes may be absent from a complex aerosol mixture at hand, thereby increasing  $\mathcal{W}_B$ . Second, higher loadings may lead to broader tails of certain absorption profiles, thereby decreasing  $\mathcal{W}_B$ . Section 2.3.1 details a method to determine these bounds.

In Step 4, we subtract final baselines from transformed segments and stitch the baseline-corrected segments together. In the overlapping region between 2000 and 1820  $\text{cm}^{-1}$ , we use the mean absorbance in the final result. The absorbance between the rightmost background region down to 1500  $\text{cm}^{-1}$  is set to zero.

## 2.3 Selection of model parameters

The problem of selecting model parameters,  $W_1$ – $W_4$  and EDF, carries key implications for the quality of fitted baselines. Our goal is to select model parameters to reproduce the structure of sample-specific PTFE variations while minimizing physically unrealistic FT-IR features, such as negative absorbance from PM spectra or absorbance from blank spectra. Referencing an extensive set of baseline-corrected ambient and blank samples (described in Sect. 2.4), we identify two common physical expectations, to which generated baseline should conform: (1) non-negative analyte absorbance and (2) near-zero blank absorbance.

### 2.3.1 Determining bounds for analyte and background regions

We determine  $W_1$  iteratively for each value of the smoothing parameter to satisfy a non-negativity constraint near the boundaries. An initial (conservative) estimate of  $W_1 = 3720 \text{ cm}^{-1}$  is congruent with our understanding of the absence of absorption bands over the subdomain between 4000 and 3720  $\text{cm}^{-1}$  (Pavia et al., 2008); yet, smaller contributions from certain functional groups, such as alcohol OH, increase the likelihood of negative background absorbance if  $W_1$  remains underspecified. Therefore, we begin with the initial estimate (gray baseline in Fig. 2 Step 3) and iteratively decrease  $W_1$  until the non-negativity constraint is satisfied or until  $W_1$  reaches  $W_2$ . We set  $W_2$  to 2220  $\text{cm}^{-1}$ , which universally marks the start of the  $\text{CO}_2$  absorbance band (Pavia et al., 2008). Similarly, we set  $W_3$  to 1820  $\text{cm}^{-1}$ , which universally marks the start of the carbonyl absorbance band observed in all PM samples.

To accommodate the specifications of individual samples,  $W_4$  is determined as a wavenumber  $\tilde{\nu}$ , for which  $a_j^*$  attains its minimum over the set of candidate values between 1520 and 1600  $\text{cm}^{-1}$ :

$$W_4 = \operatorname{argmin}_j \left\{ a_j^* : \tilde{\nu}_j \in [1520, 1600] \right\}, \quad (7)$$

where  $a_j^*$  are transformed absorbances from Step 2. To minimize the interference from the neighboring alkane peak, starting to absorb around 1510  $\text{cm}^{-1}$  (Pavia et al., 2008), we limit the lower background region to a single wavenumber adjacent to  $W_4$ ,  $W_4 - 1(\Delta\tilde{\nu})$ .

### 2.3.2 Selection of EDF

To parameterize the influence of EDF on the quality of fitted baselines via the two expectations, we derive two EDF-optimizing metrics: (1) a negative absorbance fraction for ambient samples and (2) total normalized absolute blank absorbance for blank filters. We summarize the metrics in Table 4.

The negative absorbance fraction (NAF) represents the contribution of negative analyte absorbance,  $\|a_A\|_1$ , to the



**Table 4.** Relationship between fitted baseline characteristics as a result of varying EDF and EDF-optimizing metrics to represent these characteristics

Segment		Physical criterion	Sample type	Wavenumber range ( $\text{cm}^{-1}$ )	Representation
1	1	Near-zero blank absorbance	Blank	$[4000, 2500], [2200, W_2]$	Total normalized absolute blank absorbance, $\ a_B\ _1^*$
	2	Non-negative analyte absorbance	Ambient	$[W_1, 2500]$	Negative absorbance fraction, NAF
2	1	Near-zero absorbance	Blank	$[2000, 1500]$	Total normalized absolute blank absorbance, $\ a_B\ _1^*$
	2	Non-negative analyte absorbance	Ambient	$[W_3, W_4]$	Negative absorbance fraction, NAF

total analyte absorbance,  $\|a_A\|_1$ :

$$\text{NAF} = \frac{\|a_A - \hat{a}_A\|_1}{\|a_A\|_1} \times 100\%,$$

where  $\|\cdot\|_1$  denotes the 1-norm magnitude of a vector (summation of all absolute values of vector elements). NAF is calculated across the entire wavenumber range in the analyte part of in a given segment, excluding the  $\text{CO}_2$  absorbance band.

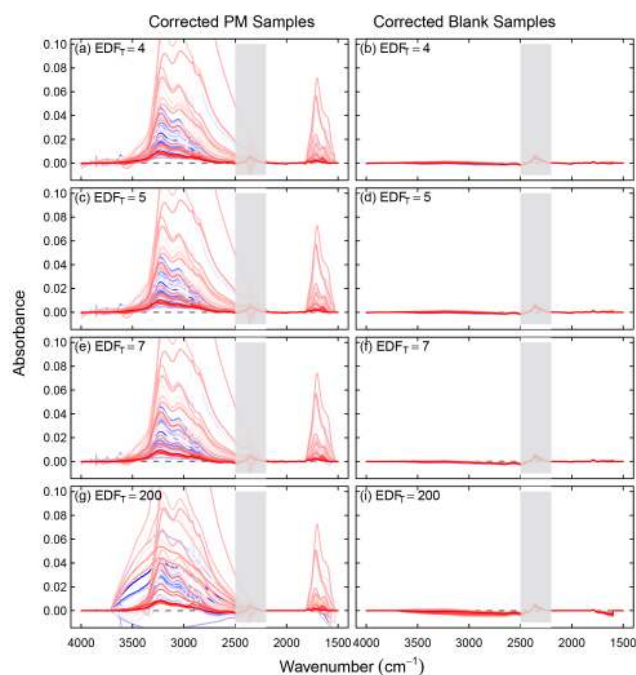
Total normalized absolute blank absorbance,  $\|a_B\|_1^*$ , quantifies the model's departure from the true result, zero absorbance, per wavelength in a given segment. It is calculated as a 1-norm magnitude of blank absorbances,  $\|a_B\|_1$ , normalized by the number of wavenumbers in the corresponding wavenumber range (Table 4),  $n_{\bar{\nu}}$ :

$$\|a_B\|_1^* = \frac{\|a_B\|_1}{n_{\bar{\nu}}}.$$

$\|a_B\|_1^*$  is calculated across the entire wavenumber range in a particular segment excluding the  $\text{CO}_2$  absorbance band. We select  $\text{EDF}^*$  from a range of  $\text{EDF}_T$  by evaluating minima from both  $\|a_B\|_1^*$  and NAF. To that end, Figs. 3 and 4 in Sect. 3.1 present a qualitative and quantitative evaluation for varying  $\text{EDF}_T$  together with  $\text{EDF}^*$  selection.

## 2.4 Experimental data

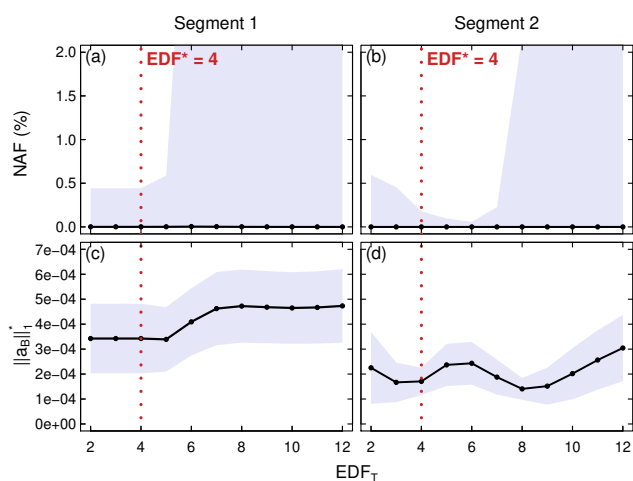
We apply smoothing splines baseline correction to 794 particulate matter ( $\leq 2.5 \mu\text{m}$  in diameter,  $\text{PM}_{2.5}$ ) samples collected on PTFE filters and 54 blank PTFE filters. The particulate matter samples were collected at IMPROVE sites on every third day in 2011. IMPROVE absorption spectra had been used in a previous studies (Ruthenburg et al., 2014; Dillner and Takahama, 2015a, b) which detail the mechanics of FT-IR spectra collection. More important for this study is the level of spectral preparation applied prior to the background correction. Following the practice established in Dillner and Takahama (2015a, b) we use unmodified spectra in which values interpolated during the zero-filling process were removed. Prior to applying the smoothing splines baseline, we truncate the original wavenumber domain between



**Figure 3.** 54 randomly selected ambient samples (left) and 54 blank samples (right) corrected by varying  $\text{EDF}_T$ . Each spectrum is color-differentiated. The  $\text{CO}_2$  absorption band between 2500 and 2220  $\text{cm}^{-1}$  not associated with PM composition is shaded in color. The x axis ranges from 4000 to 1500  $\text{cm}^{-1}$  in both left and right panels.

