

---

# Categories and Functional Units: An Infinite Hierarchical Model for Brain Activations

## Supplementary Materials

---

**Danial Lashkari      Ramesh Sridharan      Polina Golland**  
Computer Science and Artificial Intelligence Laboratory  
Massachusetts Institute of Technology  
Cambridge, MA 02139  
`{danial, rameshvs, polina}@csail.mit.edu`

In the Supplementary Materials, we derive the Gibbs free energy cost function for variational inference, and derive the update rules for inference using our variational approximation.

### 1 Joint Probability Distribution

Based on the generative model described in the paper, we can write the full joint distribution of all the observed and unobserved variables. For each variable, we use  $\omega^\cdot$  to denote the natural parameters of the distribution for that variable. For example, the variable  $e_{jih}$  is associated with natural parameters  $\omega_{jh}^{e,1}$  and  $\omega_{jh}^{e,2}$ .

#### 1.1 fMRI model

The fMRI response across different sessions within the same subject can be highly variable. Therefore, we model the variables  $a$ , the activation height, and  $\lambda$ , the voxel-level noise parameter, individually for each of the sessions in the experiment (due to space constraints in the original paper, we omitted these implementation details there). We use the index  $m$  to denote session  $m$ , giving  $a_{jm}$  and  $\lambda_{jm}$  for  $m \in \{1, \dots, M\}$ . We are also provided with binary indicators  $D_{jmt}$  for each subject  $j$ , which specifies whether time points  $t$  belonged to session  $m$ . For notational convenience, we let  $\lambda_{jit} = \prod_m \lambda_{jm}^{D_{jmt}}$ , and similarly let  $a_{jit} = \prod_m a_{jm}^{D_{jmt}}$ . Given the fMRI model parameters, we can then write the likelihood of the observed data  $y_{jit}$ :

$$p(\mathbf{y}|\mathbf{x}, \mathbf{a}, \boldsymbol{\lambda}, e) = \prod_{j,i,t} \sqrt{\frac{\lambda_{jit}}{2\pi}} \exp \left\{ -\frac{\lambda_{jit}}{2} \left( y_{jit} - \sum_h e_{jih} F_{ht} - a_{jit} \sum_s G_{jst} x_{jis} \right)^2 \right\} \quad (1)$$

For the nuisance parameters  $e$ , we have

$$p(e_{jih}) = \text{Normal}(\mu_{jh}^e, \sigma_{jh}^e) \quad (2)$$

$$\propto \exp \left\{ -\omega_{jh}^{e,2} e_{jih}^2 + \omega_{jh}^{e,1} e_{jih} \right\}, \quad (3)$$

where  $\omega_{jh}^{e,2} = \frac{1}{2}(\sigma_{jh}^e)^{-1}$  and  $\omega_{jh}^{e,1} = \frac{1}{2}\mu_h^e(\sigma_{jh}^e)^{-1}$ .

For the noise precision variables  $\boldsymbol{\lambda}$ , we have:

$$p(\lambda_{jm}) = \text{Gamma}(\kappa_{jm}, \theta_{jm}) \quad (4)$$

$$\propto \exp \left\{ \omega_{jm}^{\lambda,1} \log(\lambda_{jm}) - \omega_{jm}^{\lambda,2} \lambda_{jm} \right\}, \quad (5)$$

where  $\omega_{jm}^{\lambda,1} = \kappa_{jm} - 1$  and  $\omega_{jm}^{\lambda,2} = \theta_{jm}$ .

The distribution over the activation heights  $a$  is given by

$$p(a_{jm}) = \text{Normal}_+(\mu_{jm}^a, \sigma_{jm}^a) \quad (6)$$

$$\propto \exp \left\{ -\omega_{jm}^{a,2} a_{jm}^2 + \omega_{jm}^{a,1} a_{jm} \right\}, a_{jm} \geq 0 \quad (7)$$

We have that  $\omega_{jm}^{a,2} = \frac{1}{2}(\sigma_{jm}^a)^{-1}$  and  $\omega_{jm}^{a,1} = \frac{1}{2}\mu_{jm}^a (\sigma_{jm}^a)^{-1}$ .

