

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) (Crouse et al., 2012). 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

41 implications for oxidative stress and cellular damage (Akhtar et al., 2010; Ayres et al., 2008; Leni et al., 2020; Li
42 et al., 2008; Lodovici & Bigagli, 2011; Mudway et al., 2020; Nelin et al., 2012; Rao et al., 2018). The quantification
43 of the PM capacity to oxidize a biological media is called oxidative potential (OP) (Bates et al., 2019; Daellenbach
44 et al., 2020; Dominutti et al., 2023). Various acellular assays of OP have been introduced, differentiating ROS
45 generation mechanisms of PM (Calas et al., 2018; Dominutti et al., 2023). Dithiothreitol (DTT) and ascorbic acid
46 (AA) assays are two of the 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 [PM₁₀](#) OP SA techniques by comparing eight regression techniques:
75 multiple linear regression fitted by ~~ordinary least squares (OLS)~~, weighted least squares (WLS), positive least
76 squares (PLS), Ridge regression (Ridge), Least Absolute Shrinkage and Selection Operator (Lasso), generalized
77 linear model (GLM), random forest (RF), and multilayer perceptron (MLP). These techniques are applied to
78 apportion [PM₁₀](#)-OP_{AA} and [PM₁₀](#)-OP_{DTT} to [PM₁₀](#) sources at six sites in France. The [PM₁₀](#) SA outputs have been
79 published previously in Weber et al., (2021), using a harmonized PMF methodology based on one year of sampling
80 with similar chemical analyses for a large set of chemical tracers. The results of the [PM₁₀](#) OP SA models are
81 compared with regard to the estimated intrinsic [PM₁₀](#) OP of each source, the out-of-sample accuracy of the
82 apportionment, and the assumptions inherent in each model. The most appropriate model at each site is compared
83 with OLS to quantify the difference between choosing a model based on data characteristics vs. using the most
84 common approach. Finally, this study provides guidelines for selecting the most suitable model in the strategy for
85 OP contribution regarding sources of [PM₁₀](#). This holds particular significance in the context of the implementation

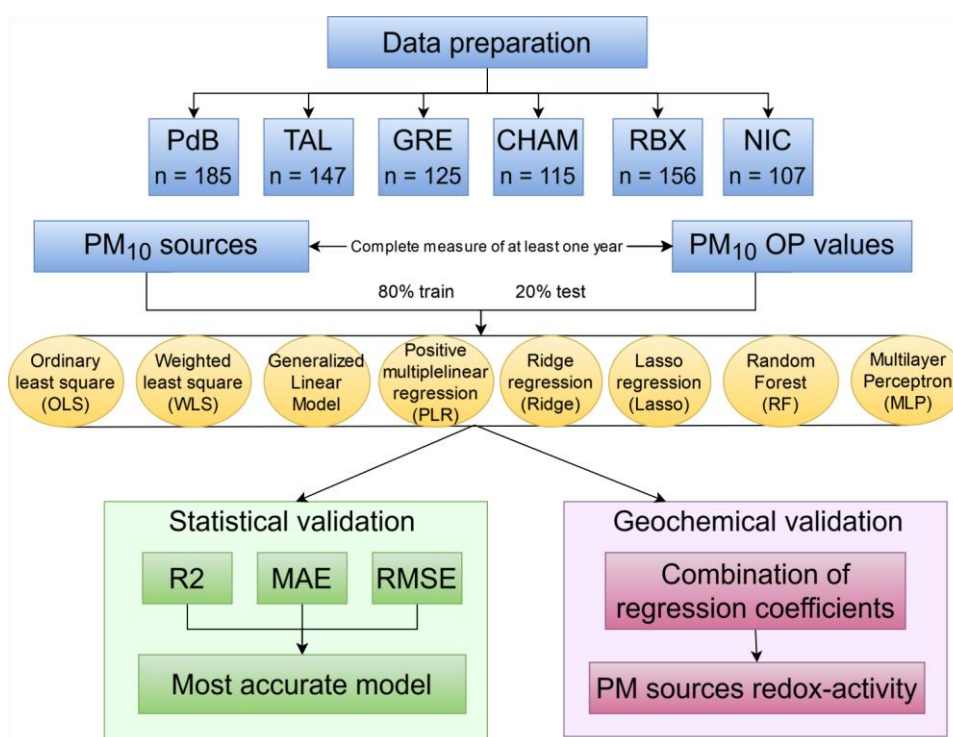
86 of OP monitoring as a novel air quality metric as foreseen in research programs (such RI-Urbans) and in the process
 87 of the revision of the European Directive 2008/50/CE.

88 2. Methodology

89 2.1. General organisation of this work

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

108



109

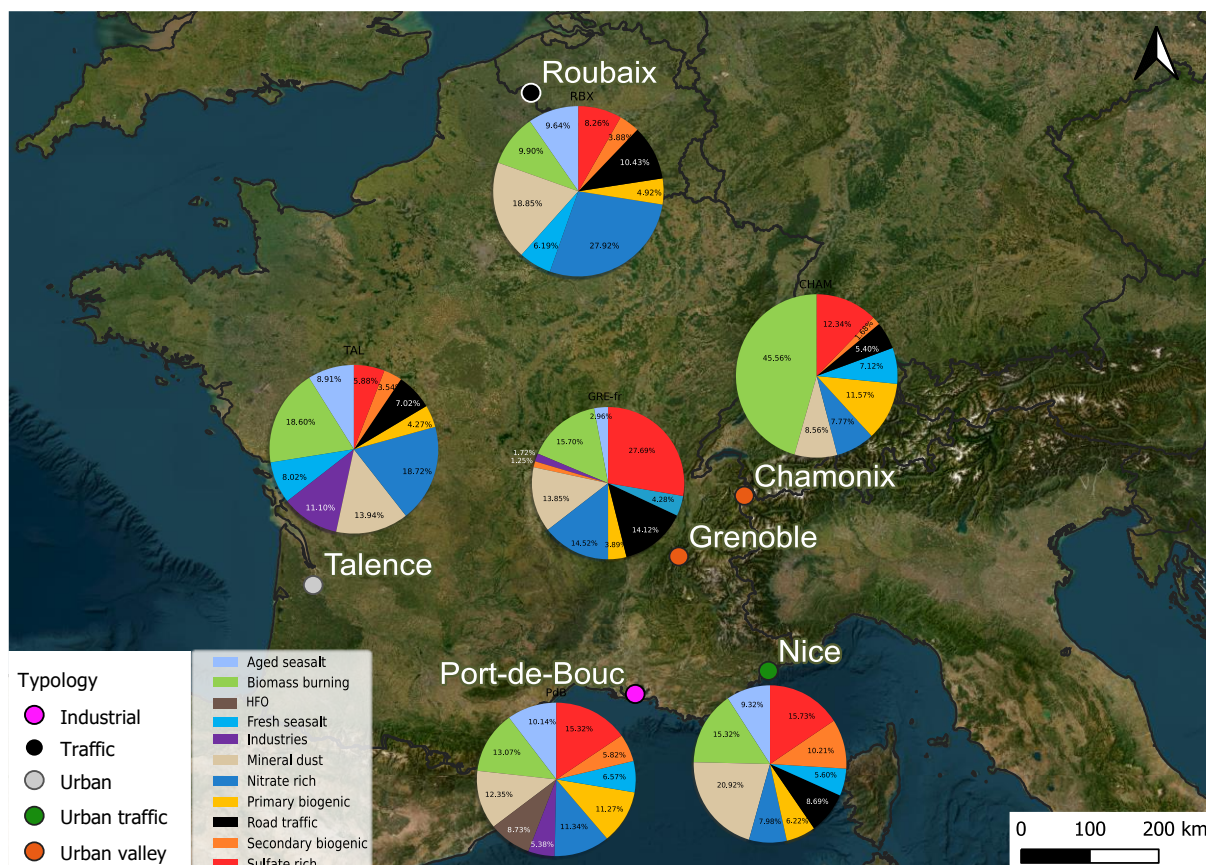
110 **Figure 1. Workflow of the comparison of PM_{10} OP sources apportionment methodology**

2.2. Study sites and PM₁₀ sources

Six French sites are selected in this work for their different typologies: Roubaix and Nice (traffic sites within urban areas), Port-de-Bouc (industrial hotspot), Talence (urban background site), Grenoble and Chamonix (urban background sites in Alpine Valley). At each site, sampling was conducted over at least one year to capture the complete annual evolution of PM₁₀ and its components. These sites and sampling series were previously used and described by Weber et al. (2019).

In brief, daily filter samples were collected on pre-heated Pallflex quartz fibre filters every third day through high-volume sampling (DA80, Digitel). These filters were analyzed to determine PM's chemical species and OP activities. Further details regarding the chemical species and PM₁₀-OP analyses methodology can be found in Weber et al. (Weber et al., 2019, 2021). Briefly, the elemental carbon (EC) and organic carbon (OC) were analyzed using the EUSAAR2 thermo-optical protocol with a Sunset Lab analyser. Major ionic components (Cl⁻, NO₃⁻, SO₄²⁻, NH₄⁺, Na⁺, K⁺, Mg²⁺, Ca²⁺) and methanesulfonic acid (MSA) were measured by ion chromatography (IC). Anhydro-sugars and saccharides (including levoglucosan, mannosan, arabitol, sorbitol, and mannitol) were analysed by high-performance liquid chromatography with pulsed amperometry detection (HPLC-PAD). Major 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 by inductively coupled plasma atomic emission spectroscopy or mass spectrometry (ICP-AES or ICP-MS). Furthermore, colocated PM₁₀ measurements were conducted automatically at each site using the Tapered Element Oscillating Microbalance equipped with a Filter Dynamics Measurement System (TEOM-FDMS).

We used the PM₁₀ sources identified by Weber et al. (2019), who performed a separate PMF for each site using a harmonized approach for all sites (same chemical species and measurement methods, same procedure to estimate uncertainties, same constraints on the preliminary solutions). Table 1 provides a data description, including the sampling duration, the number of samples collected, and the identified PM₁₀ sources at each site, while Figure 2 presents the localisation of the sites in France, together with the respective proportion of each PM₁₀ source at each 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

oxidase, thereby replicating in vivo interactions between PM₁₀ and biological oxidants (Calas et al., 2018; Cho et al., 2005). Both assays measured the consumption of AA or DTT during the assay, i.e., the rate of the transfer of electrons from AA or DTT to oxygen. The assays were conducted with 96-well plates of UV-transparent quality (CELLSTAR, Greiner-Bio), and absorption measurements were acquired using a TECAN spectrophotometer, Infinite M200 Pro, at the wavelengths of 265nm for the AA assay and 412nm for the DTT assay (Calas et al., 2017, 2018, 2019). Each sample extraction was subjected to four analyses; the [PM₁₀ OP activities](#) in this study represent the mean and the analysis uncertainty is the standard deviation of these four [PM₁₀ OP](#) analyses. After analysis, the [PM₁₀ OP](#) activities of each sample were blank-subtracted using lab and field blanks, and normalized using the air sampling volumes and the mass concentration. The resulting OP_v represents the [PM₁₀ OP activities](#) due to PM₁₀ per cubic meter of air (nmol min⁻¹ m⁻³). [To simplify the denotation of PM₁₀ OP, OP is used to represent PM₁₀ OP throughout this article.](#)

2.4. Collinearity and heteroscedasticity tests

The result of a regression model strongly depends on the characteristics of the dataset because each model makes assumptions about the data. Two critical assumptions in OLS regression analysis are that (1) there is little collinearity between independent variables (the PM₁₀ sources in this study), and (2) the variance of the regression residuals is constant (called homoscedasticity). These assumptions should be tested in different ways.

2.4.1. Collinearity

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

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

In this equation, p is the number of PM₁₀ sources, R^2 is the coefficient of determination of a multiple linear regression model between the i^{th} source and the other sources. VIF values of a PM₁₀ source present a range between 1, and ∞ . The higher the VIF values, the greater the collinearity between this PM₁₀ source and the other ones. A VIF value between 5 and 10 is commonly interpreted as moderate collinearity, while values greater than 10 indicate high collinearity (Craney & Surles, 2002).

2.4.2. Heteroscedasticity

Heteroscedasticity occurs when the variance of regression residuals is not constant but varies for different values of the dependent variable. In this case, the estimated standard errors of the regression coefficients are not reliable. The Goldfeld–Quandt test was developed by Goldfeld & Quandt (1965) to evaluate residual variance in a regression model. To implement the Goldfeld–Quandt test, an OLS regression was performed between OP and PM₁₀ sources to identify the residual of OP prediction. Next, the PM₁₀ sources and residual corresponding are divided into three segments: the upper segment is the group with higher PM₁₀ sources concentration, the lower segment is the group with lower PM₁₀ sources concentration, and the middle segment, constituting 10% of the moderate PM₁₀ concentration, is excluded. A subsequent regression analysis is then conducted on the two remaining subgroups to determine the ratio of residual sums of squares. Finally, an F-test is conducted on this ratio to assess whether the variances are the same, with a p-value below 0.05 interpreted as evidence of heteroscedasticity.

