







GAMMA_FLOW: Guided Analysis of Multi-label spectra by Matrix factorization for Lightweight Operational Workflows

Viola Rädle¹, Tilman Hartwig¹, Benjamin Oesen¹, Julius Vogt², Eike Gericke², Emily Alice Kröger², and Martin Baron²

¹ Application Lab for AI and Big Data, German Environmental Agency, Leipzig, Germany^{ROR} ² Federal Office for Radiation Protection, Berlin, Germany^{ROR} ¶ Corresponding author

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

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

Editor: [Open Journals](#) 

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

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

Summary

Most radioactive sources can be identified by measuring their emitted radiation (X-rays and gamma rays), and visualizing them as a spectrum. In nuclear security applications, the resulting gamma spectra have to be analyzed in real-time as immediate reaction and decision making may be required. However, the manual recognition of isotopes present in a spectrum constitutes a strenuous, error-prone task that depends upon expert knowledge. Hence, this raises the need for algorithms assisting in the initial categorization and recognizability of measured gamma spectra.

The delineated use case brings along several requirements:

- As mobile, room temperature detectors are often deployed in nuclear security applications, the produced spectra typically exhibit a rather low energy resolution. In addition, a high temporal resolution is required (usually around one spectrum per second), leading to a low acquisition time and a low signal-to-noise ratio. Hence, the model must be robust and be able to handle noisy data.
- For some radioactive sources, acquisition of training spectra may be challenging. Instead, spectra of those isotopes are simulated using Monte Carlo N-Particle (MCNP) code ([Kulesza et al., 2022](#)). In this process, energy deposition in a detector material is simulated, yielding spectra that can be used for model training. However, simulated spectra and measured spectra from real-world sources may differ, which may be a constraint for model performance. On this account, preliminary data exploration is crucial to assess the similarity of spectral data from different detectors and to evaluate potential data limitations.
- Lastly, not only the correct classification of single-label test spectra (stemming from one isotope) is necessary, but also the decomposition of linear combinations of various isotopes (multi-label spectra). Hence, classification approaches like k-nearest-neighbours that solely depend on the similarity between training and test spectra are not applicable.

This paper presents `gamma_flow`, a python package that includes the

- classification of test spectra to predict their constituents
- denoising of test spectra for better recognizability
- outlier detection to evaluate the model's applicability to test spectra

It is based on a dimensionality reduction model that constitutes a novel, supervised approach to non-negative matrix factorization (NMF). More explicitly, the spectral data matrix is decomposed into the product of two low-rank matrices denoted as the scores (spectral data in latent space) and the loadings (transformation matrix or latent components). The loadings matrix is predefined and consists of the mean spectra of the training isotopes. Hence, by design, the scores axes correspond to the share of an isotope in a spectrum, resulting in an interpretable latent space.

As a result, the classification of a test spectrum can be read directly from its (normalized) scores. In particular, shares of individual isotopes in a multi-label spectrum can be identified. This leads to an explainable quantitative prediction of the spectral constituents.

The scores can be transformed back into spectral space by applying the inverse model. This inverse transformation rids the test spectrum of noise and results in a smooth, easily recognizable denoised spectrum.

If a test spectrum of an isotope is unknown to the model (i.e. this isotope was not included in model training), it can still be projected into latent space. However, when the latent space information (scores) are decompressed, the resulting denoised spectrum does not resemble the original spectrum any more. Some original features may not be captured while new peaks may have been fabricated. This can be quantified by calculating the cosine similarity between the original and the denoised spectrum, which can serve as an indicator of a test spectrum to be an outlier.

Statement of need

In many research fields, spectral measurements help to assess material properties. In this context, an area of interest for many researchers is the classification (automated labelling) of the measured spectra. Proprietary spectral analysis software, however, are often limited in their functionality and adaptability (Lam, 2011; Nasereddin & Shakib, 2023). In addition, the underlying mechanisms are usually not revealed and may act as a black-box system to the user (El Amri et al., 2022). On top of that, a spectral comparison is typically only possible for spectra of pure substances (Cowger et al., 2021). However, there may be a need to decompose multi-label spectra (linear combinations of different substances) and identify their constituents.

`gamma_flow` is a Python package that can assist researchers in the classification, denoising and outlier detection of spectra. It includes data preprocessing, data exploration, model training and testing as well as an exploratory section on outlier detection. Making use of matrix decomposition methods, the designed model is lean and performant. Training and inference do not require special hardware or extensive computational power. This allows real-time application on ordinary laboratory computers and easy implementation into the measurement routine.

