

# GAMMA\_FLOW: Guided Analysis of Multi-label spectra by Matrix Factorization for Lightweight Operational Workflows

Viola Rädle<sup>a,\*</sup>, Tilman Hartwig<sup>a</sup>, Benjamin Oesen<sup>a</sup>, Emily Alice Kröger<sup>b</sup>, Julius Vogt<sup>b</sup>, Eike Gericke<sup>b</sup>, Martin Baron<sup>b</sup>

<sup>a</sup>*Application Lab for AI and Big Data, German Environmental Agency, Leipzig,  
Germany*

<sup>b</sup>*Federal Office for Radiation Protection, Berlin, Germany*

---

## Abstract

GAMMA\_FLOW is an open-source Python package for real-time analysis of spectral data. It supports classification, denoising, decomposition, and outlier detection of both single- and multi-component spectra. Instead of relying on large, computationally intensive models, it employs a supervised approach to non-negative matrix factorization (NMF) for dimensionality reduction. This ensures a fast, efficient, and adaptable analysis while reducing computational costs. GAMMA\_FLOW achieves classification accuracies above 90% and enables reliable automated spectral interpretation. Originally developed for gamma-ray spectra, it is applicable to any type of one-dimensional spectral data. As an open and flexible alternative to proprietary software, it supports various applications in research and industry.

*Keywords:* Python, Gamma spectroscopy, Non-negative Matrix Factorization, Classification, Denoising, Spectral Deconvolution

*PACS:* 07.05.Kf, 29.30.Kv, 02.50.Sk

*2020 MSC:* 15A23, 62H25, 65D10

---

\*Corresponding author: raedle.htwk - AT - web.de

## Metadata

Nr.	Code metadata description	Metadata
C1	Current code version	0.9.0
C2	Permanent link to code/repository used for this code version	<a href="https://github.com/Schlavioner/gamma_flow">https://github.com/Schlavioner/gamma_flow</a>
C3	Permanent link to Reproducible Capsule	-
C4	Legal Code License	BSD 3-Clause "New" or "Revised" License
C5	Code versioning system used	git
C6	Software code languages, tools, and services used	Python
C7	Compilation requirements, operating environments & dependencies	Python $\geq$ 3.12, matplotlib, numpy, pandas, scikit_learn, scipy, seaborn, openpyxl
C8	If available Link to developer documentation/manual	<a href="https://github.com/Schlavioner/gamma_flow/blob/main/README.md">https://github.com/Schlavioner/gamma_flow/blob/main/README.md</a> , <a href="https://github.com/Schlavioner/gamma_flow/tree/main/documentation/HTML-documentation">https://github.com/Schlavioner/gamma_flow/tree/main/documentation/HTML-documentation</a>
C9	Support email for questions	raedle.htwk - AT - web.de

Table 1: Code metadata (mandatory)

Nr.	(Executable) software meta-data description	Please fill in this column
S1	Current software version	0.9.0
S2	Permanent link to executables of this version	<a href="https://github.com/Schlavioner/gamma_flow">https://github.com/ Schlavioner/gamma_flow</a>
S3	Permanent link to Reproducible Capsule	-
S4	Legal Software License	BSD 3-Clause "New" or "Revised" License
S5	Computing platforms/Operating Systems	Linux, Microsoft Windows
S6	Installation requirements & dependencies	Python $\geq$ 3.12, matplotlib, numpy, pandas, scikit_learn, scipy, seaborn, openpyxl
S7	If available, link to user manual - if formally published include a reference to the publication in the reference list	<a href="https://github.com/Schlavioner/gamma_flow/blob/main/README.md">https://github.com/ Schlavioner/gamma_flow/blob/ main/README.md</a>
S8	Support email for questions	raedle.htwk - AT - web.de

Table 2: Software metadata (optional)

## 1. Motivation and significance

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 on 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 [1]. In this process, energy deposition in a detector material is simulated, yielding spectra that can be used for model training. However, as not all physical effects are included, simulated spectra and measured spectra from real-world sources may differ, which is a constraint for model performance. On this account, preliminary data exploration is crucial to assess their similarity 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 designed for the real-time analysis of one-dimensional spectra. It was developed in response to the

practical challenges described above and supports the following tasks:

- classification of test spectra to identify their constituents,
- denoising to enhance visibility and reduce measurement noise,
- outlier detection to evaluate the model’s applicability to unknown spectra.

