

Generative Framework

A generative modelling framework for individual brain organization across a number of different data. The Model is partitioned into a model that determines the probability of the spatial arrangement of regions in each subject s , $p(\mathbf{U}^{(s)}; \theta_A)$ and the probability of observing a set of data at each given brain location. We introduce the Markov property that the observations are mutually independent, given the spatial arrangement.

$$p(\mathbf{Y}^{(s)}|\mathbf{U}^{(s)}; \theta_E) = \prod_i p(\mathbf{y}_i^{(s)}|\mathbf{u}_i^{(s)}; \theta_E) \quad (1)$$

Inference and learning

We will learn the model, by maximizing the ELBO (Evidence lower bound). For clarity, I am dropping the index for the subject (s) for now,

$$\begin{aligned} \log p(\mathbf{Y}|\theta) &= \log \sum_{\mathbf{U}} p(\mathbf{Y}, \mathbf{U}|\theta) \\ &= \log \sum_{\mathbf{U}} q(\mathbf{U}) \frac{p(\mathbf{Y}, \mathbf{U}|\theta)}{q(\mathbf{U})} \\ &\geq \sum_{\mathbf{U}} q(\mathbf{U}) \log \frac{p(\mathbf{Y}, \mathbf{U}|\theta)}{q(\mathbf{U})} \quad (\text{Jensen's inequality}) \\ &= \langle \log p(\mathbf{Y}, \mathbf{U}|\theta) - \log q(\mathbf{U}) \rangle_q \triangleq \mathcal{L}(q, \theta) - \log \langle q(\mathbf{U}) \rangle_q \end{aligned}$$

Given the markov property, we can break the expected complete log likelihood into two pieces, one containing the parameters for the arrangement model and one containing the parameters for the emission model.

$$\begin{aligned} \langle \log p(\mathbf{Y}, \mathbf{U}|\theta) \rangle_q &= \langle \log(p(\mathbf{Y}|\mathbf{U}; \theta_E)p(\mathbf{U}|\theta_A)) \rangle_q \\ &= \langle \log p(\mathbf{Y}|\mathbf{U}; \theta_E) \rangle_q + \langle \log p(\mathbf{U}|\theta_A) \rangle_q \\ &\triangleq \mathcal{L}_E + \mathcal{L}_A \end{aligned}$$

We will refer to the first term as the expected emission log-likelihood and the second term as the expected arrangement log-likelihood. We can estimate the parameters of the emission model from maximizing the expected emission log-likelihood, and we can estimate the parameters of the arrangement model by maximizing the expected arrangement log-likelihood.

Arrangement models

This is a generative Potts model of brain activity data. The main idea is that the brain consists of K regions, each with a specific activity profile \mathbf{v}_k for a specific task set. The model consists of a arrangement model that tells us how the K regions are arranged in a specific subject s , and an emission model that provides a probability of the measured data, given the individual arrangement of regions.

Independent Arrangement model

This is the simplest spatial arrangement model - it simply learns the probability at each location that the node is part of cluster K . These probabilities are simply learned as the parameters $\pi_{ik} = p(u_i = k)$, or after a re-parameterization in log space: $\eta_{ik} = \log \pi_{ik}$. Vice versa (not assuming that the etas are correctly scaled):

$$\pi_{ik} = \frac{\exp(\eta_{ik})}{\sum_j \exp(\eta_{ij})} \quad (2)$$

This independent arrangement model can be estimated using the EM-algorithm.

In the Estep, we are integrating the evidence from the data and the prior:

$$p(u_i = k|\mathbf{y}_i) = \langle u_{ik} \rangle = \frac{\exp(\log(p(\mathbf{y}_i|\mathbf{u}_i = k)) + \eta_{ik})}{\sum_j \exp(\log(p(\mathbf{y}_i|\mathbf{u}_i = j)) + \eta_{ij})} \quad (3)$$

Or in vector notation:

$$\langle \mathbf{u}_i \rangle = \text{softmax}(\log(p(\mathbf{y}_i | \mathbf{u}_i)) + \boldsymbol{\eta}_i)$$

For the M-step, we use the derivative of the expected arrangement log-likelihood in respect to the parameters $\boldsymbol{\eta}$:

$$\begin{aligned} \mathcal{L}_A &= \sum_i \sum_k \langle u_{i,k} \rangle \log(\pi_{ik}) \\ &= \sum_i \sum_k \langle u_{ik} \rangle (\eta_{ik} - \log \sum_j \exp(\eta_{ij})) \\ &= \sum_i \sum_k \langle u_{ik} \rangle \eta_{ik} - \sum_i \log \sum_j \exp(\eta_{i,j}) \end{aligned} \quad (4)$$

