Unsupervised Connectome Learning in fMRI State Mixtures

Louis-Alexandre Leger¹ Ilaria Ricchi¹ Dimitri Van De Ville¹

Abstract

Gaussian Mixture Models are universal in their ability to approximate continuous densities and mathematical tractability has made them very popular for unsupervised learning. In graph learning, recent advances propose using the learned statistics of mixture models to infer graph structures and this methodology has proven great success in learning connectomes for different brain states from an fMRI signal. Modeling these fMRI signals as low dimensional signals of brain region activity, we can infer graphs during specific subject activities such as reading, writing, conversation or playing music. While connectome learning has provided strong insights into brain states, understanding the spinal cord's graph structure remains limited. In this document, we show the inferred state connectomes of the spinal cord with the unsupervised Graph Laplacian Mixture Model. We find laplacians highly correlated with prior knowledge despite a poor signal and show the limits of spinal cord fMRI by parametrizing the feasibility of unsupervised clustering with multivariate Gaussians. Our main findings underscore the importance of having a combination of sufficiently linearly separable or orthogonal data for Gaussian mixture modeling and for the spinal cord data an additional 1.2 euclidean distance between means or $\frac{\pi}{18}$ rotation of covariance matrices.

1. Introduction

Neurons and axons are the nodes and edges of the graph structure that integrates information responsible for perception, cognition, behavior and motor function. Understanding the connectome or wiring diagram of our central nervous system is essential for treating neurological disorders like Alzheimer's disease (Yu et al., 2021) and restoring spinal cord sensorimotor functions (Zhang et al., 2022). Recent advances in diffusion tensor imaging have allowed researchers to identify full structure of fly brains and human brain regions (Seung, 2024). The study of neural graphs is presented in two separate modalities: the structural connectome which relates to the physical wiring of neurons and the functional connectome which relates to the active graph during a subject's state. Functional Magnetic Resonance Imaging or fMRI is one of the main tools for studying functional connectomes correlated with different states such as resting, watching a movie or reading. Integrating neuronal activity for different brain regions and computing the pearson correlation matrix, we can infer a functional connectome for a time-series of fMRI volumes. In a time series where a subject's state is dynamically changing from rest to others, unsupervised learning techniques such as Graph Laplacian Mixture Model (Maretic & Frossard, 2020) have been shown to produce high quality inferred graphs (Ricchi et al., 2022) correlated with the different states of subjects.

From the central nervous system a lot of studies focus on the brain's functional graphs (Preti et al., 2017) and few target the spinal cord. The spinal cord presents unique set of constraints for graph inference, including its relatively small size, high density of interconnected neurons, and its structural organization into distinct functional regions, levels and relative contribution to function with the brain. The spatial and temporal resolution of fMRI pose as a challenge for extracting meaningful graph structures. In this work, we investigate the feasibility of unsupervised connectome learning in spinal cord fMRI using GLMM. We analyze clustering performance in both real and synthetic data and identify key limitations due to separability and covariance properties.

2. Methods

2.1. Data Acquisition

Data for this project originated from spinal cord fMRI experiments (Kinany et al., 2022) involving bilateral wrist adduction movements of 8 blocks of 18s interspersed with resting periods of equal duration. fMRI acquisition was run for an TR = 2.5s and covered four spinal levels C6, C7, C8 and T1. Regions of interest were identified with 7 PAM50 parcellations (De Leener et al., 2017): Gray Matter Dorsal, Ventral and Intermediate regions and White Matter Corticospinal Tract, Fasciculus Cuneatus, Fasciculus Gracilis and Spinothalamic Tract regions. With both left and right regions this results in a total of 56 regions of interest (ROIs). The fMRI acquisition was done for 15 subjects on two separate runs, which leads to design matrix X of shape

(3750, 56), representing activity across time and regions.



Figure 1. Illustration of Activity Paradigm followed during fMRI acquistion.