Originally developed for gamma spectroscopy, GAMMA\_FLOW is applicable to various domains including material science [2], chemistry [3], and environmental monitoring [4]. It facilitates automated analysis in settings where fast interpretation of spectral data is essential. By integrating supervised dimensionality reduction with classical signal processing techniques, it extends prior works such as those of Bilton et al. [5] or Bandstra et al. [6], contributing to reproducible and interpretable spectral analysis pipelines.

## 2. Software description

GAMMA\_FLOW is based on a dimensionality reduction model that constitutes a 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.

For spectral denoising, the scores are 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.

### *2.1. Software architecture*

GAMMA\_FLOW consists of three jupyter notebooks that are executed consecutively, as depicted in Figure 1. Each imports functions from a designated Python file, e.g. all functions called in `02_model.ipynb` are found in `tools_model.py`. In addition, the Python files `util.py`, `globals.py` and `plotting.py` provide elementary functions for all three notebooks.

### *2.2. Software functionalities*

In this section, the functionality of the software is outlined, with an emphasis on the mathematical structure of the model. To this end, the procedures realized in the three jupyter notebooks `01_preprocessing.ipynb`, `02_model.ipynb` and `03_outlier.ipynb` are delineated. In the software repository, a sample dataset of gamma spectra with exemplary results is provided to give users a first insight into the software.

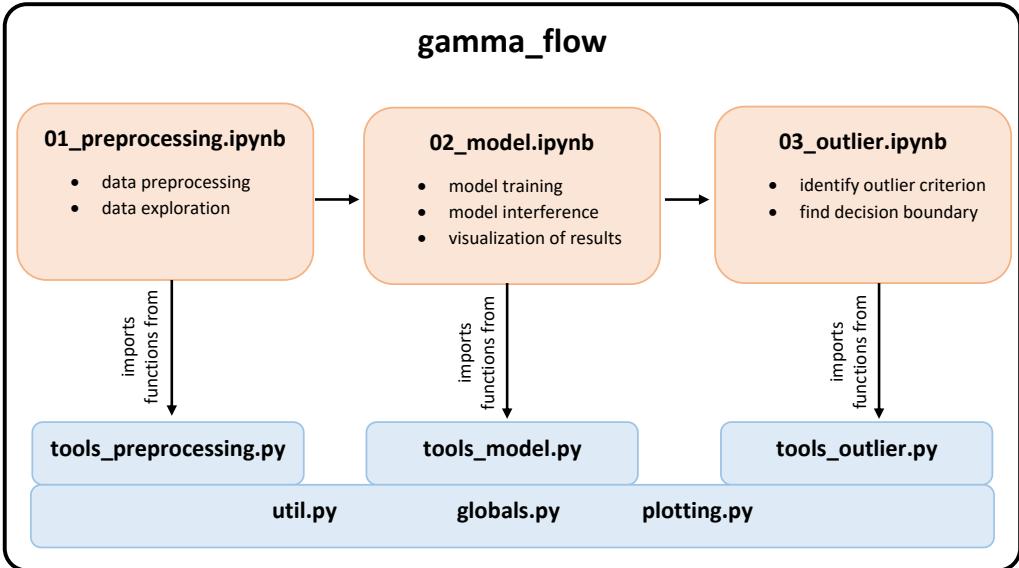


Figure 1: Software architecture of GAMMA\_FLOW: The jupyter notebooks `01_preprocessing.ipynb`, `02_model.ipynb` and `03_outlier.ipynb` are executed sequentially, using functions from different Python files.

### 2.2.1. Preprocessing and data exploration

The notebook `01_preprocessing.ipynb` harmonizes 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 (for class balance).
- 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.

### 2.2.2. 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 model** presented in this paper performs a matrix decomposition of spectral data to achieve dimensionality reduction. More precisely, the original spectra matrix  $\mathbf{X} \in \mathbb{R}_{\geq 0}^{n \times m}$ , with  $n$  spectra and  $m$  channels, is reconstructed by two low-rank matrices:

$$\mathbf{X} \approx \mathbf{S}\mathbf{L}^T$$

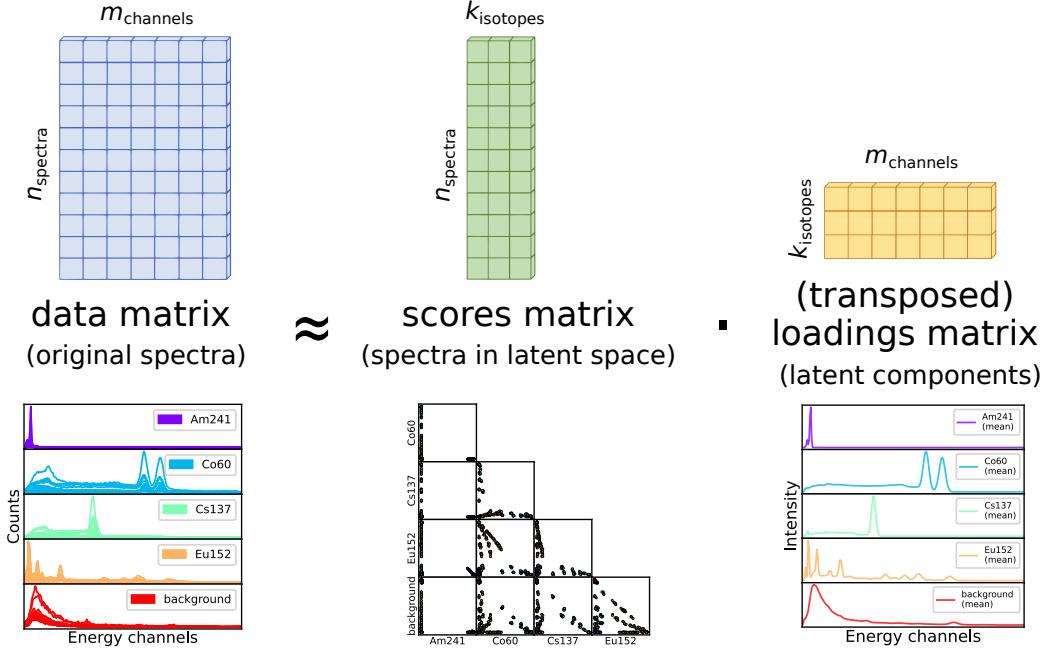


Figure 2: Decomposition of spectral data into scores and loadings, illustrated for five different isotopes. The loadings matrix contains the latent components, corresponding to the mean spectrum of each isotope. The scores matrix provides the representation of each spectrum in the five-dimensional latent space. When normalized, the scores indicate the predicted contribution of each isotope to the measured spectrum.

where  $\mathbf{S} \in \mathbb{R}_{\geq 0}^{n \times k}$  is the scores matrix and  $\mathbf{L} \in \mathbb{R}_{\geq 0}^{m \times k}$  the loadings matrix, with  $k = k_{\text{isotopes}}$  latent components.

As shown in Figure 2, each spectrum is projected into a low-dimensional space spanned by the isotope-specific mean spectra. This corresponds to a supervised variant of Non-negative Matrix Factorization (NMF) [7, 5], where the loadings matrix  $\mathbf{L}$  is predefined instead of being learned. The latent components thus retain physical interpretability, each representing the average spectral signature of one isotope.

To compute the scores  $\mathbf{S}$ , a non-negative least squares (NNLS) problem is

solved for each spectrum  $\mathbf{x}_i$  (row vector of  $\mathbf{X}$ ):

$$\mathbf{s}_i = \operatorname{argmin}_{\mathbf{s} \geq 0} \|\mathbf{x}_i - \mathbf{s}\mathbf{L}^T\|_2^2$$

This constrained optimization is carried out using the `scipy.optimize.nnls` implementation [8] based on the Lawson-Hanson active-set algorithm [9], which solves a sequence of unconstrained least squares problems by iteratively updating an active set of non-negativity constraints. The algorithm is guaranteed to converge to the global minimum of the convex problem for each spectrum and ensures strict non-negativity of the resulting coefficients. This enables an interpretable **spectral decomposition**, where the normalized scores vector  $\mathbf{s}_{i,\text{norm}} = \mathbf{s}_i / \sum \mathbf{s}_i$  directly reveals the relative contribution of each isotope. **Denoised spectra**, on the other hand, are computed by transforming the (non-normalized) scores back into spectral space by multiplication with the loadings matrix:  $\mathbf{x}_i = \mathbf{s}_i \mathbf{L}^T$ . This removes noise orthogonal to the latent space.

To assess **model performance**, the model is trained using spectral data from the specified detectors `dets_tr` and isotopes `isotopes_tr` and inferred (i.e., scores are calculated) on three different test datasets:

1. validation data/holdout data from same detector as used in training (simulated data: each spectrum including only one isotope or pure background)
2. test data from different detector (measured single-label data: each spectrum including one isotope and background)
3. multi-label test data from different detector (measured multi-label data: 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.