The Variance Inflation Factor (VIF) and the Goldfeld–Quandt test were performed in Python 3.9, using the 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 ($\text{nmol min}^{-1} \text{m}^{-3}$), X are the PM₁₀ source contributions ($\mu\text{g m}^{-3}$), y is the observed OP
197 ($\text{nmol min}^{-1} \text{m}^{-3}$), and e denotes the residuals ($\text{nmol min}^{-1} \text{m}^{-3}$). 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 ($\text{nmol min}^{-1} \text{m}^{-3}$), β_i represents the regression coefficient (intrinsic OP, nmol
208 $\text{min}^{-1} \mu\text{g}^{-1}$) of source i , x_i is the concentration of source i ($\mu\text{g m}^{-3}$), 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.

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

$$233 \text{ Minimize: } \frac{1}{2m} \sum_{i=1}^m w_i (y_i - \hat{y}_i)^2 + \lambda * \sum_{j=1}^p \beta_j^2 \text{ (Eq9)}$$

234 where λ is the parameter representing the amount of shrinkage, the larger λ , the greater the shrinkage. The
 235 hyperparameter tuning was implemented with different values of λ (5, 1, 0.5, 0.1, 0.01, 0.005, 0.001, 0.0005,
 236 0.0001). The best λ for every site varied from 0.005 to 0.01 and in this study, 0.01 was selected. Ridge can be
 237 applied with weighting to account for heteroscedasticity.

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

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

$$245 \text{ Minimize: } \sum_{i=1}^m (y_i - \hat{y}_i)^2 + \lambda * \sum_{j=1}^p |\beta_j| \text{ (Eq10)}$$

$$246 \text{ Minimize: } \frac{1}{2m} \sum_{i=1}^m w_i (y_i - \hat{y}_i)^2 + \lambda * \sum_{j=1}^p |\beta_j| \text{ (Eq11)}$$

247

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

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

251 Generalized linear models, as introduced by McCullagh (1989), provide a framework for regression analysis that
 252 can contain non-normal error distributions and capture non-linear relationships between OP activities and PM₁₀
 253 sources. GLM allows for error variance that is a function of the predicted value, hence accounting for
 254 heteroskedasticity. Key assumptions underlying GLM include (1) independence, (2) the non-normal distribution
 255 of OP, and (3) the relationship between the PM₁₀ sources and the transformed OP (logarithm in this study) is linear.
 256 The mathematical expression for GLM can be represented as follows:

257
$$\log(\hat{y}) = \beta_0 + \sum_0^p \beta_i * x_i \text{ (Eq102)}$$

258 where β_0 denotes the intercept, β_i represents the regression coefficient of source i, and x_i is the concentration of
259 source i.

260 **Random forest (RF):**

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

271 RF does not assume a specific equation to express the relationship between OP activities and PM₁₀ sources, with
272 the result that intrinsic OP could not be computed in this regression model. Nevertheless, RF can estimate the
273 relative importance of each PM₁₀ source in OP prediction. This study estimated the permutation importance of
274 each PM₁₀ source as the mean increase in the mean squared error of predicted OP when the values of the PM₁₀
275 source were permuted.

276 **Multilayer perception (MLP):**

277 MLP is an artificial neural network that consists of multiple layers of interconnected nodes or neurons organized
278 in a feedforward structure (Akhtar et al., 2018; Boudlard & Wellekens, 1989; Chianese et al., 2018). These layers
279 include an input layer (PM₁₀ sources), one or several hidden layers, and an output layer (OP_{AA} or OP_{DTT} activities).
280 In MLP, the neurons in the hidden layers are linked with the previous neurons by the connection weight, where
281 every neuron is independent and has a different weight. The output of each neuron depends on its inputs and an
282 activation function, which, if non-linear, allows the model to capture non-linear relationships. The implementation
283 of MLP includes three steps: (1) forward pass to training model: the input is passed to the model, multiplied with
284 an initial weight, add bias at every layer, then calculate output of the model. (2) error calculation: after applying
285 step 1, the output of the model and the observed data are used to calculate the error. (3) backward pass: the error
286 is propagated back through the network, and then the weights are adjusted to minimize overall error. These 3 steps
287 are repeated until the error is minimized.

288 The choice of hyperparameters to ensure the MLP model's robustness is processed by hyperparameter tuning [using
289 5-fold cross-validation as shown in section S1.2 of the supplement. Thanks to hyperparameter tuning, the two
290 hidden layers and a logistic sigmoid activation function were selected in this study to capture the non-linear
291 relationships between OP activities and PM₁₀ sources, and shown in section S1.2 of the supplement. Thanks to
292 hyperparameter tuning, the two hidden layers and a logistic sigmoid activation function were selected in this study
293 to capture the non-linear relationships between OP activities and PM sources.](#)

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

296 **2.6. Performance of the models**

297 The performance metrics R-square (R^2), mean absolute error (MAE), and root mean square error (RMSE) were
 298 used to assess the goodness of fit of models as described by Kuhn & Johnson (2013). R^2 quantifies the model's
 299 ability to explain the variance in the data. R^2 equal to 1 indicates a perfect fit. RMSE represents the aggregation of
 300 the individual differences between predicted OP and measured OP, while MAE assesses the average magnitude of
 301 errors between them. Lower RMSE and MAE values indicate a better fit, with a perfectly fitting model yielding
 302 an RMSE or MAE of 0. Eq132, Eq143, and Eq154, respectively, define R^2 , MAE, RMSE. These indicators are
 303 computed for the training and testing data of each sampling iteration and averaged across the 500 sampling
 304 iterations.

$$305 \quad R^2 = 1 - \frac{\text{Sum of Squared Residuals}}{\text{Total Sum of Squares}} = 1 - \frac{\sum_{i=0}^m (y_i - \hat{y}_i)^2}{\sum_{i=0}^m (y_i - \bar{y})^2} \quad (\text{Eq123})$$

$$306 \quad MAE = \frac{\sum_{i=0}^m |y_i - \hat{y}_i|}{m} \quad (\text{Eq14})$$

$$307 \quad RMSE = \sqrt{\frac{\sum_{i=0}^m (y_i - \hat{y}_i)^2}{m}} \quad (\text{Eq15})$$

308

309 3. Result and discussion

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

314 3.1. Dataset characteristics

315 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
 316 in Figure 3, illustrating variations in annual average OP activities and PM_{10} source contributions by sites. Most
 317 sites, including traffic and industrial ones, show higher OP_{DTT} activities than OP_{AA} . Conversely, for the alpine
 318 valley sites, CHAM presents higher OP_{AA} than OP_{DTT} , while GRE-fr experiences similar levels [of \$\text{OP}_{\text{AA}}\$ and](#)
 319 [\$\text{OP}_{\text{DTT}}\$ of the 2-OPs](#). Additionally, the average OP activities in every site are not proportional to the average PM
 320 concentration. For instance, CHAM and NIC had lower PM_{10} concentrations but higher OP activities than other
 321 sites, while TAL showed high PM_{10} concentrations but relatively lower OP activities.

322 The variations observed in the levels of PM_{10} and OP across six sites can be attributed to distinctions in identified
 323 sources and their respective contributions. These disparities are contingent upon the unique typologies of each site,
 324 which are discussed in Weber et al., 2021. Further, we can observe a significant seasonality in the OP activities
 325 (Table S.1). Strong seasonality of OP in Alpine valley sites has been addressed in previous studies (Borlaza et al.,
 326 2021; Dominutti et al., 2023; Weber et al., 2018, 2021), with thermal inversions during winter increasing pollutants
 327 concentrations and OP activities compared to summer. Conversely, OP activities in cold and warm periods in other
 328 sites are not significantly different.

329 The PM_{10} sources and their repartition vary among sites (Figure 3) because of the difference in typology and local
 330 activities. For instance, in the industrial site (PdB), two specific sources are identified: shipping emissions (HFO)
 331 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
 332 TAL also appears to be influenced by industrial sources ($2.34 \mu\text{g m}^{-3}$), which might, however, be partly due to
 333 biases induced by the application of the harmonized receptor model protocol (Weber et al., 2019). Note that the
 334 application of a site-specific PMF procedure for this site leads to a much lower contribution of this source category
 335 but relatively similar contributions of other sources (Favez, 2017). GRE-fr, an urban background site in an alpine

336 valley, presents significant long-range transport sources, with secondary sulfate contributing $3.90 \mu\text{g m}^{-3}$ followed
 337 by biomass burning at $2.21 \mu\text{g m}^{-3}$. As expected, biomass burning is an abundant source in CHAM, accounting for
 338 $7.28 \mu\text{g m}^{-3}$ of the PM contribution, while the traffic sites RBX and NIC displayed high contributions of traffic
 339 sources (at $2.43 \mu\text{g m}^{-3}$ and $1.45 \mu\text{g m}^{-3}$ respectively).

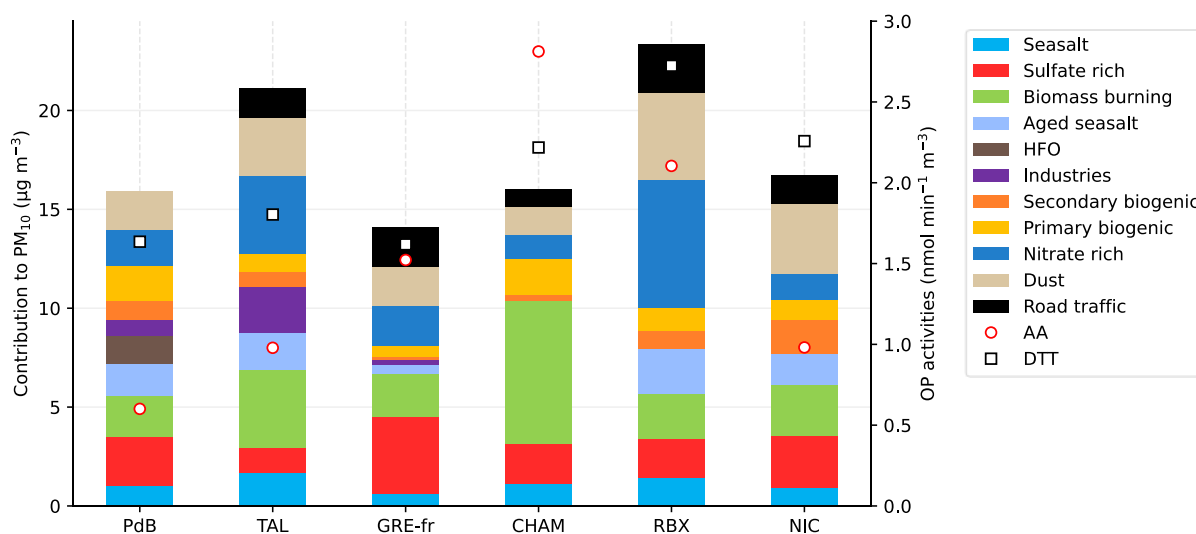
340 The presence of multicollinearity and homoscedasticity were tested to assess the data characteristic of every site.
 341 The only site with evidence of collinearity was NIC, where the VIF of the traffic source was equal to 5.0. For all
 342 other sites, VIF values are below 5, indicating limited collinearity among sources. This is expected, as the PMF
 343 analysis is constrained to avoid collinearity between sources. VIF values for each site can be found in Table S.2.