2.2. Method: Graph Laplacian Mixture Model

The Graph Laplacian Mixture Model (GLMM) is a probabilistic framework for clustering high-dimensional signals that naturally reside on different graphs. This model extends the traditional Gaussian Mixture Model (GMM) by incorporating graph structure. Given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, W)$ with N nodes, an undirected weighted adjacency matrix W, and laplacian L we define a graph signal (Kalofolias, 2016), (Kalofolias & Perraudin, 2019) $x \in \mathbb{R}^N$ as a function over the M nodes. A key assumption in graph signal processing is smoothness, which penalizes signal variations across strongly connected nodes:

$$x^{T}Lx = \frac{1}{2} \sum_{i,j} W_{ij} (x_{i} - x_{j})^{2}.$$
 (1)

A more general representation considers a graph filter g(L), modeling a kernel process:

$$x = \mu + g(L)w, \quad w \sim \mathcal{N}(0, I), \tag{2}$$

resulting in:

$$x \sim \mathcal{N}(\mu, g^2(L)).$$
 (3)

For smooth signals, this is often approximated using the pseudo-inverse of the Laplacian:

$$g(L) = \sqrt{L^{\dagger}}.$$
 (4)

The model as a mixture of multivariate gaussians can be formalized as follows with $\gamma_{m,k}$ the posterior probability that signal x_m belongs to cluster k

$$p(x_m) = \sum_{k=1}^{K} \alpha_k \mathcal{N}(x_m | \mu_k, g_k^2(L_k)).$$
(5)

$$\gamma_{m,k} = \frac{\alpha_k \mathcal{N}(x_m | \mu_k, g_k^2(L_k))}{\sum_{l=1}^K \alpha_l \mathcal{N}(x_m | \mu_l, g_l^2(L_l))}.$$
 (6)

We estimate the model parameters via an Expectation-Maximization (EM) algorithm. We iteratively update the posterior probabilities $\gamma_{m,k}$ starting for some seeded parameters and update the parameters α_k , μ_k , and L_k by maximizing the expected log-likelihood:

$$\sum_{m=1}^{M} \sum_{k=1}^{K} \gamma_{m,k} \left(\ln \alpha_k + \ln \mathcal{N}(x_m | \mu_k, g_k^2(L_k)) + \ln p(L_k) \right)$$
(7)

To infer the graph Laplacian L_k , we impose structural priors that enforce sparsity and smoothness constraints:

$$\ln p(L_k) = -\beta_{1,k} \text{tr}(\mathbf{1}^T \ln \text{diag}(L_k)) + \beta_{2,k} \|L_k\|_F^2.$$
(8)

 β_1 strengthens connectivity by promoting high node degrees. β_2 penalizes excessive off-diagonal weights, encouraging sparsity. For the rest of document, we use surrogate parame-

ters
$$\Delta = \sqrt{\frac{\beta_1}{\beta_2}}$$
 and $\theta = \sqrt{\frac{1}{\beta_1 \beta_2}}$

2.3. Evaluation Metrics

To evaluate the performance of our GLMM clustering, we employed both classical clustering metrics and a γ -F1-score that accounts for the overlap between the derived clusters and the activity paradigm labels (i.e., whether the subject was resting or not).

The γ -F1-score is computed as:

$$TP = \max_{k \in K} \gamma_k > \frac{1}{2} \cap A \tag{9}$$

F1 score =
$$\frac{\text{TP}}{\text{TP} + \frac{1}{2}(\text{FP} + \text{FN})}$$
(10)

where TP, FP and FN correspond to the true positives, false positives, and false negatives in the binary classification of rest vs. activity states based on the clustering results $A \in \{0, 1\}^N$.

Additionally, we utilized the Silhouette Score measuring cluster separation:

$$S(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}$$
(11)

where a(i) represents the average intra-cluster distance, and b(i) is the lowest average inter-cluster distance for a given data point.

2.4. Synthetic Data Generation

We parametrise the feasability of the spinal cord's clustering by parametrising the data generation process. We choose two hyperparameters R and Θ to create multiple clusters



Figure 2. Illustration of Synthetic Data Generation Parameters R and θ to simulate Spinal Cord Clustering Feasibility

with different means and orientations and then compare how well a clustering algorithm can identify these clusters. R is the parameter regulating the linear separability between clusters, which is the euclidean distance between cluster means μ_k . θ is the parameter that characterizes the rotation angle between clusters, the orthogonality between covariances or more concretely the angle between covariance matrix eigenvectors defined via Givens rotations or angles between flats.

