

## H14-182

## IMPROVING SOURCE APPORTIONMENT WITH RECEPTOR MODELS TO FOSTER THE IMPLEMENTATION OF AIR QUALITY DIRECTIVE

Federico Karagulian, Claudio A. Belis

European Commission, JRC, Unit Climate Change and Air Quality, Ispra, Italy

**Abstract:** Receptor Models identify pollution sources by solving a mass balance equation using measured chemical composition of samples in combination with known source profiles. In the last ten years, the number of receptor modelling studies performed on filed air quality data exponentially increased. Among these, 39% are performed with Positive Matrix Factorization (PMF, ME), 23% with Principal Component Analysis (PCA, APCA), 13% with Chemical Mass Balance (CMB), 10% with Factor Analysis (FA, APCFA) and 10% with the APEG model.

This approach has been extensively used in North America and South-eastern Asia mainly on particulate matter. Also in Europe it found wide acceptance and contributed to the identification of sources in support of remediation measures design.

With the aim of harmonizing the activity in receptor modelling in Europe and supporting the implementation of Directive 2008/50/EC, an evaluation of the current state of this methodology in Europe was carried out highlighting the following needs: improving data collection, harmonizing analytical protocols by setting up common criteria, promoting advanced tools, establishing criteria for the assessment and, comparing receptor model performances.

In this study we discuss the sources of uncertainty in the input data and the common approaches used to express them when preparing model input. We also analyze the contribution to the uncertainty deriving from critical model steps and the methodologies used to identify and reduce them. On the basis of the previous discussion we propose a sketch of Common QA/QC protocol.

Improving comparability and reliability of receptor models can be achieved by performing inter-comparison exercises. A community-wide inter-comparison organized and evaluated by the JRC within the framework of FAIRMODE is currently in progress. We describe the methodology for the evaluation and comparison of receptor models used in this exercise.

**Key words:** source apportionment, inter-comparison, receptor models, FAIRMODE

## INTRODUCTION

Receptor Models (RM) are used to identify the causes of pollution by analyzing concentrations and other parameters measured at one or more specific sites (receptor). In principle, RM are based on statistical analysis. At the first step they do not consider physical and chemical processes but evolved hybrid models can process additional information to constrain results. Most typical hybrid models are those using wind direction and speed or air masses trajectories to infer the geographical provenience of pollutants.

RM are independent from Emission Inventories and can work without complex meteorological and chemical processors

For that reason RM require low computational intensity if compared to Chemical Transport Models. Receptor Models rely on the mass conservation principle. However, a number of species derived from chemical transformations can be accounted for by the model using empirical equations.

Receptor models are most commonly used to apportion speciated Particulate Matter (PM) but it has been also used on VOCs mixtures, PAHs, inorganic gases and particulate size distribution. RM is appropriate for urban and regional scales and when combined with meteorological data are suitable to study medium to long range transport.

The mass balance equation that describes the basic rationale of RM is given by:

$$x_{ij} = \sum_{p=1}^P g_{ik} f_{kj} + e_{ij} \quad (1)$$

where  $x_{ij}$  is the concentration of the  $j^{\text{th}}$  species in the  $i^{\text{th}}$  sample,  $g_{ik}$  is the contribution of  $k^{\text{th}}$  source to  $i^{\text{th}}$  sample, and  $f_{kj}$  is the concentration of the  $j^{\text{th}}$  species in the  $k^{\text{th}}$  source

## RECEPTOR MODELS USED IN EUROPE

Receptor oriented source apportionment ranges from techniques based on elementary mathematical calculations and basic physical assumptions to complex models with sophisticated equation systems, pre- and post- processing routines and more or less user-friendly interfaces.

The potentials of RM are evidenced by the dramatic increase in the number of publications in scientific literature dealing with this topic in the last decade and the increasing ready-to-use available tools.