### 2.2.3. Outlier analysis

The notebook `03_outlier.ipynb` provides an exploratory approach to outlier 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 metadata are analyzed to distinguish known from unknown spectra. Using a decision tree, the most informative feature for this task is identified based on permutation feature importance, i.e., by measuring how much the performance drops when feature values are randomly shuffled [10]. Next, a decision boundary is derived for this feature, by

- a) using the condition of the first split in the decision tree, as depicted in `03_outlier.ipynb`
- b) fitting a logistic regression (sigmoid function) to the outlier data, with the most informative feature as `x` and outlier score as `y` (by setting known spectra to 0 and outliers to 1). The decision boundary can be read from the fit parameters.
- c) setting a manual threshold by considering accuracy, precision and recall of outlier identification. This is done by plotting those quantities against the most informative feature and visually evaluating their optimal trade-off for this task. Further metrics, like True Positive Rate or False Positive Rate, can be calculated in the corresponding notebook.

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

The excellent model performance on the noisy gamma dataset at hand proves the model's applicability to noisy spectra. In addition, we have verified that adding Gaussian noise to the spectra has only a minor effect on classification and denoising results. In the context of outlier detection, higher ratios of noise result in slightly diminished cosine similarities, shifting the decision boundary towards smaller values.

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

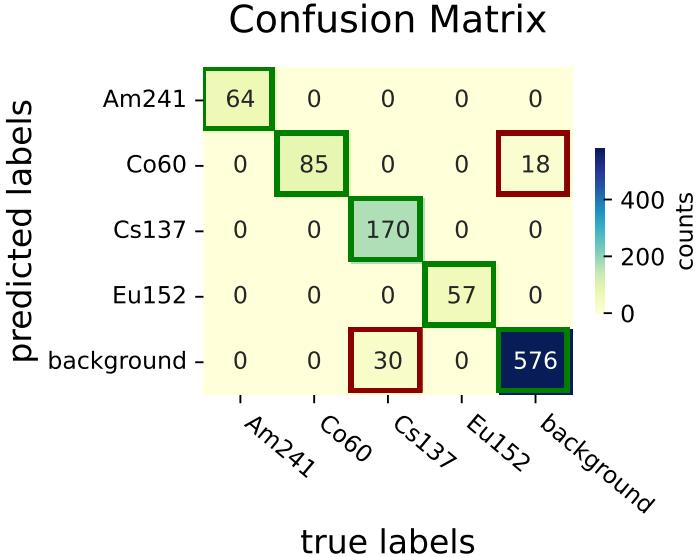


Figure 3: Confusion matrix of single-label classification, with an accuracy of 94.8%. While the model was trained on simulated spectra, measured spectra were used as test data. Misclassifications only occur between an isotope and background, not between different isotopes.

### 3. Illustrative example

The major functions of GAMMA\\_FLOW include classification, denoising, decomposition, and outlier detection of spectra. In this section, these capabilities are demonstrated using an illustrative example.

A model is trained based on simulated spectra of the isotopes Americium-241 ( $^{241}\text{Am}$ ), Cobalt-60 ( $^{60}\text{Co}$ ), Caesium-137 ( $^{137}\text{Cs}$ ), Europium-152 ( $^{152}\text{Eu}$ ), all free of background radiation, as well as pure measured background spectra. The model is then inference on various test datasets to assess its performance in different scenarios (for more detail, see Section 2.2.2). The following example focuses on a measured test dataset where spectra contain either a single isotope combined with background or pure background only. No background subtraction was applied.

The **classification** results for this test dataset can be seen in Figure 3, where the predicted labels are plotted against the ground truth as a confusion matrix. Correct classifications are highlighted with green borders while misclassifications are framed in red. Despite the difference between training and test spectra (different detector and presence of background in test data), an overall accuracy of 94.8% is achieved. Notably, all misclassifications involve confusion between an isotope and pure background — never between two isotopes. Such cases typically occur when the isotope signal is near the detection threshold and the model’s prediction deviates from the manually assigned label. Please refer to the software documentation in the code repository for further visualizations on model training and classification results, such as the loadings, the scores (visualized as scatter matrix), the mean predicted scores by isotope, misclassified example spectra and misclassification statistics.

As described in Section 2.2.2, **denoised spectra** are obtained via inverse transformation using the predefined loadings. The resulting reconstructed spectra are significantly smoother, facilitating the distinction between isotope and random features. A denoised example spectrum of  $^{137}\text{Cs}$  and background is depicted in Figure 4. As can be seen in the annotations, the predictions are correct, with an additional information on the contribution of the components: The model predicts 16% share of  $^{137}\text{Cs}$  and 83% background. The reconstruction of the original spectrum by the denoised spectrum is quantified by their cosine similarity of 0.99 and the explained variance ratio of 98.14%, which both indicate a successful denoising process.

The options for **outlier detection** are analyzed as described in Section 2.2.3.