Note that the distribution  $w \sim \text{Normal}_+(\rho, \eta^{-1})$  is a member of the exponential family of distributions and has the following properties:

$$p(w) = \sqrt{\frac{2\lambda}{\pi}} \left[ 1 + \text{erf} \left( \sqrt{\frac{\rho}{2}} \eta \right) \right]^{-1} e^{-\rho(w-\eta)^2/2}, \quad (8)$$

$$E[w] = \eta + \sqrt{\frac{2}{\pi\lambda}} \left[ 1 + \text{erf} \left( \sqrt{\frac{\rho}{2}} \eta \right) \right]^{-1} e^{-\rho\eta^2/2}, \quad (9)$$

$$E[w^2] = \eta^2 + \rho^{-1} + \eta \sqrt{\frac{2}{\pi\lambda}} \left[ 1 + \text{erf} \left( \sqrt{\frac{\rho}{2}} \eta \right) \right]^{-1} e^{-\rho\eta^2/2} \quad (10)$$

## 1.2 Nonparametric Hierarchical Co-clustering Model

The voxel activations  $x_{jis}$  are binary, with prior probability  $\phi_{k,l}$  given according to cluster memberships. We have that  $\phi \sim \text{Beta}(\omega^{\phi,1}, \omega^{\phi,2})$ , giving us the joint density of  $\boldsymbol{x}$  and  $\boldsymbol{\phi}$  conditioned on the cluster memberships  $\boldsymbol{z}$  and  $\boldsymbol{c}$ :

$$p(\boldsymbol{x}, \boldsymbol{\phi} | \boldsymbol{z}, \boldsymbol{c}) = \prod_{j,k,l} \left[ \frac{\Gamma(\omega^{\phi,1} + \omega^{\phi,2})}{\Gamma(\omega^{\phi,1})\Gamma(\omega^{\phi,2})} \phi_{k,l}^{\omega^{\phi,1}-1+\sum_{i,s} x_{jis}\delta(z_{ji}, k)\delta(c_s, l)} \cdot (1 - \phi_{k,l})^{\omega^{\phi,2}-1+\sum_{i,s}(1-x_{jis})\delta(z_{ji}, k)\delta(c_s, l)} \right]$$

Our model assumes a Dirichlet Process prior over the stimulus category memberships, using the stick-breaking construction for the prior  $\rho$  [1]:

$$p(\boldsymbol{c} | \boldsymbol{\rho}) = \prod_s \prod_l \rho_l^{\delta(c_s, l)} \quad (11)$$

$$\rho_l = u_l \prod_{l'=1}^{l-1} (1 - u_{l'}) \quad (12)$$

$$p(\boldsymbol{u}) = \prod_l (\chi - 1)(1 - u_l)^{\chi-1} \quad (13)$$

We assume a hierarchical Dirichlet process prior over the functional unit memberships, with subject-level weights  $\boldsymbol{\beta}$ . We will use a collapsed variational inference scheme [2], and therefore marginalize over these weights:

$$p(\boldsymbol{z} | \boldsymbol{\pi}, \alpha) = \int_{\boldsymbol{\beta}} p(\boldsymbol{z} | \boldsymbol{\beta}) p(\boldsymbol{\beta} | \boldsymbol{\pi}, \alpha), \quad (14)$$

$$= \prod_{j=1}^J \left[ \frac{\Gamma(\alpha)}{\Gamma(\alpha+N_j)} \prod_{k=1}^K \frac{\Gamma(\alpha\pi_k + n_{jk})}{\Gamma(\alpha\pi_k)} \right], \quad (15)$$

