
Supplementary Material: Multiple Clustering Views from Multiple Uncertain Experts

Yale Chang Junxiang Chen Michael H. Cho
Peter J. Castaldi Edwin K. Silverman Jennifer G. Dy

1 Parameter Settings

All approaches need to specify K , the number of clusters.

When applying our MCVV to the synthetic and benchmark datasets:

- If there is only one ground-truth expert view, we set K to be the true number of clusters in that view.
- If there are more than one ground-truth expert view, we set K to be the maximal value among the true number of clusters in all expert views.

For SemiCrowd, ITML, MPCKMeans, CSPA, since we only apply them to one expert view, we set K to be the true number of clusters in that view.

For COPD data, we set $K = 4$ for all approaches according to a recent study on COPD [3].

Besides the number of clusters, the parameters that are specific to each approach are set as follows.

1.1 Proposed Approach: MCVV

For variational inference in our approach, we use the following parameter settings

1. G , the number of components in truncated Dirichlet Process, is set to be $M/2$, where M is the total number of experts. In this way, we try to enforce the constraint that on average, there should be at least two experts in each view. In all experiments, the number of expert views recovered by our approach is smaller than $G = M/2$. Therefore, the value of G we use is large enough to discover the true number of expert views.
2. For the parameters of prior distributions:
 - $p(\alpha_m), p(\beta_m)$: set the parameters of prior Beta distributions to be $(10, 1)$ to incorporate the prior knowledge that each expert’s accuracy parameters should be far away from 0.5 (random guess) and close to 1. The choice can be illustrated from Figure 1. Under this setting, there is very small probability that the accuracy parameters can be close to 0.5.

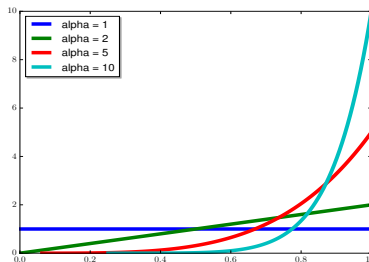


Figure 1: Probability density function of Beta distribution $Beta(\alpha, 1)$, $\alpha = 10$ can effectively make the accuracy parameters have very small probability to be close to 0.5.

- $p(\nu_g)$: set its concentration parameter $\gamma = 1$, the recovered number of experts is stable across a range of different γ values (from 0.5 to 10).
- $p(W_{ij}^{(g)}), p(b_i^{(g)})$: set the mean and standard deviation of prior Gaussian distributions to be 0 and 1 respectively.

3. Initializations of variational parameters:

- the parameters of variational Beta distribution are initialized to be equal to those of the prior distribution;
- the means of variational Gaussian distributions are initialized by randomly drawing samples from Gaussian distribution $\mathcal{N}(0, \frac{1}{D})$, the standard deviations of variational Gaussian distributions are initialized to be 0.001. This initialization strategy is similar to Xavier initialization [5] in order to avoid gradient saturation;
- $\eta_{i,:}^{(g)}$ are initialized using $W^{(g)}, b^{(g)}$ according to the discriminative clustering model described in the main paper;
- $\phi_{m,:}$ are initialized by sampling from a Dirichlet distribution with all its parameters being equal to 1, introducing most randomness for the initialization.

4. The number of random initializations for optimization is set to be 50 and the results are stable across different runs in all our experiments.

5. Cluster samples in the testing set:

After learning our model using the training set, we can obtain the variational distributions of weight $q(W^{(g)})$ and offset $q(b^{(g)})$. We can cluster x_t , a sample from the testing set by integrating out $W^{(g)}, b^{(g)}$ through Monte Carlo approximation to the integration:

$$p(Z_t^{(g)} = k | x_t) = \int \int p(Z_t^{(g)} | W^{(g)}, b^{(g)}; x_t) q(W^{(g)}) q(b^{(g)}) dW^{(g)} db^{(g)} \quad (1)$$

$$\approx \frac{1}{L} \sum_{l=1}^L p(Z_t^{(g)} | \widehat{W}^{(g)}_l, \widehat{b}^{(g)}_l; x_t) \quad (2)$$

where $\widehat{W}^{(g)}_l, \widehat{b}^{(g)}_l$ are the l -th sample from $q(W^{(g)})$ and $q(b^{(g)})$ respectively. We set $L = 100$ in all experiments. x_t is assigned to the cluster corresponding to the largest probability:

$$\hat{y}_t = \arg \max_{k=1, \dots, K} p(Z_t^{(g)} = k | x_t) \quad (3)$$

where \hat{y}_t is the predicted cluster label for x_t .