To determine the means and covariance matrices of our gaussian mixture. We place the first cluster's mean at the origin,

$$\mu_0 = 0.$$

For the remaining clusters, we draw random directions uniformly by sampling from a standard normal and normalizing:

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, I_d), \quad \mu_k = \frac{\mathbf{z}}{\|\mathbf{z}\|} \cdot \mathbf{R}$$

Ensuring the mean vector lies on the surface of a hypersphere for a given radius \mathbf{R} .

For the covariance matrices Σ_k , we first generate a random matrix $A \in \mathbb{R}^{d \times d}$ with entries drawn from a normal distribution, and then form

$$\Sigma = AA^{+} + \epsilon I.$$

We also scale Σ to emulate standardized datasets with maximal variance = 1:

$$\Sigma \leftarrow \frac{\Sigma}{\max(\operatorname{diag}(\Sigma))}.$$

Then to systematically rotate the eigenvectors of a given covariance matrix Σ , we compute the eigen-decomposition $\Sigma = U\Lambda U^{\top}$, where Λ is the diagonal matrix of eigenvalues, and U is an orthonormal matrix of eigenvectors. We then apply a product of Givens rotations to U. A Givens rotation in the (i, j)-plane by an angle θ has the general block form (in 2D):

$$G(i, j, \theta) = \begin{bmatrix} 1 & \cdots & 0 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & & \vdots & & \vdots \\ 0 & \cdots & c & \cdots & -s & \cdots & 0 \\ \vdots & & \vdots & \ddots & \vdots & & \vdots \\ 0 & \cdots & s & \cdots & c & \cdots & 0 \\ \vdots & & \vdots & & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & 0 & \cdots & 1 \end{bmatrix}$$

Where $c = \cos(\theta)$, $s = \sin(\theta)$ and in higher dimensions, we multiply these Givens rotations into U successively, for selected pairs (i, j). If we initialize the null rotation as

$$M = I_{d \times d},$$

then for each pair (i, j) we define the corresponding Givens rotation matrix G and update

$$M \leftarrow M G.$$

The final rotated eigenvector matrix is $U_{rot} = U M$. The rotated covariance matrix is then:

$$\Sigma_{\rm rot} = U_{\rm rot} \Lambda U_{\rm rot}^+$$

Lastly, having determined:

- μ_k : the mean of cluster k,
- Σ_k : the covariance of cluster k (possibly rotated),

we sample points

$$\mathbf{x}_i \sim \mathcal{N}(\mu_k, \Sigma_k), \quad i = 1, \dots, n_k.$$

All generated samples are then combined into a single dataset. An integer label is associated with each sample indicating its cluster membership and we shuffle the data points to emulate a synthetic random activity paradigm.

2.4.1. Computing Spinal Cord θ

To compute an observed data generation's θ statistic, this is not as trivial as it's R parameter which is just the euclidean distance between the mean activation of the data points of each state. For θ of two covariance matrices $\Sigma_a, \Sigma_b \in \mathbb{R}^{d \times d}$ of active and rest state with eigenvalue decompositions

$$\Sigma_a = V_a \Lambda_a V_a^{\top}, \quad \Sigma_b = V_b \Lambda_b V_b^{\top}$$

For each $k \in \{1, 2, ...\}$, form the subspaces

 $U_a = [v_{a,d}, v_{a,d-1}, \dots, v_{a,d-k+1}],$ $U_b = [v_{b,d}, v_{b,d-1}, \dots, v_{b,d-k+1}].$

Then compute

$$M = U_a^\top U_b$$

and let

$$M = P S, Q^{\top}$$

be its singular value decomposition, where $S = diag(s_1, \ldots, s_k)$. The principal angles θ_i satisfy

$$\theta_i = \arccos(s_i).$$

If all $s_i = 1$, the subspaces coincide, and the iteration is stopped. The final output is the mean of these angles across all non-collapsed iterations:

Mean Principal Angle
$$=rac{1}{k}\sum heta_i.$$

3. Experiments and Results

3.1. Data Analysis

Applying the GLMM to the spinal cord data and I observed very poor results. Unsatisfying results are caputred by the probabilities of time points averaged across patients not aligning with the activity paradigm. I would get low F1-scores ranging from 0.4 to 0.5. The observed and recurring behavior of the GLMM with K = 2 was the assignment of ~ 95% points to 1 cluster and extremely unbalanced clusters.