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

357 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

358

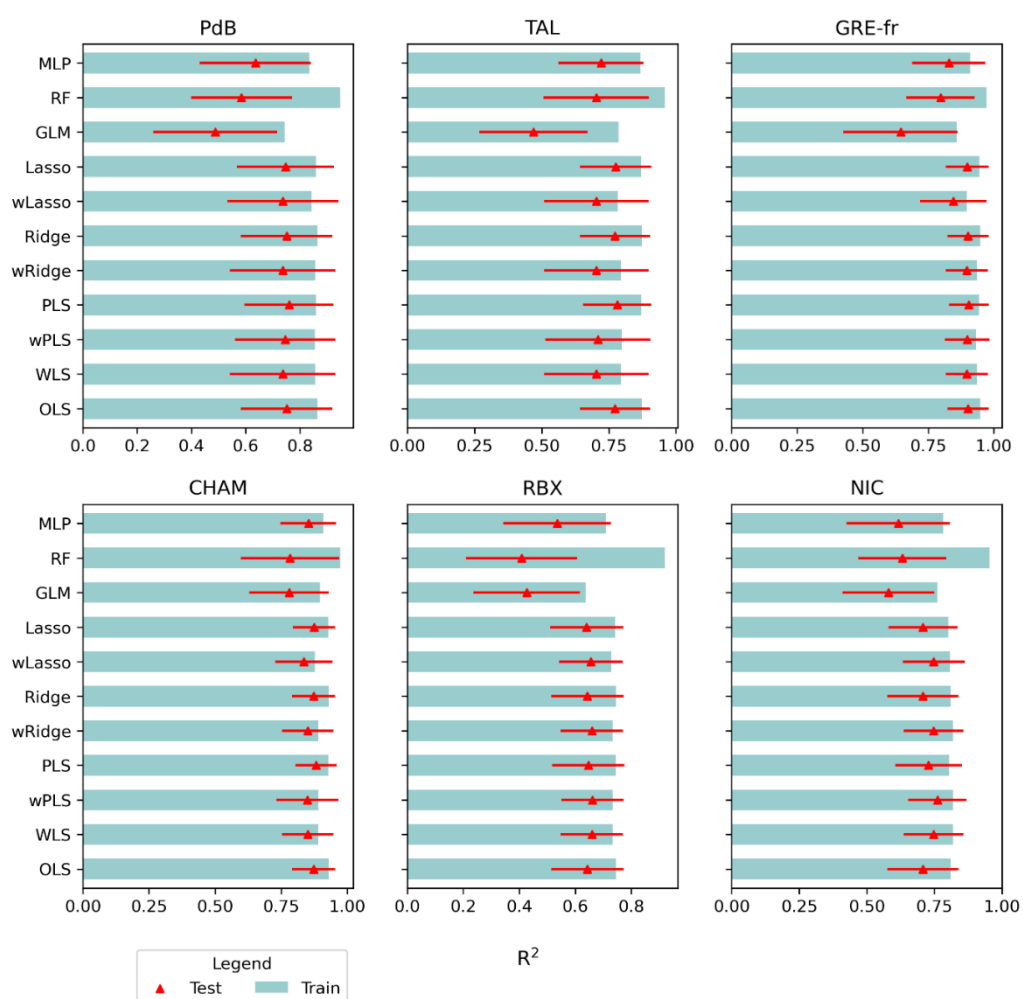


359

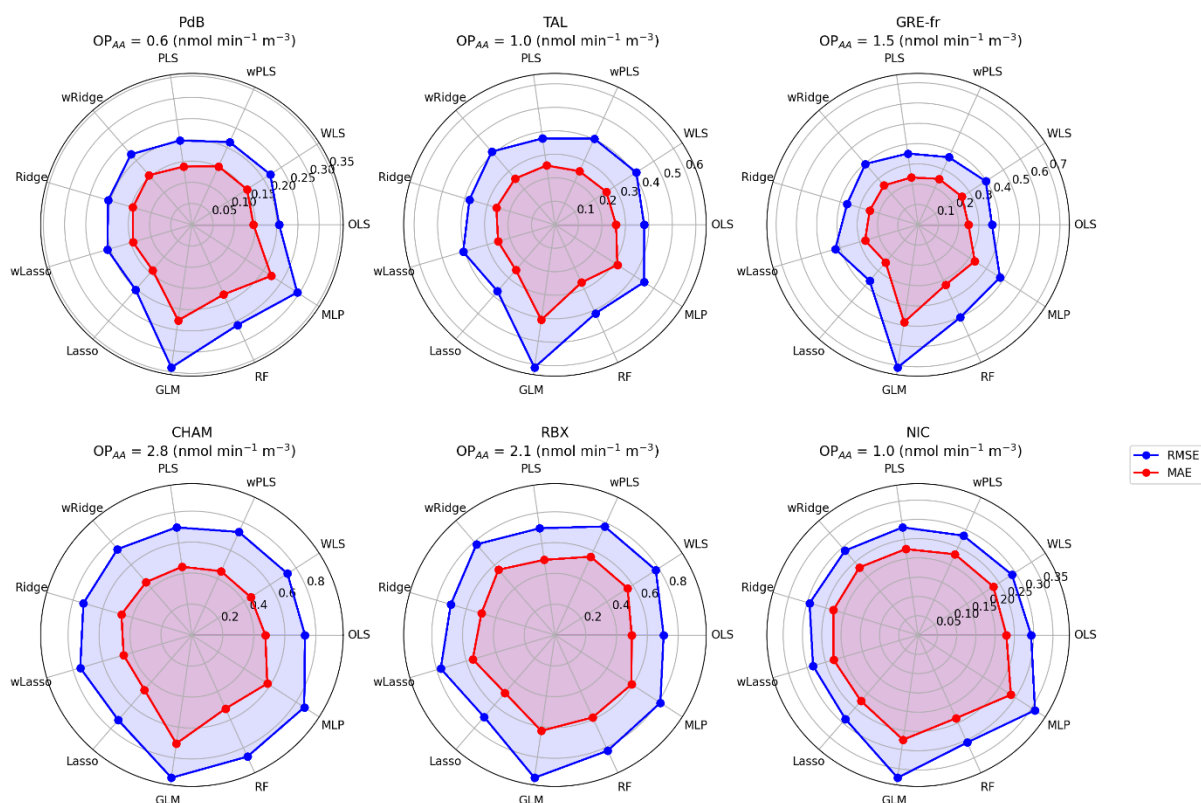
360 **Figure 3. The contribution of sources to PM₁₀ and the OP activities in 6 sites. The left y-axis and bar show**
 361 **the contribution of PM sources in $\mu\text{g m}^{-3}$. The right y-axis, circles and squares showed the mean OP_v**
 362 **activities in $\text{nmol min}^{-1} \text{m}^{-3}$, with red circle for OP_{AA} and black square for OP_{DTT}.**

363 3.2. The performances of regression models

364 The 11 regression models, with or without weighting for some of them, were tested by comparing their
 365 performance metrics between the measured and reconstructed OPs. For each run ($n = 500$ iterations), the R^2 ,
 366 RMSE, and MAE were computed for the testing and training dataset, resulting in 500 values for each performance
 367 metric. Figure 4 presents the mean R^2 values of the training data sets, the mean and the standard deviation of the
 368 testing datasets of the OP_{AA} models across the 500 sampling iterations, and Figure 5 presents the mean RMSE and
 369 MAE. The same result pattern was found for OP_{DTT}, as presented in the tables S.3, S.4, S.5, Supplement. The
 370 WLS, wPLS, wRidge, and wLasso models incorporated weighting, while the OLS, PLS, Ridge, Lasso, GLM, RF,
 371 and MLP models were unweighted.



372
 373 **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**
 374 **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**
 375 **data. The y-axis represents the models, and the x-axis denotes the R^2 values.**



377

378 **Figure 5. The MAE and RMSE of 11 OP_{AA} models in every site for the testing data. Blue and red lines**
 379 **present the RMSE and the MAE, respectively. The values in the figure are the mean of RMSE and MAE of**
 380 **500 iterations.**

381 OP predictions across all sites are statistically validated, with testing R^2 values observed in RBX, NIC, PdB, TAL,
 382 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
 383 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
 384 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
 385 exhibits for all sites the lowest R^2 values and the highest RMSE (Table S.3, S.4, S.5, Supplement). These results
 386 strongly suggest that the relationship between OP_{AA} and PM_{10} sources is not log-linear.

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

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

404 The WLS, OLS, wPLs, wRidge, and wLasso models show more robust performances with fewer differences
 405 between the training and testing data. At most sites, there is very little difference between the R^2 , RMSE, and MAE
 406 of OLS and Ridge, with or without weighting, and often PLS and Lasso as well. This consistency is observed even
 407 in the collinearity case of NIC, where $VIF = 5$. The difference between these models is a maximum of 0.06 in R^2 ,
 408 0.01 in MAE and 0.1 in RMSE, indicating that these models work well for OP prediction. Nevertheless, it is worth
 409 noting that every model exhibits different assumptions that have to be respected. The assumption violations may
 410 lead to unreliable regression coefficients (intrinsic OP) even though the prediction is good (Cohen et al., 2002;
 411 Williams et al., 2013).

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

421 **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

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

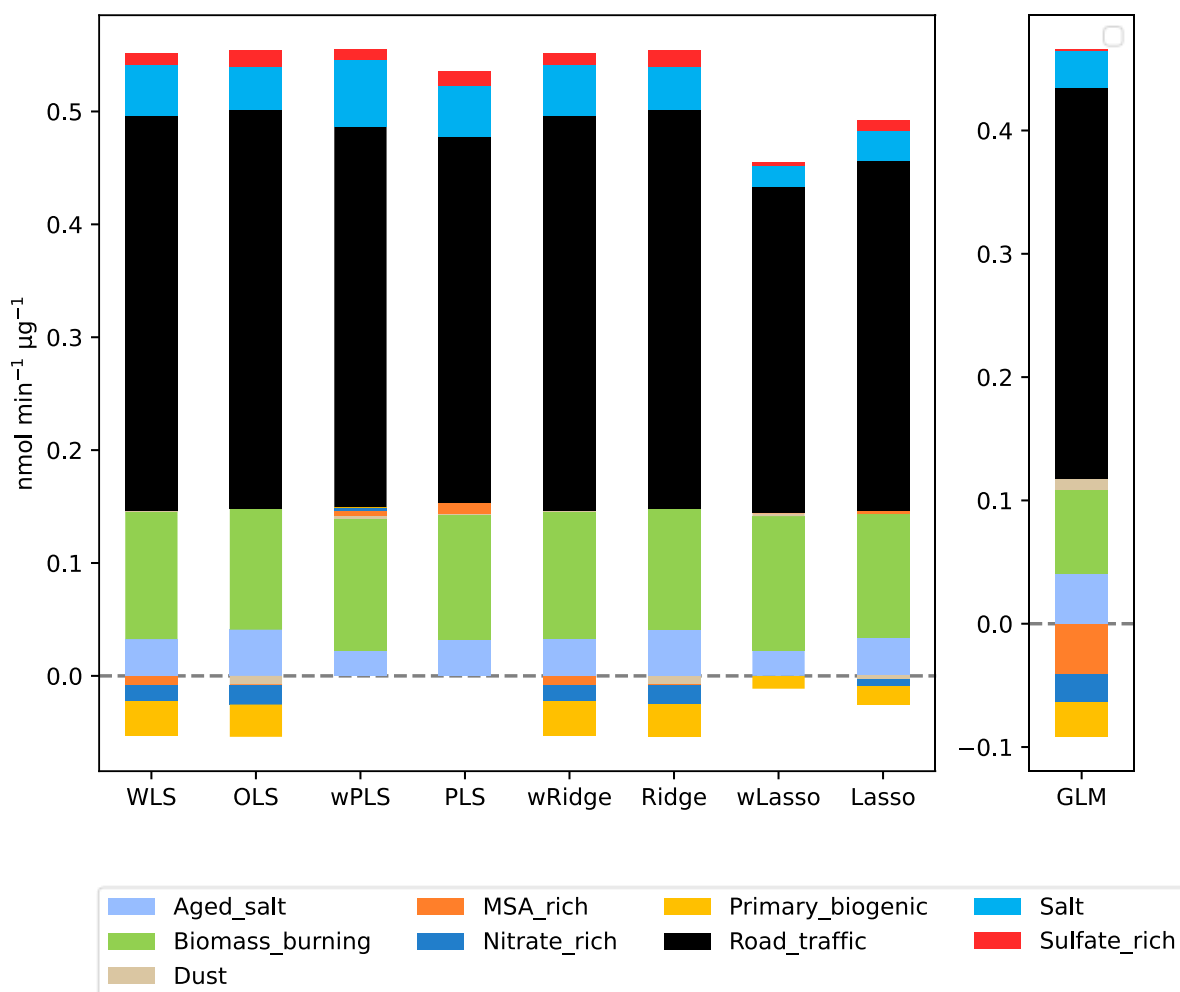
	PdB	TAL	GRE-fr	CHAM	RBX	NIC
Collinearity	No	No	No	No	No	Yes
Heteroscedasticity	No	Yes	No	Yes	No	No
Best model by characteristic	OLS/ PLS	WLS/ wPLS	OLS/ PLS	WLS/ wPLS	OLS/ PLS	Ridge/ Lasso
Best by error	OLS	wPLS	PLS	wPLS	PLS	Ridge

