

1 Applying receptor models Unmix and PMF on real data set of elements in PM 2 for sources evaluation

3 Đorđević Dragana¹, Petrović Srđan², Relić Dubravka³, Mihajlidi-Zelić Aleksandra³

4 ¹ICTM – Centre of Chemistry, University of Belgrade, Studentski trg 12-16, 11000 Belgrade,
5 Serbia (dragadj@chem.bg.ac.rs, phone +381 11 333 68 93, fax +381 11 263 60 61).

6 ²ICTM – Centre of Catalysis, University of Belgrade, Studentski trg 12-16, 11000 Belgrade,
7 Serbia

8 ³Faculty of Chemistry, University of Belgrade, Studentski trg 12-16, 11000 Belgrade, Serbia

9

10 Abstract

11 Two advanced receptor modeling techniques Unmix and PMF were applied to a data set of daily
12 measurements of 11 elements in particulate matters (PM) of 252 samples. Samples were
13 collected every sixth day as a 24h sample in the 5 year period (1995 – 2000) in the coastal part of
14 the Herceg-Novi town (Montenegro) of the sea costal side region (Southeast Adriatic Sea). In the
15 vicinity of the sampling site road traffic is a permanent.

16 The application of the receptor models to find the emission sources in the reverse order, using
17 data set of pollutants concentrations measured on the receptor, is not enough to get satisfactory
18 real solution relying only on the results of the applied models even if used the state-of-the-art
19 models such as Unmix and PMF. In this work we applied Unmix and PMF on dataset which
20 already modeled by PCA and EF in order to show how many solutions could be find and how
21 many errors could be made as well as we harmonized these advanced models to find the most
22 realistic solution. The model Unmix has the ability to suggest the solution by self-modeling
23 while PMF model can be adjusted to calculate the solution for the number of emission sources
24 that we have set. Unmix found thirteen solutions in total for several combinations of species, but
25 four solutions satisfy its criteria: $\text{Min } R^2 > 0.8$ and $\text{Min } S/N > 2$. The PMF model has given 3
26 possible solutions and by further analysis the best solution of four sources was selected. F-peak
27 refinement enabled finding a more realistic solution. We noticed that for the species with many
28 missing values but, their presence is not desirable because of its harmfulness such as cadmium in
29 this work the knowledge of emission sources is very important. Due to their limitations Unmix

30 and PMF is not able to give the solution for such cases. Other simple model applied together
31 with advanced models could help to solve similar problems.

32

33 **Key words:** Modeling, Unmix, PMF, real data set

34 **Introduction**

35 A state-of-the-art multivariate receptor models are applied in the diverse fields of
36 environmentrics, chemometrics, geology and remote sensing. Multivariate receptor modeling is a
37 term applied in the field of air quality for the solution of the general linear mixture problem. For
38 conservative chemical species, i.e. those that do not undergo reactions in the atmosphere, the
39 principle of mass balance is applied. The mass balance for species i can be written as:

$$40 C_{ij} = \sum_{k=1}^N a_{ik} S_{kj} \quad i = 1, \dots, m, j = 1, \dots, n \quad (1)$$

41 In this equation, C_{ij} is the observed concentration of species i in sample k , S_{kj} is the total amount
42 of particulate mass from source k in sample j and a_{ik} it the composition fraction of species i from
43 the source k . In air quality studies, the units of C_{ij} are usually micrograms per cubic meter. Thus,
44 since a_{ik} is a dimensionless mass fraction, the units of S_k are also micrograms per cubic meter.
45 Eq. (1) is the physical basis of all receptor models. C_{ij} is subject to random error and a_{ik} to
46 random variations (Henry, 2002).

47 Unmix seeks to solve the general mixture problem where the data is assumed to be a linear
48 combination of an unknown number of sources of unknown composition, which contribute an
49 unknown amount to each sample. Unmix assumes that the data and the compositions of the
50 sources are all strictly positive (because of the effects of errors, small values less than zero are
51 allowed in order to reduce the bias in the results). Unmix further assumes that for each source
52 there are some samples that contain little or no contribution from that source. For a given
53 selection of species, Unmix estimates the number of sources, the source composition, and source
54 contributions to each sample. The usual analytical approach to fitting the model in Eq. (1) is to
55 find the values of a_{ik} and S_{kj} that minimize the weighted mean square error F (Henry, 2002) of
56 the model:

$$57 F = \sum_{i=1}^m \sum_{j=1}^n (w_{ij} C_{ij} - \sum_{k=1}^N a_{ik} S_{kj})^2 \quad i = 1, \dots, m, j = 1, \dots, n \quad (2)$$

58 Unmix diagnostic edges plots are used to show how well-defined one or more edge is by the
59 data. If the edge plots show that all the edges are straight and well defined, then the Unmix

60 results should be more reliable and should be preferred over the PMF results (Henry and
61 Christensen, 2010).

62 The General Mixture Problem and the special case of multivariate receptor modeling are ill
63 posed problems. There are simply more unknowns than equations and thus there are many wildly
64 different solutions that are all equally good in the least-squares sense. In a statistical way these
65 problems are not identifiable. One approach to ill-posed problems is to impose conditions that
66 add additional equations, which then define more realistic solutions to be closer to unique
67 solution. The non-negativity conditions as additional conditions are imposed by the physical
68 nature of the problem (Henry, 2001). Source composition and contributions must be non-
69 negative but non-negativity conditions alone are not sufficient to give a unique solution. More
70 constraints are needed (Henry, 1987). Under certain, rather mild conditions, the data themselves
71 can provide the needed constraints (Henry 1997). This is how Unmix works.

