



1

2

3

4

5

6

7

8

**RECEPTOR MODELLING OF BOTH PARTICLE COMPOSITION
AND SIZE DISTRIBUTION FROM A BACKGROUND SITE IN
LONDON, UK – THE TWO STEP APPROACH**

10

11

12

13

David C.S. Beddows and Roy M. Harrison*†

14

15

16

17

**National Centre for Atmospheric Science
School of Geography, Earth and Environmental Sciences
University of Birmingham
Edgbaston, Birmingham B15 2TT
United Kingdom**

18

19

20

21

22

23

24

25

26

27

*To whom correspondence should be addressed.

Tele: +44 121 414 3494; Fax: +44 121 414 3708; Email: r.m.harrison@bham.ac.uk

†Also at: Department of Environmental Sciences / Center of Excellence in Environmental Studies, King Abdulaziz University, PO Box 80203, Jeddah, 21589, Saudi Arabia



28 **ABSTRACT**

29 Some air pollution datasets contain multiple variables with a range of measurement units,
30 and combined analysis by Positive Matrix Factorization (PMF) is problematic, but can offer
31 benefits from the greater information content. In this work, a novel method is devised and
32 the source apportionment of a mixed unit data set (PM₁₀ mass and Number Size Distribution
33 NSD) is achieved using a novel two-step approach to PMF. In the first step the PM₁₀ data
34 is PMF analysed using a source apportionment approach in order to provide a solution which
35 best describes the environment and conditions considered. The time series G values (and
36 errors) of the PM₁₀ solution are then taken forward into the second step where they are
37 combined with the NSD data and analysed in a second PMF analysis. This results in
38 apportioned NSD data associated with the PM₁₀ factors. We exemplify this approach using
39 data reported in the study of Beddows et al. (2015), producing one solution which unifies
40 the two separate solutions for PM₁₀ and NSD data datasets together. We also show how
41 regression of the NSD size bins and the G time series can be used to elaborate the solution
42 by identifying NSD factors (such as nucleation) not influencing the PM₁₀ mass.

43 **Keywords:** PM₁₀; London; PMF; source apportionment; receptor modelling

44



45 **1. INTRODUCTION**

46 It is unquestionable that worldwide, the scientific vista of air quality is expanding; whether it
47 is the increasing number of observatories or the refinement of information mined from the
48 increasing sophistication of measurements often incorporated in campaign work. The
49 number of metrics being measured has increased from simple measurements of PM mass
50 and gas concentrations, and we can now probe the composition of the PM mass and the
51 size distributions with mass spectrometers, mobility analysers and optical devices.

52

53 Studies using PMF as a tool for source apportionment of particle mass using
54 multicomponent chemical analysis data are published almost daily using datasets from
55 around the world. However, they do not always provide consistent outcomes (Pant and
56 Harrison, 2012), and one means by which source resolution and identification can be
57 improved is by inclusion of auxiliary data, such as gaseous pollutants (Thimmaiah et al.,
58 2009), particle number count (Masiol et al., 2017) or particle size distribution (Beddows et
59 al., 2015; Ogulei et al., 2006; Leoni et al., 2018). However, while combining, for example,
60 particle chemical composition and size distribution data in a single PMF analysis may assist
61 source resolution, it does not allow quantitative attribution of either particle mass or particle
62 number to the source factors.

63

64 Comero et al. (2009) alluded to the problem of including more than one metric with different
65 units when citing Hopke (1991). In order to obtain a physically realistic PMF solution some
66 natural constraints must be satisfied, one being, "*Only for chemical elements or compounds,*
67 *where the unit of measurement are the same, the sum of the predicted elemental mass*
68 *contributions for each source must be less than or equal to total measured mass for each*
69 *element; the whole is greater than or equal to the sum of its parts (only in the case of*
70 *chemical elements or compounds)*". This underlies the necessity to have a consistency of



71 units throughout the input dataset in order to make a quantified apportionment. To exemplify
72 this point, in Harrison et al. (2011), NSD data (merged SMPS and APS data) was analysed
73 with PMF using auxiliary data (meteorology, gas concentration, traffic counts and speed).
74 The study used particle size distribution data collected at the Marylebone Road supersite in
75 London in the autumn of 2007 and put forward a 10 factor solution comprised of roadside
76 and background particle source factors. The size distribution profiles, bivariate plots and
77 diurnal cycles were presented but the contributions of each factor were limited to percentage
78 contributions simply because of the mixed units which were inputted into the analysis, and
79 there can be no confidence as to whether the sources are apportioned by units of number
80 concentration ($1/\text{cm}^3$) or any of the other units used in the auxiliary data. Chan et al. (2011)
81 identified this “*as a matter of debate within the community concerned*” when considering the
82 use of multiple types of composition data for source apportionment. They considered
83 extracting more source information from an aerosol composition dataset by including data
84 on other air pollutants and wind data in the analysis of a small but comprehensive dataset
85 from a 24-hourly sampling programme carried out during June 2001 in an industrial area in
86 Brisbane. They chose multiple types of composition data (aerosols, VOCs and major
87 gaseous pollutants) and wind data in source apportionment of air pollutants and found it to
88 result in better defined source factors and better fit diagnostics, compared to when non-
89 combined data were used. Likewise, Wang et al. (2017) report an improvement in source
90 profiles when coupling the PMF model with ^{14}C data to constrain the PMF run as *a priori*
91 information.

