

# A Variational Approach for Bayesian Density Regression

Eric Chuu

ERICCHUU@STAT.TAMU.EDU

*Department of Statistics*

*Texas A&M University*

*College Station, TX 77840, USA*

## Abstract

In the Bayesian density regression problem, a mixture of experts model is often used because of the high flexibility in estimating conditional densities. In this paper, we discuss the case when covariate dependent weights are used in the approximating mixture density. Under this framework, however, traditional Bayesian methods result in computational difficulties when the dimension of the covariates is large. In order to remedy this problem and to provide a method for faster inference, we propose using a variational approximation to estimate the conditional density. We also discuss upper bounds for approximating quantities that lack a closed form so that a coordinate ascent algorithm is viable.

**Keywords:** Bayesian Density Regression, Variational Bayes, Mixture Models

## 1. Introduction

In the Bayesian density regression problem, we observe data  $(y_n, x_n)_{n=1}^N$ , and the goal is the estimate the conditional density of  $y | x$ . A common approach for doing this is to model the density using a mixture of gaussians, such as the following,

$$f(y | x) = \sum_k \pi_k \mathcal{N}(y | \mu_k(x), \tau_k^{-1}) \quad (1)$$

While the representation of the density using predictor-independent weights yields less expensive computation, this approach often lacks flexibility to make it useful in practice and results in a reliance on have too many mixture components. As a result, there have been many proposed models that consider predictor-dependent weights. Some examples include using a kernel stick-breaking process (**dunsonpark:08**) or logit stick-breaking prior (**durante:17**) to generate the covariate-dependent weights. In the former method, the increased flexibility comes at heavy computational cost, and in the latter method, the process from which the weights are generated does not allow for intuitive inference. In our proposed model, the covariates enter the weights through a logistic link function so that we can naturally perform inference on the coefficients. More specifically, we can model

$$f(y | x) = \sum_k^K \pi_k(x) \mathcal{N}(y | \mu_k(x), \tau_k^{-1}) \quad (2)$$

where  $\pi_k \propto \exp(x^\top \gamma_k)$ . In order to perform fast inference on the model parameters, we adopt a variational approach to obtain an approximating distribution to the true posterior.

Using this covariate-dependent setup, however, introduces a problem in the traditional coordinate ascent algorithm such that it prevents closed form updates of the variational distributions. In section 2, we give an overview of the priors used in the problem. In section 3, we consider the family of variational distributions that we use to approximate the true posterior. We then propose a way to obtain closed form updates by considering an upper bound on the problematic quantity, and in section 4, we formulate the complete algorithm. Finally, we discuss potential shortcomings of the proposed algorithm and other bounds that could be used to obtain more accurate approximating distributions.

## 2. Notation and Prior Specification

For the set of observed data, we denote  $\mathbf{y} = \{y_1, \dots, y_N\}$ ,  $\mathbf{X} = \{x_1, \dots, x_N\}$ , where each  $x_n \in \mathbb{R}^D$ . Let  $\boldsymbol{\beta} = \{\beta_1, \dots, \beta_K\}$  and  $\boldsymbol{\gamma} = \{\gamma_1, \dots, \gamma_K\}$  denote the  $D$ -dimensional coefficient vectors in the mixture weights and in gaussian mixture components, respectively. Finally, let  $\boldsymbol{\tau} = \{\tau_1, \dots, \tau_K\}$  denote the precision parameters. The mixture density can then be written explicitly as

$$p(y_n | x_n, \boldsymbol{\beta}, \boldsymbol{\tau}) = \sum_k \pi_k(x_n) \mathcal{N}(y_n | x_n^\top \beta_k, \tau_k^{-1}), \quad \pi_k(x_n) = \frac{\exp\{x_n^\top \gamma_k\}}{\sum_{j=1}^K \exp\{x_n^\top \gamma_j\}} \quad (3)$$

We can simplify the form of the density by introducing the set of latent variables  $\mathbf{Z} = \{z_1, \dots, z_N\}$ , where  $z_n \in \mathbb{R}^K$  and  $z_{nk} = 1$  if and only if  $y_n$  belongs to the  $k$ -th cluster so that  $\sum_k z_{nk} = 1$ . Conditioning on this additional variable  $\mathbf{Z}$ , we have the following density,

$$p(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\tau}, \mathbf{Z}) = \prod_n \prod_k \mathcal{N}(y_n | x_n^\top \beta_k, \tau_k^{-1})^{z_{nk}} \quad (4)$$

$$p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\gamma}) = \prod_n \prod_k \pi_k(x_n)^{z_{nk}} = \prod_n \prod_k \left( \frac{\exp\{x_n^\top \gamma_k\}}{\sum_{j=1}^K \exp\{x_n^\top \gamma_j\}} \right)^{z_{nk}} \quad (5)$$

Next, we introduce the priors over the parameters  $\boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\gamma}$ , where we simplify the calculations by considering conjugate priors. For ease of computation, we consider an independent standard normal prior on  $\gamma_k$ 's, given by

$$p(\boldsymbol{\gamma}) = \prod_k p(\gamma_k) = \prod_k \mathcal{N}(\gamma_k | 0, \mathbf{I}_D) \quad (6)$$

For  $(\boldsymbol{\beta}, \boldsymbol{\tau})$ , we consider an independent normal-gamma prior, given by

$$p(\boldsymbol{\beta}, \boldsymbol{\tau}) = p(\boldsymbol{\beta} | \boldsymbol{\tau}) = \prod_k \prod_k \mathcal{N}(\beta_k | m_0, (\tau_k \Lambda_0)^{-1}) \text{Ga}(\tau_k | a_0, b_0) \quad (7)$$

Note that in the case where the mixture weights are covariate-independent, a Dirichlet prior is typically used for the mixing weights,  $\pi_1, \dots, \pi_K$ . In this case, however, the mixing weights are fully specified by  $\mathbf{X}$  and  $\boldsymbol{\gamma}$ , so we need only place a prior on  $\boldsymbol{\gamma}$ .

### 3. Variational Distribution

At this point, the variational parameters of interest are  $\boldsymbol{\theta} = (\mathbf{Z}, \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\gamma})$ . The log of the joint distribution of these random variables (written up to constants) is given by

$$\begin{aligned}
\ln p(\mathbf{y}, \mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\tau}, \mathbf{Z}) &= \ln \left\{ p(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\tau}, \mathbf{Z}, \boldsymbol{\gamma}) p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\gamma}) p(\boldsymbol{\gamma}) p(\boldsymbol{\beta}, \boldsymbol{\tau}) \right\} \\
&= \sum_n \sum_k z_{nk} \left\{ -\frac{1}{2} \ln(2\pi) + \frac{1}{2} \ln \tau_k - \frac{\tau_k}{2} (y_n - x_n^\top \boldsymbol{\beta}_k)^2 \right\} \\
&\quad + \sum_n \sum_k z_{nk} \left\{ x_n^\top \boldsymbol{\gamma}_k - \ln \sum_{j=1}^K \exp\{x_n^\top \boldsymbol{\gamma}_j\} \right\} + \sum_k \left\{ -\frac{D}{2} \ln(2\pi) - \frac{1}{2} \boldsymbol{\gamma}_k^\top \boldsymbol{\gamma}_k \right\} \\
&\quad + \sum_k \left\{ -\frac{D}{2} \ln(2\pi) + \frac{D}{2} \ln \tau_k + \ln |\Lambda_0| - \frac{\tau_k}{2} (\boldsymbol{\beta}_k - m_0)^\top \Lambda_0 (\boldsymbol{\beta}_k - m_0) \right\} \\
&\quad + \sum_k \left\{ (a_0 - 1) \ln \tau_k - b_0 \tau_k \right\}
\end{aligned} \tag{8}$$

We consider the following variational distribution used to approximate the posterior distribution of the parameters outlined previously.

$$q(\mathbf{Z}, \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\gamma}) = q(\mathbf{Z})q(\boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\gamma}) \tag{9}$$

#### 3.1 Coordinate Ascent Updates

As is standard in variational algorithms, we now seek the sequential updates of the factors in (9). For a particular variational parameter, the optimal distribution is found by taking the expectation of the joint distribution of all the random variables with respect to all of the *other* variational parameters, excluding the one of interest (**bishop:06**). Proceeding this way, we arrive at the following update equation for  $q(\mathbf{Z})$ ,

$$\ln q^*(\mathbf{Z}) = \mathbb{E}_{-q(\mathbf{Z})} \left[ \ln \{ p(\mathbf{y}, \mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\tau}, \mathbf{Z}, \boldsymbol{\gamma}) \} \right] \tag{10}$$