So the derivative is

$$\frac{\partial \mathcal{L}_A}{\partial \eta_{ik}} = \langle u_{ik} \rangle - \frac{\partial}{\partial \eta_{ik}} \log \sum_j \exp(\eta_{ij}) \quad (5)$$

$$= \langle u_{ik} \rangle - \frac{1}{\sum_j \exp(\eta_{ij})} \frac{\partial}{\partial \eta_{ik}} \sum_j \exp(\eta_{ij}) \quad (6)$$

$$= \langle u_{ik} \rangle - \frac{\exp(\eta_{ik})}{\sum_j \exp(\eta_{ij})} \quad (7)$$

$$= \langle u_{ik} \rangle - \pi_{ik} \quad (8)$$

We can also get the same result directly by the application of chain rule: For a good introduction, see: [<https://eli.thegreenplace.net/2016/the-softmax-function-and-its-derivative/>].

$$\frac{\partial \pi_k}{\partial \eta_k} = \pi_k (\delta_{kj} - \pi_j) \quad (9)$$

Simple Potts model

The brain is sampled in P vertices (or voxels). Individual maps are aligned using anatomical normalization, such that each vertex refers to a (roughly) corresponding region in each individual brain. The assignment of each brain location to a specific parcel in subject s is expressed as the random variable $u_i^{(s)}$.

Across individual brains, we have the overall probability of a specific brain location being part of parcel k .

$$p(u_i = k) = \pi_{ki} \quad (10)$$

The spatial interdependence of brain locations is expressed as a Potts model. In this model, the overall probability of a specific assignment of brain locations to parcels (the vector \mathbf{u}) is expressed as the product of the overall prior and the product of all possible pairwise potentials (ψ_{ij}).

$$p(\mathbf{u}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_i \pi_{u_i, i} \prod_{i \neq j} \psi_{ij}(u_i, u_j) \quad (11)$$

Each local potential is defined by an exponential over all other that are connected to node i , i.e. nodes with connectivity weights of $w_{ji} = w_{ij} > 0$.

$$\psi_{ij} = \exp(\boldsymbol{\theta}_w \mathbf{u}_i^T \mathbf{u}_j w_{ij}) \quad (12)$$

Where we have introduced a one-hot encoding of u_i with a K vector of indicator variables \mathbf{u}_i , such that $\mathbf{u}_i^T \mathbf{u}_j = 1$ if $u_i = u_j$ and 0 otherwise.

The spatial co-dependence across the entire brain is therefore expressed with the pairwise weights w that encode how likely two nodes belong to the same parcel. The temperature parameter $\boldsymbol{\theta}_w$ determines how strong this co-dependence overall influences the local probabilities (relative to the prior). We can use this notation to express local co-dependencies by using a graph, where we define

$$w_{ij} = \begin{cases} 1; & \text{if } i \text{ and } j \text{ are neighbours} \\ 0; & \text{otherwise} \end{cases} \quad (13)$$

This formulation would enforce local smoothness of the map. However, we could also express in these potential more medium range potentials (two specific parietal and premotor areas likely belong to the same parcel), as well as cross-hemispheric symmetry. Given this, the matrix \mathbf{W} could be simply derived from the underlying grid or be learned to reflect known brain-connectivity.

The expected arrangement log-likelihood therefore becomes:

$$\mathcal{L}_A = \sum_i \langle \mathbf{u}_i \rangle^T \log \boldsymbol{\pi}_i + \theta_w \sum_i \sum_j w_{ij} \langle \mathbf{u}_i^T \mathbf{u}_j \rangle - \log Z$$

Inference using stochastic maximum likelihood / contrastive divergence

We can approximate the gradient of the parameters using a contrastive divergence-type algorithm. We view the arrangement log likelihood as a sum of the unnormalized part ($\tilde{\mathcal{L}}_A$) and the log partition function. For each parameter θ we then follow the gradient

$$\begin{aligned} \nabla_{\theta} \mathcal{L}_A &= \nabla_{\theta} \tilde{\mathcal{L}}_A - \nabla_{\theta} \log Z \\ &= \nabla_{\theta} \tilde{\mathcal{L}}_A - \mathbb{E}_p[\nabla_{\theta} \tilde{\mathcal{L}}_A] \\ &= \nabla_{\theta} \langle \log \tilde{p}(\mathbf{U}|\theta) \rangle_q - \nabla_{\theta} \langle \log \tilde{p}(\mathbf{U}|\theta) \rangle_p \end{aligned}$$

Thus, we can use the gradient of the unnormalized expected log-likelihood (given a distribution $q(\mathbf{U}) = p(\mathbf{U}|\mathbf{Y}; \theta)$), minus the gradient of the unnormalized expected log-likelihood in respect to the expectation under the model parameters, without seeing the data, $q(\mathbf{U}) = p(\mathbf{U}|\mathbf{Y}; \theta)$. This motivates the use of sampling / approximate inference for both of these steps. See Deep Learning (18.1).