where  $n_{jk} = \sum_{i=1}^{N_j} \delta(z_{ji}, k)$  and  $K$  is the number of non-empty functional units in the configuration. To provide conjugacy with the Dirichlet prior over the group-level functional unit weights  $\boldsymbol{\pi}$ , we prefer the terms in Equation (15) involving this weight to appear as powers of  $\pi_k$ . However, the current form of the conditional distribution makes the computation of the posterior over  $\boldsymbol{\pi}$  hard. Now, note that for  $0 \leq r \leq n$ , we have  $\sum_{r=0}^n [n]_r \vartheta^r = \Gamma(\vartheta + n)/\Gamma(\vartheta)$ , where  $[n]_r$  are unsigned Stirling numbers of the first kind [3]. The collapsed variational approach uses this fact and the properties

of the Beta distribution to add an auxiliary variable  $\mathbf{r} = \{r_{ji}\}$  to the model:

$$p(\mathbf{z}, \mathbf{r}, | \pi, \alpha) \propto \prod_{j=1}^J \prod_{k=1}^K \binom{n_{jk}}{r_{jk}} (\alpha \pi_k)^{r_{jk}}, \quad (16)$$

where  $r_{jk} \in \{0, 1, \dots, n_{ji}\}$ . If we marginalize the distribution (16) over the auxiliary variable, we obtain the expression in (15).

## 2 Minimization of the Gibbs Free Energy

Let  $\mathbf{h} = \{\mathbf{x}, \mathbf{z}, \mathbf{c}, \mathbf{a}, \mathbf{e}, \boldsymbol{\lambda}, \mathbf{r}, \boldsymbol{\phi}, \pi, \rho\}$  denote the set of all unobserved variables. In the framework of variational inference, we approximate the model posterior on  $\mathbf{h}$  given the observed data  $p(\mathbf{h}|\mathbf{y})$  by a distribution  $q(\mathbf{h})$ . The approximation is performed through minimization of the Gibbs free energy function

$$\mathcal{F}[q] = E[\log q(\mathbf{h})] - E[\log p(\mathbf{y}, \mathbf{h})].$$

Here, and in the remainder of the paper,  $E[\cdot]$  and  $V[\cdot]$  indicate expected value and variance with respect to the distribution  $q$ . We assume a distribution  $q$  of the form:

$$q(\mathbf{h}) = q(\mathbf{r}|\mathbf{z}) \left[ \prod_k q(v_k) \right] \left[ \prod_l q(u_l) \right] \left[ \prod_{k,l} q(\phi_{k,l}) \right] \left[ \prod_s q(c_s) \right] \cdot \\ \prod_{j,i} \left( q(a_{ji}) q(\lambda_{ji}) q(z_{ji}) \left[ \prod_s q(x_{jis}) \right] \left[ \prod_h q(e_{jih}) \right] \right). \quad (17)$$

where we explicitly account for the dependency of the auxiliary variables on the functional unit memberships. Including this structure maintains the quality of approximation despite the introduction of the auxiliary variables [4].

### 2.1 Auxiliary variables

Assuming that all the other components of the posterior  $q$  are constant, we have:

$$\mathcal{F}[q(\mathbf{r}, | \mathbf{z})] = E_{\mathbf{z}} \left[ \sum_{\mathbf{r}} q(\mathbf{r}|\mathbf{z}) \left( \log q(\mathbf{r}|\mathbf{z}) + \right. \right. \\ \left. \left. - \sum_{j,k} \left\{ \log \binom{n_{jk}}{r_{jk}} + r_{jk} E[\log(\alpha \pi_k)] \right\} \right) \right] + \text{const.} \quad (18)$$

From this, we find that the optimal posterior on the auxiliary variables takes the form

$$q^*(\mathbf{r}|\mathbf{z}) = \prod_j \prod_k q(r_{jk}|\mathbf{z}). \quad (19)$$