72 Based on the multivariate factor analysis and the results in factor profiles and contributions,
73 Paatero and Tapper (Paatero and Tapper, 1993; Paatero and Tapper, 1994; Paatero, 1997)
74 established the advanced factor analysis method - positive matrix factorization (PMF). Several
75 features are incorporated in this model:

- 76 - weights data points by their analytical uncertainties,
- 77 - constrains factor loadings and factor scores to non-negative values and thereby minimizes
78 the ambiguity caused by rotating the factors,
- 79 - uses weighted least-squares fits for data,
- 80 - expresses factor loadings in mass units, which allows factors to be used directly as source
81 signatures,
- 82 - provides uncertainties for factor loadings and factor scores.

83 In PMF, the matrix \mathbf{X} ($n \times m$) includes measured mass concentrations, and is represented
84 as the sum of the product of \mathbf{G} ($n \times p$) and \mathbf{F} ($p \times m$) matrices and the residual matrix \mathbf{E} ($n \times m$),
85 where \mathbf{n} is the number of samples, \mathbf{m} is the number of chemical species, and \mathbf{p} is the number of
86 independent source types. This model can give a solution that can be displayed in matrix form:

$$87 \quad \mathbf{X} = \mathbf{G} \cdot \mathbf{F} + \mathbf{E} \quad (3)$$

88 The object function Q that is to be minimized is defined as:

$$89 \quad Q = \sum_{i=1}^n \sum_{j=1}^m (\varepsilon_{ij}/u_{ij})^2 \quad (4)$$

90 where u_{ij} is the uncertainty of the species j in a sample i and residuals ϵ_{ij} i.e. the portion of the
91 measured concentration.

92 In addition, non-negativity constraints should be fulfilled, meaning that all the elements
93 in G and F are to be non-negative. The main process of the PMF is minimizing the Q -value
94 which is defined in the Eq. (4) as the sum of square of the residuals (ϵ_{ij}) weighted inversely with
95 uncertainty (u_{ij}) of the data point (Polissar et al., 1998; Lee and Hopke, 2006).

96 The solution of Eq. (4) is obtained by iteration until convergence is reached.

97 Bootstrapping is an advanced analysis that examines the stability of solutions of the
98 tested models. The bootstrap method is essentially based on resampling methods in which “new”
99 data sets that are consistent with the original data are generated. Each “new” data set (which is
100 essentially a subset of the original database), is decomposed into profile and contribution
101 matrices, and the resulting profile and contribution matrices are compared with the base run
102 (Eberly, 2005), giving the distribution for each species to evaluate the stability of the solution.

103 Numerous studies employing both the PMF and Unmix models have been done in recent years
104 (Pekney et al., 2006; Poirot et al., 2001; Kim et al., 2004; Chen et al., 2007 in: Hegg et al., 2010).

105 Paatero’s positive matrix factorization (PMF) approach weights the data by the inverse of the
106 measurement error for each observation. A major advantage of this approach is that the missing
107 data can be included as observations with a large error. However, the minimization of F is still
108 an ill-posed problem, or in other words, the model is not identifiable. Even the inclusion of the
109 non-negative constraints does not provide an identifiable model. Paatero addresses this problem,
110 which he named rotationally indeterminacy, by adding one or more user-selected parameters.
111 Park et al., (2002) have used modern constrained minimization methods on F along with specific
112 conditions, e.g. each source composition must have at least one species absent from that source.
113 Finally, Paatero has generalized F in a natural way to include the estimation of even more
114 unknown parameters associated with spatial variations (Henry, 2002).

115 Multivariate source apportionment models, Unmix and positive matrix factorization (PMF),
116 often produce nearly the same source apportionment, however some investigations have shown
117 that this is not always the case (Henry and Christensen, 2010). These models do not specify a
118 minimum number of samples, but the stability of their solutions increases with the number of
119 samples (Chen et al., 2007). In this study, we calculated sources composition and sources
120 contributions of elements in PM using real data base.

121 The main aim of this study is to show that a simple application of the most advanced
122 mathematical models may leads to erroneous conclusions because each of these models can
123 provide a larger number of mathematically correct solutions. Which solutions are really true
124 cannot be known only on the basis of the results obtained by modeling, even using models such
125 as Unmix and PMF. Our goal was to apply these state-of-the-art models, respecting their criteria,
126 on data-base previously submitted to other models; Principal Component Analysis (PCA) and
127 Enrichment Factors (EA) to compare, to be able to finding the most accurate solution relying on
128 Unmix self-modeling and PMF application to adjust and confirm the solutions found by Unmix.