I would observe this same low signal-to-noise ratio for different subsets, slices and transformations of the input data Xsuch as taking different subject sub-groups, selecting different subsets of spinal cord dimensions or regions of interest and taking the iCAPs (Karahanoğlu & Van De Ville, 2015) or the derivative of the signal to find perhaps first-order clusters. I also did shallow grid searches to optimize the 3 main hyper-parameters $\{K, \Delta, \Theta\}$ where Θ relates to the average degree of the graph, arriving at similar conclusions.

Following these results, I did some exploratory data analysis to verify the presence of a signal in our data. We can see on Fig. 3 (A and B) traditional functional connectivity matrices and can see visual differences in the correlations of regions of interest between the two. I also computed the pearson correlation of every dimension separately with the activity paradigm observed the p-value of this statistic. In 3.C we can see the distributions data for each dimension ordered in decreasing order of significant correlation with the acitivity paradigm. In our data, the C6 Left Ventral Gray Matter Region is the most significantly correlated with the activity paradgim with a p-value = 6.7×10^{-5} . If we reduce the dimensions across level and left and right regions, and sort the dimensions in increasing order of significance, we get the following 1. Where we see very coherent results where the Ventral and Intermediate Gray Matter regions are the most significantly correlated and the dorsal Gray Matter



Figure 3. Empirical Functional Connectivity and Data Distribution. A Pearson Correlation Matrix of time points of all patients in wrist rotation. B Pearson Correlation Matrix of time points of all patient resting. C Visualization of Data Distribution per dimension colored with the activity paradigm where blue is rest and orange is active time points.

is the least even though it passes the classical significance $\alpha = 5 \times 10^{-2}$ threshold but not the Bonferroni correction.

ROI	V	CST	Ι	FG	FC	SL	D
P-Value	0.0	$5.2 imes 10^{-5}$	1.02×10^{-4}	$2.8 imes 10^{-4}$	1.01×10^{-2}	$1.2 imes 10^{-2}$	$4.1 imes 10^{-2}$

Table 1. **P-Values of Pearson Correlation of Spinal Cord Regions of Interest with Activity Paradgim**, (V: Ventral, I: Interneurons, D: Dorsal Gray matters)

Even though these results were statistically significant, visually the two distributions of fMRI signals for active and rest states (Fig.4) show the spinal cord data organised as two superposed gaussians.



Figure 4. **Organization of spinal cord data for different states**A. Histogram of BOLD activity of Gray Ventral Matter B. Scatter of Ventral Gray Matter versus White CST

3.2. GLMM and Low Signal

Following these results, I conducted a fine-grained gridsearch of our GLMM's hyper-parameters for the ranges of K = [2,7], $\Theta = [1,55]$ and log range for delta $log(\Delta) = [-3,2]$ as delta is scaling ratio of two other hyper-parameters.



Figure 5. Optimization of GLMM Parameters A. Gridsearch visualization colored with the F1-score of each GLMM run for each setting. B. $\gamma_{m,k}$ or probability estimates of cluster K=1 for the best hyper-parameters $K = 6, \Delta = 0.25, \Theta = 44$.

I fitted a linear regression model onto the results of this grid search or cube in Fig.5 with the hyper-parameters as the predicting or independent variables and the F1-score as the dependent variable. I found significant but coefficients for ever hyper parameter close to 0.

$$F1 \approx \beta_0 + \sum_{i=1}^n \beta_i \cdot h_i \approx \beta_0, \quad \text{where } \beta_i \approx 0 \,\forall i$$

This is not a definitive proof of the unfeasability of our task as our hyper-parameters and F1-score could be nonlinearily explainable, however learning the null plane shows there is no observable linear trend between our clustering working well with respect to the activity paradigm and our hyper-parameters.