1.2 Meta Clustering

After computing the similarity matrix between multiple experts, we apply spectral clustering [6] to assign experts to different views. The number of views are automatically determined using the eigen-gap heuristic [7].

1.3 SemiCrowd

We follow parameter settings recommended by the author in [9]. In particular, d_0, d_1 , two thresholds used to filter out uncertain sample pairs in the average similarity matrix, are set to be 0 and 0.8 respectively. Parameters of the ℓ_1 regularized matrix completion algorithm is set according to heuristics in the literature [8].

1.4 ITML

As is suggested by the author [4], we set lower and upper bounds associated with the constraint terms to be the 5th and 95th percentiles of the observed distribution of distances between pairs of points within the dataset. Other parameters related to the convergence of the algorithm are set to be default.

1.5 MPCKMeans

We specify the constraints according to the instructions listed on the author’s website and directly run the author’s Java implementation [1]

1.6 CSPA

After computing the average similarity matrix, we use spectral clustering [6] to obtain the cluster labels.

2 Variational Inference

2.1 Variational Distribution

Denote $h = \{\alpha_{1:M}, \beta_{1:M}, c_{1:M}, \nu_{1:G}, Z^{(1:G)}, W^{(1:G)}, b^{(1:G)}\}$ as the collection of latent variables, where G is the number of components of the truncated Dirichlet Process. We assume the variational distribution has the following factorization formula:

$$q(\alpha_{1:M}, \beta_{1:M}, c_{1:M}, \nu_{1:G}, Z^{(1:G)}, W^{(1:G)}, b^{(1:G)}) \quad (4)$$

$$= q(\alpha_{1:M}) \cdot q(\beta_{1:M}) \cdot q(c_{1:M}) \cdot q(\nu_{1:G}) \cdot q(Z^{(1:G)}) \cdot q(W^{(1:G)}) \cdot q(b^{(1:G)}) \quad (5)$$

$$= \prod_{m=1}^M [q(\alpha_m) \cdot q(\beta_m) \cdot q(c_m)] \prod_{g=1}^G [q(\nu_g) q(Z^{(g)}) q(W^{(g)}) q(b^{(g)})] \quad (6)$$

$$= \prod_{m=1}^M [q(\alpha_m) \cdot q(\beta_m) \cdot q(c_m)] \prod_{g=1}^G \left[q(\nu_g) \prod_{i=1}^n q(Z_i^{(g)}) \prod_{i=1}^d \prod_{j=1}^K q(W_{ij}^{(g)}) \prod_{i=1}^K q(b_i^{(g)}) \right] \quad (7)$$

where the marginal variational distribution of each random variable is

- $q(\alpha_m) = \text{Beta}(\tau_{\alpha^1}^{(m)}, \tau_{\alpha^2}^{(m)});$
- $q(\beta_m) = \text{Beta}(\tau_{\beta^1}^{(m)}, \tau_{\beta^2}^{(m)});$
- $q(c_m) = \text{Cat}(\phi_{m,:});$
- $q(\nu_g) = \text{Beta}(\tau_{\nu^1}^{(g)}, \tau_{\nu^2}^{(g)});$
- $q(Z_i^{(g)}) = \text{Cat}(\eta_{i,:}^{(g)});$
- $q(W_{ij}^{(g)}) = \mathcal{N}(\mu_{W_{ij}}^{(g)}, \sigma_{W_{ij}}^{(g)2});$
- $q(b_i^{(g)}) = \mathcal{N}(\mu_{b_i}^{(g)}, \sigma_{b_i}^{(g)2}).$

We use θ to denote all the variational parameters, which consist of the following

- For $m = 1 \cdots M$, consider $\{\tau_{\alpha^1}^{(m)}, \tau_{\alpha^2}^{(m)}, \tau_{\beta^1}^{(m)}, \tau_{\beta^2}^{(m)}, \phi_{m,:} \in \Delta^{G-1}\}$, where Δ^{G-1} is a simplex in G -dimensional space.
- For $g = 1 \cdots G$, consider $\tau_{\nu^1}^{(g)}, \tau_{\nu^2}^{(g)}, \eta_{i,:}^{(g)} \in \Delta^{K-1} (i = 1 \cdots n), \mu_{W_{ij}}^{(g)}, \sigma_{W_{ij}}^{(g)2} (i = 1 \cdots d, j = 1 \cdots K), \mu_{b_i}^{(g)}, \sigma_{b_i}^{(g)2} (i = 1 \cdots K)$

Besides the simplex constraints on the parameters of categorical distribution, both the parameters of Beta distribution and the standard deviation of Gaussian distribution should have positive constraints.