129

130 **Materials and Methods**

131 The sampling site is situated only 10 meters away from the coast of the Adriatic Sea. Samples of
132 PM were subjected to gravimetric analysis for determination of total mass concentrations and
133 subsequently to elemental analysis for Fe, Mn, Ti, Pb, Cr, Cu, Cd, Co, Ni, Hg and Se. Suspended
134 particles were collected using a high-volume Aerosol Sampler, AQUERO model 400XT
135 sampling system with inlet for the total suspended particles, on boron-silicate fiberglass filters
136 every sixth day as a 24h sample in the period of 1995 - 2000. The sampler was located in the
137 town of Herceg-Novı (Fig. 8) 18⁰33" N, 42⁰27", Montenegro (Fig. 1). The meteorological station
138 is part of the MED POL program. The nearest road is located about 100 m north of the
139 meteorological station. There are no significant grassy areas around the meteorological station,
140 and there is no considerable construction work in progress. The terrain surrounding the receptor
141 is rocky with some small areas of soil (Đorđević et al., 2004). Filters were digested with HNO₃
142 (ultra pure). A Flame Atomic Absorption Spectrometry (F-AAS), Varian AAS–Spectr AA 55
143 instrument, was used to measure the concentrations of Cd, Co, Cr, Cu, Ni, Pb, Ti, Fe and Mn.
144 The concentrations of Hg and Se were determined by the hydride vapor AAS method (HV-AAS)
145 (Đorđević et al., 2005). The maximum expanded uncertainty of measurements for all elements
146 was about 5%.

147 The real data set of 11 trace elements in particulate matter (PM) obtained in 252 observations
148 was analyzed by Unmix 6.0 and PMF 3.0. The applied Unmix and PMF models were available
149 on the EPA Internet site (www.epa.gov).

150 Unmix and PMF used in this study do not limit the number of factors. The following initial
151 operations were subjected to the Unmix model data: *Suggest Exclusion, Initial Species,*
152 *Additional Species* including SAFER and Initial Points. PM was chosen for the total and for the
153 normalization. The data was screened using the signal-to-noise ratio (Min S/N ratio) criteria
154 higher than 2, estimated by Unmix. Only the component with S/N value greater than 2 will be
155 used for sources estimation. The agreement between the true and estimated source contribution
156 (Min R² greater than 0.8) was considered as well (Henry, 2003, EPA/600/R-07/089).

157 Applying of PMF model the procedure of Polissar et al. (1998) was used in this study to
158 calculate uncertainties in the species concentrations. Briefly, for the data below detection limit
159 (DL), the concentrations were replaced with the value DL/2 and the uncertainty was set as $\frac{5}{6}DL$.

160 For the missing data, concentrations were replaced by the geometric mean and the respective
161 uncertainty was set at four times of this mean concentration. At the first set up all elements are
162 labeled as Strong, since (the signal/noise ratio) $S/N > 2$ for all of them. Based on input data
163 statistics, residuals show bimodal distribution in the case of Ni, Mn, and Hg, so their
164 uncertainties are increased labeling them as Weak. Selenium is excluded from the model because
165 of a very small contribution and the correlation factor, while for cadmium more than 50% of
166 samples are below the detection limit. The Q value represents the goodness-of-fit and assesses
167 how the model fits the experimental data. Q_{true} is calculated taking into account all data points
168 while Q_{robust} is calculated accounting for outlier points. Data with scaled residuals above 4 are
169 regarded as outlier points. Evaluation of the validity of a solution is possible by using the G-
170 space scatter plot. Scatter plot of one versus the other factor may indicate the existence of a
171 rotational ambiguity. Namely, if the points on this graph fill the entire solution space evenly then
172 the edges of the Scatter Plot correspond to axes. If this is not a case it is indication that there is
173 rotational ambiguity and should be considered the possible rotation of the solution, using the
174 function Fpeak. The F-peak functions is used to rotate the data set, make fine tuning and
175 improvement of the model in the case of data with high noise (positive values F-peak) or clean
176 data (negative values F-peak). Normally, the default settings give satisfactory results, but in
177 some cases subsequent adjustments are needed. To ensure the robustness of statistics, 300
178 bootstrap runs were performed, while the default value of the minimum correlation (R-Value) of
179 0.60 was used.

180

181

182 Results and Discussion

183 We applied the Unmix and PMF models on dataset from our previous work regarding trace
184 elements in the PM (Đorđević et al., 2005). Fig. 2 shows the comparison of measured and the
185 predicted concentrations of trace elements in PM through time series and Min R^2 . Model Unmix
186 did not calculate R^2 values for Cd, Co, Hg and Se and neither satisfactory solution included these
187 elements since these variables contain a large number of missing values and outliers. Min R^2
188 values are given in table 1.

189 From statistical parameters displayed for each species, after input data and the following
190 operations: *Suggest exclusion*, *Influential points*, *Initial species*, *Additional species* and *SAFER*,
191 Unmix finds six combinations of species that give any kind of solution (Table 2). Min S/N for
192 each principal component and Min R^2 of all combinations of elements estimated by Unmix was
193 selected as good solutions that are in accordance with the Unmix criteria (Henry, 2003). Thirteen
194 solutions in total were found, but four solutions satisfy the above criteria, signed in bold in Table
195 2. The standard deviation of variable (σ) is the criterion for evaluation whether the variable
196 eligible for modeling or not. The sigma-based parameters (Significant/Strong Species in Sources)
197 for each of satisfactory solution are also given in Table 2.

198 Taking into account the calculated good solutions presented in Table 2, the Edges plots were
199 done for these solutions (Fig. 3). The source profile of the solutions chosen according to the
200 criteria $S/N > 2$ and $R^2 > 2$ are given in Fig. 4.

