

# Unveiling the optimal regression model for source apportionment of the oxidative potential of PM<sub>10</sub>

Vy Dinh Ngoc Thuy<sup>1</sup>, Jean-Luc Jaffrezo<sup>1</sup>, Ian Hough<sup>1</sup>, Pamela A. Dominutti<sup>1</sup>, Guillaume Salque Moreton<sup>2</sup>, Grégory Gille<sup>3</sup>, Florie Francony<sup>4</sup>, Arabelle Patron-Anquez<sup>5</sup>, Olivier Favez<sup>6,7</sup>, Gaëlle Uzu<sup>1</sup>

<sup>1</sup> Université Grenoble Alpes, CNRS, IRD, INP-G, INRAE, IGE (UMR 5001), F-38000 Grenoble, France

<sup>2</sup> Atmo AuRA, 69500 Bron, France

<sup>3</sup>Atmo Sud, 13006 Marseille, France

<sup>4</sup>Atmo Nouvelle Aquitaine, 33692 Merignac, France

<sup>5</sup>Atmo Hauts de France, 59044 Lille, France

<sup>6</sup>INERIS, Parc Technologique Alata, BP 2, 60550 Verneuil-en-Halatte, France

<sup>7</sup> Laboratoire central de surveillance de la qualité de l'air (LCSQA), 60550 Verneuil-en-Halatte, France

*Correspondance to: gaelle.uzu@ird.fr*

## Abstract

The capacity of particulate matter (PM) to generate reactive oxygen species (ROS) in vivo leading to oxidative stress, is thought to be a main pathway for the health effect of PM inhalation. Exogenous ROS from PM can be assessed by acellular oxidative potential (OP) measurements as a proxy of the induction of oxidative stress in the lungs. Here, we investigate the importance of OP apportionment methods on OP repartition by PM<sub>10</sub> sources in different types of environments. PM<sub>10</sub> sources derived from receptor models (e.g. EPA PMF) are coupled with regression models expressing the associations between PM<sub>10</sub> sources and PM<sub>10</sub> OP measured by ascorbic acid (OP<sub>AA</sub>) and dithiothreitol assay (OP<sub>DTT</sub>). These relationships are compared for eight regression techniques: Ordinary Least Squares, Weighted Least Squares, Positive Least Squares, Ridge, Lasso, Generalized Linear Model, Random Forest, and Multilayer Perceptron. The models are evaluated on one year of PM<sub>10</sub> samples and chemical analyses at each of six sites of different typologies in France to assess the possible impact of PM source variability on PM<sub>10</sub> OP apportionment. PM<sub>10</sub> source-specific OP<sub>DTT</sub> and OP<sub>AA</sub> and out-of-sample apportionment accuracy vary substantially by model, highlighting the importance of model selection depending on the datasets. Recommendations for the selection of the most accurate model are provided, encompassing considerations such as multicollinearity and homoscedasticity.

Key words: Oxidative potential, source apportionment, OP apportionment.

## 1. Introduction

Ambient particulate matter (PM) is one of the key contributors to atmospheric pollution and is responsible for approximately 7 million premature deaths worldwide yearly (WHO, 2021). Many epidemiological studies have linked PM exposure to adverse health effects including (i) acute effects studies using time series and related studies to evaluate the immediate impact of PM exposure (Bell et al., 2004; Dominici, 2004; Pope and Dockery, 2006; Peng et al., 2009) and (ii) cohort studies aiming to evaluate the long-term effects of chronic PM exposure (Pelucchi et al., 2009; Crouse et al., 2012, 2015; Beelen et al., 2014; Ayres et al., 2008; Yu et al., 2021). These studies mainly focused on the association with PM mass concentrations. However, various research shows that the impacts of PM also depend on other factors such as chemical composition, size distribution, particle morphology, and biological mechanisms (Brook et al., 2010). PM's capacity to generate reactive oxygen species (ROS) in vivo has recently been introduced as a pivotal indicator of PM biological mechanism, with direct implications for oxidative stress

41 and cellular damage (Li et al., 2008; Lodovici and Bigagli, 2011; Mudway et al., 2020; Nelin et al., 2012; Rao et  
42 al., 2018; Ayres et al., 2008; Akhtar et al., 2010; Leni et al., 2020). The quantification of the PM capacity to oxidize  
43 a biological media is called oxidative potential (OP) (Bates et al., 2019; Daellenbach et al., 2020; Dominutti et al.,  
44 2023). Various acellular assays of OP have been introduced, differentiating ROS generation mechanisms of PM  
45 (Dominutti et al., 2023; Calas et al., 2018). Dithiothreitol (DTT) and ascorbic acid (AA) assays are two of the  
46 commonly used ones in the literature (Liu and Ng, 2023).

47 The relationship between PM chemical components and OP activities may identify which components are the most  
48 prone to generate ROS (Calas et al., 2019; Godri et al., 2011; Yang et al., 2014; Janssen et al., 2014; Szigeti et al.,  
49 2016; Crobeddu et al., 2017; Szigeti et al., 2015; Calas et al., 2018). However, this research pathway struggles  
50 with the co-variation between measured and unmeasured PM components (Calas et al., 2018; Weber et al., 2018).  
51 An alternative approach is to examine the association between OP and sources of PM obtained using receptor  
52 models such as chemical mass balance, positive matrix factorization (PMF), or principal components analysis.  
53 PMF is the most popular method for its ability to quantify PM source contributions without extensive prior  
54 information on specific sources at the site studied (Belis et al., 2013; Viana et al., 2008; Paatero and Tappert, 1994;  
55 Brown et al., 2015; Paatero and Hopke, 2009).

56 Regression analysis is the most common and effective way to estimate the redox activity of receptor model-derived  
57 PM sources (Bates et al., 2015; Deng et al., 2022; Li et al., 2023; Liu et al., 2018; Shangguan et al., 2022; Verma  
58 et al., 2014; Wang et al., 2020a; Yu et al., 2019). Generally, this is achieved by regression analyses to characterize  
59 the relationship between OP activities ( $\text{nmol min}^{-1} \text{m}^{-3}$ ) and PM sources contribution ( $\mu\text{g m}^{-3}$ ). This approach  
60 provides the OP activities attributed to each microgram of each source ( $\text{nmol min}^{-1} \mu\text{g}^{-1}$ ), denoted as intrinsic OP,  
61 which can be used to calculate the contribution of each source for each observation day. Numerous regression  
62 models can be used for such OP source apportionment (SA), with multiple linear regression fitted by ordinary least  
63 squares (OLS) being the most common regression technique (Bates et al., 2015; Deng et al., 2022; Li et al., 2023;  
64 Liu et al., 2018; Shangguan et al., 2022; Verma et al., 2014; Wang et al., 2020b; Yu et al., 2019). Further, some  
65 studies exclude sources with negative intrinsic OP, assuming that negative OP activities are geochemically  
66 nonsensical (Bates et al., 2018; Weber et al., 2018). Additionally, weighted least square can be used to introduce  
67 a weighting term, usually using the OP analysis uncertainties to take into account the measurement uncertainties  
68 of the OP assays (Borlaza et al., 2021; Daellenbach et al., 2020; Dominutti et al., 2023; Fadel et al., 2023; in 't  
69 Veld et al., 2023b; Weber et al., 2021). Finally, non-linear models, such as multilayer perceptron, have been used  
70 to try to capture possible non-linearities between OP activities and PM sources (Borlaza et al., 2021; Elangasinghe  
71 et al., 2014; D. Wang et al., 2023). However, no study to date has compared the performance and applicability of  
72 these various regression models. Each model implies different assumptions which should be carefully considered  
73 when selecting a given model.

74 This study aims to evaluate the variability in  $\text{PM}_{10}$  OP SA techniques by comparing eight regression techniques:  
75 multiple linear regression fitted by OLS, weighted least squares (WLS), positive least squares (PLS), Ridge  
76 regression (Ridge), Least Absolute Shrinkage and Selection Operator (Lasso), generalized linear model (GLM),  
77 random forest (RF), and multilayer perceptron (MLP). These techniques are applied to apportion  $\text{PM}_{10}$   $\text{OP}_{\text{AA}}$  and  
78  $\text{PM}_{10}$   $\text{OP}_{\text{DTT}}$  to  $\text{PM}_{10}$  sources at six sites in France. The  $\text{PM}_{10}$  SA outputs have been published previously in Weber  
79 et al. (2021), using a harmonized PMF methodology based on one year of sampling with similar chemical analyses  
80 for a large set of chemical tracers. The results of the  $\text{PM}_{10}$  OP SA models are compared with regard to the estimated  
81 intrinsic  $\text{PM}_{10}$  OP of each source, the out-of-sample accuracy of the apportionment, and the assumptions inherent  
82 in each model. The most appropriate model at each site is compared with OLS to quantify the difference between  
83 choosing a model based on data characteristics vs. using the most common approach. Finally, this study provides  
84 guidelines for selecting the most suitable model in the strategy for OP contribution regarding sources of  $\text{PM}_{10}$ .  
85 This holds particular significance in the context of the implementation of OP monitoring as a novel air quality

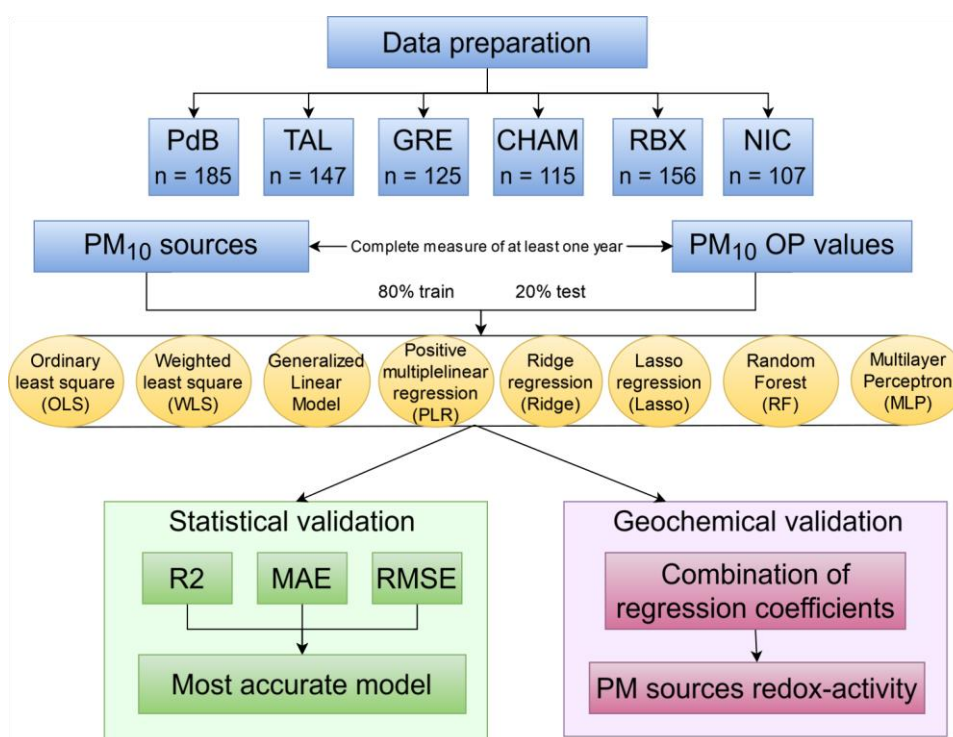
86 metric as foreseen in research programs (such RI-Urbans) and in the process of the revision of the European  
87 Directive 2008/50/CE.

## 88 2. Methodology

### 89 2.1. General organisation of this work

90 Figure 1 illustrates the general workflow of this work. Sections 2.2, 2.3, and 2.4 describe the methods used to  
91 analyse the temporal evolution of PM<sub>10</sub> sources and PM<sub>10</sub> OP, identify collinearity among PM<sub>10</sub> sources, and  
92 examine homoscedasticity in the relationship between PM<sub>10</sub> OP and PM<sub>10</sub> sources. Section 2.5 describes the eight  
93 regression techniques (OLS, WLS, PLS, Ridge, Lasso, GLM, RF, and MLP), used for PM<sub>10</sub> OP SA. Each  
94 technique is applied to each site separately using PM<sub>10</sub> OP<sub>v</sub> (nmol min<sup>-1</sup> m<sup>-3</sup>) as the dependent variable and PM<sub>10</sub>  
95 sources (µg m<sup>-3</sup>) as independent variables. The coefficient of the regression called the intrinsic PM<sub>10</sub> OP of the  
96 source (nmol min<sup>-1</sup> µg<sup>-1</sup>), represents the capacity of each µg of PM<sub>10</sub> from the given source to generate oxidative  
97 stress; the higher the intrinsic PM<sub>10</sub> OP of a source, the more redox-active. Each model is trained on a randomly  
98 selected (without replacement) 80% subsample of the dataset and validated on the remaining 20%. This process is  
99 repeated 500 times to estimate uncertainty, a method particularly needed for sources with strong seasonality. For  
100 WLS, PLS, Ridge, and Lasso models, PM<sub>10</sub> OP analytical errors were used as a weighting, implying that the PM<sub>10</sub>  
101 OP with the high analysis uncertainties has less influence on the model. These 8 regression techniques were applied  
102 to find the relationship between PM<sub>10</sub> OP and PM<sub>10</sub> sources, however, PLS, Ridge, and Lasso were performed 2  
103 times, with and without weighting, consequently, there are 11 results of regression techniques that will be  
104 presented. Section 2.6 describes the statistical validation of the models using root mean square error (RMSE),  
105 mean absolute error (MAE), R-square (R<sup>2</sup>). The geochemical validation is based on the regression coefficient (the  
106 intrinsic PM<sub>10</sub> OP) of each source. These are calculated separately for the training and testing data and averaged  
107 across the 500 sampling iterations.

108



109

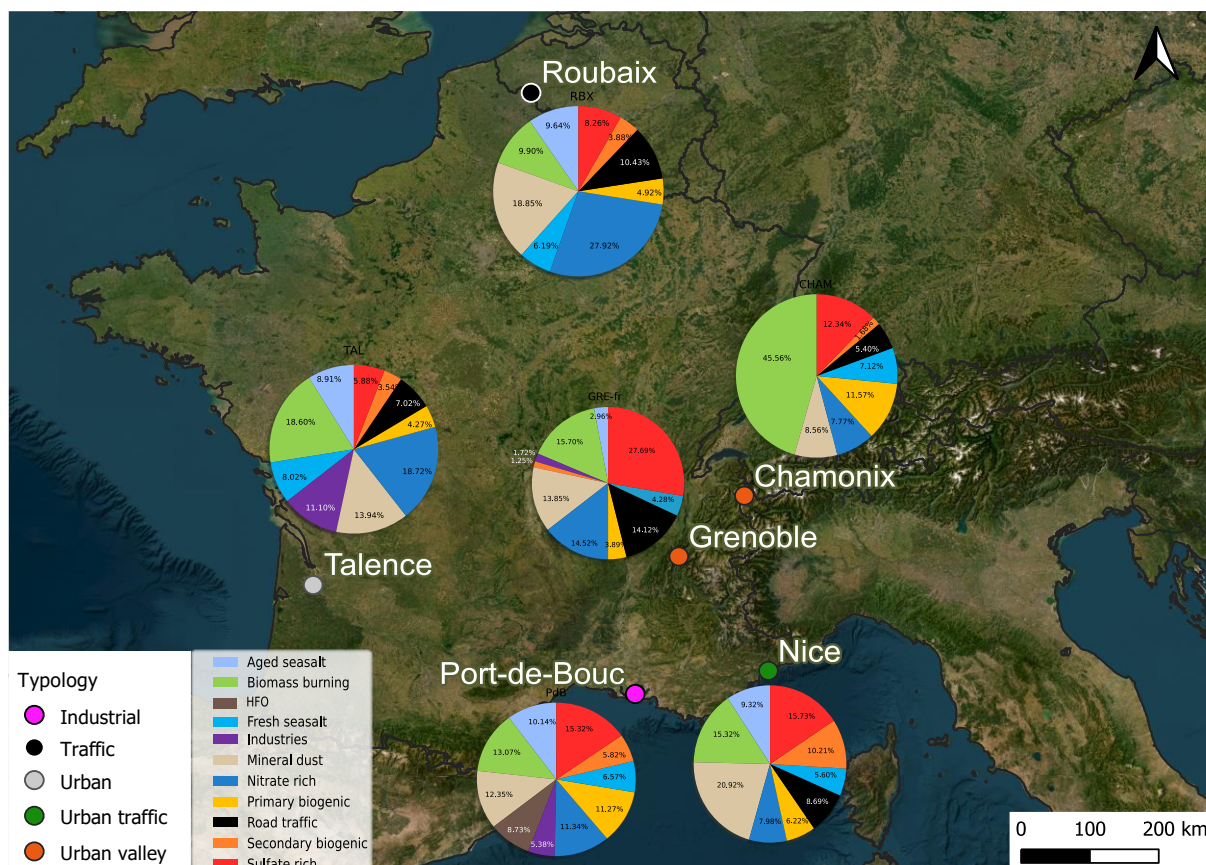
110 **Figure 1. Workflow of the comparison of PM<sub>10</sub> OP sources apportionment methodology**

## 111 2.2. Study sites and PM<sub>10</sub> sources

112 Six French sites are selected in this work for their different typologies: Roubaix and Nice (traffic sites within urban  
113 areas), Port-de-Bouc (industrial hotspot), Talence (urban background site), Grenoble and Chamonix (urban  
114 background sites in Alpine Valley). At each site, sampling was conducted over at least one year to capture the  
115 complete annual evolution of PM<sub>10</sub> and its components. These sites and sampling series were previously used and  
116 described by Weber et al. (2019).