2.2 Evidence Lower Bound

Given variational distribution $q(h; \theta)$, the log-likelihood $\log p(S^{(1:M)})$ can be decomposed as

$$\log p(S^{(1:M)}) = \mathcal{L}_{q(h;\theta)} + KL [q(h; \theta) | p(h|S^{(1:M)})] \quad (8)$$

$$\geq \mathcal{L}_{q(h;\theta)} \quad (9)$$

$$= \mathbb{E}_{q(h;\theta)} [\log p(h, S^{(1:M)}) - \log q(h; \theta)] \quad (10)$$

$$= \mathbb{E}_{q(h;\theta)} [\log p(S^{(1:M)}|h) + \log p(h) - \log q(h; \theta)] \quad (11)$$

$$= \mathbb{E}_{q(h;\theta)} [\log p(S^{(1:M)}|h)] - KL [q(h; \theta) | p(h)] \quad (12)$$

The first term prefers the variational distribution $q(h; \theta)$ that maximizes the expected conditional likelihood, the second term forces the variational distribution to be close to the prior distribution.

The overall objective of variational inference is to maximize the evidence lower bound (ELBO) $\mathcal{L}_{q(h; \theta)}$ w.r.t. θ .

Let $f(h; \theta) = \log p(S^{(1:M)} | h) + \log p(h) - \log q(h; \theta)$, the evidence lower bound can be expressed as

$$\mathcal{L}_{q(h; \theta)} = \mathbb{E}_{q(h; \theta)} [f(h; \theta)] \quad (13)$$

We first compute the three terms in $f(h; \theta)$ and then evaluate their expectations w.r.t. the variational distribution $q(h; \theta)$.

2.3 ELBO: The First Term $\log p(S^{(1:M)} | h)$

Consider the first term:

$$\log p(S^{(1:M)} | h) = \log \left(\prod_{m=1}^M \prod_{(i,j) \in E^{(m)}} p(S_{ij}^{(m)} | h) \right) \quad (14)$$

$$= \sum_{m=1}^M \sum_{(i,j) \in E^{(m)}} \log p(S_{ij}^{(m)} | h) \quad (15)$$

$$= \sum_{m=1}^M \left\{ \left(\sum_{(i,j) \in E^{(m)}} 1 \right) \log \beta_m + \left(\sum_{(i,j) \in E^{(m)}} S_{ij}^{(m)} \right) \log \left(\frac{1 - \beta_m}{\beta_m} \right) \right. \quad (16)$$

$$+ \left(\sum_{(i,j) \in E^{(m)}} S_{ij}^{(m)} A_{ij}^{(m)} \right) \log \left[\frac{\alpha_m \beta_m}{(1 - \alpha_m)(1 - \beta_m)} \right]$$

$$\left. + \left(\sum_{(i,j) \in E^{(m)}} A_{ij}^{(m)} \right) \log \left(\frac{1 - \alpha_m}{\beta_m} \right) \right\}$$

To compute $\mathbb{E}_{q(h; \theta)} [\log p(S^{(1:M)} | h)]$, the following formulas can be used:

$$\mathbb{E}_{q(h; \theta)} [\log \alpha_m] = \mathbb{E}_{q(\alpha_m)} [\log \alpha_m] = \psi(\tau_{\alpha_1}^{(m)}) - \psi(\tau_{\alpha_1}^{(m)} + \tau_{\alpha_2}^{(m)}) \quad (17)$$

$$\mathbb{E}_{q(h; \theta)} [\log(1 - \alpha_m)] = \mathbb{E}_{q(\alpha_m)} [\log(1 - \alpha_m)] = \psi(\tau_{\alpha_2}^{(m)}) - \psi(\tau_{\alpha_1}^{(m)} + \tau_{\alpha_2}^{(m)}) \quad (18)$$

$$\mathbb{E}_{q(h; \theta)} [\log \beta_m] = \mathbb{E}_{q(\beta_m)} [\log \beta_m] = \psi(\tau_{\beta_1}^{(m)}) - \psi(\tau_{\beta_1}^{(m)} + \tau_{\beta_2}^{(m)}) \quad (19)$$