201 In the first solution (combination of species Mn-Ti-Pb-Cr-Cu, 3 Sources Solution) the second
202 and third source are well defined by many points on the y-axis while source 1 has just a few
203 points on the x-axis. Pb is strong in the first source and this source can be attributed to traffic. In
204 the second source Cr and Cu are strong and Ti and Mn are significant. This source can be re-
205 suspension of elements previously settled from anthropogenic sources. In the third source neither
206 element is strong or significant.

207 The second satisfactory solution is for Fe-Mn-Ti-Pb combination of elements it also found 3
208 sources (Table 2, Fig. 4b) and does not show good accumulation of points on the x and y axes
209 (Fig. 3). This solution has the best values of Min R^2 and Min S/N compared to all combinations.
210 The first and the third source contain Pb which is a tracer for traffic. In the third source Pb is

211 strong, and it is reasonable to associate this source with traffic, while the first source could be
212 local re-suspension. The second source in this combination could be a long range transport of
213 Saharan dust since it contains crustal elements.

214 The third solution (combination of species Fe-Mn-Ti-Pb-Cr-Cu-Ni, 4 Sources Solution) shows
215 the edges on the y-axis defined by many points for the third and the fourth source, but the x-axis
216 has just a few points (Fig. 3).

217 In the fourth solution (combination of species Fe-Mn-Ti-Pb-Cr-Cu, 3 Sources Solution) good
218 accumulation of points are on the y-axis, for sources 2 and 3 while the x-axis has just a few
219 points for source 1 (Fig. 3).

220 In the third and the fourth solution, the sources where Pb is strong can be attributed to traffic;
221 namely, source 3 for Fe-Mn-Ti-Pb-Cr-Cu-Ni combination (Table 2, Fig. 4c) and source 1 for Fe-
222 Mn-Ti-Pb-Cr-Cu combination (Table 2, Fig. 4d). Another source in which Pb is present as
223 significant but not strong could be re-suspension. Source 4 for Fe-Mn-Ti-Pb-Cr-Cu-Ni
224 combination and source 2 for Fe-Mn-Ti-Pb-Cr-Cu combination could be attributed to re-
225 suspension, probably from various locations depending on wind directions. Factors containing
226 Cr and Ni indicate the existence of an anthropogenic emission source in the region (Đorđević et
227 al., 2005).

228 In our previous work (Đorđević et al., 2005) we applied the PCA method on this data set and 4
229 significant groups of sources contributions were found. The following contribution sources were
230 identified: re-suspension combined with re-suspended Saharan dust that had previously settled
231 (Fe, Mn, Ti) and settled combustion products mostly from traffic, and probably some local
232 stationary source (Cu, Pb). The remaining three factors represent the following combinations F2
233 by Cr and Ni, F3 by Cd and Se and F4 by Hg and Co.

234 The EF model revealed that in the region of the investigated receptor, the main contribution
235 source of Fe, Mn and Ti is the process of local re-suspension and that local re-suspension has no
236 influence on the content of Se in the atmospheric aerosol. The re-suspension is the dominant
237 emission source of Cd from the south-southeast direction from the nearby peninsula (Luštica) but
238 this source is not permanent (Đorđević et al., 2005).

239 The application of positive matrix factorization (PMF) to solve the number and profile of the
240 sources applied to the same database resulted in obtaining possible solutions for 3, 4 and 5
241 sources. For 6 or more sources the model does not find the convergence of the functions Q,

242 which implies that the model did not find any minima. Varying simulation conditions did not
243 contribute to significant improvement, even when the uncertainty is significantly increased.
244 Therefore possible solution should be sought among three possible cases.

245 Table 3 shows the categories of elements and the R^2 values for each of the three possible
246 solutions.

247 Each of the possible solutions obtained by PMF analysis will be considered. Fig. 5 shows F peak
248 strengths for 3 sources solution (Fig. 5a), 4 sources solution (Fig. 5b) and 5 sources solution (Fig.
249 5c).

250 3 Sources Solution: The relatively good correlation was obtained only for Cr and Pb, while
251 bimodal distribution is still present in the case of Co, Ti and Fe. Also, significant outliers are
252 present in the model. In addition, G-Space plots show considerable rotational ambiguity between
253 the sources 1 and 3.

254 Rotational ambiguity, which was found between the sources 1 and 3, decreases when the value
255 of Strength factor reaches -1.2 (Fig. 5a). This is mostly reflected in the increase of Ti
256 concentration in the source 2. However, such large values for F_{peak} is unlikely because the
257 quality of the fit decreases rapidly. The usually dataset rotations are generally much smaller and
258 they are close to the basic solution.

259 However, a small degree of correlation between the model and database indicates that the model
260 with three sources is insufficient to adequately describe a number of sources.

261 In this case, only Co, Ni and Fe show relatively good interquartile range of about 20%, while
262 other species show considerable variation and therefore represent a less stable solution. This is
263 especially pronounced in the case of Hg, Cr and Mn. Also, in some cases (Hg, Ti) base run
264 values are not within the interquartile range in the bootstrapping of results. This is probably a
265 consequence of assuming the model with only three sources. Profiles of sources are given in Fig.
266 6.