Sampling from the prior or posterior distribution

The problem is that the two expectations under the prior (p) and the posterior (q) distribution of the model cannot be easily be computed. We can evaluate the prior probability of a parcellation $p(\mathbf{U})$ or the posterior distribution $p(\mathbf{U}|\mathbf{Y})$ up to a constant of proportionality, with for example

$$p(\mathbf{U}|\mathbf{Y}; \theta) = \frac{1}{Z(\theta)} \prod_i \mu_{u_i, i} \prod_{i \neq j} \psi_{ij}(u_i, u_j) \prod_i p(\mathbf{y}_i | u_i) \quad (14)$$

Calculating the normalization constant $Z(\theta)$ (partition function, Zustandssumme, or sum over states) would involve summing this probability over all possible states, which for P brain locations and K parcels is K^P , which is intractable.

However, the conditional probability for each node, given all the other nodes, can be easily computed. Here the normalization constant is just the sum of the potential functions over the K possible states for this node

$$p(u_i | u_{j \neq i}, \mathbf{y}_i; \theta) = \frac{1}{Z(\theta)} \mu_{u_i, i} p(\mathbf{y}_i | u_i) \prod_{i \neq j} \psi_{ij}(u_i, u_j) \quad (15)$$

With Gibbs sampling, we start with a pattern $\mathbf{u}^{(0)}$ and then update $u_1^{(1)}$ by sampling from $p(u_1 | u_2^{(0)} \dots u_P^{(0)})$. We then sample $u_2^{(1)}$ by sampling from $p(u_2 | u_1^{(1)}, u_3^{(0)} \dots u_P^{(0)})$ and so on, until we have sampled each node once. Then we return to the beginning and restart the process. After some burn-in period, the samples will come from desired overall distribution. If we want to sample from the prior, rather than from the posterior, we simply drop the $p(\mathbf{y}_i | u_i)$ term from the conditional probability above.

Gradient for different parametrization of the Potts model

For the edge-energy parameters θ_w we clearly want to use the natural parametrization with the derivate:

$$\frac{\partial \tilde{\mathcal{L}}_A}{\partial \theta_w} = \sum_i \sum_j w_{ij} \langle \mathbf{u}_i^T \mathbf{u}_j \rangle \quad (16)$$

For the prior probability of each parcel k at each location i (π_{ik}) we have a number of options.

First, we can use the probabilities themselves as parameters:

$$\frac{\partial \tilde{\mathcal{L}}_A}{\partial \pi_{ik}} = \frac{\langle u_{ik} \rangle}{\pi_{ik}} \quad (17)$$

This is unconstrained (that is probabilities do not need to sum to 1), and the normalization would happen through the partition function.

Secondly, we can use a re-parameterization in log space, which is more natural: $\eta_{ik} = \log \pi_{ik}$. In this case the derivative of the non-normalized part just becomes:

$$\frac{\partial \tilde{\mathcal{L}}_A}{\partial \eta_{ik}} = \langle u_{ik} \rangle \quad (18)$$

Finally, we can implement the constraint that the probabilities at each location sum to one by the following re-parametrization:

$$\begin{aligned} \pi_{iK} &= 1 - \sum_{k=1}^{K-1} \pi_{ik} \\ \eta_{ik} &= \log\left(\frac{\pi_{ik}}{\pi_{iK}}\right) = \log \pi_{ik} - \log\left(1 - \sum_{k=1}^{K-1} \pi_{ik}\right) \\ \pi_{ik} &= \frac{\exp(\eta_{ik})}{1 + \sum_{k=1}^{K-1} \exp(\eta_{ik})} \\ \pi_{iK} &= \frac{1}{1 + \sum_{k=1}^{K-1} \exp(\eta_{ik})} \end{aligned}$$

In the implementation, we can achieve this parametrization easily by defining a non-flexible parameter $\eta_{iK} \triangleq 0$. Then we can treat the last probability like all the other ones.