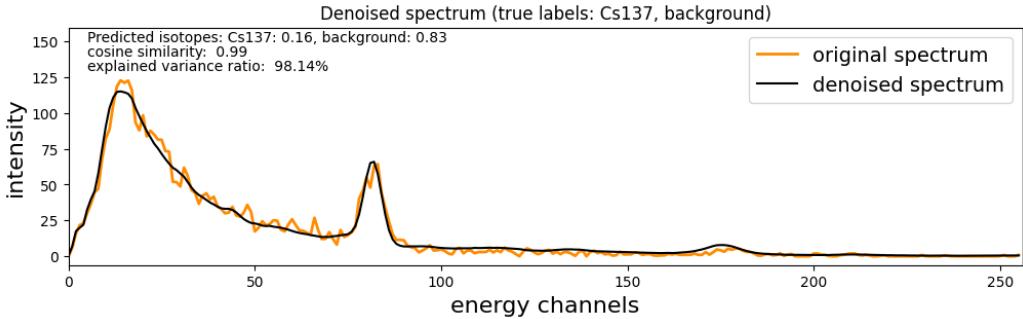


Figure 4: Comparison between the original and the reconstructed denoised spectrum of  $^{137}\text{Cs}$  and background. The denoised spectrum is smoother, leading to an easy and fast interpretability.

After the most informative feature is determined (in our case the cosine similarity between original and denoised spectrum) a decision boundary is derived. This can be done manually by considering accuracy, precision and recall, as visualized in Figure 5. Those quantities yield different information on the separation of known and unknown spectra:

- **Accuracy** quantifies the correct classifications as known and unknown spectra: 
$$\frac{\text{True Positives} + \text{True Negatives}}{\text{True Positives} + \text{False Positives} + \text{True Negatives} + \text{False Negatives}}$$
- **Precision** quantifies how many of the predicted outliers are actually unknown spectra: 
$$\frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$
- **Recall** quantifies how many of the unknown spectra are found: 
$$\frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

In this example, a threshold between 0.5 and 0.9 would be an adequate choice.

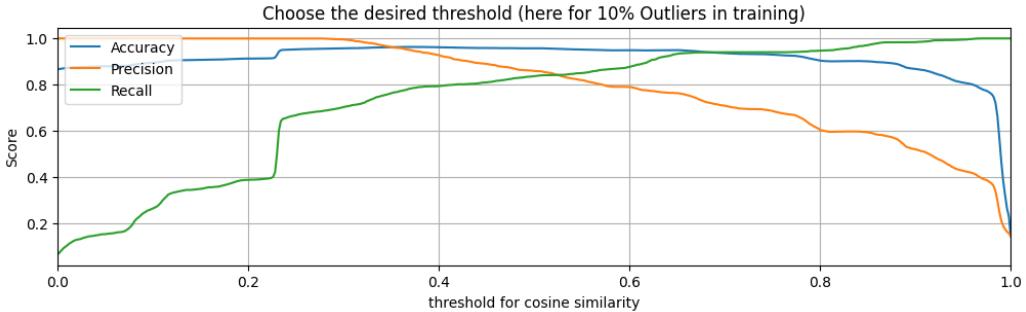


Figure 5: For outlier detection, the decision boundary can be derived based on accuracy, precision and recall of the outliers dataset.

## 4. Impact and Discussion

### 4.1. Applications

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 [11, 12]. In addition, the underlying mechanisms are usually not revealed and may act as a black-box system to the user [13]. On top of that, a spectral comparison is typically only possible for spectra of pure probes [14]. However, there may be a need to decompound multi-label spectra (linear combinations of different substances) and identify their constituents.

GAMMA\_FLOW closes this gap and provides a lean and performant model. 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 [15, 16]
- **mass spectrometry** for protein identification in snake venom [17, 18]
- **Raman spectroscopy** for analysis of complex pharmaceutical mixtures and detection of dilution products like lactose [19]
- **UV-Vis spectroscopy** for detection of pesticides in surface waters [20, 4]
- **stellar spectroscopy** to infer the chemical composition of stars [21]

#### *4.2. Comparison to Related Work*

Clustering and classification of spectral data have been explored through various NMF-based methods and have proven to excel at grouping unlabelled spectra [22]. To target robustness or interpretability, Robust and Sparse NMF [23, 24, 25, 26] mitigate noise and outliers through explicit error modelling or sparsity regularization. Meanwhile, Graph-regularized and Orthogonal NMF [24, 27, 28, 29] preserve local data structure or reduce component redundancy. While effective for clustering, these methods either alter small but relevant isotope contributions or impose constraints (e.g., orthogonality) that conflict with overlapping peaks in different isotope spectra. GAMMA\_FLOW instead uses a cosine similarity check for outlier detection while preserving all physically relevant peaks.