Under  $q^*$ , we have for the auxiliary variable  $\mathbf{r}$ :

$$q(r_{jk}|\mathbf{z}) = \frac{\Gamma(\tilde{\omega}_{jk}^r)}{\Gamma(\tilde{\omega}_{jk}^r + n_{jk})} \binom{n_{jk}}{r_{jk}} (\tilde{\omega}_{jk}^r)^{r_{jk}}. \quad (20)$$

This distribution corresponds to the probability mass function for a random variable that describes the number of tables that  $n_{jk}$  customers occupy in a Chinese Restaurant Process with parameter  $\tilde{\omega}_{jk}^r$  [3]. Here, the optimal value of the parameter  $\tilde{\omega}_{jk}^r$  is given by

$$\begin{aligned} \log \tilde{\omega}_{jk}^r &= E[\log(\alpha \pi_k)] \\ &= \log \alpha + E[\log v_k] + \sum_{k' < k} E[\log(1 - v_{k'})]. \end{aligned} \quad (21)$$

Note that as a distribution parameterized by  $\log \tilde{\omega}_{jk}^r$ , Equation (20) defines a member of the exponential family of distributions. The expected value of the distribution is therefore computed by:

$$E[r_{jk}|\mathbf{z}] = \frac{\partial}{\partial \log \tilde{\omega}_{jk}^r} \log \frac{\Gamma(\tilde{\omega}_{jk}^r + n_{jk})}{\Gamma(\tilde{\omega}_{jk}^r)} \quad (22)$$

$$= \tilde{\omega}_{jk}^r \Psi(\tilde{\omega}_{jk}^r + n_{jk}) - \tilde{\omega}_{jk}^r \Psi(\tilde{\omega}_{jk}^r), \quad (23)$$

where  $\Psi(\omega) = \frac{\partial}{\partial \omega} \log \Gamma(\omega)$ . This expression is useful when we need the expected values of the auxiliary variable  $r_{jk}$  for updating the other components of the distribution:

$$E[r_{jk}] = E_{\mathbf{z}}[E_{\mathbf{r}}[r_{jk}|\mathbf{z}]] = \tilde{\omega}_{jk}^r E_{\mathbf{z}}[\Psi(\tilde{\omega}_{jk}^r + n_{jk}) - \Psi(\tilde{\omega}_{jk}^r)]. \quad (24)$$

Under  $q(\mathbf{z})$ , each variable  $n_{jk}$  is the sum of  $N_j$  independent Bernoulli random variables  $\delta(z_{ji}, k)$  for  $1 \leq i \leq N_j$  with the probability of success  $q(z_{ji} = k)$ . Therefore, as suggested in [2], we can use the Central Limit Theorem and approximate this term using a Gaussian distribution for  $n_{jk} > 0$ . Due to the independence of these Bernoulli variables, we have

$$\Pr(n_{jk} > 0) = 1 - \prod_{i=1}^{N_j} q(z_{ji} \neq k), \quad (25)$$

$$E[n_{jk}] = E[n_{jk}|n_{jk} > 0] \Pr(n_{jk} > 0), \quad (26)$$

$$E[n_{jk}^2] = E[n_{jk}^2|n_{jk} > 0] \Pr(n_{jk} > 0), \quad (27)$$

which we can use to easily compute  $E^+[n_{jk}] = E[n_{jk}|n_{jk} > 0]$  and  $V^+[n_{jk}] = V[n_{jk}|n_{jk} > 0]$ . Now, we can calculate  $E[r_{jk}]$  using Equation (24) by noting that:

$$E_{\mathbf{z}}[\Psi(\tilde{\omega}_{jk}^r + n_{jk}) - \Psi(\tilde{\omega}_{jk}^r)] \approx \Pr(n_{jk} > 0) \left[ \Psi(\tilde{\omega}_{jk}^r + E^+[n_{jk}]) - \Psi(\tilde{\omega}_{jk}^r) + \frac{V^+[n_{jk}]}{2} \Psi''(\tilde{\omega}_{jk}^r + E^+[n_{jk}]) \right]. \quad (28)$$