With this constrained parameterization, we can rewrite the unnormalized part of the expected log-likelihood as:

$$\begin{aligned} \tilde{\mathcal{L}}_A &= \sum_i \sum_k^{K-1} \langle u_{ik} \rangle \log \pi_{ik} + [1 - \sum_k^{K-1} \langle u_{ik} \rangle] \log \pi_{iK} + C \\ &= \sum_i \sum_k^{K-1} \langle u_{ik} \rangle (\log \pi_{ik} - \log \pi_{iK}) + \log \pi_{iK} + C \\ &= \sum_i \sum_k^{K-1} \langle u_{ik} \rangle \eta_{ik} - \log\left(1 + \sum_{k=1}^{K-1} \exp(\eta_{ik})\right) + C \end{aligned}$$

where C is the part of the normalized log-likelihood that does not depend on π . Taking derivative in respect to η_{ik} yields:

$$\begin{aligned} \frac{\partial \tilde{\mathcal{L}}_A}{\partial \eta_{ik}} &= \langle u_{ik} \rangle - \frac{\exp(\eta_{ik})}{1 + \sum_k^{K-1} \exp(\eta_{ik})} \\ &= \langle u_{ik} \rangle - \pi_{ik} \end{aligned}$$

So in this parameterization in the iid case, $Z = 1$ and we don't need the negative step. In general, however, we cannot simply set the above derivative to zero and solve it, as the parameter θ_w will also have an influences on $\langle u_{ik} \rangle$.

Emission models

Given the Markov property, the emission models specify the log probability of the observed data as a function of \mathbf{u} .

$$\log p(\mathbf{Y}|\mathbf{U}; \theta_E) = \sum_i \log p(\mathbf{y}_i|\mathbf{u}_i; \theta_E) \quad (19)$$

Furthermore, assuming that \mathbf{u}_i is a one-hot encoded indicator variable (parcellation model), we can write the expected emission log-likelihood as:

$$\langle \log p(\mathbf{Y}|\mathbf{U}; \theta_E) \rangle = \sum_i \sum_k \langle u_i^{(k)} \rangle \log p(\mathbf{y}_i|u_i = k; \theta_E) \quad (20)$$

In the E-step the emission model simply passes $p(\mathbf{y}_i|\mathbf{u}_i; \theta_E)$ as a message to the arrangement model. In the M-step, $q(\mathbf{u}_i) = \langle \mathbf{u}_i \rangle$ is passed back, and the emission model optimizes the above quantity in respect to θ_E .

Emission model 1: Mixture of Gaussians

Under the Gaussian mixture model, we model the emissions as a Gaussian with a parcel-specific mean (\mathbf{v}_k), and with equal isotropic variance across parcels and observations:

$$p(\mathbf{y}_i | u^{(k)}; \theta_E) = \frac{1}{(2\pi)^{N/2} (\sigma^2)^{N/2}} \exp\left\{-\frac{1}{2\sigma^2} (\mathbf{y}_i - \mathbf{v}_k)^T (\mathbf{y}_i - \mathbf{v}_k)\right\} \quad (21)$$

The expected emission log-likelihood therefore is:

$$\begin{aligned} \mathcal{L}_E &= \sum_i \sum_k \langle u_i^{(k)} \rangle_q \left[-\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y}_i - \mathbf{v}_k)^T (\mathbf{y}_i - \mathbf{v}_k) \right] \\ &= -\frac{PN}{2} \log(2\pi) - \frac{PN}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_i \sum_k \langle u_i^{(k)} \rangle_q [(\mathbf{y}_i - \mathbf{v}_k)^T (\mathbf{y}_i - \mathbf{v}_k)] \end{aligned}$$

Now, with the above expected emission log likelihood by hand, we can update the parameters $\theta_E = \{\mathbf{v}_1, \dots, \mathbf{v}_K, \sigma^2\}$ in the **M** step.

1. Updating \mathbf{v}_k , we take derivative of *expected emission log likelihood* with respect to \mathbf{v}_k and set it to 0:

$$\frac{\partial \mathcal{L}_E}{\partial \mathbf{v}_k} = \frac{1}{\sigma^2} \sum_i \langle u_i^{(k)} \rangle_q (\mathbf{y}_i - \mathbf{v}_k) = 0 \quad (22)$$

Thus, we get the updated \mathbf{v}_k in current **M** step as,

$$\mathbf{v}_k^{(t)} = \frac{\sum_i \langle u_i^{(k)} \rangle_q^{(t)} \mathbf{y}_i}{\sum_i \langle u_i^{(k)} \rangle_q^{(t)}} \quad (23)$$