92

93 The potential for an improved factor solution obtained by mixing data types in PMF provides
94 a motivation in the community to develop a methodology which can overcome the
95 aforementioned difficulties. In this study, we present such a method for analysing



96 simultaneously collected PM₁₀ composition and NSD data. In the work of Beddows et al.
97 (2015), both particle composition and number size distribution (NSD) data from a
98 background site in London (2011 and 2012) was analysed using Positive Matrix
99 Factorization. As part of the methodology development, it was concluded that it was
100 preferable not to combine these two data types in the analysis but to conduct separate PMF
101 analyses for PM₁₀ mass and particle number. This yielded a 6 factor solution for the PM₁₀
102 data (Diffuse Urban; Marine; Secondary; Non-Exhaust Traffic / Crustal (NET/Crustal); Fuel
103 Oil; and Traffic. Factors described as Diffuse Urban; Secondary; and Traffic were identified
104 in the 4 factor solution for the NSD data. A further factor was the Nucleation factor. When
105 combining the PM₁₀ and NSD data in a single PMF analysis, Diffuse Urban; Nucleation;
106 Secondary; Aged Marine and Traffic Factors were identified but the factors were not as
107 clearly separated from each other as the factors derived from the separate datasets. For
108 example, Fuel Oil was now mixed in with Marine and called Aged Marine. This is
109 summarized in Figure 1. However in the analysis, it would still be useful to obtain a number
110 size distribution for each of the 6 PM₁₀ factors and/or a chemical composition for the 4 NSD
111 factors.

112

113 In this work, we present a continuation of the analysis of Beddows et al. (2015) describing
114 a two-step methodology in which we use the first step to analyse a primary dataset (PM₁₀;
115 units: $\mu\text{g}/\text{m}^3$) and then combine the output with a second dataset (NSD; units: $1/\text{cm}^3$). The
116 first step identifies sources and apportions their contribution to mass. Then in the second
117 step, PM₁₀ factors are augmented by number size distribution factors. We show that a more
118 complete picture of the sources can be obtained using a 2-step (PMF-PMF) analysis.
119 Furthermore, we also consider linear regression as a second step in a PMF-LR analysis to
120 show how this can reveal hidden factors.



121 **2. EXPERIMENTAL**

122 With a population of 8.5 million in 2014 (ONS, 2017), the UK city of London is the focus of
123 study in this work where the London *North Kensington* (NK) Site ($LAT = 51^{\circ} : 31' : 15.780''$
124 N and $LONG = 0^{\circ} : 12' : 48.571''$ W) was considered. NK is part of both the London Air
125 Quality Network and the national Automatic Urban and Rural Network and is owned and
126 part-funded by the Royal Borough of Kensington and Chelsea. The facility is located within
127 a self contained cabin within the grounds of Sion Manning School. The nearest road, St.
128 Charles Square, is a quiet residential street approximately 5 metres from the monitoring site
129 and the surrounding area is mainly residential. The nearest heavily trafficked roads are the
130 B450 (~100 m East) and the very busy A40 (~400 m South). For a detailed overview of the
131 air pollution climate at North Kensington, the reader is referred to Bigi and Harrison (2010).

132

133 **2.1 Data**

134 For this study, the same datasets considered, and PMF analysis outputs generated, by
135 Beddows et al. (2015) were used. For this, 24h air samples were taken daily over a two
136 year period (2011 and 2012) using a Thermo Partisol 2025 sampler fitted with a PM₁₀ size
137 selective inlet, and alongside, Number Size Distribution (NSD) data were collected
138 continuously every ¼ hour using a Scanning Mobility Particle Sizer (SMPS) consisting of a
139 CPC (TSI model 3775) combined with an electrostatic classifier (TSI model 3080) in air dried
140 according to the EUSAAR protocol (Wiedensohler et al., 2012). The particle sizes covered
141 were 51 size bins ranging from 16.55 nm to 604.3 nm. Analysis of this data resulted in PMF
142 source profiles F and source time series G of the PM₁₀ and NSD data sets which were
143 carried forward into this work. Further details of the data, collection methods, coverage and
144 first analysis are given in Beddows et al. (2015).

145



146 2.2 Proxy Data

147 Besides the PM₁₀ mass, estimates of PM mass can be derived using the NSD assuming
148 spherical particles of a fixed density. For the SMPS settings, a particle size range between
149 16 and 604 nm is collected which can be used to estimate a PM_{0.6} value using equation 1.

$$PM_{0.6} = \rho_{eff} \times \frac{\pi}{6} \sum_{sizeBins} d^3 \quad (1)$$

150

151 Where ρ_{eff} is set to 2 g/cm³ for a Diffuse Urban background (based upon 1.8-2.5 g/cm³ for
152 an urban background aerosol; Beddows et al., 2010).

153

154 Figure S1 plots the total apportioned PM₁₀ mass against the PM_{0.6} estimates and shows that
155 although the SMPS does not account for the whole mass, it does track with the total PM,
156 with a fitted gradient of 0.65, i.e. accounting for 65% of the mass. To account for the particles
157 greater than 600 nm in the PMF analysis, a proxy was used created by using the difference
158 between the total daily apportioned PM mass in the step 1 of the PMF analysis and the mass
159 estimated from the SMPS data. This difference was then converted back into a number and
160 added to the NSD matrix of counts as $PN_{0.6-10}$ to improve the match of the NSD matrix to
161 the PM₁₀.

162

163 2.3 Methods

164 2.3.1 PMF

165 Positive Matrix Factorization (PMF) is a well-established multivariate data analysis method
166 used in the field of aerosol science. PMF can be described as a least-squares formulation
167 of factor analysis developed by Paatero (Paatero and Tapper, 1994). It assumes that the



168 ambient aerosol concentration X (represented by $n \times m$ matrix of n observations and m PM₁₀
169 constituents or NSD size bins), measured at one or more sites can be explained by the
170 product of a source profile matrix F and source contribution matrix G whose elements are
171 given by equation 1:

$$x_{ij} = \sum_{k=1}^p g_{ik} \cdot f_{kj} + e_{ij} \quad i=1 \dots n; j=1 \dots m \quad (1)$$

172 where the j^{th} PM constituent (element, size bin, or auxiliary measurement) on the i^{th}
173 observation (i.e. hour) is represented by x_{ij} . The term g_{ik} is the contribution of the k^{th} factor
174 to the receptor on the i^{th} hour, f_{kj} is the fraction of j^{th} PM constituent in the k^{th} factor, and e_{ij}
175 is the residual for the j^{th} measurement on the i^{th} hour. The residuals (i.e. difference between
176 measured and reconstructed concentrations) are accounted for in matrix E and the two
177 matrices G and F are obtained by an iterative algorithm which minimises the object function
178 Q (see equation 2).