Lastly, based on the auxiliary variable  $r$ , we find that the optimal posterior over the system weight stick-breaking parameters is given by  $v_k \sim \text{Beta}(\tilde{\omega}_k^{v,1}, \tilde{\omega}_k^{v,2})$ , with parameters as follows:

$$\tilde{\omega}_k^{v,1} = 1 + \sum_j E[r_{jk}] \quad (29)$$

$$\tilde{\omega}_k^{v,2} = \gamma + \sum_{j,k'>k} E[r_{jk'}] \quad (30)$$

## 2.2 fMRI model variables

The free energy terms corresponding to  $e$  are

$$\begin{aligned} \mathcal{F}[q(e)] = \int_e q(e) \left( \log q(e) + e_{jih}^2 \omega_{jh}^{e,2} - e_{jih} \omega_{jh}^{e,1} - \sum_{j,i,t,h} \frac{E[\lambda_{jit}]}{2} \left[ e_{jih}^2 F_{ht}^2 + \right. \right. \\ \left. \left. - e_{jih} F_{ht} \left( y_{jit} - \sum_{h' \neq h} e_{jih'} F_{h't} - E[a_{jit}] \sum_s E[x_{jis}] G_{jst} \right) \right] \right) + \text{const.} \end{aligned} \quad (31)$$

Recall that we assume a factored form for  $q(e) = \prod_{j,i,h} q(e_{jih})$ . Minimizing with respect to this distribution gives  $q(e_{jih}) \propto \exp \left\{ -e_{jih}^2 \tilde{\omega}_{jh}^{e,2} + e_{jih} \tilde{\omega}_{jh}^{e,1} \right\}$ , with the parameters given by

$$\tilde{\omega}_{jh}^{e,2} = \omega_{jh}^{e,2} + \sum_t \frac{\lambda_{jht}}{2} F_{ht}^2 \quad (32)$$

$$\tilde{\omega}_{jh}^{e,1} = \omega_{jh}^{e,1} + \sum_t \lambda_{jht} F_{ht} \left( y_{jht} - \sum_{h' \neq h} E[e_{jih'}] F_{h't} - E[a_{jht}] \sum_s E[x_{jis}] G_{jst} \right) \quad (33)$$

For the activation heights  $a$ , we find

$$\begin{aligned} \mathcal{F}[q(a)] = \int_a q(a) \left( \log q(a) + a_{jim}^2 \omega_{jm}^{a,2} - a_{jim} \omega_{jm}^{a,1} + \right. \\ \left. \sum_{j,i,m} \frac{E[\lambda_{jim}]}{2} \sum_t D_{jmt} \left[ a_{jim}^2 \sum_{s,s'} E[x_{jis} x_{jis'}] G_{jst} G_{js't} + \right. \right. \\ \left. \left. - a_{jim} \sum_s G_{jst} E[x_{jis}] \left( y_{jim} - \sum_h E[e_{jih}] F_{ht} \right) \right] \right) + \text{const.} \end{aligned} \quad (34)$$

Assuming factorizability, minimizing gives  $q(a_{jim}) \propto \exp\left\{-a_{jim}^2\tilde{\omega}_{jim}^{a,2} + a_{jim}\tilde{\omega}_{jim}^{a,1}\right\}$ ,  $a \geq 0$ , with parameters given by

$$\tilde{\omega}_{jim}^{a,2} = \omega_{jm}^{a,2} + \sum_t D_{jmt} \frac{\lambda_{jim}}{2} \sum_{s,s',t} G_{jst} G_{js't} E[x_{jis} x_{jis'}] \quad (35)$$

$$\tilde{\omega}_{jim}^{a,1} = \omega_{jm}^{a,1} + \sum_t D_{jmt} \sum_s G_{jst} E[x_{jis}] \left( y_{jim} - \sum_h E[e_{jih}] F_{ht} \right) \quad (36)$$