267 4 Sources Solution: The model with four sources shows a significantly better correlation with
268 measured concentrations of elements. Although the agreement of time series for Ti and Cr is
269 excellent ($R^2 > 0.95$), and for Pb satisfactory ($R^2 = 0.70$), in the case of other elements there are
270 still episodes with very high concentrations that this model cannot fit. It should be noted that Cu
271 shows very good agreement between the predicted and observed concentrations, but the

272 existence of outliers have reduced the correlation to 0.34. A small degree of correlation in the
273 case of Co is the result of a significant number of measurements below the detection limit.

274 Bimodal distribution is still present in the case of Ni and Hg. G-Space plot shows that there is a
275 rotational ambiguity between sources 1 - 3 2 - 3, 3 - 4.

276 For a model with four sources, rotational ambiguity disappears when the F-Peak strength reaches
277 -0.8 (Fig. 5b). This rotation is mostly reflected in the increase of Ti content in source 1, and
278 largely in sources 2 and 3. On the other hand, this may just mean that the content of titanium in
279 this solution is divided among several sources. As in the case of a solution with three sources, a
280 significant rotation of the dataset ($F_{peak} = -0.8$) is less likely. It is necessary to consider these
281 results carefully and determine whether there is justification for it to be included in further
282 solving of the composition of the sources.

283 Interquartile range of solutions obtained by bootstrapping in the case of Fe, Pb, Cu and Cr are
284 about 20%, while in the case of other species this range is much higher indicating the instability
285 of the solution. Base run values which are not within the interquartile range in the bootstrapping
286 of the results are, in the case of Cu, Mn, Pb and especially Hg, calculated by the model only in
287 the fourth source.

288 5 Sources Solution: The model and data from the database show agreement (R^2) over 90% for
289 Cr, Ti, Fe and Pb, while just over 50% for Mn. The model also fits Cu real data very well, and
290 the correlation of 0.47 is caused by significant outliers that are related to individual episodes of
291 Cu emissions. In spite of the increased uncertainty Mn, Ni, Co and Hg show a lack of fit.

292 G-Space plot only shows some rotational ambiguity in the case of sources: 2 - 5, 3 - 5 and 4 - 5.
293 The F-Peak in the range -2 to 2 (Fig. 5c) showed the most impact on the sources of Cu,
294 especially at higher strength values, while the ambiguity between the sources mentioned above
295 still exist. The Peak F-curve is generally symmetrical in the examined interval.

296 In the case of five sources there are also unmapped results, suggesting a reduced stability of the
297 solution. In general, the most stable solutions are obtained for those elements that are present in
298 the source with the highest percentage. In these cases, the distribution of solutions obtained by
299 bootstrapping lie in the range of 15% of the concentration calculated in the base run. This is the
300 case for Fe, Pb, Ti. A slightly worse result of the bootstrap analysis is obtained for Cu, Mn, Cr,
301 Co and Ni (bootstrapping distribution of solutions equal or higher than 20%). The least stable
302 solution is for Hg with considerable dispersion in the bootstrap analysis solutions.

303 When discussing the number and origin of pollution sources, it is preferred to take into account
304 the real situation on the field. In this case the following sources that contribute to the overall PM
305 deposition can clearly be predicted: marine aerosols, traffic, re-suspension from the ground,
306 probably some local stationary source, as a shipyard located in the vicinity. Based on these
307 obvious sources, PMF analysis solution with only three sources is exempt from further
308 consideration.

309 In the case of PMF solution of five sources it may be noted that source No. 5 (Fig. 8), in which
310 Co, Cu, Ni and Mn are present, can be described rather as a splitting factor than as a separate
311 source. The most realistic solution that is imposed upon a detailed analysis is the solution with
312 four sources (Fig. 7).

313 Identification of sources was carried out and it agrees with the results of the Enrichment Factors
314 analysis well (Đorđević et al., 2005). The F-peak profiles shown in Fig. 7 in rotation of data set
315 for -0.8, increase the contents of Ti, in the case of sources 2 and 4.

316 Source 1 has been identified as re-suspension in combination with the long-range transport of
317 Saharan dust. The prevailing wind directions are over open sea (Đorđević et al., 2005).

318 Source 2 is attributed to the re-suspension, indicated in our previous work. Titanium found by F-
319 peak is in better accordance with the EF analysis (Đorđević et al., 2005).

320 Source 3 corresponds to the composition of the particles that come from some anthropogenic
321 source.

322 Source 4 with the highest content of Pb, is characteristic for urban traffic. F-peak is increasing
323 the value for Ti which is in better agreement with the traffic profile.

324

325 **Conclusion**

326 In this work we applied state-of-the-art mathematical models Unmix and PMF on database
327 previously modeled using more simple models (PCA and EF) to be compared. In this study we
328 have shown that only application Unmix and PMF for sources apportionment is not guarantee to
329 obtain the unique realistic solution. Thirteen solutions in total were found, but four solutions
330 satisfy the Unmix criteria: three solutions with three sources and one with four sources. In terms
331 of modeling all four solutions found by Unmix are satisfactory. PMF model has given three
332 possible solutions: one with three sources, one with four sources and one with five sources. By
333 further analysis of the results of PMF model the best solution with four sources was selected. F-

334 peak refinement was enabled to find a more realistic solution. Also we have shown that due to
335 their limitations Unmix and PMF were unable to calculate Cd and Se in used database, due to
336 large number of missing values. For example, although the presence of cadmium in terms of
337 concentration is negligible and there are many missing values the knowledge of emission sources
338 is very important regarding its harmfulness. The simple model of EF applied could help to solve
339 similar problem. For obtaining the best results using Unmix and PMF models our
340 recommendation is to start modeling by Unmix relying on its self-modeling to estimate all
341 possible types of sources and then apply PMF for confirmation. For the species that are
342 important and that cannot be modeled by advanced models like Unmix and PMF should be apply
343 other, even, the simple model.