We adopt the convention that the expectation with respect to a negative subscript indicates an expectation taken with respect to the other variational parameters. Ignoring terms that are not functionally dependent on  $\mathbf{Z}$ , we can exponentiate both sides of (10) to obtain the optimal solution for  $q(\mathbf{Z})$ ,

$$q^*(\mathbf{Z}) = \prod_n \prod_k r_{nk}^{z_{nk}}, \quad r_{nk} = \frac{\rho_{nk}}{\sum_j \rho_{nj}} \tag{11}$$

For the details in this calculation, refer to Appendix A. The form of this discrete distribution gives us  $\mathbb{E}[z_{nk}] = r_{nk}$ . If we consider the quantity  $\ln \rho_{nk}$ , defined in (12) and discussed in more detail in Appendix A, we note that the exact computation of  $\ln \rho_{nk}$  involves four expectations taken with respect to the variational distribution  $q(\boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\gamma})$ .

$$\begin{aligned} \ln \rho_{nk} = & -\frac{1}{2} \ln(2\pi) + \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\tau})}[\ln \tau_k] - \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\beta}, \boldsymbol{\tau})}[\tau_k (y_n - x_n^\top \beta_k)^2] \\ & + x_n^\top \mathbb{E}_{q(\boldsymbol{\gamma})}[\gamma_k] - \mathbb{E}_{q(\boldsymbol{\gamma})} \left[ \ln \sum_j \exp\{x_n^\top \gamma_j\} \right] \end{aligned} \quad (12)$$

We briefly consider each of these expectations. Since we used conjugate families, we know that  $\mathbb{E}_{q(\boldsymbol{\tau})}[\ln \tau_k]$ ,  $\mathbb{E}_{q(\boldsymbol{\beta}, \boldsymbol{\tau})}[\tau_k (y_n - x_n^\top \beta_k)^2]$ , and  $\mathbb{E}_{q(\boldsymbol{\gamma})}[\gamma_k]$  will have closed form expressions. The remaining expectation,  $\mathbb{E}_{q(\boldsymbol{\gamma})}[\ln \sum_j \exp\{x_n^\top \gamma_j\}]$  presents a problem in that there lacks a closed form expression. Therefore, in order to complete this update, we use the following upper bound to approximate this quantity,

$$\mathbb{E}_{q(\boldsymbol{\gamma})} \left[ \ln \sum_j \exp\{x_n^\top \gamma_j\} \right] \approx \alpha_n + \varphi_n \quad (13)$$

where  $\varphi_n = \sum_j^K \frac{1}{2} (x_n^\top \mu_j - \alpha_n + \xi_{nj}) + \log(1 + e^{\xi_{nj}})$ . The details for the approximation and how to find  $\alpha_n$  and  $\xi_{nj}$  can be found in Appendix B. Using this in place of the problematic expectation in (12), we are able to obtain  $r_{nk}$  in closed form.

The remaining variational distribution can be found by considering the expectation of the joint density in (8) taken with respect to  $q(\mathbf{Z})$ , as derived above. The resulting variational distribution can be written up to constants as shown below,

$$\begin{aligned} \ln q^*(\boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\gamma}) = & \sum_k \sum_n \mathbb{E}_{q(\mathbf{Z})}[z_{nk}] \ln \mathcal{N}(y_n | x_n^\top \beta_k, \tau_k^{-1}) + \sum_k \ln p(\beta_k, \tau_k) \\ & + \sum_k \sum_n \mathbb{E}_{q(\mathbf{Z})}[z_{nk}] \left( x_n^\top \gamma_k - \ln \sum_{j=1}^K \exp\{x_n^\top \gamma_j\} \right) + \sum_k \ln \mathcal{N}(\gamma_k | 0, \mathbf{I}_D) \end{aligned} \quad (14)$$

Having just deriving the form for  $q(\mathbf{Z})$ , we can write out closed forms for the two expectation terms above. Noting in the summations above that the optimal distribution for  $q(\boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\gamma})$  has a sum involving only the  $\gamma_k$ 's in the second line and a sum involving only the  $(\beta_k, \tau_k)$ 's in the first line of (14), we can deduce  $q(\boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\gamma}) = \prod_k q(\beta_k, \tau_k) q(\gamma_k)$  and find the optimal distributions separately. With this factorization in mind, we obtain the following updates for the remaining variational parameters,

$$q^*(\beta_k | \tau_k) = \mathcal{N}(\beta_k | m_k, (\tau_k \mathbf{Q}_k)^{-1}) \quad (15)$$

$$q^*(\tau_k) = \text{Ga}(\tau_k | a_k, b_k) \quad (16)$$

$$q^*(\gamma_k) = \mathcal{N}(\gamma_k | \mu_k, \mathbf{V}_k^{-1}) \quad (17)$$

for  $k = 1, \dots, K$ . The derivation and parameter definitions for (15) and (16) can be found in Appendix C, and the details for (17) can be found in Appendix B. Now that we have

closed form updates for the variational distributions  $q(\mathbf{Z})$  and  $q(\boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\gamma})$ , we can cycle through a two-step update procedure. In the first step (variational E-step), we compute the expectation of each of the  $z_{nk}$ 's using the variational distributions (15), (16), and (17). Then in the second step (variational M-step), we derive new optimal distributions using results from the E-step. We alternate between these two steps until the variational lower bound converges, as discussed in the section 3.2.

### 3.2 Evidence Lower Bound (ELBO)

In order to evaluate the convergence of the coordinate ascent algorithm, we can calculate the evidence lower bound using the updated variational parameters at the end of each iteration. Since the ELBO is monotonic increasing, we continue the coordinate ascent until the change in the ELBO between iterations falls below a predetermined tolerance. Note that the expectations taken below are with respect to the optimal variational distributions defined in the previous section.

$$\begin{aligned}\mathcal{L}(q) &= \sum_{\mathbf{z}} \int \int \int q(\boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\gamma}, \mathbf{Z}) \ln \left\{ \frac{p(\mathbf{y}, \mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\gamma}, \mathbf{Z})}{q(\boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\gamma}, \mathbf{Z})} \right\} d\boldsymbol{\beta} d\boldsymbol{\tau} d\boldsymbol{\gamma} \\ &= \mathbb{E}[\ln p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\tau}, \mathbf{Z})] + \mathbb{E}[\ln p(\mathbf{Z} \mid \mathbf{X}, \boldsymbol{\gamma})] + \mathbb{E}[\ln p(\boldsymbol{\gamma})] + \mathbb{E}[\ln p(\boldsymbol{\beta}, \boldsymbol{\tau})] \\ &\quad - \mathbb{E}[\ln q(\mathbf{Z})] - \mathbb{E}[\ln q(\boldsymbol{\beta}, \boldsymbol{\tau})] - \mathbb{E}[\ln q(\boldsymbol{\gamma})]\end{aligned}\tag{18}$$

Details for each of the expectations can be found in Appendix D.

$$\begin{aligned}\mathbb{E}[\ln p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\tau}, \mathbf{Z})] &= -\frac{1}{2} \sum_n \sum_k r_{nk} \left\{ \ln(2\pi) - \psi(a_k) + \psi(b_k) + \frac{a_k}{b_k} (y_n - x_n^\top m_k)^2 + x_n^\top \mathbf{Q}_k^{-1} x_n \right\} \\ \mathbb{E}[\ln p(\mathbf{Z} \mid \mathbf{X}, \boldsymbol{\gamma})] &= \sum_n \sum_k r_{nk} (x_n^\top \mu_k - \alpha_n - \varphi_n) \\ \mathbb{E}[\ln p(\boldsymbol{\gamma})] &= -\frac{K \cdot D}{2} \ln(2\pi) - \frac{1}{2} \sum_k \mu_k^\top \mu_k \\ \mathbb{E}[\ln p(\boldsymbol{\beta}, \boldsymbol{\tau})] &= -K \left( \frac{D}{2} \ln(2\pi) - \ln |\Lambda_0| - a_0 \ln b_0 + \ln \Gamma(\alpha_0) \right) + \left( a_0 + \frac{D}{2} - 1 \right) \sum_k \psi(a_k) - \psi(b_k) \\ &\quad - \frac{1}{2} \sum_k \left\{ \frac{a_k}{b_k} \left[ (m_k - m_0)^\top \Lambda_0 (m_k - m_0) + b_0 \right] + \text{tr}(\Lambda_0 \mathbf{Q}_k^{-1}) \right\} \\ \mathbb{E}[\ln q(\mathbf{Z})] &= \sum_n \sum_k r_{nk} \ln r_{nk} \\ \mathbb{E}[\ln q(\boldsymbol{\beta}, \boldsymbol{\tau})] &= \sum_k \left( \frac{D}{2} + a_k - 1 \right) [\psi(a_k) - \psi(b_k)] + a_k (\ln b_k - 1) - \ln \Gamma(a_k) + \ln |\mathbf{Q}_k| \\ &\quad - \frac{KD}{2} (\ln(2\pi) + 1) \\ \mathbb{E}[\ln q(\boldsymbol{\gamma})] &= -\frac{KD}{2} (\ln(2\pi) + 1) + \sum_k \ln |\mathbf{V}_k|\end{aligned}$$

#### 4. Algorithm

Using the updates discussed in the previous section, we can formalize the variational algorithm below.