2. Updating σ^2 , we take derivative of *expected emission log likelihood* $\mathcal{L}(q, \theta)$ with respect to σ^2 and set it to 0:

$$\frac{\partial \mathcal{L}_E}{\partial \sigma^2} = -\frac{PN}{2\sigma^2} + \sum_i \sum_k \langle u_i^{(k)} \rangle_q \left[\frac{1}{2\sigma^4} (\mathbf{y}_i - \mathbf{v}_k^{(t)})^T (\mathbf{y}_i - \mathbf{v}_k^{(t)}) \right] = 0 \quad (24)$$

Thus, we get the updated σ^2 for parcel k in the current **M** step as,

$$\sigma^{2(t)} = \frac{1}{PN} \sum_i \sum_k \langle u_i^{(k)} \rangle_q^{(t)} (\mathbf{y}_i - \mathbf{v}_k^{(t)})^T (\mathbf{y}_i - \mathbf{v}_k^{(t)}) \quad (25)$$

where P is the total number of voxels i .

The updated parameters $\theta_k^{(t)}$ from current **M**-step will be passed to the next **E**-step ($t + 1$) until convergence.

Emission model 2: Mixture of Gaussians with Exponential signal strength

The emission model should depend on the type of data that is measured. A common application is that the data measured at location i are the task activation in N tasks, arranged in the $N \times 1$ data vector \mathbf{y}_i . The averaged expected response for each of the parcels is \mathbf{v}_k . One issue of the functional activation is that the signal-to-noise ratio (SNR) can be quite different across different participants, and voxels, with many voxels having relatively low SNR. We model this signal to noise for each brain location (and subject) as $s_i \sim \text{exponential}(\beta_s)$. Therefore the probability model for gamma is defined as:

$$p(s_i | \theta) = \beta e^{-\beta s_i} \quad (26)$$

Overall, the expected signal at each brain location is then

$$\mathbf{E}(\mathbf{y}_i) = \mathbf{u}_i^T \mathbf{V} \mathbf{s}_i \quad (27)$$

Finally, relative to the signal, we assume that the noise is distributed i.i.d Gaussian with:

$$\boldsymbol{\epsilon}_i \sim \text{Normal}(0, \mathbf{I}_K \theta_{\sigma s}) \quad (28)$$

Here, the proposal distribution $q(u_i^{(k)}, s_i | \mathbf{y}_i)$ is now a multivariate distribution across u_i and s_i . Thus, the *expected emission log likelihood* $\mathcal{L}_E(q, \theta)$ is defined as:

$$\begin{aligned}
\mathcal{L}_E &= \left\langle \sum_i \log p(\mathbf{y}_i, s_i | u_i; \theta_E) \right\rangle_q \\
&= \sum_i \sum_k \langle u_i^{(k)} [-\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y}_i - \mathbf{v}_k s_i)^T (\mathbf{y}_i - \mathbf{v}_k s_i)] \rangle_q \\
&+ \sum_i \sum_k \langle u_i^{(k)} [\log \beta - \beta s_i] \rangle_q \\
&= -\frac{NP}{2} \log(2\pi) - \frac{NP}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_i \sum_k \langle u_i^{(k)} (\mathbf{y}_i - \mathbf{v}_k s_i)^T (\mathbf{y}_i - \mathbf{v}_k s_i) \rangle_q \\
&+ P \log \beta - \sum_i \sum_k \beta \langle u_i^{(k)} s_i \rangle_q \\
&= -\frac{NP}{2} \log(2\pi) - \frac{NP}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_i \mathbf{y}_i^T \mathbf{y}_i - \frac{1}{2\sigma^2} \sum_i \sum_k (-2\mathbf{y}_i^T \mathbf{v}_k \langle u_i^{(k)} s_i \rangle_q + \mathbf{v}_k^T \mathbf{v}_k \langle u_i^{(k)} s_i^2 \rangle_q) \\
&+ \log \beta - \beta \sum_i \sum_k \langle u_i^{(k)} s_i \rangle_q
\end{aligned}$$

Now, we can update the parameters θ of the Gaussians/Exponential mixture in the \mathbf{M} step. The parameters of the gaussian mixture model are $\theta_E = \{\mathbf{v}_1, \dots, \sigma^2, \beta\}$.