117 In brief, daily filter samples were collected on pre-heated Pallflex quartz fibre filters every third day through high-  
118 volume sampling (DA80, Digitel). These filters were analyzed to determine PM's chemical species and OP  
119 activities. Further details regarding the chemical species and PM<sub>10</sub> OP analyses methodology can be found in  
120 Weber et al. (Weber et al., 2019, 2021). Briefly, the elemental carbon (EC) and organic carbon (OC) were analyzed  
121 using the EUSAAR2 thermo-optical protocol with a Sunset Lab analyser. Major ionic components (Cl<sup>-</sup>, NO<sub>3</sub><sup>-</sup>,  
122 SO<sub>4</sub><sup>2-</sup>, NH<sub>4</sub><sup>+</sup>, Na<sup>+</sup>, K<sup>+</sup>, Mg<sup>2+</sup>, Ca<sup>2+</sup>) and methanesulfonic acid (MSA) were measured by ion chromatography (IC).  
123 Anhydro-sugars and saccharides (including levoglucosan, mannosan, arabitol, sorbitol, and mannitol) were  
124 analysed by high-performance liquid chromatography with pulsed amperometry detection (HPLC-PAD). Major  
125 and trace elements (Al, Ca, Fe, K, As, Ba, Cd, Co, Cu, La, Mn, Mo, Ni, Pb, Rb, Sb, Sr, V, and Zn) were determined  
126 by inductively coupled plasma atomic emission spectroscopy or mass spectrometry (ICP-AES or ICP-MS).  
127 Furthermore, colocated PM<sub>10</sub> measurements were conducted automatically at each site using the Tapered Element  
128 Oscillating Microbalance equipped with a Filter Dynamics Measurement System (TEOM-FDMS).

129 We used the PM<sub>10</sub> sources identified by Weber et al. (2019), who performed a separate PMF for each site using a  
130 harmonized approach for all sites (same chemical species and measurement methods, same procedure to estimate  
131 uncertainties, same constraints on the preliminary solutions). Table 1 provides a data description, including the  
132 sampling duration, the number of samples collected, and the identified PM<sub>10</sub> sources at each site, while Figure 2  
133 presents the localisation of the sites in France, together with the respective proportion of each PM<sub>10</sub> source at each  
134 site.



135  
 136 **Figure 2. The location of the selected sites for this study. The small colored dots represent the typology of**  
 137 **sites. The pie charts are the PM<sub>10</sub> source apportionment for each site with the colors identifying the PM<sub>10</sub>**  
 138 **sources. Background photography from ESRI satellite.**

139 Table 1. Data description

	<b>PdB</b>	<b>TAL</b>	<b>GRE-fr</b>	<b>CHAM</b>	<b>RBX</b>	<b>NIC</b>
<b>Name</b>	Port de Bouc	Talence	Grenoble	Chamonix	Roubaix	Nice
<b>N of samples</b>	185	147	125	115	156	107
<b>Sampling dates</b>	2014-06 to 2016-06	2012-02 to 2013-04	2017-02 to 2018-03	2013-11 to 2014-10	2013-01 to 2014-05	2014-07 to 2015-05
<b>N of sources</b>	10	10	10	8	9	9

140  
 141 **2.3. OP analysis**  
 142 PM<sub>10</sub> OP assays were performed on PM<sub>10</sub> extracted from the filters using simulated lung fluid, as detailed in Calas  
 143 et al. (2017, 2018). The AA assay involved ascorbic acid, a natural antioxidant in the lungs inhibiting lipid and  
 144 protein oxidation in the lining fluid, using the method presented by Kelly & Mudway (2003) and further described  
 145 by Calas et al. (2018). Conversely, the DTT assay used dithiothreitol (DTT) as a chemical surrogate for cellular  
 146 reducing agents, specifically nicotinamide adenine dinucleotide and nicotinamide adenine dinucleotide phosphate

147 oxidase, thereby replicating in vivo interactions between PM<sub>10</sub> and biological oxidants (Cho et al., 2005; Calas et  
148 al., 2018). Both assays measured the consumption of AA or DTT during the assay, i.e., the rate of the transfer of  
149 electrons from AA or DTT to oxygen. The assays were conducted with 96-well plates of UV-transparent quality  
150 (CELLSTAR, Greiner-Bio), and absorption measurements were acquired using a TECAN spectrophotometer,  
151 Infinite M200 Pro, at the wavelengths of 265nm for the AA assay and 412nm for the DTT assay (Calas et al.,  
152 2017, 2018, 2019). Each sample extraction was subjected to four analyses; the PM<sub>10</sub> OP in this study represents  
153 the mean and the analysis uncertainty is the standard deviation of these four PM<sub>10</sub> OP analyses. After analysis, the  
154 PM<sub>10</sub> OP activities of each sample were blank-subtracted using lab and field blanks, and normalized using the air  
155 sampling volumes and the mass concentration. The resulting OP<sub>v</sub> represents the PM<sub>10</sub> OP due to PM<sub>10</sub> per cubic  
156 meter of air (nmol min<sup>-1</sup> m<sup>-3</sup>). To simplify the denotation of PM<sub>10</sub> OP, OP is used to represent PM<sub>10</sub> OP throughout  
157 this article.

## 158 **2.4. Collinearity and heteroscedasticity tests**

159 The result of a regression model strongly depends on the characteristics of the dataset because each model makes  
160 assumptions about the data. Two critical assumptions in OLS regression analysis are that (1) there is little  
161 collinearity between independent variables (the PM<sub>10</sub> sources in this study), and (2) the variance of the regression  
162 residuals is constant (called homoscedasticity). These assumptions should be tested in different ways.

### 163 **2.4.1. Collinearity**

164 Collinearity occurs when one or more of the independent variables is close to a linear combination of the other  
165 independent variables. When collinearity is present, small changes in the data can cause large changes in estimated  
166 coefficients, and the estimated standard errors of the coefficients are large. Variance Inflation Factor (VIF) is an  
167 indicator of the collinearity between the independent variables (Craney & Surles, 2002; O'Brien, 2007; Rosenblad,  
168 2011). VIF of a specific source is calculated as:

$$169 \quad VIF_i = \frac{1}{1 - R_i^2}, i = 1, \dots, p - 1 \text{ (Eq1)}$$

170 In this equation,  $p$  is the number of PM<sub>10</sub> sources,  $R^2$  is the coefficient of determination of a multiple linear  
171 regression model between the  $i^{\text{th}}$  source and the other sources. VIF values of a PM<sub>10</sub> source present a range between  
172 1, and  $\infty$ . The higher the VIF values, the greater the collinearity between this PM<sub>10</sub> source and the other ones. A  
173 VIF value between 5 and 10 is commonly interpreted as moderate collinearity, while values greater than 10 indicate  
174 high collinearity (Craney and Surles, 2002).

### 175 **2.4.2. Heteroscedasticity**

176 Heteroscedasticity occurs when the variance of regression residuals is not constant but varies for different values  
177 of the dependent variable. In this case, the estimated standard errors of the regression coefficients are not reliable.  
178 The Goldfeld–Quandt test was developed by Goldfeld & Quandt (1965) to evaluate residual variance in a  
179 regression model. To implement the Goldfeld–Quandt test, an OLS regression was performed between OP and  
180 PM<sub>10</sub> sources to identify the residual of OP prediction. Next, the PM<sub>10</sub> sources and residual corresponding are  
181 divided into three segments: the upper segment is the group with higher PM<sub>10</sub> sources concentration, the lower  
182 segment is the group with lower PM<sub>10</sub> sources concentration, and the middle segment, constituting 10% of the  
183 moderate PM<sub>10</sub> concentration, is excluded. A subsequent regression analysis is then conducted on the two  
184 remaining subgroups to determine the ratio of residual sums of squares. Finally, an F-test is conducted on this ratio  
185 to assess whether the variances are the same, with a p-value below 0.05 interpreted as evidence of  
186 heteroscedasticity.

187 The Variance Inflation Factor (VIF) and the Goldfeld–Quandt test were performed in Python 3.9, using the  
188 statsmodels 0.14.0 package (Seabold and Perktold, 2010).

## 189 2.5. Regression models

190 The fundamental principle of regression models in this study is to use the PM<sub>10</sub> sources to predict OP activities by  
191 identifying the parameters (coefficients and residuals) that minimize an error term (Hastie, 2009). A simple  
192 regression model can be represented by Eq. 2, which defines the estimated function of the regression model, and  
193 Eq. 3, which estimates the residuals.

$$194 \hat{y} = f(X) + e \text{ (Eq2)}$$

$$195 e = y - \hat{y} \text{ (Eq3)}$$

196 Here,  $\hat{y}$  is the estimated OP (nmol min<sup>-1</sup> m<sup>-3</sup>),  $X$  are the PM<sub>10</sub> source contributions (μg m<sup>-3</sup>),  $y$  is the observed OP  
197 (nmol min<sup>-1</sup> m<sup>-3</sup>), and  $e$  denotes the residuals (nmol min<sup>-1</sup> m<sup>-3</sup>). Each model has certain assumptions and a  
198 minimization term, as presented below.

### 199 Ordinary least squares (OLS):

200 OLS is a linear regression technique that minimizes the residual sum of squares. This model is based on several  
201 assumptions: (1) **Linearity**: The relationship between OP and PM<sub>10</sub> sources is linear. (2) **Independence**: The  
202 PM<sub>10</sub> sources must be independent, with no collinearity. (3) **Homoscedasticity**: The variance of residuals is  
203 constant across all values of PM<sub>10</sub> sources. (4) **Normality**: The residuals are normally distributed. In the OLS  
204 model, the estimated equation and objective to minimize are defined as follows:

$$205 \hat{y} = \beta_0 + \sum_{i=1}^p \beta_i * x_i \text{ (Eq4)}$$

$$206 \text{Minimize: } \sum_{i=1}^m (y_i - \hat{y}_i)^2 \text{ (Eq5)}$$

207 Here, the  $\beta_0$  denotes the intercept (nmol min<sup>-1</sup> m<sup>-3</sup>),  $\beta_i$  represents the regression coefficient (intrinsic OP, nmol  
208 min<sup>-1</sup> μg<sup>-1</sup>) of source  $i$ ,  $x_i$  is the concentration of source  $i$  (μg m<sup>-3</sup>),  $p$  is the number of PM<sub>10</sub> sources, and  $m$  is the  
209 number of observations.

### 210 Weighted least square (WLS):

211 The assumptions and the minimization term in WLS closely align with those in OLS. The only difference is that  
212 WLS accounts for heteroscedasticity by introducing a weighting term for individual OP observations, whose  
213 variance is assumed to be related to the variance of the residuals. The estimation equation in WLS is the same as  
214 that of OLS, but the objective to minimize is expressed as:

$$215 \text{Minimize: } \sum_{i=1}^m (y_i - \hat{y}_i)^2 * w_i \text{ (Eq6)}$$

$$216 w_i = \frac{1}{SD_i^2}$$

217 With  $w_i$  being the weight assigned to each observation, and  $SD_i$  is the OP analysis variance of each observation.

### 218 Positive least square (PLS):

219 The assumptions for PLS primarily include linearity, independence, and normality. PLS can be applied with  
220 weighting, if there is heteroscedasticity in the data. PLS extends OLS with the constraint that the regression  
221 coefficients must be non-negative. The estimation equation and the error term, PLS, are similar to OLS (without  
222 weighting) and WLS (applying weighting). To ensure the positivity of coefficients, a specific condition must be  
223 met:

224

$$\beta_i \geq 0, \forall i \text{ in PM sources (Eq7)}$$

225 **Ridge:**

226 Shrinkage methods such as Ridge regression try to produce a more interpretable model or reduce error in the  
 227 presence of collinearity by selecting a subset of the independent variables. Ridge regression is introduced by Hoerl  
 228 & Kennard (1970), which incorporates a penalty term that shrinks the coefficients towards zero. The Ridge  
 229 regression minimizes the residual sum of squares plus a penalty term proportional to the sum of squares of the  
 230 coefficients (L2 regularization) as shown in Eq 8 and Eq 9. Consequently, Ridge regression reduces the influence  
 231 of a PM<sub>10</sub> source that exhibits minimal impact on OP prediction without excluding it from the model.

$$232 \text{ Minimize: } \sum_{i=1}^m (y_i - \hat{y}_i)^2 + \lambda * \sum_{j=1}^p \beta_j^2 \text{ (Eq8)}$$

233 *Minimize:*  $\frac{1}{2m} \sum_{i=1}^m w_i (y_i - \hat{y}_i)^2 + \lambda * \sum_{j=1}^p \beta_j^2$  (Eq9) where  $\lambda$  is the parameter representing the amount of  
 234 shrinkage, the larger  $\lambda$ , the greater the shrinkage. The hyperparameter tuning was implemented with different  
 235 values of  $\lambda$  (5, 1, 0.5, 0.1, 0.01, 0.005, 0.001, 0.0005, 0.0001). The best  $\lambda$  for every site varied from 0.005 to 0.01  
 236 and in this study, 0.01 was selected. Ridge can be applied with weighting to account for heteroscedasticity.

237 **Least Absolute Shrinkage and Selection Operator (Lasso):**

238 Lasso (Tibshirani, 1996) is a shrinkage method that uses a penalty term proportional to the sum of the absolute  
 239 regression coefficients (L1 regularization). This penalty term shrinks the coefficients of a source with a low impact  
 240 on OP prediction to zero, effectively removing it from the model. This results in a sparse model that may be easier  
 241 to interpret and may reduce error on out-of-sample data. However, Lasso is more sensitive to outliers than ridge  
 242 regression and is less stable when data are collinear. Lasso can be applied with weighting to account for  
 243 heteroscedasticity.

$$244 \text{ Minimize: } \sum_{i=1}^m (y_i - \hat{y}_i)^2 + \lambda * \sum_{j=1}^p |\beta_j| \text{ (Eq10)}$$

$$245 \text{ Minimize: } \frac{1}{2m} \sum_{i=1}^m w_i (y_i - \hat{y}_i)^2 + \lambda * \sum_{j=1}^p |\beta_j| \text{ (Eq11)}$$

246

247 Similar to Ridge,  $\lambda$  is the parameter representing the amount of shrinkage.  $\lambda$  is selected as 0.01 in this study by  
 248 running the hyperparameter tuning using the same values as for Ridge.

249 **Generalized linear model (GLM):**

250 Generalized linear models, as introduced by McCullagh (1989), provide a framework for regression analysis that  
 251 can contain non-normal error distributions and capture non-linear relationships between OP activities and PM<sub>10</sub>  
 252 sources. GLM allows for error variance that is a function of the predicted value, hence accounting for  
 253 heteroskedasticity. Key assumptions underlying GLM include (1) independence, (2) the non-normal distribution  
 254 of OP, and (3) the relationship between the PM<sub>10</sub> sources and the transformed OP (logarithm in this study) is linear.  
 255 The mathematical expression for GLM can be represented as follows:

$$256 \log(\hat{y}) = \beta_0 + \sum_0^p \beta_i * x_i \text{ (Eq12)}$$



257 where  $\beta_0$  denotes the intercept,  $\beta_i$  represents the regression coefficient of source  $i$ , and  $x_i$  is the concentration of  
258 source  $i$ .

### 259 **Random forest (RF):**

260 RF, an ensemble learning method introduced by Breiman (2001), combines multiple decision trees to make  
261 predictions. In the reference implementation, each tree is grown on a bootstrap sample of the data and a random  
262 subset of the available features is evaluated at each node to choose the best split. The predictions of all trees are  
263 averaged to give the forest's final prediction. RF is customizable via hyperparameters such as the number of trees,  
264 the size of the bootstrap sample, and the number of features to evaluate at each node. The hyperparameters tuning  
265 used 5-fold cross-validation on the training data for hyperparameter tuning. The training dataset was separated  
266 into 5 parts: 4 parts were used for training, and the remaining was used for validation. This process was repeated  
267 5 times, and the hyperparameter value producing the lowest mean RMSE across the 5 parts was selected. The  
268 hyperparameters tuning is shown in section S1.1 Supplement.

269 RF does not assume a specific equation to express the relationship between OP activities and PM<sub>10</sub> sources, with  
270 the result that intrinsic OP could not be computed in this regression model. Nevertheless, RF can estimate the  
271 relative importance of each PM<sub>10</sub> source in OP prediction. This study estimated the permutation importance of  
272 each PM<sub>10</sub> source as the mean increase in the mean squared error of predicted OP when the values of the PM<sub>10</sub>  
273 source were permuted.

### 274 **Multilayer perception (MLP):**

275 MLP is an artificial neural network that consists of multiple layers of interconnected nodes or neurons organized  
276 in a feedforward structure (Akhtar et al., 2018; Chianese et al., 2018; Bourlard and Wellekens, 1989). These layers  
277 include an input layer (PM<sub>10</sub> sources), one or several hidden layers, and an output layer (OP<sub>AA</sub> or OP<sub>DTT</sub> activities).  
278 In MLP, the neurons in the hidden layers are linked with the previous neurons by the connection weight, where  
279 every neuron is independent and has a different weight. The output of each neuron depends on its inputs and an  
280 activation function, which, if non-linear, allows the model to capture non-linear relationships. The implementation  
281 of MLP includes three steps: (1) forward pass to training model: the input is passed to the model, multiplied with  
282 an initial weight, add bias at every layer, then calculate output of the model. (2) error calculation: after applying  
283 step 1, the output of the model and the observed data are used to calculate the error. (3) backward pass: the error  
284 is propagated back through the network, and then the weights are adjusted to minimize overall error. These 3 steps  
285 are repeated until the error is minimized.

