

Unveiling the optimal regression model for source apportionment of the oxidative potential of PM₁₀

Vy N.T. Dinh¹, Jean-Luc Jaffrezo¹, Ian Hough¹, Pamela A. Dominutti¹, Guillaume Salque Moreton², Grégory Gille³, Florie Francony⁴, Arabelle Patron-Anquez⁵, Olivier Favez^{6,7}, Gaëlle Uzu¹

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

² Atmo AuRA, 69500 Bron, France

³Atmo Sud, 13006 Marseille, France

⁴Atmo Nouvelle Aquitaine, 33692 Merignac, France

⁵Atmo Hauts de France, 59044 Lille, France

⁶INERIS, Parc Technologique Alata, BP 2, 60550 Verneuil-en-Halatte, France

⁷ 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 PM₁₀ sources in different types of environments. PM₁₀ sources derived from receptor models (e.g. EPA PMF) are coupled with regression models expressing the associations between PM₁₀ sources and PM₁₀ OP measured by ascorbic acid (OP_{AA}) and dithiothreitol assay (OP_{DTT}). 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 PM₁₀ samples and chemical analyses at each of six sites of different typologies in France to assess the possible impact of PM source variability on PM₁₀ OP apportionment. PM₁₀ source-specific OP_{DTT} and OP_{AA} 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; Peng et al., 2009; Pope & Dockery, 2006) and (ii) cohort studies aiming to evaluate the long-term effects of chronic PM exposure (Ayres et al., 2008; Beelen et al., 2014; Crouse et al., 2012, 2015; Pelucchi et al., 2009; 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

41 stress and cellular damage (Akhtar et al., 2010; Ayres et al., 2008; Leni et al., 2020; Li et al., 2008; Lodovici &
42 Bigagli, 2011; Mudway et al., 2020; Nelin et al., 2012; Rao et al., 2018). The quantification of the PM capacity to
43 oxidize a biological media is called oxidative potential (OP) (Bates et al., 2019; Daellenbach et al., 2020; Dominutti
44 et al., 2023). Various acellular assays of OP have been introduced, differentiating ROS generation mechanisms of
45 PM (Calas et al., 2018; Dominutti et al., 2023). Dithiothreitol (DTT) and ascorbic acid (AA) assays are two of the
46 commonly used ones in the literature (Liu & 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., 2018, 2019; Crobeddu et al., 2017; Godri et al., 2011; Janssen et al., 2014;
49 Szigeti et al., 2015, 2016; Yang et al., 2014). However, this research pathway struggles with the co-variation
50 between measured and unmeasured PM components (Calas et al., 2018; Weber et al., 2018). An alternative
51 approach is to examine the association between OP and sources of PM obtained using receptor models such as
52 chemical mass balance, positive matrix factorization (PMF), or principal components analysis. PMF is the most
53 popular method for its ability to quantify PM source contributions without extensive prior information on specific
54 sources at the site studied (Belis et al., 2013; Brown et al., 2015; Paatero & Hopke, 2009; Paatero & Tappert, 1994;
55 Viana et al., 2008).

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; J. Wang et al., 2020; 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; Y. Wang et al., 2020; 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 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 PM_{10} OP_{AA} and
78 PM_{10} OP_{DTT} to PM_{10} sources at six sites in France. The 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 PM_{10} OP SA models are compared with regard to the estimated
81 intrinsic 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 PM_{10} .
85 This holds particular significance in the context of the implementation of OP monitoring as a novel air quality

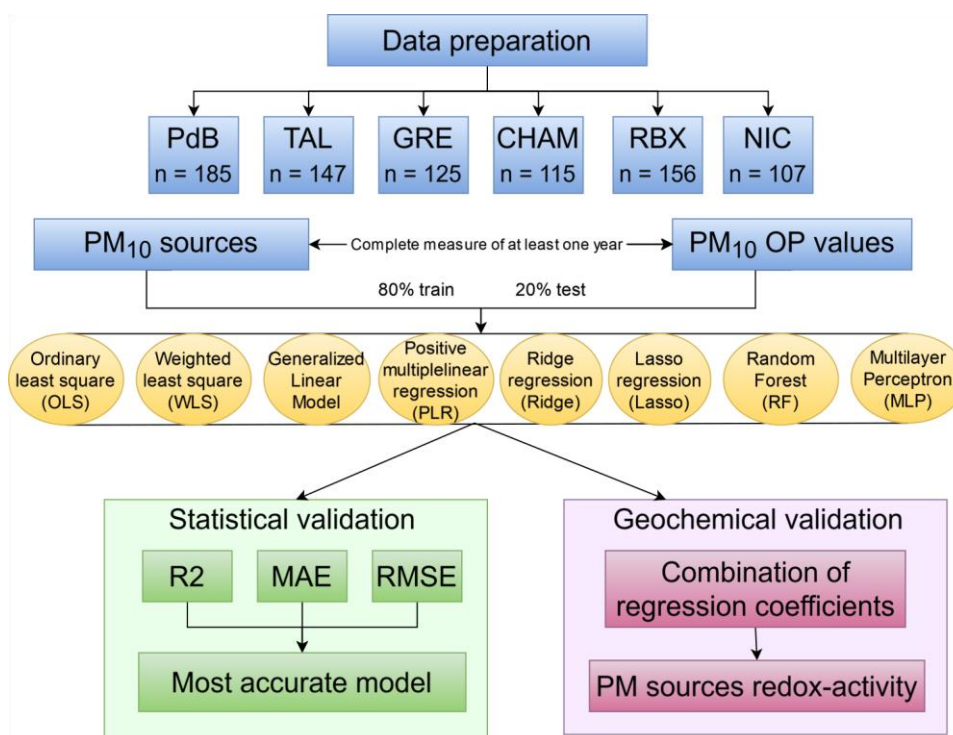
86 metric as foreseen in research programs (such RI-Urbans) and in the process of the revision of the European
87 Directive 2008/50/CE.

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 PM₁₀ sources and PM₁₀ OP, identify collinearity among PM₁₀ sources, and
92 examine homoscedasticity in the relationship between PM₁₀ OP and PM₁₀ sources. Section 2.5 describes the eight
93 regression techniques (OLS, WLS, PLS, Ridge, Lasso, GLM, RF, and MLP), used for PM₁₀ OP SA. Each
94 technique is applied to each site separately using PM₁₀ OP_v (nmol min⁻¹ m⁻³) as the dependent variable and PM₁₀
95 sources (µg m⁻³) as independent variables. The coefficient of the regression called the intrinsic PM₁₀ OP of the
96 source (nmol min⁻¹ µg⁻¹), represents the capacity of each µg of PM₁₀ from the given source to generate oxidative
97 stress; the higher the intrinsic PM₁₀ 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, PM₁₀ OP analytical errors were used as a weighting, implying that the PM₁₀
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 PM₁₀ OP and PM₁₀ 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 (R²). The geochemical validation is based on the regression coefficient (the
106 intrinsic PM₁₀ OP) of each source. These are calculated separately for the training and testing data and averaged
107 across the 500 sampling iterations.

108



109

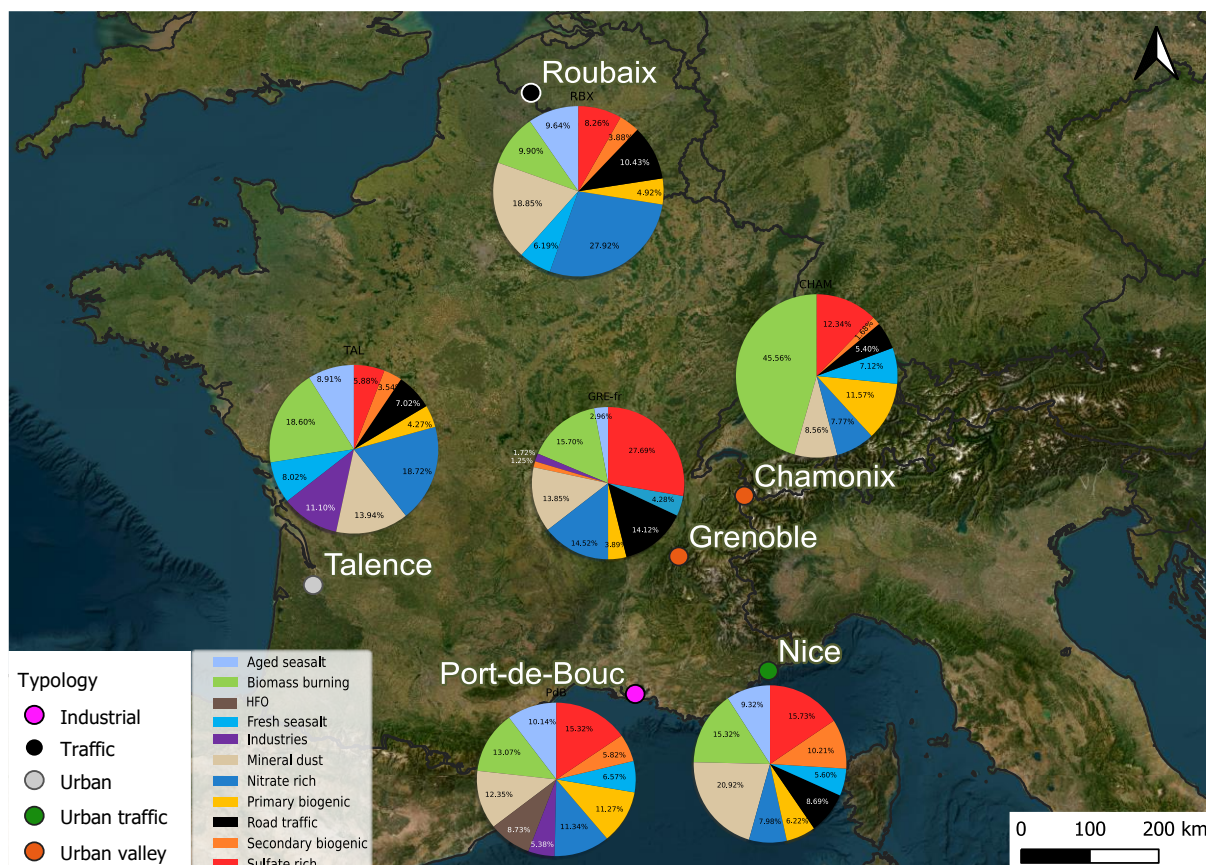
110 **Figure 1. Workflow of the comparison of PM₁₀ OP sources apportionment methodology**