Subspace-based extensions [30, 31] learn low-dimensional manifolds or incorporate subspace priors to improve robustness. In contrast, GAMMA\_FLOW defines the latent space a priori from reference spectra, enabling supervised and physically interpretable unmixing.

Finally, supervised NMF approaches [32, 33, 34] often learn latent components and coefficients jointly via classification loss. Meanwhile, GAMMA\_FLOW fixes the basis from labelled training data and estimates coefficients via non-negative least squares, which reduces computational cost while ensuring physical interpretability.

#### 4.3. Limitations

The applicability of GAMMA\_FLOW is limited by its **reliance on a predefined basis**, making it suitable only when relevant isotopes are known and labeled spectra are available. If only mixtures are present, an initial NMF decomposition is required to obtain training spectra.

Furthermore, the **absence of orthogonality constraints** makes performance strongly dataset-dependent: similar isotope spectra can be confused, which may **limit scalability**. Preliminary tests successfully reduced 256 channels to 18 components, but the upper limit remains dataset-specific. Therefore, the correlation between latent components should be assessed based on their cosine similarity matrix (data exploration in `01_preprocessing.ipynb`).

In addition, no intrinsic robustness exists against **mislabeled training spectra**, requiring careful expert validation. **Preprocessing choices**, especially rebinning parameters, can also affect results: parameters far outside the calibration range may cause data loss or resolution degradation, though this step can be skipped or adapted for other channel-energy relationships.

Broader **applicability to non-gamma spectra** is possible, but the fixed-basis approach assumes that loadings match the resolution and peak shapes of the target data. Substantially broader or asymmetric peaks may reduce reconstruction accuracy and classification performance, but this is mitigated by deriving the loadings from spectra acquired under comparable conditions.

In theory, the non-negativity constraint only applies to the scores matrix, enabling GAMMA\_FLOW to analyze spectra containing both positive and negative signals, e.g. difference spectra. However, as those datasets are uncommon in most spectroscopic applications, the algorithm has not yet been tested for these use cases.

## 5. Conclusions

In this paper, we introduced GAMMA\_FLOW, an open-source, Python-based framework for the analysis of one-dimensional spectra. The approach combines established data science techniques in a resource-efficient way, enabling real-time analysis at low computational cost. Based on a supervised variant of non-negative matrix factorization, it constructs a fast and interpretable model for classification, decomposition, denoising, and outlier detection of spectral data. The method is applicable to both pure (single-label) and mixed (multi-label) spectra.

For a user-friendly implementation, the software is organized into three jupyter notebooks. Each notebook combines explanatory text with code, guiding users step-by-step through data preprocessing and exploration, model training and evaluation, and outlier detection. This enables users not only to apply the model to spectral data but also develop a deep understanding of the underlying processes - without requiring advanced mathematical or programming expertise.

GAMMA\_FLOW bridges the gap between complex, often proprietary spectral analysis tools and the limitations of manual interpretation in laboratory workflows. By offering an open, transparent, and extensible framework, it empowers researchers to move beyond traditional black-box approaches. With its modular architecture and general applicability, it opens new av-

venues for innovation across disciplines — from nuclear safety and materials science to environmental monitoring and beyond. As such, GAMMA\_FLOW has the potential to accelerate scientific discovery wherever spectral data plays a role.

## Acknowledgements

We gratefully acknowledge the support provided by the Federal Ministry for the Environment, Climate Action, Nature Conservation and Nuclear Safety (BMUKN), 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

- [1] J. A. Kulesza, T. R. Adams, J. C. Armstrong, S. R. Bolding, F. B. Brown, J. S. Bull, T. P. Burke, A. R. Clark, R. A. A. Forster III, J. F. Giron, T. S. Grieve, C. J. Josey, R. L. Martz, G. W. McKinney, E. J. Pearson, M. E. Rising, C. J. C. Solomon Jr., S. Swaminarayan, T. J. Trahan, S. C. Wilson, A. J. Zukaitis, MCNP® Code Version 6.3.0 Theory & User Manual, Tech. Rep. LA-UR-22-30006, Los Alamos National Laboratory (LANL), Los Alamos, NM (United States) (Sep. 2022). doi:10.2172/1889957.
- [2] D. A. Zatsepин, A. F. Zatsepин, Ultraviolet Photoelectron Spectroscopy – Materials Science Technique, in: Spectroscopy for Materials Char-

