
Bayesian models for Large-scale Hierarchical Classification (Supplementary Material)

Anonymous Author(s)

Affiliation

Address

email

1 VARIATIONAL INFERENCE FOR MODELS

1.1 Details of Variational Inference for HBLR-M2

Generally in Bayesian Inference, computing the posterior distribution of parameters in closed form solution might not be possible. In such situations, one often resorts to computing an approximate posterior which is ‘close’ to the true posterior distribution. For HBLR-M2, the posterior of the model parameters is given by,

$$p(\mathbf{W}, \alpha | \mathbf{D}) \propto p(\mathbf{D} | \mathbf{W}, \alpha) p(\mathbf{W}, \alpha)$$

$$p(\mathbf{D} | \mathbf{W}, \alpha) = \prod_{(x,t) \in \mathbf{D}} \frac{\exp(w_t^\top x)}{\sum_{t' \in T} \exp(w_{t'}^\top x)} \quad (1)$$

$$p(\mathbf{W}, \alpha) = \prod_{y \in Y \setminus T} \prod_{i=1}^d p(\alpha_y^{(i)} | a_y^{(i)}, b_y^{(i)}) \prod_{y \in Y} p(w_y | w_{\pi(y)}, \Sigma_{\pi(y)})$$

$$= \prod_{y \in Y \setminus T} \prod_{i=1}^d \Gamma(\alpha_y^{(i)} | a_y^{(i)}, b_y^{(i)}) \prod_{y \in Y} \mathcal{N}(w_y | w_{\pi(y)}, \Sigma_{\pi(y)}) \quad (2)$$

The posterior has a logistic likelihood term with the \mathbf{W} in (1) and with a Gamma and Normal prior over α , \mathbf{W} in (2). The convolution between a normal-gamma prior and logistic likelihood cannot be computed in closed form; therefore one has to resort to approximate methods to calculate the posterior.

Variational methods try to compute an approximate posterior having a simplified factored form which is closest in KL divergence to the true posterior. They rely on the following bound for the log-marginal probability of D . For any distribution $q(\mathbf{W}, \alpha)$,

$$P(D) = \int q(\mathbf{W}, \alpha) \log \frac{p(\mathbf{W}, \alpha, D)}{q(\mathbf{W}, \alpha)} d\mathbf{W} d\alpha - \int q(\mathbf{W}, \alpha) \log \frac{p(\mathbf{W}, \alpha | D)}{q(\mathbf{W}, \alpha)} d\mathbf{W} d\alpha \quad (3)$$

$$\geq \int q(\mathbf{W}, \alpha) \log \frac{p(\mathbf{W}, \alpha, D)}{q(\mathbf{W}, \alpha)} d\mathbf{W} d\alpha$$

$$= \int q(\mathbf{W}, \alpha) \log p(\mathbf{W}, \alpha, D) d\mathbf{W} d\alpha - \int q(\mathbf{W}, \alpha) \log q(\mathbf{W}, \alpha) d\mathbf{W} d\alpha \quad (4)$$

$$= E_q [\log p(\mathbf{W}, \alpha, D)] + H(q) \quad (5)$$

The (3) can be easily verified by combining the RHS terms.

In order to compute such a q , we start by assuming a simplified factored form using independent distributions for each parameter. Note that this does not neglect the dependence between the various parameters in the model; it just finds the suitable factored form which best approximates the dependencies.

$$\begin{aligned} q(\mathbf{W}, \alpha) &= \prod_{y \in Y \setminus T} q(\alpha_y) \prod_{y \in Y} q(w_y) \\ q(w_y) &\propto \mathcal{N}(\cdot | \mu_y, \Psi_y) \\ q(\alpha_y) &= \prod_{i=1}^d q(\alpha_y^{(i)}) \propto \prod_{i=1}^d \Gamma(\cdot | \tau_y^{(i)}, v_y^{(i)}) \end{aligned}$$

Here $q(\mathbf{W}, \alpha) \equiv q(\mathbf{W} | \mu, \Psi) q(\alpha | \tau, v)$ where τ, v, μ, Σ are variational parameters which we can optimize one at a time to maximize (4).