344

345 **Acknowledgment**

346 The authors gratefully acknowledge the financial support of the Ministry of Education and
347 Science of the Republic of Serbia, which supported this research within the project 172001. The
348 authors are gratefully acknowledged to Saša Savić as well, for language improving.

349

350 **References**

351 Chen, L.-W. A., Watson, J. G., Chow J. C., and Magliano, K. L., 2007. Quantifying PM_{2.5}
352 Source Contributions for the San Joaquin Valley with Multivariate Receptor Models, Environ.
353 Sci. Technol., 41, 2818-2826.

354 Đorđević, D., Vukmirović, Z., Tosić, I., and Unkasević, M.,2004. Contribution of dust transport
355 and resuspension to particulate matter levels in the Mediterranean atmosphere, Atmos.
356 Environ., 38, 3637–3645.

357 Đorđević, D., Mihajlidi-Zelić, A., and Relić, D.,2005.Differentiation of the contribution of local
358 resuspension from that of regional and remote sources on trace elements content in the
359 atmospheric aerosol in the Mediterranean area, Atmos. Environ., 39, 6271–6281.

360 Eberly, S.,2005. EPA PMF 1.1 user's guide. Prepared by the U.S. Environmental Protection
361 Agency, National Exposure Research Laboratory, Research Triangle Park, NC,
362 June,http://www.epa.gov/heads/products/pmf/users_guide_old.pdf

363 Norris, G. A., Vedantham R., and Duvall R. M., 2007. EPA UNMIX 6.0 USER GUIDE. U.S.
364 Environmental Protection Agency, Washington, DC, EPA/600/R-07/089 (NTIS PB2007-
365 112630).

366 Hegg, D. A., Covert, D. S., Jonsson, H. H., and Woods, R. K., 2010. The contribution of
367 anthropogenic aerosols to aerosol light-scattering and CCN activity in the California coastal
368 zone, *Atmos. Chem. Phys.*, 10, 7341-7351.

369 Henry, R. C.,1987. Current Factor Analysis Receptor Models are Ill-Posed, *Atmos. Environ.*, 21,
370 1815-1820.

371 Henry, R. C.,1997. History and Fundamentals of Multivariate Air Quality Receptor Models,
372 *Chemometr. Intell. Lab.*, 37, 525-530.

373 Henry, R. C. UNMIX Version 2.4 Manual; U.S. Environmental Protection Agency: Research
374 Triangle Park, NC, June 2001.

375 Henry, R. C., 2002. Multivariate receptor models-current practice and future trends, *Chemometr.*
376 *Intell. Lab.*, 60, 43-48.

377 Henry, R. C., 2003. Multivariate receptor modeling by N-dimensional edge detection,
378 *Chemometr. Intell. Lab.*, 65, 179-189.

379 Henry, R. C. and Christensen, E. R.,2010. Selecting an Appropriate Multivariate Source
380 Apportionment Model Result, *Environ. Sci. Technol.*, 44, 2474-2481.

381 Kim, E., Hopke, P. K., Larson, T. V., and Covert, D. S.,2004. Analysis of ambient particle size
382 distributions using Unmix and positive matrix factorization, *Environ. Sci. Technol.*, 38, 202-
383 209.

384 Lee, J. H. and Hopke, P. K., 2006. Apportioning sources of PM_{2.5} in St. Luis, MO using
385 speciation trends network data, *Atmos. Environ.*, 40, S360–S377.

386 Park, E. S., Spiegelman, C. H., and Henry, R. C.,2002. Bilinear estimation of pollution source
387 profiles and amounts by using multivariate receptor models, *Environmetrics*, 13, 775-798.

388 Paatero, P., and Tapper, U.,1993. Analysis of different modes of factor analysis as least square fit
389 problem,*Chemometr. Intell. Lab.*, 18, 183–194.

390 Paatero, P. and Tapper,U.,1994. Positive Matrix Factorization: a non-negative factor model with
391 optimal utilization of error estimates of data values, *Environmetrics*, 5, 111–126.

392 Paatero, P.,1997. Least squares formulation of robust non-negative factor analysis, *Chemometr.*
393 *Intell. Lab.*, 37, 23-35.

394 Pekney, N. J., Davidson, C. I., Robinson, A., Zhou, L., Hopke, P., Eatough, D., and Rogge, W.
395 F.,2006. Major source categories for PM_{2.5} in Pittsburgh using PMF and UNMIX, *Aerosol*
396 *Sci. Tech.*, 40, 910-924.

397 Poirot, R. L., Wishinski, P. R., Hopke, P. K., and Polissar, A. V.,2001. Comparative application
398 of multiple receptor methods to identify aerosol sources in Northern Vermont, *Environ. Sci.*
399 *Technol.*, 35, 4622-4636.