Figure 1 show the steady increase in the number of source apportionment studies using RM in Europe in the last decade. The highest increase rate occurred in 2005 and 2010 which coincides with the entry into force of the limit value for  $PM_{10}$  (Dir 1999/30/EC) and the target value for  $PM_{2.5}$  (Dir 2008/50/EC), respectively.

A survey on the use of receptor models for PM source apportionment in Europe between 2001 and 2010 was carried out. The survey examined 180 studies in 18 countries where 11 different types of model categories were recorded (Figure 1). The mass fractions considered span from  $PM_{10}$  to  $PM_1$ . About 65% of the studies are carried out in urban background sites, 18%

are source oriented sites and 17% are rural and remote sites. The most common model is the Positive Matrix Factorization (PMF), which is sometimes solved using ME platform (80 records). This kind of model became widespread especially after 2005 when the US-EPA made available user friendly on line versions. The success of this model is also linked to its use in the elaboration of AMS data, mostly oriented to the apportionment of the PM<sub>1</sub> organic particulate (9 studies). PCA and Absolute PCA (APCS) family of models (50 records) is also very popular and dominated during the first part of the last decade. CMB and the traditional factor analysis models (FA and APCFA), are used in 28 and 21 records, respectively. There are conceptual models based on a number of empirical assumptions that have been used in specific geographic areas like APEG (UK) and Lenschow (Germany) and others, such as UNMIX and COPREM that despite their good potentials have received little attention of European experts.

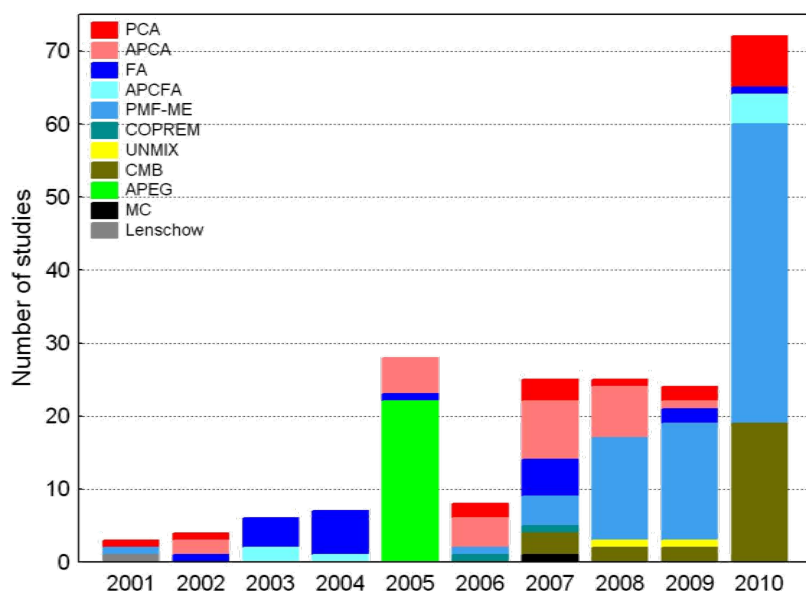


Figure 1. Time trend of PM source apportionment studies using RM in Europe between 2001 and 2010.

APEG receptor model was developed by the “Air Pollution Experts Group” in UK and only used for extensive (preliminary screening) source apportionment studies on several cities in that country. The Lenschow model was developed to estimate broad categories of sources combining different type of data from large cities and their surrounding area.

Spain is the country with the highest number of source apportionment studies with receptor models, followed by Italy and UK. Germany, France, Switzerland, Finland, Poland, Turkey and Ireland present between 6 and 8 studies each. Only 1 record is available for Norway and Czech Republic (Figure 3).