179

180 Given, the data and uncertainty to matrices for the model, equation 1 is optimised in the
181 PMF2 algorithm by minimising the Q value (equation 2),

$$Q = \sum_{i=1}^n \sum_{j=1}^m \left(\frac{e_{ij}}{s_{ij}} \right)^2 \quad (2)$$

182

183 where s_{ij} is the uncertainty in the j^{th} measurement for hour i .

184

185 It may be seen from equation (2) that PMF is a weighted technique; the value of Q , and
186 hence the model fit, is determined by the input variables with the lowest values of
187 uncertainty, S_{ij} . Input variables with high uncertainty have little effect upon the value of Q ,



188 even when their residuals are large. This can be used to the advantage of the operator.
189 When apportioning total PM mass in a conventional one-stage PMF, the total PM
190 concentrations are normally input with artificially high uncertainty, so that they are essentially
191 passive in the PMF analysis and do not influence its outcome. By doing so, the chemical
192 composition data determine the apportionment of PM mass to the source-related factors
193 identified by the PMF. In this work, the primary aim was to define a particle size distribution
194 associated with each factor derived from the PMF of PM₁₀ composition. Consequently, in
195 the second stage of the PMF, large uncertainties were input for the particle number data,
196 combined with realistic uncertainties for the PM G-values, so that the latter determine the
197 outcome of (“drive”) the PMF analysis.

198

199 In this work, the Q value is outputted by PMF2 and compared to a theoretical value Q_{theory}
200 which is approximately the difference between the product of the dimensions of X and the
201 product of the number of factors and the sum of dimensions of X (i.e. $n \times m - p(n + m)$) $p \times$
202 $x \ m$. For a given number of factors, the whole uncertainty matrix is scaled by a factor X_{scale}
203 until the ratio between Q and Q_{theory} is approximately one ($rQ \text{ value} = Q/Q_{\text{theory}} \approx 1$).

204

205 **2.3.2. Application of PMF**

206 The two step method is shown schematically in Figure 2. In the current example, it uses the
207 PMF output of Beddows et al. (2015) as a starting point and assumes that a PMF analysis
208 of the PM₁₀ chemical composition dataset (Step One) has already been carried out and dealt
209 with as in the previous study. In this current work, a second step which takes the output
210 from the first step and uses it as an input for the second step is developed. This is done by
211 using the G1 time series from the PMF analysis of PM₁₀ and combining this with secondary



212 data, (i.e. NSD data). The uncertainties of the G1 matrix are transferred from the output of
213 the first step and entered as input uncertainties for the second step. For the NSD data, the
214 uncertainties are taken as X times the NSD values in order to be large and ensure that the
215 PMF is driven by the G1 matrix (see Figure 2). The value of X was optimised in Cran R so
216 that the ratio of $Q/Q_{\text{theory}} \sim 1$.

217

218 **2.3.3 Fkey**

219 Fkey is a feature in for incorporating a priori information into a PMF analysis and is used in
220 the second step of the PMF-PMF analysis. It is used to “pull” elements of the source profiles
221 to zero. This method uses a matrix that indicates the location of suspected zeros in source
222 profiles or contributions (Figure 3). Since here it is concerned with the profiles, this
223 information is given in the form of integer values in an Fkey. The greater the certainty that
224 an element of a source profile is zero, the larger the integer value that is specified. In this
225 case, in the second step, it is certain that only one PM G score from one of the sources will
226 be strong, e.g. the traffic source will be the only contributing to the PM G value in the Traffic
227 NSD profile, and likewise for the other sources: Diffuse Urban; Secondary; Marine; Fuel Oil;
228 and NET & Crustal (Figure 3).

229

230 **2.4 Regression**

231 The output of the regression of a dependent variable Y regressed against independent
232 variables X1, X2, X3, ... Xn is n gradients and one intercept. When n = 1 it yields a line,
233 when n = 2 it is a fitted plane. But when n > 2 or in this case n = 6, it is a multidimensional
234 fitted model. Each of the n gradients show how Y varies with the n X values given that the
235 other X values are fixed and the intercept provides a bias value. If Y is allowed to take on



236 each value of the NSD size bin and X variables are set to the 6 G time series from the first
237 step of PMF analysis, then it can be seen how the NSD are correlated to the 6 G time series
238 and infer an associated NSD for each of the factors derived in the first step of the PMF-LR
239 analysis.

240

241 As an alternative second step, each size bin within the NSD was regressed with the six G1
242 time series, Equation 3.

$$NSD[a, j] = \sum_{a=1G1...1G6} grad_{a,j} + int_j \quad (3)$$

243

244 This results in a 7 by 51 matrix of values. Each column represents a size bin of the NSD
245 data and each row represents the gradients associated with 6 of the factors (giving an
246 indication of how each size bin correlates with each of the 6 factors) and an intercept. When
247 $grad_{a,j}$ is plotted against the size bin, 6 plots showing the dependence of each size bin on
248 each of the 6 PM₁₀ factors are produced. It is also assumed that these will be comparable
249 to the actual source profile. Similarly, the int_j values are expected to give a background
250 value, possibly noise. However, this method can extract information known as a remainder
251 factor, shown later in this paper.

252

253 **2.5 Peak Fitting**

254 If it is assumed that the factors derived from the daily NSD data are the same as those
255 present in the hourly data, i.e. the factors are conserved when averaging the data from
256 hourly to daily data before PMF analysis, then daily NSD profiles can be fitted to the hourly
257 NSD spectra to recover a diurnal cycle for the factors. Given the i^{th} number size distribution,



258 NSD_i , the difference $D_{i,j,k}$ (equation 3), between the k^{th} element and the linear superposition
259 of the k^{th} element of the seven factors $f_{j,k}$ is minimised.

$$D_{i,j,k} = f(x) = \begin{cases} NSD_{i,k} - \sum_{j=1}^7 a_j \times f_{j,k}, & a_j \geq 0 \\ 1 \times 10^{10}, & a_j < 0 \end{cases} \quad (3)$$

260

261 The Cran R package Non-Linear Minimization (nlm) (R Core Team, 2018) was used to
262 minimise equation 3. A non-negative constraint is placed in the function. If a negative value
263 is returned by any of the a_j values then $D_{i,j,k}$ returns an excessively large value. Furthermore,
264 in order to extract an apportionment to number concentration ($1/\text{cm}^3$) the fitted values were
265 scaled using a factor SA_j . Six values were derived for SA_j by regressing the total particle
266 number (total hourly SMPS) against each of the fitted values a_j .

