

ESAT: Environmental Source Apportionment Toolkit Python package

Deron Smith¹, Michael Cyterski¹, John M Johnston¹, Kurt Wolfe¹,
and Rajbir Parmar¹

¹ United States Environmental Protection Agency, Office of Research and Development, Center for Environmental Measurement and Modeling

DOI: [10.21105/joss.07316](https://doi.org/10.21105/joss.07316)

Software

- [Review](#)
- [Repository](#)
- [Archive](#)

Editor: [Mengqi Zhao](#)

Reviewers:

- [@gutabeshu](#)
- [@ifoxfoot](#)
- [@niravlekinwala](#)

Submitted: 09 September 2024

Published: 03 December 2024

License

Authors of papers retain copyright
and release the work under a
Creative Commons Attribution 4.0
International License ([CC BY 4.0](#)).

Summary

Source apportionment is an important tool in environmental science where sample or sensor data are often the product of many, often unknown, contributing sources. Source apportionment is used to understand the relative contributions of air sources (Bhandari et al., 2022) like vehicle emissions, industrial activities, or dust; as well as particulate matter pollution and to identify relative contributions of point sources and non-point sources in water bodies such as lakes, rivers, and estuaries (Jiang et al., 2019; Mamun & An, 2021). Using non-negative matrix factorization (NMF), source apportionment models estimate potential source profiles and contributions providing a cost-efficient method for further strategic data collection or modeling.

Environmental Source Apportionment Toolkit (ESAT) is an open-source Python package that provides a flexible and transparent workflow for source apportionment using NMF algorithms, developed to replace the EPA's Positive Matrix Factorization version 5 (PMF5) application (EPA, 2014; Pentti Paatero, 1999). ESAT recreates the source apportionment workflow of PMF5 including pre- and post-processing analytical tools, batch modeling, uncertainty estimations and customized constraints. ESAT offers a simulator for generating datasets from synthetic profiles and contributions, allowing for model output evaluation. The synthetic profiles can be randomly generated, use a pre-defined set of profiles, or be a combination of the two. The random synthetic contributions can follow specified curves and value ranges. By running ESAT using the synthetic datasets, users are able to accurately assess ESAT's ability to find a solution that recreates the original synthetic profiles and contributions.

Statement of Need

The EPA's PMF5, released in 2014, provides a widely-used source apportionment modeling and analysis workflow that is no longer supported and relies on the proprietary Multilinear Engine v2 (ME2). ESAT has been developed as a replacement to PMF5, and has been designed for increased flexibility, documentation and transparency.

The Python API and CLI of ESAT provides a programmatic interface that can recreate the PMF5 workflow. The matrix factorization algorithms in ESAT have been written in Rust for runtime optimization. The ESAT API and CLI provides a flexible way to create source apportionment workflows and novel research applications. ESAT was developed for environmental research, though it's not limited to that domain, as matrix factorization is used in many different fields.

Algorithms

Source apportionment algorithms use a loss function to quantify the difference between the input data matrix (V) and the product of a factor contribution matrix (W) and a factor profile matrix (H), weighted by an uncertainty matrix (U) (Pentti Paatero & Tapper, 1994). The goal is to find factor matrices that best reproduce the input matrix, while constraining all, or most of, the factor elements to be non-negative. The solution, W and H , can be used to calculate the residuals and overall model loss. ESAT has two NMF algorithms for updating the profile and contribution matrices: least-squares NMF (LS-NMF) (Wang et al., 2006) and weighted-semi NMF (WS-NMF) (Ding et al., 2008; Melo & Wainer, 2012).

The loss function used in ESAT, and PMF5, is a variation of squared-error loss, where data uncertainty is taken into consideration (both in the loss function and in the matrix update equations):

$$Q = \sum_{i=1}^n \sum_{j=1}^m \left[\frac{V_{ij} - \sum_{k=1}^K W_{ik} H_{kj}}{U_{ij}} \right]^2$$

here V is the input data matrix of features (columns= M) by samples (rows= N), U is the uncertainty matrix of the input data matrix, W is the factor contribution matrix of samples by factors= K , H is the factor profile of factors by features.

The ESAT versions of NMF algorithms convert the uncertainty U into weights defined as $Uw = \frac{1}{U^2}$. The update equations for LS-NMF then become:

$$H_{t+1} = H_t \circ \frac{W_t(V \circ Uw)}{W_t((W_t H_t) \circ Uw)}$$

$$W_{t+1} = W_t \circ \frac{(V \circ Uw) H_{t+1}}{((W_t H_{t+1}) \circ Uw) H_{t+1}}$$

while the update equations for WS-NMF:

$$W_{t+1,i} = (H^T U w_i^d H)^{-1} (H^T U w_i^d V_i)$$

$$H_{t+1,i} = H_{t,i} \sqrt{\frac{((V^T U w) W_{t+1})_i^+ + [H_t (W_{t+1}^T U w W)^-]_i}{((V^T U w) W_{t+1})_i^- + [H_t (W_{t+1}^T U w W)^+]_i}}$$

where $W^- = \frac{(|W| - W)}{2.0}$ and $W^+ = \frac{(|W| + W)}{2.0}$.