Interestingly, if we fit linear models to optimize for the Silhouette Score or how balanced the resulting clusters are (how equally distributed the data points are in each cluster with respect to K), we increase the explainability of our model from $R^2 = 0.02$ to $R^2 = 0.12$ and $R^2 = 0.52$ respectively.

In Figure 6, we see the learned parameters for the best GLMM run with an F1-score of 0.65. Despite this low signal, we see a very coherent adjacency matrix and learned activation means 6.A 6.C. With a strong intra-spinal-level connectivity in adjacency matrix and high activity of the ventral gray ventral matter regions in the means.

Lastly, I also tried K-Graphs (Araghi et al., 2019) algorithm and initializing the GLMM with the Functional Connectivity matrices and observed similar unsatisfactory results and the algorithm drifting away from the solution.



Figure 6. Learned Cluster K = 1 (Active/Wrist Rotation) Statistics for Best GLMM Hyper Parameters A. Learned Adjacency Matrix B. Learned Mean Activations C. Reduced Learned Means

3.3. Synthetic Data and Data Generation

Given the challenges observed in real data, we explore synthetic data modeling to analyze clustering feasibility. Figure 4 shows the spinal cord data between active and rest states being organised a two superposed gaussians. If this were truly the case, the task of learning unsupervised Laplacians and means describing the two states would be trivially impossible. To verify this, I generated synthetic data clusters with two generation process parameters R and θ to simulate an fMRI signal from the spinal cord and observe the GLMM's performance on the task.

I observed a phase transition of the feasability of the GLMM 7 varying generating two synthetic data clusters with ranges of parameters R = [0, 3] and $\theta = [0, 60]$ where theta is in degrees for the orthogonality of the covariance matrices. I ran each setting 5 times and averaged the F1-score of the GLMM solving the synthetic clustering tasks. We see a rather sharp feasability transition in when the GLMM can solve the clustering task and cannot. Furthermore, we unfortunately see that computing the same statistics for our real data, the spinal cord is close to feasability but corresponds to a F1-score 0.57. This matrix gives a feasability look up table for unsupervised clustering tasks.

4. Discussion & Conclusion

If we can model the spinal cord's data distribution between active and rest as a gaussian mixture parametrised with R and θ , we see that we a get a very small R = 0.45 and $\theta = 30$ where independently a minimum R > 2.4 or $\theta > 40$ is needed. The latter conditions intuitively makes sense as the 99% confidence interval of a gaussian distribution is approximately 3σ or 3 standard deviations away and 40 degrees is close to $\frac{\pi}{4}$ in radians. In higher dimensions, the notion of angles between eigenvectors, flats and orthogonality between covariance matrices seems to be less related



Figure 7. Heat Map of two-dimensional GLMM Feasability Transition and Data Generation Parameters of Spinal Cord

to an angle and more to an overlap between distributions. The transformation for one multivariate normal distribution to overlap another is a simple rescaling of the covariance matrix which can also be given by our Givens rotation. I think these results provide convincing arguments to justify the failure of the GLMM on the spinal cord's data.

In conclusion, I think the organization of the spinal chord's functional connectome we dervied from the GLMM of wrist activity 6 is correct despite the weak signal and shows logical results of high intra-level connectivity and activity being dominated by ventral and intermediate gray matter regions. I also think that with synthetic data, we showed the difficulty of this task comes from both the data and the GLMM's modeling capabilities. The assumptions on the data being made that spinal cord activity corresponds to the translation and the rotation of a high dimensional ellipsoid of the resting state.

Our findings suggest that future research should focus on improving data quality of spinal cord acquisitions with a stronger difference in means between states, explore other unsupervised methodologies with different assumptions.

Another synthetic feasability transition that would be interesting to explore, would be the one of data sufficiency. I ran all my synthetic data generations with 4000 sampled points to emulate the spinal cord's data sufficiency, however low data regimes are noisy. Exploring methods of mean and covariance estimation for small data samples and noise reduction method, that are well document perhaps in finance like Ledoit-Wolf's shrinking or more recent approaches may be interesting. As in the spinal cord, we are trying to estimate a gaussian mixture of two cluster's statistics (two means and laplacians) in d = 56 dimensions which is $O(d^2)$ with $3750 > 3136(= 56^2) \sim O(d^2)$ data points, and we would alteast need another order of magnitude to estimate our statistics more precisely.