423

424 3.3. Effect of the choice of a model on intrinsic OP

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

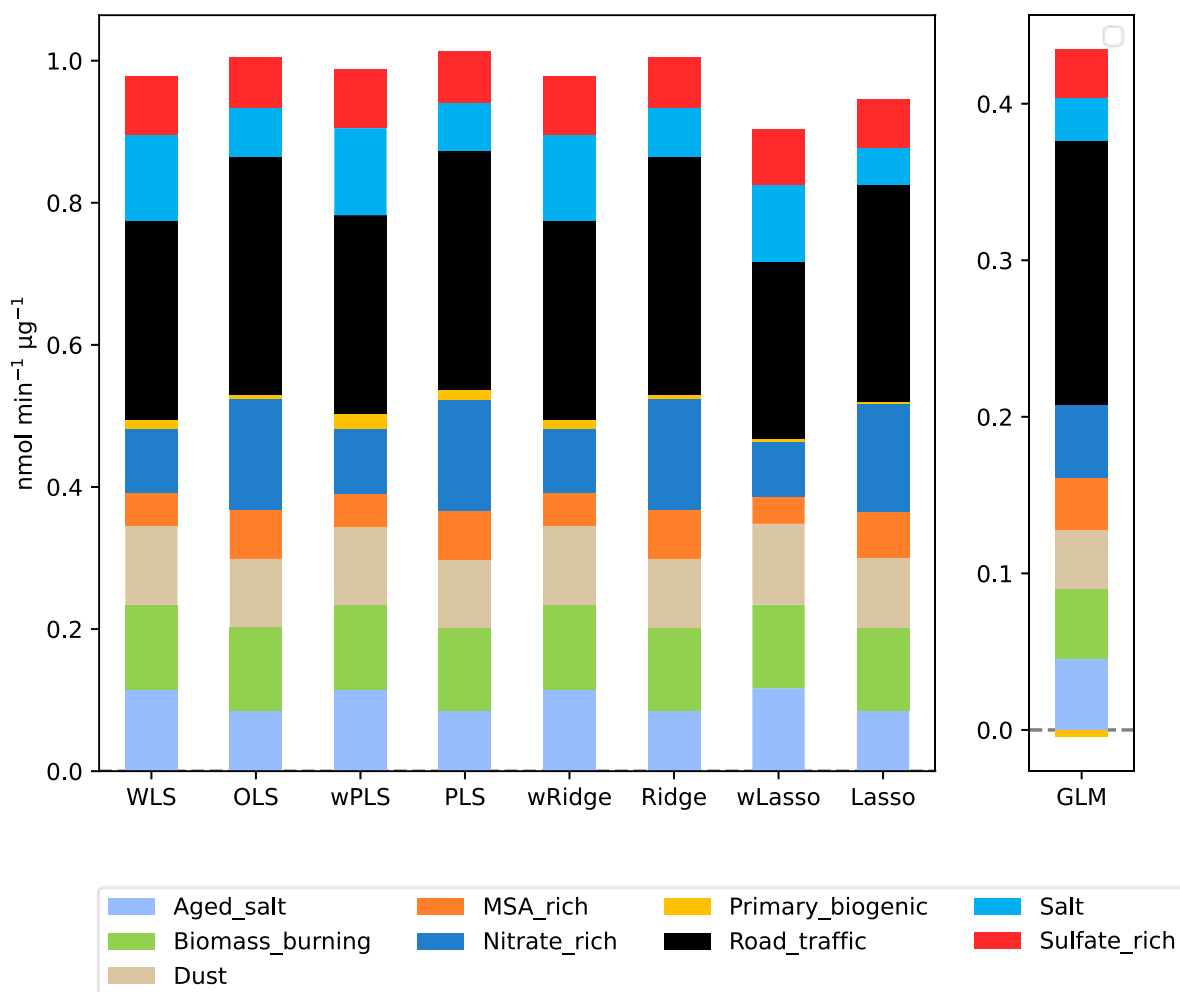
437 Intrinsic OP_{AA} of PM_{10} sources at NIC is the same between WLS and wRidge and between the OLS and Ridge,
438 revealing that the moderate collinearity of the road traffic source did not affect the estimated intrinsic OP_{AA} . PLS
439 sets the intrinsic OP_{AA} of some sources to zero, therefore producing slightly different results. Lasso regression sets
440 the intrinsic OP_{AA} of some sources to zero and shrinks the estimates for all other sources toward zero. GLM
441 produces intrinsic OP_{AA} values that represent a multiplicative change on the log scale, so they are not directly
442 comparable to the other models. However, the direction and importance of the sources are similar to the other
443 models. Whatever the model, road traffic appears as the source with the highest intrinsic OP_{AA} , followed by
444 biomass burning, aged salt, salt and sulfate-rich sources, in NIC. Traffic and biomass burning sources have been
445 similarly recognized as significant contributors to OP_{AA} in prior studies (Borlaza et al., 2021; Dominutti et al.,
446 2023; Stevanović et al., 2023). The intrinsic OP of the dominant sources is stable, indicating that all these models
447 could give the same information about the intrinsic OP of the main sources. Conversely, the differences are larger
448 between models for the sources with small to very small intrinsic OP (MSA rich, primary biogenic, nitrate-rich,
449 dust), whose intrinsic OP varies from positive to negative among models.



450

451 **Figure 6. Intrinsic OP_{AA} values of the different PM₁₀ sources at Nice were obtained with the different models.**

452 The OP_{DTT} intrinsic values in NIC (Figure 7) display minimal variation among the WLS, wPLS. This consistency
 453 is linked to the absence of negative intrinsic values. On the other hand, even though there is the presence of
 454 moderate collinearity, wRidge still has the same result as WLS and wPLS. In line with the OP_{AA} results, the wLasso
 455 and GLM models exhibit distinct responses compared to the other models. The intrinsic OP_{DTT} of all sources varies
 456 depending on the presence or absence of weighting. While the WLS models tend to amplify the influence of some
 457 sources (aged sea salt, primary biogenic, sea salt, and sulfate-rich), the OLS reduces the intrinsic OP_{DTT} of these
 458 sources. Conversely, MSA-rich, nitrate, and road traffic sources undergo less influence in WLS but higher in OLS.
 459 Different from OP_{AA}, OP_{DTT} prediction shows more variation among models, highlighting the effect of choosing
 460 a model on evaluating the intrinsic OP_{DTT} of PM₁₀ sources.



461

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

464 The comparison of intrinsic OP among regression models in NIC demonstrated that OP_{DTT} and OP_{AA} intrinsic
 465 values exhibit variation across different models with and without weighting, illustrating that the choice of the
 466 model significantly influences the values obtained for intrinsic OP of PM₁₀ sources (A similar pattern is observed
 467 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
 468 in $\frac{OP}{intrinsic_OP}$ across models, a comparison between the best-performing and most commonly used models
 469 (OLS) is presented in the following section to elucidate the advantage of choosing a model based on data
 470 characteristics (section 3.4).

471 3.4. Comparisons between the best site-specific model and OLS

472 In this section, the intrinsic OP of the best model is selected for each site as discussed in Section 3.2, and the
 473 intrinsic values of each source are compared to the ones returned by the OLS model. The OLS model is used as a
 474 representative of usual practices that do not consider the database characteristics_ (Williams et al., 2013). Each
 475 PM₁₀ source's average intrinsic OP value is calculated from all the 500 bootstrapping iterations for all sites where
 476 that particular source is identified. Intrinsic OP values obtained in this way from the best model ([the best model](#)
 477 [presented in Table 3 for OP_{AA} and Table 4 for OP_{DTT}](#)) encompassing all six sites are called **intrinsic OP of the**
 478 **best model**, and the intrinsic OP values derived from the OLS from all six sites are called **intrinsic OP of the**
 479 **reference model**.

480 A meaningful comparison of the two series of intrinsic values requires two conditions. First, intrinsic OP values
481 should be consistent across all sites. While recognizing that intrinsic OP values depend on diverse factors, we
482 assumed the sites share fairly uniform PM₁₀ chemical source profiles in France. This is demonstrated by evaluating
483 the Pearson distance and standardized identity distance similarity indicators of the source chemical profiles (Belis
484 et al., 2015; Weber et al., 2019), and Figure S.13 indicates consistent profiles of sources for the 6 sites.
485 Consequently, we could expect to observe minimal divergence in intrinsic OP values among these sites. Second,
486 we postulate that negative intrinsic OP values are possible since previous studies have reported that total PM₁₀
487 intrinsic OP can be modulated due to the synergetic/antagonistic effects involving, for example, soluble copper,
488 quinones, and bacteria (Borlaza et al., 2021; Pietrogrande et al., 2022; Samake et al., 2017; S. Wang et al., 2018;
489 Xiong et al., 2017). Samake et al. (2017) demonstrated that the presence of bacterial cells in aerosol decreases the
490 redox activity of Cu and 1,4-naphthoquinone, with a maximum decreasing of 60% compared to the oxidative
491 reactivity considered individually. Pietrogrande et al. (2022) indicated that the mixture of Cu, Fe, 9,10-
492 phenanthrene quinone and 1,2-naphthoquinone reduces the rate consumption of AA and DTT, up to 50%
493 depending on the quantity of each chemical. Wang et al. (2018) reported that the mixing of Cu and naphthalene
494 secondary organic aerosol (SOA) and phenanthrene SOA only got half of DTT rate consumption compared to the
495 consumption when considered separately. Xiong et al. (2017) showed the presence of antagonists in the interaction
496 of Fe and quinones, nevertheless, much lower than those in the other studies (under 10%). These references
497 reported that the antagonistic effects of a mixture can significantly reduce the consumption rate of OP_{DTT} and
498 OP_{AA}, and this impact varies widely from 10% to 60% depending on the type of chemical species and the quantity
499 of each species in the mixture. These last studies showed that the impact of synergistic and antagonistic effects
500 cannot exceed 60% of the intrinsic OP value when assessed independently for each chemical. Consequently, we
501 consider here that the intrinsic OP value of an individual site for a given source could be negative only within a
502 range of at most 60% of the mean combined intrinsic OP value of this source across all sites. Negative intrinsic
503 OP exceeding this criterion may result from the mathematical construction of the model. The comparison of
504 intrinsic OP_{AA} of the best and reference model is presented in 3.4.1 and that of OP_{DTT} is shown in 3.4.2.

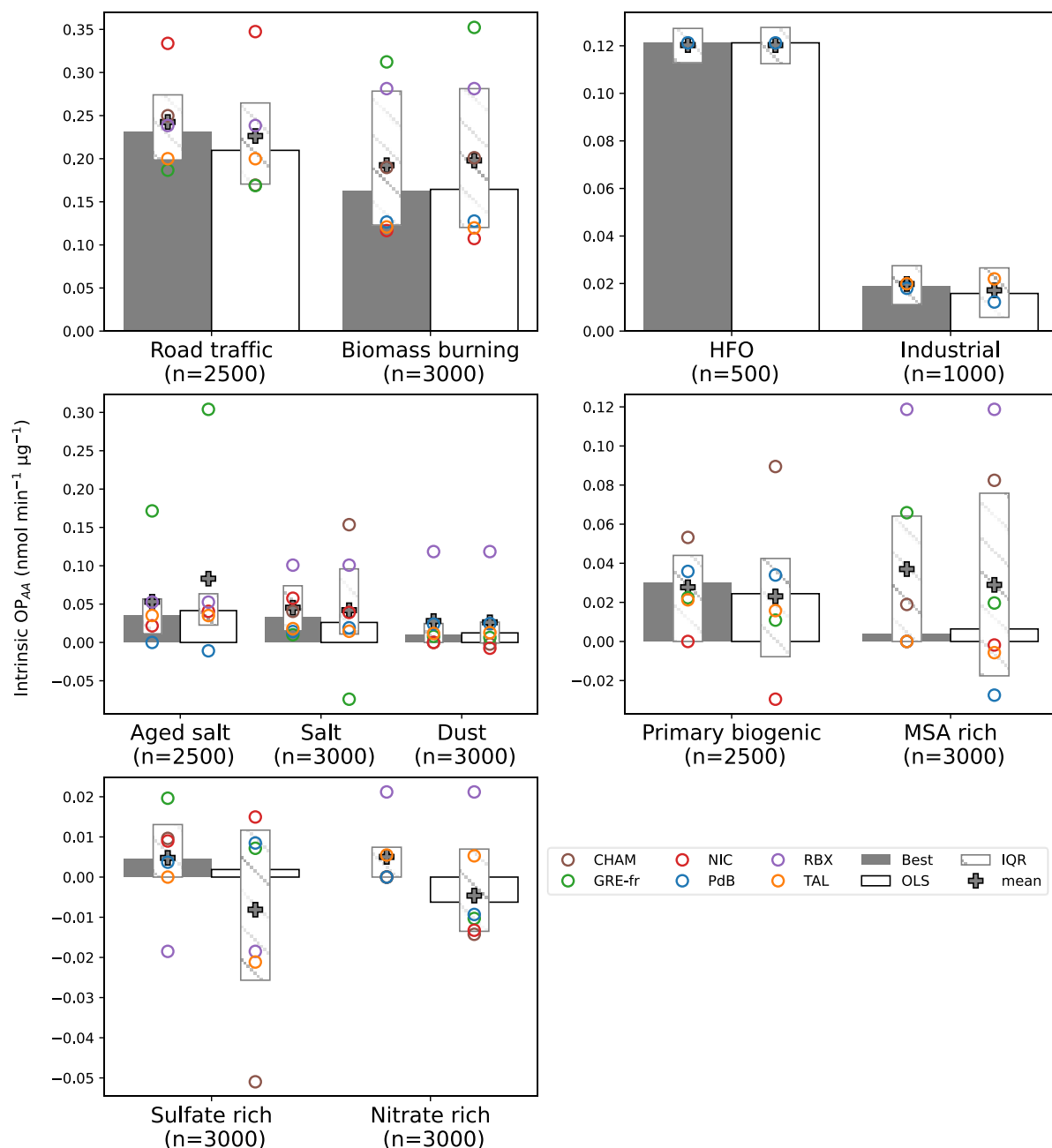
505 3.4.1. OP_{AA} activities

506 The results of the comparison of OP_{AA} intrinsic values (Figure 8 and Table S.8) show that the anthropogenic
507 sources get the highest intrinsic OP values in both the best and reference models. Among these sources, road traffic
508 appears as the most prominent potent fraction, followed by biomass burning, HFO, and industrial. These results
509 are aligned with prior research (Calas et al., 2019; Daellenbach et al., 2020; Dominutti et al., 2023; Fadel et al.,
510 2023; Fang et al., 2016; in 't Veld et al., 2023; Weber et al., 2018; Zhang et al., 2020) which has highlighted the
511 sensitivity of OP_{AA} to concentrations of metals, black carbon, and organic carbon. The differences between the
512 best and reference models were insignificant for these sources, demonstrating that **the best and reference models**
513 **consistently captured similar patterns for the most critical sources of OP activities.**

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

523 When considering sources with low intrinsic OP, the variability can be larger between the two methods. As an
524 example, for the sulfate-rich sources, the median intrinsic OP values were positive (0.002 nmol min⁻¹ μg⁻¹), while

525 the mean intrinsic OP values were negative ($-0.008 \text{ nmol min}^{-1} \mu\text{g}^{-1}$). The mean intrinsic OP in the best model
526 exhibited fewer negative values in individual sites than in the reference model (for aged salt, salt, primary biogenic,
527 MSA rich, sulfate-rich and nitrate-rich). In addition, the best model showed the less disparate intrinsic OP among
528 individual sites; for instance, the aged salt sources in GRE-fr and the primary biogenic and salt sources in CHAM,
529 highlighting the advantage of considering the data in model selection. For example, Furthermore, the best model
530 displayed an intrinsic OP meaningful in terms of geochemicals, which showed in the source of salt, primary
531 biogenic, sulfate-rich. For instance, in the reference model, The mean-average intrinsic OP values of the primary
532 biogenic source revealed a negative intrinsic OP in NIC ($-0.03 \text{ nmol min}^{-1} \mu\text{g}^{-1}$), the intrinsic OP of salt in GRE-
533 ft ($-0.07 \text{ nmol min}^{-1} \mu\text{g}^{-1}$) as well as the sulfate-rich source in CHAM ($-0.05 \text{ nmol min}^{-1} \mu\text{g}^{-1}$). This negative value
534 represented represented a 100% reduction compared to the mean intrinsic OP of all sites. Moreover, In the OLS
535 model, the negative intrinsic OP was observed in NIC (Primary biogenic), and some extreme values in GRE-fr
536 (aged salt, salt), CHAM (salt, primary biogenic, MSA-rich), NIC (where heteroscedasticity was presented) in the
537 OLS model, underscores that the model assumptions on data characteristics proving false could impact the
538 accuracy of OP prediction. Consequently, these results highlight the advantage of considering the data in model
539 selection.