267

268 3. RESULTS AND DISCUSSION

269 The aim of this work is to take the results from the first step of a PMF analysis where a
270 successful source apportionment study has been completed and then complement the
271 results with a second step to derive further information about the sources. This can be done
272 using a second PMF analysis or a regression.

273

274 3.1 2-Step PMF-PMF Analysis

275 Figure 4 presents our results from the second PMF analysis of a combined dataset. The
276 G1 time series and uncertainties from the first PMF analysis of PM_{10} data are carried over
277 into the second step where they are combined with the NSD data for PMF analysis. The
278 uncertainties of the NSD data are taken as an optimised multiple of the NSD values
279 themselves. Also in order to maintain the solution from step 1 in step 2 the Fkey matrix is



280 applied to pull elements in the source matrix to zero as described. This ensures that PMF
281 analysis of the NSD data is driven by the G1 time series. This results in a 6 factor solution
282 in which there are unique contributions from one of the G1 scores and an associated NSD
283 source profile, and it is notable that they are surprisingly similar to those calculated for the
284 just-NSD and PM₁₀+NSD data in Beddows et al. (2015). The Diffuse Urban factor has a
285 modal-diameter just below 0.1 μm which is comparable to the NSD factor in the just-NSD
286 analysis. Marine is comparable to the Aged Marine factor derived from the PM₁₀+NSD
287 analysis. The Secondary factor is again the factor with the largest modal diameter (between
288 0.4 and 0.5 μm) and traffic has as expected a modal diameter between 30 and 40 nm. Fuel
289 Oil is interesting as it appears to be a combination of a nucleation factor and a mode
290 comparable to diesel exhaust seen in the Traffic factor.

291

292 **3.2 2-Step PMF-LR Analysis**

293 Figure S2 shows the results of the linear regression of the NSD data plotted against the
294 PM₁₀ G1 scores and again what is remarkable is the similarity between these correlation
295 plots and both the factors derived in Beddows et al. (2015) and those from the 2-step PMF-
296 PMF analysis. This analysis was carried out using daily averaged data. To obtain hourly
297 information and thus obtain the diurnal patterns, the resulting correlation factors were re-
298 fitted to the original NSD data. On inspection of these source profiles and diurnal plots, the
299 negative values make interpretation a struggle reinforcing one of the 4 conditions (Hopke,
300 1991) in the analysis if it is to make sense. We can however fit non-negative gradients using
301 non-negative regression. However, the surprising consequence of applying this constraint
302 is that the same profiles are derived but they are clipped so that all negative values are
303 replaced by zero values – hence, information is lost by doing this. One interpretation is that
304 these are particle sinks but this contradicts the PMF-PMF findings and hence it is concluded



305 that the PMF-LR analysis only serves as an indication of how the PM₁₀ factors are
306 augmented by the NSD data. If all profiles are shifted to above the zero line then
307 comparisons to the PMF-PMF data can be made. However, what is interesting to note in
308 this result is the intercept NSD which is comparable in profile and diurnal pattern to the
309 nucleation mode identified in Beddows et al. (2015). This is a seventh factor in addition to
310 the 6 PM₁₀ factors and suggests that although the PMF analysis of the PM₁₀ data alone
311 misses a Nucleation factor, this can be recovered in a second analysis as a remainder or
312 bias in the data. Furthermore, this result indicates that the composition of the Nucleation
313 NSD factor has no link to the chemical PM₁₀ composition and cannot be used to infer a
314 composition.

315

316 Returning to the PMF-PMF analysis and extending the analysis from 6 factors to 7 factors
317 and adding an extra row in the Fkey matrix which pulls all of the G1 scores to zero in the
318 solution, the same 6 factor solution is obtained with the additional 7th factor (Figure 5 and
319 Figure S3). As expected, this seventh factor is a Nucleation factor by separating out of the
320 fuel oil factor a nucleation mode leaving a mode with a modal-diameter between 50 and 60
321 nm. In the results of Beddows et al. (2015), the Nucleation factor was only seen when
322 applying PMF to the just-NSD and PM₁₀+NSD data, and in the PM₁₀+NSD results, Fuel Oil
323 was not separated and appeared to be smeared across all 5 factors. A seven factor solution
324 to PMF of the PM₁₀ chemical composition data did not reveal this factor, presumably
325 because the mass associated with nucleation mode particles is too small to affect
326 composition significantly.

327

328 Another interesting observation is that although only 4 factors were derived from the PMF
329 analysis of NSD data alone (Diffuse Urban; Secondary; Traffic and Nucleation), when extra



330 information is included from the PMF analysis of the PM₁₀ data, more information can be
331 extracted from the PMF analysis of the NSD data in the form of the Marine; Fuel Oil and
332 NET & Crustal factors. The Nucleation factor is only revealed when performing a regression
333 between the NSD size bins and the G scores of the PMF analysis which leads to increasing
334 the factor number from 6 to 7 which yields the Nucleation profile. It is also reassuring that
335 the bivariate plots for of the 7 factors (discussed in the next section) correspond to the
336 bivariate plots given in Beddows et al. (2015).

337

338 3.3 Diurnal and Bivariate Plots

339 The original PMF was carried out on daily PM₁₀ data and in order to make diurnal and
340 bivariate plots, a higher time resolution is required. It is assumed that the factors derived in
341 the hourly NSD data are the same as those derived from the daily averaged data, i.e. the
342 factors are conserved when averaging the data from hourly to daily data before PMF
343 analysis. Then the hourly NSD data can be fit with the PMF profiles derived from the daily
344 data. Figure 6 shows the resulting diurnal profiles.