---

##### Algorithm 1 CAVI for Conditional Density Estimation

---

**Result:** An approximating distribution to the true posterior of  $\theta$

**Input:** Data  $y_{1:N}, x_{1:N}$ , number of components  $K$ , prior mean and precision for  $\beta_{1:K}$ , prior shape, rate parameters for precision parameters  $\tau_{1:K}$

**Output:** A variational density  $q(\mathbf{Z}, \beta, \tau, \gamma) = q(\mathbf{Z})q(\beta, \tau, \gamma) = q(\mathbf{Z}) \prod_k q(\beta_k, \tau_k)q(\gamma_k)$ ;

**Initialize:** Variational parameters  $\mathbf{m}_{1:K}, \mathbf{V}_{1:K}, \boldsymbol{\mu}_{1:K}, \mathbf{Q}_{1:K}, a_{1:K}, b_{1:K}, \xi_{1:N, 1:K}, \alpha_{1:N}$

**while** the ELBO has not converged **do**

**for**  $n \in \{1, \dots, N\}$  **do**

**for**  $k \in \{1, \dots, K\}$  **do**

      Set  $r_{nk} \propto \exp \left\{ -\frac{1}{2} \ln(2\pi) + \frac{1}{2} \mathbb{E}[\ln \tau_k] - \frac{1}{2} \mathbb{E}[\tau_k (y_n - x_n^\top \beta_k)^2] + x_n^\top \mathbb{E}[\gamma_k] \right. \\ \left. - \mathbb{E} \left[ \ln \left( \sum_{j=1}^K \exp \{ x_n^\top \gamma_j \} \right) \right] \right\}$

**end**

**end**

**for**  $n \in \{1, \dots, N\}$  **do**

**for**  $k \in \{1, \dots, K\}$  **do**

      Set  $\xi_{nk} \leftarrow \sqrt{(x_n^\top \mu_k - \alpha_n)^2 + x_n^\top \mathbf{Q}_k^{-1} x_n}$

**end**

    Set  $\alpha_n \leftarrow \frac{\frac{1}{2} \left( \frac{K}{2} - 1 \right) + \sum_k \lambda(\xi_{nk}) \mu_k^\top x_n}{\sum_k \lambda(\xi_{nk})}$

**end**

**for**  $k \in \{1, \dots, K\}$  **do**

    Set  $\mathbf{Q}_k \leftarrow \mathbf{I}_D + 2 \sum_n r_{nk} \lambda(\xi_{nk}) x_n x_n^\top$

    Set  $\eta_k \leftarrow \sum_n r_{nk} \left[ \frac{1}{2} + 2 \lambda(\xi_{nk}) \alpha_n \right] x_n$

    Set  $\mu_k \leftarrow \mathbf{Q}_k^{-1} \eta_k$

    Set  $\mathbf{V}_k \leftarrow \sum_n r_{nk} x_n x_n^\top + \Lambda_0$

    Set  $\zeta_k \leftarrow \sum_n r_{nk} y_n x_n + \Lambda_0 m_0$

    Set  $m_k \leftarrow \mathbf{V}_k^{-1} \zeta_k$

    Set  $a_k \leftarrow a_0 + N_k$

    Set  $b_k \leftarrow b_0 + \frac{1}{2} \left[ \sum_n r_{nk} y_n^2 + m_0^\top \Lambda_0 m_0 - \zeta_k^\top \mathbf{V}_k^{-1} \zeta_k \right]$

**end**

  Compute ELBO using updated parameters

**end**

**return**  $q(\mathbf{Z}, \beta, \tau, \gamma)$

---

## 5. Conclusion

In this project, we focused on one primary approximation to the difficult-to-compute expectation of the log sum of exponentials,  $\mathbb{E}_{q(\gamma)}[\ln \sum_j \exp\{x_n^\top \gamma_j\}]$ . Alternatively, we also considered a simpler bound in section ??, which only requires us to compute moment generating functions. As has been noted in previous work and simulations, the former bound provides better approximations when the variance of  $q(\gamma)$  is extremely large because the approximation is asymptotically optimal (**bouchard:07**), but in most other cases, the performance deteriorated (**Depraetere:17**). Although both of these methods require additional variational parameters, this increase in model complexity is offset by the faster inference that variational Bayesian methods provide. As a result, even though each iteration involves  $2K$  ( $D \times D$ )-matrix inversions to update the variational distributions for  $\beta_k$  and  $\gamma_k$ , we expect the number of iterations in the coordinate ascent algorithm to be far fewer than if we use a traditional Gibbs sampling scheme.

While we only considered two candidate approximations in this project, there are numerous ways to approach the problem. Although in a slightly different context, Depraetere and Vandebroek (2017) provide several approximations to the same quantity of interest, all of which admit closed form updates so that variational Bayesian methods are tractable.

## 6. Variable Selection for Gaussian Component

### 6.1 Notation and Prior Specification

Recall the previous setup, where we denote  $\mathbf{y} = \{y_1, \dots, y_N\}$ ,  $\mathbf{X} = \{x_1, \dots, x_N\}$ , where each  $x_n \in \mathbb{R}^D$ . Let  $\boldsymbol{\beta} = \{\beta_1, \dots, \beta_K\}$  and  $\boldsymbol{\gamma} = \{\gamma_1, \dots, \gamma_K\}$  denote the  $D$ -dimensional coefficient vectors in the mixture weights and in gaussian mixture components, respectively. Finally, let  $\boldsymbol{\tau} = \{\tau_1, \dots, \tau_K\}$  denote the precision parameters for each of the gaussian components in the mixture density. The mixture density can then be written explicitly as

$$p(y_n | x_n, \boldsymbol{\beta}, \boldsymbol{\tau}) = \sum_k \pi_k(x_n) \mathcal{N}(y_n | x_n^\top \beta_k, \tau_k^{-1}), \quad \pi_k(x_n) = \frac{\exp\{x_n^\top \gamma_k\}}{\sum_{j=1}^K \exp\{x_n^\top \gamma_j\}} \quad (19)$$

Conditioning on the same set of latent variables,  $\mathbf{Z} = \{z_1, \dots, z_N\}$ , used indicate the membership for each of the response variables, we obtain the mixture density in product form.

$$p(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\tau}, \mathbf{Z}) = \prod_n \prod_k \mathcal{N}(y_n | x_n^\top \beta_k, \tau_k^{-1})^{z_{nk}} \quad (20)$$

$$p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\gamma}) = \prod_n \prod_k \pi_k(x_n)^{z_{nk}} = \prod_n \prod_k \left( \frac{\exp\{x_n^\top \gamma_k\}}{\sum_{j=1}^K \exp\{x_n^\top \gamma_j\}} \right)^{z_{nk}} \quad (21)$$

For  $\boldsymbol{\gamma}$ , we consider an independent standard normal prior on the  $\gamma_k$ 's, given by

$$p(\boldsymbol{\gamma}) = \prod_k p(\gamma_k) = \prod_k \mathcal{N}(\gamma_k | 0, \mathbf{I}_D) \quad (22)$$

Instead of the normal-gamma prior that we used in the standard setup, we place a gamma prior on  $\tau_k$ 's:

$$p(\boldsymbol{\tau}) = \prod_k p(\tau_k) = \prod_k \text{Ga}(\tau_k | a_0, b_0) \quad (23)$$

Finally, to incorporate variable selection we consider the *rows* of  $\boldsymbol{\beta}$ , as follows:

$$\boldsymbol{\beta} = \begin{bmatrix} | & | & & | \\ \beta_1 & \beta_2 & \dots & \beta_K \\ | & | & & | \end{bmatrix} = \begin{bmatrix} - & \tilde{\beta}_1 & - \\ \vdots & & \\ - & \tilde{\beta}_D & - \end{bmatrix}$$

Using the above formulation,  $\tilde{\beta}_d = (\beta_{1d}, \beta_{2d}, \dots, \beta_{Kd})^\top \in \mathbb{R}^K$  represents the  $d$ -th row of the coefficient matrix. In other words,  $\tilde{\beta}_d$  represents the coefficient vector for the  $d$ -th predictor across all  $K$  clusters. We place independent spike and slab priors on each *row* vectors of  $\boldsymbol{\beta}$ ,

$$\tilde{\beta}_d \sim \pi \mathcal{N}(\tilde{\beta}_d | 0, \xi_0^{-1} \cdot \mathbf{I}_K) + (1 - \pi) \delta_0(\tilde{\beta}_d) \quad (24)$$

for  $\tilde{\beta}_d$  for  $d = 1, \dots, D$ .