The terms relating to the noise precisions  $\boldsymbol{\lambda}$  are computed as

$$\begin{aligned} \mathcal{F}[q(\boldsymbol{\lambda})] = & \int_{\boldsymbol{\lambda}} q(\boldsymbol{\lambda}) \left[ \log q(\boldsymbol{\lambda}) - \log(\lambda_{jim}) \omega_{jm}^{\lambda,1} - \lambda_{jim} \omega_{jm}^{\lambda,2} + \right. \\ & - \sum_{j,i,m,t} D_{jmt} \left( -\frac{\log(\lambda_{jim})}{2} + \frac{\lambda_{jim}}{2} \left[ y_{jim}^2 + (E[a_{jim}^2] \sum_{s,s'} E[x_{jis} x_{jis'}] G_{jst} G_{js't}) + \right. \right. \\ & \left. \left. \sum_h (E[e_{jih}^2] F_{ht}^2 + \sum_{h' \neq h} E[e_{jih}] E[e_{jih'}] F_{ht} F_{h't}) + \right. \right. \\ & \left. \left. - y_{jim} (\sum_h E[e_{jih}] F_{ht} + E[a_{jim}] \sum_s G_{jst} E[x_{jis}]) + \right. \right. \\ & \left. \left. \left( \sum_h E[e_{jih}] F_{ht} \right) \left( E[a_{jim}] \sum_s G_{jst} E[x_{jis}] \right) \right] \right) + \text{const.} \quad (37) \end{aligned}$$

For a factorizable  $q(\boldsymbol{\lambda}) = \prod_{jim} q(\lambda_{jim})$ , minimization with respect to  $q(\lambda_{jim})$  shows that  $q(\lambda_{jim}) \propto \exp\left\{\log(\lambda_{jim})\tilde{\omega}_{jim}^{\lambda,1} - \lambda_{jim}\tilde{\omega}_{jim}^{\lambda,2}\right\}$ , where the parameters are given by

$$\begin{aligned} \tilde{\omega}_{jim}^{\lambda,2} &= \omega_{jm}^{\lambda,2} + \sum_t \frac{D_{jmt}}{2} \\ \tilde{\omega}_{jim}^{\lambda,1} &= \omega_{jm}^{\lambda,1} + \left( y_{jim}^2 + \sum_h (E[e_{jih}^2] F_{ht}^2 + \sum_{h' \neq h} E[e_{jih}] E[e_{jih'}] F_{ht} F_{h't}) + \right. \\ &\quad \left. (E[a_{jim}^2] \sum_{s,s'} E[x_{jis} x_{jis'}] G_{jst} G_{js't}) + \right. \\ &\quad \left. - y_{jim} (\sum_h E[e_{jih}] F_{ht} + E[a_{jim}] \sum_s G_{jst} E[x_{jis}]) + \right. \\ &\quad \left. \left( \sum_h E[e_{jih}] F_{ht} \right) \left( E[a_{jim}] \sum_s G_{jst} E[x_{jis}] \right) \right) \quad (38) \end{aligned}$$

### 2.3 Voxel activation variables

The Gibbs free energy as a function only of the posterior over voxel activation variables  $\mathbf{x}$  can be computed as below. For notational convenience, we let  $\Phi_{jis} = \sum_{k,l} E[\log(\phi_{k,l})] \delta(z_{ji}, k) \delta(c_s, l)$  and  $\bar{\Phi}_{jis} = \sum_{k,l} E[\log(1 - \phi_{k,l})] \delta(z_{ji}, k) \delta(c_s, l)$ .