286 The choice of hyperparameters to ensure the MLP model's robustness is processed by hyperparameter tuning using  
287 5-fold cross-validation as shown in section S1.2 of the supplement. Thanks to hyperparameter tuning, the two  
288 hidden layers and a logistic sigmoid activation function were selected in this study to capture the non-linear  
289 relationships between OP activities and PM<sub>10</sub> sources.

290 All regression models were performed using the Python package statsmodels 0.14.0 (Seabold and Perktold, 2010)  
291 and scikit-learn 1.3.1 (Pedregosa et al., 2011).

### 292 **Performance of the models**

293 The performance metrics R-square ( $R^2$ ), mean absolute error (MAE), and root mean square error (RMSE) were  
294 used to assess the goodness of fit of models as described by Kuhn & Johnson (2013).  $R^2$  quantifies the model's  
295 ability to explain the variance in the data.  $R^2$  equal to 1 indicates a perfect fit. RMSE represents the aggregation of  
296 the individual differences between predicted OP and measured OP, while MAE assesses the average magnitude of  
297 errors between them. Lower RMSE and MAE values indicate a better fit, with a perfectly fitting model yielding  
298 an RMSE or MAE of 0. Eq13, Eq14, and Eq15, respectively, define  $R^2$ , MAE, RMSE. These indicators are

309 computed for the training and testing data of each sampling iteration and averaged across the 500 sampling  
300 iterations.

$$301 \quad R^2 = 1 - \frac{\text{Sum of Squared Residuals}}{\text{Total Sum of Squares}} = 1 - \frac{\sum_{i=0}^m (y_i - \hat{y}_i)^2}{\sum_{i=0}^m (y_i - \bar{y})^2} \quad (Eq13)$$

$$302 \quad MAE = \frac{\sum_{i=0}^m |y_i - \hat{y}_i|}{m} \quad (Eq14)$$

$$303 \quad RMSE = \sqrt{\frac{\sum_{i=0}^m (y_i - \hat{y}_i)^2}{m}} \quad (Eq15)$$

304

### 305 **3. Result and discussion**

306 Assessments of collinearity and homoscedasticity are addressed in Section 3.1. Model performance, including key  
307 performance metrics and identification of the optimal model, is detailed in Section 3.2. Section 3.3 compares the  
308 intrinsic OP estimated by the different models. Section 3.4 compares intrinsic OP between the combined best-fit  
309 and reference models. Lastly, Section 3.5 proposes recommendations for selecting an appropriate model.

#### 310 **3.1. Dataset characteristics**

311 The contributions of identified sources ( $\mu\text{g m}^{-3}$ ) and the  $\text{OP}_v$  activities ( $\text{nmol min}^{-1} \text{m}^{-3}$ ) in each site are presented  
312 in Figure 3, illustrating variations in annual average OP activities and  $\text{PM}_{10}$  source contributions by sites. Most  
313 sites, including traffic and industrial ones, show higher  $\text{OP}_{\text{DTT}}$  activities than  $\text{OP}_{\text{AA}}$ . Conversely, for the alpine  
314 valley sites, CHAM presents higher  $\text{OP}_{\text{AA}}$  than  $\text{OP}_{\text{DTT}}$ , while GRE-fr experiences similar levels of  $\text{OP}_{\text{AA}}$  and  
315  $\text{OP}_{\text{DTT}}$ . Additionally, the average OP activities in every site are not proportional to the average PM concentration.  
316 For instance, CHAM and NIC had lower  $\text{PM}_{10}$  concentrations but higher OP activities than other sites, while TAL  
317 showed high  $\text{PM}_{10}$  concentrations but relatively lower OP activities.

318 The variations observed in the levels of  $\text{PM}_{10}$  and OP across six sites can be attributed to distinctions in identified  
319 sources and their respective contributions. These disparities are contingent upon the unique typologies of each site,  
320 which are discussed in Weber et al., 2021. Further, we can observe a significant seasonality in the OP activities  
321 (Table S.1). Strong seasonality of OP in Alpine valley sites has been addressed in previous studies (Borlaza et al.,  
322 2021; Dominutti et al., 2023; Weber et al., 2018, 2021), with thermal inversions during winter increasing pollutants  
323 concentrations and OP activities compared to summer. Conversely, OP activities in cold and warm periods in other  
324 sites are not significantly different.

325 The  $\text{PM}_{10}$  sources and their repartition vary among sites (Figure 3) because of the difference in typology and local  
326 activities. For instance, in the industrial site (PdB), two specific sources are identified: shipping emissions (HFO)  
327 with an annual mean contribution of  $1.39 \mu\text{g m}^{-3}$  and industrial sources at  $0.86 \mu\text{g m}^{-3}$ . The urban background site  
328 TAL also appears to be influenced by industrial sources ( $2.34 \mu\text{g m}^{-3}$ ), which might, however, be partly due to  
329 biases induced by the application of the harmonized receptor model protocol (Weber et al., 2019). Note that the  
330 application of a site-specific PMF procedure for this site leads to a much lower contribution of this source category  
331 but relatively similar contributions of other sources (Favez, 2017). GRE-fr, an urban background site in an alpine  
332 valley, presents significant long-range transport sources, with secondary sulfate contributing  $3.90 \mu\text{g m}^{-3}$  followed  
333 by biomass burning at  $2.21 \mu\text{g m}^{-3}$ . As expected, biomass burning is an abundant source in CHAM, accounting for  
334  $7.28 \mu\text{g m}^{-3}$  of the PM contribution, while the traffic sites RBX and NIC displayed high contributions of traffic  
335 sources (at  $2.43 \mu\text{g m}^{-3}$  and  $1.45 \mu\text{g m}^{-3}$  respectively).

336 The presence of multicollinearity and homoscedasticity were tested to assess the data characteristic of every site.  
337 The only site with evidence of collinearity was NIC, where the VIF of the traffic source was equal to 5.0. For all

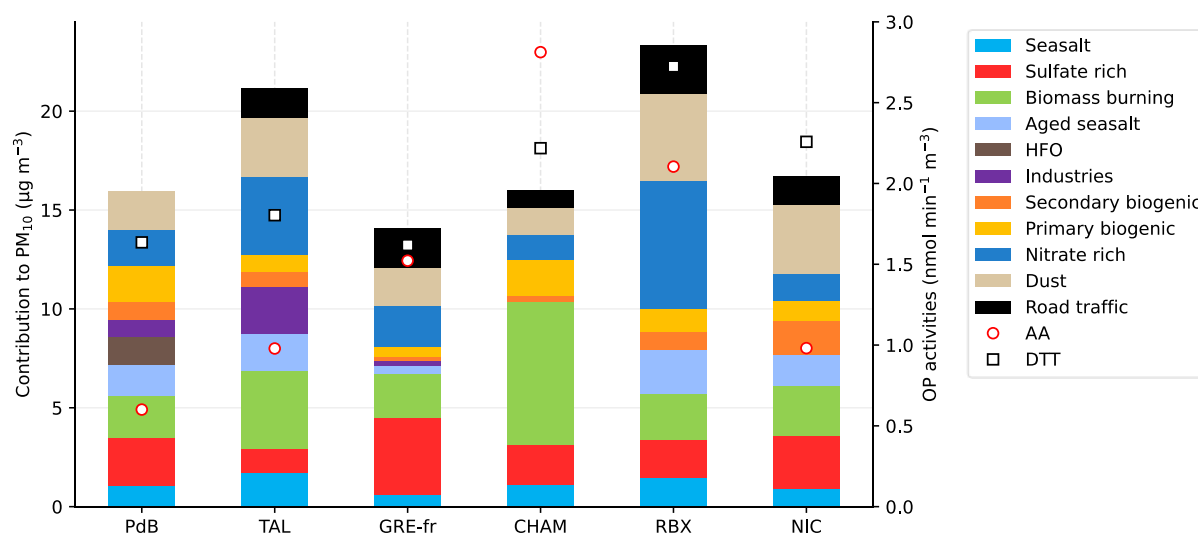
338 other sites, VIF values are below 5, indicating limited collinearity among sources. This is expected, as the PMF  
 339 analysis is constrained to avoid collinearity between sources. VIF values for each site can be found in Table S.2.

340 The presence of heteroscedasticity is commonly found when the dependent variable (or OP in this study) exhibits  
 341 a large difference between the minimum and maximum values or when the error variance varies proportionally  
 342 with an independent variable (PM<sub>10</sub> sources). The heteroscedasticity was assessed by applying the Goldfeld–  
 343 Quandt test. Table 2 presents the p-values of the Goldfeld–Quandt test, indicating homoscedasticity of OP  
 344 prediction when  $p > 0.05$ . This test reveals that heteroscedasticity was detected in CHAM, GRE-fr, NIC for OP<sub>AA</sub>  
 345 and in CHAM and TAL for OP<sub>DTT</sub> (Table 2). We observed a large difference between the cold and warm periods  
 346 for both OP<sub>AA</sub> and OP<sub>DTT</sub> in CHAM, similar to what was seen for OP<sub>AA</sub> in GRE-fr (Table S1), which can be the  
 347 reason for the presence of heteroscedasticity. For NIC and TAL, there is an insignificant difference between the  
 348 cold and warm periods, which indicates the presence of heteroscedasticity may be because of the relationship  
 349 between the PM<sub>10</sub> sources and error variance. When heteroscedasticity is detected, unweighted regression for OP  
 350 prediction according to sources may not accurately reflect the uncertainty of each source's intrinsic OP. The  
 351 scatterplots representing the relationship between the regression analysis residuals and the fitted values (for  
 352 observed OP) are available in Figures S.1 and S.2, Supplement.

353 Table 2. The p-value of the Goldfeld–Quandt heteroscedasticity test

	PdB	TAL	GRE-fr	CHAM	RBX	NIC
AA	0.15	0.78	<< 0.001	<< 0.001	0.44	0.002
DTT	0.59	<< 0.001	0.189	<< 0.001	0.56	0.91

354



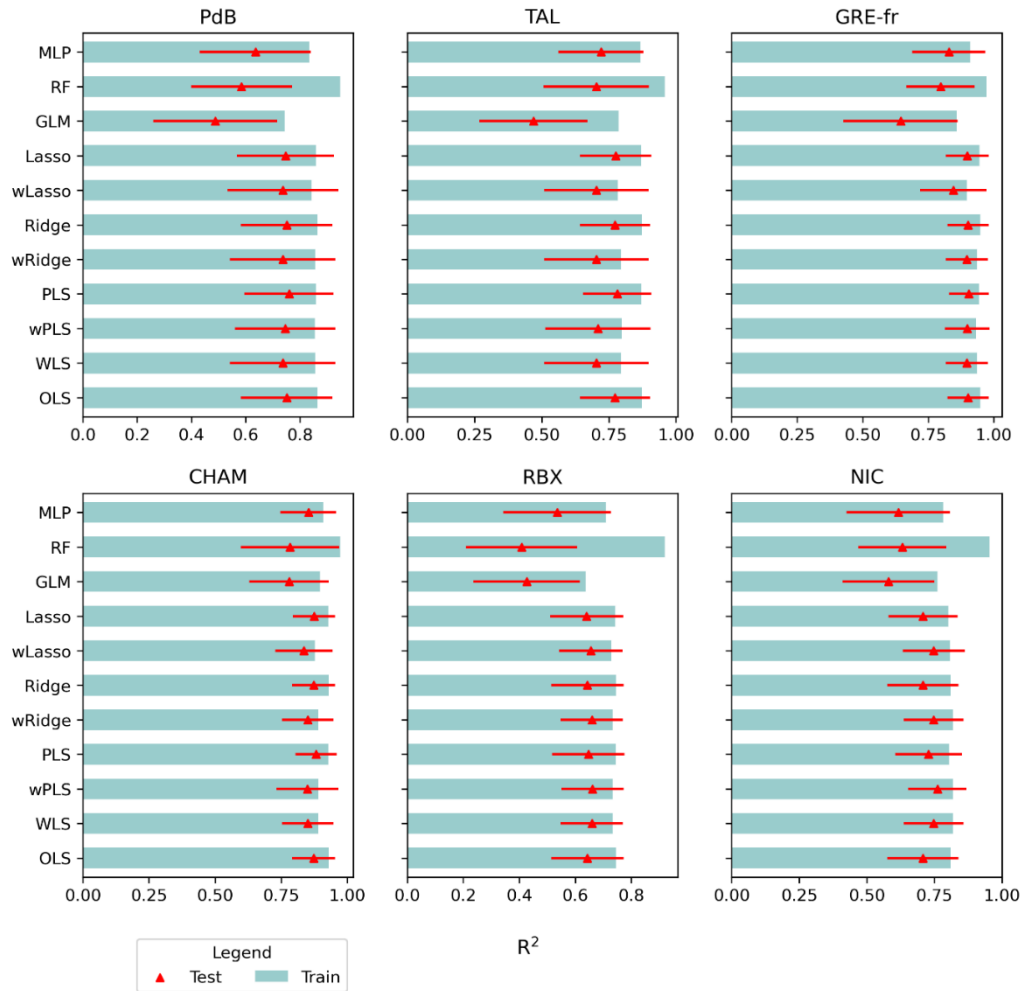
355

356 **Figure 3. The contribution of sources to PM<sub>10</sub> and the OP activities in 6 sites. The left y-axis and bar show**  
 357 **the contribution of PM sources in  $\mu\text{g m}^{-3}$ . The right y-axis, circles and squares showed the mean OP<sub>v</sub>**  
 358 **activities in  $\text{nmol min}^{-1} \text{m}^{-3}$ , with red circle for OP<sub>AA</sub> and black square for OP<sub>DTT</sub>.**

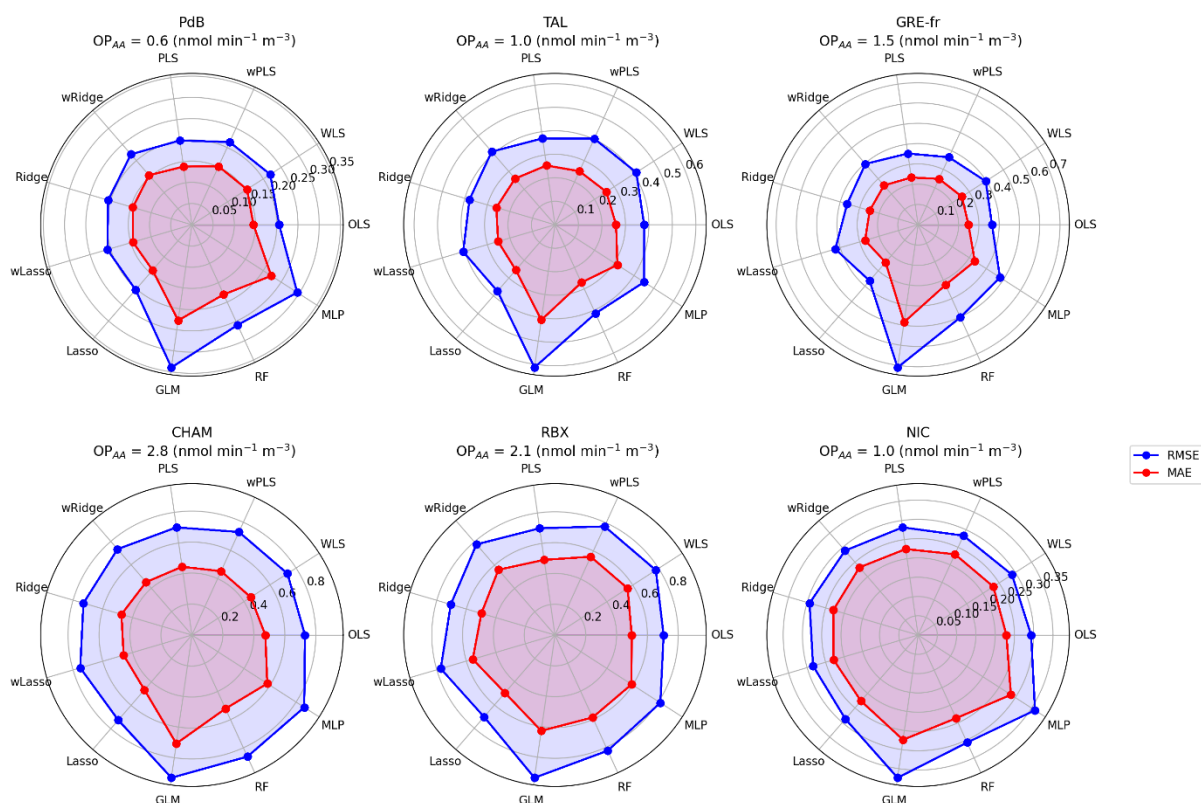
### 359 3.2. The performances of regression models

360 The 11 regression models, with or without the weighting for some of them, were tested by comparing their  
 361 performance metrics between the measured and reconstructed OPs. For each run ( $n = 500$  iterations), the  $R^2$ ,  
 362 RMSE, and MAE were computed for the testing and training dataset, resulting in 500 values for each performance

363 metric. Figure 4 presents the mean  $R^2$  values of the training data sets, the mean and the standard deviation of the  
 364 testing datasets of the  $OP_{AA}$  models across the 500 sampling iterations, and Figure 5 presents the mean RMSE and  
 365 MAE. The same result pattern was found for  $OP_{DTT}$ , as presented in the tables S.3, S.4, S.5, Supplement. The  
 366 WLS, wPLS, wRidge, and wLasso models incorporated weighting, while the OLS, PLS, Ridge, Lasso, GLM, RF,  
 367 and MLP models were unweighted.



