

The Spiked Mixture Model: A New Clustering Approach for Imaging Mass Spectrometry Offers Advanced Resilience to Noise

The Spiked Mixture Model (SMM)

 $\alpha \mathbf{x_1} + \varepsilon$ with probability π_1 $\alpha \mathbf{x_2} + \varepsilon$ with probability π_2 $\mathbf{y} =$ $\alpha \mathbf{x}_{\mathbf{K}} + \varepsilon$ with probability $\pi_{\mathbf{K}}$

$$\begin{aligned} \alpha &\sim \mathcal{N}(\mathbf{0}, \mathbf{1}), \ \varepsilon &\sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d) \\ \sum_{k=1}^{K} \pi_k = \mathbf{1}, \ \mathbf{x}_1, \dots, \mathbf{x}_K \in \mathcal{R}^d \end{aligned}$$

Imaging mass spectrometry (IMS):

- IMS is a molecular imaging technique that combines spatial mapping with mass spectral analysis.
- IMS offers detailed chemical maps of organic tissue samples, measuring the spatial distributions of hundreds of molecular species concurrently.
- Each pixel records a full mass spectrum or m/z profile, and prior labeling of compounds is not required.

Challenges:

- 1. IMS experiments can acquire vast amounts of spatially resolved data.
- 2. IMS data can be noisy.

Modeling IMS data using the Spiked Mixture Model:

- ullet Every observation or pixel ${f y}$ in an IMS dataset is modeled as a randomly scaled representative of a biological signature, subpopulation, or spike $\mathbf{x}_{\mathbf{K}}$.
- α is the random scaling factor of observation y:
- ε is the random noise of observation y:
- z is a latent categorical variable indicating which spike x_K is underlying y: $\mathbf{z} \sim \operatorname{cat}(\pi)$

Recovering signals underlying noisy IMS measurements (and clustering them):

ullet Given noisy observations $\{y_1,\ldots y_N\}$, find f K underlying subpopulations x_1,\ldots,x_K

i.e., find parameters of interest: $\theta = \{x_1, \dots, x_K, \pi_1, \dots, \pi_K, \sigma^2\}$

that deliver a maximum likelihood estimate (MLE) for the given data: $\theta_{MLE} = \arg \max \sum \log (p(y_i \mid \theta))$

Custom EM algorithm for SMM-based recovery and clustering:

- MLE is intractable in this case.
- We developed an Expectation-Maximization (EM) approach, customized for SMM, to find a candidate MLE that: • is guaranteed to converge to a local maximum; and
- does not require matrix inversion of large matrices.

Our custom EM algorithm is described in https://doi.org/10.48550/arXiv.2501.01840.

- $\theta = \{x_1, ..., x_K, \pi_1, ..., \pi_K, \sigma^2\}$
- 1 Delft Center for Systems and Control, Delft University of Technology, Delft, Netherlands
- 2 Karlsruhe Institute of Technology, Institute of Industrial Information Technology, Karlsruhe, Germany 6 Department of Biochemistry, Vanderbilt University, Nashville, TN, USA
- 3 Mass Spectrometry Research Center, Vanderbilt University, Nashville, TN, USA
- 4 Chemical Physical Biology Program, Vanderbilt University, Nashville, TN, USA





SMM substantially outperforms GMM.

Conclusions:



that might lie hidden in the noise.

- 5 Department of Cell and Developmental Biology, Vanderbilt University, Nashville, TN, USA
- 7 Department of Chemistry, Vanderbilt University, Nashville, TN, USA

$$lpha \sim \mathcal{N}(\mathbf{0}, \mathbf{1}), \, \varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I_d})$$

le x_K is underlying y: $\mathbf{z} \sim \operatorname{cat}(\pi)$

Under low-noise conditions, SMM performs comparably to GMM, but in high-noise regimes, where robust computational methods are most needed,

Overall, the SMM offers a means of making algorithms aware of structures and perturbations that naturally arise in MS measurements, enabling advanced noise resilience and aiding discovery of biological patterns

Jeffrey M. Spraggins^{3,5,6,7}, Lukasz G. Migas¹, Raf Van de Plas^{1,3,6}

Results – Clustering of Rat Brain IMS dataset

Goal: Use SMM to recover underlying molecular signatures from real-world noisy IMS measurements of a rat brain tissue section, and to implicitly cluster the mass spectra, i.e., segment the tissue section into areas of similar chemical composition. Furthermore, we compare SMM's clustering results to ones given by traditional methods such as GMM and k-means clustering.

Rat brain IMS dataset:

A transverse rat brain section was measured using a timsTOF FleX (Bruker Daltonics) in QTOF mode across m/z 400-2,000, using a 10- μ m pixel size and yielding 572,832 individual spectra. Spectral alignment was performed to correct for drift along the m/z domain. After alignment and calibration, an average mass spectrum based on all pixels in the dataset was computed. The average spectrum was peak-picked, 843 peaks were detected, and their corresponding ion intensities were retrieved.

Comparison of clustering results by SMM, GMM, k-Means Clustering:



Spectral perspective

Research was supported by the National Institutes of Health (NIH)'s NIDDK (U54DK134302 and U01DK133766), NEI (U54EY032442), NIAID (R01AI138581 and R01AI145992), NIA (R01AG078803), NCI (U01CA294527), and by the Chan Zuckerberg Initiative DAF (2021-240339 and 2022-309518). The content is solely the responsibility of the authors and does not necessarily represent the official views of the NIH. The authors declare no competing financial interest.

Paul-Louis Delacour¹, Sander Wahls², Lauren N. Emmerson^{3,4}, Madeline E. Colley^{3,5},





🔁 Cerebral cortex laye

- Estimated responsibility variables segment the tissue according to molecular content. SMM retrieves histological patterns missed by other methods.
- SMM found subdivisions of cerebral cortex layers (e.g., molecular, granular, pyramidal, and multiform) that are known to align parallel to the surface of the brain, but that were missed by GMM.
- SMM delivered sharper delineation of anatomical structures (e.g., white matter, molecular layer, and granule cell layer) and exhibited less susceptibility to noise than k-means clustering.