- acterization, John Wiley & Sons, Ltd, 2021, Ch. 13, pp. 383–403.  
doi:10.1002/9781119698029.ch13.
- [3] M. M. Roessler, E. Salvadori, Principles and applications of EPR spectroscopy in the chemical sciences, *Chemical Society Reviews* 47 (8) (2018) 2534–2553. doi:10.1039/C6CS00565A.
  - [4] X. Qi, Y. Lian, L. Xie, Y. Wang, Z. Lu, Water quality detection based on UV-Vis and NIR spectroscopy: A review, *Applied Spectroscopy Reviews* 59 (8) (2024) 1036–1060. doi:10.1080/05704928.2023.2294458.
  - [5] K. J. Bilton, T. H. Joshi, M. S. Bandstra, J. C. Curtis, B. J. Quiter, R. J. Cooper, K. Vetter, Non-negative Matrix Factorization of Gamma-Ray Spectra for Background Modeling, Detection, and Source Identification, *IEEE Transactions on Nuclear Science* 66 (5) (2019) 827–837. doi:10.1109/TNS.2019.2907267.
  - [6] M. S. Bandstra, T. H. Y. Joshi, K. J. Bilton, A. Zoglauer, B. J. Quiter, Modeling Aerial Gamma-Ray Backgrounds Using Non-negative Matrix Factorization, *IEEE Transactions on Nuclear Science* 67 (5) (2020) 777–790. doi:10.1109/TNS.2020.2978798.
  - [7] P. Shreeves, Dimensionality reduction techniques with applications in Raman spectroscopy - UBC Library Open Collections, Ph.D. thesis, University of British Columbia (2020).
  - [8] SciPy, Official code documentation of `scipy.optimize.nnls`, API Reference, URL: <https://docs.scipy.org/doc/scipy-1.16.0/reference/generated/scipy.optimize.nnls.html> , SciPy v1.16.0 Manual.

- [9] C. L. Lawson, R. J. Hanson, Solving Least Squares Problems, Society for Industrial and Applied Mathematics, 1995. doi:10.1137/1.9781611971217.
- [10] scikit-learn, 5.2. Permutation feature importance, [https://scikit-learn/stable/modules/permuation\\_importance.html](https://scikit-learn/stable/modules/permuation_importance.html) (2025).
- [11] H. Lam, Building and Searching Tandem Mass Spectral Libraries for Peptide Identification, *Molecular & Cellular Proteomics* 10 (12) (Dec. 2011). doi:10.1074/mcp.R111.008565.
- [12] J. Nasereddin, M. Shakib, Ira: A free and open-source Fourier transform infrared (FTIR) data analysis widget for pharmaceutical applications, *Analytical Letters* 56 (16) (2023) 2637–2648. doi:10.1080/00032719.2023.2180516.
- [13] L. El Amri, A. Chetaine, H. Amsil, B. El Mokhtari, H. Bounouira, A. Didi, A. Benchrif, K. Laraki, H. Marah, New open-source software for gamma-ray spectra analysis, *Applied Radiation and Isotopes* 185 (2022) 110227. doi:10.1016/j.apradiso.2022.110227.
- [14] W. Cowger, Z. Steinmetz, A. Gray, K. Munno, J. Lynch, H. Hapich, S. Primpke, H. De Frond, C. Rochman, O. Herodotou, Microplastic Spectral Classification Needs an Open Source Community: Open Specy to the Rescue!, *Analytical Chemistry* 93 (21) (2021) 7543–7548. doi:10.1021/acs.analchem.1c00123.
- [15] B. Ferreiro, R. Leardi, E. Farinini, J. M. Andrade, Supervised classification combined with genetic algorithm variable selection for a fast identification of polymeric microdebris using infrared reflectance, Ma-

rine Pollution Bulletin 195 (2023) 115540. doi:10.1016/j.marpolbul.2023.115540.

- [16] Q. T. Whiting, K. F. O'Connor, P. M. Potter, S. R. Al-Abed, 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) (2022) 8353–8364. doi:10.1007/s00216-022-04371-2.
- [17] A. Zelanis, D. A. Silva, E. S. Kitano, T. Liberato, I. Fukushima, S. M. T. Serrano, A. K. Tashima, 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 (2019) 100599. doi:10.1016/j.cbd.2019.100599.
- [18] N. Ç. Yasemin, A. Izzet, K. M. Ali, K. Selçuk, S. Bekir, Identification of Snake Venoms According to their Protein Content Using the MALDI-TOF-MS Method, Analytical Chemistry Letters 11 (2) (2021) 153–167. doi:10.1080/22297928.2021.1894974.
- [19] X. Fu, L.-m. Zhong, Y.-b. Cao, H. Chen, F. Lu, Quantitative analysis of excipient dominated drug formulations by Raman spectroscopy combined with deep learning, Analytical Methods 13 (1) (2021) 64–68. doi:10.1039/D0AY01874K.
- [20] Y. Guo, C. Liu, R. Ye, Q. Duan, Advances on Water Quality Detection by UV-Vis Spectroscopy, Applied Sciences 10 (19) (2020) 6874. doi:10.3390/app10196874.
- [21] D. F. Gray, The Observation and Analysis of Stellar Photospheres, <https://www.cambridge.org/highereducation/books/the-observation-and-analysis-of-stellar-photospheres>