368  
 369 **Figure 4. The  $R^2$  of 11  $OP_{AA}$  models in 6 sites. The mean  $R^2$  of training data is shown in a blue bar, the mean**  
 370  **$R^2$  of testing data is shown by a red triangle, and the red bar is the standard deviation of the  $R^2$  of the testing**  
 371 **data. The y-axis represents the models, and the x-axis denotes the  $R^2$  values.**



373

374 **Figure 5. The MAE and RMSE of 11  $OP_{AA}$  models in every site for the testing data. Blue and red lines**  
 375 **present the RMSE and the MAE, respectively. The values in the figure are the mean of RMSE and MAE of**  
 376 **500 iterations.**

377  $OP$  predictions across all sites are statistically validated, with testing  $R^2$  values observed in RBX, NIC, PdB, TAL,  
 378 CHAM, and GRE-fr being 0.66, 0.76, 0.76, 0.78, 0.87, 0.90, respectively. The lowest mean test set RMSE values  
 379 are 0.70, 0.28, 0.21, 0.37, 0.70, 0.31  $\text{nmol min}^{-1} \text{m}^{-3}$ , respectively, for the same sites. The lowest mean test set  
 380 MAE values are 0.49, 0.23, 0.14, 0.25, 0.45, and 0.21  $\text{nmol min}^{-1} \text{m}^{-3}$ , respectively. Notably, the GLM model  
 381 exhibits for all sites the lowest  $R^2$  values and the highest RMSE (Table S.3, S.4, S.5, Supplement). These results  
 382 strongly suggest that the relationship between  $OP_{AA}$  and  $PM_{10}$  sources is not log-linear.

383 Differences in MAE, RMSE, and  $R^2$  between the training and testing database for RF and MLP are significant  
 384 across the sites. Notably, RF displays a large difference in  $R^2$ , with a gap of up to 0.6 in RBX ( $R^2$  training: 0.92,  
 385  $R^2$  testing: 0.27). Similar gaps were found in RMSE and MAE. RF consistently performed best on the training set,  
 386 characterized by the highest  $R^2$  and the lowest MAE and RMSE values, but had lower set test  $R^2$  values than the  
 387 other models (except GLM). Conversely, MLP exhibited training  $R^2$  values comparable to other models but lower  
 388 test  $R^2$ . These findings suggest overfitting: the flexible algorithms identify relationships in the training data that  
 389 do not generalize to the testing data. This observation may be attributed to the limitations of data coverage, possibly  
 390 failing to fully represent the underlying relationships, leading to poor performance in testing datasets (Matsuki et  
 391 al., 2016; Benkendorf and Hawkins, 2020; Stockwell and Peterson, 2002; Wisz et al., 2008; Hernandez et al.,  
 392 2006; Hawkins, 2004; Raudys and Jain, 1991). Pearce and Ferrier (2000) recommended that the minimum number  
 393 of samples for robust performance should be over 250 for GLM model, while (Raudys and Jain, 1991) showed  
 394 that the minimum number of sample are based on the complexity of the model and the number of predictors.  
 395 Additionally, Harrell (2016) suggested that the number of predictors (PM sources) should be below the number of  
 396 samples divided by 15, a threshold not reached in this analysis. For example, in NIC, the minimum number of  
 397 samples should be 135 for the training set (9 PM sources x 15), while in total, we have only 107 samples. Therefore,

398 we can also recommend that, for optimal performance of RF, and MLP, the number of samples and PM sources  
 399 should satisfy these thresholds.

400 The WLS, OLS, wPLs, wRidge, and wLasso models show more robust performances with fewer differences  
 401 between the training and testing data. At most sites, there is very little difference between the  $R^2$ , RMSE, and MAE  
 402 of OLS and Ridge, with or without weighting, and often PLS and Lasso as well. This consistency is observed even  
 403 in the collinearity case of NIC, where  $VIF = 5$ . The difference between these models is a maximum of 0.06 in  $R^2$ ,  
 404 0.01 in MAE and 0.1 in RMSE, indicating that these models work well for OP prediction. Nevertheless, it is worth  
 405 noting that every model exhibits different assumptions that have to be respected. The assumption violations may  
 406 lead to unreliable regression coefficients (intrinsic OP) even though the prediction is good (Williams et al., 2013;  
 407 Cohen et al., 2002).

408 The best model for each site was selected based on both data characteristics (collinearity and heteroscedasticity)  
 409 and testing data performance. For sites with collinearity, the Ridge, Lasso were considered most appropriate. For  
 410 sites with heteroscedasticity, models with weights were considered the most appropriate. For sites with neither  
 411 collinearity nor heteroskedasticity, OLS and PLS were considered most appropriate. Tables 3 and 4 present the  
 412 best  $OP_{AA}$  and  $OP_{DTT}$  prediction models for each site. It follows that the best model is not necessarily the same one  
 413 for both series of OP for a given site. As a rule, the model that exhibits the best performance metrics (the best  
 414 model by error in Table 3 for  $OP_{AA}$  and Table 4 for  $OP_{DTT}$ ) is suited to the best model chosen by data  
 415 characteristics; therefore, choosing a model according to data characteristics help to more reliable in OP  
 416 predictions.

417 **Table 3. Criteria to select the best model for  $OP_{AA}$**

	<b>PdB</b>	<b>TAL</b>	<b>GRE-fr</b>	<b>CHAM</b>	<b>RBX</b>	<b>NIC</b>
<b>Collinearity</b>	No	No	No	No	No	Yes
<b>Heteroscedasticity</b>	No	No	Yes	Yes	No	Yes
<b>Best model by characteristic</b>	OLS/ PLS	OLS/ PLS	WLS/ wPLS	WLS/ wPLS	OLS/ PLS	wRidge/ wLasso
<b>Best by error</b>	PLS	PLS	wPLS	wPLS	OLS	wRidge

418 **Table 4. Criteria to select the best model for  $OP_{DTT}$**

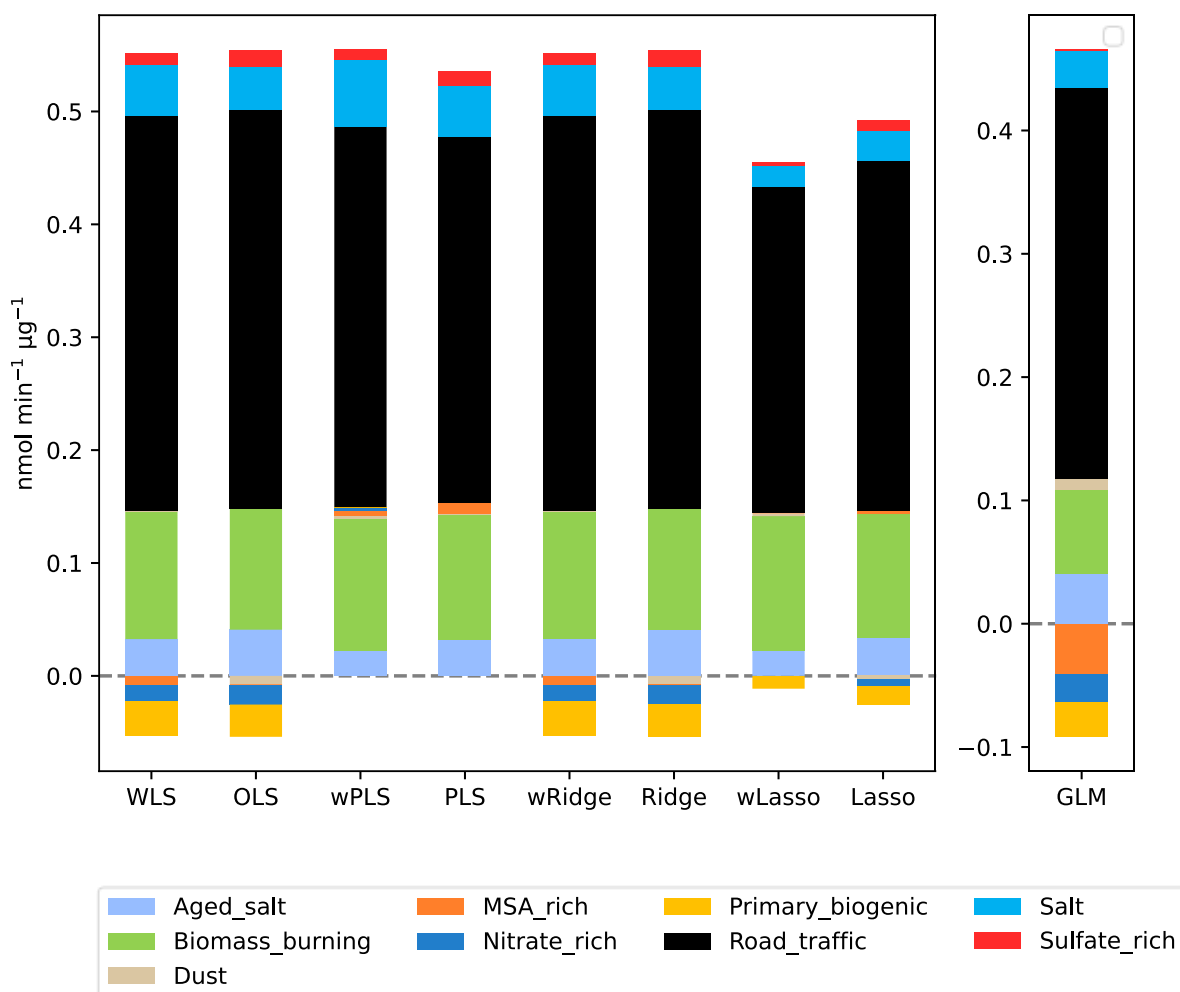
	<b>PdB</b>	<b>TAL</b>	<b>GRE-fr</b>	<b>CHAM</b>	<b>RBX</b>	<b>NIC</b>
<b>Collinearity</b>	No	No	No	No	No	Yes
<b>Heteroscedasticity</b>	No	Yes	No	Yes	No	No
<b>Best model by characteristic</b>	OLS/ PLS	WLS/ wPLS	OLS/ PLS	WLS/ wPLS	OLS/ PLS	Ridge/ Lasso
<b>Best by error</b>	OLS	wPLS	PLS	wPLS	PLS	Ridge

419

### 420 3.3. Effect of the choice of a model on intrinsic OP

421 It is particularly important to try to define the best way of calculating the more accurate PM sources intrinsic OP  
422 and the contribution of sources to OP, since these values are fundamental inputs in all the works of large-scale  
423 modelling of OP with chemical transport models (CTM) (Daellenbach et al., 2020; Vida et al., 2024). Figures 6  
424 and 7 show the variations of intrinsic OP for all the models, focusing on the results of NIC as an example. The  
425 evaluation of the 5 other sites is presented in Fig S.3 to Fig S.7 for  $OP_{AA}$  and Fig S.8 to S.12 for  $OP_{DTT}$ . The  
426 differences in equations, error term minimizations, and assumptions can explain the differences in intrinsic OP per  
427  $\mu\text{g}$  of source among the eight regression models. While the  $R^2$ , RMSE, and MAE values are similar among models  
428 (except for GLM, RF, and MLP), the intrinsic OP values significantly differ between the models with and without  
429 weighting and between the linear and non-linear regression models. The average intrinsic OP of 500 iterations is  
430 discussed in this section since these values are usually used to calculate the contribution of the  $PM_{10}$  source to OP  
431 in prior studies (Borlaza et al., 2021; Dominutti et al., 2023; Weber et al., 2018). The mean and standard deviation  
432 of intrinsic  $OP_{AA}$  and  $OP_{DTT}$  for the 6 sites are shown in Table S.6 and S.7, respectively.

433 Intrinsic  $OP_{AA}$  of  $PM_{10}$  sources at NIC is the same between WLS and wRidge and between the OLS and Ridge,  
434 revealing that the moderate collinearity of the road traffic source did not affect the estimated intrinsic  $OP_{AA}$ . PLS  
435 sets the intrinsic  $OP_{AA}$  of some sources to zero, therefore producing slightly different results. Lasso regression sets  
436 the intrinsic  $OP_{AA}$  of some sources to zero and shrinks the estimates for all other sources toward zero. GLM  
437 produces intrinsic  $OP_{AA}$  values that represent a multiplicative change on the log scale, so they are not directly  
438 comparable to the other models. However, the direction and importance of the sources are similar to the other  
439 models. Whatever the model, road traffic appears as the source with the highest intrinsic  $OP_{AA}$ , followed by  
440 biomass burning, aged salt, salt and sulfate-rich sources, in NIC. Traffic and biomass burning sources have been  
441 similarly recognized as significant contributors to  $OP_{AA}$  in prior studies (Borlaza et al., 2021; Dominutti et al.,  
442 2023; Stevanović et al., 2023). The intrinsic OP of the dominant sources is stable, indicating that all these models  
443 could give the same information about the intrinsic OP of the main sources. Conversely, the differences are larger  
444 between models for the sources with small to very small intrinsic OP (MSA rich, primary biogenic, nitrate-rich,  
445 dust), whose intrinsic OP varies from positive to negative among models.

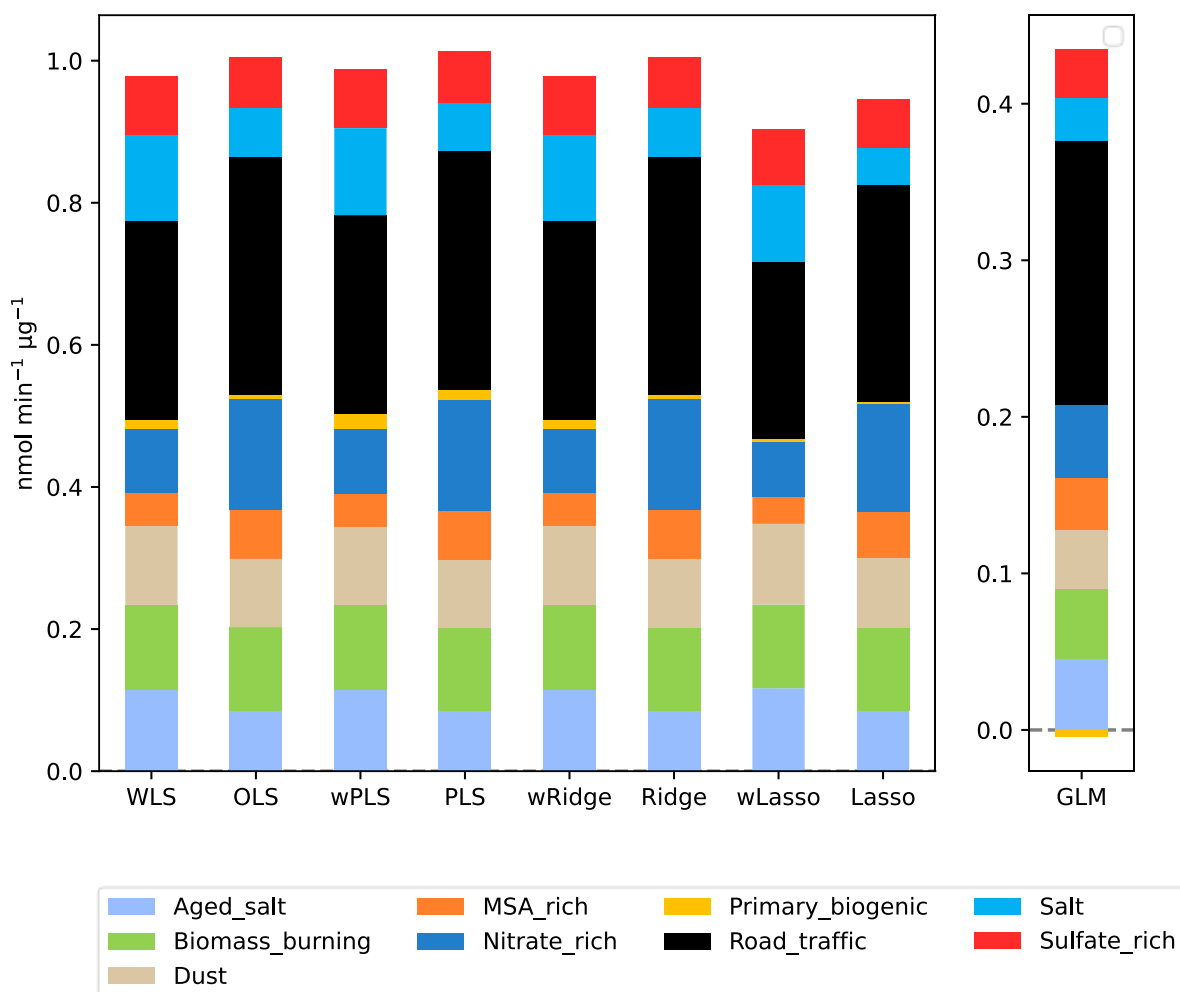


446

447 **Figure 6. Intrinsic OP<sub>AA</sub> values of the different PM<sub>10</sub> sources at Nice were obtained with the different models.**

448 The OP<sub>DTT</sub> intrinsic values in NIC (Figure 7) display minimal variation among the WLS, wPLS. This consistency  
 449 is linked to the absence of negative intrinsic values. On the other hand, even though there is the presence of  
 450 moderate collinearity, wRidge still has the same result as WLS and wPLS. In line with the OP<sub>AA</sub> results, the wLasso  
 451 and GLM models exhibit distinct responses compared to the other models. The intrinsic OP<sub>DTT</sub> of all sources varies  
 452 depending on the presence or absence of weighting. While the WLS models tend to amplify the influence of some  
 453 sources (aged sea salt, primary biogenic, sea salt, and sulfate-rich), the OLS reduces the intrinsic OP<sub>DTT</sub> of these  
 454 sources. Conversely, MSA-rich, nitrate, and road traffic sources undergo less influence in WLS but higher in OLS.  
 455 Different from OP<sub>AA</sub>, OP<sub>DTT</sub> prediction shows more variation among models, highlighting the effect of choosing  
 456 a model on evaluating the intrinsic OP<sub>DTT</sub> of PM<sub>10</sub> sources.