345

346 The diurnal trends of the fitted peaks show the values required in equation 3 to fit the 7 daily
347 NSD factors to the hourly NSD data. These have been scaled in these plots according to
348 the integral of the NSD factor measured in 1/cm³. The nucleation diurnal trend behaves as
349 expected rising to a maximum during the day and then falling back down to a minimum at
350 night. This corresponds to the intensity of the sun during the day and the increased
351 likelihood of nucleation on clean days when there is sufficient precursor material to form
352 particles with a low particle condensation sink. Marine is also high during the day
353 presumably due to higher wind speeds. Diffuse Urban, NET and Crustal, and Traffic all



354 follow a trend which is synchronised to the daily cycle of anthropogenic activity and traffic
355 as influenced by greater atmospheric stability at night. Secondary also follows a similar
356 anthropogenic cycle and would be expected to be strongest at night. Fuel Oil is highest
357 during the evening and night and may correspond to home heating rather than marine
358 activity. The particle size distributions associated with the Marine and NET and Crustal
359 sources are of limited value as these sources are dominated by coarse particles, beyond
360 the range of the SMPS data.

361

362 The hourly contributions are aggregated into daily values and plotted as bivariate plots in
363 Figure 7 to assist comparison with the daily plot in Beddows et al. (2015). In that work, the
364 same PMF analysis of the NSD data yielded 4 factors which are represented here again in
365 the bivariate plots. The similarity of both of the polar and annular plots for each of the 4
366 factors justifies our aforementioned factor-fitting assumption. The Secondary and Diffuse
367 Urban are background sources with strongest contributions in the evening and morning.
368 Traffic is strongest for all wind speeds from the East which makes sense since North
369 Kensington is to the West of the city centre of London where traffic is expected to be most
370 dense. Nucleation is also seen to be strongest for those wind direction from the West which
371 are expected to be cleaner, and have a lower condensation sink. NET & Crustal and Fuel
372 Oil are similar to Diffuse Urban suggesting a similar predominant source location in the
373 centre of London. Marine is observed to be strongest for elevated wind speeds for all wind
374 directions which is consistent with the expected strong contribution for all high wind speeds
375 from the South West, as observed in the daily polar plots in Beddows et al. (2015).

376

377



378 **3.4 Composition of Hidden Factor**

379 The Nucleation factor was extracted from the two-step PMF-PMF analysis when forcing the
380 condition of no PM₁₀ contribution through G1 to G6. It might be reasonable to suggest that
381 if the two-step PMF-PMF analysis is repeated and the order of analysis of PM₁₀ and NSD
382 datasets reversed that it would be possible to derive the chemical conditions within the
383 atmosphere which were conducive to nucleation. Ideally, for this the chemical data would
384 be more informed with regards to the composition of the particles below 100 nm. However,
385 when using the PM₁₀ data the Nucleation factor was associated with marine air with strong
386 contributions to Na, Cl and Mg (Figure S4). There are also traces of V, Cr, Ni and a high
387 PM level which are all associated with marine air. This is explained by an association with
388 the south-westerly wind sector which brings strong winds and marine aerosol rather than
389 reflecting the composition of the nucleation particles themselves. Secondary shows a strong
390 association with ammonium, nitrate and sulphate but there are also traces of organics, Al,
391 Cd, Mn, Pb, Ti and Zn and high PM_{2.5} and PM₁₀. Diffuse Urban makes the smallest
392 contribution to PM but shows strong elemental carbon, wood smoke, Ba, Cr, Fe, Mo, Sb, V
393 and Zn; indications of recreational wood burning and brake dust. Traffic has strong
394 associations with Ba, Al, Ca, Cu, Mn, Ti and Zn which have sources in tyre and brake dust
395 and resuspension.

396

397 **4. CONCLUSIONS**

398 It is recommended when applying PMF to atmospheric PM data that only metrics with the
399 same unit are input in order to make a meaningful quantitative apportionment. However,
400 the inclusion of meteorological and particle number data has proved to give a clearer
401 separation of factors. Mixed unit datasets limit the PMF to a qualitative analysis and the
402 quantitative step of apportioning the sources to a mass or number concentration has to be



403 omitted. This problem is overcome in this work by using a novel Two-Step PMF approach.
404 In the first step the PM₁₀ data is PMF analysed using the standard approach without the
405 inclusion of additional data. An appropriate solution is derived using the methods described
406 in the literature in order to give an initial separation of source factors. The time series G
407 (and errors) of the PM₁₀ solution are then taken forward into the second step where they are
408 combined with the NSD data. The PMF analysis is then repeated using the combined and
409 mixed unit G time series and NSD dataset. In order to ensure that unique factors are
410 obtained for the G scores, Fkey is used to pull off diagonal values to zero thus driving the
411 NSD data. This ensures that the NSD factors are specific to the PM₁₀ solution. This results
412 in 6 PM₁₀ factors which are not only apportioned in mass but are augmented by the NSD
413 data. Comparisons of both the factor profiles, diurnal trends and bivariate plots to those of
414 Beddows et al. (2015), show that this technique produces one solution linking the two
415 separated solutions for PM₁₀ and NSD data datasets together. This generates confidence
416 that the NSD and PM₁₀ factors ascribed to one source are in fact attributable to that same
417 source.

418

419 Hence, the process starts with a dataset which produces a solution which is sensitive to
420 mass but the factors more sensitive to number can be accessed using a second step.
421 Furthermore, by exploring a higher number of factors, NSD factors which are insensitive to
422 PM₁₀ mass can be identified as in the case of the Nucleation factor. This information can
423 also be extracted using a linear regression PMF-LR where the size bins of the NSD data are
424 regressed against the PM₁₀ PMF time series. For this dataset, the Nucleation factor profile
425 is identified as an intercept within the fitted model leading to an increase in the number of
426 PMF factors from 6 to 7.