4000 and 420  $\text{cm}^{-1}$  to capture the subdomain between 4000 and 1500  $\text{cm}^{-1}$  (1944 wavenumbers). As a reference, the same subdomain is used in the polynomial method (Takahama et al., 2013). In contrast to Takahama et al. (2013), we do not apply smoothing to remove water vapor interference and carbon dioxide to minimize the number of preprocessing steps.





**Figure 4.** Median NAF in Segment 1 (a) and Segment 2 (b) calculated from 794 ambient samples (black points). Lower and upper bounds of shaded areas denote 3rd and 97th percentiles. Mean  $\|a_B\|_1^*$  for  $2 \leq \text{EDF} \leq 12$  in Segment 1 (c) and Segment 2 (d), calculated from 54 IMPROVE 2011 laboratory blank samples (black points). Shaded areas denote 3 standard deviations from the mean. In all panels, the black line is drawn to capture the overall trend. While we select the interval  $2 \leq \text{EDF}_T \leq 12$  specifically to highlight each metric's minima, we present results from the entire interval  $2 \leq \text{EDF}_T \leq n$  for completeness in Fig. S2.

## 2.5 Applications for model evaluation

### 2.5.1 Cluster analysis

Cluster analysis with FT-IR measurements generates natural categories for PM samples based on spectral similarity. These categories can represent mixture classes of chemically complex aerosols, and their association with meteorological and collocated measurements has been shown to provide complementary information for source apportionment (Takahama et al., 2011; Corrigan et al., 2013). For this purpose, each spectrum is SSB-corrected to isolate the analyte contribution to the IR absorbance, normalized by its 2-norm magnitude to emphasize variation in relative composition rather than absolute concentration, and grouped according to the hierarchical clustering algorithm of Ward (1963). There are inherent differences in the vapor artifacts between the PB-corrected and SSB-corrected spectra that are not critical for the algorithms used for quantification of functional groups, or TOR organic and elemental carbon but influence clusters formed from the naïve clustering approach described above. As the PB-corrected signal requires differencing the IR spectrum of the PTFE before and after sample collection, water vapor and CO<sub>2</sub> signals remaining in the PB-corrected spectra represent differences in concentrations present in the chamber during both scans, whereas SSB-corrected spectra only contain the amount present in the latter. Therefore, regions where these artifacts are present ( $\tilde{\nu} > 3600 \text{ cm}^{-1}$  and

$\tilde{\nu} < 2400 \text{ cm}^{-1}$  in Segment 1) are excluded from the normalization and clustering, though some water vapor artifact overlapping with analyte absorption remains in Segment 2. In addition, seven samples with specific features or low signal-to-noise ratios are removed from the set prior to the clustering as they are not well discriminated by the algorithm, or influences the grouping of the rest of the spectra.

### 2.5.2 Peak fitting

We apply the peak-fitting algorithm based on parameter constraints described by Takahama et al. (2013) to both SSB- and PB-corrected spectra and evaluate the differences between two baseline correction methods by comparing peak areas. Peak areas correspond to integrated absorbances from line shapes fitted for alcohol COH, carboxylic COH, alkane CH, carbonyl CO, and amine NH. We examine the comparability and implications of replacing the PB correction approach with SSB correction in future analyses of this type.

### 2.5.3 Prediction of TOR organic carbon (OC) and elemental carbon (EC)

Dillner and Takahama (2015a, b) recently demonstrated that collocated PTFE samples analyzed by FT-IR and quartz fiber filters analyzed by TOR can be used to build calibration models that predict TOR-equivalent OC and EC concentrations from new FT-IR spectra. One of several calibration models with accuracy and precision on a par with TOR precision can be constructed when the concentration range and composition of carbonaceous samples in the calibration set approximately resemble those in the test (challenge) set. For this work, we use an identical procedure as described by Dillner and Takahama (2015a, b) for building calibration and test sets from 794 IMPROVE 2011 samples chronologically stratified within each site. The spectra are SSB-corrected and calibration and test samples are drawn to contain two-thirds and one-third of the entire set, respectively. Only TOR OC and EC predictions necessitate dividing the data set into calibration and test subsets; the previous two applications, clustering and peak fitting, are applied to the entire data set.

## 3 Results

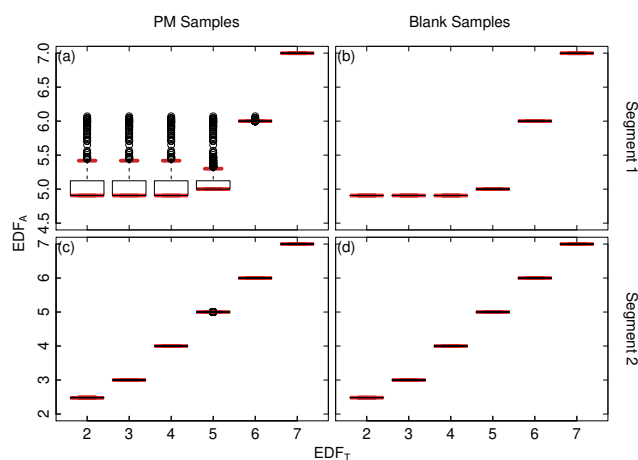
At the physical level, we evaluate the feasibility of our model by selecting the optimal smoothing parameters in Sect. 3.1 and by presenting the sample-specific bounds for analyte and background regions in Sect. 3.2. At the application level, we begin our evaluation of smoothing splines baseline-corrected spectra with visual and cluster analysis in Sect. 3.3, followed by functional group quantification analysis in Sect. 3.4, and predicted TOR OC and EC analysis in Sect. 3.5.

### 3.1 EDF selection

Qualitatively, in Fig. 3 we compare the behavior of PM (left panel) and blank samples (right panel) using varying  $\text{EDF}_T$  (4, 5, 7, and 200, from top to bottom). In this analysis we used all 54 blank samples and randomly sampled 54 out of 794 PM samples to keep the counts equal and allow for representative cross-comparison. The trend from top to bottom shows both PM and blank samples exhibit increasing sensitivity to the amount of smoothing applied. With increasing  $\text{EDF}_T$ , baseline-corrected ambient spectra begin to exhibit negative analyte absorbance (left column). Simultaneously, baseline-corrected blanks in the region at  $3700\text{--}2500$  and  $1820\text{--}1600\text{ cm}^{-1}$  begin to depart from our target, zero absorbance (right column).