Denote  $\boldsymbol{\omega} = \{\omega_1, \dots, \omega_D\}$ , where  $\omega_d \stackrel{i.i.d}{\sim} \text{Ber}(\omega_d | \pi)$ , we have the following reparametrization of the spike and slab prior on  $\boldsymbol{\beta}$ :

$$p(\boldsymbol{\beta}, \boldsymbol{\omega}) = \prod_d p(\tilde{\beta}_d, \omega_d) = \prod_d \mathcal{N}(\tilde{\beta}_d | 0, \xi_0^{-1} \cdot \mathbf{I}_K)^{\omega_d} \pi^{\omega_d} (1 - \pi)^{1 - \omega_d} \quad (25)$$



## 6.2 Joint Likelihood

Using the priors outlined in the previous section, we have the following expression for the joint likelihood.

$$\begin{aligned}
\ln p(\mathbf{y}, \mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\tau}, \mathbf{Z}) &= \ln \left\{ p(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\tau}, \mathbf{Z}, \boldsymbol{\gamma}) p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\gamma}) p(\boldsymbol{\gamma}) p(\tilde{\boldsymbol{\beta}}_d, \omega_d) p(\boldsymbol{\tau}) \right\} \\
&= \sum_n \sum_k z_{nk} \left\{ -\frac{1}{2} \ln(2\pi) + \frac{1}{2} \ln \tau_k - \frac{\tau_k}{2} (y_n - x_n^\top \boldsymbol{\beta}_k)^2 \right\} \\
&\quad + \sum_n \sum_k z_{nk} \left\{ x_n^\top \boldsymbol{\gamma}_k - \ln \sum_{j=1}^K \exp\{x_n^\top \boldsymbol{\gamma}_j\} \right\} + \sum_k \left\{ -\frac{D}{2} \ln(2\pi) - \frac{1}{2} \boldsymbol{\gamma}_k^\top \boldsymbol{\gamma}_k \right\} \\
&\quad + \sum_d \omega_d \left\{ -\frac{K}{2} \ln(2\pi) + \frac{K}{2} \ln \xi_0 - \frac{\xi_0}{2} \tilde{\boldsymbol{\beta}}_d^\top \tilde{\boldsymbol{\beta}}_d \right\} + \omega_d \ln \pi + (1 - \omega_d) \ln(1 - \pi) \\
&\quad + \sum_k \left\{ (a_0 - 1) \ln \tau_k - b_0 \tau_k \right\}
\end{aligned} \tag{26}$$

## 6.3 Approximating Distribution

The variational parameters of interest are  $\boldsymbol{\theta} = (\mathbf{Z}, \boldsymbol{\tau}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\omega})$ . We consider an approximating distribution of the following form

$$q(\boldsymbol{\theta}) = q(\mathbf{Z}, \boldsymbol{\tau}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\omega}) = q(\mathbf{Z}) \cdot q(\boldsymbol{\tau}, \boldsymbol{\gamma}) \cdot \prod_d q(\tilde{\boldsymbol{\beta}}_d, \omega_d) \tag{27}$$

In fact, it follows easily (without assuming anything about the functional forms) that the right hand side of (27) has the following induced factorization,

$$q(\mathbf{Z}, \boldsymbol{\tau}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\omega}) = \prod_n \prod_k q(z_{nk}) \cdot \prod_k q(\tau_k) \cdot \prod_k q(\boldsymbol{\gamma}_k) \cdot \prod_d q(\tilde{\boldsymbol{\beta}}_d, \omega_d) \tag{28}$$

From our choice of conjugate families in the prior distributions, each of variational parameters, with the exception of  $q(\boldsymbol{\gamma})$ , can be updated in closed form. Proceeding in standard fashion by taking the expectation with respect to the variational distribution, we can derive the functional form for each of the approximating distributions and their corresponding coordinate ascent updates.

## 6.4 Coordinate Ascent Algorithm: Component-wise VB

While the updates for  $\{\mathbf{Z}, \boldsymbol{\tau}, \boldsymbol{\gamma}\}$  are performed in a similar manner as the standard algorithm, the variable selection part of the algorithm iterates over  $D$  dimensions, updating the features sequentially.

**Update  $q(z_{nk})$ .** As shown in the previous case, we know the optimal distribution for  $\mathbf{Z}$  is of the form

$$q(z_{nk}) = r_{nk}^{z_{nk}}, \quad r_{nk} = \frac{\rho_{nk}}{\sum_j \rho_{nj}}$$

Note that  $\rho_{nk}$  involves evaluating a number of expectations taken with respect to various variational distribution, each of which is updated sequentially during each iteration of the coordinate ascent algorithm. In particular,

$$\begin{aligned} \ln \rho_{nk} = & -\frac{1}{2} \ln(2\pi) + \frac{1}{2} \mathbb{E}_{q(\tau)}[\ln \tau_k] - \frac{1}{2} \mathbb{E}_{q(\tau)}[\tau_k] \cdot \mathbb{E}_{q(\beta, \omega)}[(y_n - x_n^\top \beta_k)^2] \\ & + x_n^\top \mathbb{E}_{q(\gamma)}[\gamma_k] - \mathbb{E}_{q(\gamma)} \left[ \ln \sum_j^K \exp\{x_n^\top \gamma_j\} \right] \end{aligned} \quad (29)$$

so that  $\mathbb{E}[z_{nk}] = r_{nk}$ . Explicit expressions for each of the expectations, which can either be computed in closed form or approximated, can be found in the supplemental section.

**Update**  $q(\tau_k)$ . For  $k \in \{1, \dots, K\}$ ,  $q(\tau_k) = \text{Ga}(\tau_k \mid a_k, b_k)$ . The shape and rate parameters can be updated with

$$a_k = a_0 + \frac{N_k}{2} - 1 \quad (30)$$

$$b_k = b_0 + \frac{1}{2} \sum_n (y_n - x_n^\top \tilde{m}_k)^2 + x_n^\top \text{Var}(\beta_k) x_n \quad (31)$$

where  $\tilde{m}_k = \mathbb{E}(\beta_k)$  and  $\text{Var}(\beta_k) =$ .