540

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

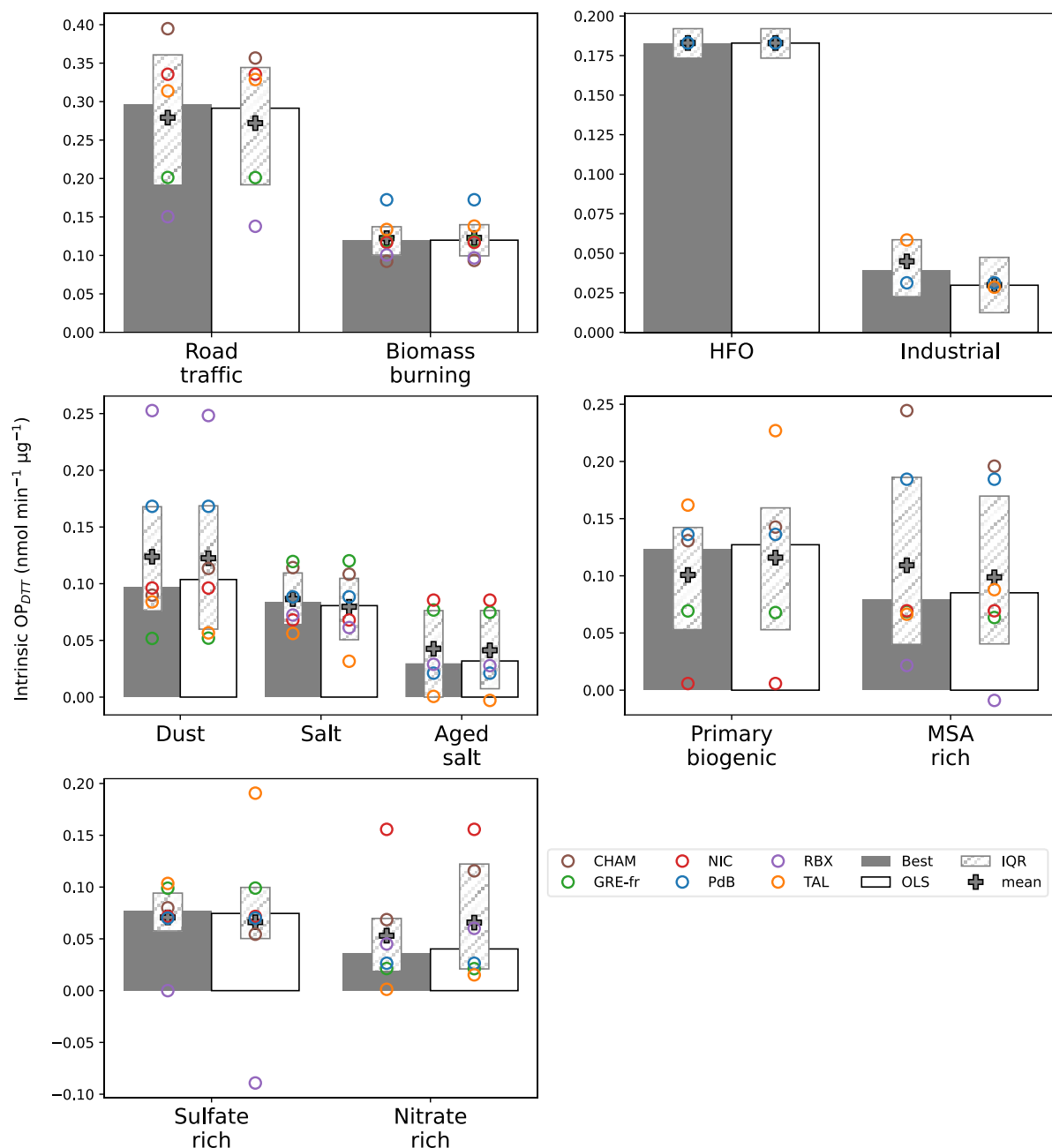
549 The detailed comparison of intrinsic OP_{AA} between the best and reference models is categorized into four groups
 550 and discussed in detail in section S9. These groups include (1) anthropogenic sources without nitrate and sulfate
 551 (road traffic, biomass burning, HFO, industrial), (2) natural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

553 3.4.2. OP_{DTT} activities

554 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
555 reference model (OLS). Except for the road traffic, industrial, and MSA-rich, the IQR is slightly higher in the best
556 model (Figure 9 and Table S.9). In the two models, the mean intrinsic OP is essentially unchanged, where the
557 traffic is the most critical source (0.27 ± 0.10), followed by HFO (0.18 ± 0.01), biomass burning (0.12 ± 0.03), dust
558 (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:
559 0.09 ± 0.09). ~~The remaining sources, such as sea salt, sulfate rich, industrial, and nitrate rich, show a negligible
560 contribution to OP_{DTT} with an intrinsic OP_{DTT} from 0.02 to 0.08.~~ The minimum difference between the two models
561 ~~in the dominant sources~~ again confirms the conclusion in the OP_{AA} comparison, demonstrating **the similar pattern
562 of the best and the reference model in the most crucial sources of OP**. For both best and reference, OP_{DTT}
563 activities showed sensitivity to more sources than OP_{AA} , as discussed in previous studies (Borlaza et al., 2021;
564 Calas et al., 2019; Dominutti et al., 2023; Fadel et al., 2023). ~~The traffic, HFO and biomass burning sources
565 highlighted in are the most contributing to OP_{DTT} activities. The primary biogenic highly contributes to OP_{DTT} in
566 both models, likely reflecting the sensitivity of OP_{DTT} to organic compounds, as mentioned in . The intrinsic OP
567 of dust and MSA rich have been shown to vary in the literature, indicating the different effects of the compositions
568 to generate DTT ROS.~~

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

582 ~~For the sulfate rich source, the reference model showed a negative intrinsic OP_{DTT} in RBX ($-0.09 \text{ nmol min}^{-1} \mu\text{g}^{-1}$)
583 $^{-1}$), while the best model showed an intrinsic of 0. On the other hand, the reference model presents that the intrinsic
584 OP in TAL is $0.20 \text{ nmol min}^{-1} \mu\text{g}^{-1}$, far from the mean of all sites ($0.07 \text{ nmol min}^{-1} \mu\text{g}^{-1}$). The best model,
585 conversely, shows a more consistent intrinsic OP in TAL compared to the other sites. A similar result was also
586 found in primary biogenic sources, where the reference model overestimates the intrinsic OP of this source in TAL
587 compared to the other sites. The reason is that heteroscedasticity was detected in TAL, which does not satisfy the
588 assumption of OLS.~~



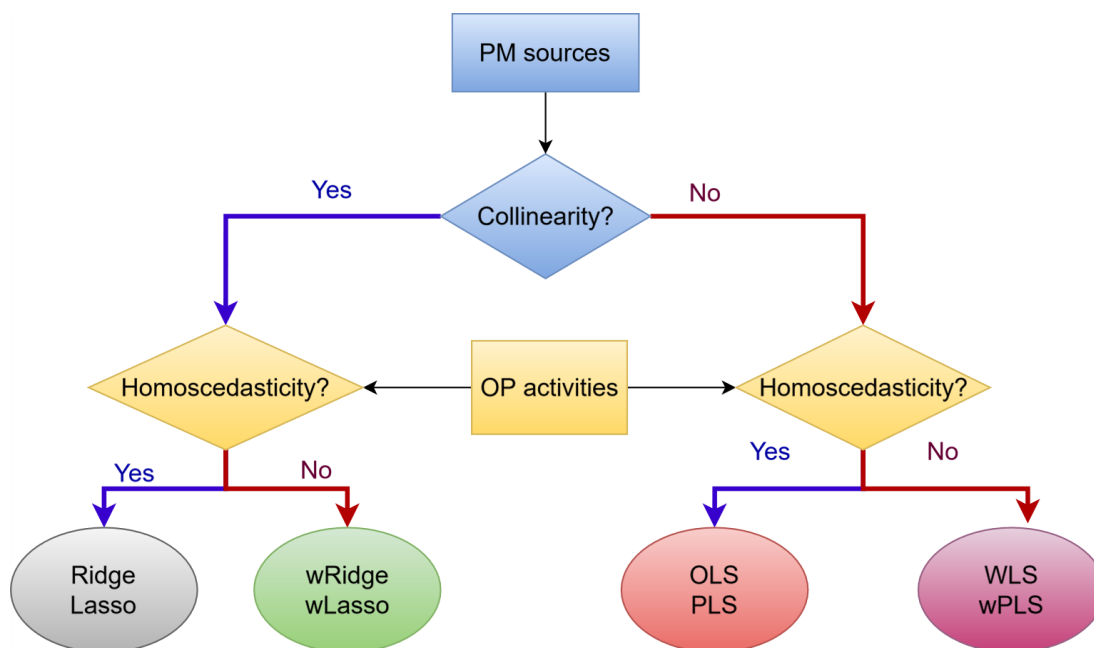
589

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

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

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

603 **3.5. Guidelines for the selection of regression model for OP SA.**



604

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

606

607 Our results have highlighted the benefits of choosing a model that matches the characteristics of the data to improve
608 the robustness of OP SA method. For this reason, this section develops a workflow to help make model selection
609 decisions. Before selecting a regression for OP SA, the first question is whether the PM_{10} sources are collinear and
610 the second is whether the residual variance of the regression between OP and PM_{10} mass is constant. These two
611 questions represent the characteristics of PM_{10} sources and OP activities, which vary according to the study site.

612 For data exhibiting collinearity between sources and generating a residual variance that varies according to the
613 value of the PM_{10} sources, weighted regularisation regression can help to reduce collinearity and to match the
614 model assumption about the residual. On the other hand, the unweighted Ridge and Lasso are introduced for data
615 showing collinearity and homoscedasticity. Additionally, data with no collinearity are suitable for OLS and
616 unweighted PLS in the case of homoscedasticity, while WLS, weighted PLS are used for data with
617 heteroscedasticity.

618 If the number of predictors (PM_{10} sources) is below the number of samples divided by 15, RF and MLP can also
619 be employed to capture possible non-linear relationships between the OP and PM_{10} sources. However, cross-
620 validation must be used to ensure that there is no over-fitting. In addition, these models do not estimate intrinsic
621 OP ($nmol\ min^{-1}\mu g^{-1}$) but only the importance of each PM_{10} source to the OP prediction. This is a large drawback
622 since the intrinsic OP of sources is a must for the modelling effort of OP with CTM. However, RF and MLP could
623 be useful for OP prediction in the case of larger datasets generated by online instruments.