The provided example dataset contains gamma spectra of several measured and simulated isotopes as well as pure background spectra. While this package was developed in need of an analysis tool for gamma spectra, it is suitable for any one-dimensional spectra.

Exemplary applications encompass

- **Infrared spectroscopy** for the assessment of the polymer composition of microplastics in water (Ferreiro et al., 2023; Whiting et al., 2022)
- **mass spectrometry** for protein identification in snake venom (Yasemin et al., 2021; Zelanis et al., 2019)
- **Raman spectroscopy** for analysis of complex pharmaceutical mixtures and detection of dilution products like lactose (Fu et al., 2021)

- **UV-Vis spectroscopy** for detection of pesticides in surface waters (Guo et al., 2020; Qi et al., 2024)
- **stellar spectroscopy** to infer the chemical composition of stars (Gray, 2021)

Methodology and structure

This python package consists of three jupyter notebooks that are executed consecutively. In this section, their functionality and is outlined, with an emphasis on the mathematical structure of the model.

1. Preprocessing and data exploration

The notebook 01_preprocessing.ipynb synchronizes spectral data and provides a framework of visualizations for data exploration. All functions called in this notebook are found in tools_preprocessing.py.

During **preprocessing**, the following steps are performed:

- Spectral data files are converted from .xslm/.spe data to .npy format and saved.
- Spectra of different energy calibrations are rebinned to a standard energy calibration.
- Spectral data are aggregated by label classes and detectors. Thus, it is possible to collect data from different files and formats.
- Optional: The spectra per isotope are limited to a maximum number.
- The preprocessed spectra are saved as .npy files.

Data exploration involves the following visualizations:

- For each label class (e.g. for each isotope), the mean spectra are calculated detector-wise and compared quantitatively by the cosine similarity.
- For each label class, example spectra are chosen randomly and plotted to provide an overview over the data.
- The cosine similarity is calculated and visualized as a matrix for all label classes and detectors. This helps to assess whether the model can handle spectra from different detectors.

2. Model training and testing

The notebook 02_model.ipynb trains and tests a dimensionality reduction model that allows for denoising, classification and outlier detection of test spectra. All functions called in this notebook are found in tools_model.py.

The dimensionality reduction model presented in this paper comprises a matrix decomposition of spectral data. More precisely, the original spectra matrix X is reconstructed by two low-rank matrices S and L :

$$X \approx SL^T$$

with S : scores matrix (spectra in latent space)

L : loadings matrix (transformation matrix or latent components)

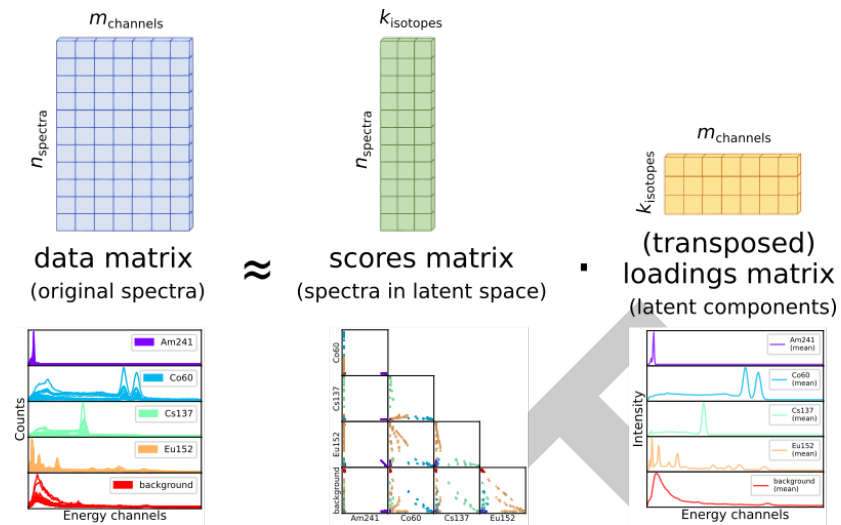


Figure 1: Matrix decomposition of spectral data.

As illustrated in Figure 1, original spectral data can be compressed into k_{isotopes} dimensions. To ensure a conclusive assignment of the latent space axes to the isotopes (i.e. one axis stands for of one isotope), the loadings matrix is predefined as the mean spectra of the k_{isotopes} isotopes.

During model training, mean spectra for all isotopes are calculated. The scores are then derived by non-negative least squares fit of the original spectra to the loadings matrix. Thus, the components of the normalized scores vectors directly reveal the contributions of the individual isotopes. Denoised spectra, on the other hand, are computed by transforming the non-normalized scores back into spectral space (i.e. by multiplication of with the loadings matrix).