**Update**  $q(\gamma_k)$ . The update for  $\gamma$  remains the same as before:

$$q(\gamma_k) = \mathcal{N}(\gamma_k \mid \mu_k, V_k^{-1}) \quad (32)$$

**Update**  $q(\tilde{\beta}_d, \omega_d)$ . We do this sequentially by first deriving the update for  $q(\tilde{\beta}_d \mid \omega_d = 1)$  and then for  $q(\omega_d)$ . It is helpful to rewrite the first term in the log joint likelihood (mixture density) in terms of  $\tilde{\beta}_d$ . Conditioning on  $\omega_d = 1$ , we have the variational distribution with equality written up to constants (in  $\tilde{\beta}_d$ ),

$$\ln q(\tilde{\beta}_d \mid \omega_d = 1) = -\frac{1}{2} \tilde{\beta}_d^\top U_d \tilde{\beta}_d + \tilde{\beta}_d^\top \eta_d - \frac{\xi_0}{2} \tilde{\beta}_d^\top \tilde{\beta}_d \quad (33)$$

$$= -\frac{1}{2} \tilde{\beta}_d^\top [U_d + \xi_0 \mathbf{I}_K] \tilde{\beta}_d + \tilde{\beta}_d^\top \eta_d \quad (34)$$

Note that we have used the following result (equality written up to constants in  $\tilde{\beta}_d$ ).

$$E_{-q(\beta)} \left[ \frac{1}{2} \sum_n \sum_k -z_{nk} \tau_k (y_n - x_n^\top \beta_k)^2 \right] = -\frac{1}{2} \sum_d \tilde{\beta}_d^\top U_d \tilde{\beta}_d + \sum_d \tilde{\beta}_d^\top \eta_d \quad (35)$$

From (34), we see that

$$q(\tilde{\beta}_d) = \mathcal{N}(\tilde{\beta}_d \mid m_d, Q_d^{-1})$$

where  $m_d = \mathbf{Q}_d^{-1} \eta_d$ ,  $\mathbf{Q}_d = \mathbf{U}_d + \xi_0 \mathbf{I}_K$ , and  $\eta_d = \zeta_d - \frac{1}{2} \sum_{j \neq d}^D \mathbf{R}_{dj} \omega_j m_j$ . We have made use of the following matrices

$$\mathbf{R}_{dj} = \begin{bmatrix} \frac{a_1}{b_1} \sum_n r_{n1} x_{nd} x_{nj} & & \\ & \ddots & \\ & & \frac{a_K}{b_K} \sum_n r_{nK} x_{nd} x_{nj} \end{bmatrix} \quad \zeta_d = \begin{bmatrix} \frac{a_1}{b_1} \sum_n r_{n1} x_{nd} y_n \\ \vdots \\ \frac{a_K}{b_K} \sum_n r_{nK} x_{nd} y_n \end{bmatrix}$$

$$\mathbf{U}_d = \begin{bmatrix} \frac{a_1}{b_1} \sum_n r_{n1} x_{nd}^2 & & \\ & \ddots & \\ & & \frac{a_K}{b_K} \sum_n r_{nK} x_{nd}^2 \end{bmatrix}$$

From (34), we have

$$q^*(\tilde{\beta}_d) = \mathcal{N}(\tilde{\beta}_d \mid m_d, \mathbf{Q}_d^{-1})$$

where  $m_d = \mathbf{Q}_d^{-1} \eta_d$ ,  $\mathbf{Q}_d = \mathbf{U}_d + \xi_0 \mathbf{I}_K$ , and  $\eta_d = \zeta_d - \frac{1}{2} \sum_{j \neq d}^D \mathbf{R}_{dj} \omega_j m_j$ . We have made use of the following matrices

$$\mathbf{U}_d = \text{diag} \left( \frac{a_1}{b_1} \sum_n r_{n1} x_{nd} x_{nj}, \dots, \frac{a_K}{b_K} \sum_n r_{nK} x_{nd} x_{nj} \right) \quad (36)$$

$$\mathbf{R}_{dj} = \text{diag} \left( \frac{a_1}{b_1} \sum_n r_{n1} x_{nd} x_{nj}, \dots, \frac{a_K}{b_K} \sum_n r_{nK} x_{nd} x_{nj} \right) \quad (37)$$

$$\zeta_d = \left( \frac{a_1}{b_1} \sum_n r_{n1} x_{nd} y_n, \dots, \frac{a_K}{b_K} \sum_n r_{nK} x_{nd} y_n \right)^\top \quad (38)$$

Note that we restricted the approximating distribution for  $p(\boldsymbol{\beta}, \boldsymbol{\omega})$  to be of the form

$$q(\boldsymbol{\beta}, \boldsymbol{\omega}) = \prod_d q(\tilde{\beta}_d, \omega_d)$$

where each of the  $d$  factors then inherits a spike and slab density of the form

$$q(\tilde{\beta}_d, \omega_d) = \begin{cases} \lambda_d \mathcal{N}(\tilde{\beta}_d \mid m_d, \mathbf{Q}_d^{-1}), & \text{if } \omega_d = 1 \\ (1 - \lambda_d) \delta_0(\tilde{\beta}_d) & \text{if } \omega_d = 0 \end{cases}$$

The variational parameters  $(m_d, \mathbf{Q}_d, \lambda_d)$  can be updated using the following by differentiating the variational lower bound with respect to each of the parameters, setting the resulting partial derivatives to zero and solving for each parameter. Doing so yields following:

$$\text{Cov}(\tilde{\beta}_d \mid \omega_d = 1) \approx \mathbf{Q}_d^{-1} = (\mathbf{U}_d + \xi_0 \mathbf{I}_K)^{-1} \quad (39)$$

$$\mathbb{E}[\tilde{\beta}_d \mid \omega_d = 1] \approx m_d = \mathbf{Q}_d^{-1} \eta_d \quad (40)$$

$$\frac{p(\omega_d = 1 \mid \mathbf{X}, \mathbf{y}, \boldsymbol{\theta})}{p(\omega_d = 0 \mid \mathbf{X}, \mathbf{y}, \boldsymbol{\theta})} \approx \frac{\lambda_d}{1 - \lambda_d} = \frac{\pi}{1 - \pi} \cdot \log \xi_0^{\frac{K}{2}} \cdot \exp \left\{ \frac{1}{2} m_d^\top \eta_d \right\} \quad (41)$$

## Appendix A. Variational Update for $q(\mathbf{Z})$

Taking the expectation with respect to the other variational parameters, we can derive the following variational distribution for  $\mathbf{Z}$ ,

$$\begin{aligned}\ln q^*(\mathbf{Z}) &= \sum_n \sum_k z_{nk} \left\{ -\frac{1}{2} \ln(2\pi) + \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\tau})}[\ln \tau_k] - \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\beta}, \boldsymbol{\tau})}[\tau_k (y_n - x_n^\top \beta_k)^2] \right. \\ &\quad \left. + x_n^\top \mathbb{E}_{q(\boldsymbol{\gamma})}[\boldsymbol{\gamma}_k] - \mathbb{E}_{q(\boldsymbol{\gamma})} \left[ \ln \sum_j \exp\{x_n^\top \boldsymbol{\gamma}_j\} \right] \right\} \\ &= \sum_n \sum_k z_{nk} \ln \rho_{nk}\end{aligned}$$

where we have defined

$$\begin{aligned}\ln \rho_{nk} &= -\frac{1}{2} \ln(2\pi) + \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\tau})}[\ln \tau_k] - \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\beta}, \boldsymbol{\tau})}[\tau_k (y_n - x_n^\top \beta_k)^2] \\ &\quad + x_n^\top \mathbb{E}_{q(\boldsymbol{\gamma})}[\boldsymbol{\gamma}_k] - \mathbb{E}_{q(\boldsymbol{\gamma})} \left[ \ln \sum_j \exp\{x_n^\top \boldsymbol{\gamma}_j\} \right]\end{aligned}\tag{42}$$

Exponentiating and normalizing, we have

$$q^*(\mathbf{Z}) = \prod_n \prod_k r_{nk}^{z_{nk}}, \quad r_{nk} = \frac{\rho_{nk}}{\sum_j \rho_{nj}}\tag{43}$$

For the discrete distribution  $q^*(\mathbf{Z})$  given in (43) above, we have  $E[z_{nk}] = r_{nk}$ . Note, however, that in order to compute the expectation in closed form, we need an expression for the four expectations involved in the quantity  $\ln \rho_{nk}$ , as defined in (42).