Error Estimation

An important part of the source apportionment workflow is quantifying potential model error. ESAT offers the error estimation methods that were developed and made available in PMF5 (Brown et al., 2015; P. Paatero et al., 2014).

The displacement method (DISP) determines the amount that a source profile feature, a single value in the H matrix, must increase and decrease to cause targeted changes to the loss value. One or more features can be selected in the DISP uncertainty analysis. The bootstrap method (BS) uses block bootstrap resampling with replacement to create datasets with the original dimensions of the input, where the order of the samples has been modified in blocks of a specified size. The BS method then calculates a new model from the bootstrap dataset, and original initialization, to evaluate how the profiles and concentrations change as a result of

sample reordering. The bootstrap-displacement method (BS-DISP) is the combination of the two techniques, where DISP is run for each bootstrap model on one or more features.

These error estimation methods address different uncertainty aspects: DISP targets rotational uncertainty, BS addresses random errors and sample variability, and BS-DISP provides the most comprehensive understanding of how the uncertainty impacts a source apportionment solution.

Acknowledgements

This paper has been reviewed in accordance with EPA policy and approved for publication. ESAT development has been funded by U.S. EPA. Mention of any trade names, products, or services does not convey, and should not be interpreted as conveying, official EPA approval, endorsement, or recommendation. The views expressed in this paper are those of the authors and do not necessarily represent the views or policies of the U.S. EPA.

References

- Bhandari, S., Arub, Z., Habib, G., Apte, J. S., & Hildebrandt Ruiz, L. (2022). Source apportionment resolved by time of day for improved deconvolution of primary source contributions to air pollution. *Atmospheric Measurement Techniques*, 15(20), 6051–6074. <https://doi.org/10.5194/amt-15-6051-2022>
- Brown, S. G., Eberly, S., Paatero, P., & Norris, G. A. (2015). Methods for estimating uncertainty in PMF solutions: Examples with ambient air and water quality data and guidance on reporting PMF results. *Science of the Total Environment*, 518, 626–635. <https://doi.org/10.1016/j.scitotenv.2015.01.022>
- Ding, C. H., Li, T., & Jordan, M. I. (2008). Convex and semi-nonnegative matrix factorizations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(1), 45–55. <https://doi.org/10.1109/TPAMI.2008.277>
- EPA, U. S. (2014). *Positive matrix factorization model for environmental data analyses*. <https://www.epa.gov/air-research/positive-matrix-factorization-model-environmental-data-analyses>
- Jiang, J., Khan, A. U., & Shi, B. (2019). Application of positive matrix factorization to identify potential sources of water quality deterioration of huaihe river, china. *Applied Water Science*, 9(63, 3). <https://doi.org/10.1007/s13201-019-0938-4>
- Mamun, M., & An, K.-G. (2021). Application of multivariate statistical techniques and water quality index for the assessment of water quality and apportionment of pollution sources in the yeongsan river, south korea. *International Journal of Environmental Research and Public Health*, 18(16). <https://doi.org/10.3390/ijerph18168268>
- Melo, E. V. de, & Wainer, J. (2012). *Semi-NMF and weighted semi-NMF algorithms comparison*.
- Paatero, Pentti. (1999). The multilinear engine—a table-driven, least squares program for solving multilinear problems, including the n-way parallel factor analysis model. *Journal of Computational and Graphical Statistics*, 8(4), 854–888. <https://doi.org/10.1080/10618600.1999.10474853>
- Paatero, P., Eberly, S., Brown, S. G., & Norris, G. A. (2014). Methods for estimating uncertainty in factor analytic solutions. *Atmospheric Measurement Techniques*, 7(3), 781–797. <https://doi.org/10.5194/amt-7-781-2014>
- Paatero, Pentti, & Tapper, U. (1994). Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics*, 5(2),

111–126. <https://doi.org/10.1002/env.3170050203>

Wang, G., Kossenkov, A. V., & Ochs, M. F. (2006). LS-NMF: A modified non-negative matrix factorization algorithm utilizing uncertainty estimates. *BMC Bioinformatics*, 7, 1–10. <https://doi.org/10.1186/1471-2105-7-175>