For example, to optimize w_y , we differentiate $P(D)$ w.r.t $q(w_y)$,

$$\begin{aligned} P(D) &= \int q(w_y | \mu_y, \Psi_y) [q(\mathbf{W}^{-w_y}, \alpha) \log p(\mathbf{W}, \alpha, D) d\mathbf{W} d\alpha] - \int q(w_y | \mu_y, \Psi_y) \log q(w_y | \mu_y, \Psi_y) dw_y \\ \log q^*(w_y | \mu_y, \Psi_y) &= E_{q^{-w_y}} [\log p(\mathbf{W}, \alpha, D) d\mathbf{W} d\alpha] + \text{constant} \end{aligned} \quad (6)$$

where $-w_y$ denotes all parameters other than w_y . Similary,

$$\log q^*(\alpha_y | \tau_y, v_y) = E_{q^{-\alpha_y}} [\log p(\mathbf{W}, \alpha, D) d\mathbf{W} d\alpha] + \text{constant} \quad (7)$$

In (7) and (6), to update the parameters we need to calculate the expected log-likelihood under the variational posterior i.e. $E_q[\log p(\mathbf{W}, \alpha, D)]$.

$$E_q[\log p(\mathbf{W}, \alpha, D)] = E_q[\log P(D | \mathbf{W}, \alpha)] + E_q[\log P(\mathbf{W}, \alpha)]$$

Where,

$$\begin{aligned} E_q[\log P(\mathbf{W}, \alpha)] &= \sum_{Y \in Y} -\frac{1}{2} E_q \left[(w_y - w_{\pi(y)})^\top \Sigma_{\pi(y)}^{-1} (w_y - w_{\pi(y)}) \right] \\ &\quad - \sum_{y \in Y} \frac{1}{2} E_q [\log |\Sigma_y|] - \sum_{y \in Y \setminus T} \sum_{i=1}^d E_q [\alpha_y^{(i)}] b_y + E_q [(\alpha_y^{(i)} - 1) \log(\alpha_y^{(i)})] \end{aligned} \quad (8)$$

$$E_q[\log P(D | \mathbf{W}, \alpha)] = \sum_{(x,t) \in D} \mu_y^\top x - \sum_{(x,t) \in D} E_q \left[\log \left(\sum_{y \in T} \exp(w_y^\top x) \right) \right] \quad (9)$$

The RHS term in the above equation is not directly computable. Therefore we use a suitable lower-bound proposed in [1], to bound the log-likelihood term from below.

$$\log \left(\sum_{y \in T} \exp(w_y^\top x) \right) \leq \beta_x + \sum_{y \in T} \frac{(w_y^\top x - \beta_x - \xi_{xy})}{2} + \sum_{y \in T} \lambda(\xi_{xy}) \left((w_y^\top x - \beta_x)^2 - \xi_{xy}^2 \right) \quad (10)$$

$$+ \sum_{y \in T} \log(1 + \xi_{xy}^2) \quad (11)$$

$$= -\left(\frac{|T|}{2} - 1\right) \beta_x + \sum_{y \in T} \frac{w_y^\top x}{2} - \sum_{y \in T} \frac{\xi_{xy}}{2} \quad (12)$$

$$+ \sum_{y \in T} \lambda(\xi_{xy}) (w_y^\top (xx^\top) w_y - 2\beta_x (w_y^\top x) \lambda(\xi_{xy}) + \beta_x^2 - \xi_{xy}^2) + \sum_{y \in T} \log(1 + e^{\xi_{xy}}) \quad (13)$$

Here we have introduced variational parameter β_x and ξ_{xy} for every $x \in D, y \in Y$. In order to get the tightest possible bound, we can optimize over these variational parameters. We discuss this optimization later in the text.

$$E_q \left[\log \left(\sum_{y \in T} \exp(w_y^\top x) \right) \right] = \sum_{y \in Y} \left(\frac{1}{2} - 2\lambda(\xi_{xy})\beta_x \right) \mu_y^\top x + \sum_{y \in Y} \mu_y^\top (2\lambda(\xi_{xy})xx^\top) \mu_y - \left(\frac{|T|}{2} - 1 \right) \beta_x - \sum_{y \in Y} \left(\lambda(\xi_{xy})\beta_x^2 - \frac{\xi_{xy}}{2} - \frac{\xi_{xy}^2}{2} + \log(1 + e^{\xi_{xy}}) \right) \quad (14)$$

1.1.1 Optimizing $q^*(w_y)$

Combining the above equation with (8) and (6), the update for parameter $w_y, y \in T$ can be written as

$$\log q^*(w_y | \mu_y, \Psi_y) = \left(\sum_{(x,t) \in D} \left(I(t=y) - \frac{1}{2} + 2\lambda(\xi_{xy})\beta_x \right) x + E_{q^{-w_y}} \left[\Sigma_{\pi(y)}^{-1} \right] \mu_y - w_y \right)^\top \left(\sum_{(x,t) \in D} 2\lambda(\xi_{xy})xx^\top + E_{q^{-w_y}} \left[\Sigma_{\pi(y)}^{-1} \right] \right) \left(\sum_{(x,t) \in D} \left(I(t=y) - \frac{1}{2} + 2\lambda(\xi_{xy})\beta_x \right) x + E_{q^{-w_y}} \left[\Sigma_{\pi(y)}^{-1} \right] \mu_y - w_y \right) + constant$$