From the results derived in Appendix C, we know that  $q^*(\tau_k) = \text{Ga}(\tau_k \mid a_k, b_k)$ . We can then compute the following expectation with respect to  $q^*(\boldsymbol{\tau})$ .

$$E_{q(\boldsymbol{\tau})}[\ln \tau_k] = \psi(a_k) - \psi(b_k)\tag{44}$$

Again from Appendix C, we can then compute the following expectation with respect to  $q^*(\beta_k, \tau_k)$ .

$$\begin{aligned}\mathbb{E}_{q(\boldsymbol{\beta}, \boldsymbol{\tau})}[\tau_k (y_n - x_n^\top \beta_k)^2] &= \mathbb{E} \left[ \tau_k \left( y_n + m_k^\top x_n x_n^\top m_k + \text{tr} \left( x_n x_n^\top (\tau_k \mathbf{Q}_k)^{-1} \right) - 2y_n x_n^\top m_k \right) \right] \\ &= \frac{a_k}{b_k} (y_n^2 + m_k^\top x_n x_n^\top m_k - 2y_n x_n^\top m_k) + \text{tr} (x_n x_n^\top \mathbf{Q}_k^{-1}) \\ &= \frac{a_k}{b_k} (y_n - m_k^\top x_n)^2 + x_n^\top \mathbf{Q}_k^{-1} x_n\end{aligned}\tag{45}$$

From the expression derived in (50) of Appendix B, we have  $q^*(\gamma_k) = \mathcal{N}(\gamma_k | \mu_k, V_k^{-1})$ , then we have

$$\mathbb{E}_{q(\gamma_k)}[\gamma_k] = \mu_k \quad (46)$$

Using the bound discussed in Appendix B, equation (49), we can then compute the following expectation with respect to  $q^*(\gamma)$ .

$$\begin{aligned} & \mathbb{E}_{q(\gamma)} \left[ \ln \sum_j^K \exp\{x_n^\top \gamma_j\} \right] \\ & \approx \mathbb{E}_{q(\gamma)} \left[ \alpha_n + \sum_{j=1}^K \frac{x_n^\top \gamma_j - \alpha_n - \xi_{nj}}{2} + \lambda(\xi_{nj}) \left( (x_n^\top \gamma_j - \alpha_n)^2 - \xi_{nj}^2 \right) + \log(1 + e^{\xi_{nj}}) \right] \\ & = \alpha_n + \sum_j^K \frac{1}{2} (x_n^\top \mu_j - \alpha_n - \xi_{nj}) + \lambda(\xi_{nj}) \left( (x_n^\top \mu_j - \alpha_n)^2 - \xi_{nj}^2 + x_n^\top V_j^{-1} x_n \right) + \log(1 + e^{\xi_{nj}}) \\ & = \alpha_n + \sum_j^K \frac{1}{2} (x_n^\top \mu_j - \alpha_n - \xi_{nj}) + \log(1 + e^{\xi_{nj}}) \end{aligned} \quad (47)$$

since  $(x_n^\top \mu_j - \alpha_n)^2 - \xi_{nj}^2 = -x_n^\top V_j^{-1} x_n$ . Gathering the results in (44), (45), (46), and (47), and substituting these into (42), we can compute  $\mathbb{E}[z_{nk}] = r_{nk}$  in closed form.

## Appendix B. Variational Updates for $q(\gamma_k)$

For the variational distribution for  $\gamma_k$ , we first note the following bound (**bouchard:07**),  $\sum_{j=1}^K e^{t_j} \leq \prod_{j=1}^K (1 + e^{t_j})$ . Setting  $t_j = x_n^\top \gamma_j - \alpha_n$  and then taking log, we have the following bound:

$$\log \left( \sum_{j=1}^K \exp\{x_n^\top \gamma_j\} \right) \leq \alpha_n + \sum_{j=1}^K \log(1 + \exp\{x_n^\top \gamma_j - \alpha_n\}) \quad (48)$$

We can further bound this by using the following tangential bound (**jj:2001**),

$$\log(1 + e^x) \leq \frac{x - \xi}{2} + \frac{1}{4\xi} \tanh\left(\frac{\xi}{2}\right) (x^2 - \xi^2) + \log(1 + e^\xi)$$

then we arrive at the following bound:

$$\log \left( \sum_{j=1}^K \exp\{x_n^\top \gamma_j\} \right) \leq \alpha_n + \sum_{j=1}^K \frac{x_n^\top \gamma_j - \alpha_n - \xi_{nj}}{2} + \lambda(\xi_{nj}) \left( (x_n^\top \gamma_j - \alpha_n)^2 - \xi_{nj}^2 \right) + \log(1 + e^{\xi_{nj}}) \quad (49)$$

where  $\lambda(\xi) = \frac{1}{4\xi} \tanh\left(\frac{\xi}{2}\right)$ . Then we can substitute this back into  $\ln q^*(\gamma_k)$  to obtain an approximation for the left hand side of (48), thus allowing us to obtain a closed form for the

variational distribution. Note that all of the equalities above are written up to constants.

$$\begin{aligned}
\ln q^*(\gamma_k) &= -\frac{1}{2}\gamma_k^\top \gamma_k + \sum_n r_{nk} x_n^\top \gamma_k - \sum_n r_{nk} \ln \left( \sum_j \exp\{x_n^\top \gamma_j\} \right) \\
&\approx -\frac{1}{2}\gamma_k^\top \gamma_k + \gamma_k^\top \sum_n r_{nk} x_n \\
&\quad - \sum_n r_{nk} \left\{ \alpha_n + \sum_{j=1}^K \frac{x_n^\top \gamma_j - \alpha_n - \xi_{nj}}{2} + \lambda(\xi_{nj}) ((x_n^\top \gamma_j - \alpha_n)^2 - \xi_{nj}^2) + \log(1 + e^{\xi_{nj}}) \right\} \\
&= -\frac{1}{2}\gamma_k^\top \gamma_k + \gamma_k^\top \sum_n r_{nk} x_n - \sum_n r_{nk} \left\{ \frac{1}{2}\gamma_k^\top x_n + \lambda(\xi_{nj}) (\gamma_k^\top x_n x_n^\top \gamma_k - 2\alpha_n \gamma_k^\top x_n) \right\} \\
&= -\frac{1}{2}\gamma_k^\top \left( \mathbf{I}_D + 2 \sum_n r_{nk} \lambda(\xi_{nk}) x_n x_n^\top \right) \gamma_k + \gamma_k' \left( \sum_n r_{nk} \left( \frac{1}{2} + 2\lambda(\xi_{nk}) \alpha_n x_n \right) \right)
\end{aligned}$$

Exponentiating, we can recover  $q^*(\gamma_k) = \mathcal{N}(\gamma_k | \mu_k, \mathbf{V}_k^{-1})$ , where

$$\begin{aligned}
\mu_k &= \mathbf{V}_k^{-1} \eta_k \\
\eta_k &= \sum_n r_{nk} \left( \frac{1}{2} + 2\lambda(\xi_{nj}) \alpha_n \right) x_n \\
\mathbf{V}_k &= \mathbf{I}_D + 2 \sum_n r_{nk} \lambda(\xi_{nk}) x_n x_n^\top
\end{aligned} \tag{50}$$

The additional parameters introduced in the two upper bounds can be updated using the following equations (**Depraetere:17**),

$$\xi_{nk} = \sqrt{(\mu_k^\top x_n - \alpha_n)^2 + x_n^\top \mathbf{V}_k^{-1} x_n} \quad \forall k, n$$

$$\alpha_n = \frac{\frac{1}{2} \left( \frac{K}{2} - 1 \right) + \sum_{j=1}^K \lambda(\xi_{nj}) \mu_j^\top x_n}{\sum_{j=1}^K \lambda(\xi_{nj})} \quad \forall n$$

### Appendix C. Variational Updates for $q(\beta_k, \tau_k)$

Using results from Appendix A, we can write the following expression for the joint variational distribution of  $(\beta_k, \tau_k)$ ,

$$\begin{aligned}
\ln q^*(\beta_k, \tau_k) &= \sum_n -\frac{1}{2} r_{nk} \tau_k (y_n^2 + \beta_k^\top x_n x_n^\top \beta_k - 2y_n \beta_k^\top x_n) + \frac{r_{nk}}{2} \ln \tau_k + \frac{D}{2} \ln \tau_k \\
&\quad - \frac{\tau_k}{2} (\beta_k^\top \Lambda_0 \beta_k + m_0^\top \Lambda_0 m_0 - 2\beta_k^\top \Lambda_0 m_0) + (a_0 - 1) \ln \tau_k - b_0 \tau_k
\end{aligned} \tag{51}$$

We first consider terms on the right hand side of (51) that depend on  $\beta_k$  to find  $\ln q^*(\beta_k | \tau_k)$ , giving

$$\ln q^*(\beta_k | \tau_k) = -\frac{\tau_k}{2} \beta_k^\top \left[ \sum_n r_{nk} x_n x_n^\top + \Lambda_0 \right] \beta_k + \tau_k \beta_k^\top \left[ \sum_n r_{nk} y_n x_n + \Lambda_0 m_0 \right] \quad (52)$$

$$q^*(\beta_k | \tau_k) = \mathcal{N}(\beta_k | m_k, (\tau_k \mathbf{Q}_k)^{-1}) \quad (53)$$

$$\mathbf{Q}_k = \sum_n r_{nk} x_n x_n^\top + \Lambda_0$$

$$\zeta_k = \sum_n r_{nk} y_n x_n + \Lambda_0 m_0 \quad (54)$$

$$m_k = \mathbf{Q}_k^{-1} \zeta_k$$

Then we can make use of the relation  $\ln q^*(\tau_k) = \ln q^*(\beta_k, \tau_k) - \ln q^*(\beta_k | \tau_k)$ , where the quantities on the right hand side come from (51) and (53). Note that equality below is written up to constants, keeping only terms involving  $\tau_k$ .