111 2.2. Study sites and PM₁₀ 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 PM₁₀ 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 PM₁₀ 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⁻, NO₃⁻,
122 SO₄²⁻, NH₄⁺, Na⁺, K⁺, Mg²⁺, Ca²⁺) 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 PM₁₀ 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 PM₁₀ 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 PM₁₀ sources at each site, while Figure 2
133 presents the localisation of the sites in France, together with the respective proportion of each PM₁₀ 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 PM₁₀ source apportionment for each site with the colors identifying the PM₁₀**
 138 **sources. Background photography from ESRI satellite.**

139 Table 1. Data description

	PdB	TAL	GRE-fr	CHAM	RBX	NIC
Name	Port de Bouc	Talence	Grenoble	Chamonix	Roubaix	Nice
N of samples	185	147	125	115	156	107
Sampling dates	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
N of sources	10	10	10	8	9	9

140
 141 **2.3. OP analysis**
 142 PM₁₀ OP assays were performed on PM₁₀ 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 PM₁₀ and biological oxidants (Calas et al., 2018; Cho et
148 al., 2005). 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 PM₁₀ OP in this study represents
153 the mean and the analysis uncertainty is the standard deviation of these four PM₁₀ OP analyses. After analysis, the
154 PM₁₀ 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_v represents the PM₁₀ OP due to PM₁₀ per cubic
156 meter of air (nmol min⁻¹ m⁻³). To simplify the denotation of PM₁₀ OP, OP is used to represent PM₁₀ 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 PM₁₀ 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 PM₁₀ sources, R^2 is the coefficient of determination of a multiple linear
171 regression model between the i^{th} source and the other sources. VIF values of a PM₁₀ source present a range between
172 1, and ∞ . The higher the VIF values, the greater the collinearity between this PM₁₀ 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 & 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 PM₁₀ sources to identify the residual of OP prediction. Next, the PM₁₀ sources and residual corresponding are
181 divided into three segments: the upper segment is the group with higher PM₁₀ sources concentration, the lower
182 segment is the group with lower PM₁₀ sources concentration, and the middle segment, constituting 10% of the
183 moderate PM₁₀ 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 & Perktold, 2010).

189 2.5. Regression models

190 The fundamental principle of regression models in this study is to use the PM₁₀ 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⁻¹ m⁻³), X are the PM₁₀ source contributions (μg m⁻³), y is the observed OP
197 (nmol min⁻¹ m⁻³), and e denotes the residuals (nmol min⁻¹ m⁻³). 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 PM₁₀ sources is linear. (2) **Independence**: The
202 PM₁₀ sources must be independent, with no collinearity. (3) **Homoscedasticity**: The variance of residuals is
203 constant across all values of PM₁₀ sources. (4) **Normality**: The residuals are normally distributed. In the OLS
204 model, the estimated equation and objective to minimize are defined as follows:

$$205 \hat{y} = \beta_0 + \sum_1^p \beta_i * x_i \text{ (Eq4)}$$

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

207 Here, the β_0 denotes the intercept (nmol min⁻¹ m⁻³), β_i represents the regression coefficient (intrinsic OP, nmol
208 min⁻¹ μg⁻¹) of source i , x_i is the concentration of source i (μg m⁻³), p is the number of PM₁₀ 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 PM₁₀ 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 λ is the parameter representing the amount of
 234 shrinkage, the larger λ , the greater the shrinkage. The hyperparameter tuning was implemented with different
 235 values of λ (5, 1, 0.5, 0.1, 0.01, 0.005, 0.001, 0.0005, 0.0001). The best λ 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, λ is the parameter representing the amount of shrinkage. λ 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 PM₁₀
 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 PM₁₀ 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 β_0 denotes the intercept, β_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 PM₁₀ 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 PM₁₀ source in OP prediction. This study estimated the permutation importance of
272 each PM₁₀ source as the mean increase in the mean squared error of predicted OP when the values of the PM₁₀
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; Boursard & Wellekens, 1989; Chianese et al., 2018). These layers
277 include an input layer (PM₁₀ sources), one or several hidden layers, and an output layer (OP_{AA} or OP_{DTT} 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 PM₁₀ sources.

290 All regression models were performed using the Python package statsmodels 0.14.0 (Seabold & 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

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

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

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

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

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 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 PM_{10} source contributions by sites. Most
313 sites, including traffic and industrial ones, show higher OP_{DTT} activities than OP_{AA} . Conversely, for the alpine
314 valley sites, CHAM presents higher OP_{AA} than OP_{DTT} , while GRE-fr experiences similar levels of OP_{AA} and
315 OP_{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 PM_{10} concentrations but higher OP activities than other sites, while TAL
317 showed high PM_{10} concentrations but relatively lower OP activities.

318 The variations observed in the levels of 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 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

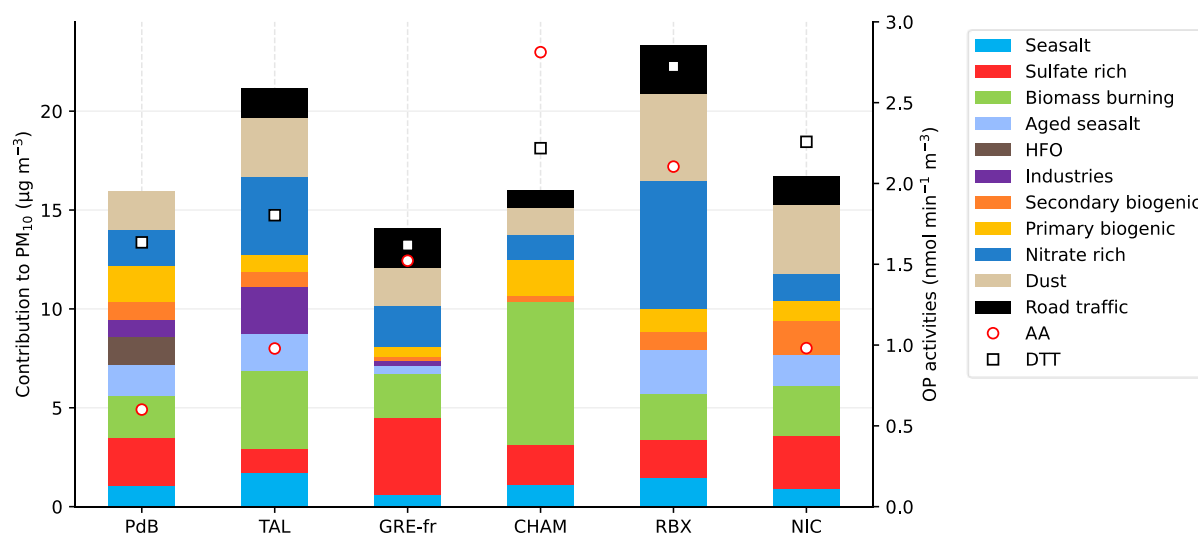
338 other sites, VIF values are below 5, indicating limited collinearity among sources. This is expected, as the PMF
 339 analysis is constrained to avoid collinearity between sources. VIF values for each site can be found in Table S.2.

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 (PM₁₀ 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_{AA}
 345 and in CHAM and TAL for OP_{DTT} (Table 2). We observed a large difference between the cold and warm periods
 346 for both OP_{AA} and OP_{DTT} in CHAM, similar to what was seen for OP_{AA} 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 PM₁₀ 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



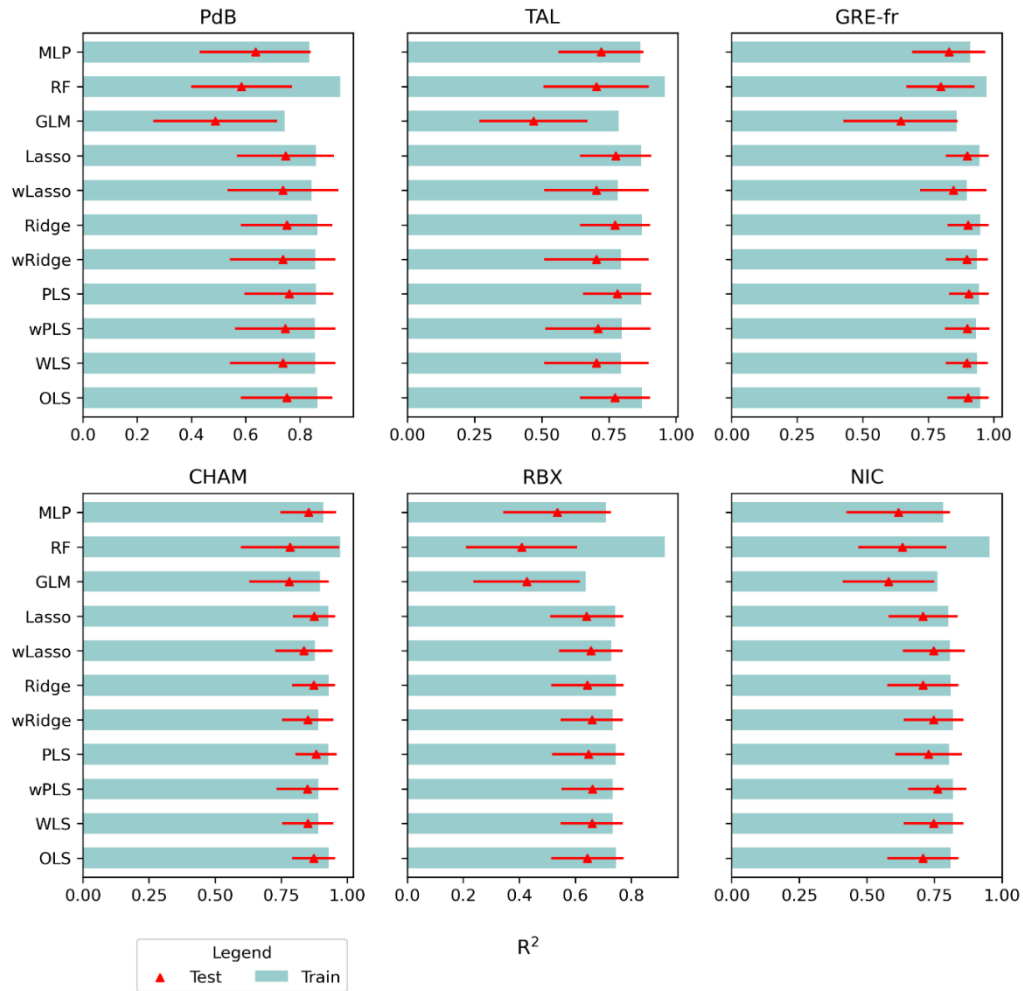
355

356 **Figure 3. The contribution of sources to PM₁₀ 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_v**
 358 **activities in $\text{nmol min}^{-1} \text{m}^{-3}$, with red circle for OP_{AA} and black square for OP_{DTT}.**