Quantitatively, in Fig. 4 we evaluate the impact of  $\text{EDF}_T$  on negative absorbance fraction metric, NAF, (top panel) and total normalized absolute blank absorbance metric,  $\|\mathbf{a}_B\|_1^*$ , (bottom panel) in segments 1 and 2 (left and right panel). Horizontal panels share the same  $x$  axis and vertical panels share the same  $y$  axis to allow for representative cross-comparison. Therefore, each plot in the matrix in Fig. 4 corresponds to a unique condition in terms of a metric and segment. Starting from Fig. 4 A (top left), we find that any  $\text{EDF}_T$  between 2 and 4 minimizes median NAF and its variance simultaneously: median  $\text{NAF} \approx 0.0\%$  and variance =  $0.44\%$ . Moving down to Fig. 4c (bottom left), we look at the effect of  $\text{EDF}_T$  on blank absorbance in Segment 1. We find that any  $\text{EDF}_T$  between 2 and 4 generates very low  $\|\mathbf{a}_B\|_1^*$ : mean  $\|\mathbf{a}_B\|_1^* = 3.42 \times 10^{-4}$  and  $3\sigma$  (the extent of shaded areas) =  $2.79 \times 10^{-4}$ . Technically, the minimum variance in  $\|\mathbf{a}_B\|_1^*$  occurs for  $\text{EDF}_T = 5$  but the difference is less than  $1.5\%$ . Of the two metrics, we prefer to minimize NAF over  $\|\mathbf{a}_B\|_1^*$  as NAF represents a more robust metric (the sample size is an order of magnitude greater and in future applications the choice of  $\text{EDF}_T$  will likely affect disproportionately more PM samples than blank samples). To finalize the choice of  $\text{EDF}_T$  from  $2 \leq \text{EDF}_T \leq 4$ , we now consider how these  $\text{EDF}_T$  values compare to  $\text{EDF}_A$  obtained by the smoothing splines algorithm from Eq. (6). We plot the distributions of  $\text{EDF}_A$  given  $\text{EDF}_T$  in Segment 1 using all 794 PM samples and 54 blank samples in Fig. 5a and b.

The extensive number of knots to form bases for fitting splines (that is, wavenumbers in observed absorbances used for fitting:  $x_j$  for which  $w_j \neq 0$  from Eq. 1) creates limitations on minimum achievable EDF. This is particularly acute when  $\text{EDF}_T$  is low ( $< 7$  in Segment 1 and  $< 3$  in Segment 2). For instance, if we apply baselines with  $\text{EDF}_T = 4$  in Segment 1 (Fig. 5a and b), the distribution of  $\text{EDF}_A$  will span between 4.9 and 6.1 depending on the number of basis-forming knots (Fig. S3). However, applying baselines with target  $\text{EDF} < 4$  will lead to identical  $\text{EDF}_A$  results, confirming that the set of  $\text{EDF}_A$  between 4.9 and 6.1 is indeed the minimum achievable EDF in the search domain. Therefore,



**Figure 5.** Box-and-whisker plots representing distributions of  $\text{EDF}_A$  for a given  $\text{EDF}_T$  used in Segment 1 (a, b) and Segment 2 (c, d) in both PM ( $n = 794$ ) and blank samples ( $n = 54$ ). Median and whiskers in each box-and-whisker plot are highlighted in red.

out of  $\text{EDF}_T$  candidates for  $\text{EDF}^*$  we choose 4 as it represents the actual, true parameters most accurately; given  $\text{EDF}^* = 4$  we obtain  $\text{EDF}_A \in ([4.9, 6.1])$  for PM samples and  $[4.9, 4.9]$  for blank samples.

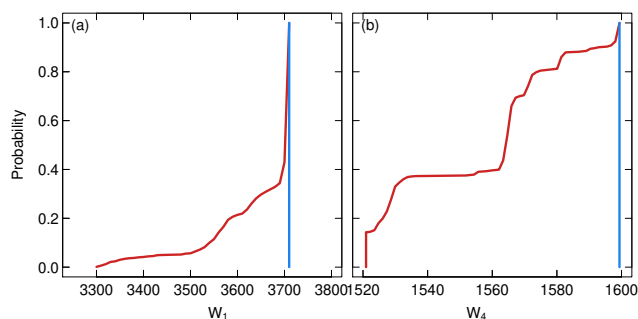
In Fig. 4b (top right) we start by limiting the evaluation in Segment 2 to  $\text{EDF}_T$  for which NAF variance is greater than  $0.22\%$  (roughly a half of the value from the best-fit model in Segment 1). This leaves us with  $4 \leq \text{EDF}_T \leq 7$ . Out of this subset, we find selecting 4 as  $\text{EDF}_T$  minimizes  $\|\mathbf{a}_B\|_1^*$  in Fig. 4d; mean  $\|\mathbf{a}_B\|_1^* = 1.71 \times 10^{-4}$ ,  $3\sigma = 1.06 \times 10^{-4}$ . Additionally, and importantly, 4 represents the most parsimonious solutions without visually distorting the blank baseline and shape of the PM peaks (Fig. 3). By selecting  $\text{EDF}^* = \text{EDF}_T = 4$ , now the actual EDF parameters match the target EDF parameter (Fig. 5a and d).

### 3.2 $W_1$ and $W_4$

Figure 6 presents empirical cumulative distributions' functions of  $W_1$  and  $W_4$  from PM and blank samples. Distribution of  $W_1$  in PM samples spans values between  $3300$  and  $3710\text{ cm}^{-1}$ , with  $50\%$  of samples having  $W_1 > 3700\text{ cm}^{-1}$ , reflecting sample-specific PM mixture composition (illustration of spectra in Fig. 3a).  $W_1$  in blank samples was determined to be  $3710\text{ cm}^{-1}$  (Fig. 3b). Distribution of  $W_4$  in PM samples spans values between  $1520$  and  $1600\text{ cm}^{-1}$ , reflecting sample-specific ammonium absorbance width (Fig. 3a).  $W_4$  in blank samples was determined to be  $1600\text{ cm}^{-1}$ , which is consistent with our physical expectation about zero amine absorbance (Fig. 3b).

### 3.3 Cluster analysis

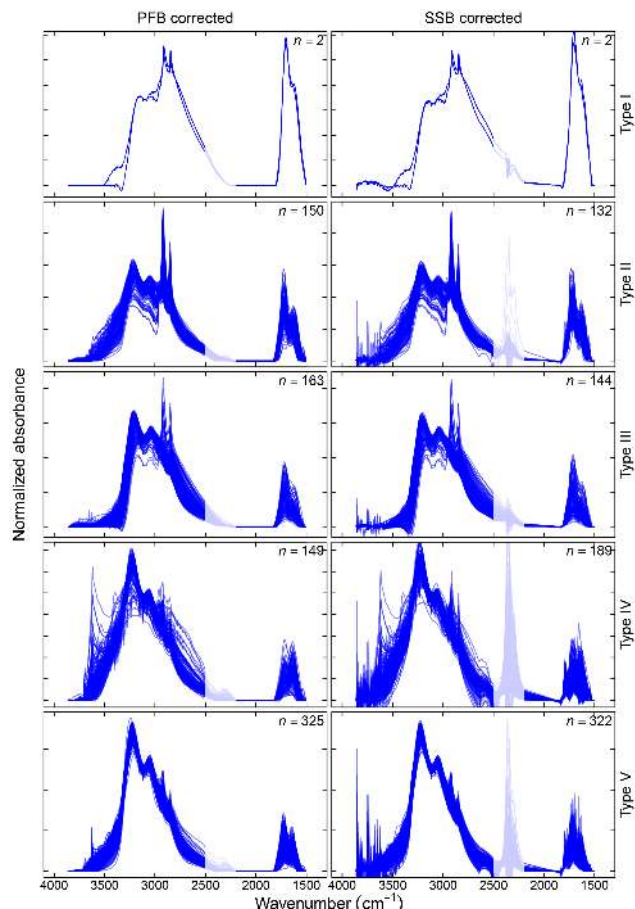
The number of samples from SSB-corrected spectra not sharing the same relative labeling as those from PB-corrected



**Figure 6.** Empirical cumulative distribution functions representing distributions of  $W_1$  and  $W_4$  in PM samples ( $n = 794$ ) in red and blank samples ( $n = 54$ ) in blue.