427



428 **5. ACKNOWLEDGEMENTS**

429 The National Centre for Atmospheric Science is funded by the U.K. Natural Environment

430 Research Council. Figures were produced using CRAN R and Openair (R Core Team, 2016;

431 Carslaw and Ropkins, 2012).

432

433

434 **REFERENCES**

- 435 Beddows, D. C. S., Harrison, R. M., Green, D. C. and Fuller, G. W.: Receptor modelling of
436 both particle composition and size distribution from a background site in London, UK, Atmos.
437 Chem. Phys., 15, 10107-10125, 2015.
- 438
439 Beddows, D. C. S., Dall'Osto, M. and Harrison, R. M.: An enhanced procedure for the
440 merging of atmospheric particle size distribution data measured using electrical mobility and
441 time-of-flight analysers, Aerosol Sci. Technol., 44, 930-938 (2010).
- 442
443 Bigi, A. and Harrison, R. M.: Analysis of the air pollution climate at a central urban
444 background site, Atmos. Environ., 44, 2004-2012, 2010.
- 445
446 Carslaw, D. C. and Ropkins, K.: openair - an R package for air quality data analysis, Environ.
447 Model Softw. 27-28, 52-61, 2012.
- 448
449 Chan, Y.-C., Hawas, O., Hawker, D., Vowles, P., Cohen, D. D., Stelcer, E., Simpson, R.,
450 Golding, G. and Christensen E.: Using multiple type composition data and wind data in PMF
451 analysis to apportion and locate sources of air pollutants, Atmos. Environ., 45, 439-449,
452 2011.
- 453
454 Comero, S., Capitani, L. and Gawlik, B. M.: Positive Matrix Factorisation (PMF) An
455 introduction to the chemometric evaluation of environmental monitoring data using PMF,
456 JRC Scientific and Technical Reports, 2009.
457 [https://pdfs.semanticscholar.org/24d3/7d67993228d05513f877b007a16d9a677ff0.pdf?_ga](https://pdfs.semanticscholar.org/24d3/7d67993228d05513f877b007a16d9a677ff0.pdf?_ga=2.219748825.1591505292.1532007606-1458741279.1523619082)
458 [=2.219748825.1591505292.1532007606-1458741279.1523619082](https://pdfs.semanticscholar.org/24d3/7d67993228d05513f877b007a16d9a677ff0.pdf?_ga=2.219748825.1591505292.1532007606-1458741279.1523619082).
- 459
460 Harrison, R. M., Beddows, D. C. S. and Dall'Osto, M.: PMF analysis of wide-range particle
461 size spectra collected on a major highway, Environ.Sci.Technol., 45, 5522-5528, 2011.
- 462
463 Hopke, P. K.: A guide to Positive Matrix Factorization, J. Neuroscience, 2, 1-16, 1991.
- 464
465 Leoni, C., Pokorna, P., Hovorka, J., Masio, M., Topinka J., Zhao, Y., Krumal, K., Cliff, S.,
466 Mikuska, P. and Hopke, P. K.: Source apportionment of aerosol particles at a European air
467 pollution hot spot using particle number size distributions and chemical composition, Environ.
468 Pollut., 234, 45-154, 2018.
- 469
470 Masiol, M., Hopke, P. K., Felton, H. D., Frank, B. P., Rattigan, O. V., Wurth, M. J. and
471 LaDuke, G. H.: Source apportionment of PM_{2.5} chemically speciated mass and particle
472 number concentrations in New York City, Atmos. Environ., 148, 215-229, 2017.
- 473
474 Ogulei, D., Hopke, P. K., Zhou, L., Pancras, J. P., Nair, N., and Ondov, J.M.: Source
475 apportionment of Baltimore aerosol from combined size distribution and chemical
476 composition data, Atmos. Environ., 40, S396-S410, 2006.
- 477
478 Pant, P. and Harrison, R. M.: Critical review of receptor modelling for particulate matter: A
479 case study of India, Atmos. Environ., 49, 1-12, 2012.
- 480
481 R Core Team. R: A language and environment for statistical computing. R Foundation for
482 Statistical Computing, Vienna, Austria, 2018. Available at: <https://www.r-project.org/>.



483 R Core Team. R: A language and environment for statistical computing. R Foundation for
484 Statistical Computing, Vienna, Austria, 2016. Available at: <https://www.R-project.org/>.
485
486 Thimmaiah, D., Hovorka, J. and Hopke, P. K.: Source apportionment of winter submicron
487 Prague aerosols from combined particle number size distribution and gaseous composition
488 data, *Aerosol Air Qual. Res.*, 9, 209-236, 2009.
489
490 Wang, X., Zong, Z., Tian, C., Chen, Y., Luo, C., Li, J., Zhang, G. and Luo, Y.: Combining
491 Positive Matrix Factorization and radiocarbon measurements for source apportionment of
492 PM_{2.5} from a national background site in north China, *Sci. Rep.*, 7, 10648, 2017. doi:
493 10.1038/s41598-017-10762-8.
494
495 Wiedensohler, A., Birmili, W., Nowak, A., Sonntag, A., Weinhold, K., Merkel, M., Wehner,
496 B., Tuch, T., Pfeifer, S., Fiebig, M., Fjaraa, A. M., Asmi, E., Sellegri, K., Depuy, R., Venzac,
497 H., Villani, P., Laj, P., Aalto, P., Ogren, J. A., Swietlicki, E., Williams, P., Roldin, P., Quincey,
498 P., Hüglin, C., Fierz-Schmidhauser, R., Gysel, M., Weingartner, E., Riccobono, F., Santos,
499 S., Gruning, C., Faloon, K., Beddows, D., Harrison, R. M., Monahan, C., Jennings, S. G.,
500 O'Dowd, C. D., Marinoni, A., Horn, H.-G., Keck, L., Jiang, J., Scheckman, J., McMurry, P.
501 H., Deng, Z., Zhao, C. S., Moerman, M., Henzing, B., de Leeuw, G., Loschau, G. and Bastian
502 S.: Mobility particle size spectrometers: Harmonization of technical standards and data
503 structure to facilitate high quality long-term observations of atmospheric particle number size
504 distributions, *Atmos. Meas. Tech.*, 5, 657-685, 2012.
505
506



507 **FIGURE LEGENDS:**

508

509 **Figure 1.** Venn Diagram showing the summary of the findings of Beddows et al. (2015)
510 applying PMF to PM₁₀-only, NSD-only and PM₁₀+NSD datasets.

511

512 **Figure 2.** Flow diagram showing the flow of data through the 2-step PMF-PMF analysis.