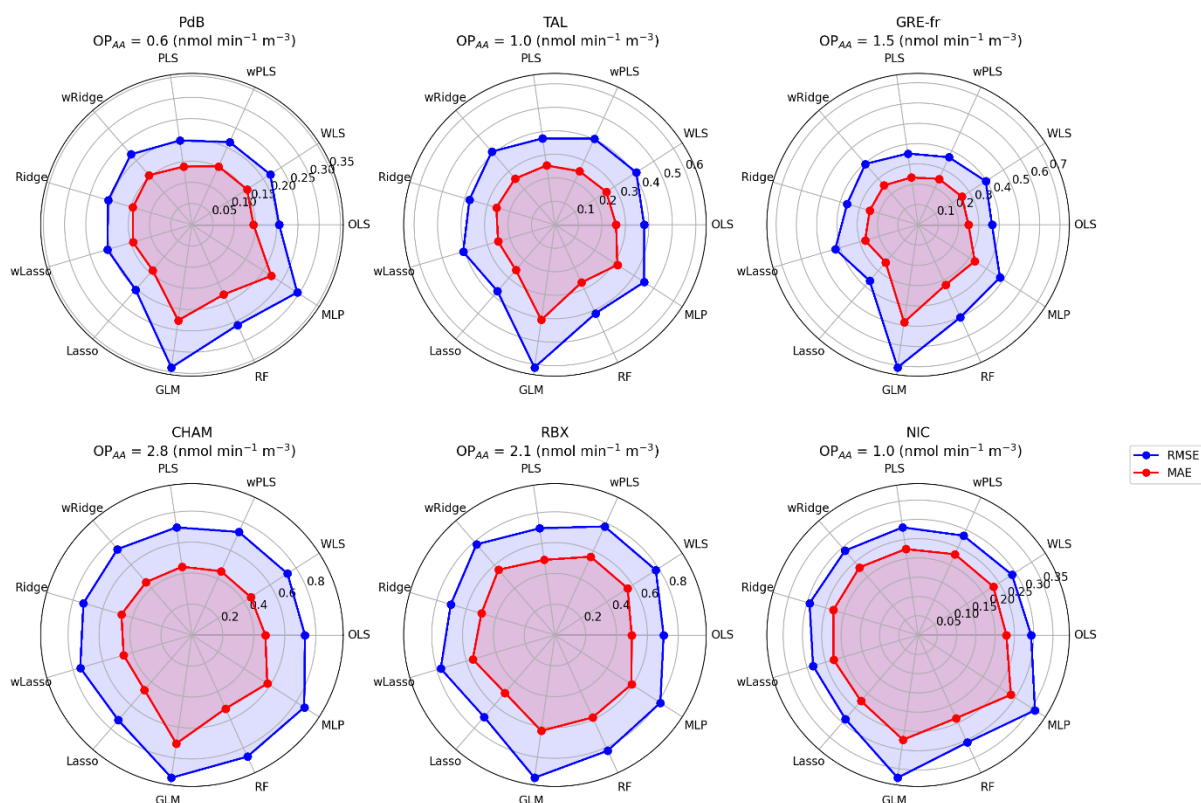
359 3.2. The performances of regression models

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 (Benkendorf
 391 & Hawkins, 2020; Hawkins, 2004; Hernandez et al., 2006; Matsuki et al., 2016; Raudys & Jain, 1991; Stockwell
 392 & Peterson, 2002; Wisz et al., 2008). Pearce and Ferrier (2000) recommended that the minimum number of
 393 samples for robust performance should be over 250 for GLM model, while (Raudys & Jain, 1991) showed that the
 394 minimum number of sample are based on the complexity of the model and the number of predictors. Additionally,
 395 Harrell (2016) suggested that the number of predictors (PM sources) should be below the number of samples
 396 divided by 15, a threshold not reached in this analysis. For example, in NIC, the minimum number of samples
 397 should be 135 for the training set (9 PM sources x 15), while in total, we have only 107 samples. Therefore, we

398 can also recommend that, for optimal performance of RF, and MLP, the number of samples and PM sources should
 399 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 (Cohen et al., 2002;
 407 Williams et al., 2013).

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 heteroscedasticity, 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}**

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

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

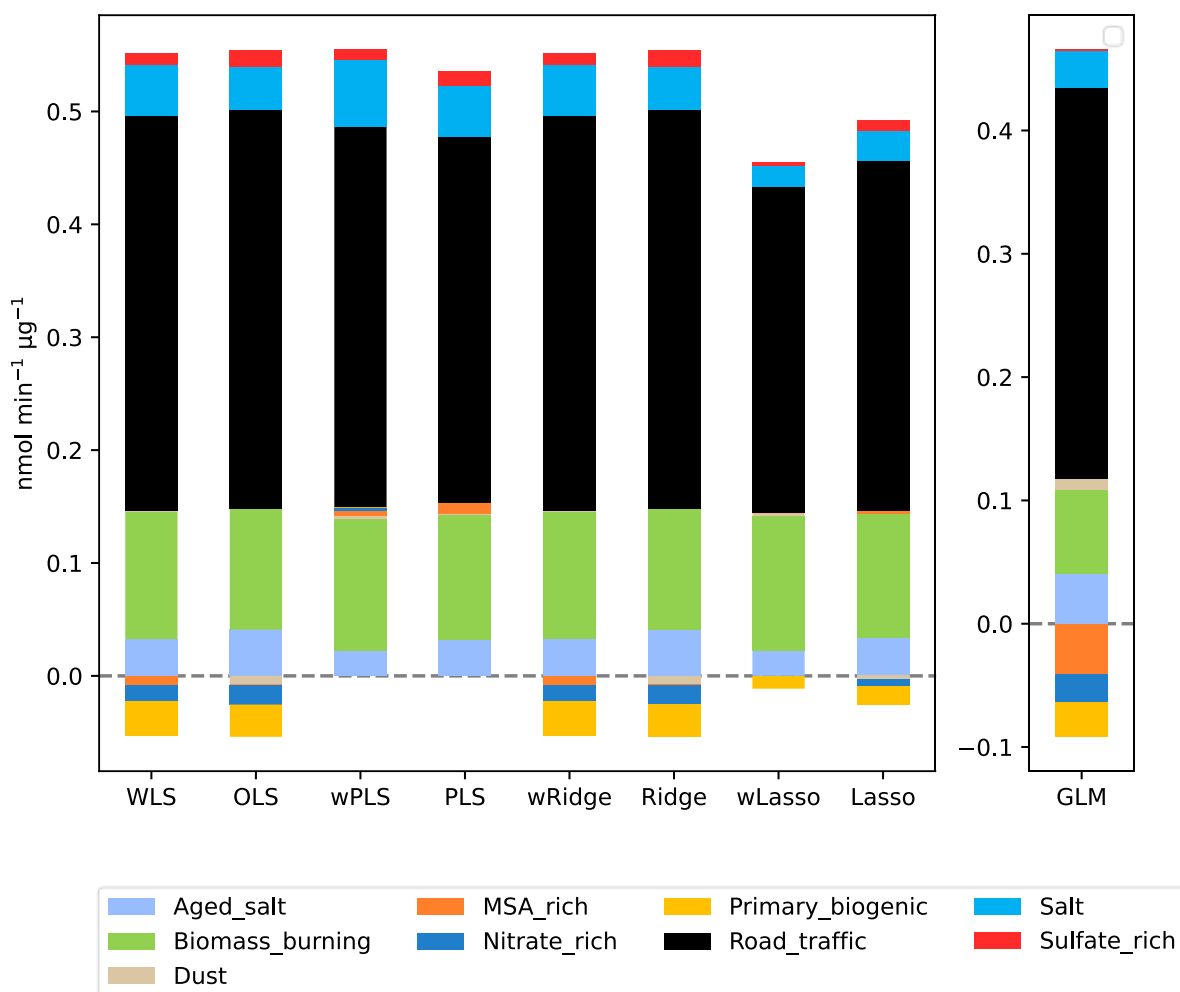
	PdB	TAL	GRE-fr	CHAM	RBX	NIC
Collinearity	No	No	No	No	No	Yes
Heteroscedasticity	No	Yes	No	Yes	No	No
Best model by characteristic	OLS/ PLS	WLS/ wPLS	OLS/ PLS	WLS/ wPLS	OLS/ PLS	Ridge/ Lasso
Best by error	OLS	wPLS	PLS	wPLS	PLS	Ridge

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 μg 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.