400 Polissar, A. V.,Hopke, P.K.,Paatero, P., Malm, W. C., andSisler,J. F.,1998. Atmospheric aerosol
401 over Alaska 2. Elemental composition and sources, *J. Geophys. Res.*, 103, 19045–19057

402 U.S. Environmental Protection Agency. EPA Unmix 6.0 Model:
403 <http://www.epa.gov/heads/products/Unmix/Unmix.htm>, access: June 2007.

404 U.S. Environmental Protection Agency. EPA Positive Matrix Factorisation 3.0:
405 http://www.epa.gov/heads/products/pmf/pmf_registration.htm, access: July 2008.

406

407 **Tables**

408 Table 1. R^2 values obtained by Unmix of measured and the predicted concentrations of PM and
409 trace elements in PM

	PM	Fe	Mn	Ti	Pb	Cr	Cu	Ni
R^2	0.46	0.75	0.92	0.66	0.83	0.40	0.99	0.00

Table 2. All combination of elements for solutions obtained by calculation by Unmix

Combination of species	Number of sources	Min R ²	Min S/N	Significant/Strong Species in Sources (sigma-based)
Mn-Ti-Pb-Cr-Cu	3	0.84	2.49	Source 1: *Strong – Pb; Source 2: *Strong - Cr, Cu, **Significant - Ti, Mn; Source 3: *Strong – None, **Significant - None
	4	0.89	1.94	
	5	0.90	1.59	
Cr-Cu-Pb-Ti-Mn-Se-Cd-Co-Fe	3	0.68	2.41	
	4	0.76	2.13	
Cu-Ti-Fe-Mn-Pb-Cr-Hg-Se	2	0.56	2.29	
Fe-Mn-Ti-Pb	3	0.90	2.85	Source 1: *Strong – None, **Significant - PM, Pb, Fe Source 2: *Strong – None **Significant - PM, Fe, Mn; Source 3: *Strong - Pb, Mn, **Significant - Ti, Fe
	3	0.76	2.20	
Fe-Mn-Ti-Pb-Cr-Cu-Ni	4	0.83	2.18	Source 1: *Strong - Cr, Ni, **Significant – None Source 2: *Strong – None, **Significant – Cu; Source 3: *Strong – Pb, **Significant - PM, Cr, Cu, Ti, Fe, Mn; Source 4: *Strong – Ti, **Significant - PM, Cr, Pb, Fe
	5	0.89	1.67	
	3	0.83	2.57	Source 1: *Strong - Pb, **Significant – Cu; Source 2: *Strong – None, **Significant - PM, Cr, Pb, Ti, Fe; Source 3: *Strong – None, **Significant - PM, Cr, Cu, Ti, Fe, Mn
Fe-Mn-Ti-Pb-Cr-Cu	4	0.88	1.97	
	5	0.90	1.62	

* Source Composition ≥ 1 sigma

**Source Composition ≥ 2 sigma

Table 3. R^2 values obtained by PMF of measured and the predicted concentrations

Species	Category	R^2		
		3 sources	4 sources	5 sources
Cr	Strong	0.617	0.980	0.998
Ti	Strong	0.381	0.962	0.959
Fe	Strong	0.401	0.472	0.942
Pb	Strong	0.695	0.701	0.905
Mn	Weak	0.394	0.340	0.528
Cu	Strong	0.345	0.337	0.473
Ni	Weak	0.027	0.031	0.025
Co	Weak	0.006	0.013	0.018
Hg	Weak	0.002	0.005	0.003

Figure Legends

Fig.1 Sampling site and prevailing wind directions

Fig. 2 Predicted and measured concentrations

Fig. 3 Edge plots for chosen solutions that satisfy the conditions of Min S/N and Min R^2

Fig. 4 Source profiles for selected solutions that are in accordance with the Unmix criteria

Fig. 5 F-peak analysis for three a), four b) and five c) source solutions. The red mark represents the value of F-peak Strength, at which the rotational ambiguity disappears.

Fig. 6 Profiles in the case of three sources solutions. Comparison of base run profile and F-peak run profile with the strength of -1.2 (disappearance of rotational ambiguity).

Fig. 7 Profiles in the case of three sources solutions. Comparison of base run profile and F-peak run profile with the strength of -0.8 (disappearance of rotational ambiguity).

Fig. 8 Profiles in the case of three sources solutions. Comparison of base run profile and F-peak run profile with the strength of -2.0 where it can be seen that F-peak Strength does not affect the existing rotational ambiguity.

Interactive comment on

“Applying receptor models Unmix and PMF on real data set of elements in PM for sources evaluation of the sea coastal side region (Southeast Adriatic Sea)” by D. Đorđević et al.

Anonymous Referee #2

Received and published: 5 November 2013

In the study presented by Dordevic et al., a characterization of the possible sources of TSP (?) was carried out in a coastal area (Southeast Adriatic Sea). An attempt to identify the PM sources is made using two receptor modelling techniques (UNMIX and PMF). The results obtained by both techniques are compared with the results obtained using other models like Enrichment Factors (EF) and Principal Component Analysis (PCA). From my point of view the main problem of the article, is that it is not clear what the objective of the manuscript is. It seems that both models (PMF and Unmix) are used to find a solution similar to that found in the previous article written by the authors (Dordevic et al 2005). If that is the aim then the article does not add anything new to the scientific knowledge which already exists on this subject. If however the overall objective is to present a new attempt to identify the PM sources in the region, the article should be re-written in a different way. In this case, the results of the best solution of both models can then be compared with the result obtained using the other models.