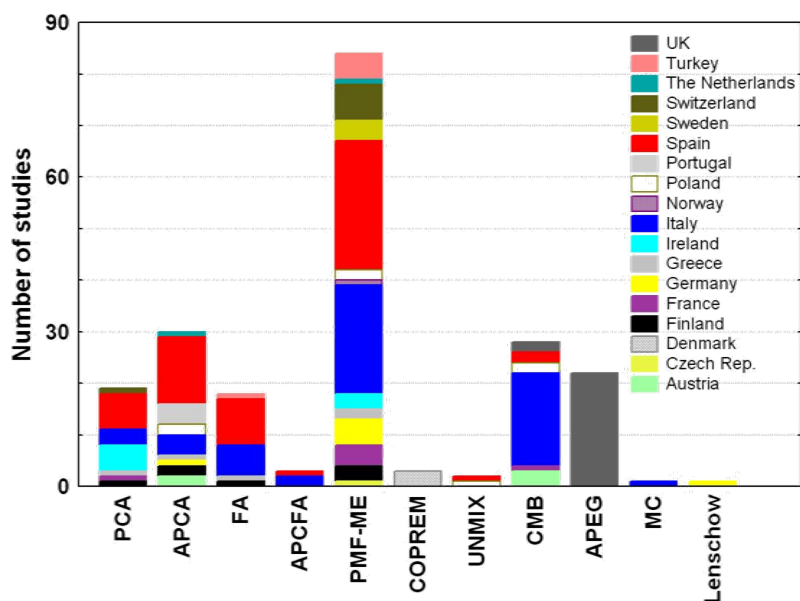


Figure 2. Number of European RM studies published between 2001 and 2010 grouped by model type.

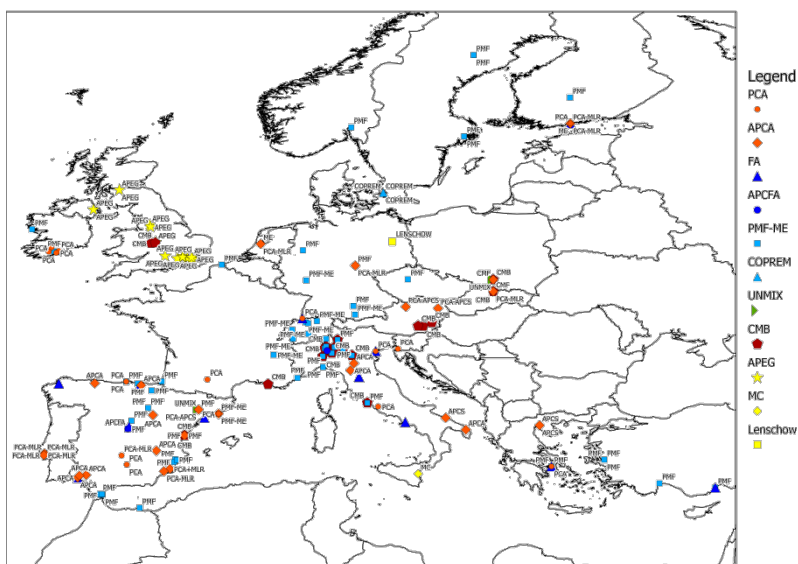


Figure 3. Geographic distribution of RM studies reporting model types between 2001 and 2010 in Europe

### SOURCES OF UNCERTAINTY

Interpreting the results of a source apportionment study and comparing results from different studies in different sites or in the same site with different models requires a proper uncertainty estimation and expression. In addition, there are models like PMF that weight data on the basis of their uncertainty, so that, an appropriate uncertainty estimation is a prerequisite for model execution. Comparison of results requires harmonization of the procedures according to international accepted standards like the Guide for the Expression of Uncertainty in Measurements (GUM). The first step to accomplish a proper estimation of the uncertainty is the identification of the main sources of uncertainty and the choice of the methodology to assess them. In RM the uncertainty derives from both inaccuracy in the input data and model assumptions and ambiguities.