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

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

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

638

639 Limitations and perspectives of the study:

- 640 - This study compares eight regression models but is not exhaustive; further research could add more
641 regression techniques to evaluate result variations across models. The potential techniques that could be
642 applied for OP SA are gradient boosting techniques for resolving regression models, or supervised
643 machine learning techniques which allows the investigation of linear and non-linear regression
644 relationships. However, the consistently strong performance of ordinary linear regression across six
645 locations in France suggests that there may be little to gain from applying more complex models in areas
646 with similar PM₁₀ sources.
- 647 - PMF coupled with a regression model remains a popular approach for OP SA. Notably, the uncertainties
648 in PMF are typically addressed in chemical profiles, but not in contributions. Incorporating uncertainty
649 from variations in contribution into models could enhance their robustness compared to relying only on
650 absolute PMF results.
- 651 - ~~Observations ranged between 100 and 200 samples at each site, which may be insufficient to obtain fair~~
652 ~~performance of GLM, decision trees and neural network models. Such a number of samples is sufficient~~
653 ~~to address SA through PMF model for offline analyses. Therefore, such study outlines well the limitations~~
654 ~~of GLM, RF, MLP for such types of datasets. Future investigations should be performed in an extended~~
655 ~~dataset, such as long term or real time measurement data, to investigate the performance of such machine~~
656 ~~learning algorithms.~~ Observations ranged between 100 and 200 samples at each site, which may be
657 insufficient to obtain a fair performance of GLM, decision trees and neural network models even though
658 this number of samples is sufficient to address SA through the PMF model for offline analyses. Therefore,
659 this study outlines well the limitations of GLM, RF, and MLP for offline datasets. Future investigations
660 should be performed in an extended dataset, such as long-term or real-time measurement data, to
661 investigate the performance of machine learning algorithms.
- 662 - This study only focused on the two most popular OP assays of PM₁₀ (OP_{DTT} and OP_{AA}). However, there
663 are actually various OP assays, such as OP_{DCFH}, OP_{OH}, OP_{FOX}, OP_{GSH}, OP_{ESR} and different sizes of PM
664 (PM₁, PM_{2.5}, PM₅). Further research should include more OP assays, which can be helpful in evaluating
665 the performance of various regression models for different OP and different PM sizes.
- 666 - This study used the analytical uncertainty as the weighting for the weighted model. However, the
667 weighting can be selected based on different ways, as reported by Montgomery et al. (2012): (1) Prior
668 information from the theoretical model, (2) Using the residual extracted from the OLS model, (3) The
669 selecting of weighting based on the uncertainty of instrument if the dependent variable measured by a
670 different method and (4) If the dependent variable is the average of different observations, the weighting
671 selected based on the error of these observations.

673 **4. Conclusion**

674 The results of the OP SA marked an important milestone as they were revealed for the first time through the use
675 of eight regression models, including OLS, WLS, PLS, GLM, Ridge, Lasso, RF and MLP. This in-depth analysis
676 was carried out on a complete set of data collected from six sites with different characteristics. The approach of
677 selecting a suitable model for each site based on specific data characteristics resulted in a more consistent intrinsic
678 OP across sites, in stark contrast to the variation observed when using the basic OLS model. The revelations of the
679 study have provided concrete recommendations for the judicious selection of an appropriate regression model
680 based on the unique characteristics of the dataset. These guidelines should help to improve the accuracy of OP
681 assessments and contribute to the refinement of air quality assessment methods. In addition, the implications of
682 this research extend to the implementation of OP monitoring as a new measure of air quality, particularly on
683 European supersites. As this initiative aligns with the ongoing revision process of the European Directive
684 2008/50/CE, the study's findings assume a pivotal role in shaping the methodologies underpinning air quality
685 assessments at a broader regulatory level.

686 **Code availability**

687 The software code could be made available by contacting the corresponding author upon request.

688 **Data availability**

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

690 **Author contributions**

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

695 **Competing interests**

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

697 **Acknowledgments**

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

704 **Financial support**

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

712 **Reference**

- 713 Akhtar, A., Islamia, J. M., Masood, S., Islamia, J. M., Masood, A., & Islamia, J. M. (2018). *Prediction and Analysis*
714 *of Pollution Levels in Delhi Using Multilayer Perceptron*. June. <https://doi.org/10.1007/978-981-10-3223-3>
- 715 Akhtar, McWhinney, R. D., Rastogi, N., Abbatt, J. P. D., Evans, G. J., & Scott, J. A. (2010). Cytotoxic and
716 proinflammatory effects of ambient and source-related particulate matter (PM) in relation to the production
717 of reactive oxygen species (ROS) and cytokine adsorption by particles. *Inhalation Toxicology*, 22(SUPPL.
718 2), 37–47. <https://doi.org/10.3109/08958378.2010.518377>
- 719 Alleman, L. Y., Lamaison, L., Perdrix, E., Robache, A., & Galloo, J. C. (2010). PM10 metal concentrations and
720 source identification using positive matrix factorization and wind sectoring in a French industrial zone.
721 *Atmospheric Research*, 96(4), 612–625. <https://doi.org/10.1016/j.atmosres.2010.02.008>
- 722 Ayres, J. G., Borm, P., Cassee, F. R., Castranova, V., Donaldson, K., Ghio, A., Harrison, R. M., Hider, R., Kelly,
723 F., Kooter, I. M., Marano, F., Maynard, R. L., Mudway, I., Nel, A., Sioutas, C., Smith, S., Baeza-Squiban,
724 A., Cho, A., Duggan, S., & Froines, J. (2008). Evaluating the toxicity of airborne particulate matter and
725 nanoparticles by measuring oxidative stress potential - A workshop report and consensus statement.
726 *Inhalation Toxicology*, 20(1), 75–99. <https://doi.org/10.1080/08958370701665517>
- 727 Bates, J. T., Fang, T., Verma, V., Zeng, L., Weber, R. J., Tolbert, P. E., Abrams, J. Y., Sarnat, S. E., Klein, M.,
728 Mulholland, J. A., & Russell, A. G. (2019). Review of Acellular Assays of Ambient Particulate Matter
729 Oxidative Potential: Methods and Relationships with Composition, Sources, and Health Effects.
730 *Environmental Science and Technology*, 53(8), 4003–4019. <https://doi.org/10.1021/acs.est.8b03430>
- 731 Bates, J. T., Weber, R. J., Abrams, J., Verma, V., Fang, T., Klein, M., Strickland, M. J., Sarnat, S. E., Chang, H.
732 H., Mulholland, J. A., Tolbert, P. E., & Russell, A. G. (2015). Reactive Oxygen Species Generation Linked
733 to Sources of Atmospheric Particulate Matter and Cardiorespiratory Effects. *Environmental Science and*
734 *Technology*, 49(22), 13605–13612. <https://doi.org/10.1021/acs.est.5b02967>
- 735 Bates, J. T., Weber, R. J., Verma, V., Fang, T., Ivey, C., Liu, C., Sarnat, S. E., Chang, H. H., Mulholland, J. A., &
736 Russell, A. (2018). Source impact modeling of spatiotemporal trends in PM2.5 oxidative potential across
737 the eastern United States. *Atmospheric Environment*, 193(August), 158–167.
738 <https://doi.org/10.1016/j.atmosenv.2018.08.055>
- 739 Beelen, R., Stafoggia, M., Raaschou-Nielsen, O., Andersen, Z. J., Xun, W. W., Katsouyanni, K., Dimakopoulou,
740 K., Brunekreef, B., Weinmayr, G., Hoffmann, B., Wolf, K., Samoli, E., Houthuijs, D., Nieuwenhuijsen, M.,
741 Oudin, A., Forsberg, B., Olsson, D., Salomaa, V., Lanki, T., ... Hoek, G. (2014). Long-term exposure to air
742 pollution and cardiovascular mortality: An analysis of 22 European cohorts. *Epidemiology*, 25(3), 368–378.
743 <https://doi.org/10.1097/EDE.0000000000000076>
- 744 Belis, C. A., Karagulian, F., Amato, F., Almeida, M., Artaxo, P., Beddows, D. C. S., Bernardoni, V., Bove, M. C.,
745 Carbone, S., Cesari, D., Contini, D., Cuccia, E., Diapouli, E., Eleftheriadis, K., Favez, O., El Haddad, I.,
746 Harrison, R. M., Hellebust, S., Hovorka, J., ... Hopke, P. K. (2015). A new methodology to assess the
747 performance and uncertainty of source apportionment models II: The results of two European
748 intercomparison exercises. *Atmospheric Environment*, 123, 240–250.
749 <https://doi.org/10.1016/j.atmosenv.2015.10.068>
- 750 Belis, C. A., Karagulian, F., Larsen, B. R., & Hopke, P. K. (2013). Critical review and meta-analysis of ambient
751 particulate matter source apportionment using receptor models in Europe. In *Atmospheric Environment* (Vol.
752 69, pp. 94–108). <https://doi.org/10.1016/j.atmosenv.2012.11.009>
- 753 Bell, M. L., Samet, J. M., & Dominici, F. (2004). Time-series studies of particulate matter. *Annual Review of*
754 *Public Health*, 25, 247–280. <https://doi.org/10.1146/annurev.publhealth.25.102802.124329>
- 755 Benkendorf, D. J., & Hawkins, C. P. (2020). Effects of sample size and network depth on a deep learning approach
756 to species distribution modeling. *Ecological Informatics*, 60(February).
757 <https://doi.org/10.1016/j.ecoinf.2020.101137>
- 758 Borlaza. (2021). Disparities in particulate matter (PM10) origins and oxidative potential at a city scale (Grenoble,
759 France) - Part 2: Sources of PM10 oxidative potential using multiple linear regression analysis and the
760 predictive applicability of multilayer perceptron n. *Atmospheric Chemistry and Physics*, 21(12), 9719–9739.
761 <https://doi.org/10.5194/acp-21-9719-2021>
- 762 Borlaza, L., Weber, S., Jaffrezo, J. L., Houdier, S., Slama, R., Rieux, C., Albinet, A., Micallef, S., Trébluchon, C.,

- 763 & Uzu, G. (2021). Disparities in particulate matter (PM10) origins and oxidative potential at a city scale
764 (Grenoble, France) - Part 2: Sources of PM10 oxidative potential using multiple linear regression analysis
765 and the predictive applicability of multilayer perceptron n. *Atmospheric Chemistry and Physics*, 21(12),
766 9719–9739. <https://doi.org/10.5194/acp-21-9719-2021>
- 767 Borlaza, L., Weber, S., Uzu, G., Jacob, V., Cañete, T., Micallef, S., Trébuchon, C., Slama, R., Favez, O., &
768 Jaffrezo, J.-L. (2021). Disparities in particulate matter (PM10) origins and oxidative potential at a city scale
769 (Grenoble, France) - Part 1: Source apportionment at three neighbouring sites. *Atmospheric Chemistry and*
770 *Physics*, 21(7), 5415–5437. <https://doi.org/10.5194/acp-21-5415-2021>
- 771 Bourlard, H., & Wellekens, C. J. (1989). Speech pattern discrimination and multilayer perceptrons. *Computer*
772 *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)
- 773 Breiman, L. (2001). RFRSF: Employee Turnover Prediction Based on Random Forests and Survival Analysis.
774 *Machine Learning*, 12343 LNCS, 503–515. https://doi.org/10.1007/978-3-030-62008-0_35
- 775 Brook, R. D., Rajagopalan, S., Pope, C. A., Brook, J. R., Bhatnagar, A., Diez-Roux, A. V., Holguin, F., Hong, Y.,
776 Luepker, R. V., Mittleman, M. A., Peters, A., Siscovick, D., Smith, S. C., Whitsel, L., & Kaufman, J. D.
777 (2010). Particulate matter air pollution and cardiovascular disease: An update to the scientific statement from
778 the american heart association. *Circulation*, 121(21), 2331–2378.
779 <https://doi.org/10.1161/CIR.0b013e3181d8ce1>
- 780 Brown, S. G., Eberly, S., Paatero, P., & Norris, G. A. (2015). Methods for estimating uncertainty in PMF solutions:
781 Examples with ambient air and water quality data and guidance on reporting PMF results. *Science of the*
782 *Total Environment*, 518–519, 626–635. <https://doi.org/10.1016/j.scitotenv.2015.01.022>
- 783 Calas, A., Uzu, G., Besombes, J. L., Martins, J. M. F., Redaelli, M., Weber, S., Charron, A., Albinet, A., Chevrier,
784 F., Brulfert, G., Mesbah, B., Favez, O., & Jaffrezo, J. L. (2019). Seasonal variations and chemical predictors
785 of oxidative potential (OP) of particulate matter (PM), for seven urban French sites. *Atmosphere*, 10(11).
786 <https://doi.org/10.3390/atmos10110698>
- 787 Calas, A., Uzu, G., Kelly, F. J., Houdier, S., Martins, J. M. F., Thomas, F., Molton, F., Charron, A., Dunster, C.,
788 Oliete, A., Jacob, V., Besombes, J. L., Chevrier, F., & Jaffrezo, J. L. (2018). Comparison between five
789 acellular oxidative potential measurement assays performed with detailed chemistry on PM10 samples from
790 the city of Chamonix (France). *Atmospheric Chemistry and Physics*, 18(11), 7863–7875.
791 <https://doi.org/10.5194/acp-18-7863-2018>
- 792 Calas, A., Uzu, G., Martins, J. M. F., Voisin, Di., Spadini, L., Lacroix, T., & Jaffrezo, J. L. (2017). The importance
793 of simulated lung fluid (SLF) extractions for a more relevant evaluation of the oxidative potential of
794 particulate matter. *Scientific Reports*, 7(1), 1–12. <https://doi.org/10.1038/s41598-017-11979-3>
- 795 Chianese, E., Camastra, F., & Ciaramella, A. (2018). *Spatio-temporal learning in predicting ambient particulate*
796 *matter concentration by multi-layer perceptron Spatio-temporal Learning in Predicting Ambient Particulate*
797 *Matter Concentration by Multi-Layer*. December. <https://doi.org/10.1016/j.econinf.2018.12.001>
- 798 Cho, A., Sioutas, C., Miguel, A. H., Kumagai, Y., Schmitz, D. A., Singh, M., Eiguren-Fernandez, A., & Froines,
799 J. R. (2005). Redox activity of airborne particulate matter at different sites in the Los Angeles Basin.
800 *Environmental Research*, 99(1), 40–47. <https://doi.org/10.1016/j.envres.2005.01.003>
- 801 Cohen, J., Cohen, P., West, S. G., & Aiken, L. S. (2002). *Applied multiple regression/correlation analysis for the*
802 *behavioral sciences*. Routledge. [https://doi.org/https://doi-org.sid2nomade-](https://doi.org/https://doi-org.sid2nomade-1.grenet.fr/10.4324/9780203774441)
803 [1.grenet.fr/10.4324/9780203774441](https://doi.org/https://doi-org.sid2nomade-1.grenet.fr/10.4324/9780203774441)
- 804 Craney, T. A., & Surles, J. G. (2002). Model-dependent variance inflation factor cutoff values. *Quality*
805 *Engineering*, 14(3), 391–403. <https://doi.org/10.1081/QEN-120001878>
- 806 Crobeddu, B., Aragao-Santiago, L., Bui, L. C., Boland, S., & Baeza Squiban, A. (2017). Oxidative potential of
807 particulate matter 2.5 as predictive indicator of cellular stress. *Environmental Pollution*, 230, 125–133.
808 <https://doi.org/10.1016/j.envpol.2017.06.051>
- 809 Crouse, D. L., Peters, P. A., Hystad, P., Brook, J. R., van Donkelaar, A., Martin, R. V., Villeneuve, P. J., Jerrett,
810 M., Goldberg, M. S., Arden Pope, C., Brauer, M., Brook, R. D., Robichaud, A., Menard, R., & Burnett, R.
811 T. (2015). Ambient PM2.5, O3, and NO2 exposures and associations with mortality over 16 years of follow-
812 up in the canadian census health and environment cohort (CanCHEC). *Environmental Health Perspectives*,
813 123(11), 1180–1186. <https://doi.org/10.1289/ehp.1409276>