$$\mathbb{E}_{q(h; \theta)} [\log(1 - \beta_m)] = \mathbb{E}_{q(\beta_m)} [\log(1 - \beta_m)] = \psi(\tau_{\beta_2}^{(m)}) - \psi(\tau_{\beta_1}^{(m)} + \tau_{\beta_2}^{(m)}) \quad (20)$$

$$\mathbb{E}_{q(h; \theta)} [A_{ij}^{(m)}] = \mathbb{E}_{q(c_m, Z^{(c_m)})} [\mathbb{I}(Z_i^{(c_m)} = Z_j^{(c_m)})] \quad (21)$$

$$= \mathbb{E}_{q(c_m)} [\mathbb{E}_{q(Z^{(c_m)})} [\mathbb{I}(Z_i^{(c_m)} = Z_j^{(c_m)})]] \quad (22)$$

$$= \mathbb{E}_{q(c_m)} \left[\sum_{k=1}^K \eta_{ik}^{(c_m)} \eta_{jk}^{(c_m)} \right] \quad (23)$$

$$= \sum_{g=1}^G \left[\phi_{mg} \sum_{k=1}^K \eta_{ik}^g \eta_{jk}^g \right] \quad (24)$$

2.4 ELBO: The Second Term $\log p(h)$

Consider the second term:

$$\log p(h) = \log \left(\prod_{m=1}^M \left(p(\alpha_m) p(\beta_m) p(c_m | \nu_{1:\infty}) \right) \prod_{g=1}^{\infty} \left(p(\nu_g) p(Z^{(g)} | W^{(g)}, b^{(g)}) p(W^{(g)}) p(b^{(g)}) \right) \right) \quad (25)$$

$$\begin{aligned} &= \sum_{m=1}^M \left(\log p(\alpha_m) + \log p(\beta_m) + \log p(c_m | \nu_{1:\infty}) \right) \\ &+ \sum_{g=1}^{\infty} \left(\log p(\nu_g) + \log p(Z^{(g)} | W^{(g)}, b^{(g)}) + \log p(W^{(g)}) + \log p(b^{(g)}) \right) \end{aligned} \quad (26)$$

where

$$\log p(\alpha_m) = (\tau_{\alpha^{10}}^{(m)} - 1) \log \alpha_m + (\tau_{\alpha^{20}}^{(m)} - 1) \log(1 - \alpha_m) - \log \left(\frac{\Gamma(\tau_{\alpha^{10}}^{(m)}) \Gamma(\tau_{\alpha^{20}}^{(m)})}{\Gamma(\tau_{\alpha^{10}}^{(m)} + \tau_{\alpha^{20}}^{(m)})} \right) \quad (27)$$

$$\log p(\beta_m) = (\tau_{\beta^{10}}^{(m)} - 1) \log \beta_m + (\tau_{\beta^{20}}^{(m)} - 1) \log(1 - \beta_m) - \log \left(\frac{\Gamma(\tau_{\beta^{10}}^{(m)}) \Gamma(\tau_{\beta^{20}}^{(m)})}{\Gamma(\tau_{\beta^{10}}^{(m)} + \tau_{\beta^{20}}^{(m)})} \right) \quad (28)$$

$$\log p(c_m | \nu_{1:\infty}) = \sum_{g=1}^{\infty} c_{m,g} \left(\log \nu_g + \sum_{j=1}^{g-1} \log(1 - \nu_j) \right) \quad (29)$$

$$\log p(\nu_g) = (\gamma - 1) \log(1 - \nu_g) - \log \left(\frac{\Gamma(\gamma)}{\Gamma(1 + \gamma)} \right) \quad (30)$$

$$= (\gamma - 1) \log(1 - \nu_g) + \log \gamma \quad (31)$$

$$\log p(Z^{(g)} | W^{(g)}, b^{(g)}) = \log \prod_{i=1}^N p(Z_i^{(g)} | W^{(g)}, b^{(g)}) \quad (32)$$

$$= \log \prod_{i=1}^N \prod_{k=1}^K p(Z_i^{(g)} = k | W^{(g)}, b^{(g)})^{\mathbb{I}[Z_i^{(g)} = k]} \quad (33)$$

$$= \sum_{i=1}^N \sum_{k=1}^K \mathbb{I}[Z_i^{(g)} = k] \log p(Z_i^{(g)} = k | W^{(g)}, b^{(g)}) \quad (34)$$