$$\begin{aligned} \ln q^*(\tau_k) = (a_0 + N_k - 1) \ln \tau_k - \tau_k & \left\{ b_0 + \frac{1}{2} \left( \sum_n r_{nk} y_n^2 + m_0^\top \Lambda_0 m_0 - m_k^\top \mathbf{Q}_k m_k \right) \right. \\ & + \frac{1}{2} \beta_k^\top \left( \sum_n r_{nk} x_n x_n^\top + \Lambda_0 - \mathbf{Q}_k \right) \beta_k \\ & \left. - 2 \beta_k^\top \left( \sum_n r_{nk} y_n x_n + \Lambda_0 m_0 - \mathbf{Q}_k m_k \right) \right\} \end{aligned} \quad (55)$$

Exponentiating, we arrive at the following distribution

$$q^*(\tau_k) = \text{Ga}(\tau_k | a_k, b_k) \quad (56)$$

where we have defined

$$\begin{aligned} a_k &= a_0 + \frac{N_k}{2} \\ b_k &= b_0 + \frac{1}{2} \sum_n r_{nk} y_n^2 + m_0^\top \Lambda_0 m_0 - \zeta_k^\top \mathbf{Q}_k^{-1} \zeta_k \end{aligned} \quad (57)$$

The expression for  $b_k$  arises by noting that the three following simplifications for the summation terms in the coefficient of  $\tau_k$  in (55),

$$\begin{aligned} \sum_n r_{nk} y_n^2 + m_0^\top \Lambda_0 m_0 - m_k^\top \mathbf{Q}_k m_k &= \sum_n r_{nk} y_n^2 + m_0^\top \Lambda_0 m_0 - \zeta_k^\top \mathbf{Q}_k^{-1} \zeta_k \\ \sum_n r_{nk} x_n x_n^\top + \Lambda_0 - \mathbf{Q}_k &= 0 \\ \sum_n r_{nk} y_n x_n + \Lambda_0 m_0 - \mathbf{Q}_k m_k &= 0 \end{aligned}$$

where the first equality holds by expanding  $m_k^\top Q_k m_k = \zeta_k^\top (Q_k^{-1})^\top Q_k Q_k^{-1} \zeta_k = b_k^\top Q_k^{-1} \zeta_k$ . The second quality holds by recalling the definition of  $Q_k$  in (54), and the third equality holds by observing from (54) that  $Q_k m_k = \zeta_k$

## Appendix D. Variational Lower Bound

As seen in (18) of section 3.2, we need to calculate seven expectations (taken with respect to the variational distribution  $q(\mathbf{Z}, \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\gamma})$ ). Below, we compute each of these expectations in detail, making extensive use of the variational distributions derived in Appendix A, B, and C.

$$\begin{aligned}
\mathbb{E}[\ln p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\tau}, \mathbf{Z})] &= \sum_n \sum_k \mathbb{E}[z_{nk}] \mathbb{E} \left[ \ln \mathcal{N}(y_n \mid x_n^\top \beta_k, \tau_k^{-1}) \right] \\
&= -\frac{1}{2} \sum_n \sum_k r_{nk} \left\{ \ln(2\pi) - \mathbb{E}[\ln \tau_k] + \mathbb{E} \left[ \tau_k (y_n - x_n^\top \beta_k)^2 \right] \right\} \\
&= -\frac{1}{2} \sum_n \sum_k r_{nk} \left\{ \ln(2\pi) - (\psi(a_k) - \psi(b_k)) + x_n^\top Q_k^{-1} x_n \right. \\
&\quad \left. + \frac{a_k}{b_k} (y_n - x_n^\top m_k)^2 \right\}
\end{aligned} \tag{58}$$

$$\begin{aligned}
\mathbb{E}[\ln p(\mathbf{Z} \mid \mathbf{X}, \boldsymbol{\gamma})] &= \sum_n \sum_k \mathbb{E}[z_{nk}] \left( \mathbb{E}[x_n^\top \gamma_k] - \mathbb{E} \left[ \ln \sum_j \exp \{ x_n^\top \gamma_j \} \right] \right) \\
&\approx \sum_n \sum_k r_{nk} (x_n^\top \mu_k - \alpha_n - \varphi_n)
\end{aligned} \tag{59}$$

where  $\varphi_n = \sum_{j=1}^K \frac{1}{2} (x_n^\top \mu_j - \alpha_n - \xi_{nj}) + \log(1 + e^{\xi_{nj}})$ . Here, we make use of the result in (47), where we take the expectation of the upper bound previously derived.

$$\begin{aligned}
\mathbb{E}[\ln p(\boldsymbol{\gamma})] &= \sum_k \mathbb{E} \left[ \ln \mathcal{N}(\gamma_k \mid 0, \mathbf{I}_D) \right] \\
&= -\frac{K \cdot D}{2} \ln(2\pi) - \frac{1}{2} \sum_k \mu_k^\top \mu_k + \text{tr}(\mathbf{V}_k^{-1})
\end{aligned} \tag{60}$$



$$\begin{aligned}
\mathbb{E}[\ln p(\boldsymbol{\beta}, \boldsymbol{\tau})] &= \sum_k \mathbb{E} \left[ \ln \mathcal{N}(\beta_k \mid m_0, (\tau_k \Lambda_0)^{-1}) \right] + \mathbb{E} \left[ \ln \text{Ga}(\tau_k \mid a_0, b_0) \right] \\
&= \left( a_0 + \frac{D}{2} - 1 \right) \sum_k \psi(a_k) - \psi(b_k) \\
&\quad - \frac{1}{2} \sum_k \left\{ \frac{a_k}{b_k} \left[ (m_k - m_0)^\top \Lambda_0 (m_k - m_0) + b_0 \right] + \text{tr}(\Lambda_0 \mathbf{Q}_k^{-1}) \right\} \\
&\quad - K \left( \frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\Lambda_0| - a_0 \ln b_0 + \ln \Gamma(\alpha_0) \right)
\end{aligned} \tag{61}$$

where we make use of  $\mathbb{E}[\tau_k(\beta_k - m_0)\Lambda_0(\beta_k - m_0)^\top] = \frac{a_k}{b_k}(m_k - m_0)^\top \Lambda_0 (m_k - m_0) + \text{tr}(\Lambda_0 \mathbf{V}_k^{-1})$ . The other expectations can be calculated using results derived in Appendix A. In the following expectation, we make use of the established result that  $E[z_{nk}] = r_{nk}$ .

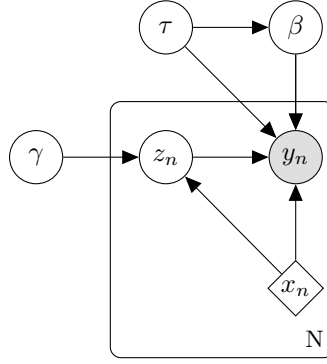
$$\mathbb{E}[\ln q(\mathbf{Z})] = \sum_n \sum_k r_{nk} \ln r_{nk} \tag{62}$$

$$\begin{aligned}
\mathbb{E}[\ln q(\boldsymbol{\beta}, \boldsymbol{\tau})] &= \sum_k \mathbb{E} \left[ \ln \mathcal{N}(\beta_k \mid m_k, (\tau_k \mathbf{Q}_k)^{-1}) \right] + \mathbb{E} \left[ \ln \text{Ga}(\tau_k \mid a_k, b_k) \right] \\
&= \sum_k \left( \frac{D}{2} + a_k - 1 \right) [\psi(a_k) - \psi(b_k)] + a_k \ln b_k - a_k - \ln \Gamma(a_k) + \frac{1}{2} \ln |\mathbf{Q}_k| \\
&\quad - \frac{KD}{2} (\ln(2\pi) + 1)
\end{aligned} \tag{63}$$

where we make use of  $\mathbb{E}[\tau_k(\beta_k - m_k)^\top \mathbf{Q}_k(\beta_k - m_k)] = D$ . The other expectations can be calculated using results derived in Appendix A.

$$\begin{aligned}
\mathbb{E}[\ln q(\boldsymbol{\gamma})] &= \sum_k \mathbb{E} \left[ \ln \mathcal{N}(\gamma_k \mid \mu_k, \mathbf{V}_k^{-1}) \right] \\
&= -\frac{KD}{2} (\ln(2\pi) + 1) + \frac{1}{2} \sum_k \ln |\mathbf{V}_k|
\end{aligned} \tag{64}$$

since  $\mathbb{E}[(\gamma_k - \mu_k)^\top \mathbf{V}_k(\gamma_k - \mu_k)] = D$ .