$$\begin{aligned} \mathcal{F}[q(\mathbf{x})] = & \sum_{\mathbf{x}} q(\mathbf{x}) \left[ \log q(\mathbf{x}) - \sum_{jis} \left[ (1 - x_{jis}) \bar{\Phi}_{jis} + \right. \right. \\ & x_{jis} \left( \Phi_{jis} + \sum_m \frac{E[\lambda_{jim}] E[a_{jim}]}{2} \sum_t D_{jmt} G_{jst} (y_{jim} - \sum_h e_{jih} F_{ht}) + \right. \\ & \left. \left. \sum_m \frac{E[\lambda_{jim}] E[a_{jim}^2]}{2} \sum_t D_{jmt} \sum_{s'} G_{jst} G_{js't} x_{jis'} \right) \right] \right] \quad (39) \end{aligned}$$

Minimization of this function with respect to  $q(\mathbf{x})$  yields

$$q(x_{jis} = 1) \propto \exp \left\{ \Phi_{jis} + \sum_m E[\lambda_{jim}] \sum_t D_{jmt} \left[ E[a_{jim}] G_{jst} (y_{jit} - \sum_h e_{jih} F_{ht}) + \right. \right. \\ \left. \left. - \frac{1}{2} E[a_{jim}^2] (G_{jst}^2 + \sum_{s' \neq s} 2G_{jst} G_{js't} x_{jis'}) \right] \right\} \quad (40)$$

$$q(x_{jis} = 0) \propto \exp \left\{ \hat{\Phi}_{jis} \right\} \quad (41)$$

## 2.4 System membership and activation probabilities

First, note that with the optimal posterior over the auxiliary variables  $q^*$ , defined in Equation (19), we have:

$$E[\log q^*(\mathbf{r}|\mathbf{z}) - \log p(\mathbf{z}, \mathbf{r} | \pi, \alpha)] = \sum_j \left( \log \Gamma(\alpha + N_j) - \log \Gamma(\alpha) \right) \\ + \sum_{jk} E_{\mathbf{z}} \left[ \log \Gamma(\tilde{\omega}_{jk}^r) - \log \Gamma(\tilde{\omega}_{jk}^r + n_{jk}) \right]. \quad (42)$$

The Gibbs free energy as a function of posterior over a single membership variable  $q(z_{ji})$  then becomes:

$$\mathcal{F}[q(z_{ji})] = \sum_k q(z_{ji} = k) \log q(z_{ji} = k) - \sum_k E_{\mathbf{z}} \left[ \log \Gamma(\tilde{\omega}_{jk}^r + n_{jk}) \right] + \\ - \sum_k q(z_{ji} = k) \sum_{l,s} \left[ q(x_{jis} = 1) q(c_s = l) E[\log \phi_{k,l}] + q(x_{jis} = 0) q(c_s = l) E[\log(1 - \phi_{k,l})] \right] \quad (43)$$

We can simplify the second term on the right hand side of Equation (43) as:

$$E_{\mathbf{z}} \left[ \log \Gamma(\tilde{\omega}_{r_{ji}} + n_{jk}) \right] = E_{\mathbf{z}} \left[ \delta(z_{ji}, k) \log(\tilde{\omega}_{jk}^r + n_{jk}^{-ji}) + \log \Gamma(\tilde{\omega}_{jk}^r + n_{jk}^{-ji}) \right], \quad (44)$$

$$= q(z_{ji} = k) E_{\mathbf{z}^{-ji}} [\log(\tilde{\omega}_{jk}^r + n_{jk}^{-ji})] + E_{\mathbf{z}^{-ji}} [\log \Gamma(\tilde{\omega}_{jk}^r + n_{jk}^{-ji})], \quad (45)$$

where  $n_{jk}^{-ji}$  and  $\mathbf{z}^{-ji}$  indicate the exclusion of voxel  $i$  in subject  $j$  and only the first term is a function of  $q(z_{ji})$ . Now, minimizing Equation (43) yields the following update for membership variables:

$$q(z_{ji} = k) \propto \exp \left\{ E_{\mathbf{z}^{-ji}} [\log(\tilde{\omega}_{jk}^r + n_{jk}^{-ji})] + \right. \\ \left. \sum_{l,s} \left( q(x_{jis} = 1) q(c_s = l) E[\log \phi_{k,l}] + q(x_{jis} = 0) q(c_s = l) E[\log(1 - \phi_{k,l})] \right) \right\},$$