457

458 **Figure 7. The variations of the intrinsic OP<sub>DTT</sub> of the different PM<sub>10</sub> sources at Nice were obtained with the**  
 459 **different models.**

460 The comparison of intrinsic OP among regression models in NIC demonstrated that OP<sub>DTT</sub> and OP<sub>AA</sub> intrinsic  
 461 values exhibit variation across different models with and without weighting, illustrating that the choice of the  
 462 model significantly influences the values obtained for intrinsic OP of PM<sub>10</sub> sources (A similar pattern is observed  
 463 for all other sites and shown in Fig S.3 to Fig S.7 for OP<sub>AA</sub> and Fig S.8 to S.12 for OP<sub>DTT</sub>). Because of the difference  
 464 in intrinsic OP across models, a comparison between the best-performing and most commonly used models (OLS)  
 465 is presented in the following section to elucidate the advantage of choosing a model based on data characteristics  
 466 (section 3.4).

### 467 3.4. Comparisons between the best site-specific model and OLS

468 In this section, the intrinsic OP of the best model is selected for each site as discussed in Section 3.2, and the  
 469 intrinsic values of each source are compared to the ones returned by the OLS model. The OLS model is used as a  
 470 representative of usual practices that do not consider the database characteristics (Williams et al., 2013). Each  
 471 PM<sub>10</sub> source's average intrinsic OP value is calculated from all the 500 bootstrapping iterations for all sites where  
 472 that particular source is identified. Intrinsic OP values obtained in this way from the best model (the best model  
 473 presented in Table 3 for OP<sub>AA</sub> and Table 4 for OP<sub>DTT</sub>) encompassing all six sites are called **intrinsic OP of the**  
 474 **best model**, and the intrinsic OP values derived from the OLS from all six sites are called **intrinsic OP of the**  
 475 **reference model**.

476 A meaningful comparison of the two series of intrinsic values requires two conditions. First, intrinsic OP should  
477 be consistent across all sites. While recognizing that intrinsic OP values depend on diverse factors, we assumed  
478 the sites share fairly uniform PM<sub>10</sub> chemical source profiles in France. This is demonstrated by evaluating the  
479 Pearson distance and standardized identity distance similarity indicators of the source chemical profiles (Belis et  
480 al., 2015; Weber et al., 2019), and Figure S.13 indicates consistent profiles of sources for the 6 sites. Consequently,  
481 we could expect to observe minimal divergence in intrinsic OP values among these sites. Second, we postulate  
482 that negative intrinsic OP values are possible since previous studies have reported that total PM<sub>10</sub> intrinsic OP can  
483 be modulated due to the synergetic/antagonistic effects involving, for example, soluble copper, quinones, and  
484 bacteria (Borlaza et al., 2021; Pietrogrande et al., 2022; Samake et al., 2017; S. Wang et al., 2018; Xiong et al.,  
485 2017). Samake et al. (2017) demonstrated that the presence of bacterial cells in aerosol decreases the redox activity  
486 of Cu and 1,4-naphthoquinone, with a maximum decreasing of 60% compared to the oxidative reactivity  
487 considered individually. Pietrogrande et al. (2022) indicated that the mixture of Cu, Fe, 9,10-phenanthrene quinone  
488 and 1,2-naphthoquinone reduces the rate consumption of AA and DTT, up to 50% depending on the quantity of  
489 each chemical. Wang et al. (2018) reported that the mixing of Cu and naphthalene secondary organic aerosol  
490 (SOA) and phenanthrene SOA only got half of DTT rate consumption compared to the consumption when  
491 considered separately. Xiong et al. (2017) showed the presence of antagonists in the interaction of Fe and quinones,  
492 nevertheless, much lower than those in the other studies (under 10%). These references reported that the  
493 antagonistic effects of a mixture can significantly reduce the consumption rate of OP<sub>DTT</sub> and OP<sub>AA</sub>, and this impact  
494 varies widely from 10% to 60% depending on the type of chemical species and the quantity of each species in the  
495 mixture. Consequently, we consider here that the intrinsic OP value of an individual site for a given source could  
496 be negative only within a range of at most 60% of the mean combined intrinsic OP value of this source across all  
497 sites. Negative intrinsic OP exceeding this criterion may result from the mathematical construction of the model.  
498 The comparison of intrinsic OP<sub>AA</sub> of the best and reference model is presented in 3.4.1 and that of OP<sub>DTT</sub> is shown  
499 in 3.4.2.

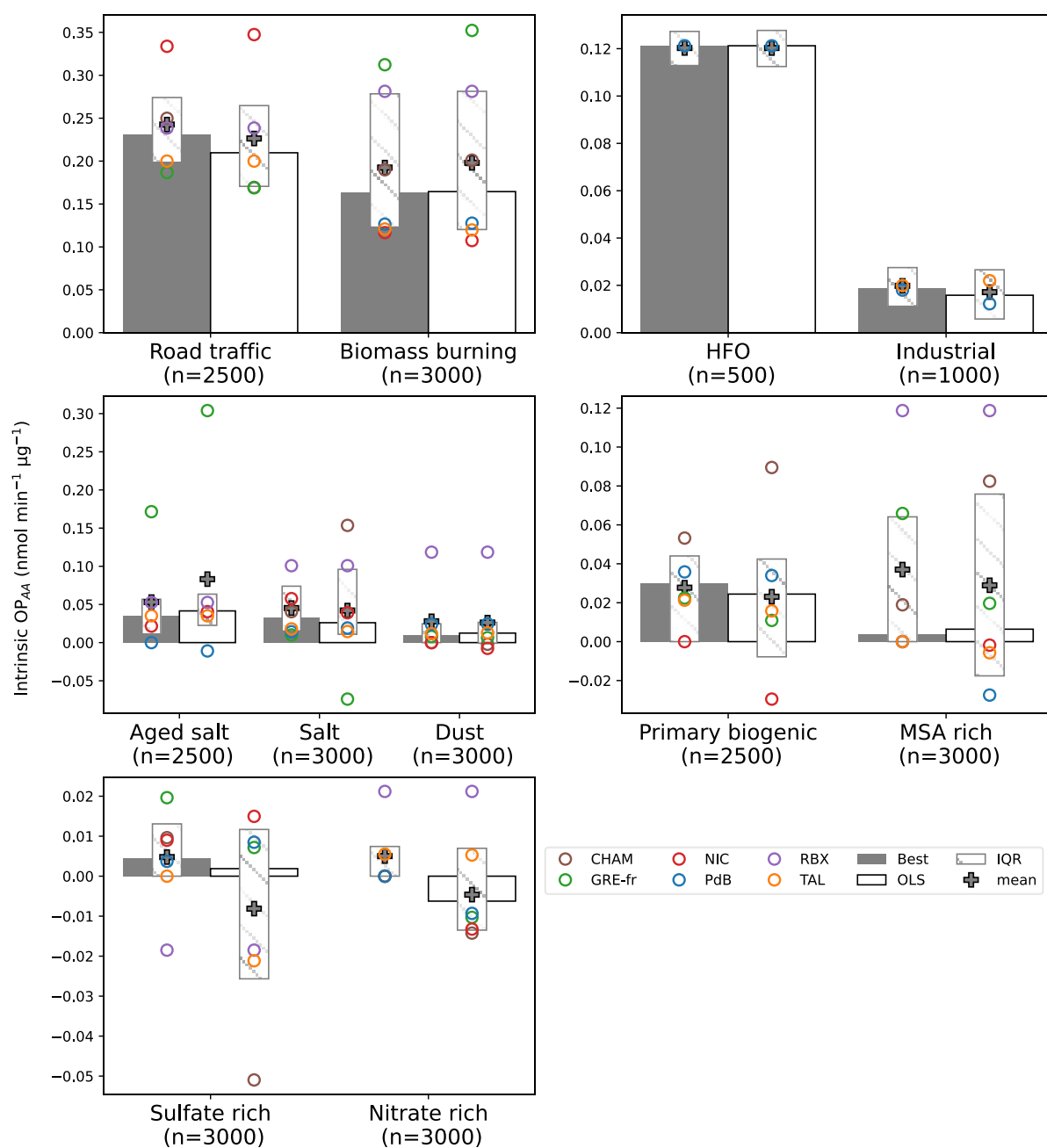
#### 500 3.4.1. OP<sub>AA</sub> activities

501 The results of the comparison of OP<sub>AA</sub> intrinsic values (Figure 8 and Table S.8) show that the anthropogenic  
502 sources get the highest intrinsic OP values in both the best and reference models. Among these sources, road traffic  
503 appears as the most prominent potent fraction, followed by biomass burning, HFO, and industrial. These results  
504 are aligned with prior research (Calas et al., 2019; Daellenbach et al., 2020; Dominutti et al., 2023; Fadel et al.,  
505 2023; Fang et al., 2016; in 't Veld et al., 2023; Weber et al., 2018; Zhang et al., 2020) which has highlighted the  
506 sensitivity of OP<sub>AA</sub> to concentrations of metals, black carbon, and organic carbon. The differences between the  
507 best and reference models were insignificant for these sources, demonstrating that **the best and reference models**  
508 **consistently captured similar patterns for the most critical sources of OP activities.**

509 However, the interquartile ranges (IQR) of the intrinsic OP values are consistently narrower for the best models  
510 across all sources, accounting for less divergence in intrinsic OP values across sites. Moreover, the median intrinsic  
511 OP values obtained from the best model closely approximated the mean values, indicating the absence of extreme  
512 intrinsic OP values. For instance, in the case of road traffic, the mean and median values were 0.24 and 0.23 nmol  
513 min<sup>-1</sup> μg<sup>-1</sup>, respectively. Conversely, the reference model exhibited a large difference between the mean and  
514 median values, implying lower consistency across sites and sampling iterations. The same result was observed in  
515 biomass burning source, in which the median and mean intrinsic OP in the best model had fewer discrepancies.  
516 Further, the biomass burning intrinsic OP in GRE-fr of the best model is more consistent with those in other sites  
517 (best: 0.30 nmol min<sup>-1</sup> μg<sup>-1</sup>, reference: 0.35 nmol min<sup>-1</sup> μg<sup>-1</sup>).

518 When considering sources with low intrinsic OP, the variability can be larger between the two methods. As an  
519 example, for the sulfate-rich sources, the median intrinsic OP values were positive (0.002 nmol min<sup>-1</sup> μg<sup>-1</sup>), while  
520 the mean intrinsic OP values were negative (-0.008 nmol min<sup>-1</sup> μg<sup>-1</sup>). The mean intrinsic OP in the best model

521 exhibited fewer negative values in individual sites than in the reference model (for aged salt, salt, primary biogenic,  
 522 MSA rich, sulfate-rich and nitrate-rich). In addition, the best model showed the less disparate intrinsic OP among  
 523 individual sites for instance, the aged salt sources in GRE-fr and the primary biogenic and salt sources in CHAM.  
 524 Furthermore, the best model displayed an intrinsic OP meaningful in terms of geochemical, which showed in the  
 525 source of salt, primary biogenic, sulfate-rich. For instance, in the reference model, the average intrinsic OP of the  
 526 primary biogenic in NIC (-0.03 nmol min<sup>-1</sup> μg<sup>-1</sup>), the intrinsic OP of salt in GRE-ft (-0.07 nmol min<sup>-1</sup> μg<sup>-1</sup>) as well  
 527 as the sulfate-rich source in CHAM (-0.05 nmol min<sup>-1</sup> μg<sup>-1</sup>) represented a 100% reduction compared to the mean  
 528 intrinsic OP of all sites. Moreover, the negative intrinsic OP was observed in NIC (Primary biogenic), and some  
 529 extreme values in GRE-fr (aged salt, salt), CHAM (salt, primary biogenic, MSA-rich) (where heteroscedasticity  
 530 was presented) in the OLS model, underscores that the model assumptions on data characteristics proving false  
 531 could impact the accuracy of OP prediction. Consequently, these results highlight the advantage of considering  
 532 the data in model selection.



533

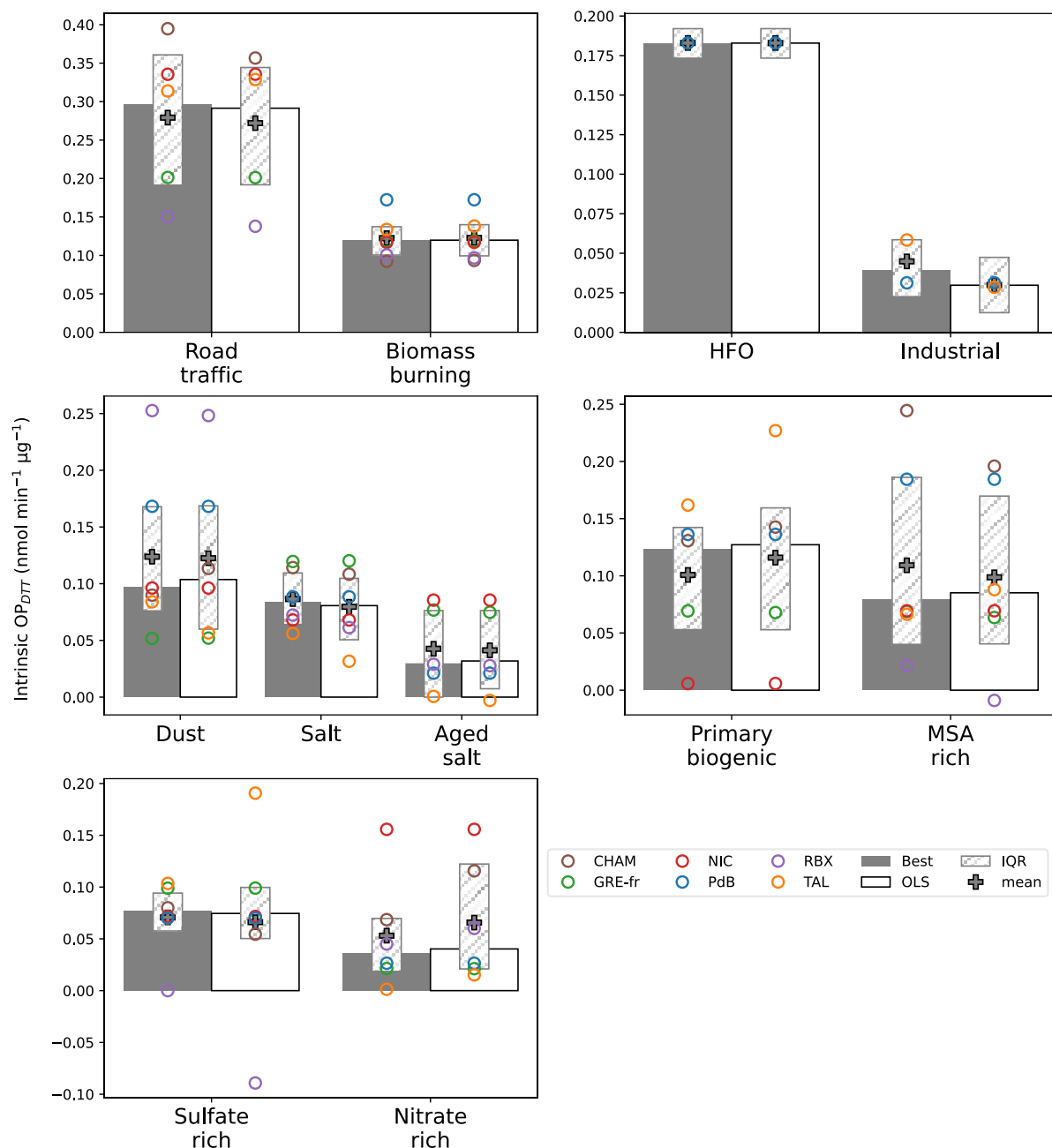
534 **Figure 8. Intrinsic  $OP_{AA}$  estimated by the best and the reference methods in the 6 sites. The y-axis represents**  
535 **the intrinsic OP values in  $nmol\ min^{-1}\ \mu g^{-1}$ , the x-axis represents the sources. The grey bars are the median**  
536 **intrinsic OP values of the best models in the 6 sites ( $n = 500$  bootstrapping \* number of sites where the given**  
537 **source is detected) for each source. The white bars are the same median intrinsic OP values for the reference**  
538 **(OLS) model. The grey plus symbol represents the mean of intrinsic OP values. The hatched bars are the**  
539 **interquartile ranges of the intrinsic OP values. The dots represent the mean intrinsic OP of all sites,**  
540 **including grey – Chamonix, green – Grenoble, red – Nice, blue – Port-de-Bouc, purple – Roubaix, and**  
541 **orange-Talence.**

542 The detailed comparison of intrinsic  $OP_{AA}$  between the best and reference models is categorized into four groups  
543 and discussed in detail in section S9. These groups include (1) anthropogenic sources without nitrate and sulfate  
544 (road traffic, biomass burning, HFO, industrial), (2) natural inorganic sources (aged sea salt, sea salt, dust), (3)  
545 biogenic sources (primary biogenic, MSA rich), and (4) nitrate and sulfate-rich sources.