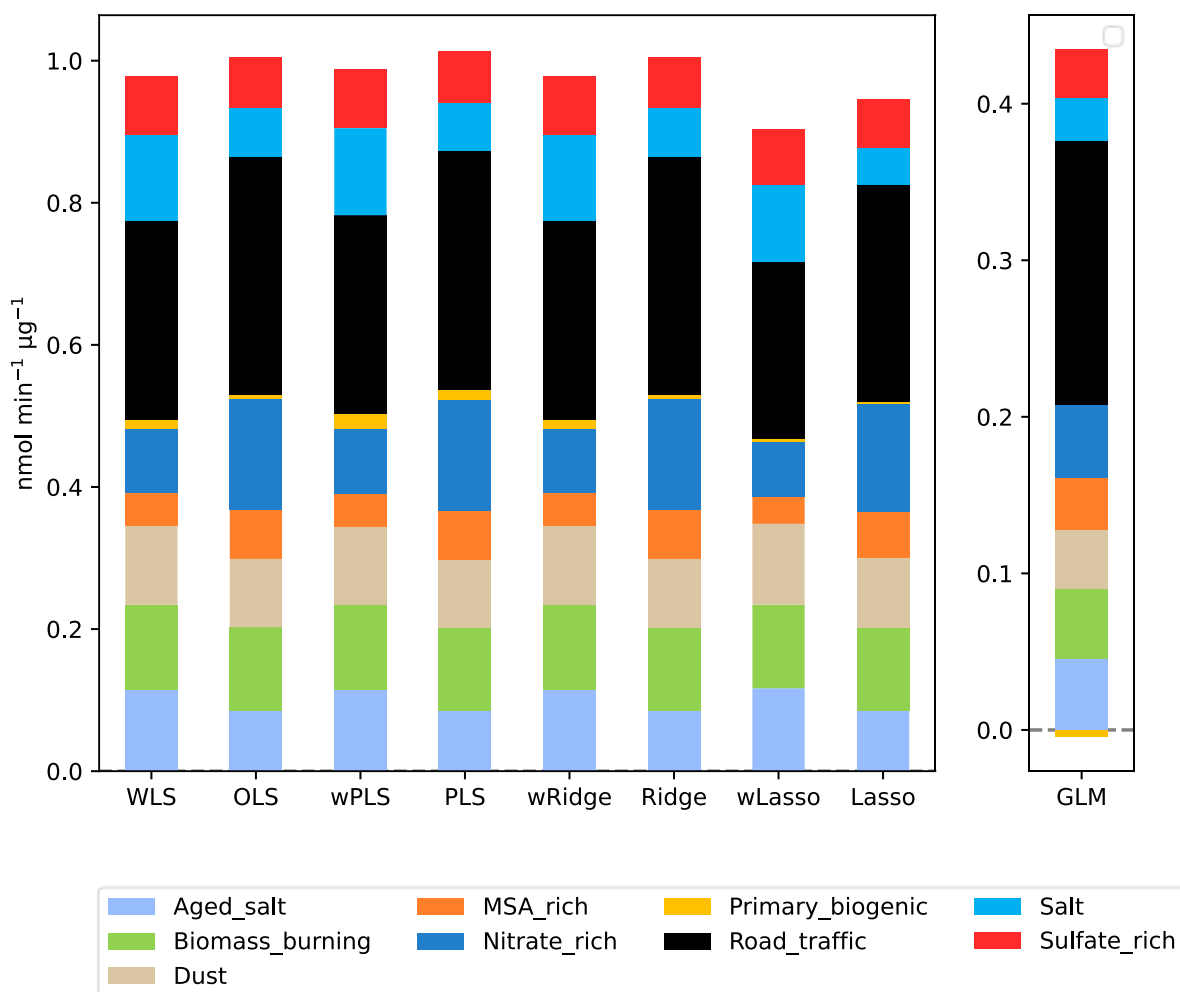
433 Intrinsic OP_{AA} of PM_{10} sources at NIC is the same between WLS and wRidge and between the OLS and Ridge,
434 revealing that the moderate collinearity of the road traffic source did not affect the estimated intrinsic OP_{AA} . PLS
435 sets the intrinsic OP_{AA} of some sources to zero, therefore 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_{AA} values of the different PM₁₀ sources at Nice were obtained with the different models.**

448 The OP_{DTT} 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_{AA} results, the wLasso
 451 and GLM models exhibit distinct responses compared to the other models. The intrinsic OP_{DTT} 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_{DTT} 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_{AA}, OP_{DTT} prediction shows more variation among models, highlighting the effect of choosing
 456 a model on evaluating the intrinsic OP_{DTT} of PM₁₀ sources.



457

458 **Figure 7. The variations of the intrinsic OP_{DTT} of the different PM₁₀ sources at Nice were obtained with the**
 459 **different models.**

460 The comparison of intrinsic OP among regression models in NIC demonstrated that OP_{DTT} and OP_{AA} 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 PM₁₀ sources (A similar pattern is observed
 463 for all other sites and shown in Fig S.3 to Fig S.7 for OP_{AA} and Fig S.8 to S.12 for OP_{DTT}). 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 PM₁₀ 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_{AA} and Table 4 for OP_{DTT}) 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 PM₁₀ 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 PM₁₀ 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_{DTT} and OP_{AA}, 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_{AA} of the best and reference model is presented in 3.4.1 and that of OP_{DTT} is shown
499 in 3.4.2.

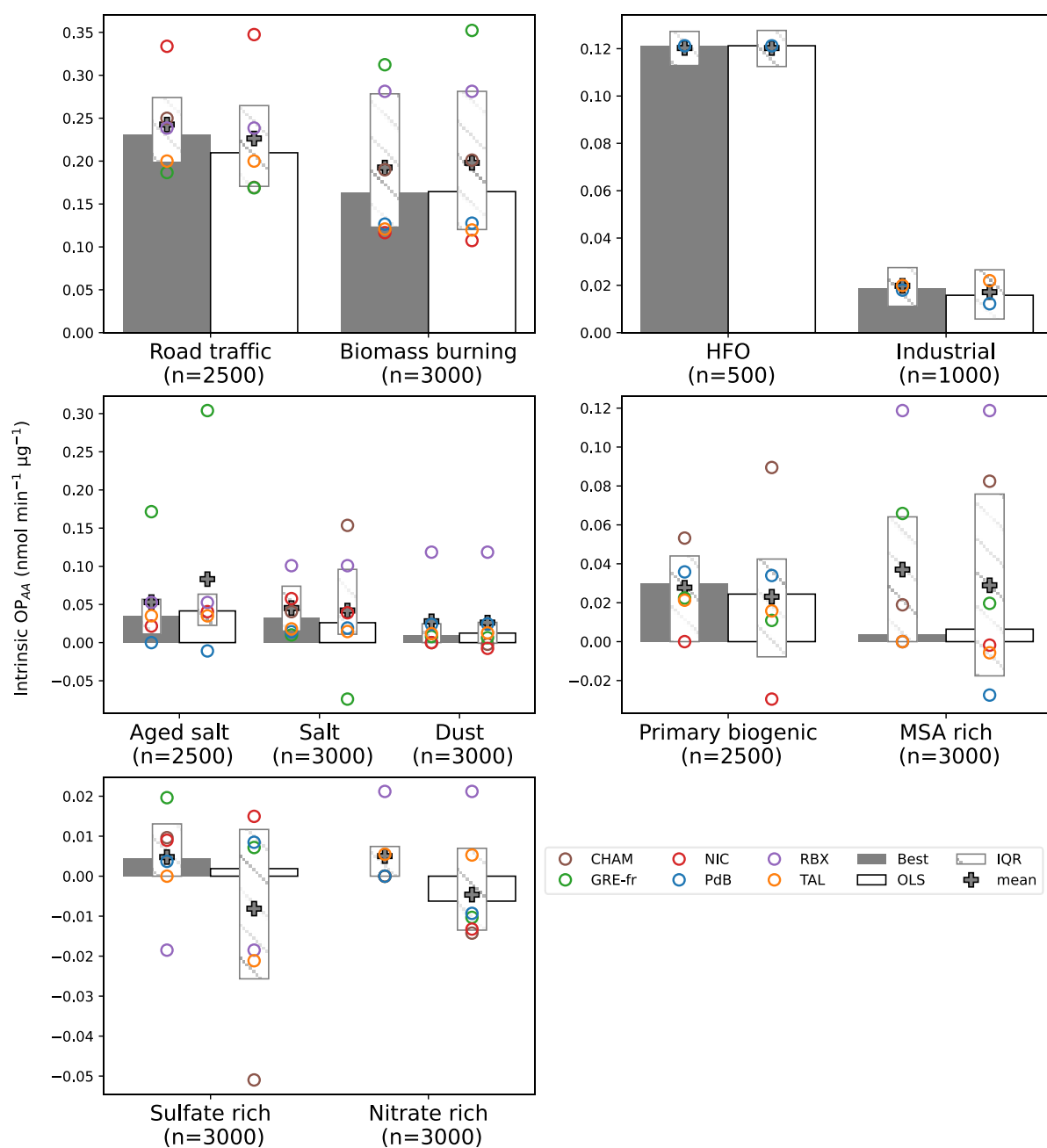
500 3.4.1. OP_{AA} activities

501 The results of the comparison of OP_{AA} 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_{AA} 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⁻¹ μg⁻¹, 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⁻¹ μg⁻¹, reference: 0.35 nmol min⁻¹ μg⁻¹).

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⁻¹ μg⁻¹), while
520 the mean intrinsic OP values were negative (-0.008 nmol min⁻¹ μg⁻¹). 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⁻¹ μg⁻¹), the intrinsic OP of salt in GRE-ft (-0.07 nmol min⁻¹ μg⁻¹) as well
 527 as the sulfate-rich source in CHAM (-0.05 nmol min⁻¹ μg⁻¹) 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

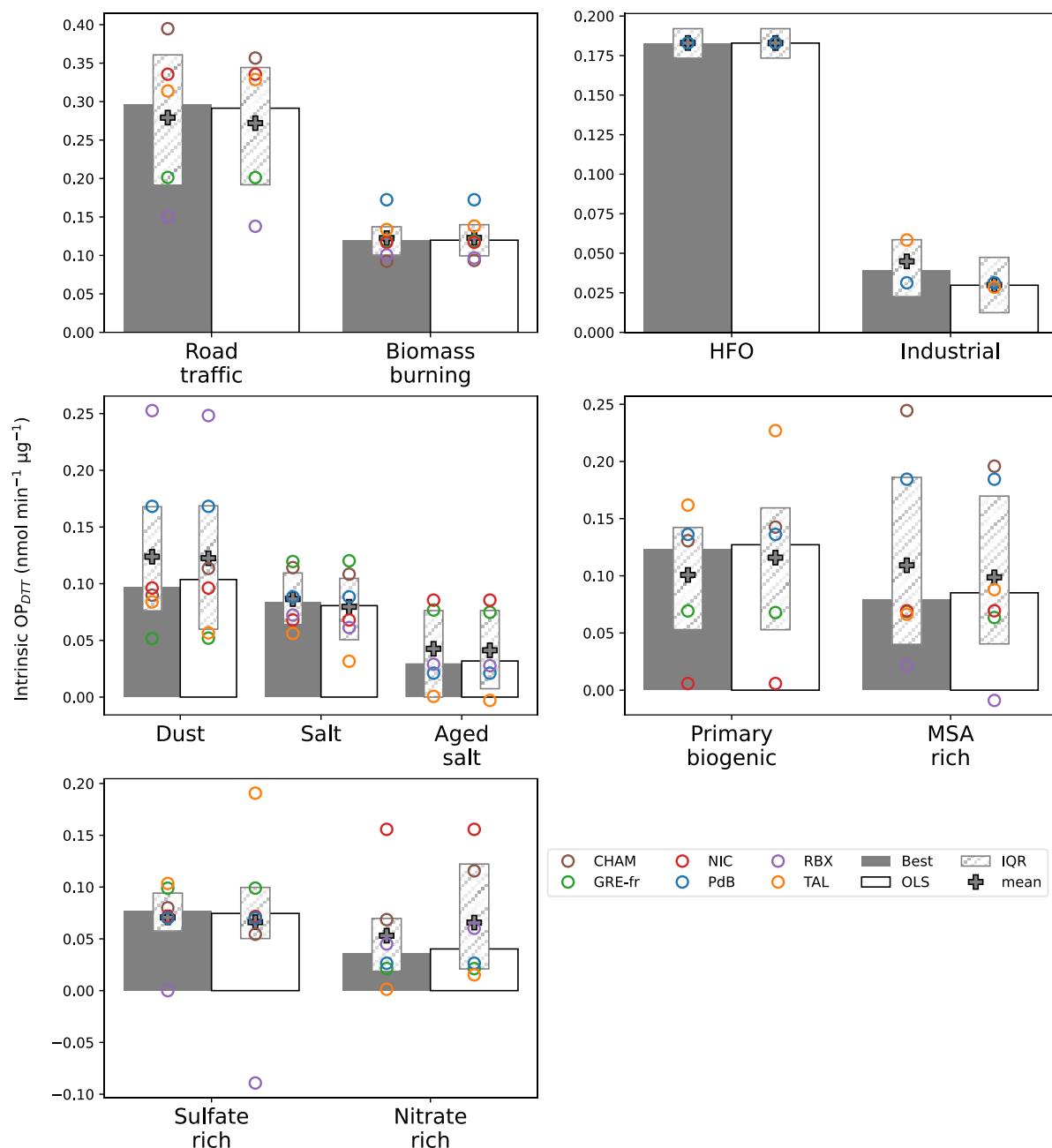
534 **Figure 8. Intrinsic OP_{AA} estimated by the best and the reference methods in the 6 sites. The y-axis represents**
535 **the intrinsic OP values in $nmol\ min^{-1}\ \mu g^{-1}$, the x-axis represents the sources. The grey bars are the median**
536 **intrinsic OP values of the best models in the 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 ± 0.10), followed by HFO (0.18 ± 0.01), biomass burning (0.12 ± 0.03), dust
551 (0.12 ± 0.07), primary biogenic (best: 0.10 ± 0.06 , reference: 0.12 ± 0.08) and MSA rich (best: 0.11 ± 0.09 , reference:
552 0.09 ± 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.