- 814 Crouse, D. L., Peters, P. A., van Donkelaar, A., Goldberg, M. S., Villeneuve, P. J., Brion, O., Khan, S., Atari, D.
815 O., Jerrett, M., Pope, C. A., Brauer, M., Brook, J. R., Martin, R. V., Stieb, D., & Burnett, R. T. (2012). Risk
816 of nonaccidental and cardiovascular mortality in relation to long-term exposure to low concentrations of fine
817 particulate matter: A Canadian national-level cohort study. *Environmental Health Perspectives*, *120*(5), 708–
818 714. <https://doi.org/10.1289/ehp.1104049>
- 819 Daellenbach, K. R., Uzu, G., Jiang, J., Cassagnes, L.-E., Leni, Z., Vlachou, A., Stefenelli, G., Canonaco, F., Weber,
820 S., Segers, A., & Sources, al. (2020). Sources of particulate-matter air pollution and its oxidative potential
821 in Europe of particulate-matter air pollution and its oxidative potential in Europe. *Nature*, *587*(7834).
822 <https://doi.org/10.1038/s41586-020-2902-8i>
- 823 Deng, M., Chen, D., Zhang, G., & Cheng, H. (2022). Policy-driven variations in oxidation potential and source
824 apportionment of PM_{2.5} in Wuhan, central China. *Science of the Total Environment*, *853*(May), 158255.
825 <https://doi.org/10.1016/j.scitotenv.2022.158255>
- 826 Dominici, F. (2004). Time-series analysis of air pollution and mortality: a statistical review. *Research Report*
827 (*Health Effects Institute*), *123*, 3–27.
- 828 Dominutti, P. A., Borlaza, L., Sauvain, J. J., Ngoc Thuy, V. D., Houdier, S., Suarez, G., Jaffrezo, J. L., Tobin, S.,
829 Trébuchon, C., Socquet, S., Moussu, E., Mary, G., & Uzu, G. (2023). Source apportionment of oxidative
830 potential depends on the choice of the assay: insights into 5 protocols comparison and implications for
831 mitigation measures. *Environmental Science: Atmospheres*. <https://doi.org/10.1039/d3ea00007a>
- 832 Elangasinghe, M. A., Singhal, N., Dirks, K. N., & Salmond, J. A. (2014). Development of an ANN-based air
833 pollution forecasting system with explicit knowledge through sensitivity analysis. *Atmospheric Pollution*
834 *Research*, *5*(4), 696–708. <https://doi.org/10.5094/APR.2014.079>
- 835 Fadel, M., Courcot, D., Delmaire, G., Roussel, G., Afif, C., & Ledoux, F. (2023). Source apportionment of PM_{2.5}
836 oxidative potential in an East Mediterranean site. *Science of the Total Environment*, *900*(July).
837 <https://doi.org/10.1016/j.scitotenv.2023.165843>
- 838 Fang, T., Verma, V., T Bates, J., Abrams, J., Klein, M., Strickland, J. M., Sarnat, E. S., Chang, H. H., Mulholland,
839 A. J., Tolbert, E. P., Russell, G. A., & Weber, J. R. (2016). Oxidative potential of ambient water-soluble
840 PM_{2.5} in the southeastern United States: Contrasts in sources and health associations between ascorbic acid
841 (AA) and dithiothreitol (DTT) assays. *Atmospheric Chemistry and Physics*, *16*(6), 3865–3879.
842 <https://doi.org/10.5194/acp-16-3865-2016>
- 843 Favez, O. (2017). *Traitement harmonisé de jeux de données multi-sites pour l'étude des sources de PM par*
844 *Positive Matrix Factorization*.
- 845 Godri, K. J., Harrison, R. M., Evans, T., Baker, T., Dunster, C., Mudway, I. S., & Kelly, F. J. (2011). Increased
846 oxidative burden associated with traffic component of ambient particulate matter at roadside and Urban
847 background schools sites in London. *PLoS ONE*, *6*(7). <https://doi.org/10.1371/journal.pone.0021961>
- 848 Goldfeld, S. M., & Quandt, R. E. (1965). Some Tests for Homoscedasticity Author (s): Stephen M . Goldfeld and
849 Richard E . Quandt Source : Journal of the American Statistical Association , Jun ., 1965 , Vol . 60 , No .
850 310 Published by : Taylor & Francis , Ltd . on behalf of the American Statis. *Journal of the American*
851 *Statistical Association*, *60*(310), 539–547.
- 852 Harrell. (2016). Regression Modeling Strategies. *Technometrics*, *45*(2), 170–170.
853 <https://doi.org/10.1198/tech.2003.s158>
- 854 Hastie, T. et. all. (2009). Springer Series in Statistics The Elements of Statistical Learning. *The Mathematical*
855 *Intelligencer*, *27*(2), 83–85. <http://www.springerlink.com/index/D7X7KX6772HQ2135.pdf>
- 856 Hawkins, D. M. (2004). The Problem of Overfitting. *Journal of Chemical Information and Computer Sciences*,
857 *44*(1), 1–12. <https://doi.org/10.1021/ci0342472>
- 858 Hernandez, P. A., Graham, C. H., Master, L. L., & Albert, D. L. (2006). The effect of sample size and species
859 characteristics on performance of different species distribution modeling methods. *Ecography*, *29*(5), 773–
860 785. <https://doi.org/10.1111/j.0906-7590.2006.04700.x>
- 861 Hoerl, A. E., & Kennard, R. W. (1970). Ridge Regression: Applications to Nonorthogonal Problems.
862 *Technometrics*, *12*(1), 69. <https://doi.org/10.2307/1267352>

- 863 in 't Veld, M., Pandolfi, M., Amato, F., Pérez, N., Reche, C., Dominutti, P., Jaffrezo, J., Alastuey, A., Querol, X.,
864 & Uzu, G. (2023). Discovering oxidative potential (OP) drivers of atmospheric PM₁₀, PM_{2.5}, and PM₁
865 simultaneously in North-Eastern Spain. *Science of the Total Environment*, 857(August 2022).
866 <https://doi.org/10.1016/j.scitotenv.2022.159386>
- 867 Janssen, N. A. H., Yang, A., Strak, M., Steenhof, M., Hellack, B., Gerlofs-Nijland, M. E., Kuhlbusch, T., Kelly,
868 F., Harrison, R., Brunekreef, B., Hoek, G., & Cassee, F. (2014). Oxidative potential of particulate matter
869 collected at sites with different source characteristics. *Science of the Total Environment*, 472, 572–581.
870 <https://doi.org/10.1016/j.scitotenv.2013.11.099>
- 871 Kelly, F. J., & Mudway, I. S. (2003). Protein oxidation at the air-lung interface. *Amino Acids*, 25(3–4), 375–396.
872 <https://doi.org/10.1007/s00726-003-0024-x>
- 873 Kuhn, M., & Johnson, K. (2013). Applied predictive modeling. In *Applied Predictive Modeling*.
874 <https://doi.org/10.1007/978-1-4614-6849-3>
- 875 Leni, Z., Cassagnes, L. E., Daellenbach, K. R., Haddad, I. El, Vlachou, A., Uzu, G., Prévôt, A. S. H., Jaffrezo, J.
876 L., Baumlin, N., Salathe, M., Baltensperger, U., Dommen, J., & Geiser, M. (2020). Oxidative stress-induced
877 inflammation in susceptible airways by anthropogenic aerosol. *PLoS ONE*, 15(11 November).
878 <https://doi.org/10.1371/journal.pone.0233425>
- 879 Li, J., Zhao, S., Xiao, S., Li, X., Wu, S., Zhang, J., & Schwab, J. J. (2023). Source apportionment of water-soluble
880 oxidative potential of PM_{2.5} in a port city of Xiamen, Southeast China. *Atmospheric Environment*,
881 314(June), 120122. <https://doi.org/10.1016/j.atmosenv.2023.120122>
- 882 Li, Xia, T., & Nel, A. E. (2008). The role of oxidative stress in ambient particulate matter-induced lung diseases
883 and its implications in the toxicity of engineered nanoparticles. *Free Radical Biology and Medicine*, 44(9),
884 1689–1699. <https://doi.org/10.1016/j.freeradbiomed.2008.01.028>
- 885 Liu, & Ng. (2023). Toxicity of Atmospheric Aerosols: Methodologies & Assays. *American Chemical Society*.
886 <https://doi.org/DOI:10.1021/acscinfocus.7e7012>
- 887 Liu, W. J., Xu, Y. S., Liu, W. X., Liu, Q. Y., Yu, S. Y., Liu, Y., Wang, X., & Tao, S. (2018). Oxidative potential
888 of ambient PM_{2.5} in the coastal cities of the Bohai Sea, northern China: Seasonal variation and source
889 apportionment. *Environmental Pollution*, 236, 514–528. <https://doi.org/10.1016/j.envpol.2018.01.116>
- 890 Lodovici, M., & Bigagli, E. (2011). Oxidative stress and air pollution exposure. *Journal of Toxicology*, 2011.
891 <https://doi.org/10.1155/2011/487074>
- 892 Matsuki, K., Kuperman, V., & Van Dyke, J. A. (2016). The Random Forests statistical technique: An examination
893 of its value for the study of reading. *Scientific Studies of Reading*, 20(1), 20–33.
894 <https://doi.org/10.1080/10888438.2015.1107073>
- 895 McCullagh. (1989). Generalized linear models. In *Statistical Models in S* (pp. 195–247).
896 <https://doi.org/10.1201/9780203738535>
- 897 Montgomery C, D., Peck A, E., & Vining, G. G. (2012). *Introducing To Linear Regression Analysis (5th ed.)*.
- 898 Mudway, I. S., Kelly, F. J., & Holgate, S. T. (2020). Oxidative stress in air pollution research. In *Free Radical*
899 *Biology and Medicine* (Vol. 151, pp. 2–6). Elsevier Inc.
900 <https://doi.org/10.1016/j.freeradbiomed.2020.04.031>
- 901 Nelin, T. D., Joseph, A. M., Gorr, M. W., & Wold, L. E. (2012). Direct and indirect effects of particulate matter
902 on the cardiovascular system. *Toxicology Letters*, 208(3), 293–299.
903 <https://doi.org/10.1016/j.toxlet.2011.11.008>
- 904 O'Brien, R. M. (2007). A caution regarding rules of thumb for variance inflation factors. *Quality and Quantity*,
905 41(5), 673–690. <https://doi.org/10.1007/s11135-006-9018-6>
- 906 Paatero, P., & Hopke, P. K. (2009). Rotational tools for factor analytic models. *Journal of Chemometrics*, 23(2),
907 91–100. <https://doi.org/10.1002/cem.1197>
- 908 Paatero, P., & Tappert, U. (1994). Positive matrix factorization: A non-negative factor model with optimal
909 utilization of error estimates of data values. In *Environmetrics* (Vol. 5).
910 <https://doi.org/https://doi.org/10.1002/env.3170050203>