#### 546 3.4.2. $OP_{DTT}$ activities

547 Similar to  $OP_{AA}$ , for  $OP_{DTT}$  the IQR of the best model is narrower for most of the sources than the IQR of the  
548 reference model (OLS). Except for the road traffic, industrial, and MSA-rich, the IQR is slightly higher in the best  
549 model (Figure 9 and Table S.9). In the two models, the mean intrinsic OP is essentially unchanged, where the  
550 traffic is the most critical source ( $0.27\pm 0.10$ ), followed by HFO ( $0.18\pm 0.01$ ), biomass burning ( $0.12\pm 0.03$ ), dust  
551 ( $0.12\pm 0.07$ ), primary biogenic (best:  $0.10\pm 0.06$ , reference:  $0.12\pm 0.08$ ) and MSA rich (best:  $0.11\pm 0.09$ , reference:  
552  $0.09\pm 0.09$ ). The minimum difference between the two models in the dominant sources again confirms the  
553 conclusion in the  $OP_{AA}$  comparison, demonstrating **the similar pattern of the best and the reference model in**  
554 **the most crucial sources of OP**. For both best and reference,  $OP_{DTT}$  activities showed sensitivity to more sources  
555 than  $OP_{AA}$ , as discussed in previous studies (Borlaza et al., 2021; Calas et al., 2019; Dominutti et al., 2023; Fadel  
556 et al., 2023).

557 While the best and reference models give the same mean intrinsic  $OP_{DTT}$  of all sites, the mean  $OP_{DTT}$  at each  
558 individual site can vary substantially between the two models. The best model exhibited the positive intrinsic OP  
559 for all sources, while the reference model displayed negative intrinsic OP in RBX (MSA-rich and sulfate-rich).  
560 Especially in the case of sulfate-rich in RBX, the negative intrinsic OP in the reference model passed the threshold  
561 of negative value, which presented a 110% reduction compared to the mean intrinsic OP of all sites. This is also  
562 found in the  $OP_{AA}$  comparison, which confirmed that the best model generates a geochemical meaningful OP  
563 intrinsic. In addition, the best model exhibited consistent intrinsic OP across sites, especially for the source of dust,  
564 salt, primary biogenic, sulfate-rich in TAL (heteroscedasticity is presented in this site), where intrinsic OP in TAL  
565 in the best model is more similar to the other sites. For instance, the reference model presented that the intrinsic  
566 OP in TAL is  $0.20\ nmol\ min^{-1}\ \mu g^{-1}$ , far from the mean of all sites ( $0.07\ nmol\ min^{-1}\ \mu g^{-1}$ ). We observed the same  
567 for OP intrinsic of nitrate-rich source in CHAM (where the heteroscedasticity is detected), which displayed the  
568 less dissimilar of CHAM with the other site in the best model. This again validates the conclusion in  $OP_{AA}$   
569 comparison, demonstrating that respecting model assumption is essential to obtain a robust OP SA result.



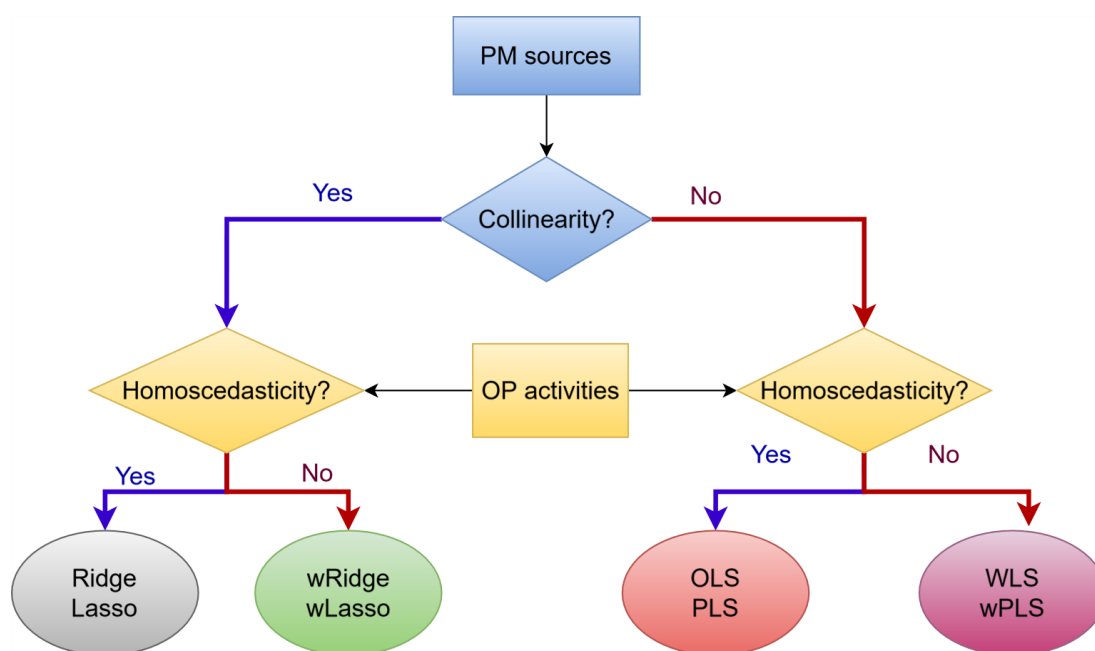
570

571 **Figure 9. Intrinsic  $OP_{DTT}$  was estimated by the best and the reference methods in the 6 sites. The y-axis**  
 572 **represents the intrinsic OP values in  $\text{nmol min}^{-1} \mu\text{g}^{-1}$ , the x-axis represents the sources. The grey bars are**  
 573 **the median intrinsic OP values of the best models in the 6 sites ( $n = 500$  bootstrapping \* number of sites**  
 574 **where the given source is detected) for each source. The white bars are the same median intrinsic OP values**  
 575 **for the reference (OLS) model. The grey plus symbol represents the mean of intrinsic OP values. The**  
 576 **hatched bars are the interquartile ranges of the Intrinsic OP values. The dots represent the mean intrinsic**  
 577 **OP of all sites, including grey – Chamonix, green – Grenoble, red – Nice, blue – Port-de-Bouc, purple –**  
 578 **Roubaix, and orange-Talence.**

579 The comparison of intrinsic OP between the best models and the reference model highlights the importance of  
 580 considering the database characteristics when selecting a model for OP SA. For all the datasets studied here, using  
 581 the best model for each site delivered more robust results with reduced uncertainty, reduced differences in intrinsic

582 OP across sites, and provided a more geochemically meaningful intrinsic OP. The recommendation for selecting  
583 a model based on the characteristics of the database is presented in section 3.5.

584 **3.5. Guidelines for the selection of regression model for OP SA.**



585

586 **Figure 10. Workflow in model selection considering the characteristics of data**

587

588 Our results have highlighted the benefits of choosing a model that matches the characteristics of the data to improve  
589 the robustness of OP SA method. For this reason, this section develops a workflow to help make model selection  
590 decisions. Before selecting a regression for OP SA, the first question is whether the PM<sub>10</sub> sources are collinear and  
591 the second is whether the residual variance of the regression between OP and PM<sub>10</sub> mass is constant. These two  
592 questions represent the characteristics of PM<sub>10</sub> sources and OP activities, which vary according to the study site.

593 For data exhibiting collinearity between sources and generating a residual variance that varies according to the  
594 value of the PM<sub>10</sub> sources, weighted regularisation regression can help to reduce collinearity and to match the  
595 model assumption about the residual. On the other hand, the unweighted Ridge and Lasso are introduced for data  
596 showing collinearity and homoscedasticity. Additionally, data with no collinearity are suitable for OLS and  
597 unweighted PLS in the case of homoscedasticity, while WLS, weighted PLS are used for data with  
598 heteroscedasticity.

599 If the number of predictors (PM<sub>10</sub> sources) is below the number of samples divided by 15, RF and MLP can also  
600 be employed to capture possible non-linear relationships between the OP and PM<sub>10</sub> sources. However, cross-  
601 validation must be used to ensure that there is no over-fitting. In addition, these models do not estimate intrinsic  
602 OP (nmol min<sup>-1</sup>μg<sup>-1</sup>) but only the importance of each PM<sub>10</sub> source to the OP prediction. This is a large drawback  
603 since the intrinsic OP of sources is a must for the modelling effort of OP with CTM. However, RF and MLP could  
604 be useful for OP prediction in the case of larger datasets generated by online instruments.

605 For each data characteristic there is more than one model that suits. Out-of-sample performance metrics should be  
606 employed to identify the most accurate of these models.

607 Finally, these techniques of OP apportionment could not be well performed with uncertain PMF-derived sources.  
608 The PMF results sometimes do not adequately represent PM mass concentration for several reasons, such as the  
609 lack of a trace species to identify a source, an insufficient sample size, the source contribution being too small to

610 be identified (under 1%), or collinearity matters. The important information could be missed because of these  
611 problems in PMF implementation, which is apprehended by the model's low accuracy. Our study did not encounter  
612 this problem since the PMF is harmonized and performed according to European recommendations which could  
613 well perform the regression technique and allow to obtain a very satisfactory successive OP modelled in  
614 comparison to observations after regression techniques ( $R^2$  from 0.7 to 0.9). However, this problem could  
615 potentially happen, and for these cases, we could recommend either subtracting the total source contribution from  
616 PM mass concentration to get a part that PMF cannot simulate. The information in this part may contain vital  
617 sources. Alternatively, it is possible to re-execute the PMF to validate the result and ensure the robustness of the  
618 chemical profile and the contribution of sources.

619

620 Limitations and perspectives of the study:

- 621 - This study compares eight regression models but is not exhaustive; further research could add more  
622 regression techniques to evaluate result variations across models. The potential techniques that could be  
623 applied for OP SA are gradient boosting techniques for resolving regression models, or supervised  
624 machine learning techniques which allows the investigation of linear and non-linear regression  
625 relationships. However, the consistently strong performance of ordinary linear regression across six  
626 locations in France suggests that there may be little to gain from applying more complex models in areas  
627 with similar  $PM_{10}$  sources.
- 628 - PMF coupled with a regression model remains a popular approach for OP SA. Notably, the uncertainties  
629 in PMF are typically addressed in chemical profiles, but not in contributions. Incorporating uncertainty  
630 from variations in contribution into models could enhance their robustness compared to relying only on  
631 absolute PMF results.
- 632 - Observations ranged between 100 and 200 samples at each site, which may be insufficient to obtain a fair  
633 performance of GLM, decision trees and neural network models even though this number of samples is  
634 sufficient to address SA through the PMF model for offline analyses. Therefore, this study outlines well  
635 the limitations of GLM, RF, and MLP for offline datasets. Future investigations should be performed in  
636 an extended dataset, such as long-term or real-time measurement data, to investigate the performance of  
637 machine learning algorithms.
- 638 - This study only focused on the two most popular OP assays of  $PM_{10}$  ( $OP_{DTT}$  and  $OP_{AA}$ ). However, there  
639 are actually various OP assays, such as  $OP_{DCFH}$ ,  $OP_{OH}$ ,  $OP_{FOX}$ ,  $OP_{GSH}$ ,  $OP_{ESR}$  and different sizes of PM  
640 ( $PM_1$ ,  $PM_{2.5}$ ,  $PM_5$ ). Further research should include more OP assays, which can be helpful in evaluating  
641 the performance of various regression models for different OP and different PM sizes.
- 642 - This study used the analytical uncertainty as the weighting for the weighted model. However, the  
643 weighting can be selected based on different ways, as reported by Montgomery et al. (2012): (1) Prior  
644 information from the theoretical model, (2) Using the residual extracted from the OLS model, (3) The  
645 selecting of weighting based on the uncertainty of instrument if the dependent variable measured by a  
646 different method and (4) If the dependent variable is the average of different observations, the weighting  
647 selected based on the error of these observations.

648

#### 649 **4. Conclusion**

650 The results of the OP SA marked an important milestone as they were revealed for the first time through the use  
651 of eight regression models, including OLS, WLS, PLS, GLM, Ridge, Lasso, RF and MLP. This in-depth analysis  
652 was carried out on a complete set of data collected from six sites with different characteristics. The approach of  
653 selecting a suitable model for each site based on specific data characteristics resulted in a more consistent intrinsic

654 OP across sites, in stark contrast to the variation observed when using the basic OLS model. The revelations of the  
655 study have provided concrete recommendations for the judicious selection of an appropriate regression model  
656 based on the unique characteristics of the dataset. These guidelines should help to improve the accuracy of OP  
657 assessments and contribute to the refinement of air quality assessment methods. In addition, the implications of  
658 this research extend to the implementation of OP monitoring as a new measure of air quality, particularly on  
659 European supersites. As this initiative aligns with the ongoing revision process of the European Directive  
660 2008/50/CE, the study's findings assume a pivotal role in shaping the methodologies underpinning air quality  
661 assessments at a broader regulatory level.

#### 662 **Code availability**

663 The code is available in <https://doi.org/10.5281/zenodo.11071884> (Dinh Ngoc, 2024).

#### 664 **Data availability**

665 The datasets could be made available upon request by contacting the corresponding author.

#### 666 **Author contributions**

667 VDNT performed the data analysis for the OP source apportionment setup. GU, JLJ mentoring, supervision, and  
668 validation of the methodology and results. IH, PD, and VDNT worked on the result visualization. OF, JLJ, and  
669 GU acquired fundings for the original PM sampling and analysis. VDNT wrote the original draft. All authors  
670 reviewed and edited the manuscript.

#### 671 **Competing interests**

672 The authors declare that they have no conflict of interest.

#### 673 **Acknowledgments**

674 The authors would like to express their sincere gratitude to many people of the Air-O-Sol analytical platform at  
675 IGE (including S. Darfeuil, R. Elazzouzi, and T Madhbi), to R. Aujay (Ineris) for sample management at TAL and  
676 RBX, to L. Alleman (IMT Nord-Europe) and N. Bonnaire (LSCE) for part of the chemical analyses for some sites,  
677 and to all the personnel within the AASQA in charge of the sites for their contribution in conducting the dedicated  
678 sample collection. The authors would like to thank S. Weber for running the PMF model in his previous  
679 professional life.

#### 680 **Financial support**

681 The PhD grant of VDNT was funded by grant PR-PRE-2021, UGA-UGA 2022-16 FUGA-Fondation Air Liquide,  
682 and ANR ABS (ANR-21-CE01-0021-01). Analytical work on OP was funded through ANR GET OP STAND  
683 (ANR-19-CE34-0002), MOBILAIR and ACME IDEX projects at UGA (ANR-15-IDEX-02). The sampling and  
684 chemical analyses performed at TAL, GRE, RBX, PdB and NIC sites have been partly funded by the French  
685 Ministry of Environment in the frame of the CARA program. The present work was also supported by European  
686 Union's Horizon 2020 research and innovation program under grant agreement 101036245 (RI-URBANS) for the  
687 Post-doc salary of Pamela Dominutti.