$$= \sum_{i=1}^N \sum_{k=1}^K \mathbb{I}[Z_i^{(g)} = k] \left((w_k^{(g)T} x_i + b_k^{(g)}) - \log \left(\sum_{j=1}^K \exp(w_j^{(g)T} x_i + b_j^{(g)}) \right) \right) \quad (35)$$

$$\log p(W^{(g)}) = \sum_{i=1}^d \sum_{j=1}^K \left(-\log(\sqrt{2\pi} \sigma_{W_{ij}^{(g)}}^{(g)}) - \frac{(W_{ij}^{(g)} - \mu_{W_{ij}^{(g)}}^{(g)})^2}{2\sigma_{W_{ij}^{(g)}}^{(g)2}} \right) \quad (36)$$

$$\log p(b^{(g)}) = \sum_{i=1}^K \left(-\log(\sqrt{2\pi} \sigma_{b_i^{(g)}}^{(g)}) - \frac{(b_i^{(g)} - \mu_{b_i^{(g)}}^{(g)})^2}{2\sigma_{b_i^{(g)}}^{(g)2}} \right) \quad (37)$$

To compute $\mathbb{E}_{q(h;\theta)}[\log p(h)]$, the following formulas can be used besides the formulas used to compute

$\mathbb{E}_{q(h;\theta)}[\log p(S^{(1:M)} | h)]:$

$$\mathbb{E}_{q(h;\theta)}[\log \nu_g] = \mathbb{E}_{q(\nu_g)}[\log \nu_g] = \psi(\tau_{\nu^1}^{(g)}) - \psi(\tau_{\nu^1}^{(g)} + \tau_{\nu^2}^{(g)}) \quad (38)$$

$$\mathbb{E}_{q(h;\theta)}[\log(1 - \nu_g)] = \mathbb{E}_{q(\nu_g)}[\log(1 - \nu_g)] = \psi(\tau_{\nu^2}^{(g)}) - \psi(\tau_{\nu^1}^{(g)} + \tau_{\nu^2}^{(g)}) \quad (39)$$

$$\mathbb{E}_{q(h;\theta)}[c_{mg}] = \mathbb{E}_{q(c_m)}[c_{mg}] = \phi_{mg} \quad (40)$$

$$\mathbb{E}_{q(h;\theta)}[\mathbb{I}[Z_i^{(g)} = k]] = \eta_{ik}^{(g)} \quad (41)$$

$$\mathbb{E}_{q(h;\theta)}[w_k^{(g)}] = \mathbb{E}_{q(W^{(g)})}[w_k^{(g)}] = \mu_{W.k}^{(g)} \quad (42)$$

$$\mathbb{E}_{q(h;\theta)}[(W_{ij}^{(g)} - \mu_{W_{ij}0}^{(g)})^2] = \sigma_{W_{ij}}^{(g)2} + (\mu_{W_{ij}}^{(g)} - \mu_{W_{ij}0}^{(g)})^2 \quad (43)$$

$$\mathbb{E}_{q(h;\theta)}[(b_i^{(g)} - \mu_{b_i0}^{(g)})^2] = \sigma_{b_i}^{(g)2} + (\mu_{b_i}^{(g)} - \mu_{b_i0}^{(g)})^2 \quad (44)$$

2.4.1 Upper Bound of Log-sum Function

The computation of $\mathbb{E}_{q(h;\theta)} \left[\log \left(\sum_{j=1}^K \exp(w_j^{(g)T} x_i + b_j) \right) \right]$ does not have closed form and can only be approximated through sampling. However, we can use the upper bound of log-sum function based on the log concavity to derive a lower bound of ELBO [2].

$$\log \left(\sum_{j=1}^K \exp(w_j^{(g)T} x_i + b_j^{(g)}) \right) \leq r_i^{(g)} \sum_{j=1}^K \exp(w_j^{(g)T} x_i + b_j^{(g)}) - \log(r_i^{(g)}) - 1 \quad (45)$$