and-analysis-of-stellar-photospheres/67B340445C56F4421BCBA0AFFAAFDEE0  
(Dec. 2021). doi:10.1017/9781009082136.

- [22] A. Mirzal, Statistical Analysis of Microarray Data Clustering using NMF, Spectral Clustering, Kmeans, and GMM, *IEEE/ACM Transactions on Computational Biology and Bioinformatics* 19 (2) (2022) 1173–1192. doi:10.1109/TCBB.2020.3025486.
- [23] L. Zhang, Z. Chen, M. Zheng, X. He, Robust non-negative matrix factorization, *Frontiers of Electrical and Electronic Engineering in China* 6 (2) (2011) 192–200. doi:10.1007/s11460-011-0128-0.
- [24] Y. Chen, G. Qu, J. Zhao, Orthogonal graph regularized non-negative matrix factorization under sparse constraints for clustering, *Expert Systems with Applications* 249 (2024) 123797. doi:10.1016/j.eswa.2024.123797.
- [25] K. Qu, Z. Li, A Fast Sparse NMF Optimization Algorithm for Hyperspectral Unmixing, *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* 17 (2024) 1885–1902. doi:10.1109/JSTARS.2023.3341583.
- [26] L. Mella, A. Lal, F. Angaroni, D. Maspero, R. Piazza, A. Sidow, M. Antoniotti, A. Graudenzi, D. Ramazzotti, SparseSignatures: An R package using LASSO-regularized non-negative matrix factorization to identify mutational signatures from human tumor samples, *STAR Protocols* 3 (3) (2022) 101513. doi:10.1016/j.xpro.2022.101513.
- [27] Z. Liu, F. Zhu, H. Xiong, X. Chen, D. Pelusi, A. V. Vasilakos, Graph regularized discriminative nonnegative matrix factorization, *Engineering*

Applications of Artificial Intelligence 139 (2025) 109629. doi:10.1016/j.engappai.2024.109629.

- [28] W. Li, J. Zhao, Y. Chen, Orthogonal-Constrained Graph Non-Negative Matrix Factorization for Clustering, Symmetry 17 (7) (2025) 1154. doi:10.3390/sym17071154.
- [29] S. Wang, T.-H. Chang, Y. Cui, J.-S. Pang, Clustering by Orthogonal NMF Model and Non-Convex Penalty Optimization, IEEE Transactions on Signal Processing 69 (2021) 5273–5288. doi:10.1109/TSP.2021.3102106.
- [30] N. Naseri, M. Eftekhari, F. Saberi-Movahed, M. Radjabalipour, L. A. Belanche, A similarity measure based on subspace distance for spectral clustering, Neurocomputing 620 (2025) 129187. doi:10.1016/j.neucom.2024.129187.
- [31] L. Zhou, X. Zhang, J. Wang, X. Bai, L. Tong, L. Zhang, J. Zhou, E. Hancock, Subspace Structure Regularized Nonnegative Matrix Factorization for Hyperspectral Unmixing, IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing 13 (2020) 4257–4270. doi:10.1109/JSTARS.2020.3011257.
- [32] J. Leuschner, M. Schmidt, P. Fernsel, D. Lachmund, T. Boskamp, P. Maass, Supervised non-negative matrix factorization methods for MALDI imaging applications, Bioinformatics 35 (11) (2019) 1940–1947. doi:10.1093/bioinformatics/bty909.
- [33] H. Lee, J. Yoo, S. Choi, Semi-Supervised Nonnegative Matrix Factorization, IEEE Signal Processing Letters 17 (1) (2010) 4–7. doi:10.1109/LSP.2009.2027163.

- [34] V. Bisot, R. Serizel, S. Essid, G. Richard, Supervised nonnegative matrix factorization for acoustic scene classification, in: IEEE International Evaluation Campaign on Detection and Classification of Acousitc Scenes and Events (DCASE 2016), Budapest, Hungary, 2016.