688 **Reference**

- 689 Akhtar, McWhinney, R. D., Rastogi, N., Abbatt, J. P. D., Evans, G. J., and Scott, J. A.: Cytotoxic and  
690 proinflammatory effects of ambient and source-related particulate matter (PM) in relation to the production of  
691 reactive oxygen species (ROS) and cytokine adsorption by particles, *Inhal. Toxicol.*, 22, 37–47,  
692 <https://doi.org/10.3109/08958378.2010.518377>, 2010.
- 693 Akhtar, A., Islamia, J. M., Masood, S., Islamia, J. M., Masood, A., and Islamia, J. M.: Prediction and Analysis of  
694 Pollution Levels in Delhi Using Multilayer Perceptron, <https://doi.org/10.1007/978-981-10-3223-3>, 2018.
- 695 Alleman, L. Y., Lamaison, L., Perdrix, E., Robache, A., and Galloo, J. C.: PM10 metal concentrations and source  
696 identification using positive matrix factorization and wind sectoring in a French industrial zone, *Atmos. Res.*, 96,  
697 612–625, <https://doi.org/10.1016/j.atmosres.2010.02.008>, 2010.
- 698 Ayres, J. G., Borm, P., Cassee, F. R., Castranova, V., Donaldson, K., Ghio, A., Harrison, R. M., Hider, R., Kelly,  
699 F., Kooter, I. M., Marano, F., Maynard, R. L., Mudway, I., Nel, A., Sioutas, C., Smith, S., Baeza-Squiban, A.,  
700 Cho, A., Duggan, S., and Froines, J.: Evaluating the toxicity of airborne particulate matter and nanoparticles by  
701 measuring oxidative stress potential - A workshop report and consensus statement, *Inhal. Toxicol.*, 20, 75–99,  
702 <https://doi.org/10.1080/08958370701665517>, 2008.
- 703 Bates, J. T., Weber, R. J., Abrams, J., Verma, V., Fang, T., Klein, M., Strickland, M. J., Sarnat, S. E., Chang, H.  
704 H., Mulholland, J. A., Tolbert, P. E., and Russell, A. G.: Reactive Oxygen Species Generation Linked to Sources  
705 of Atmospheric Particulate Matter and Cardiorespiratory Effects, *Environ. Sci. Technol.*, 49, 13605–13612,  
706 <https://doi.org/10.1021/acs.est.5b02967>, 2015.
- 707 Bates, J. T., Weber, R. J., Verma, V., Fang, T., Ivey, C., Liu, C., Sarnat, S. E., Chang, H. H., Mulholland, J. A.,  
708 and Russell, A.: Source impact modeling of spatiotemporal trends in PM2.5 oxidative potential across the eastern  
709 United States, *Atmos. Environ.*, 193, 158–167, <https://doi.org/10.1016/j.atmosenv.2018.08.055>, 2018.
- 710 Bates, J. T., Fang, T., Verma, V., Zeng, L., Weber, R. J., Tolbert, P. E., Abrams, J. Y., Sarnat, S. E., Klein, M.,  
711 Mulholland, J. A., and Russell, A. G.: Review of Acellular Assays of Ambient Particulate Matter Oxidative  
712 Potential: Methods and Relationships with Composition, Sources, and Health Effects, *Environ. Sci. Technol.*, 53,  
713 4003–4019, <https://doi.org/10.1021/acs.est.8b03430>, 2019.
- 714 Beelen, R., Stafoggia, M., Raaschou-Nielsen, O., Andersen, Z. J., Xun, W. W., Katsouyanni, K., Dimakopoulou,  
715 K., Brunekreef, B., Weinmayr, G., Hoffmann, B., Wolf, K., Samoli, E., Houthuijs, D., Nieuwenhuijsen, M., Oudin,  
716 A., Forsberg, B., Olsson, D., Salomaa, V., Lanki, T., Yli-Tuomi, T., Oftedal, B., Aamodt, G., Nafstad, P., De Faire,  
717 U., Pedersen, N. L., Östenson, C. G., Fratiglioni, L., Penell, J., Korek, M., Pyko, A., Eriksen, K. T., Tjønneland,  
718 A., Becker, T., Eeftens, M., Bots, M., Meliefste, K., Wang, M., Bueno-De-Mesquita, B., Sugiri, D., Krämer, U.,  
719 Heinrich, J., De Hoogh, K., Key, T., Peters, A., Cyrus, J., Concin, H., Nagel, G., Ineichen, A., Schaffner, E.,  
720 Probst-Hensch, N., Dratva, J., Ducret-Stich, R., Vilier, A., Clavel-Chapelon, F., Stempfelet, M., Gioni, S., Krogh,  
721 V., Tsai, M. Y., Marcon, A., Ricceri, F., Sacerdote, C., Galassi, C., Migliore, E., Ranzi, A., Cesaroni, G., Badaloni,  
722 C., Forastiere, F., Tamayo, I., Amiano, P., Dorronsoro, M., Katsoulis, M., Trichopoulou, A., Vineis, P., and Hoek,  
723 G.: Long-term exposure to air pollution and cardiovascular mortality: An analysis of 22 European cohorts,  
724 *Epidemiology*, 25, 368–378, <https://doi.org/10.1097/EDE.000000000000076>, 2014.
- 725 Belis, C. A., Karagulian, F., Larsen, B. R., and Hopke, P. K.: Critical review and meta-analysis of ambient  
726 particulate matter source apportionment using receptor models in Europe,  
727 <https://doi.org/10.1016/j.atmosenv.2012.11.009>, April 2013.
- 728 Belis, C. A., Karagulian, F., Amato, F., Almeida, M., Artaxo, P., Beddows, D. C. S., Bernardoni, V., Bove, M. C.,  
729 Carbone, S., Cesari, D., Contini, D., Cuccia, E., Diapouli, E., Eleftheriadis, K., Favez, O., El Haddad, I., Harrison,  
730 R. M., Hellebust, S., Hovorka, J., Jang, E., Jorquera, H., Kammermeier, T., Karl, M., Lucarelli, F., Mooibroek, D.,  
731 Nava, S., Nøjgaard, J. K., Paatero, P., Pandolfi, M., Perrone, M. G., Petit, J. E., Pietrodangelo, A., Pokorná, P.,  
732 Prati, P., Prevot, A. S. H., Quass, U., Querol, X., Saraga, D., Sciare, J., Sfetsos, A., Valli, G., Vecchi, R., Vestenius,  
733 M., Yubero, E., and Hopke, P. K.: A new methodology to assess the performance and uncertainty of source  
734 apportionment models II: The results of two European intercomparison exercises, *Atmos. Environ.*, 123, 240–250,  
735 <https://doi.org/10.1016/j.atmosenv.2015.10.068>, 2015.
- 736 Bell, M. L., Samet, J. M., and Dominici, F.: Time-series studies of particulate matter, *Annu. Rev. Public Health*,  
737 25, 247–280, <https://doi.org/10.1146/annurev.publhealth.25.102802.124329>, 2004.
- 738 Benkendorf, D. J. and Hawkins, C. P.: Effects of sample size and network depth on a deep learning approach to  
739 species distribution modeling, *Ecol. Inform.*, 60, <https://doi.org/10.1016/j.ecoinf.2020.101137>, 2020.

740 Borlaza: Disparities in particulate matter (PM10) origins and oxidative potential at a city scale (Grenoble, France)  
741 - Part 2: Sources of PM10 oxidative potential using multiple linear regression analysis and the predictive  
742 applicability of multilayer perceptron n, *Atmos. Chem. Phys.*, 21, 9719–9739, [https://doi.org/10.5194/acp-21-](https://doi.org/10.5194/acp-21-9719-2021)  
743 9719-2021, 2021.

744 Borlaza, L., Weber, S., Uzu, G., Jacob, V., Cañete, T., Micallef, S., Trébuchon, C., Slama, R., Favez, O., and  
745 Jaffrezo, J.-L.: Disparities in particulate matter (PM10) origins and oxidative potential at a city scale (Grenoble,  
746 France) - Part 1: Source apportionment at three neighbouring sites, *Atmos. Chem. Phys.*, 21, 5415–5437,  
747 <https://doi.org/10.5194/acp-21-5415-2021>, 2021a.

748 Borlaza, L., Weber, S., Jaffrezo, J. L., Houdier, S., Slama, R., Rieux, C., Albinet, A., Micallef, S., Trébluchon, C.,  
749 and Uzu, G.: Disparities in particulate matter (PM10) origins and oxidative potential at a city scale (Grenoble,  
750 France) - Part 2: Sources of PM10 oxidative potential using multiple linear regression analysis and the predictive  
751 applicability of multilayer perceptron n, *Atmos. Chem. Phys.*, 21, 9719–9739, [https://doi.org/10.5194/acp-21-](https://doi.org/10.5194/acp-21-9719-2021)  
752 9719-2021, 2021b.

753 Boulard, H. and Wellekens, C. J.: Speech pattern discrimination and multilayer perceptrons, *Comput. Speech*  
754 *Lang.*, 3, 1–19, [https://doi.org/https://dx.doi.org/10.1016/0885-2308\(89\)90011-9](https://doi.org/https://dx.doi.org/10.1016/0885-2308(89)90011-9), 1989.

755 Breiman, L.: RFRSF: Employee Turnover Prediction Based on Random Forests and Survival Analysis, *Mach.*  
756 *Learn.*, 12343 LNCS, 503–515, [https://doi.org/10.1007/978-3-030-62008-0\\_35](https://doi.org/10.1007/978-3-030-62008-0_35), 2001.

757 Brook, R. D., Rajagopalan, S., Pope, C. A., Brook, J. R., Bhatnagar, A., Diez-Roux, A. V., Holguin, F., Hong, Y.,  
758 Luepker, R. V., Mittleman, M. A., Peters, A., Siscovick, D., Smith, S. C., Whitsel, L., and Kaufman, J. D.:  
759 Particulate matter air pollution and cardiovascular disease: An update to the scientific statement from the american  
760 heart association, *Circulation*, 121, 2331–2378, <https://doi.org/10.1161/CIR.0b013e3181d8e1>, 2010.

761 Brown, S. G., Eberly, S., Paatero, P., and Norris, G. A.: Methods for estimating uncertainty in PMF solutions:  
762 Examples with ambient air and water quality data and guidance on reporting PMF results, *Sci. Total Environ.*,  
763 518–519, 626–635, <https://doi.org/10.1016/j.scitotenv.2015.01.022>, 2015.

764 Calas, A., Uzu, G., Martins, J. M. F., Voisin, Di., Spadini, L., Lacroix, T., and Jaffrezo, J. L.: The importance of  
765 simulated lung fluid (SLF) extractions for a more relevant evaluation of the oxidative potential of particulate  
766 matter, *Sci. Rep.*, 7, 1–12, <https://doi.org/10.1038/s41598-017-11979-3>, 2017.

767 Calas, A., Uzu, G., Kelly, F. J., Houdier, S., Martins, J. M. F., Thomas, F., Molton, F., Charron, A., Dunster, C.,  
768 Oliete, A., Jacob, V., Besombes, J. L., Chevrier, F., and Jaffrezo, J. L.: Comparison between five acellular  
769 oxidative potential measurement assays performed with detailed chemistry on PM10 samples from the city of  
770 Chamonix (France), *Atmos. Chem. Phys.*, 18, 7863–7875, <https://doi.org/10.5194/acp-18-7863-2018>, 2018.

771 Calas, A., Uzu, G., Besombes, J. L., Martins, J. M. F., Redaelli, M., Weber, S., Charron, A., Albinet, A., Chevrier,  
772 F., Brulfert, G., Mesbah, B., Favez, O., and Jaffrezo, J. L.: Seasonal variations and chemical predictors of oxidative  
773 potential (OP) of particulate matter (PM), for seven urban French sites, *Atmosphere (Basel)*, 10,  
774 <https://doi.org/10.3390/atmos10110698>, 2019.

775 Chianese, E., Camastra, F., and Ciaramella, A.: Spatio-temporal learning in predicting ambient particulate matter  
776 concentration by multi-layer perceptron Spatio-temporal Learning in Predicting Ambient Particulate Matter  
777 Concentration by Multi-Layer, <https://doi.org/10.1016/j.ecoinf.2018.12.001>, 2018.

778 Cho, A., Sioutas, C., Miguel, A. H., Kumagai, Y., Schmitz, D. A., Singh, M., Eiguren-Fernandez, A., and Froines,  
779 J. R.: Redox activity of airborne particulate matter at different sites in the Los Angeles Basin, *Environ. Res.*, 99,  
780 40–47, <https://doi.org/10.1016/j.envres.2005.01.003>, 2005.

781 Cohen, J., Cohen, P., West, S. G., and Aiken, L. S.: Applied multiple regression/correlation analysis for the  
782 behavioral sciences, Routledge, 536 pp., [https://doi.org/https://doi-org.sid2nomade-](https://doi.org/https://doi-org.sid2nomade-1.grenet.fr/10.4324/9780203774441)  
783 1.grenet.fr/10.4324/9780203774441, 2002.

784 Craney, T. A. and Surles, J. G.: Model-dependent variance inflation factor cutoff values, *Qual. Eng.*, 14, 391–403,  
785 <https://doi.org/10.1081/QEN-120001878>, 2002.

786 Crobeddu, B., Aragao-Santiago, L., Bui, L. C., Boland, S., and Baeza Squiban, A.: Oxidative potential of  
787 particulate matter 2.5 as predictive indicator of cellular stress, *Environ. Pollut.*, 230, 125–133,  
788 <https://doi.org/10.1016/j.envpol.2017.06.051>, 2017.

789 Crouse, D. L., Peters, P. A., van Donkelaar, A., Goldberg, M. S., Villeneuve, P. J., Brion, O., Khan, S., Atari, D.

790 O., Jerrett, M., Pope, C. A., Brauer, M., Brook, J. R., Martin, R. V., Stieb, D., and Burnett, R. T.: Risk of  
791 nonaccidental and cardiovascular mortality in relation to long-term exposure to low concentrations of fine  
792 particulate matter: A canadian national-level cohort study, *Environ. Health Perspect.*, 120, 708–714,  
793 <https://doi.org/10.1289/ehp.1104049>, 2012.

794 Crouse, D. L., Peters, P. A., Hystad, P., Brook, J. R., van Donkelaar, A., Martin, R. V., Villeneuve, P. J., Jerrett,  
795 M., Goldberg, M. S., Arden Pope, C., Brauer, M., Brook, R. D., Robichaud, A., Menard, R., and Burnett, R. T.:  
796 Ambient PM<sub>2.5</sub>, O<sub>3</sub>, and NO<sub>2</sub> exposures and associations with mortality over 16 years of follow-up in the  
797 canadian census health and environment cohort (CanCHEC), *Environ. Health Perspect.*, 123, 1180–1186,  
798 <https://doi.org/10.1289/ehp.1409276>, 2015.

799 Daellenbach, K. R., Uzu, G., Jiang, J., Cassagnes, L.-E., Leni, Z., Vlachou, A., Stefenelli, G., Canonaco, F., Weber,  
800 S., Segers, A., and Sources, al: Sources of particulate-matter air pollution and its oxidative potential in Europe of  
801 particulate-matter air pollution and its oxidative potential in Europe, *Nature*, 587, <https://doi.org/10.1038/s41586-020-2902-8i>, 2020.

803 Deng, M., Chen, D., Zhang, G., and Cheng, H.: Policy-driven variations in oxidation potential and source  
804 apportionment of PM<sub>2.5</sub> in Wuhan, central China, *Sci. Total Environ.*, 853, 158255,  
805 <https://doi.org/10.1016/j.scitotenv.2022.158255>, 2022.

806 Dinh Ngoc, T. V.: Regression techniques applied to particulate matter oxidative potential source apportionment.  
807 In *Unveiling the optimal regression model for source apportionment of the oxidative potential of PM*,  
808 <https://doi.org/https://doi.org/10.5281/zenodo.11071884>, 2024.

809 Dominici, F.: Time-series analysis of air pollution and mortality: a statistical review., *Res. Rep. Health. Eff. Inst.*,  
810 3–27, 2004.

811 Dominutti, P. A., Borlaza, L., Sauvain, J. J., Ngoc Thuy, V. D., Houdier, S., Suarez, G., Jaffrezo, J. L., Tobin, S.,  
812 Trébuchon, C., Socquet, S., Moussu, E., Mary, G., and Uzu, G.: Source apportionment of oxidative potential  
813 depends on the choice of the assay: insights into 5 protocols comparison and implications for mitigation measures,  
814 *Environ. Sci. Atmos.*, <https://doi.org/10.1039/d3ea00007a>, 2023.

815 Elangasinghe, M. A., Singhal, N., Dirks, K. N., and Salmond, J. A.: Development of an ANN-based air pollution  
816 forecasting system with explicit knowledge through sensitivity analysis, *Atmos. Pollut. Res.*, 5, 696–708,  
817 <https://doi.org/10.5094/APR.2014.079>, 2014.

818 Fadel, M., Courcot, D., Delmaire, G., Roussel, G., Afif, C., and Ledoux, F.: Source apportionment of PM<sub>2.5</sub>  
819 oxidative potential in an East Mediterranean site, *Sci. Total Environ.*, 900,  
820 <https://doi.org/10.1016/j.scitotenv.2023.165843>, 2023.

821 Fang, T., Verma, V., T Bates, J., Abrams, J., Klein, M., Strickland, J. M., Sarnat, E. S., Chang, H. H., Mulholland,  
822 A. J., Tolbert, E. P., Russell, G. A., and Weber, J. R.: Oxidative potential of ambient water-soluble PM<sub>2.5</sub> in the  
823 southeastern United States: Contrasts in sources and health associations between ascorbic acid (AA) and  
824 dithiothreitol (DTT) assays, *Atmos. Chem. Phys.*, 16, 3865–3879, <https://doi.org/10.5194/acp-16-3865-2016>,  
825 2016.

826 Favez, O.: Traitement harmonisé de jeux de données multi-sites pour l'étude des sources de PM par Positive Matrix  
827 Factorization, 2017.

828 Godri, K. J., Harrison, R. M., Evans, T., Baker, T., Dunster, C., Mudway, I. S., and Kelly, F. J.: Increased oxidative  
829 burden associated with traffic component of ambient particulate matter at roadside and Urban background schools  
830 sites in London, *PLoS One*, 6, <https://doi.org/10.1371/journal.pone.0021961>, 2011.

831 Goldfeld, S. M. and Quandt, R. E.: Some Tests for Homoscedasticity Author ( s ): Stephen M . Goldfeld and  
832 Richard E . Quandt Source : *Journal of the American Statistical Association* , Jun ., 1965 , Vol . 60 , No . 310  
833 Published by : Taylor & Francis , Ltd . on behalf of the American Statis, *J. Am. Stat. Assoc.*, 60, 539–547, 1965.

834 Harrell: Regression Modeling Strategies, *Technometrics*, 45, 170–170, <https://doi.org/10.1198/tech.2003.s158>,  
835 2016.

836 Hastie, T. et. all.: Springer Series in Statistics The Elements of Statistical Learning, *Math. Intell.*, 27, 83–85, 2009.

837 Hawkins, D. M.: The Problem of Overfitting, *J. Chem. Inf. Comput. Sci.*, 44, 1–12,  
838 <https://doi.org/10.1021/ci0342472>, 2004.