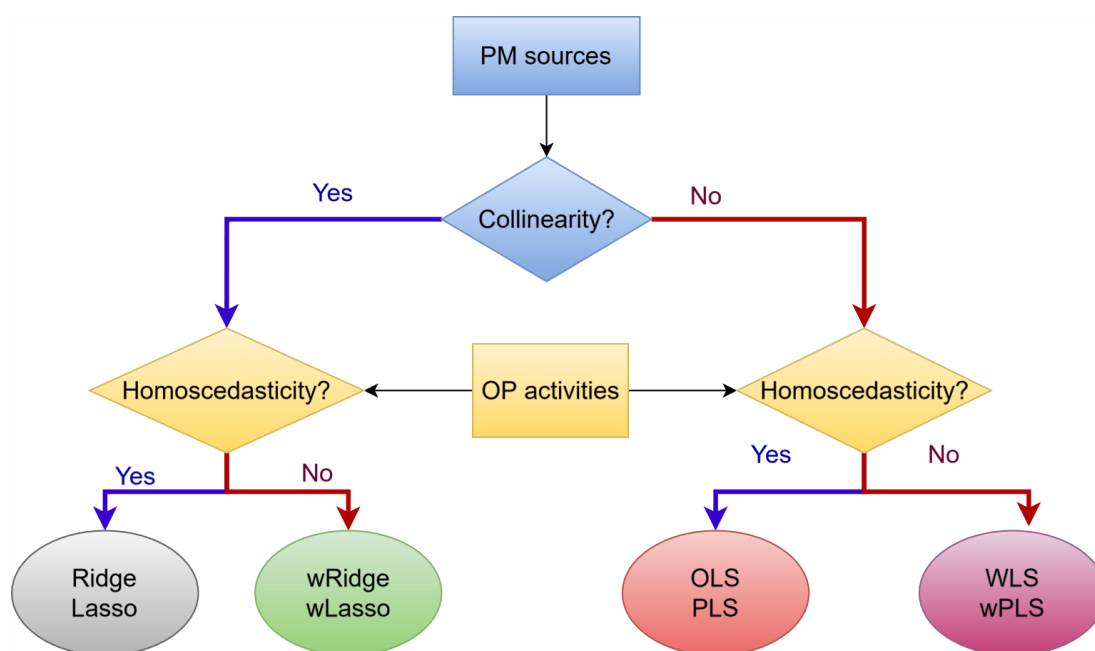
570

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 PM₁₀ sources are collinear and
591 the second is whether the residual variance of the regression between OP and PM₁₀ mass is constant. These two
592 questions represent the characteristics of PM₁₀ 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 PM₁₀ 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 (PM₁₀ 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 PM₁₀ 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⁻¹μg⁻¹) but only the importance of each PM₁₀ 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 software code could be made available by contacting the corresponding author upon request.

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. Darfeuill, 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, A., Islamia, J. M., Masood, S., Islamia, J. M., Masood, A., & Islamia, J. M. (2018). *Prediction and Analysis*
690 *of Pollution Levels in Delhi Using Multilayer Perceptron*. June. <https://doi.org/10.1007/978-981-10-3223-3>
- 691 Akhtar, McWhinney, R. D., Rastogi, N., Abbatt, J. P. D., Evans, G. J., & Scott, J. A. (2010). Cytotoxic and
692 proinflammatory effects of ambient and source-related particulate matter (PM) in relation to the production
693 of reactive oxygen species (ROS) and cytokine adsorption by particles. *Inhalation Toxicology*, 22(SUPPL.
694 2), 37–47. <https://doi.org/10.3109/08958378.2010.518377>
- 695 Alleman, L. Y., Lamaison, L., Perdrix, E., Robache, A., & Galloo, J. C. (2010). PM10 metal concentrations and
696 source identification using positive matrix factorization and wind sectoring in a French industrial zone.
697 *Atmospheric Research*, 96(4), 612–625. <https://doi.org/10.1016/j.atmosres.2010.02.008>
- 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,
700 A., Cho, A., Duggan, S., & Froines, J. (2008). Evaluating the toxicity of airborne particulate matter and
701 nanoparticles by measuring oxidative stress potential - A workshop report and consensus statement.
702 *Inhalation Toxicology*, 20(1), 75–99. <https://doi.org/10.1080/08958370701665517>
- 703 Bates, J. T., Fang, T., Verma, V., Zeng, L., Weber, R. J., Tolbert, P. E., Abrams, J. Y., Sarnat, S. E., Klein, M.,
704 Mulholland, J. A., & Russell, A. G. (2019). Review of Acellular Assays of Ambient Particulate Matter
705 Oxidative Potential: Methods and Relationships with Composition, Sources, and Health Effects.
706 *Environmental Science and Technology*, 53(8), 4003–4019. <https://doi.org/10.1021/acs.est.8b03430>
- 707 Bates, J. T., Weber, R. J., Abrams, J., Verma, V., Fang, T., Klein, M., Strickland, M. J., Sarnat, S. E., Chang, H.
708 H., Mulholland, J. A., Tolbert, P. E., & Russell, A. G. (2015). Reactive Oxygen Species Generation Linked
709 to Sources of Atmospheric Particulate Matter and Cardiorespiratory Effects. *Environmental Science and*
710 *Technology*, 49(22), 13605–13612. <https://doi.org/10.1021/acs.est.5b02967>
- 711 Bates, J. T., Weber, R. J., Verma, V., Fang, T., Ivey, C., Liu, C., Sarnat, S. E., Chang, H. H., Mulholland, J. A., &
712 Russell, A. (2018). Source impact modeling of spatiotemporal trends in PM_{2.5} oxidative potential across
713 the eastern United States. *Atmospheric Environment*, 193(August), 158–167.
714 <https://doi.org/10.1016/j.atmosenv.2018.08.055>
- 715 Beelen, R., Stafoggia, M., Raaschou-Nielsen, O., Andersen, Z. J., Xun, W. W., Katsouyanni, K., Dimakopoulou,
716 K., Brunekreef, B., Weinmayr, G., Hoffmann, B., Wolf, K., Samoli, E., Houthuijs, D., Nieuwenhuijsen, M.,
717 Oudin, A., Forsberg, B., Olsson, D., Salomaa, V., Lanki, T., ... Hoek, G. (2014). Long-term exposure to air
718 pollution and cardiovascular mortality: An analysis of 22 European cohorts. *Epidemiology*, 25(3), 368–378.
719 <https://doi.org/10.1097/EDE.0000000000000076>
- 720 Belis, C. A., Karagulian, F., Amato, F., Almeida, M., Artaxo, P., Beddows, D. C. S., Bernardoni, V., Bove, M. C.,
721 Carbone, S., Cesari, D., Contini, D., Cuccia, E., Diapouli, E., Eleftheriadis, K., Favez, O., El Haddad, I.,
722 Harrison, R. M., Hellebust, S., Hovorka, J., ... Hopke, P. K. (2015). A new methodology to assess the
723 performance and uncertainty of source apportionment models II: The results of two European
724 intercomparison exercises. *Atmospheric Environment*, 123, 240–250.
725 <https://doi.org/10.1016/j.atmosenv.2015.10.068>
- 726 Belis, C. A., Karagulian, F., Larsen, B. R., & Hopke, P. K. (2013). Critical review and meta-analysis of ambient
727 particulate matter source apportionment using receptor models in Europe. In *Atmospheric Environment* (Vol.
728 69, pp. 94–108). <https://doi.org/10.1016/j.atmosenv.2012.11.009>
- 729 Bell, M. L., Samet, J. M., & Dominici, F. (2004). Time-series studies of particulate matter. *Annual Review of*
730 *Public Health*, 25, 247–280. <https://doi.org/10.1146/annurev.publhealth.25.102802.124329>
- 731 Benkendorf, D. J., & Hawkins, C. P. (2020). Effects of sample size and network depth on a deep learning approach
732 to species distribution modeling. *Ecological Informatics*, 60(February).
733 <https://doi.org/10.1016/j.ecoinf.2020.101137>
- 734 Borlaza. (2021). Disparities in particulate matter (PM10) origins and oxidative potential at a city scale (Grenoble,
735 France) - Part 2: Sources of PM10 oxidative potential using multiple linear regression analysis and the
736 predictive applicability of multilayer perceptron n. *Atmospheric Chemistry and Physics*, 21(12), 9719–9739.
737 <https://doi.org/10.5194/acp-21-9719-2021>
- 738 Borlaza, L., Weber, S., Jaffrezo, J. L., Houdier, S., Slama, R., Rieux, C., Albinet, A., Micallef, S., Trébluchon, C.,