1. We start with updating the \mathbf{v}_k (Note: the updates only consider a single subject). We take the derivative of *expected emission log likelihood* \mathcal{L}_E with respect to \mathbf{v}_k and make it equals to 0 as following:

$$\frac{\partial \mathcal{L}_E}{\partial \mathbf{v}_k} = -\frac{1}{\sigma^2} \sum_i -\mathbf{y}_i^T \langle u_i^{(k)} s_i \rangle_q + \mathbf{v}_k^T \langle u_i^{(k)} s_i^2 \rangle_q = 0 \quad (29)$$

Thus, we get the updated \mathbf{v}_k in current \mathbf{M} step as,

$$\mathbf{v}_k^{(t)} = \frac{\sum_i \langle u_i^{(k)} s_i \rangle_q \mathbf{y}_i}{\sum_i \langle u_i^{(k)} s_i^2 \rangle_q} \quad (30)$$

2. Updating σ^2 , we take derivative of with respect to σ^2 and set it equals to 0 as following:

$$\frac{\partial \mathcal{L}_E}{\partial \sigma^2} = -\frac{NP}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_i \mathbf{y}_i^T \mathbf{y}_i + \frac{1}{2\sigma^4} \sum_i \sum_k (-2\mathbf{y}_i^T \mathbf{v}_k \langle u_i^{(k)} s_i \rangle_q + \mathbf{v}_k^T \mathbf{v}_k \langle u_i^{(k)} s_i^2 \rangle_q) = 0 \quad (31)$$

Thus, we get the updated σ^2 for parcel k in the current \mathbf{M} step as,

$$\sigma^{2(t)} = \frac{1}{NP} \left(\sum_i \mathbf{y}_i^T \mathbf{y}_i + \sum_i \sum_k (-2\mathbf{y}_i^T \mathbf{v}_k \langle u_i^{(k)} s_i \rangle_q + \mathbf{v}_k^T \mathbf{v}_k \langle u_i^{(k)} s_i^2 \rangle_q) \right) \quad (32)$$

3. Updating β , we take derivative of $\mathcal{L}_E(q, \theta)$ with respect to β and set it equal to 0 as following:

$$\begin{aligned}
\frac{\partial \mathcal{L}_E}{\partial \beta} &= \frac{\partial [P \log \beta - \beta \sum_i \sum_k \langle u_i^{(k)} s_i \rangle_q]}{\partial \beta} \\
&= \frac{P}{\beta} - \sum_i \sum_k \langle u_i^{(k)} s_i \rangle_q = 0
\end{aligned}$$

Thus, we get the updated β_k in current \mathbf{M} step as,

$$\beta_k^{(t)} = \frac{P}{\sum_i \sum_k \langle u_i^{(k)} s_i \rangle_q} \quad (33)$$

Emission model 2b: Mixture of Gaussians with Gamma signal strength

The emission model should depend on the type of data that is measured. A common application is that the data measured at location i are the task activation in N tasks, arranged in the $N \times 1$ data vector \mathbf{y}_i . The averaged expected response for each of the parcels is \mathbf{v}_k . One issue of the functional activation is that the signal-to-noise ratio (SNR) can be quite different across different participants, and voxels, with many voxels

having relatively low SNR. We model this signal to noise for each brain location (and subject) as $s_i \sim \text{Gamma}(\theta_\alpha, \theta_{\beta s})$. Therefore the probability model for gamma is defined as:

$$p(s_i|\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} s_i^{\alpha-1} e^{-\beta s_i} \quad (34)$$

Overall, the expected signal at each brain location is then

$$\mathbf{E}(\mathbf{y}_i) = \mathbf{u}_i^T \mathbf{V} \mathbf{s}_i \quad (35)$$

Finally, relative to the signal, we assume that the noise is distributed i.i.d Gaussian with:

$$\boldsymbol{\epsilon}_i \sim \text{Normal}(0, \mathbf{I}_K \theta_{\sigma s}) \quad (36)$$

Here, the proposal distribution $q(u_i^{(k)}, s_i|\mathbf{y}_i)$ is now a multivariate distribution across u_i and s_i . Thus, the *expected emission log likelihood* $\mathcal{L}_E(q, \theta)$ is defined as:

$$\begin{aligned} \mathcal{L}_E &= \left\langle \sum_i \log p(\mathbf{y}_i, s_i | u_i; \theta_E) \right\rangle_q \\ &= \sum_i \sum_k \langle u_i^{(k)} [-\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y}_i - \mathbf{v}_k s_i)^T (\mathbf{y}_i - \mathbf{v}_k s_i)] \rangle_q \\ &\quad + \sum_i \sum_k \langle u_i^{(k)} [\alpha \log \beta - \log \Gamma(\alpha) + (\alpha - 1) \log s_i - \beta s_i] \rangle_q \\ &= -\frac{NP}{2} \log(2\pi) - \frac{NP}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_i \sum_k \langle u_i^{(k)} (\mathbf{y}_i - \mathbf{v}_k s_i)^T (\mathbf{y}_i - \mathbf{v}_k s_i) \rangle_q \\ &\quad + P\alpha \log \beta - P \log \Gamma(\alpha) + \sum_i \sum_k \langle u_i^{(k)} (\alpha - 1) \log s_i - u_i^{(k)} \beta s_i \rangle_q \\ &= -\frac{NP}{2} \log(2\pi) - \frac{NP}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_i \mathbf{y}_i^T \mathbf{y}_i - \frac{1}{2\sigma^2} \sum_i \sum_k (-2\mathbf{y}_i^T \mathbf{v}_k \langle u_i^{(k)} s_i \rangle_q + \mathbf{v}_k^T \mathbf{v}_k \langle u_i^{(k)} s_i^2 \rangle_q) \\ &\quad + P\alpha \log \beta - P \log \Gamma(\alpha) + (\alpha - 1) \sum_i \sum_k \langle u_i^{(k)} \log s_i \rangle_q - \beta \sum_i \sum_k \langle u_i^{(k)} s_i \rangle_q \end{aligned}$$

Now, we can update the parameters θ of the Gaussians/Gamma mixture in the M step. The parameters of the gaussian mixture model are $\theta_E = \{\mathbf{v}_1, \dots, \sigma^2, \alpha, \beta\}$.

1. We start with updating the \mathbf{v}_k (Note: the updates only consider a single subject). We take the derivative of *expected emission log likelihood* \mathcal{L}_E with respect to \mathbf{v}_k and make it equals to 0 as following:

$$\frac{\partial \mathcal{L}_E}{\partial \mathbf{v}_k} = -\frac{1}{\sigma^2} \sum_i -\mathbf{y}_i^T \langle u_i^{(k)} s_i \rangle_q + \mathbf{v}_k^T \langle u_i^{(k)} s_i^2 \rangle_q = 0 \quad (37)$$

Thus, we get the updated \mathbf{v}_k in current M step as,

$$\mathbf{v}_k^{(t)} = \frac{\sum_i \langle u_i^{(k)} s_i \rangle_q \mathbf{y}_i}{\sum_i \langle u_i^{(k)} s_i^2 \rangle_q} \quad (38)$$

2. Updating σ^2 , we take derivative of with respect to σ^2 and set it equals to 0 as following:

$$\frac{\partial \mathcal{L}_E}{\partial \sigma^2} = -\frac{NP}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_i \mathbf{y}_i^T \mathbf{y}_i + \frac{1}{2\sigma^4} \sum_i \sum_k (-2\mathbf{y}_i^T \mathbf{v}_k \langle u_i^{(k)} s_i \rangle_q + \mathbf{v}_k^T \mathbf{v}_k \langle u_i^{(k)} s_i^2 \rangle_q) = 0 \quad (39)$$

Thus, we get the updated σ^2 for parcel k in the current M step as,

$$\sigma^{2(t)} = \frac{1}{NP} (\sum_i \mathbf{y}_i^T \mathbf{y}_i + \sum_i \sum_k (-2\mathbf{y}_i^T \mathbf{v}_k \langle u_i^{(k)} s_i \rangle_q + \mathbf{v}_k^T \mathbf{v}_k \langle u_i^{(k)} s_i^2 \rangle_q)) \quad (40)$$

3. Updating β , we take derivative of $\mathcal{L}_E(q, \theta)$ with respect to β and set it equal to 0 as following:

$$\begin{aligned} \frac{\partial \mathcal{L}_E}{\partial \beta} &= \frac{\partial [P\alpha \log \beta - \beta \sum_i \sum_k \langle u_i^{(k)} s_i \rangle_q]}{\partial \beta} \\ &= \frac{P\alpha}{\beta} - \sum_i \sum_k \langle u_i^{(k)} s_i \rangle_q = 0 \end{aligned}$$

Thus, we get the updated β_k in current M step as,

$$\beta_k^{(t)} = \frac{P\alpha_k^{(t)}}{\sum_i \sum_k \langle u_i^{(k)} s_i \rangle_q} \quad (41)$$

4. Updating α_k is comparatively hard since we cannot derive closed-form, we take derivative of $\mathcal{L}_E(q, \theta)$ with respect to α and make it equals to 0 as following:

$$\begin{aligned} \frac{\partial \mathcal{L}_E}{\partial \alpha} &= \frac{\partial [P\alpha \log \beta - P \log \Gamma(\alpha) + (\alpha - 1) \sum_i \sum_k \langle u_i^{(k)} \log s_i \rangle_q]}{\partial \alpha} \\ &= P \log \beta - P \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} + \sum_i \sum_k \langle u_i^{(k)} \log s_i \rangle_q = 0 \end{aligned}$$