### D.1 Commonly Used Expectations

Two useful formulations of the approximating density:  $q(\tilde{\beta}_d, \omega_d)$

$$q(\tilde{\beta}_d | w_d) = \omega_d \mathcal{N}(\tilde{\beta}_d | m_d, \mathbf{Q}_d^{-1}) + (1 - \omega_d) \delta_0(\tilde{\beta}_d) \quad (65)$$

$$q(\tilde{\beta}_d) = \lambda_d \mathcal{N}(\tilde{\beta}_d | m_d, \mathbf{Q}_d^{-1}) + (1 - \lambda_d) \delta_0(\tilde{\beta}_d) \quad (66)$$

$$\mathbb{E}[\tilde{\beta}_d] = \lambda_d m_d$$

$$\mathbb{E}[\tilde{\beta}_d | w_d] = w_d m_d$$

$$\mathbb{E}[\tilde{\beta}_d \tilde{\beta}_d^\top | w_d] = w_d (\mathbf{Q}_d^{-1} + m_d m_d^\top)$$

$$\text{Var}(\tilde{\beta}_d | w_d) = w_d \mathbf{Q}_d^{-1} + w_d(1 - w_d) m_d m_d^\top$$

$$\mathbb{E}[\tilde{\beta}_d - m_d | w_d] = (w_d - 1) m_d$$

$$\mathbb{E}[\tilde{\beta}_d^\top \mathbf{U}_d \tilde{\beta}_d] = \lambda_d \text{tr}(\mathbf{U}_d \mathbf{Q}_d^{-1}) + \lambda_d m_d^\top \mathbf{U}_d m_d$$

$$\mathbb{E}[(\tilde{\beta}_d - m_d)^\top \mathbf{Q}_d (\tilde{\beta}_d - m_d) | w_d] = w_d \cdot K$$

$$\mathbb{E}[w_d (\tilde{\beta}_d - m_d)^\top \mathbf{Q}_d (\tilde{\beta}_d - m_d)] = K \cdot \lambda_d$$

$$\mathbb{E}[\tilde{\beta}_d \tilde{\beta}_d^\top] = \lambda_d (\mathbf{Q}_d^{-1} + m_d m_d^\top)$$

$$\text{Var}(\tilde{\beta}_d) = \lambda_d \mathbf{Q}_d^{-1} + \lambda_d(1 - \lambda_d) m_d m_d^\top$$



## D.2 Expectation Calculations Details

$$\begin{aligned}
\mathbb{E}[\tilde{\beta}_d] &= \int \tilde{\beta}_d \cdot \lambda_d \cdot \mathcal{N}(\tilde{\beta}_d | m_d, Q_d^{-1}) d\beta_d = \lambda_d m_d \\
\mathbb{E}[\tilde{\beta}_d | w_d] &= \int \tilde{\beta}_d \cdot w_d \cdot \mathcal{N}(\tilde{\beta}_d | m_d, Q_d^{-1}) d\tilde{\beta}_d = w_d m_d \\
\mathbb{E}[\tilde{\beta}_d \tilde{\beta}_d^\top | w_d] &= \int \tilde{\beta}_d \tilde{\beta}_d^\top \cdot w_d \mathcal{N}(\tilde{\beta}_d | m_d, Q_d^{-1}) d\tilde{\beta}_d = w_d (Q_d^{-1} + m_d m_d^\top) \\
\text{Var}(\tilde{\beta}_d | w_d) &= \mathbb{E}[(\tilde{\beta}_d - \mathbb{E}(\tilde{\beta}_d))(\tilde{\beta}_d - \mathbb{E}(\tilde{\beta}_d))^\top | w_d] \\
&= \mathbb{E}[\tilde{\beta}_d \tilde{\beta}_d^\top | w_d] + \mathbb{E}[\tilde{\beta}_d | w_d] \mathbb{E}[\tilde{\beta}_d | w_d]^\top - 2\mathbb{E}[\tilde{\beta}_d | w_d] \mathbb{E}[\tilde{\beta}_d | w_d]^\top \\
&= w_d (m_d m_d^\top + Q_d^{-1}) + w_d^2 m_d m_d^\top - 2w_d^2 m_d m_d^\top \\
&= w_d Q_d^{-1} + w_d(1 - w_d) m_d m_d^\top \\
\mathbb{E}[\tilde{\beta}_d - m_d | w_d] &= w_d m_d - m_d = (w_d - 1) m_d \\
\mathbb{E}[\tilde{\beta}_d^\top U_d \tilde{\beta}_d] &= \mathbb{E}(\tilde{\beta}_d)^\top U_d \mathbb{E}(\tilde{\beta}_d) + \text{tr}(U_d \text{Var}(\tilde{\beta}_d)) \\
&= \lambda_d^2 m_d^\top U_d m_d + \text{tr}(U_d (\lambda_d Q_d^{-1} + \lambda_d(1 - \lambda_d) m_d m_d^\top)) \\
&= \lambda_d^2 m_d^\top U_d m_d + \lambda_d \text{tr}(U_d Q_d^{-1}) + \lambda_d(1 - \lambda_d) m_d^\top U_d m_d \\
&= \lambda_d \text{tr}(U_d Q_d^{-1}) + \lambda_d m_d^\top U_d m_d \\
\mathbb{E}[(\tilde{\beta}_d - m_d)^\top Q_d (\tilde{\beta}_d - m_d) | w_d] &= \mathbb{E}[(\tilde{\beta}_d - m_d)^\top | w_d]^\top Q_d \mathbb{E}[(\tilde{\beta}_d - m_d) | w_d] + \text{tr}(Q_d \text{Var}(\tilde{\beta}_d | w_d)) \\
&= (w_d - 1)^2 m_d^\top Q_d m_d + w_d \text{tr}(Q_d Q_d^{-1}) + w_d(1 - w_d) m_d^\top Q_d m_d \\
&= (w_d^2 - 2w_d + 1) m_d^\top Q_d m_d + w_d \text{tr}(I_K) + w_d m_d^\top Q_d m_d - w_d^2 m_d^\top Q_d m_d \\
&= K \cdot w_d + m_d^\top Q_d m_d - w_d m_d^\top Q_d m_d \\
\mathbb{E}[w_d (\tilde{\beta}_d - m_d)^\top Q_d (\tilde{\beta}_d - m_d)] &= \mathbb{E}[w_d \mathbb{E}((\tilde{\beta}_d - m_d)^\top Q_d (\tilde{\beta}_d - m_d) | w_d)] \\
&= \mathbb{E}[w_d (K \cdot w_d + m_d^\top Q_d m_d - w_d m_d^\top Q_d m_d)] \\
&= \mathbb{E}[K \cdot w_d^2 + w_d m_d^\top Q_d m_d - w_d^2 m_d^\top Q_d m_d] \\
&= K \cdot \lambda_d + \lambda_d m_d^\top Q_d m_d - \lambda_d m_d^\top Q_d m_d \\
&= K \cdot \lambda_d \\
\mathbb{E}[\tilde{\beta}_d \tilde{\beta}_d^\top] &= \int \tilde{\beta}_d \tilde{\beta}_d^\top \cdot \lambda_d \cdot \mathcal{N}(\tilde{\beta}_d | m_d, Q_d^{-1}) d\beta_d = \lambda_d (Q_d^{-1} + m_d m_d^\top) \\
\text{Var}(\tilde{\beta}_d) &= \mathbb{E}[\tilde{\beta}_d \tilde{\beta}_d^\top] - \mathbb{E}[\tilde{\beta}_d] \mathbb{E}[\tilde{\beta}_d]^\top \\
&= \lambda_d (Q_d^{-1} + m_d m_d^\top) - \lambda_d^2 m_d m_d^\top
\end{aligned}$$