513

514 **Figure 3.** Entries in the Fkey matrix used in step 2 of the PMF-PMF analysis. An extremely
515 strong value of 24 was chosen for Fkey.

516

517 **Figure 4.** Second step PMF result. PM₁₀ G score driven PMF with quadruple NSD
518 uncertainties. Also Fkey applied to G score part of F factors to pull off-diagonal elements to
519 zero.

520

521 **Figure 5.** Second step PMF result. 7th factor. PM₁₀ G score driven PMF with quadruple
522 NSD uncertainties. Also Fkey applied to G score part of F factors to pull off-diagonal
523 elements to zero.

524

525 **Figure 6.** Diurnal cycles derived from the hourly NSD data fitted by the daily Factor profiles.
526 Left-left column – diurnal trends of the fitted peaks; left-middle column – bivariate plot of the
527 hourly fitted peaks; middle-right – annular plot of the hourly fitted peaks; right-right – bivariate
528 plot of the daily averaged fitted peaks, plotted using the Openair program.

529

530

531



532

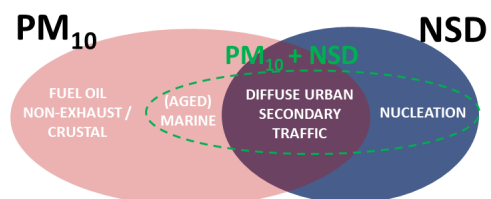


Figure 1. Venn Diagram showing the summary of the findings of Beddows et al. (2015) applying PMF to PM_{10} -only, NSD-only and $PM_{10}+NSD$ datasets.

533

534



535

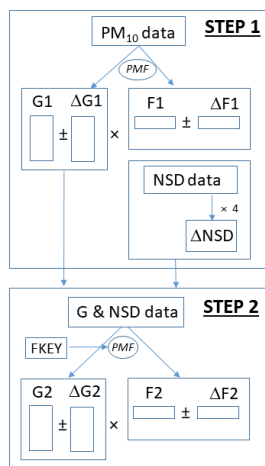


Figure 2. Flow diagram showing the flow of data through the 2-step PMF-PMF analysis.

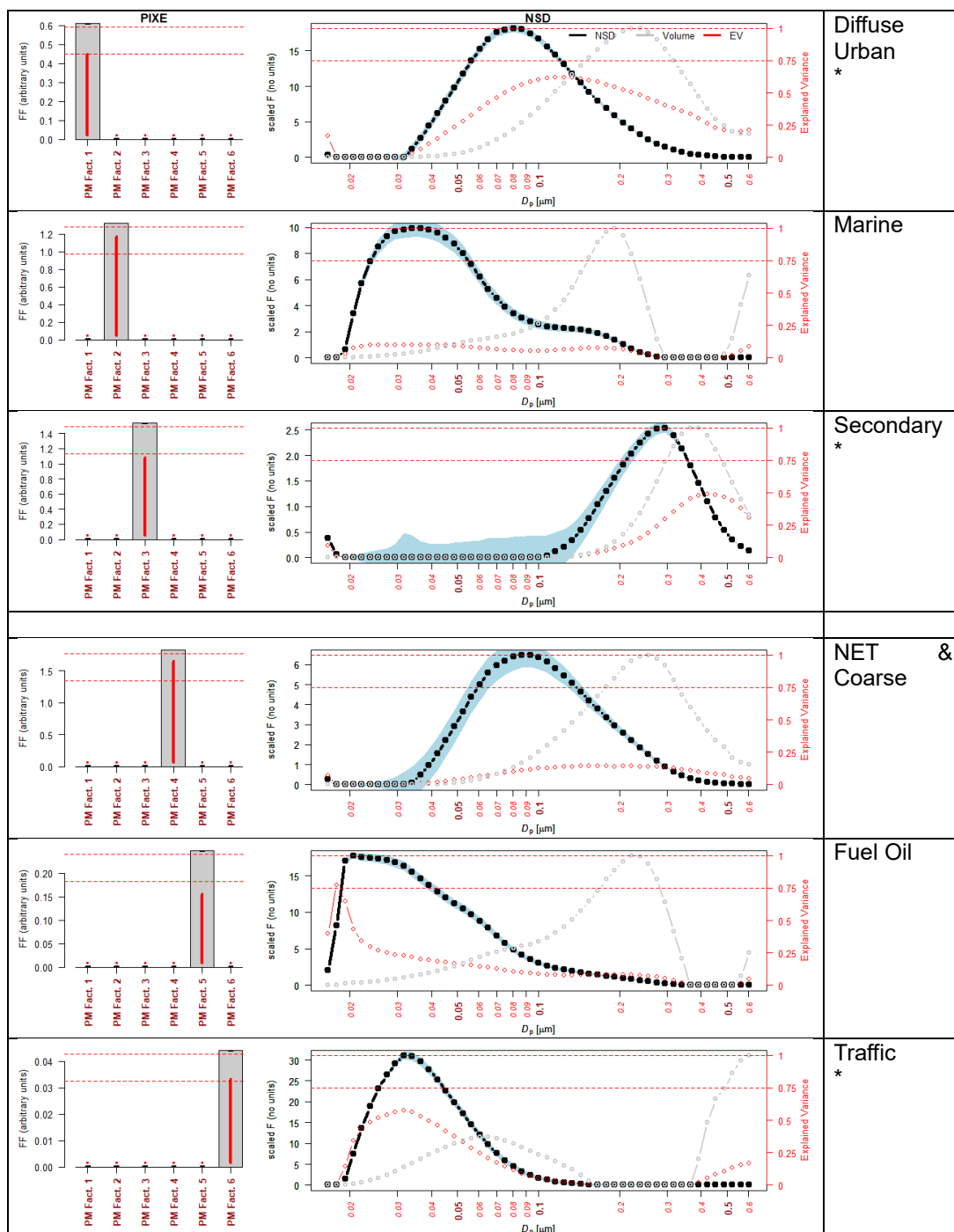
536
537



538
 539

		<i>G</i> Time Series from Step 1						Number Size Distribution (nm)				
		¹ G1	¹ G2	¹ G3	¹ G4	¹ G5	¹ G6	16.6	17.8	19.2	⋮	604
Factors from Step 2	² F ₁	0	FKEY	FKEY	FKEY	FKEY	FKEY	0	0	0		0
	² F ₂	FKEY	0	FKEY	FKEY	FKEY	FKEY	0	0	0		0
	² F ₃	FKEY	FKEY	0	FKEY	FKEY	FKEY	0	0	0		0
	² F ₄	FKEY	FKEY	FKEY	0	FKEY	FKEY	0	0	0		0
	² F ₅	FKEY	FKEY	FKEY	FKEY	0	FKEY	0	0	0		0
	² F ₆	FKEY	FKEY	FKEY	FKEY	FKEY	0	0	0	0		0