#### Input Data uncertainty

Before starting model runs it is necessary to perform a preliminary data screening by accomplishing basic statistic tests and plots (e.g. average/median, standard deviation/interquartile distance, maximum and minimum values). Identification of outliers and anomalous values is of primary importance to keep uncertainty under control. This can be accomplished with Grubbs or other outlier tests. Criteria to express values below the minimum detection limit (MDL) need to be established at this point of the analysis. The most common choices are to set them to zero, to MDL/2 or to the MDL value itself. Certain models do not treat missing values; therefore, it is necessary to estimate them in order to avoid empty cells in the input data matrix. Clearly, missing values and values below MDL have an impact on the quality of results since the uncertainty of calculated values is higher than the one of measured ones.

Likely the most relevant component of input data uncertainty is the analytical uncertainty (Amato *et al.*, 2009). The estimation of this component is carried out by the analyst according to the analytical method specifications. However, it often happens to deal with databases in which single entry uncertainties are unavailable or inconsistent. In these cases uncertainties can be estimated using equation based approaches which rely on species MDL, empirical constants ( $k$ ), species concentration ( $C$ ) and/or coefficient of variation ( $CV$ ).

Analytical uncertainty can be estimated by linear regression where  $\sigma_a$  is the uncertainty of the analytical procedure and  $m$  is the mass of the analyte while  $\sigma_0$  and  $\alpha$  are fitting parameters (Anttila *et al.*, 1995)

$$\sigma_a^2 = \sigma_0^2 + (\alpha m)^2 \quad (2)$$

In the estimated fractional uncertainties (EFU) method (Kim and Hopke, 2005) the error structures ( $s_{ij}$ ) are calculated using the following equation:

$$s_{ij} = [MDL] / 3 + kx_{ij} \quad (3)$$

No empirical constants but  $MDL$ , and coefficient of variation ( $CV$ ) are used in the analytical uncertainty equations used in a recent work (Chow *et al.*, 2007).

$$\sigma_{i,t}^a = \sqrt{MDL_i^2 + (CV_i \times C_{i,t})^2} \quad (4)$$

Also sampling contributes to the uncertainty of measured values due to sampling volume uncertainty, selective effect and other artifacts caused by the inlet, losses due to sample transport and conservation. These contributions can be assessed by

field tests and comparison with reference instrumentation. Other authors (Amato *et al.*, 2009) incorporated the sampling uncertainty into the expanded uncertainty by considering the sampled volume ( $V_i$ ) and a coefficient ( $\beta$ ) to account for additional uncertainty sources.

$$\sigma_{ij}^2 = \sqrt{\frac{\sigma_A^2}{V_i^2} + (\beta x_{ij})^2} \quad (5)$$

By definition receptor models reconstruct the source contribution at a specific point. However, the ultimate interest of these studies is to assess the influence of pollution on the population, in a more or less homogeneous area, where the monitoring site is located. Therefore, the generalization of results depends on the representativeness of the monitoring site. To improve the geographic representativeness of the study (also in terms of source distribution) it is possible to combine more sites in a single study. Alternatively, the representativeness of the monitoring site can be evaluated using geostatistical techniques. Similarly, the time representativeness of the model outputs depends on the amount and distribution of collected samples. Studies oriented to observe average levels use samples collected throughout the year while studies interested in specific kind of events may concentrate the data collection in specific seasons. Data time resolution, and consequently uncertainty, vary according to the kind of processes under study and the available resources. In order to obtain measurements representative of the whole year with a limited amount of samples to analyze it is common to collect data every three to six days. On the other hand, modern analytical techniques made available high time resolution data. This kind of information makes it possible to observe short lived patterns but the effort to process such a big amount of data imposes a limitation to the length of monitoring campaigns and consequently the representativeness of the time series.

The output of RM like CMB which use source profiles as input data is seriously affected by their uncertainty. To prevent problems of collinearity sources with similar chemical composition need to be combined into source categories. The selection of sources to include in the final input has to be representative of the local sources in the study area and of their variability in time. The suitability of source profiles could be checked using techniques like edge identification and ratio-ratio plots (Robinson *et al.*, 2006). In general, the use of local source profiles improves the model performance indicators.