In this manuscript were not only made an attempt to identify the PM sources using two receptor modeling techniques (UNMIX and PMF) and these results compared with the results obtained using other models like Enrichment Factors (EF) and Principal Component Analysis (PCA) but our main goal was to show how many mistakes could be made during the modeling even using the state-of-the-art models such as Unmix and PMF. We corrected the manuscript in this manner including the correction of Abstract, Main goals, Discussion and the Conclusion as well as the title of the Manuscript.

Another thing that concerns me is the use of only 11 elements for the determination of the possible sources affecting the study region. Some of the 11 sources are normally below the MDL (Co, Hg, Cd, Se), so the analysis is mainly done with 7 elements only. I am not sure that the results obtained are meaningful. This should be commented throughout article. I have serious doubts that this article is suitable for publication in this journal. I recommend a major revision of the article including a complete restructuring of the contents.

This data base was chosen because of measurements over longer period of time allow the greater accuracy than one obtained in a shorter time period. We also, modeled other data sets with large number of variables but we didn't get

satisfactory solutions, even we didn't get any solution due to a large number of excluded variables by the model because of a lot of outliers or missing values. So this data set has got the most representative solutions although with 11 elements.

Specific comments: The abstract should be rewritten. The location of the measurements, duration of the sampling campaign, sampling time, PM inlet used should be included. The remainder of the abstract should be completed once the aim of article is clarified. The same can be applied to the introduction. Both models are described in quite a lot of detail. That can be replace for some references.

The Abstract has rewritten and the location of the measurements, duration of the sampling campaign, sampling time added. PM inlet have included in the Material and Methods section in the part P4948L6. In the manuscript it is highlighted in yellow.

The Introduction has shortened.

Part of the introduction, like the Polissar criteria, should be in the experimental section.

It is moved to Material and Methods part along with another the paragraph which belong to this section

The procedure of Polissar et al. (1998) was used in this study to calculate uncertainties in the species concentrations. Briefly, for the data below detection limit (DL), the concentrations were replaced with the value $DL/2$ and the uncertainty was set as $\frac{5}{6}DL$. For the missing data, concentrations were replaced by the geometric mean and the respective uncertainty was set at four times of this mean concentration.

...

At the first set up all elements are labeled as Strong, since (the signal/noise ratio) $S/N > 2$ for all of them. Based on input data statistics, residuals show bi modal distribution in the case of Ni, Mn, and Hg, so their uncertainties are increased labeling them as Weak. Selenium is excluded from the model because of a very small contribution and the correlation factor, while for cadmium more than 50% of samples are below the detection limit.

Depending on the aim, the rest of the introduction should be changed to include more important references. A detailed main goal should appear at the end of this section. The materials and methods section should be expanded. Detailed information about the sampling point, the PM inlet used, sampling period, etc... are needed. A description of how the models are configured is also needed (Initial operations of the Unmix model,

data preparation for the PMF model, S/N ratio, Min R², IM, IS, G-space plots, ...). A table with the data (geometric mean, max, min, standard deviation, number of samples under the mdl) should be added as well.

It is not usual to use the values of FPEAK of -0.8, -1 and -1.2 to eliminate the rotational ambiguity. An explanation for that should be added because it is not evident how you have used Fpeak and the G-space plots. What should be done is to use G-space plots of the base run to identify possible rotations in the solution.

It is corrected:

In addition, G-Space plots show considerable rotational ambiguity between the sources 1 and 3. ... However, such large values for Fpeak are unlikely because the quality of the fit decreases rapidly. The usually dataset rotations are generally much smaller and they are close to the basic solution.

...

G-Space plot shows that there is a rotational ambiguity between sources 1 - 3 2 - 3, 3 - 4. As in the case of a solution with three sources, a significant rotation of the dataset (F peak = -0.8) is less likely.

Corresponding G-space plots of F peak solution factors should be examined to see if any edges viewed in the base runs are more or less evident in the F peak runs. Additionally, profiles and contributions should be examined for species that deviate from the base run to ensure that they are reasonable.

It was observed that all the basic solutions cases of 3, 4 and 5 sources, the analysis of G-space plots there are couples of factors which indicate rotational ambiguity. It is improved – highlighted in the yellow in the new Manuscript. In the new part Results and Discussion in the part of PMF modeling for results for 3, 4 and 5 sources now is existing. These cases were analyzed using F-peak function and the solutions were analyzed. It is improved – highlighted in the yellow in the new Manuscript. In the new part Results and Discussion in the part of PMF modeling for results for 3, 4 and 5 sources now is existing.

The conclusion should also be rewritten once the objective of the article has been clarified. The conclusion is not a summary of the results and discussion sections. It should answer the objective proposed in the introduction. About the figures: Figure 1 is very small and the units are not present. Figure 5, 6, 7: There are some errors in the axis labels. The errors of the concentrations should be added. The caption in Figure 6 should state “four” instead of “three”.

The conclusion is corrected

The errors in the axis labels of Figures 5, 6, and 7 are corrected, now these Figures are 6, 7 and 8.

Interactive comment on Atmos. Meas. Tech. Discuss., 6, 4941, 2013.