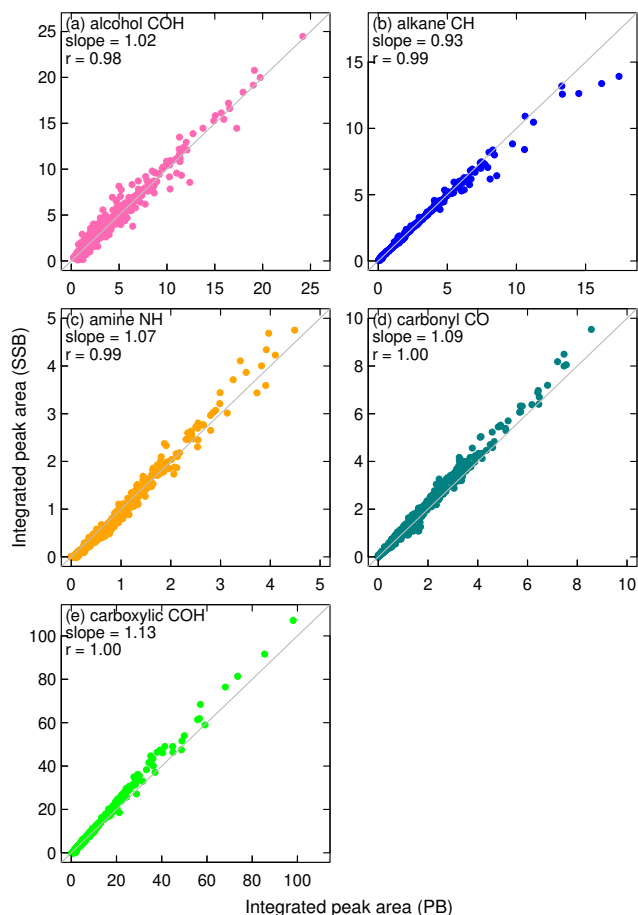
spectra varies with the total number of clusters used to partition the spectra set. Figure S1 in the Supplement shows that the discrepancy for 787 samples increases as the set is partitioned into a larger number of clusters. The difference in sample labeling varies between 5 % for two clusters and 11 % for five clusters; the increase is observed for larger number of clusters because spectra are grouped according to finer variations in their features. Feature (wavenumber) selection and advanced algorithms can lead to more robust clustering that is less sensitive to small variations in spectra (Hastie et al., 2009), but visual comparisons of spectra in the present form of aggregation can provide useful interpretations as discussed below. The inter-cluster differences will further depend on the number of clusters and the type of clustering algorithm. Since there is no absolute reference for baseline-corrected spectra, these discrepancies speak to the differences between two candidate methods.

Figure 7 shows spectra from the two baseline correction algorithms grouped into categories using the approach described in Sect. 2. Type I spectra are selected manually, and Types II–V are determined by a four-cluster solution by hierarchical clustering (with a discrepancy rate between PB and SSB of 10 %). Type I spectra display low absorbance in the alcohol COH region, visible methylene paired peaks (2920 and 2850  $\text{cm}^{-1}$ ) from CH<sub>2</sub> bonds present in vegetative detritus (Hawkins and Russell, 2010), and the largest absorbance in the carbonyl CO region (centered near 1700  $\text{cm}^{-1}$ ) compared to the rest of the sample spectra. This spectra type indicates a dominant contribution from biomass burning aerosol spectra (Hawkins and Russell, 2010; Takahama et al., 2011). These two samples were collected in St. Marks, FL, during January and February; fire burning is prescribed near this location during January through May of each year. Type II spectra also contain sharp methylene peaks but also stronger absorption above 3100  $\text{cm}^{-1}$  associated with alcohol COH and less pronounced carbonyl CO absorption. Sixty percent of the 132 SSB-corrected spectra are found in Phoenix, AZ, so this is interpreted to be associated with urban aerosol (we note that Phoenix samples may be overrepresented in



**Figure 7.** Cluster membership for polynomial and smoothing splines methods. The region between 2500 and 2200  $\text{cm}^{-1}$  is masked to indicate the region of CO<sub>2</sub> absorption not associated with aerosol composition.

this spectra set as two sampling sites out of the seven analyzed in this work are located in this city). Similar features have been found in spectra from the urban environment of Mexico City (Liu et al., 2009). Type V contains spectra for which peaks near 3200–3100  $\text{cm}^{-1}$  are most prominent, indicating the significant presence of ammonium. These features have commonly been reported in fossil fuel burning samples or factor analysis components (Hawkins and Russell, 2010; Takahama et al., 2011; Guzman-Morales et al., 2014) that have been assigned by correlation with combustion tracers (e.g., V, Cr, Ni, Zn, As) and back trajectory analyses. These aerosols presumably arise from a combination of aged background aerosol and aerosols produced locally in the presence of high oxidant concentrations of polluted environments (Liu et al., 2011). However, 87 % of the 322 SSB-corrected Type V samples are found in the five non-urban sites, suggesting that in this data set this spectroscopic signature is more indicative of aged secondary aerosol. Ammonium concentrations are often temporally correlated with oxidized organic aerosol (e.g., Jimenez et al., 2009; Lanz



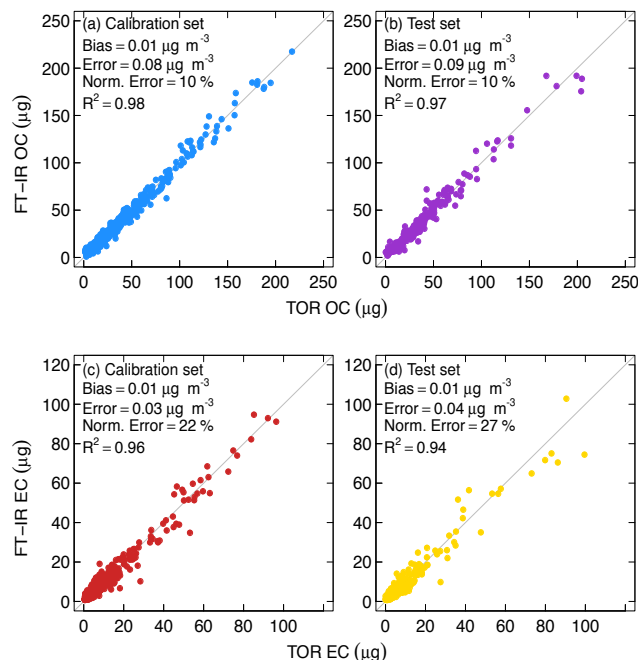
**Figure 8.** Integrated peak area corresponding to different functional groups (a–e) from polynomial baseline and smoothing splines baseline-corrected spectra. Slope magnitudes represent the slope of the regressed line. The silver line represents a one to one line.

et al., 2010) which increases in abundance toward rural areas (Zhang et al., 2007). Types III and IV share some combination of features with types I, II, and V, with the ammonium peak near  $3200\text{ cm}^{-1}$  more visible in type IV and larger contributions from methylene peaks visible in type III. The peak near  $3700\text{ cm}^{-1}$  present in several type IV spectra is suggestive of phenolic compounds also present in biogenic aerosol (Bahadur et al., 2010).

This analysis demonstrates that the new SSB correction method can generate spectra similar in profile to PB-corrected spectra used in past studies, providing a basis for further mixture analysis.

### 3.4 Peak fitting analysis

Figure 8 presents integrated absorbances for alcohol COH, carboxylic COH, alkane CH, carbonyl CO, and amine NH quantified from PB and SSB-corrected spectra. For all functional groups but carboxylic COH the discrepancy between the two methods is  $< 10\%$  (the slope of the regressed line



**Figure 9.** Predicted FT-IR OC vs. measured TOR OC using smoothing splines-corrected spectra for (a) calibration set ( $n = 517$ ) and (b) test set ( $n = 268$ ). Predicted FT-IR EC vs. measured TOR EC using smoothing splines-corrected spectra for (c) the calibration set ( $n = 501$ ) and (d) the test set ( $n = 268$ ).

$< 1 \pm 0.1$ ). The difference is on the same order of magnitude as the cluster discrepancy rate. The bias in carboxylic COH fitting is likely due to the fact that its line shape was fixed specifically to the PB-corrected spectra (Takahama et al., 2013), and is more sensitive to the absorption profile to which it is fitted than the Gaussian peaks with adjustable parameters used for fitting other functional groups. The bias in may be alleviated by rederiving the carboxylic COH line shape for the smoothing splines method, or applying an adjusted molar absorption coefficient. The bias of 13 % is on the order of variation in absorption coefficients of carboxylic COH estimated for different organic acid compounds, and also within uncertainty for an absorption coefficient estimated from the mean of these values (Takahama et al., 2013).