- 739 & Uzu, G. (2021). Disparities in particulate matter (PM10) origins and oxidative potential at a city scale
740 (Grenoble, France) - Part 2: Sources of PM10 oxidative potential using multiple linear regression analysis
741 and the predictive applicability of multilayer perceptron n. *Atmospheric Chemistry and Physics*, 21(12),
742 9719–9739. <https://doi.org/10.5194/acp-21-9719-2021>
- 743 Borlaza, L., Weber, S., Uzu, G., Jacob, V., Cañete, T., Micallef, S., Trébuchon, C., Slama, R., Favez, O., &
744 Jaffrezo, J.-L. (2021). Disparities in particulate matter (PM10) origins and oxidative potential at a city scale
745 (Grenoble, France) - Part 1: Source apportionment at three neighbouring sites. *Atmospheric Chemistry and*
746 *Physics*, 21(7), 5415–5437. <https://doi.org/10.5194/acp-21-5415-2021>
- 747 Bourlard, H., & Wellekens, C. J. (1989). Speech pattern discrimination and multilayer perceptrons. *Computer*
748 *Speech & Language*, 3(1), 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)
- 749 Breiman, L. (2001). RFRSF: Employee Turnover Prediction Based on Random Forests and Survival Analysis.
750 *Machine Learning*, 12343 LNCS, 503–515. https://doi.org/10.1007/978-3-030-62008-0_35
- 751 Brook, R. D., Rajagopalan, S., Pope, C. A., Brook, J. R., Bhatnagar, A., Diez-Roux, A. V., Holguin, F., Hong, Y.,
752 Luepker, R. V., Mittleman, M. A., Peters, A., Siscovick, D., Smith, S. C., Whitsel, L., & Kaufman, J. D.
753 (2010). Particulate matter air pollution and cardiovascular disease: An update to the scientific statement from
754 the american heart association. *Circulation*, 121(21), 2331–2378.
755 <https://doi.org/10.1161/CIR.0b013e3181d8ce1>
- 756 Brown, S. G., Eberly, S., Paatero, P., & Norris, G. A. (2015). Methods for estimating uncertainty in PMF solutions:
757 Examples with ambient air and water quality data and guidance on reporting PMF results. *Science of the*
758 *Total Environment*, 518–519, 626–635. <https://doi.org/10.1016/j.scitotenv.2015.01.022>
- 759 Calas, A., Uzu, G., Besombes, J. L., Martins, J. M. F., Redaelli, M., Weber, S., Charron, A., Albinet, A., Chevrier,
760 F., Brulfert, G., Mesbah, B., Favez, O., & Jaffrezo, J. L. (2019). Seasonal variations and chemical predictors
761 of oxidative potential (OP) of particulate matter (PM), for seven urban French sites. *Atmosphere*, 10(11).
762 <https://doi.org/10.3390/atmos10110698>
- 763 Calas, A., Uzu, G., Kelly, F. J., Houdier, S., Martins, J. M. F., Thomas, F., Molton, F., Charron, A., Dunster, C.,
764 Oliete, A., Jacob, V., Besombes, J. L., Chevrier, F., & Jaffrezo, J. L. (2018). Comparison between five
765 acellular oxidative potential measurement assays performed with detailed chemistry on PM10 samples from
766 the city of Chamonix (France). *Atmospheric Chemistry and Physics*, 18(11), 7863–7875.
767 <https://doi.org/10.5194/acp-18-7863-2018>
- 768 Calas, A., Uzu, G., Martins, J. M. F., Voisin, Di., Spadini, L., Lacroix, T., & Jaffrezo, J. L. (2017). The importance
769 of simulated lung fluid (SLF) extractions for a more relevant evaluation of the oxidative potential of
770 particulate matter. *Scientific Reports*, 7(1), 1–12. <https://doi.org/10.1038/s41598-017-11979-3>
- 771 Chianese, E., Camastra, F., & Ciaramella, A. (2018). *Spatio-temporal learning in predicting ambient particulate*
772 *matter concentration by multi-layer perceptron Spatio-temporal Learning in Predicting Ambient Particulate*
773 *Matter Concentration by Multi-Layer*. December. <https://doi.org/10.1016/j.econinf.2018.12.001>
- 774 Cho, A., Sioutas, C., Miguel, A. H., Kumagai, Y., Schmitz, D. A., Singh, M., Eiguren-Fernandez, A., & Froines,
775 J. R. (2005). Redox activity of airborne particulate matter at different sites in the Los Angeles Basin.
776 *Environmental Research*, 99(1), 40–47. <https://doi.org/10.1016/j.envres.2005.01.003>
- 777 Cohen, J., Cohen, P., West, S. G., & Aiken, L. S. (2002). *Applied multiple regression/correlation analysis for the*
778 *behavioral sciences*. Routledge. [https://doi.org/https://doi-org.sid2nomade-](https://doi.org/https://doi-org.sid2nomade-1.grenet.fr/10.4324/9780203774441)
779 [1.grenet.fr/10.4324/9780203774441](https://doi.org/https://doi-org.sid2nomade-1.grenet.fr/10.4324/9780203774441)
- 780 Craney, T. A., & Surles, J. G. (2002). Model-dependent variance inflation factor cutoff values. *Quality*
781 *Engineering*, 14(3), 391–403. <https://doi.org/10.1081/QEN-120001878>
- 782 Crobeddu, B., Aragao-Santiago, L., Bui, L. C., Boland, S., & Baeza Squiban, A. (2017). Oxidative potential of
783 particulate matter 2.5 as predictive indicator of cellular stress. *Environmental Pollution*, 230, 125–133.
784 <https://doi.org/10.1016/j.envpol.2017.06.051>
- 785 Crouse, D. L., Peters, P. A., Hystad, P., Brook, J. R., van Donkelaar, A., Martin, R. V., Villeneuve, P. J., Jerrett,
786 M., Goldberg, M. S., Arden Pope, C., Brauer, M., Brook, R. D., Robichaud, A., Menard, R., & Burnett, R.
787 T. (2015). Ambient PM2.5, O3, and NO2 exposures and associations with mortality over 16 years of follow-
788 up in the canadian census health and environment cohort (CanCHEC). *Environmental Health Perspectives*,
789 123(11), 1180–1186. <https://doi.org/10.1289/ehp.1409276>

- 790 Crouse, D. L., Peters, P. A., van Donkelaar, A., Goldberg, M. S., Villeneuve, P. J., Brion, O., Khan, S., Atari, D.
791 O., Jerrett, M., Pope, C. A., Brauer, M., Brook, J. R., Martin, R. V., Stieb, D., & Burnett, R. T. (2012). Risk
792 of nonaccidental and cardiovascular mortality in relation to long-term exposure to low concentrations of fine
793 particulate matter: A Canadian national-level cohort study. *Environmental Health Perspectives*, *120*(5), 708–
794 714. <https://doi.org/10.1289/ehp.1104049>
- 795 Daellenbach, K. R., Uzu, G., Jiang, J., Cassagnes, L.-E., Leni, Z., Vlachou, A., Stefenelli, G., Canonaco, F., Weber,
796 S., Segers, A., & Sources, al. (2020). Sources of particulate-matter air pollution and its oxidative potential
797 in Europe of particulate-matter air pollution and its oxidative potential in Europe. *Nature*, *587*(7834).
798 <https://doi.org/10.1038/s41586-020-2902-8i>
- 799 Deng, M., Chen, D., Zhang, G., & Cheng, H. (2022). Policy-driven variations in oxidation potential and source
800 apportionment of PM_{2.5} in Wuhan, central China. *Science of the Total Environment*, *853*(May), 158255.
801 <https://doi.org/10.1016/j.scitotenv.2022.158255>
- 802 Dominici, F. (2004). Time-series analysis of air pollution and mortality: a statistical review. *Research Report*
803 (*Health Effects Institute*), *123*, 3–27.
- 804 Dominutti, P. A., Borlaza, L., Sauvain, J. J., Ngoc Thuy, V. D., Houdier, S., Suarez, G., Jaffrezo, J. L., Tobin, S.,
805 Trébuchon, C., Socquet, S., Moussu, E., Mary, G., & Uzu, G. (2023). Source apportionment of oxidative
806 potential depends on the choice of the assay: insights into 5 protocols comparison and implications for
807 mitigation measures. *Environmental Science: Atmospheres*. <https://doi.org/10.1039/d3ea00007a>
- 808 Elangasinghe, M. A., Singhal, N., Dirks, K. N., & Salmond, J. A. (2014). Development of an ANN-based air
809 pollution forecasting system with explicit knowledge through sensitivity analysis. *Atmospheric Pollution*
810 *Research*, *5*(4), 696–708. <https://doi.org/10.5094/APR.2014.079>
- 811 Fadel, M., Courcot, D., Delmaire, G., Roussel, G., Afif, C., & Ledoux, F. (2023). Source apportionment of PM_{2.5}
812 oxidative potential in an East Mediterranean site. *Science of the Total Environment*, *900*(July).
813 <https://doi.org/10.1016/j.scitotenv.2023.165843>
- 814 Fang, T., Verma, V., T Bates, J., Abrams, J., Klein, M., Strickland, J. M., Sarnat, E. S., Chang, H. H., Mulholland,
815 A. J., Tolbert, E. P., Russell, G. A., & Weber, J. R. (2016). Oxidative potential of ambient water-soluble
816 PM_{2.5} in the southeastern United States: Contrasts in sources and health associations between ascorbic acid
817 (AA) and dithiothreitol (DTT) assays. *Atmospheric Chemistry and Physics*, *16*(6), 3865–3879.
818 <https://doi.org/10.5194/acp-16-3865-2016>
- 819 Favez, O. (2017). *Traitement harmonisé de jeux de données multi-sites pour l'étude des sources de PM par*
820 *Positive Matrix Factorization*.
- 821 Godri, K. J., Harrison, R. M., Evans, T., Baker, T., Dunster, C., Mudway, I. S., & Kelly, F. J. (2011). Increased
822 oxidative burden associated with traffic component of ambient particulate matter at roadside and Urban
823 background schools sites in London. *PLoS ONE*, *6*(7). <https://doi.org/10.1371/journal.pone.0021961>
- 824 Goldfeld, S. M., & Quandt, R. E. (1965). Some Tests for Homoscedasticity Author (s): Stephen M . Goldfeld and
825 Richard E . Quandt Source : Journal of the American Statistical Association , Jun ., 1965 , Vol . 60 , No .
826 310 Published by : Taylor & Francis , Ltd . on behalf of the American Statis. *Journal of the American*
827 *Statistical Association*, *60*(310), 539–547.
- 828 Harrell. (2016). Regression Modeling Strategies. *Technometrics*, *45*(2), 170–170.
829 <https://doi.org/10.1198/tech.2003.s158>
- 830 Hastie, T. et. all. (2009). Springer Series in Statistics The Elements of Statistical Learning. *The Mathematical*
831 *Intelligencer*, *27*(2), 83–85. <http://www.springerlink.com/index/D7X7KX6772HQ2135.pdf>
- 832 Hawkins, D. M. (2004). The Problem of Overfitting. *Journal of Chemical Information and Computer Sciences*,
833 *44*(1), 1–12. <https://doi.org/10.1021/ci0342472>
- 834 Hernandez, P. A., Graham, C. H., Master, L. L., & Albert, D. L. (2006). The effect of sample size and species
835 characteristics on performance of different species distribution modeling methods. *Ecography*, *29*(5), 773–
836 785. <https://doi.org/10.1111/j.0906-7590.2006.04700.x>
- 837 Hoerl, A. E., & Kennard, R. W. (1970). Ridge Regression: Applications to Nonorthogonal Problems.
838 *Technometrics*, *12*(1), 69. <https://doi.org/10.2307/1267352>