In mathematical terms, this model represents a 'supervised' approach to Non-negative Matrix Factorization (NMF) (Bilton et al., 2019; Shreeves, 2020). While dimensionality reduction is conventionally an unsupervised task as it only considers data structure (Olaya & Otman, 2022), our approach integrates labels in model training. This leads to an interpretable latent space and obviates the need for an additional classification step. While other supervised NMF approaches incorporate classification loss in model training (Bisot et al., 2016; Lee et al., 2010; Leuschner et al., 2019), our model focuses on a comprehensible construction of the latent space.

The model is trained using spectral data from the specified detectors `dets_tr` and isotopes `isotopes_tr`. Subsequently, it is inferred (i.e. scores are calculated) on three different test datasets: 1. validation data/holdout data from same detector as used in training (each spectrum including only one isotope or pure background) 2. test data from different detector (each spectrum including one isotope and background) 3. multi-label test data from different detector (each spectrum including multiple isotopes and background)

For all test datasets, spectra are classified and denoised. The results are visualized as - confusion matrix

- misclassified spectra
- denoised example spectrum
- misclassification statistics
- scores as scatter matrix
- mean scores as bar plot

This helps to assess model performance with respect to classification and denoising.

3. Outlier analysis

The notebook `03_outlier.ipynb` provides an exploratory approach to outliers detection, i.e. to identify spectra from isotopes that were not used in model training. All functions called in this notebook are found in `tools_outlier.py`.

To simulate outlier spectra, a mock dataset is generated by training a model after removing one specific isotope. The trained model is then inferenced on spectra of this unknown isotope to investigate its behaviour with outliers. First, the resulting latent space distribution and further meta data are analyzed to distinguish known from unknown spectra. Using a decision tree, the most informative feature is identified. Next, a decision boundary is derived for this feature, by

- a) using the condition of the first split in the decision tree
- b) fitting a logistic regression (sigmoid function) to the data
- c) setting a manual threshold by considering accuracy, precision and recall of outlier identification.

The derived decision boundary can then be implemented in the measurement pipeline by the user.

Apart from the jupyter notebooks and python files described above, the project includes the following python files:

- `globals.py`: global variables
- `plotting.py`: all visualizations and plotting routines
- `util.py`: basic functions that are used by all notebooks

Acknowledgements

We gratefully acknowledge the support provided by the Federal Ministry for the Environment, Nature Conservation and Nuclear Safety (BMUV), whose funding has been instrumental in enabling us to achieve our research objectives and explore new directions. We also extend our appreciation to Martin Bussick in his function as the AI coordinator. Additionally, we thank the entire AI-Lab team for their support and inspiration, with special recognition to Ruth Brodte for guidance on legal and licensing matters.

References

- Bilton, K. J., Joshi, T. H., Bandstra, M. S., Curtis, J. C., Quiter, B. J., Cooper, R. J., & Vetter, K. (2019). Non-negative Matrix Factorization of Gamma-Ray Spectra for Background Modeling, Detection, and Source Identification. *IEEE Transactions on Nuclear Science*, 66(5), 827–837. <https://doi.org/10.1109/TNS.2019.2907267>
- Bisot, V., Serizel, R., Essid, S., & Richard, G. (2016, September). Supervised nonnegative matrix factorization for acoustic scene classification. *IEEE International Evaluation Campaign on Detection and Classification of Acoustic Scenes and Events (DCASE 2016)*.
- Cowger, W., Steinmetz, Z., Gray, A., Munno, K., Lynch, J., Hapich, H., Primpke, S., De Frond, H., Rochman, C., & Herodotou, O. (2021). Microplastic Spectral Classification Needs an Open Source Community: Open Specy to the Rescue! *Analytical Chemistry*, 93(21), 7543–7548. <https://doi.org/10.1021/acs.analchem.1c00123>
- El Amri, L., Chetaine, A., Amsil, H., El Mokhtari, B., Bounouira, H., Didi, A., Benchrif, A., Laraki, K., & Marah, H. (2022). New open-source software for gamma-ray spectra analysis. *Applied Radiation and Isotopes*, 185, 110227. <https://doi.org/10.1016/j.apradiso.2022.110227>
- Ferreiro, B., Leardi, R., Farinini, E., & Andrade, J. M. (2023). Supervised classification combined with genetic algorithm variable selection for a fast identification of polymeric