Since $q^*(w_y)$ is assumed to be a normal distribution, we can directly match the sufficient statistics i.e. mean and the covariance matrix and set μ_y and Ψ_y as in the paper. Note that we have used the fact that

$$E_{q^{-w_y}} \left[\Sigma_y^{-1} \right] = \text{diag} \left(E_{q^{-w_y}} \left[\alpha_y^{(1)} \right], E_{q^{-w_y}} \left[\alpha_y^{(2)} \right], \dots, E_{q^{-w_y}} \left[\alpha_y^{(d)} \right] \right) = \text{diag} \left(\frac{\tau_y^{(1)}}{v_y^{(1)}}, \frac{\tau_y^{(2)}}{v_y^{(2)}}, \dots, \frac{\tau_y^{(d)}}{v_y^{(d)}} \right)$$

For $q^*(w_y), y \notin T$, fortunately the w_y is not present in the logistic function, therefore, we can use just (8) to match the likelihood with (6) and get the equations as in the paper.

1.1.2 Optimizing $q^*(\alpha_y)$

For optimizing $q^*(\alpha_y)$, we essentially follow the same strategy as above. Each optimize each $\alpha_y^{(i)}$ by matching (8) and (7)

$$\log q^*(\alpha_y^{(i)} | \tau_y^{(i)}, v_y^{(i)}) = - \sum_{y \in C_y} \frac{\alpha_y^{(i)}}{2} E_q \left[(w^{(i)} - w_{\pi(y)}^{(i)})^2 \right] + \sum_{y \in C_y} \frac{\alpha_y^{(i)}}{2} - \alpha_y^{(i)} b_y^{(i)} + (\alpha_y^{(i)} - 1) \log(\alpha_y^{(i)}) = -\alpha_y^{(i)} \left(b_y^{(i)} + \sum_{y \in C_y} \frac{1}{2} \left((\mu^{(i)} - \mu_{\pi(y)}^{(i)})^2 + \Psi^{(i,i)} + \Psi_{\pi(y)}^{(i,i)} \right) \right) + (\alpha_y^{(i)} + \frac{|C_y|}{2} - 1) \log(\alpha_y^{(i)}) + constant$$

Since $q^*(\alpha_y)$ is a gamma distribution, we simply match the sufficient statistics and set update the as in the paper. Note that we have used the fact that

$$-\log |\Sigma_y| = \log |\Sigma_y^{-1}| = \log \prod_{i=1}^d \alpha_y^{(i)} = \sum_{i=1}^d \log \alpha_y^{(i)}$$

1.2 Extension to HBLR-M1

Most of the derivation simply goes through changed except that $\alpha_y^{(1)} = \alpha_y^{(2)} = \dots = \alpha_y$. Although, this makes a difference only in how the α 's are updated; we present the full update equations.

If $y \in T$

$$\begin{aligned}\Psi_y^{-1} &= \sum_{(x,t) \in D} 2\lambda(\xi_{xy})xx^\top + E_{q^{-w_y}} [\Sigma_{\pi(y)}^{-1}] \\ \mu_y &= \Psi_y \left(\sum_{(x,t) \in D} (I(t=y) - \frac{1}{2} + 2\lambda(\xi_{xy})\beta_x)x + E_{q^{-w_y}} [\Sigma_{\pi(y)}^{-1}] \mu_{\pi(y)} \right)\end{aligned}$$

If $y \notin T$

$$\begin{aligned}\Psi_y^{-1} &= E_{q^{-w_y}} [\Sigma_{\pi(y)}^{-1}] + |C_y| E_{q^{-w_y}} [\Sigma_y^{-1}] \\ \mu_y &= \Psi_y \left(E_{q^{-w_y}} [\Sigma_{\pi(y)}^{-1}] \mu_{\pi(y)} + E_{q^{-w_y}} [\Sigma_y^{-1}] \sum_{c \in C_y} \mu_c \right) \\ v_y &= (b_y + \sum_{c \in C_y} \text{diag}(\Psi_y) + \text{diag}(\Psi_c) + (\mu_y - \mu_c)^\top (\mu_y - \mu_c))^\top e \\ \tau_y &= a_y + \frac{|C_y|d}{2}\end{aligned}$$

where

$$E_{q^{-w_y}} [\Sigma_y^{-1}] = \text{diag} \left(\frac{\tau_y}{v_y}, \frac{\tau_y}{v_y}, \dots, \frac{\tau_y}{v_y} \right)$$

and e is a unit vector.