### 3.5 Prediction of TOR organic and elemental carbon

Figure 9 presents performance metrics from TOR OC and TOR EC predictions obtained from SSB-corrected spectra. All fits are characterized by high coefficients of variations ( $R^2 \geq 0.94$ ) and near-zero bias ( $\leq 0.01\text{ }\mu\text{g m}^{-3}$ ), demonstrating accurate predictions. With respect to predicted TOR OC, performance metrics from the test set (Fig. 9b) are on a par with those obtained from raw spectra and PB-corrected spectra. Specifically, error ( $0.09\text{ }\mu\text{g m}^{-3}$ ) and normalized error (10 %) are on the same order as those obtained from raw spectra (error of  $0.08\text{ }\mu\text{g m}^{-3}$ , normalized error of 11 %) and

**Table 5.** MDL and precision for FT-IR OC and TOR OC.

Carbon type	Metric	TOR	FT-IR raw spectra <sup>d</sup>	FT-IR PB-corrected spectra <sup>d</sup>	FT-IR SSB-corrected spectra
OC	MDL ( $\mu\text{g m}^{-3}$ ) <sup>b, c</sup>	0.05	0.14, [0.11, 0.28]	0.11, [0.08, 0.17]	0.06, [0.04, 0.09]
	% below MDL	1.5	2.6	0.7	0.0
	Precision ( $\mu\text{g m}^{-3}$ ) <sup>b</sup>	0.14	0.12	0.21	0.06
	Mean blank ( $\mu\text{g}$ )	NR <sup>e</sup>	$0.1 \pm 1.5$	$1.9 \pm 1.2$	$0.1 \pm 0.6$
EC	MDL ( $\mu\text{g m}^{-3}$ ) <sup>b, c</sup>	0.01	0.02, [0.01, 0.02]	0.01, [0.00, 0.01]	0.01, [0.01, 0.02]
	% below MDL	3	1	2	1
	Precision ( $\mu\text{g m}^{-3}$ ) <sup>b</sup>	0.11	0.04	0.06	0.06
	Mean blank ( $\mu\text{g}$ )	NR <sup>e</sup>	$0.06 \pm 0.17$	$0.08 \pm 0.15$	$0.01 \pm 0.12$

<sup>b</sup> Concentration units of  $\mu\text{g m}^{-3}$  for MDL and precision are based on the IMPROVE volume of  $32.8 \text{ m}^3$ . <sup>c</sup> Numbers inside the interval denote 95 % confidence intervals on the estimate. <sup>d</sup> (Dillner and Takahama, 2015a, b). <sup>e</sup> Not reported.

PB-corrected spectra (error of  $0.08 \mu\text{g m}^{-3}$ , normalized error of 12 %) (Dillner and Takahama, 2015a). In Table 5 we show that applying SSB leads to a lower minimum detection limit (MDL) of  $0.06 \mu\text{g m}^{-3}$ , which leaves no samples below MDL. This is statistically different from the no baseline case, where MDL is  $0.14 \mu\text{g m}^{-3}$ . Precision ( $0.06 \mu\text{g m}^{-3}$ ) obtained from SSB-corrected spectra is on the same order as that obtained from raw ( $0.12 \mu\text{g m}^{-3}$ ) or PB-corrected spectra ( $0.21 \mu\text{g m}^{-3}$ ).

Likewise, TOR EC performance metrics from the test set (Fig. 9d) are on a par with those obtained from raw spectra and PB-corrected spectra. Specifically, error ( $0.04 \mu\text{g m}^{-3}$ ) and normalized error (27 %) are on the same order as those obtained from raw spectra (error of  $0.02 \mu\text{g m}^{-3}$ , normalized error of 21 %) and PB-corrected spectra (error of  $0.04 \mu\text{g m}^{-3}$ , normalized error of 24 %) (Dillner and Takahama, 2015b). Table 5 shows that MDL ( $0.01 \mu\text{g m}^{-3}$ ) obtained from SSB-corrected spectra is similar to MDL obtained from raw or PB-corrected spectra (all  $\leq 0.02 \mu\text{g m}^{-3}$ ).

In summary, SSB-corrected spectra OC and EC predictions from blank and ambient samples are as accurate and precise as those from raw or PB-corrected spectra. No additional bias is introduced as a result of SSB correction implementation. However, the reduction in the complexity of baseline correction is amenable for scaling up to a large number of samples. To some extent, PLS is a robust regression method and is able to effectively remove contributions to the signal which are not related to the target analyte. While individual predictions vary, we show in Fig. S4 that the quality of TOR OC and EC predictions is not statistically affected by the choice of EDF between 2 and 30.

#### 4 Conclusions

Within the past few years the guided polynomial baseline-corrected algorithm has been applied to characterize the ambient FT-IR spectra by classifying mixtures (Russell et al., 2009; Liu et al., 2009; Takahama et al., 2011; Ruthenburg

et al., 2014), quantifying organic functional groups (Takahama et al., 2013), and predicting TOR OC and EC (Dillner and Takahama, 2015a, b). Here our results demonstrate that similar estimates (cluster discrepancy rate of 10 %, functional group difference  $\leq 13$  %, and  $R^2 \geq 0.94$  %, bias  $\leq 0.01 \mu\text{g m}^{-3}$ , error  $\leq 0.04 \mu\text{g m}^{-3}$  in TOR OC and EC predictions) can be obtained using a new, automated baseline correction protocol. Contrasting with the polynomial method, this paper detailed the statistical framework, which applies nonparametric smoothing splines to model sample-specific PTFE variations, reduces the number of free parameters from four to one, and selects the parameter by minimizing two evaluation metrics: negative analyte absorbance and blank absolute absorbance. The proposed protocol unifies and simplifies many of the steps in existing techniques while eliminating the need for expert intervention in manually adjusting background regions specific to each sample. More importantly, the automated solution allows us and future users to evaluate its analytical reproducibility while minimizing reducible bias due to current default background regions or a variability in human judgement in adjusting these regions. The solution was developed as a direct response to the growing body of research on statistical applications for characterization of FT-IR atmospheric aerosol samples collected on PTFE filters and a rising interest in analyzing FT-IR samples collected by air quality monitoring networks. As a result, we anticipate that the model will enable FT-IR researchers and data analysts to quickly and reliably analyze a large amount of data. Although the exact reduction in user time may be difficult to generalize due to high variability across different users, we reason that the following approximation applies. Qualitatively, if  $N$  values are considered for each free parameter in each method, then the amount of time for expert examination of each model solution scales up with  $N^4$  for the polynomial method (due to four boundary points as free parameters) and  $N$  for the smoothing splines method (due to 1 EDF parameter). Additionally, and importantly, the evaluation metrics, which we established in this manuscript,



have been shown to sufficiently simplify the parameter selection process for users of any level of experience.

One of the important avenues for future research include implementing sample-specific EDF when the parameter choice affects model performance significantly across samples. As Fig. 3 demonstrates, the individual differences between EDF<sub>T</sub> 4 and 7 in Segment 1 are negligible; on the whole these parameters do a very similar job in minimizing the undesirable quantities (negative analyte absorbance and blank absorbance) in Fig. 4. However, we anticipate that we and other FT-IR analysts may benefit from sample-specific EDF when analyzing data sets collected under different conditions, be it a different sampling flow rate or filter type. Another line of future work may include extending this approach to the remaining part of mid-IR absorbance spectrum (1500–420 cm<sup>-1</sup>). The fingerprint region contains important functional groups (Day et al., 2010), such as organonitrates, which can benefit from an adaptive baseline correction algorithm. As demonstrated in this paper, the general strategy of (1) segmenting baseline regions of interest such that they contain a smoothly varying (or uniformly sloping) baseline and (2) using conservative estimates for background regions, and (3) using FT-IR physical criteria (such as minimal blank absorbance, non-negative analyte/background absorbance, and no baseline discontinuities) for parameter selection can provide a good starting point for these tasks.

The automated smoothing splines baseline correction method has been implemented in R package APRLssb and can be accessed at this repository: <https://bitbucket.org/stakahama/aprlssb> by contacting the corresponding author.