- microdebris using infrared reflectance. *Marine Pollution Bulletin*, 195, 115540. <https://doi.org/10.1016/j.marpolbul.2023.115540>
- Fu, X., Zhong, L., Cao, Y., Chen, H., & Lu, F. (2021). Quantitative analysis of excipient dominated drug formulations by Raman spectroscopy combined with deep learning. *Analytical Methods*, 13(1), 64–68. <https://doi.org/10.1039/D0AY01874K>
- Gray, D. F. (2021). The Observation and Analysis of Stellar Photospheres. In *Higher Education from Cambridge University Press*. <https://www.cambridge.org/highereducation/books/the-observation-and-analysis-of-stellar-photospheres/67B340445C56F4421BCBA0AF-FAAFDEE0>; Cambridge University Press. <https://doi.org/10.1017/9781009082136>
- Guo, Y., Liu, C., Ye, R., & Duan, Q. (2020). Advances on Water Quality Detection by UV-Vis Spectroscopy. *Applied Sciences*, 10(19), 6874. <https://doi.org/10.3390/app10196874>
- Kulesza, J. A., Adams, T. R., Armstrong, J. C., Bolding, S. R., Brown, F. B., Bull, J. S., Burke, T. P., Clark, A. R., Forster III, R. A. (Art), Giron, J. F., Grieve, T. S., Josey, C. J., Martz, R. L., McKinney, G. W., Pearson, E. J., Rising, M. E., Solomon Jr., C. J. (CJ), Swaminarayan, S., Trahan, T. J., ... Zukaitis, A. J. (2022). *MCNP® Code Version 6.3.0 Theory & User Manual* (LA-UR-22-30006). Los Alamos National Laboratory (LANL), Los Alamos, NM (United States). <https://doi.org/10.2172/1889957>
- Lam, H. (2011). Building and Searching Tandem Mass Spectral Libraries for Peptide Identification. *Molecular & Cellular Proteomics*, 10(12). <https://doi.org/10.1074/mcp.R111.008565>
- Lee, H., Yoo, J., & Choi, S. (2010). Semi-Supervised Nonnegative Matrix Factorization. *IEEE Signal Processing Letters*, 17(1), 4–7. <https://doi.org/10.1109/LSP.2009.2027163>
- Leuschner, J., Schmidt, M., Fernsel, P., Lachmund, D., Boskamp, T., & Maass, P. (2019). Supervised non-negative matrix factorization methods for MALDI imaging applications. *Bioinformatics*, 35(11), 1940–1947. <https://doi.org/10.1093/bioinformatics/bty909>
- Nasereddin, J., & Shakib, M. (2023). Ira: A free and open-source Fourier transform infrared (FTIR) data analysis widget for pharmaceutical applications. *Analytical Letters*, 56(16), 2637–2648. <https://doi.org/10.1080/00032719.2023.2180516>
- Olaya, J., & Otman, C. (2022). Non-negative Matrix Factorization for Dimensionality Reduction. *ITM Web of Conferences*, 48, 03006. <https://doi.org/10.1051/itmconf/20224803006>
- Qi, X., Lian, Y., Xie, L., Wang, Y., & Lu, Z. (2024). Water quality detection based on UV-Vis and NIR spectroscopy: A review. *Applied Spectroscopy Reviews*, 59(8), 1036–1060. <https://doi.org/10.1080/05704928.2023.2294458>
- Shreeves, P. (2020). *Dimensionality reduction techniques with applications in Raman spectroscopy - UBC Library Open Collections* [PhD thesis]. University of British Columbia.
- Whiting, Q. T., O'Connor, K. F., Potter, P. M., & Al-Abed, S. R. (2022). A high-throughput, automated technique for microplastics detection, quantification, and characterization in surface waters using laser direct infrared spectroscopy. *Analytical and Bioanalytical Chemistry*, 414(29), 8353–8364. <https://doi.org/10.1007/s00216-022-04371-2>
- Yasemin, N. Ç., Izzet, A., Ali, K. M., Selçuk, K., & Bekir, S. (2021). Identification of Snake Venoms According to their Protein Content Using the MALDI-TOF-MS Method. *Analytical Chemistry Letters*, 11(2), 153–167. <https://doi.org/10.1080/22297928.2021.1894974>
- Zelanis, A., Silva, D. A., Kitano, E. S., Liberato, T., Fukushima, I., Serrano, S. M. T., & Tashima, A. K. (2019). A first step towards building spectral libraries as complementary tools for snake venom proteome/peptidome studies. *Comparative Biochemistry and Physiology Part D: Genomics and Proteomics*, 31, 100599. <https://doi.org/10.1016/j.cbd.2019.100599>