We need introduce new variational parameters $r_i^{(g)} > 0$ to optimize. The expectation term can also be upper bounded as follows

$$\begin{aligned} & \mathbb{E}_{q(h;\theta)} \left[\log \left(\sum_{j=1}^K \exp(w_j^{(g)T} x_i + b_j) \right) \right] \\ & \leq \mathbb{E}_{q(h;\theta)} \left[r_i^{(g)} \sum_{j=1}^K \exp(w_j^{(g)T} x_i + b_j^{(g)}) - \log(r_i^{(g)}) - 1 \right] \end{aligned} \quad (46)$$

$$= r_i^{(g)} \sum_{j=1}^K \mathbb{E}_{q(h;\theta)} \left[\exp \left(\sum_{l=1}^d w_{lj}^{(g)} x_{il} + b_j^{(g)} \right) \right] - \log(r_i^{(g)}) - 1 \quad (47)$$

$$= r_i^{(g)} \sum_{j=1}^K \left(\mathbb{E}_{q(b_j^{(g)})} [e^{b_j^{(g)}}] \prod_{l=1}^d \mathbb{E}_{q(W_{lj}^{(g)})} [e^{w_{lj}^{(g)} x_{il}}] \right) - \log(r_i^{(g)}) - 1 \quad (48)$$

where the expectation can be evaluated using the mean of log-normal distribution.

$$\mathbb{E}_{q(b_j^{(g)})} [e^{b_j^{(g)}}] = \exp \left(\mu_{b_j}^{(g)} + \frac{\sigma_{b_j}^{(g)2}}{2} \right) \quad (49)$$

$$\mathbb{E}_{q(W_{lj}^{(g)})} [e^{w_{lj}^{(g)} x_{il}}] = \exp \left(\mu_{W_{lj}}^{(g)} x_{il} + \frac{\sigma_{W_{lj}}^{(g)2} x_{il}^2}{2} \right) \quad (50)$$

Therefore, the upper bound can be written as

$$\begin{aligned} & \mathbb{E}_{q(h;\theta)} \left[\log \left(\sum_{j=1}^K \exp(w_j^{(g)T} x_i + b_j) \right) \right] \\ & \leq r_i^{(g)} \sum_{j=1}^K \left(\exp \left(\mu_{b_j}^{(g)} + \frac{\sigma_{b_j}^{(g)2}}{2} + \sum_{l=1}^d \left(\mu_{W_{lj}}^{(g)} x_{il} + \frac{\sigma_{W_{lj}}^{(g)2} x_{il}^2}{2} \right) \right) \right) - \log(r_i^{(g)}) - 1 \end{aligned} \quad (51)$$

The upper bound becomes minimal when

$$r_i^{(g)} = \frac{1}{\sum_{j=1}^K \left(\exp \left(\mu_{b_j}^{(g)} + \frac{\sigma_{b_j}^{(g)2}}{2} + \sum_{l=1}^d \left(\mu_{W_{lj}}^{(g)} x_{il} + \frac{\sigma_{W_{lj}}^{(g)2} x_{il}^2}{2} \right) \right) \right)} \quad (52)$$

Then the upper bound becomes

$$E_{q(h;\theta)} \left[\log \left(\sum_{j=1}^K \exp(w_j^{(g)T} x_i + b_j) \right) \right] \quad (53)$$

$$\leq \log \sum_{j=1}^K \left(\exp \left(\mu_{b_j}^{(g)} + \frac{\sigma_{b_j}^{(g)2}}{2} + \sum_{l=1}^d \left(\mu_{W_{lj}}^{(g)} x_{il} + \frac{\sigma_{W_{lj}}^{(g)2} x_{il}^2}{2} \right) \right) \right) \quad (54)$$

To avoid numerical overflow when computing this bound due to the potential large exponent of the exponential function, we can use the `logsumexp()` function implemented in Scipy. The idea is to extract the maximal value in the sequence and compute its log, then each exponent in the sequence will be less than 1, leading to stable numerical behavior.

2.5 ELBO: The Third Term $\log q(h; \theta)$

Consider the third term:

$$\log q(h; \theta) = \sum_{m=1}^M \left(\log q(\alpha_m) + \log q(\beta_m) + \log q(c_m) \right) \quad (55)$$

$$+ \sum_{g=1}^G \left(\log q(\nu_g) + \sum_{i=1}^n \sum_{k=1}^K \mathbb{I}[Z_i^{(g)} = k] \log q(Z_i^{(g)} = k) + \sum_{i=1}^d \sum_{j=1}^K \log q(W_{ij}^{(g)}) + \sum_{i=1}^K \log q(b_i^{(g)}) \right) \quad (56)$$