**The Supplement related to this article is available online at doi:10.5194/amt-9-2615-2016-supplement.**

*Acknowledgements.* The authors acknowledge EPFL discretionary funding and funding from IMPROVE program and EPA (National Park Service cooperative agreement P11AC91045). We thank Matteo Reggente and Andrew Weakley for helpful conversations.

Edited by: G. Phillips

## References

- Bacsik, Z., Mink, J., and Keresztury, G.: FTIR spectroscopy of the atmosphere. I. Principles and methods, *Appl. Spectrosc. Rev.*, 39, 295–363, doi:10.1081/asr-200030192, 2004.
- Bahadur, R., Uplinger, T., Russell, L. M., Sive, B. C., Cliff, S. S., Millet, D. B., Goldstein, A., and Bates, T. S.: Phenol Groups in Northeastern US Submicrometer Aerosol Particles Produced from Seawater Sources, *Environ. Sci. Technol.*, 44, 2542–2548, doi:10.1021/es9032277, 2010.

- Baker, M. J., Trevisan, J., Bassan, P., Bhargava, R., Butler, H. J., Dorling, K. M., Fielden, P. R., Fogarty, S. W., Fullwood, N. J., Heys, K. A., Hughes, C., Lasch, P., Martin-Hirsch, P. L., Obinaju, B., Sockalingum, G. D., Sule-Suso, J., Strong, R. J., Walsh, M. J., Wood, B. R., Gardner, P., and Martin, F. L.: Using Fourier transform IR spectroscopy to analyze biological materials, *Nat. Protoc.*, 9, 1771–1791, doi:10.1038/nprot.2014.110, 2014.
- Cantoni, E. and Hastie, T.: Degrees-of-freedom tests for smoothing splines, *Biometrika*, 89, 251–263, doi:10.1093/biomet/89.2.251, 2002.
- Corrigan, A. L., Russell, L. M., Takahama, S., Äijälä, M., Ehn, M., Junninen, H., Rinne, J., Petäjä, T., Kulmala, M., Vogel, A. L., Hoffmann, T., Ebben, C. J., Geiger, F. M., Chhabra, P., Seinfeld, J. H., Worsnop, D. R., Song, W., Auld, J., and Williams, J.: Biogenic and biomass burning organic aerosol in a boreal forest at Hyytiälä, Finland, during HUMPPA-COPEC 2010, *Atmos. Chem. Phys.*, 13, 12233–12256, doi:10.5194/acp-13-12233-2013, 2013.
- Coury, C. and Dillner, A. M.: A method to quantify organic functional groups and inorganic compounds in ambient aerosols using attenuated total reflectance FTIR spectroscopy and multivariate chemometric techniques, *Atmos. Environ.*, 42, 5923–5932, doi:10.1016/j.atmosenv.2008.03.026, 2008.
- Day, D. A., Liu, S., Russell, L. M., and Ziemann, P. J.: Organonitrate group concentrations in submicron particles with high nitrate and organic fractions in coastal southern California, *Atmos. Environ.*, 44, 1970–1979, doi:10.1016/j.atmosenv.2010.02.045, 2010.
- de Rooij, J. J. and Eilers, P. H. C.: Mixture models for baseline estimation, *Chemometr. Intell. Lab.*, 117, 56–60, doi:10.1016/j.chemolab.2011.11.001, 2012.
- Dillner, A. M. and Takahama, S.: Predicting ambient aerosol thermal-optical reflectance (TOR) measurements from infrared spectra: organic carbon, *Atmos. Meas. Tech.*, 8, 1097–1109, doi:10.5194/amt-8-1097-2015, 2015a.
- Dillner, A. M. and Takahama, S.: Predicting ambient aerosol thermal-optical reflectance measurements from infrared spectra: elemental carbon, *Atmos. Meas. Tech.*, 8, 4013–4023, doi:10.5194/amt-8-4013-2015, 2015b.
- Drouet, L., Bosetti, V., and Tavoni, M.: Selection of climate policies under the uncertainties in the Fifth Assessment Report of the IPCC, *Nat. Clim. Change*, 5, 937–940, doi:10.1038/nclimate2721, 2015.
- Felten, J., Hall, H., Jaumot, J., Tauler, R., de Juan, A., and Gorzsas, A.: Vibrational spectroscopic image analysis of biological material using multivariate curve resolution-alternating least squares (MCR-ALS), *Nat. Protoc.*, 10, 217–240, doi:10.1038/nprot.2015.008, 2015.
- Fourmond, V., Hoke, K., Heering, H. A., Baffert, C., Leroux, F., Bertrand, P., and Leger, C.: SOAS: A free program to analyze electrochemical data and other one-dimensional signals, *Bioelectrochemistry*, 76, 141–147, doi:10.1016/j.bioelechem.2009.02.010, 2009.
- Frossard, A. A., Russell, L. M., Burrows, S. M., Elliott, S. M., Bates, T. S., and Quinn, P. K.: Sources and composition of submicron organic mass in marine aerosol particles, *J. Geophys. Res.-Atmos.*, 119, 12977–13003, doi:10.1002/2014jd021913, 2014.
- Gilardoni, S., Russell, L. M., Sorooshian, A., Flagan, R. C., Seinfeld, J. H., Bates, T. S., Quinn, P. K., Allan, J. D., Williams, B., Goldstein, A. H., Onasch, T. B., and Worsnop, D. R.: Re-