- 839 in 't Veld, M., Pandolfi, M., Amato, F., Pérez, N., Reche, C., Dominutti, P., Jaffrezo, J., Alastuey, A., Querol, X.,
840 & Uzu, G. (2023). Discovering oxidative potential (OP) drivers of atmospheric PM₁₀, PM_{2.5}, and PM₁
841 simultaneously in North-Eastern Spain. *Science of the Total Environment*, 857(August 2022).
842 <https://doi.org/10.1016/j.scitotenv.2022.159386>
- 843 Janssen, N. A. H., Yang, A., Strak, M., Steenhof, M., Hellack, B., Gerlofs-Nijland, M. E., Kuhlbusch, T., Kelly,
844 F., Harrison, R., Brunekreef, B., Hoek, G., & Cassee, F. (2014). Oxidative potential of particulate matter
845 collected at sites with different source characteristics. *Science of the Total Environment*, 472, 572–581.
846 <https://doi.org/10.1016/j.scitotenv.2013.11.099>
- 847 Kelly, F. J., & Mudway, I. S. (2003). Protein oxidation at the air-lung interface. *Amino Acids*, 25(3–4), 375–396.
848 <https://doi.org/10.1007/s00726-003-0024-x>
- 849 Kuhn, M., & Johnson, K. (2013). Applied predictive modeling. In *Applied Predictive Modeling*.
850 <https://doi.org/10.1007/978-1-4614-6849-3>
- 851 Leni, Z., Cassagnes, L. E., Daellenbach, K. R., Haddad, I. El, Vlachou, A., Uzu, G., Prévôt, A. S. H., Jaffrezo, J.
852 L., Baumlin, N., Salathe, M., Baltensperger, U., Dommen, J., & Geiser, M. (2020). Oxidative stress-induced
853 inflammation in susceptible airways by anthropogenic aerosol. *PLoS ONE*, 15(11 November).
854 <https://doi.org/10.1371/journal.pone.0233425>
- 855 Li, J., Zhao, S., Xiao, S., Li, X., Wu, S., Zhang, J., & Schwab, J. J. (2023). Source apportionment of water-soluble
856 oxidative potential of PM_{2.5} in a port city of Xiamen, Southeast China. *Atmospheric Environment*,
857 314(June), 120122. <https://doi.org/10.1016/j.atmosenv.2023.120122>
- 858 Li, Xia, T., & Nel, A. E. (2008). The role of oxidative stress in ambient particulate matter-induced lung diseases
859 and its implications in the toxicity of engineered nanoparticles. *Free Radical Biology and Medicine*, 44(9),
860 1689–1699. <https://doi.org/10.1016/j.freeradbiomed.2008.01.028>
- 861 Liu, & Ng. (2023). Toxicity of Atmospheric Aerosols: Methodologies & Assays. *American Chemical Society*.
862 <https://doi.org/DOI:10.1021/acscinfocus.7e7012>
- 863 Liu, W. J., Xu, Y. S., Liu, W. X., Liu, Q. Y., Yu, S. Y., Liu, Y., Wang, X., & Tao, S. (2018). Oxidative potential
864 of ambient PM_{2.5} in the coastal cities of the Bohai Sea, northern China: Seasonal variation and source
865 apportionment. *Environmental Pollution*, 236, 514–528. <https://doi.org/10.1016/j.envpol.2018.01.116>
- 866 Lodovici, M., & Bigagli, E. (2011). Oxidative stress and air pollution exposure. *Journal of Toxicology*, 2011.
867 <https://doi.org/10.1155/2011/487074>
- 868 Matsuki, K., Kuperman, V., & Van Dyke, J. A. (2016). The Random Forests statistical technique: An examination
869 of its value for the study of reading. *Scientific Studies of Reading*, 20(1), 20–33.
870 <https://doi.org/10.1080/10888438.2015.1107073>
- 871 McCullagh. (1989). Generalized linear models. In *Statistical Models in S* (pp. 195–247).
872 <https://doi.org/10.1201/9780203738535>
- 873 Montgomery C, D., Peck A, E., & Vining, G. G. (2012). *Introducing To Linear Regression Analysis (5th ed.)*.
- 874 Mudway, I. S., Kelly, F. J., & Holgate, S. T. (2020). Oxidative stress in air pollution research. In *Free Radical*
875 *Biology and Medicine* (Vol. 151, pp. 2–6). Elsevier Inc.
876 <https://doi.org/10.1016/j.freeradbiomed.2020.04.031>
- 877 Nelin, T. D., Joseph, A. M., Gorr, M. W., & Wold, L. E. (2012). Direct and indirect effects of particulate matter
878 on the cardiovascular system. *Toxicology Letters*, 208(3), 293–299.
879 <https://doi.org/10.1016/j.toxlet.2011.11.008>
- 880 O'Brien, R. M. (2007). A caution regarding rules of thumb for variance inflation factors. *Quality and Quantity*,
881 41(5), 673–690. <https://doi.org/10.1007/s11135-006-9018-6>
- 882 Paatero, P., & Hopke, P. K. (2009). Rotational tools for factor analytic models. *Journal of Chemometrics*, 23(2),
883 91–100. <https://doi.org/10.1002/cem.1197>
- 884 Paatero, P., & Tappert, U. (1994). Positive matrix factorization: A non-negative factor model with optimal
885 utilization of error estimates of data values. In *Environmetrics* (Vol. 5).
886 <https://doi.org/https://doi.org/10.1002/env.3170050203>

- 887 Pearce, J., & Ferrier, S. (2000). An evaluation of alternative algorithms for fitting species distribution models using
888 logistic regression. *Ecological Modelling*, *128*(2–3), 127–147. [https://doi.org/10.1016/S0304-](https://doi.org/10.1016/S0304-3800(99)00227-6)
889 [3800\(99\)00227-6](https://doi.org/10.1016/S0304-3800(99)00227-6)
- 890 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P.,
891 Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., & Duchesnay,
892 E. (2011). Scikit-learn: Machine Learning in {P}ython. *Journal of Machine Learning Research*, *12*, 2825–
893 2830.
- 894 Pelucchi, C., Negri, E., Gallus, S., Boffetta, P., Tramacere, I., & La Vecchia, C. (2009). Long-term particulate
895 matter exposure and mortality: A review of European epidemiological studies. *BMC Public Health*, *9*, 1–8.
896 <https://doi.org/10.1186/1471-2458-9-453>
- 897 Peng, R. D., Bell, M. L., Geyh, A. S., McDermott, A., Zeger, S. L., Samet, J. M., & Dominici, F. (2009).
898 Emergency admissions for cardiovascular and respiratory diseases and the chemical composition of fine
899 particle air pollution. *Environmental Health Perspectives*, *117*(6), 957–963.
900 <https://doi.org/10.1289/ehp.0800185>
- 901 Pietrogrande, M. C., Romanato, L., & Russo, M. (2022). Synergistic and Antagonistic Effects of Aerosol
902 Components on Its Oxidative Potential as Predictor of Particle Toxicity. *Toxics*, *10*(4).
903 <https://doi.org/10.3390/toxics10040196>
- 904 Pope, C. A., & Dockery, D. W. (2006). Health effects of fine particulate air pollution: Lines that connect. *Journal*
905 *of the Air and Waste Management Association*, *56*(6), 709–742.
906 <https://doi.org/10.1080/10473289.2006.10464485>
- 907 Rao, X., Zhong, J., Brook, R. D., & Rajagopalan, S. (2018). Effect of Particulate Matter Air Pollution on
908 Cardiovascular Oxidative Stress Pathways. *Antioxidants and Redox Signaling*, *28*(9), 797–818.
909 <https://doi.org/10.1089/ars.2017.7394>
- 910 Raudys, S. J., & Jain, A. K. (1991). Small Sample Size Effects in Statistical Pattern Recognition:
911 Recommendations for Practitioners. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*
912 (Vol. 13, Issue 3, pp. 252–264). <https://doi.org/10.1109/34.75512>
- 913 Rosenblad, A. (2011). The Concise Encyclopedia of Statistics. In *Journal of Applied Statistics* (Vol. 38, Issue 4).
914 <https://doi.org/10.1080/02664760903075614>
- 915 Samake, A., Uzu, G., Martins, J. M. F., Calas, A., Vince, E., Parat, S., & Jaffrezo, J. L. (2017). The unexpected
916 role of bioaerosols in the Oxidative Potential of PM. *Scientific Reports*, *7*(1). [https://doi.org/10.1038/s41598-](https://doi.org/10.1038/s41598-017-11178-0)
917 [017-11178-0](https://doi.org/10.1038/s41598-017-11178-0)
- 918 Seabold, S., & Perktold, J. (2010). Statsmodels: Econometric and statistical modeling with python. *9th Python in*
919 *Science Conference*.
- 920 Shangguan, Y., Zhuang, X., Querol, X., Li, B., Moreno, N., Trechera, P., Sola, P. C., Uzu, G., & Li, J. (2022).
921 Characterization of deposited dust and its respirable fractions in underground coal mines: Implications for
922 oxidative potential-driving species and source apportionment. *International Journal of Coal Geology*,
923 *258*(December 2021). <https://doi.org/10.1016/j.coal.2022.104017>
- 924 Stevanović, S., Jovanović, M. V., Jovašević-Stojanović, M. V., & Ristovski, Z. (2023). SOURCE
925 APPORTIONMENT OF OXIDATIVE POTENTIAL What We Know So Far. *Thermal Science*, *27*(3),
926 2347–2357. <https://doi.org/10.2298/TSCI221107111S>
- 927 Stockwell, D. R. B., & Peterson, A. T. (2002). Effects of sample size on accuracy of species distribution models.
928 *Ecological Modelling*, *148*(1), 1–13. [https://doi.org/10.1016/S0304-3800\(01\)00388-X](https://doi.org/10.1016/S0304-3800(01)00388-X)
- 929 Szigeti, T., Dunster, C., Cattaneo, A., Cavallo, D., Spinazzè, A., Saraga, D. E., Sakellaris, I. A., de Kluizenaar, Y.,
930 Cornelissen, E. J. M., Hänninen, O., Peltonen, M., Calzolari, G., Lucarelli, F., Mandin, C., Bartzis, J. G.,
931 Záray, G., & Kelly, F. J. (2016). Oxidative potential and chemical composition of PM_{2.5} in office buildings
932 across Europe - The OFFICAIR study. *Environment International*, *92–93*, 324–333.
933 <https://doi.org/10.1016/j.envint.2016.04.015>
- 934 Szigeti, T., Óvári, M., Dunster, C., Kelly, F. J., Lucarelli, F., & Záray, G. (2015). Changes in chemical composition
935 and oxidative potential of urban PM_{2.5} between 2010 and 2013 in Hungary. *Science of the Total*
936 *Environment*, *518–519*, 534–544. <https://doi.org/10.1016/j.scitotenv.2015.03.025>