#### Uncertainty associated with model performance

Considering that classical receptor models rely on the principle of mass conservation between source and receptor, substantial departures from this assumption due to evaporation, condensation or degradation of species constitutes a source of uncertainty. In general, RMs are able to identify secondary inorganic aerosol assuming that ammonium sulphate and nitrate mainly derive from gaseous precursors ( $\text{SO}_2$ ,  $\text{NO}_x$  and  $\text{NH}_3$ ). Moreover, the mass conservation assumption can be relaxed downweighting the volatile or reactive species. When quantitative information about the processes that precursor species undergo after emission is available, it is possible to introduce compensation equations using empirical coefficients to correct the expected amounts of products at the receptor.

In factor analysis, the number of relevant factors and their correspondence with sources is unknown and represents another source of uncertainty. Estimating the number of factors is often performed with an iterative procedure by checking the influence of the number of factors on the model performance. A number of indicators is used to guide the selection of the number of factors e.g. signal to noise, residuals,  $Q$  value (in PMF). Establishing a correspondence between the factors resulting from the analysis and the sources in the area is not straightforward and may substantially contribute to the uncertainty of the output. This step has to be accomplished by experienced users with the knowledge of the local sources and the study area geographic and meteorological patterns. The identified sources have to be sound from a physical-chemical point of view. Marker compounds are useful for specific types of sources. Comparison of the factor profiles with local source profiles is essential and may be performed using simple mathematical algorithms to identify the best fit.

Another contribution to the overall uncertainty in factor analysis is the lack of a unique solution due to the large number of unknown variables. This limitation of factor analysis, called rotational uncertainty (Paatero and Hopke, 2009) is partially removed by non-negativity constraints in PMF. In addition, an analysis of the rotations in the area close to the final solution can be performed by analyzing the  $Q/Q_{\text{zero}}$  parameter. Tests on residuals, FPEAK and analysis of edges also help to identify the best rotation. Moreover, introducing additional information about the sources (e.g. set the source to zero when it is known that the source is absent) contributes to eliminate the rotational uncertainty.

The most evolved tools (e.g. PMF2 -3, CMB 8.2, UNMIX, etc.) include routine tests to evaluate the overall model performance and provide an estimation of the uncertainty associated to the Source Contribution Estimations (SCE).

Since it is not possible to know what is the actual contribution of the sources in real-world conditions the overall model uncertainty can be evaluated by performing an intercomparison. Models are complementary and mutually reinforce. Source apportionment has shown to be more robust when different approaches are used and the consistency of results is tested (e.g. Rizzo and Scheff, 2007).

#### INTERCOMPARISON

In order to promote the harmonization and continuous improvement of source apportionment with receptor models in Europe an intercomparison has been launched. The exercise consists of comparing the results of source apportionment analyses performed by independent groups using different techniques on the same dataset. The intercomparison is expected to a) provide information about the degree of reproducibility between different approaches and scientific backgrounds and b) assess the uncertainty of the SCE.

### Intercomparison Assessment

In source apportionment studies it is not possible to validate the model outputs against measured values since the actual contributions from the sources are unknown. Different approaches have been used to compare the performance of different models on the same dataset: visual comparison of models' SCE mean and SD for each source (Favez *et al.*, 2010; Hopke *et al.*, 2006; Larsen *et al.*, 2008), correlation coefficient (Favez *et al.*, 2010; Hopke *et al.*, 2006; Sandradewi *et al.*, 2008) and regression analysis between SCE provided by different models (Rizzo and Scheff, 2007; Sandradewi *et al.*, 2008). In one case ANOVA was carried out on the mean source contributions (Hopke *et al.*, 2006).

The present intercomparison is evaluated according to international standards for proficiency testing exercises (ISO 13528). The assigned value, that is the reference against participants results are compared, is generated by the robust analysis iterative algorithm. The criterion to assess data comparability, standard deviation for proficiency assessment ( $\hat{\sigma}$ ), is derived from the Directive 2008/50/EC data quality objectives for PM modelling (equivalent to 50% of the annual mean for PM<sub>10</sub> and PM<sub>2.5</sub> total mass) or in alternative from the robust analysis. The SCE provided by participants is compared with the assigned values using the bias and the z-score algorithms.