- gional variation of organic functional groups in aerosol particles on four US east coast platforms during the International Consortium for Atmospheric Research on Transport and Transformation 2004 campaign, *J. Geophys. Res.-Atmos.*, 112, D10S27, doi:10.1029/2006JD007737, 2007.
- Goldstein, A. H. and Galbally, I. E.: Known and unexplored organic constituents in the earth's atmosphere, *Environ. Sci. Technol.*, 41, 1514–1521, doi:10.1021/es072476p, 2007.
- Golotvin, S. and Williams, A.: Improved baseline recognition and modeling of FT NMR spectra, *J. Magn. Reson.*, 146, 122–125, doi:10.1006/jmre.2000.2121, 2000.
- Griffiths, P., Shao, L., and Leytem, A.: Completely automated open-path FT-IR spectrometry, *Anal. Bioanal. Chem.*, 393, 45–50, doi:10.1007/s00216-008-2429-6, 2009.
- Griffiths, P. R. and De Haseth, J. A.: *Fourier transform infrared spectrometry*, vol. 171, John Wiley & Sons, New York, 201–205, 2007.
- Guzman-Morales, J., Frossard, A., Corrigan, A., Russell, L., Liu, S., Takahama, S., Taylor, J., Allan, J., Coe, H., Zhao, Y., and Goldstein, A.: Estimated contributions of primary and secondary organic aerosol from fossil fuel combustion during the CalNex and Cal-Mex campaigns, *Atmos. Environ.*, 88, 330–340, doi:10.1016/j.atmosenv.2013.08.047, 2014.
- Hamilton, J. F., Webb, P. J., Lewis, A. C., Hopkins, J. R., Smith, S., and Davy, P.: Partially oxidised organic components in urban aerosol using GCXGC-TOF/MS, *Atmos. Chem. Phys.*, 4, 1279–1290, doi:10.5194/acp-4-1279-2004, 2004.
- Hastie, T., Tibshirani, R., and Friedman, J.: *The elements of statistical learning: data mining, inference, and prediction*, Springer Verlag, New York, NY, USA, 2009.
- Hawkins, L. N. and Russell, L. M.: Oxidation of ketone groups in transported biomass burning aerosol from the 2008 Northern California Lightning Series fires, *Atmos. Environ.*, 44, 4142–4154, doi:10.1016/j.atmosenv.2010.07.036, 2010.
- Isaksen, I. S. A., Granier, C., Myhre, G., Berntsen, T. K., Dalsson, S. B., Gauss, M., Klimont, Z., Benestad, R., Bousquet, P., Collins, W., Cox, T., Eyring, V., Fowler, D., Fuzzi, S., Jockel, P., Laj, P., Lohmann, U., Maione, M., Monks, P., Prevot, A. S. H., Raes, F., Richter, A., Rognnerud, B., Schulz, M., Shindell, D., Stevenson, D. S., Storelvmo, T., Wang, W. C., van Weele, M., Wild, M., and Wuebbles, D.: Atmospheric composition change: Climate-Chemistry interactions, *Atmos. Environ.*, 43, 5138–5192, doi:10.1016/j.atmosenv.2009.08.003, 2009.
- Jimenez, J. L., Canagaratna, M. R., Donahue, N. M., Prevot, A. S. H., Zhang, Q., Kroll, J. H., DeCarlo, P. F., Allan, J. D., Coe, H., Ng, N. L., Aiken, A. C., Docherty, K. S., Ulbrich, I. M., Grieshop, A. P., Robinson, A. L., Duplissy, J., Smith, J. D., Wilson, K. R., Lanz, V. A., Hueglin, C., Sun, Y. L., Tian, J., Laaksonen, A., Raatikainen, T., Rautiainen, J., Vaattovaara, P., Ehn, M., Kulmala, M., Tomlinson, J. M., Collins, D. R., Cubison, M. J., Dunlea, E. J., Huffman, J. A., Onasch, T. B., Alfarra, M. R., Williams, P. I., Bower, K., Kondo, Y., Schneider, J., Drewnick, F., Borrmann, S., Weimer, S., Demerjian, K., Salcedo, D., Cottrell, L., Griffin, R., Takami, A., Miyoshi, T., Hatakeyama, S., Shimojo, A., Sun, J. Y., Zhang, Y. M., Dzepina, K., Kimmel, J. R., Sueper, D., Jayne, J. T., Herndon, S. C., Trimborn, A. M., Williams, L. R., Wood, E. C., Middlebrook, A. M., Kolb, C. E., Baltensperger, U., and Worsnop, D. R.: Evolution of Organic Aerosols in the Atmosphere, *Science*, 326, 1525–1529, doi:10.1126/science.1180353, 2009.
- Kanakidou, M., Seinfeld, J. H., Pandis, S. N., Barnes, I., Dentener, F. J., Facchini, M. C., Van Dingenen, R., Ervens, B., Nenes, A., Nielsen, C. J., Swietlicki, E., Putaud, J. P., Balkanski, Y., Fuzzi, S., Horth, J., Moortgat, G. K., Winterhalter, R., Myhre, C. E. L., Tsigaridis, K., Vignati, E., Stephanou, E. G., and Wilson, J.: Organic aerosol and global climate modelling: a review, *Atmos. Chem. Phys.*, 5, 1053–1123, doi:10.5194/acp-5-1053-2005, 2005.
- Katajamaa, M. and Oresic, M.: Data processing for mass spectrometry-based metabolomics, *J. Chromatogr. A*, 1158, 318–328, doi:10.1016/j.chroma.2007.04.021, 2007.
- Kroll, J. H., Donahue, N. M., Jimenez, J. L., Kessler, S. H., Canagaratna, M. R., Wilson, K. R., Altieri, K. E., Mazzoleni, L. R., Wozniak, A. S., Bluhm, H., Mysak, E. R., Smith, J. D., Kolb, C. E., and Worsnop, D. R.: Carbon oxidation state as a metric for describing the chemistry of atmospheric organic aerosol, *Nature Chemistry*, 3, 133–139, doi:10.1038/nchem.948, 2011.
- Kulkarni, P., Baron, P. A., and Willeke, K.: *Aerosol Measurement: Principles, Techniques, and Applications*, John Wiley & Sons, Hoboken, NJ, USA, 2011.
- Lanz, V. A., Prévôt, A. S. H., Alfarra, M. R., Weimer, S., Mohr, C., DeCarlo, P. F., Gianini, M. F. D., Hueglin, C., Schneider, J., Favez, O., D'Anna, B., George, C., and Baltensperger, U.: Characterization of aerosol chemical composition with aerosol mass spectrometry in Central Europe: an overview, *Atmos. Chem. Phys.*, 10, 10453–10471, doi:10.5194/acp-10-10453-2010, 2010.
- Liland, K. H., Almoy, T., and Mevik, B. H.: Optimal Choice of Baseline Correction for Multivariate Calibration of Spectra, *Appl. Spectrosc.*, 64, 1007–1016, 2010.
- Liland, K. H., Rukke, E. O., Olsen, E. F., and Isaksson, T.: Customized baseline correction, *Chemometr. Intell. Lab.*, 109, 51–56, doi:10.1016/j.chemolab.2011.07.005, 2011.
- Liu, C.-L., Smith, J. D., Che, D. L., Ahmed, M., Leone, S. R., and Wilson, K. R.: The direct observation of secondary radical chain chemistry in the heterogeneous reaction of chlorine atoms with submicron squalane droplets, *Phys. Chem. Chem. Phys.*, 13, 8993–9007, doi:10.1039/c1cp20236g, 2011.
- Liu, S., Takahama, S., Russell, L. M., Gilardoni, S., and Baumgardner, D.: Oxygenated organic functional groups and their sources in single and submicron organic particles in MILAGRO 2006 campaign, *Atmos. Chem. Phys.*, 9, 6849–6863, doi:10.5194/acp-9-6849-2009, 2009.
- Maria, S. F., Russell, L. M., Turpin, B. J., and Porcja, R. J.: FTIR measurements of functional groups and organic mass in aerosol samples over the Caribbean, *Atmos. Environ.*, 36, 5185–5196, doi:10.1016/s1352-2310(02)00654-4, 2002.
- Maria, S. F., Russell, L. M., Turpin, B. J., Porcja, R. J., Campos, T. L., Weber, R. J., and Huebert, B. J.: Source signatures of carbon monoxide and organic functional groups in Asian Pacific Regional Aerosol Characterization Experiment (ACE-Asia) submicron aerosol types, *J. Geophys. Res.-Atmos.*, 108, doi:10.1029/2003JD003703, 2003.
- Mazet, V., Carteret, C., Brie, D., Idier, J., and Humbert, B.: Background removal from spectra by designing and minimising a non-quadratic cost function, *Chemometr. Intell. Lab.*, 76, 121–133, doi:10.1016/j.chemolab.2004.10.003, 2005.