540 **Figure 3.** Entries in the Fkey matrix used in step 2 of the PMF-PMF analysis. An extremely
 541 strong value of 24 was chosen for Fkey.
 542

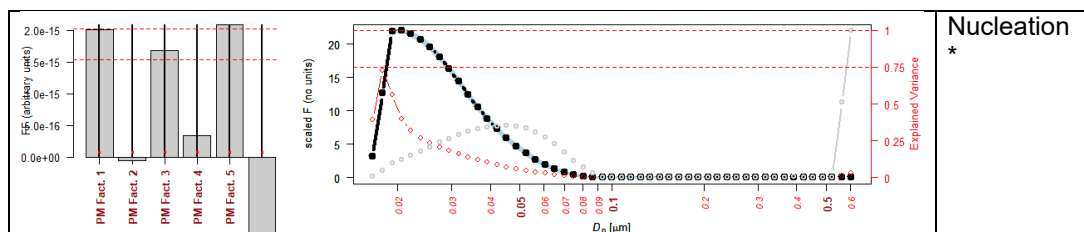


543
 544
 545
 546

Figure 4. Second step PMF result. PM10 G score driven PMF with quadruple NSD uncertainties. Also Fkey applied to G score part of F factors to pull off-diagonal elements to zero.



547
548

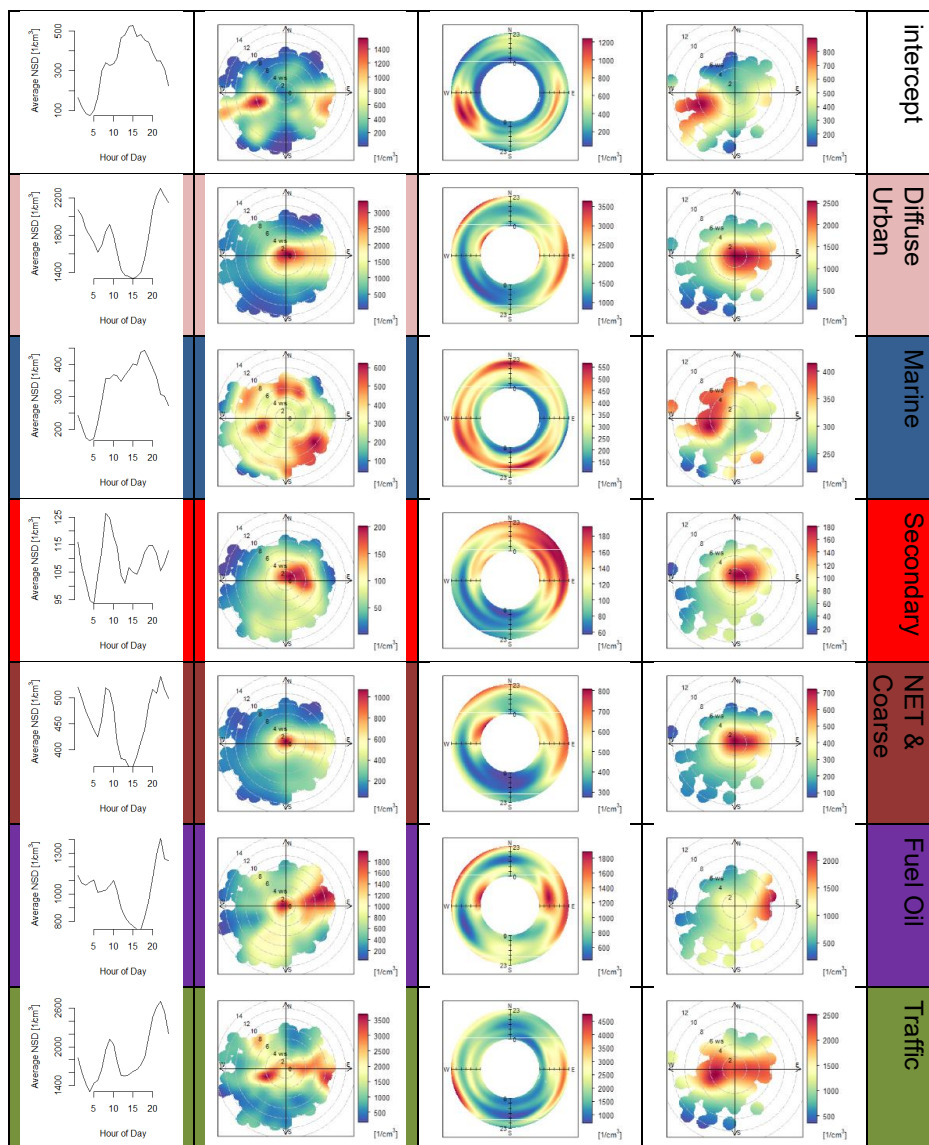


549
550
551
552
553
554
555
556
557
558
559
560
561

Figure 5. Second step PMF result. 7th factor. PM10 G score driven PMF with quadruple NSD uncertainties. Also Fkey applied to G score part of F factors to pull off-diagonal elements to zero.



562
 563



564
 565
 566
 567
 568
 569

Figure 6. Diurnal cycles derived from the hourly NSD data fitted by the daily Factor profiles. Left-left column – diurnal trends of the fitted peaks; left-middle column – bivariate plot of the hourly fitted peaks; middle-right – annular plot of the hourly fitted peaks; right-right – bivariate plot of the daily averaged fitted peaks, plotted using the Openair program.