The term $\frac{\Gamma'(\alpha)}{\Gamma(\alpha)}$ in above equation is exactly the *digamma function* and we use $F(\alpha)$ to represent. Also from (4), we know $\beta = \frac{P\alpha_k^{(t)}}{\sum_i \sum_k \langle u_i^{(k)} s_i \rangle_q}$ Thus, we get the updated α in current M step as,

$$\begin{aligned} F(\alpha)^{(t)} &= \log \frac{P\alpha_k^{(t)}}{\sum_i \sum_k \langle u_i^{(k)} s_i \rangle_q} + \frac{1}{P} \sum_i \sum_k \log \langle u_i^{(k)} s_i \rangle_q \\ &= \log P\alpha^{(t)} - \log \sum_i \sum_k \langle u_i^{(k)} s_i \rangle_q + \frac{1}{P} \sum_i \sum_k \log \langle u_i^{(k)} s_i \rangle_q \end{aligned}$$

By applying "generalized Newton" approximation form, the updated α is as follows:

$$\alpha^{(t)} \approx \frac{0.5}{\log \sum_i \sum_k \langle u_i^{(k)} s_i \rangle_q - \frac{1}{P} \sum_i \sum_k \langle u_i^{(k)} \log s_i \rangle_q} \quad (42)$$

Note that $\log \sum_i \sum_k \langle u_i^{(k)} s_i \rangle_q \geq \frac{1}{P} \sum_i \sum_k \log \langle u_i^{(k)} s_i \rangle_q$ is given by Jensen's inequality.

Emission model 3: Mixture of Von-Mises Distributions

For a N -dimensional data $\mathbf{y} \in \mathbb{R}^N$ the probability density function of von Mises-Fisher distribution is defined as following,

$$V_N(\mathbf{y} | \mathbf{v}_k, \kappa) = C_N(\kappa) \exp(\kappa \mathbf{v}_k^T \mathbf{y}) \quad (43)$$

where \mathbf{v}_k denotes the mean direction (unit vectors for each parcels), κ indicates the concentration parameter ($\kappa \geq 0$), which is joint over all parcels. $C_N(\kappa) = \frac{\kappa^{N/2-1}}{(2\pi)^{N/2} I_{N/2-1}(\kappa)}$ is the normalization constant where $I_r(\cdot)$ refers to the modified Bessel function of the r order. Thus, the *expected emission log-likelihood* of a mixture of K -classes von-Mises fisher distributions is defined as:

$$\begin{aligned} \mathcal{L}_E &= \langle \sum_i \log p(\mathbf{y}_i | \mathbf{u}_i; \theta_E) \rangle_q \\ &= P \log C_N(\kappa) + \sum_i \sum_k \langle u_i^{(k)} \rangle_q \kappa \mu^{(k)T} \mathbf{y}_i \end{aligned}$$

Now, we update the parameters θ of the von-Mises mixture in the M step by maximizing \mathcal{L}_E in respect to the parameters in von-Mises mixture $\theta_k = \{\mathbf{v}_k, \kappa\}$. (Note: the updates only consider a single subject).

1. Updating mean direction \mathbf{v}_k , we take derivative in respect to \mathbf{v}_k and set it to 0. Thus, we get the updated \mathbf{v}_k in current M step as,

$$\mathbf{v}_k^{(t)} = \frac{\bar{\mathbf{y}}_k}{r_k}, \quad \text{where } \bar{\mathbf{y}}_k = \sum_i \langle u_i^{(k)} \rangle_q \mathbf{y}_i; \quad r_k = \|\bar{\mathbf{y}}_k\|$$

2. Updating concentration parameter κ is difficult in particularly for high dimensional problems since it involves inverting ratio of two Bessel functions. Here we use approximate solutions suggested in (Banerjee et al., 2005):

$$\begin{aligned} \kappa_k^{(t)} &\approx \frac{\bar{r}_k N - \bar{r}_k^3}{1 - \bar{r}_k^2} \\ \bar{r}_k &= \frac{r_k}{\dots} \end{aligned} \quad (44)$$

$$'' \quad \sum_i \langle u_i^{(k)} \rangle_q$$

The updated parameters from current **M**-step will be passed to the **E**-step of $(t + 1)$ times for calculating the expectation.