- McClenny, W. A., Childers, J. W., Rohl, R., and Palmer, R. A.: FTIR transmission spectrometry for the nondestructive determination of ammonium and sulfate in ambient aerosols collected on Teflon filters, *Atmos. Environ.*, 19, 1891–1898, doi:10.1016/0004-6981(85)90014-9, 1985.
- Monks, P. S., Granier, C., Fuzzi, S., Stohl, A., Williams, M. L., Aki-moto, H., Amann, M., Baklanov, A., Baltensperger, U., Bey, I., Blake, N., Blake, R. S., Carslaw, K., Cooper, O. R., Dentener, F., Fowler, D., Fragkou, E., Frost, G. J., Generoso, S., Ginoux, P., Grewe, V., Guenther, A., Hansson, H. C., Henne, S., Hjorth, J., Hofzumahaus, A., Huntrieser, H., Isaksen, I. S. A., Jenkin, M. E., Kaiser, J., Kanakidou, M., Klimont, Z., Kulmala, M., Laj, P., Lawrence, M. G., Lee, J. D., Lioussis, C., Maione, M., McFiggans, G., Metzger, A., Mieville, A., Moussiopoulos, N., Orlando, J. J., O'Dowd, C. D., Palmer, P. I., Parrish, D. D., Petzold, A., Platt, U., Poschl, U., Prevot, A. S. H., Reeves, C. E., Reimann, S., Rudich, Y., Sellegri, K., Steinbrecher, R., Simpson, D., ten Brink, H., Theloke, J., van der Werf, G. R., Vautard, R., Vestreng, V., Vlachokostas, C., and von Glasow, R.: Atmospheric composition change - global and regional air quality, *Atmos. Environ.*, 43, 5268–5350, doi:10.1016/j.atmosenv.2009.08.021, 2009.
- Pavia, D., Lampman, G., Kriz, G., and Vyvyan, J.: *Introduction to spectroscopy*, Cengage Learning, Belmont, CA, USA, 2008.
- Peng, J. T., Peng, S. L., Jiang, A., Wei, J. P., Li, C. W., and Tan, J.: Asymmetric least squares for multiple spectra baseline correction, *Anal. Chim. Acta*, 683, 63–68, doi:10.1016/j.aca.2010.08.033, 2010.
- Persson, P. B., Stauss, H., Chung, O., Wittmann, U., and Unger, T.: Spectrum analysis of sympathetic-nerve activity and blood-pressure in conscious rats, *Am. J. Physiol.*, 263, H1348–H1355, 1992.
- Pouillet, J. B., Sima, D. M., Simonetti, A. W., De Neuter, B., Vanhamme, L., Lemmerling, P., and Van Huffel, S.: An automated quantitation of short echo time MRS spectra in an open source software environment: AQSES, *NMR Biomed.*, 20, 493–504, doi:10.1002/nbm.1112, 2007.
- R Core Team: R: A Language and Environment for Statistical Computing, R Foundation for Statistical Computing, Vienna, Austria, available at: <http://www.R-project.org/> (last access: 7 December 2015), 2014.
- Rinnan, A., van den Berg, F., and Engelsen, S. B.: Review of the most common pre-processing techniques for near-infrared spectra, *Trac-Trend. Anal. Chem.*, 28, 1201–1222, doi:10.1016/j.trac.2009.07.007, 2009.
- Rouh, A., Delsuc, M. A., Bertrand, G., and Lallemand, J. Y.: The use of classification in base-line correction of FT NMR spectra, *J. Magn. Reson. Ser. A*, 102, 357–359, doi:10.1006/jmra.1993.1117, 1993.
- Rowlands, C. and Elliott, S.: Automated algorithm for baseline subtraction in spectra, *J. Raman Spectrosc.*, 42, 363–369, doi:10.1002/jrs.2691, 2011a.
- Rowlands, C. J. and Elliott, S. R.: Denoising of spectra with no user input: a spline-smoothing algorithm, *J. Raman Spectrosc.*, 42, 370–376, doi:10.1002/jrs.2692, 2011b.
- Ruckstuhl, A. F., Jacobson, M. P., Field, R. W., and Dodd, J. A.: Baseline subtraction using robust local regression estimation, *J. Quant. Spectrosc. Ra.*, 68, 179–193, doi:10.1016/s0022-4073(00)00021-2, 2001.
- Russell, L. M.: Aerosol organic-mass-to-organic-carbon ratio measurements, *Environ. Sci. Technol.*, 37, 2982–2987, doi:10.1021/es026123w, 2003.
- Russell, L. M., Bahadur, R., Hawkins, L. N., Allan, J., Baumgardner, D., Quinn, P. K., and Bates, T. S.: Organic aerosol characterization by complementary measurements of chemical bonds and molecular fragments, *Atmos. Environ.*, 43, 6100–6105, doi:10.1016/j.atmosenv.2009.09.036, 2009.
- Russell, L. M., Bahadur, R., and Ziemann, P. J.: Identifying organic aerosol sources by comparing functional group composition in chamber and atmospheric particles, *P. Natl. Acad. Sci USA*, 108, 3516–3521, doi:10.1073/pnas.1006461108, 2011.
- Ruthenburg, T. C., Perlin, P. C., Liu, V., McDade, C. E., and Dillner, A. M.: Determination of organic matter and organic matter to organic carbon ratios by infrared spectroscopy with application to selected sites in the IMPROVE network, *Atmos. Environ.*, 86, 47–57, doi:10.1016/j.atmosenv.2013.12.034, 2014.
- Schirm, B. and Watzig, H.: Peak recognition imitating human judgement, *Chromatographia*, 48, 331–346, doi:10.1007/bf02467701, 1998.
- Schulze, G., Jirasek, A., Yu, M. M. L., Lim, A., Turner, R. F. B., and Blades, M. W.: Investigation of selected baseline removal techniques as candidates for automated implementation, *Appl. Spectrosc.*, 59, 545–574, doi:10.1366/0003702053945985, 2005.
- Seinfeld, J. H. and Pandis, S. N.: *Atmospheric Chemistry and Physics: From Air Pollution to Climate Change*, 2nd edn., John Wiley & Sons, New York, USA, , 2006.
- Shao, L. and Griffiths, P. R.: Automatic Baseline Correction by Wavelet Transform for Quantitative Open-Path Fourier Transform Infrared Spectroscopy, *Environ. Sci. Technol.*, 41, 7054–7059, doi:10.1021/es062188d, 2007.
- Shao, X. G., Leung, A. K. M., and Chau, F. T.: Wavelet: A new trend in chemistry, *Accounts Chem. Res.*, 36, 276–283, doi:10.1021/ar990163w, 2003.
- Surratt, J. D., Kroll, J. H., Kleindienst, T. E., Edney, E. O., Claeys, M., Sorooshian, A., Ng, N. L., Offenberg, J. H., Lewandowski, M., Jaoui, M., Flagan, R. C., and Seinfeld, J. H.: Evidence for organosulfates in secondary organic aerosol, *Environ. Sci. Technol.*, 41, 517–527, doi:10.1021/es062081q, 2007.
- Takahama, S., Schwartz, R. E., Russell, L. M., Macdonald, A. M., Sharma, S., and Leaitch, W. R.: Organic functional groups in aerosol particles from burning and non-burning forest emissions at a high-elevation mountain site, *Atmos. Chem. Phys.*, 11, 6367–6386, doi:10.5194/acp-11-6367-2011, 2011.
- Takahama, S., Johnson, A., and Russell, L. M.: Quantification of Carboxylic and Carbonyl Functional Groups in Organic Aerosol Infrared Absorbance Spectra, *Aerosol Sci. Tech.*, 47, 310–325, doi:10.1080/02786826.2012.752065, 2013.
- Trevisan, J., Angelov, P. P., Carmichael, P. L., Scott, A. D., and Martin, F. L.: Extracting biological information with computational analysis of Fourier-transform infrared (FTIR) biospectroscopy datasets: current practices to future perspectives, *Analyst*, 137, 3202–3215, doi:10.1039/c2an16300d, 2012.
- Turpin, B. J., Huntzicker, J. J., and Hering, S. V.: Investigation of organic aerosol sampling artifacts in the los angeles basin, *Atmos. Environ.*, 28, 3061–3071, doi:10.1016/1352-2310(94)00133-6, 1994.
- Ward, Jr., J.: Hierarchical grouping to optimize an objective function, *J. Am. Stat. Assoc.*, 58, 236–244, 1963.

- Weakley, A. T., Griffiths, P. R., and Aston, D. E.: Automatic Baseline Subtraction of Vibrational Spectra Using Minima Identification and Discrimination via Adaptive, Least-Squares Thresholding, *Appl. Spectrosc.*, 66, 519–529, doi:10.1366/110-06526, 2012.
- Xi, Y. and Rocke, D. M.: Baseline correction for NMR spectroscopic metabolomics data analysis, *BMC Bioinformatics*, 9, 329, doi:10.1186/1471-2105-9-324, 2008.
- Zhang, Q., Jimenez, J. L., Canagaratna, M. R., Allan, J. D., Coe, H., Ulbrich, I., Alfarra, M. R., Takami, A., Middlebrook, A. M., Sun, Y. L., Dzepina, K., Dunlea, E., Docherty, K., DeCarlo, P. F., Salcedo, D., Onasch, T., Jayne, J. T., Miyoshi, T., Shimojo, A., Hatakeyama, S., Takegawa, N., Kondo, Y., Schneider, J., Drewnick, F., Borrmann, S., Weimer, S., Demerjian, K., Williams, P., Bower, K., Bahreini, R., Cottrell, L., Griffin, R. J., Rautiainen, J., Sun, J. Y., Zhang, Y. M., and Worsnop, D. R.: Ubiquity and dominance of oxygenated species in organic aerosols in anthropogenically-influenced Northern Hemisphere midlatitudes, *Geophys. Res. Lett.*, 34, L13801, doi:10.1029/2007GL029979, 2007.
- Zhao, A. X., Tang, X. J., Li, W. D., Zhang, Z. H., and Liu, J. H.: The Piecewise Two Points Autolinear Correlated Correction Method for Fourier Transform Infrared Baseline Wander, *Spectrosc. Lett.*, 48, 274–279, doi:10.1080/00387010.2013.874530, 2015.