In order to compute the first term on the right hand side, as with the Equation (28), we use a Gaussian approximation for the distribution of  $n_{jk}$  to find:

$$E_{\mathbf{z}^{-ji}} [\log(\tilde{\omega}_{jk}^r + n_{jk}^{-ji})] \approx \log(\tilde{\omega}_{jk}^r + E[n_{jk}^{-ji}]) - \frac{V[n_{jk}^{-ji}]}{2(\tilde{\omega}_{jk}^r + E[n_{jk}^{-ji}])^2}. \quad (46)$$

For the unit/category activation variables, we find

$$\mathcal{F}[q(\phi_{k,l})] = \int_v q(\phi_{k,l}) \left( \log q(\phi_{k,l}) - \sum_k \left[ \left\{ \omega^{\phi,1} + \sum_{j,i,s} q(z_{ji} = k) q(c_s = l) q(x_{jis} = 1) \right\} \log \phi_{k,l} + \right. \right. \\ \left. \left. \left\{ \omega^{\phi,2} + \sum_{j,i,s} q(z_{ji} = k) q(c_s = l) q(x_{jis} = 0) \right\} \log(1 - \phi_{k,l}) \right] \right) + \text{const.} \quad (47)$$

The posterior distribution is thus found to be  $\phi_{k,l} \sim \text{Beta}(\tilde{\omega}_{k,l}^{\phi,1}, \tilde{\omega}_{k,l}^{\phi,2})$ , with parameters as described by:

$$\tilde{\omega}_{k,l}^{\phi,1} = \omega^{\phi,1} + \sum_{j,i,s} q(z_{ji} = k)q(c_s = l)q(x_{jis} = 1) \quad (48)$$

$$\tilde{\omega}_{k,l}^{\phi,2} = \omega^{\phi,2} + \sum_{j,i,s} q(z_{ji} = k)q(c_s = l)q(x_{jis} = 0) \quad (49)$$

The stimulus category memberships are given by

$$\begin{aligned} \log q(c_s = l) &= \sum_{i,k} \left[ q(x_{jis} = 1)q(z_{ji} = k)E[\log \phi_{k,l}] + q(x_{jis} = 0)q(z_{ji} = k)E[\log(1 - \phi_{k,l})] \right] \\ &\quad + E[\log \rho] + \text{const.} \end{aligned} \quad (50)$$

Where  $\rho$  is given as in Equation (12), and the posterior over the stick-breaking weights is  $u_l \sim \text{Beta}(\tilde{\omega}_l^{u,1}, \tilde{\omega}_l^{u,2})$ :

$$\tilde{\omega}_l^{u,1} = 1 + \sum_s q(c_s = l) \quad (51)$$

$$\tilde{\omega}_l^{u,1} = \chi + \sum_{s, l' < l} q(c_s = l') \quad (52)$$

## References

- [1] D.M. Blei and M.I. Jordan. Variational inference for Dirichlet process mixtures. *Bayesian Analysis*, 1(1):121–144, 2006.
- [2] Y.W. Teh, K. Kurihara, and M. Welling. Collapsed variational inference for HDP. *Advances in Neural Information Processing Systems*, 20:1481–1488, 2008.
- [3] C.E. Antoniak. Mixtures of Dirichlet processes with applications to Bayesian nonparametric problems. *The Annals of Statistics*, 2(6):1152–1174, 1974.
- [4] Y.W. Teh, D. Newman, and M. Welling. A collapsed variational bayesian inference algorithm for latent dirichlet allocation. *Advances in Neural Information Processing Systems*, 19:1353–1360, 2007.