where

$$\log q(\alpha_m) = (\tau_{\alpha^1}^{(m)} - 1) \log \alpha_m + (\tau_{\alpha^2}^{(m)} - 1) \log(1 - \alpha_m) - \log \left(\frac{\Gamma(\tau_{\alpha^1}^{(m)}) \Gamma(\tau_{\alpha^2}^{(m)})}{\Gamma(\tau_{\alpha^1}^{(m)} + \tau_{\alpha^2}^{(m)})} \right) \quad (57)$$

$$\log q(\beta_m) = (\tau_{\beta^1}^{(m)} - 1) \log \beta_m + (\tau_{\beta^2}^{(m)} - 1) \log(1 - \beta_m) - \log \left(\frac{\Gamma(\tau_{\beta^1}^{(m)}) \Gamma(\tau_{\beta^2}^{(m)})}{\Gamma(\tau_{\beta^1}^{(m)} + \tau_{\beta^2}^{(m)})} \right) \quad (58)$$

$$\log q(c_m) = \sum_{g=1}^G c_{mg} \log \phi_{mg} \quad (59)$$

$$\log q(\nu_g) = (\tau_{\nu^1}^{(g)} - 1) \log \nu_g + (\tau_{\nu^2}^{(g)} - 1) \log(1 - \nu_g) - \log \left(\frac{\Gamma(\tau_{\nu^1}^{(g)}) \Gamma(\tau_{\nu^2}^{(g)})}{\Gamma(\tau_{\nu^1}^{(g)} + \tau_{\nu^2}^{(g)})} \right) \quad (60)$$

$$\log q(Z_i^{(g)} = k) = \log \eta_{ik}^{(g)} \quad (61)$$

$$\log q(W_{ij}^{(g)}) = -\log(\sqrt{2\pi}\sigma_{W_{ij}}^{(g)}) - \frac{(W_{ij}^{(g)} - \mu_{W_{ij}}^{(g)})^2}{2\sigma_{W_{ij}}^{(g)2}} \quad (62)$$

$$\log q(b_i^{(g)}) = -\log(\sqrt{2\pi}\sigma_{b_i}^{(g)}) - \frac{(b_i^{(g)} - \mu_{b_i}^{(g)})^2}{2\sigma_{b_i}^{(g)2}} \quad (63)$$

The computation of $E_{q(h;\theta)}[\log q(h; \theta)]$ can be done using formulas from the previous two subsections.

3 Weights of Different Experts

As we can see from the derivation of the first-term $\log p(S^{(1:M)}|h)$ of ELBO, if we view Equation (16) as a function of $A_{ij}^{(m)}$, the weight of $S_{ij}^{(m)}$ (constraints provided by the m -th expert) is: $\log \frac{\alpha_m \beta_m}{(1-\alpha_m)(1-\beta_m)}$.

References

- [1] M. Bilenko. Java implementation of mpckmeans, 2004.
- [2] G. Bouchard. Efficient bounds for the softmax function and applications to approximate inference in hybrid models. In *NIPS 2007 workshop for approximate Bayesian inference in continuous/hybrid systems*. Citeseer, 2007.
- [3] P. J. Castaldi, J. Dy, J. Ross, Y. Chang, G. R. Washko, D. Curran-Everett, A. Williams, D. A. Lynch, B. J. Make, J. D. Crapo, et al. Cluster analysis in the copdgene study identifies subtypes of smokers with distinct patterns of airway disease and emphysema. *Thorax*, pages thoraxjnl–2013, 2014.
- [4] J. V. Davis, B. Kulis, P. Jain, S. Sra, and I. S. Dhillon. Information-theoretic metric learning. In *Proceedings of the 24th international conference on Machine learning*, pages 209–216. ACM, 2007.
- [5] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Aistats*, volume 9, pages 249–256, 2010.
- [6] A. Y. Ng, M. I. Jordan, Y. Weiss, et al. On spectral clustering: Analysis and an algorithm. *Advances in neural information processing systems*, 2:849–856, 2002.
- [7] U. Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007.
- [8] J. Yi, R. Jin, A. K. Jain, and S. Jain. Crowdclustering with sparse pairwise labels: A matrix completion approach. In *AAAI Workshop on Human Computation*, volume 2. Citeseer, 2012.
- [9] J. Yi, R. Jin, S. Jain, T. Yang, and A. K. Jain. Semi-crowdsourced clustering: Generalizing crowd labeling by robust distance metric learning. In *Advances in Neural Information Processing Systems*, pages 1772–1780, 2012.