References

- Araghi, H., Sabbaqi, M., and Babaie–Zadeh, M. K-Graphs: An Algorithm for Graph Signal Clustering and Multiple Graph Learning. *IEEE Signal Processing Letters*, 26(10):1486–1490, October 2019. ISSN 1558-2361. doi: 10.1109/LSP.2019. 2936665. URL https://ieeexplore.ieee.org/document/8809198. Conference Name: IEEE Signal Processing Letters.
- De Leener, B., Lévy, S., Dupont, S. M., Fonov, V. S., Stikov, N., Louis Collins, D., Callot, V., and Cohen-Adad, J. SCT: Spinal Cord Toolbox, an open-source software for processing spinal cord MRI data. *NeuroImage*, 145:24–43, January 2017. ISSN 10538119. doi: 10.1016/j.neuroimage.2016. 10.009. URL https://linkinghub.elsevier. com/retrieve/pii/S1053811916305560.
- Kalofolias, V. How to learn a graph from smooth signals, January 2016. URL http://arxiv.org/abs/ 1601.02513. arXiv:1601.02513 [stat].
- Kalofolias, V. and Perraudin, N. Large Scale Graph Learning from Smooth Signals, May 2019. URL http:// arxiv.org/abs/1710.05654. arXiv:1710.05654 [stat].
- Karahanoğlu, F. I. and Van De Ville, D. Transient brain activity disentangles fMRI resting-state dynamics in terms of spatially and temporally overlapping networks. *Nature Communications*, 6(1):7751, July 2015. ISSN 2041-1723. doi: 10.1038/ncomms8751. URL https://www. nature.com/articles/ncomms8751. Publisher: Nature Publishing Group.
- Kinany, N., Pirondini, E., Mattera, L., Martuzzi, R., Micera, S., and Van De Ville, D. Towards reliable spinal cord fMRI: Assessment of common imaging protocols. *NeuroImage*, 250:118964, April 2022. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2022.118964.
- Maretic, H. P. and Frossard, P. Graph Laplacian mixture model, March 2020. URL http://arxiv.org/ abs/1810.10053. arXiv:1810.10053.
- Preti, M. G., Bolton, T. A., and Van De Ville, D. The dynamic functional connectome: State-of-the-art and perspectives. *NeuroImage*, 160:41–54, October 2017. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2016.12.

061. URL https://www.sciencedirect.com/ science/article/pii/S1053811916307881.

- Ricchi, I., Tarun, A., Maretic, H. P., Frossard, P., and Van De Ville, D. Dynamics of functional network organization through graph mixture learning. *NeuroImage*, 252:119037, May 2022. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2022.119037. URL https://www.sciencedirect.com/ science/article/pii/S1053811922001665.
- Seung, H. S. Predicting visual function by interpreting a neuronal wiring diagram. *Nature*, 634(8032):113– 123, October 2024. ISSN 1476-4687. doi: 10.1038/ s41586-024-07953-5. URL https://www.nature. com/articles/s41586-024-07953-5. Publisher: Nature Publishing Group.
- Yu, M., Sporns, O., and Saykin, A. J. The human connectome in Alzheimer disease — relationship to biomarkers and genetics. *Nature reviews. Neurology*, 17(9):545– 563, September 2021. ISSN 1759-4758. doi: 10.1038/ s41582-021-00529-1. URL https://www.ncbi. nlm.nih.gov/pmc/articles/PMC8403643/.
- Zhang, L., Wang, L., Xia, H., Tan, Y., Li, C., and Fang, C. Connectomic mapping of brain-spinal cord neural networks: Future directions in assessing spinal cord injury at rest. *Neuroscience Research*, 176:9–17, March 2022. ISSN 0168-0102. doi: 10.1016/j.neures.2021.10. 008. URL https://www.sciencedirect.com/ science/article/pii/S0168010221002169.