1.3 Extension to HBLR-M3

The extension to HBLR-M3 follows on similar lines.

If $y \in T$

$$\begin{aligned}\Psi_y^{-1} &= \sum_{(x,t) \in D} 2\lambda(\xi_{xy})xx^\top + E_{q^{-w_y}} [\Sigma_y^{-1}] \\ \mu_y &= \Psi_y \left(\sum_{(x,t) \in D} (I(t=y) - \frac{1}{2} + 2\lambda(\xi_{xy})\beta_x)x + E_{q^{-w_y}} [\Sigma_y^{-1}] \mu_{\pi(y)} \right)\end{aligned}$$

If $y \notin T$

$$\begin{aligned}\Psi_y^{-1} &= E_{q^{-w_y}} [\Sigma_{\pi(y)}^{-1}] + |C_y| E_{q^{-w_y}} [\Sigma_y^{-1}] \\ \mu_y &= \Psi_y \left(E_{q^{-w_y}} [\Sigma_{\pi(y)}^{-1}] \mu_{\pi(y)} + E_{q^{-w_y}} [\Sigma_y^{-1}] \sum_{c \in C_y} \mu_c \right) \\ v_y &= (b_y + \text{diag}(\Psi_{\pi(y)}) + \text{diag}(\Psi_y) + (\mu_y - \mu_{\pi(y)})^\top (\mu_y - \mu_{\pi(y)}))^\top e \\ \tau_y &= a_y + \frac{d}{2}\end{aligned}$$

where

$$E_{q^{-w_y}} [\Sigma_y^{-1}] = \text{diag} \left(\frac{\tau_y}{v_y}, \frac{\tau_y}{v_y}, \dots, \frac{\tau_y}{v_y} \right)$$

2 Empirical Bayes

In this section, we show why Empirical Bayes route for learning the hyperparameters in our model does not work. Let us model M1 for instance. The general procedure for Empirical Bayes is to maximize the marginal likelihood w.r.t to the hyperparameters and get point estimates for them (another approach would be use MCMC, which we do not pursue in the interested of scalability).

In M1, the marginal likelihood of the data can be computed as

$$P(D|a, b) \propto \int_{\alpha} P(\alpha|a, b) \int_{\mathbf{W}} P(D|\mathbf{W})P(\mathbf{W}|\alpha)$$

As before the integral cannot be computed in closed form and hence cannot be maximized analytically. Therefore we use the variational lower-bound as a proxy which is more amenable to maximization.

$$\begin{aligned} \log P(D|a, b) &\geq E_q[\log P(D|a, b)] \\ &\geq E_q[\log P(\alpha|a, b)] + \text{terms that do not depend on } a, b \\ &\geq \sum_{Y \in Y \setminus T} E_q[\log(P(\alpha_y|a_y, b_y))] \\ &\geq \sum_{Y \in Y \setminus T} E_q \left[b_y^{a_y} \frac{1}{\Gamma(a_y)} \alpha_y^{a_y-1} e^{-b_y \alpha} \right] \end{aligned}$$

Note that the maximization can be carried out independently for all the a_y 's. This leads to a Gamma distribution type MLE of the following form,

$$(a_y^*, b_y^*) = \arg \max a_y \log(b_y) - \log \Gamma(a_y) + (a_y - 1)E_q[\log(\alpha_y)] - b_y E_q[\alpha_y]$$

Note that $E_q[\log(\alpha_y)] = \Psi(\tau_y) - \log(v_y)$ and $E_q[\alpha_y] = \frac{\tau_y}{v_y}$; both of twich can be considered as constants in the maximization. But since there is exactly only one sample, one cannot learn the a_y, b_y effectively [2].

One way to overcome this, is to assume that all the α_y 's are commonly drawn from a single a, b . This enables a larger number of samples to succesfully estimate the the single a, b . The downside is that, this commonly shrinks all the α_y 's to the expected value of the distribution $\frac{a}{b}$, which might be not a good thing.

We conducted several preliminary experiments were we tried sharing all α_y 's under single a, b as well sharing the α_y 's of sinling nodes etc. None of the models seemed to achieve competitive performance. For example, using a common a, b on the best performing model M3 on the CLEF dataset achieved a Micro- F_1 of 80.32 and Macro- F_1 of 54.82 both of which are lower than M3-var. Further investigation is required to establish how empirical can be succesfully applied.

References

- [1] G. Bouchard. Efficient bounds for the softmax function. 2007.
- [2] George Casella. Empirical bayes method - a tutorial. Technical report.