- 839 Hernandez, P. A., Graham, C. H., Master, L. L., and Albert, D. L.: The effect of sample size and species  
840 characteristics on performance of different species distribution modeling methods, *Ecography (Cop.)*, 29, 773–  
841 785, <https://doi.org/10.1111/j.0906-7590.2006.04700.x>, 2006.
- 842 Hoerl, A. E. and Kennard, R. W.: Ridge Regression: Applications to Nonorthogonal Problems, *Technometrics*,  
843 12, 69, <https://doi.org/10.2307/1267352>, 1970.
- 844 Janssen, N. A. H., Yang, A., Strak, M., Steenhof, M., Hellack, B., Gerlofs-Nijland, M. E., Kuhlbusch, T., Kelly,  
845 F., Harrison, R., Brunekreef, B., Hoek, G., and Cassee, F.: Oxidative potential of particulate matter collected at  
846 sites with different source characteristics, *Sci. Total Environ.*, 472, 572–581,  
847 <https://doi.org/10.1016/j.scitotenv.2013.11.099>, 2014.
- 848 Kelly, F. J. and Mudway, I. S.: Protein oxidation at the air-lung interface, *Amino Acids*, 25, 375–396,  
849 <https://doi.org/10.1007/s00726-003-0024-x>, 2003.
- 850 Kuhn, M. and Johnson, K.: Applied predictive modeling, 1–600 pp., <https://doi.org/10.1007/978-1-4614-6849-3>,  
851 2013.
- 852 Leni, Z., Cassagnes, L. E., Daellenbach, K. R., Haddad, I. El, Vlachou, A., Uzu, G., Prévôt, A. S. H., Jaffrezo, J.  
853 L., Baumlin, N., Salathe, M., Baltensperger, U., Dommen, J., and Geiser, M.: Oxidative stress-induced  
854 inflammation in susceptible airways by anthropogenic aerosol, *PLoS One*, 15,  
855 <https://doi.org/10.1371/journal.pone.0233425>, 2020.
- 856 Li, Xia, T., and Nel, A. E.: The role of oxidative stress in ambient particulate matter-induced lung diseases and its  
857 implications in the toxicity of engineered nanoparticles, *Free Radic. Biol. Med.*, 44, 1689–1699,  
858 <https://doi.org/10.1016/j.freeradbiomed.2008.01.028>, 2008.
- 859 Li, J., Zhao, S., Xiao, S., Li, X., Wu, S., Zhang, J., and Schwab, J. J.: Source apportionment of water-soluble  
860 oxidative potential of PM<sub>2.5</sub> in a port city of Xiamen, Southeast China, *Atmos. Environ.*, 314, 120122,  
861 <https://doi.org/10.1016/j.atmosenv.2023.120122>, 2023.
- 862 Liu and Ng: Toxicity of Atmospheric Aerosols: Methodologies & Assays, *Am. Chem. Soc.*,  
863 <https://doi.org/DOI:10.1021/acsinfocus.7e7012>, 2023.
- 864 Liu, W. J., Xu, Y. S., Liu, W. X., Liu, Q. Y., Yu, S. Y., Liu, Y., Wang, X., and Tao, S.: Oxidative potential of  
865 ambient PM<sub>2.5</sub> in the coastal cities of the Bohai Sea, northern China: Seasonal variation and source apportionment,  
866 *Environ. Pollut.*, 236, 514–528, <https://doi.org/10.1016/j.envpol.2018.01.116>, 2018.
- 867 Lodovici, M. and Bigagli, E.: Oxidative stress and air pollution exposure, *J. Toxicol.*, 2011,  
868 <https://doi.org/10.1155/2011/487074>, 2011.
- 869 Matsuki, K., Kuperman, V., and Van Dyke, J. A.: The Random Forests statistical technique: An examination of  
870 its value for the study of reading, *Sci. Stud. Read.*, 20, 20–33, <https://doi.org/10.1080/10888438.2015.1107073>,  
871 2016.
- 872 McCullagh: Generalized linear models, <https://doi.org/10.1201/9780203738535>, 1989.
- 873 Montgomery C, D., Peck A, E., and Vining, G. G.: *Introducing To Linear Regression Analysis* (5th ed.), 2012.
- 874 Mudway, I. S., Kelly, F. J., and Holgate, S. T.: Oxidative stress in air pollution research,  
875 <https://doi.org/10.1016/j.freeradbiomed.2020.04.031>, 1 May 2020.
- 876 Nelin, T. D., Joseph, A. M., Gorr, M. W., and Wold, L. E.: Direct and indirect effects of particulate matter on the  
877 cardiovascular system, *Toxicol. Lett.*, 208, 293–299, <https://doi.org/10.1016/j.toxlet.2011.11.008>, 2012.
- 878 O'Brien, R. M.: A caution regarding rules of thumb for variance inflation factors, *Qual. Quant.*, 41, 673–690,  
879 <https://doi.org/10.1007/s11135-006-9018-6>, 2007.
- 880 Paatero, P. and Hopke, P. K.: Rotational tools for factor analytic models, *J. Chemom.*, 23, 91–100,  
881 <https://doi.org/10.1002/cem.1197>, 2009.
- 882 Paatero, P. and Tappert, U.: Positive matrix factorization: A non-negative factor model with optimal utilization of  
883 error estimates of data values, *Environmetrics*, 111–126 pp.,  
884 <https://doi.org/https://doi.org/10.1002/env.3170050203>, 1994.
- 885 Pearce, J. and Ferrier, S.: An evaluation of alternative algorithms for fitting species distribution models using

886 logistic regression, *Ecol. Modell.*, 128, 127–147, [https://doi.org/10.1016/S0304-3800\(99\)00227-6](https://doi.org/10.1016/S0304-3800(99)00227-6), 2000.

887 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P.,  
888 Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E.:  
889 Scikit-learn: Machine Learning in {P}ython, *J. Mach. Learn. Res.*, 12, 2825–2830, 2011.

890 Pelucchi, C., Negri, E., Gallus, S., Boffetta, P., Tramacere, I., and La Vecchia, C.: Long-term particulate matter  
891 exposure and mortality: A review of European epidemiological studies, *BMC Public Health*, 9, 1–8,  
892 <https://doi.org/10.1186/1471-2458-9-453>, 2009.

893 Peng, R. D., Bell, M. L., Geyh, A. S., McDermott, A., Zeger, S. L., Samet, J. M., and Dominici, F.: Emergency  
894 admissions for cardiovascular and respiratory diseases and the chemical composition of fine particle air pollution,  
895 *Environ. Health Perspect.*, 117, 957–963, <https://doi.org/10.1289/ehp.0800185>, 2009.

896 Pietrogrande, M. C., Romanato, L., and Russo, M.: Synergistic and Antagonistic Effects of Aerosol Components  
897 on Its Oxidative Potential as Predictor of Particle Toxicity, *Toxics*, 10, <https://doi.org/10.3390/toxics10040196>,  
898 2022.

899 Pope, C. A. and Dockery, D. W.: Health effects of fine particulate air pollution: Lines that connect, *J. Air Waste*  
900 *Manag. Assoc.*, 56, 709–742, <https://doi.org/10.1080/10473289.2006.10464485>, 2006.

901 Rao, X., Zhong, J., Brook, R. D., and Rajagopalan, S.: Effect of Particulate Matter Air Pollution on Cardiovascular  
902 Oxidative Stress Pathways, *Antioxidants Redox Signal.*, 28, 797–818, <https://doi.org/10.1089/ars.2017.7394>,  
903 2018.

904 Raudys, S. J. and Jain, A. K.: Small Sample Size Effects in Statistical Pattern Recognition: Recommendations for  
905 Practitioners, <https://doi.org/10.1109/34.75512>, 1991.

906 Rosenblad, A.: The Concise Encyclopedia of Statistics, 867–868 pp.,  
907 <https://doi.org/10.1080/02664760903075614>, 2011.

908 Samake, A., Uzu, G., Martins, J. M. F., Calas, A., Vince, E., Parat, S., and Jaffrezo, J. L.: The unexpected role of  
909 bioaerosols in the Oxidative Potential of PM, *Sci. Rep.*, 7, <https://doi.org/10.1038/s41598-017-11178-0>, 2017.

910 Seabold, S. and Perktold, J.: Statsmodels: Econometric and statistical modeling with python, in: 9th Python in  
911 Science Conference, 2010.

912 Shangguan, Y., Zhuang, X., Querol, X., Li, B., Moreno, N., Trechera, P., Sola, P. C., Uzu, G., and Li, J.:  
913 Characterization of deposited dust and its respirable fractions in underground coal mines: Implications for  
914 oxidative potential-driving species and source apportionment, *Int. J. Coal Geol.*, 258,  
915 <https://doi.org/10.1016/j.coal.2022.104017>, 2022.

916 Stevanović, S., Jovanović, M. V., Jovašević-Stojanović, M. V., and Ristovski, Z.: SOURCE APPORTIONMENT  
917 OF OXIDATIVE POTENTIAL What We Know So Far, *Therm. Sci.*, 27, 2347–2357,  
918 <https://doi.org/10.2298/TSCI221107111S>, 2023.

919 Stockwell, D. R. B. and Peterson, A. T.: Effects of sample size on accuracy of species distribution models, *Ecol.*  
920 *Modell.*, 148, 1–13, [https://doi.org/10.1016/S0304-3800\(01\)00388-X](https://doi.org/10.1016/S0304-3800(01)00388-X), 2002.

921 Szigeti, T., Óvári, M., Dunster, C., Kelly, F. J., Lucarelli, F., and Záráy, G.: Changes in chemical composition and  
922 oxidative potential of urban PM<sub>2.5</sub> between 2010 and 2013 in Hungary, *Sci. Total Environ.*, 518–519, 534–544,  
923 <https://doi.org/10.1016/j.scitotenv.2015.03.025>, 2015.

924 Szigeti, T., Dunster, C., Cattaneo, A., Cavallo, D., Spinazzè, A., Saraga, D. E., Sakellaris, I. A., de Kluizenaar, Y.,  
925 Cornelissen, E. J. M., Hänninen, O., Peltonen, M., Calzolari, G., Lucarelli, F., Mandin, C., Bartzis, J. G., Záráy, G.,  
926 and Kelly, F. J.: Oxidative potential and chemical composition of PM<sub>2.5</sub> in office buildings across Europe - The  
927 OFFICAIR study, *Environ. Int.*, 92–93, 324–333, <https://doi.org/10.1016/j.envint.2016.04.015>, 2016.

928 Tibshirani, R.: Regression Shrinkage and Selection Via the Lasso, *J. R. Stat. Soc. Ser. B*, 58, 267–288,  
929 <https://doi.org/10.1111/j.2517-6161.1996.tb02080.x>, 1996.

930 in 't Veld, M., Pandolfi, M., Amato, F., Pérez, N., Reche, C., Dominutti, P., Jaffrezo, J., Alastuey, A., Querol, X.,  
931 and Uzu, G.: Discovering oxidative potential (OP) drivers of atmospheric PM<sub>10</sub>, PM<sub>2.5</sub>, and PM<sub>1</sub> simultaneously  
932 in North-Eastern Spain, *Sci. Total Environ.*, 857, <https://doi.org/10.1016/j.scitotenv.2022.159386>, 2023.

933 Verma, V., Fang, T., Guo, H., King, L., Bates, J. T., Peltier, R. E., Edgerton, E., Russell, A. G., and Weber, R. J.:

934 Reactive oxygen species associated with water-soluble PM<sub>2.5</sub> in the southeastern United States: Spatiotemporal  
935 trends and source apportionment, *Atmos. Chem. Phys.*, 14, 12915–12930, [https://doi.org/10.5194/acp-14-12915-](https://doi.org/10.5194/acp-14-12915-2014)  
936 2014, 2014.

937 Viana, M., Kuhlbusch, T. A. J., Querol, X., Alastuey, A., Harrison, R. M., Hopke, P. K., Winiwarter, W., Vallius,  
938 M., Szidat, S., Prévôt, A. S. H., Hueglin, C., Bloemen, H., Wählin, P., Vecchi, R., Miranda, A. I., Kasper-Giebl,  
939 A., Maenhaut, W., and Hitzenberger, R.: Source apportionment of particulate matter in Europe: A review of  
940 methods and results, <https://doi.org/10.1016/j.jaerosci.2008.05.007>, 2008.

941 Vida, M., Foret, G., Siour, G., Coman, A., Weber, S., Favez, O., Jaffrezo, J., Pontet, S., Mesbah, B., Gille, G.,  
942 Zhang, S., Chevrier, F., Pallares, C., Uzu, G., and Beekmann, M.: Oxidative potential modelling of PM<sub>10</sub>: a 2-  
943 year study over France, *ACDP*, 2024.

944 Wang, D., Yang, X., Lu, H., Li, D., Xu, H., Luo, Y., Sun, J., Hang Ho, S. S., and Shen, Z.: Oxidative potential of  
945 atmospheric brown carbon in six Chinese megacities: Seasonal variation and source apportionment, *Atmos.*  
946 *Environ.*, 309, 119909, <https://doi.org/10.1016/j.atmosenv.2023.119909>, 2023.

947 Wang, J., Jiang, H., Jiang, H., Mo, Y., Geng, X., Li, J., Mao, S., Bualert, S., Ma, S., Li, J., and Zhang, G.: Source  
948 apportionment of water-soluble oxidative potential in ambient total suspended particulate from Bangkok: Biomass  
949 burning versus fossil fuel combustion, *Atmos. Environ.*, 235, 117624,  
950 <https://doi.org/10.1016/j.atmosenv.2020.117624>, 2020a.

951 Wang, S., Ye, J., Soong, R., Wu, B., Yu, L., Simpson, A. J., and Chan, A. W. H.: Relationship between chemical  
952 composition and oxidative potential of secondary organic aerosol from polycyclic aromatic hydrocarbons, *Atmos.*  
953 *Chem. Phys.*, 18, 3987–4003, <https://doi.org/10.5194/acp-18-3987-2018>, 2018.

954 Wang, Y., Wang, M., Li, S., Sun, H., Mu, Z., Zhang, L., Li, Y., and Chen, Q.: Study on the oxidation potential of  
955 the water-soluble components of ambient PM<sub>2.5</sub> over Xi'an, China: Pollution levels, source apportionment and  
956 transport pathways, *Environ. Int.*, 136, 105515, <https://doi.org/10.1016/j.envint.2020.105515>, 2020b.

957 Weber, S., Uzu, G., Calas, A., Chevrier, F., Besombes, J. L., Charron, A., Salameh, D., Ježek, I., Močnik, G., and  
958 Jaffrezo, J. L.: An apportionment method for the oxidative potential of atmospheric particulate matter sources:  
959 Application to a one-year study in Chamonix, France, *Atmos. Chem. Phys.*, 18, 9617–9629,  
960 <https://doi.org/10.5194/acp-18-9617-2018>, 2018.

961 Weber, S., Salameh, D., Albinet, A., Alleman, L. Y., Waked, A., Besombes, J. L., Jacob, V., Guillaud, G.,  
962 Meshbah, B., Rocq, B., Hulin, A., Dominik-Sègue, M., Chrétien, E., Jaffrezo, J. L., and Favez, O.: Comparison of  
963 PM<sub>10</sub> sources profiles at 15 french sites using a harmonized constrained positive matrix factorization approach,  
964 *Atmosphere (Basel)*, 10, <https://doi.org/10.3390/atmos10060310>, 2019.

965 Weber, S., Uzu, G., Favez, O., Borlaza, L., Calas, A., Salameh, D., Chevrier, F., Allard, J., Besombes, J. L.,  
966 Albinet, A., Pontet, S., Mesbah, B., Gille, G., Zhang, S., Pallares, C., Leoz-Garziandia, E., and Jaffrezo, J. L.:  
967 Source apportionment of atmospheric PM<sub>10</sub> oxidative potential: Synthesis of 15 year-round urban datasets in  
968 France, *Atmos. Chem. Phys.*, 21, 11353–11378, <https://doi.org/10.5194/acp-21-11353-2021>, 2021.

969 WHO: WHO global air quality guidelines, 2021.

970 Williams, M., Gomez Grajales, C. A., and Kurkiewicz, D.: Assumptions of Multiple Regression: Correcting Two  
971 Misconceptions - Practical Assessment, Research & Evaluation, *Pract. Assessment, Res. Eval.*, 18, 1–16, 2013.

972 Wisz, M. S., Hijmans, R. J., Li, J., Peterson, A. T., Graham, C. H., Guisan, A., Elith, J., Dudík, M., Ferrier, S.,  
973 Huettmann, F., Leathwick, J. R., Lehmann, A., Lohmann, L., Loiselle, B. A., Manion, G., Moritz, C., Nakamura,  
974 M., Nakazawa, Y., Overton, J. M. C., Phillips, S. J., Richardson, K. S., Scachetti-Pereira, R., Schapire, R. E.,  
975 Soberón, J., Williams, S. E., and Zimmermann, N. E.: Effects of sample size on the performance of species  
976 distribution models, *Divers. Distrib.*, 14, 763–773, <https://doi.org/10.1111/j.1472-4642.2008.00482.x>, 2008.

977 Xiong, Q., Yu, H., Wang, R., Wei, J., and Verma, V.: Rethinking Dithiothreitol-Based Particulate Matter Oxidative  
978 Potential: Measuring Dithiothreitol Consumption versus Reactive Oxygen Species Generation, *Environ. Sci.*  
979 *Technol.*, 51, 6507–6514, <https://doi.org/10.1021/acs.est.7b01272>, 2017.

980 Yang, A., Jedynska, A., Hellack, B., Kooter, I., Hoek, G., Brunekreef, B., Kuhlbusch, T. A. J., Cassee, F. R., and  
981 Janssen, N. A. H.: Measurement of the oxidative potential of PM<sub>2.5</sub> and its constituents: The effect of extraction  
982 solvent and filter type, *Atmos. Environ.*, 83, 35–42, <https://doi.org/10.1016/j.atmosenv.2013.10.049>, 2014.

983 Yu, Guo, S., Xu, R., Ye, T., Li, S., Sim, M. R., Abramson, M. J., and Guo, Y.: Cohort studies of long-term exposure

- 984 to outdoor particulate matter and risks of cancer: A systematic review and meta-analysis, *Innovation*, 2, 100143,  
985 <https://doi.org/10.1016/j.xinn.2021.100143>, 2021.
- 986 Yu, S. Y., Liu, W. J., Xu, Y. S., Yi, K., Zhou, M., Tao, S., and Liu, W. X.: Characteristics and oxidative potential  
987 of atmospheric PM<sub>2.5</sub> in Beijing: Source apportionment and seasonal variation, *Sci. Total Environ.*, 650, 277–  
988 287, <https://doi.org/10.1016/j.scitotenv.2018.09.021>, 2019.
- 989 Zhang, Y., Albinet, A., Petit, J. E., Jacob, V., Chevrier, F., Gille, G., Pontet, S., Chrétien, E., Dominik-Sègue, M.,  
990 Levigoureux, G., Močnik, G., Gros, V., Jaffrezo, J. L., and Favez, O.: Substantial brown carbon emissions from  
991 wintertime residential wood burning over France, *Sci. Total Environ.*, 743,  
992 <https://doi.org/10.1016/j.scitotenv.2020.140752>, 2020.

993