- 911 Pearce, J., & Ferrier, S. (2000). An evaluation of alternative algorithms for fitting species distribution models using
 912 logistic regression. *Ecological Modelling*, *128*(2–3), 127–147. [https://doi.org/10.1016/S0304-3800\(99\)00227-6](https://doi.org/10.1016/S0304-3800(99)00227-6)
 913
- 914 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P.,
 915 Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., & Duchesnay,
 916 E. (2011). Scikit-learn: Machine Learning in {P}ython. *Journal of Machine Learning Research*, *12*, 2825–
 917 2830.
- 918 Pelucchi, C., Negri, E., Gallus, S., Boffetta, P., Tramacere, I., & La Vecchia, C. (2009). Long-term particulate
 919 matter exposure and mortality: A review of European epidemiological studies. *BMC Public Health*, *9*, 1–8.
 920 <https://doi.org/10.1186/1471-2458-9-453>
- 921 Peng, R. D., Bell, M. L., Geyh, A. S., McDermott, A., Zeger, S. L., Samet, J. M., & Dominici, F. (2009).
 922 Emergency admissions for cardiovascular and respiratory diseases and the chemical composition of fine
 923 particle air pollution. *Environmental Health Perspectives*, *117*(6), 957–963.
 924 <https://doi.org/10.1289/ehp.0800185>
- 925 Pietrogrande, M. C., Romanato, L., & Russo, M. (2022). Synergistic and Antagonistic Effects of Aerosol
 926 Components on Its Oxidative Potential as Predictor of Particle Toxicity. *Toxics*, *10*(4).
 927 <https://doi.org/10.3390/toxics10040196>
- 928 Pope, C. A., & Dockery, D. W. (2006). Health effects of fine particulate air pollution: Lines that connect. *Journal*
 929 *of the Air and Waste Management Association*, *56*(6), 709–742.
 930 <https://doi.org/10.1080/10473289.2006.10464485>
- 931 Rao, X., Zhong, J., Brook, R. D., & Rajagopalan, S. (2018). Effect of Particulate Matter Air Pollution on
 932 Cardiovascular Oxidative Stress Pathways. *Antioxidants and Redox Signaling*, *28*(9), 797–818.
 933 <https://doi.org/10.1089/ars.2017.7394>
- 934 Raudys, S. J., & Jain, A. K. (1991). Small Sample Size Effects in Statistical Pattern Recognition:
 935 Recommendations for Practitioners. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*
 936 (Vol. 13, Issue 3, pp. 252–264). <https://doi.org/10.1109/34.75512>
- 937 Rosenblad, A. (2011). The Concise Encyclopedia of Statistics. In *Journal of Applied Statistics* (Vol. 38, Issue 4).
 938 <https://doi.org/10.1080/02664760903075614>
- 939 Samake, A., Uzu, G., Martins, J. M. F., Calas, A., Vince, E., Parat, S., & Jaffrezo, J. L. (2017). The unexpected
 940 role of bioaerosols in the Oxidative Potential of PM. *Scientific Reports*, *7*(1). <https://doi.org/10.1038/s41598-017-11178-0>
 941
- 942 Seabold, S., & Perktold, J. (2010). Statsmodels: Econometric and statistical modeling with python. *9th Python in*
 943 *Science Conference*.
- 944 Shangguan, Y., Zhuang, X., Querol, X., Li, B., Moreno, N., Trechera, P., Sola, P. C., Uzu, G., & Li, J. (2022).
 945 Characterization of deposited dust and its respirable fractions in underground coal mines: Implications for
 946 oxidative potential-driving species and source apportionment. *International Journal of Coal Geology*,
 947 *258*(December 2021). <https://doi.org/10.1016/j.coal.2022.104017>
- 948 Stevanović, S., Jovanović, M. V., Jovašević-Stojanović, M. V., & Ristovski, Z. (2023). SOURCE
 949 APPORTIONMENT OF OXIDATIVE POTENTIAL What We Know So Far. *Thermal Science*, *27*(3),
 950 2347–2357. <https://doi.org/10.2298/TSCI221107111S>
- 951 Stockwell, D. R. B., & Peterson, A. T. (2002). Effects of sample size on accuracy of species distribution models.
 952 *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)
- 953 Szigeti, T., Dunster, C., Cattaneo, A., Cavallo, D., Spinazzè, A., Saraga, D. E., Sakellaris, I. A., de Kluizenaar, Y.,
 954 Cornelissen, E. J. M., Hänninen, O., Peltonen, M., Calzolari, G., Lucarelli, F., Mandin, C., Bartzis, J. G.,
 955 Záráy, G., & Kelly, F. J. (2016). Oxidative potential and chemical composition of PM_{2.5} in office buildings
 956 across Europe - The OFFICAIR study. *Environment International*, *92–93*, 324–333.
 957 <https://doi.org/10.1016/j.envint.2016.04.015>
- 958 Szigeti, T., Óvári, M., Dunster, C., Kelly, F. J., Lucarelli, F., & Záráy, G. (2015). Changes in chemical composition
 959 and oxidative potential of urban PM_{2.5} between 2010 and 2013 in Hungary. *Science of the Total*
 960 *Environment*, *518–519*, 534–544. <https://doi.org/10.1016/j.scitotenv.2015.03.025>

- 961 Tibshirani, R. (1996). Regression Shrinkage and Selection Via the Lasso. *Journal of the Royal Statistical Society:*
962 *Series B (Methodological)*, 58(1), 267–288. <https://doi.org/10.1111/j.2517-6161.1996.tb02080.x>
- 963 Verma, V., Fang, T., Guo, H., King, L., Bates, J. T., Peltier, R. E., Edgerton, E., Russell, A. G., & Weber, R. J.
964 (2014). Reactive oxygen species associated with water-soluble PM_{2.5} in the southeastern United States:
965 Spatiotemporal trends and source apportionment. *Atmospheric Chemistry and Physics*, 14(23), 12915–
966 12930. <https://doi.org/10.5194/acp-14-12915-2014>
- 967 Viana, M., Kuhlbusch, T. A. J., Querol, X., Alastuey, A., Harrison, R. M., Hopke, P. K., Winiwarter, W., Vallius,
968 M., Szidat, S., Prévôt, A. S. H., Hueglin, C., Bloemen, H., Wählin, P., Vecchi, R., Miranda, A. I., Kasper-
969 Giebl, A., Maenhaut, W., & Hitzenberger, R. (2008). Source apportionment of particulate matter in Europe:
970 A review of methods and results. In *Journal of Aerosol Science* (Vol. 39, Issue 10, pp. 827–849). Elsevier
971 Ltd. <https://doi.org/10.1016/j.jaerosci.2008.05.007>
- 972 Vida, M., Foret, G., Siour, G., Coman, A., Weber, S., Favez, O., Jaffrezo, J., Pontet, S., Mesbah, B., Gille, G.,
973 Zhang, S., Chevrier, F., Pallares, C., Uzu, G., & Beekmann, M. (2024). Oxidative potential modelling of
974 PM₁₀: a 2-year study over France. *ACDP*.
- 975 Wang, D., Yang, X., Lu, H., Li, D., Xu, H., Luo, Y., Sun, J., Hang Ho, S. S., & Shen, Z. (2023). Oxidative potential
976 of atmospheric brown carbon in six Chinese megacities: Seasonal variation and source apportionment.
977 *Atmospheric Environment*, 309(June), 119909. <https://doi.org/10.1016/j.atmosenv.2023.119909>
- 978 Wang, J., Jiang, H., Jiang, H., Mo, Y., Geng, X., Li, J., Mao, S., Bualert, S., Ma, S., Li, J., & Zhang, G. (2020).
979 Source apportionment of water-soluble oxidative potential in ambient total suspended particulate from
980 Bangkok: Biomass burning versus fossil fuel combustion. *Atmospheric Environment*, 235(May), 117624.
981 <https://doi.org/10.1016/j.atmosenv.2020.117624>
- 982 Wang, S., Ye, J., Soong, R., Wu, B., Yu, L., Simpson, A. J., & Chan, A. W. H. (2018). Relationship between
983 chemical composition and oxidative potential of secondary organic aerosol from polycyclic aromatic
984 hydrocarbons. *Atmospheric Chemistry and Physics*, 18(6), 3987–4003. <https://doi.org/10.5194/acp-18-3987-2018>
- 986 Wang, Y., Wang, M., Li, S., Sun, H., Mu, Z., Zhang, L., Li, Y., & Chen, Q. (2020). Study on the oxidation potential
987 of the water-soluble components of ambient PM_{2.5} over Xi'an, China: Pollution levels, source
988 apportionment and transport pathways. *Environment International*, 136(January), 105515.
989 <https://doi.org/10.1016/j.envint.2020.105515>
- 990 Weber, S., Salameh, D., Albinet, A., Alleman, L. Y., Waked, A., Besombes, J. L., Jacob, V., Guillaud, G.,
991 Meshbah, B., Rocq, B., Hulin, A., Dominik-Sègue, M., Chrétien, E., Jaffrezo, J. L., & Favez, O. (2019).
992 Comparison of PM₁₀ sources profiles at 15 french sites using a harmonized constrained positive matrix
993 factorization approach. *Atmosphere*, 10(6). <https://doi.org/10.3390/atmos10060310>
- 994 Weber, S., Uzu, G., Calas, A., Chevrier, F., Besombes, J. L., Charron, A., Salameh, D., Ježek, I., Močnik, G., &
995 Jaffrezo, J. L. (2018). An apportionment method for the oxidative potential of atmospheric particulate matter
996 sources: Application to a one-year study in Chamonix, France. *Atmospheric Chemistry and Physics*, 18(13),
997 9617–9629. <https://doi.org/10.5194/acp-18-9617-2018>
- 998 Weber, S., Uzu, G., Favez, O., Borlaza, L., Calas, A., Salameh, D., Chevrier, F., Allard, J., Besombes, J. L.,
999 Albinet, A., Pontet, S., Mesbah, B., Gille, G., Zhang, S., Pallares, C., Leoz-Garziandia, E., & Jaffrezo, J. L.
1000 (2021). Source apportionment of atmospheric PM₁₀ oxidative potential: Synthesis of 15 year-round urban
1001 datasets in France. *Atmospheric Chemistry and Physics*, 21(14), 11353–11378. <https://doi.org/10.5194/acp-21-11353-2021>
- 1003 WHO. (2021). *WHO global air quality guidelines*.
- 1004 Williams, M., Gomez Grajales, C. A., & Kurkiewicz, D. (2013). Assumptions of Multiple Regression: Correcting
1005 Two Misconceptions - Practical Assessment, Research & Evaluation. *Practical Assessment, Research, and*
1006 *Evaluation (PARE)*, 18(11), 1–16. <https://scholarworks.umass.edu/pare/vol18/iss1/11>
- 1007 Wisz, M. S., Hijmans, R. J., Li, J., Peterson, A. T., Graham, C. H., Guisan, A., Elith, J., Dudík, M., Ferrier, S.,
1008 Huettmann, F., Leathwick, J. R., Lehmann, A., Lohmann, L., Loiselle, B. A., Manion, G., Moritz, C.,
1009 Nakamura, M., Nakazawa, Y., Overton, J. M. C., ... Zimmermann, N. E. (2008). Effects of sample size on
1010 the performance of species distribution models. *Diversity and Distributions*, 14(5), 763–773.
1011 <https://doi.org/10.1111/j.1472-4642.2008.00482.x>

- 1012 Xiong, Q., Yu, H., Wang, R., Wei, J., & Verma, V. (2017). Rethinking Dithiothreitol-Based Particulate Matter
1013 Oxidative Potential: Measuring Dithiothreitol Consumption versus Reactive Oxygen Species Generation.
1014 *Environmental Science and Technology*, 51(11), 6507–6514. <https://doi.org/10.1021/acs.est.7b01272>
- 1015 Yang, A., Jedynska, A., Hellack, B., Kooter, I., Hoek, G., Brunekreef, B., Kuhlbusch, T. A. J., Cassee, F. R., &
1016 Janssen, N. A. H. (2014). Measurement of the oxidative potential of PM_{2.5} and its constituents: The effect
1017 of extraction solvent and filter type. *Atmospheric Environment*, 83, 35–42.
1018 <https://doi.org/10.1016/j.atmosenv.2013.10.049>
- 1019 Yu, Guo, S., Xu, R., Ye, T., Li, S., Sim, M. R., Abramson, M. J., & Guo, Y. (2021). Cohort studies of long-term
1020 exposure to outdoor particulate matter and risks of cancer: A systematic review and meta-analysis.
1021 *Innovation*, 2(3), 100143. <https://doi.org/10.1016/j.xinn.2021.100143>
- 1022 Yu, S. Y., Liu, W. J., Xu, Y. S., Yi, K., Zhou, M., Tao, S., & Liu, W. X. (2019). Characteristics and oxidative
1023 potential of atmospheric PM_{2.5} in Beijing: Source apportionment and seasonal variation. *Science of the
1024 Total Environment*, 650, 277–287. <https://doi.org/10.1016/j.scitotenv.2018.09.021>
- 1025 Zhang, Y., Albinet, A., Petit, J. E., Jacob, V., Chevrier, F., Gille, G., Pontet, S., Chrétien, E., Dominik-Sègue, M.,
1026 Levigoureux, G., Močnik, G., Gros, V., Jaffrezo, J. L., & Favez, O. (2020). Substantial brown carbon
1027 emissions from wintertime residential wood burning over France. *Science of the Total Environment*, 743.
1028 <https://doi.org/10.1016/j.scitotenv.2020.140752>
- 1029