- 937 Tibshirani, R. (1996). Regression Shrinkage and Selection Via the Lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1), 267–288. <https://doi.org/10.1111/j.2517-6161.1996.tb02080.x>
- 938
- 939 Verma, V., Fang, T., Guo, H., King, L., Bates, J. T., Peltier, R. E., Edgerton, E., Russell, A. G., & Weber, R. J. (2014). Reactive oxygen species associated with water-soluble PM_{2.5} in the southeastern United States: Spatiotemporal trends and source apportionment. *Atmospheric Chemistry and Physics*, 14(23), 12915–12930. <https://doi.org/10.5194/acp-14-12915-2014>
- 940
- 941
- 942
- 943 Viana, M., Kuhlbusch, T. A. J., Querol, X., Alastuey, A., Harrison, R. M., Hopke, P. K., Winiwarter, W., Vallius, M., Szidat, S., Prévôt, A. S. H., Hueglin, C., Bloemen, H., Wählin, P., Vecchi, R., Miranda, A. I., Kasper-Giebl, A., Maenhaut, W., & Hitzenberger, R. (2008). Source apportionment of particulate matter in Europe: A review of methods and results. In *Journal of Aerosol Science* (Vol. 39, Issue 10, pp. 827–849). Elsevier Ltd. <https://doi.org/10.1016/j.jaerosci.2008.05.007>
- 944
- 945
- 946
- 947
- 948 Vida, M., Foret, G., Siour, G., Coman, A., Weber, S., Favez, O., Jaffrezo, J., Pontet, S., Mesbah, B., Gille, G., Zhang, S., Chevrier, F., Pallares, C., Uzu, G., & Beekmann, M. (2024). Oxidative potential modelling of PM₁₀: a 2-year study over France. *ACDP*.
- 949
- 950
- 951 Wang, D., Yang, X., Lu, H., Li, D., Xu, H., Luo, Y., Sun, J., Hang Ho, S. S., & Shen, Z. (2023). Oxidative potential of atmospheric brown carbon in six Chinese megacities: Seasonal variation and source apportionment. *Atmospheric Environment*, 309(June), 119909. <https://doi.org/10.1016/j.atmosenv.2023.119909>
- 952
- 953
- 954 Wang, J., Jiang, H., Jiang, H., Mo, Y., Geng, X., Li, J., Mao, S., Bualert, S., Ma, S., Li, J., & Zhang, G. (2020). Source apportionment of water-soluble oxidative potential in ambient total suspended particulate from Bangkok: Biomass burning versus fossil fuel combustion. *Atmospheric Environment*, 235(May), 117624. <https://doi.org/10.1016/j.atmosenv.2020.117624>
- 955
- 956
- 957
- 958 Wang, S., Ye, J., Soong, R., Wu, B., Yu, L., Simpson, A. J., & Chan, A. W. H. (2018). Relationship between chemical composition and oxidative potential of secondary organic aerosol from polycyclic aromatic hydrocarbons. *Atmospheric Chemistry and Physics*, 18(6), 3987–4003. <https://doi.org/10.5194/acp-18-3987-2018>
- 959
- 960
- 961
- 962 Wang, Y., Wang, M., Li, S., Sun, H., Mu, Z., Zhang, L., Li, Y., & Chen, Q. (2020). Study on the oxidation potential of the water-soluble components of ambient PM_{2.5} over Xi'an, China: Pollution levels, source apportionment and transport pathways. *Environment International*, 136(January), 105515. <https://doi.org/10.1016/j.envint.2020.105515>
- 963
- 964
- 965
- 966 Weber, S., Salameh, D., Albinet, A., Alleman, L. Y., Waked, A., Besombes, J. L., Jacob, V., Guillaud, G., Meshbah, B., Rocq, B., Hulin, A., Dominik-Sègue, M., Chrétien, E., Jaffrezo, J. L., & Favez, O. (2019). Comparison of PM₁₀ sources profiles at 15 french sites using a harmonized constrained positive matrix factorization approach. *Atmosphere*, 10(6). <https://doi.org/10.3390/atmos10060310>
- 967
- 968
- 969
- 970 Weber, S., Uzu, G., Calas, A., Chevrier, F., Besombes, J. L., Charron, A., Salameh, D., Ježek, I., Močnik, G., & Jaffrezo, J. L. (2018). An apportionment method for the oxidative potential of atmospheric particulate matter sources: Application to a one-year study in Chamonix, France. *Atmospheric Chemistry and Physics*, 18(13), 9617–9629. <https://doi.org/10.5194/acp-18-9617-2018>
- 971
- 972
- 973
- 974 Weber, S., Uzu, G., Favez, O., Borlaza, L., Calas, A., Salameh, D., Chevrier, F., Allard, J., Besombes, J. L., Albinet, A., Pontet, S., Mesbah, B., Gille, G., Zhang, S., Pallares, C., Leoz-Garziandia, E., & Jaffrezo, J. L. (2021). Source apportionment of atmospheric PM₁₀ oxidative potential: Synthesis of 15 year-round urban datasets in France. *Atmospheric Chemistry and Physics*, 21(14), 11353–11378. <https://doi.org/10.5194/acp-21-11353-2021>
- 975
- 976
- 977
- 978
- 979 WHO. (2021). *WHO global air quality guidelines*.
- 980 Williams, M., Gomez Grajales, C. A., & Kurkiewicz, D. (2013). Assumptions of Multiple Regression: Correcting Two Misconceptions - Practical Assessment, Research & Evaluation. *Practical Assessment, Research, and Evaluation (PARE)*, 18(11), 1–16. <https://scholarworks.umass.edu/pare/vol18/iss1/11>
- 981
- 982
- 983 Wisz, M. S., Hijmans, R. J., Li, J., Peterson, A. T., Graham, C. H., Guisan, A., Elith, J., Dudík, M., Ferrier, S., Huettmann, F., Leathwick, J. R., Lehmann, A., Lohmann, L., Loiselle, B. A., Manion, G., Moritz, C., Nakamura, M., Nakazawa, Y., Overton, J. M. C., ... Zimmermann, N. E. (2008). Effects of sample size on the performance of species distribution models. *Diversity and Distributions*, 14(5), 763–773. <https://doi.org/10.1111/j.1472-4642.2008.00482.x>
- 984
- 985
- 986
- 987

- 988 Xiong, Q., Yu, H., Wang, R., Wei, J., & Verma, V. (2017). Rethinking Dithiothreitol-Based Particulate Matter
989 Oxidative Potential: Measuring Dithiothreitol Consumption versus Reactive Oxygen Species Generation.
990 *Environmental Science and Technology*, 51(11), 6507–6514. <https://doi.org/10.1021/acs.est.7b01272>
- 991 Yang, A., Jedynska, A., Hellack, B., Kooter, I., Hoek, G., Brunekreef, B., Kuhlbusch, T. A. J., Cassee, F. R., &
992 Janssen, N. A. H. (2014). Measurement of the oxidative potential of PM_{2.5} and its constituents: The effect
993 of extraction solvent and filter type. *Atmospheric Environment*, 83, 35–42.
994 <https://doi.org/10.1016/j.atmosenv.2013.10.049>
- 995 Yu, Guo, S., Xu, R., Ye, T., Li, S., Sim, M. R., Abramson, M. J., & Guo, Y. (2021). Cohort studies of long-term
996 exposure to outdoor particulate matter and risks of cancer: A systematic review and meta-analysis.
997 *Innovation*, 2(3), 100143. <https://doi.org/10.1016/j.xinn.2021.100143>
- 998 Yu, S. Y., Liu, W. J., Xu, Y. S., Yi, K., Zhou, M., Tao, S., & Liu, W. X. (2019). Characteristics and oxidative
999 potential of atmospheric PM_{2.5} in Beijing: Source apportionment and seasonal variation. *Science of the*
1000 *Total Environment*, 650, 277–287. <https://doi.org/10.1016/j.scitotenv.2018.09.021>
- 1001 Zhang, Y., Albinet, A., Petit, J. E., Jacob, V., Chevrier, F., Gille, G., Pontet, S., Chrétien, E., Dominik-Sègue, M.,
1002 Levigoureux, G., Močnik, G., Gros, V., Jaffrezo, J. L., & Favez, O. (2020). Substantial brown carbon
1003 emissions from wintertime residential wood burning over France. *Science of the Total Environment*, 743.
1004 <https://doi.org/10.1016/j.scitotenv.2020.140752>
- 1005