### SKETCH OF COMMON PROTOCOL

An extensive and accurate examination of the scientific literature in the field of source apportionment has shown that receptor models are widely used in Europe in support of the implementation of the Air Quality Directive (Douros *et al.*, 2011). However, there is a need of harmonization of the tools, procedures and criteria applied in the different countries. In order to promote and guide the discussion of the experts an outline containing items considered relevant for a Common Protocol has been drafted:

- criteria for site selection, and minimum number of samples
- guidelines for the estimation of uncertainty in input data
- input data pre-treatment
- determination of the number of factors and correspondence between factors and sources
- interpretation of model specific performance indicators and uncertainty
- sensitivity analysis (for model validation)
- comparison and integration of the outputs of different models

### REFERENCES

- Amato, F., M. Pandolfi, A. Escrig, X. Querol, A. Alastuey, J. Pey, N. Perez, and P. K. Hopke (2009), Quantifying road dust resuspension in urban environment by Multilinear Engine: A comparison with PMF2, *Atmospheric Environment*, 43(17), 2770-2780.
- Anttila, P., P. Paatero, U. Tapper, and O. Järvinen (1995), Source identification of bulk wet deposition in Finland by positive matrix factorization, *Atmospheric Environment*, 29, 1705-1718.
- Chow, J. C., J. G. Watson Jr, L. W. Antony-Chen, M. C. Oliver Chang, N. F. Robinson, D. Trimble, and S. Kohl (2007), The IMPROVE\_A temperature protocol for Thermal/Optical carbon analysis: maintaining consistency with a long-term database, *Journal of Air Waste and Management Association*, 57, 1014-1023.
- Douros, J., E. Fragkou, N. Moussiopoulos, and C. A. Belis (2011), The use and evaluation of multi-pollutant source apportionment methodologies by EU authorities and research groups, *this volume*.
- Favez, O., et al. (2010), Inter-comparison of source apportionment models for the estimation of wood burning aerosols during wintertime in an Alpine city (Grenoble, France), *Atmospheric Chemistry and Physics*, 10(12), 5295-5314.
- Hopke, P. K., et al. (2006), PM source apportionment and health effects: 1. Intercomparison of source apportionment results, *Journal of Exposure Science and Environmental Epidemiology*, 16(3), 275-286.
- Kim, Y. J., and P. K. Hopke (2005), Estimation of Organic Carbon Blank Values and Error Structures of the Speciation Trends Network Data for Source Apportionment, *J. Air & Waste Manage. Assoc.*, 55, 1190-1199.
- Larsen, B. R., H. Junninen, J. Monster, M. Viana, P. Tsakovski, R. M. Duvall, G. Norris, and X. Querol (2008), The Krakow receptor modelling intercomparison exercise, *JRC Scientific and Technical Reports, EUR 23621 EN 2008*.
- Paatero, P., and P. K. Hopke (2009), Rotational tools for factors analytic models, *Journal of Chemometrics*, 23, 91-100.
- Rizzo, M. J., and P. A. Scheff (2007), Utilizing the Chemical Mass Balance and Positive Matrix Factorization models to determine influential species and examine possible rotations in receptor modeling results, *Atmospheric Environment - Part A General Topics*, 41, 6986-6998.
- Robinson, A. L., R. Subramanian, N. M. Donahue, A. Bernardo-Bricker, and W. F. Rogge (2006), Source apportionment of molecular markers and organic aerosol-1. Polycyclic Aromatic Hydrocarbons and methodology for data visualization, *Environmental Science and Technology*, 40, 7803-7810.
- Sandradewi, J., et al. (2008), Comparison of several wood smoke markers and source apportionment methods for wood burning particulate mass, *Atmospheric Chemistry and Physics Discussions*, 8, 8091-8118.