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

Siddharth Gopal
sgopal1@andrew.cmu.edu

Yiming Yang
yiming@cs.cmu.edu

Bing Bai
bing@nec-labs.com

Alex Niculescu-Mizil
alex@nec-labs.com

1 Variational Inference for HBLR models

For the supplementary material we slightly change the notation as follows,

Define a hierarchy as a set of nodes $\mathcal{N} = \{1, 2, \dots\}$ with the parent relationship $\pi : \mathcal{N} \rightarrow \mathcal{N}$ where $\pi(n)$ is the parent of node $n \in \mathcal{N}$. Let $\mathbf{D} = \{(x_i, t_i)\}_{i=1}^N$ denote the training data where $x_i \in \mathbb{R}^d$ is an instance, $t_i \in T$ is a label, where $T \subset \mathcal{N}$ is the set of leaf nodes in the hierarchy labeled from 1 to $|T|$. We assume that each instance is assigned to one of the leaf nodes in the hierarchy. Let C_n be the set of all children of node n .

Models M1, M2 and M3 are defined as follows,

$$\begin{aligned}
 \mathbf{M1} \quad & w_{root} \sim \mathcal{N}(w_0, \Sigma_0), \quad \alpha_{root} \sim \Gamma(a_0, b_0) \\
 & w_n | w_{\pi(n)}, \Sigma_{\pi(n)} \sim \mathcal{N}(w_{\pi(n)}, \Sigma_{\pi(n)}) \quad \forall n, \quad \alpha_n \sim \Gamma(a_n, b_n) \quad \forall n \notin T \\
 & t | x, \mathbf{W} \sim \text{Categorical}(p_1(x), p_2(x), \dots, p_{|T|}(x)) \quad \forall (x, t) \in \mathbf{D} \\
 & p_i(x) = \exp(w_i^\top x) / \sum_{t' \in T} \exp(w_{t'}^\top x)
 \end{aligned} \tag{1}$$

$$\begin{aligned}
 \mathbf{M2} \quad & w_n | w_{\pi(n)}, \Sigma_{\pi(n)} \sim \mathcal{N}(w_{\pi(n)}, \Sigma_{\pi(n)}) \quad \forall n \\
 & \alpha_n^{(i)} \sim \Gamma(a_n^{(i)}, b_n^{(i)}) \quad i = 1..d, \forall n \notin T \\
 & \text{where } \Sigma_{\pi(n)}^{-1} = \text{diag}(\alpha_{\pi(n)}^{(1)}, \alpha_{\pi(n)}^{(2)}, \dots, \alpha_{\pi(n)}^{(d)})
 \end{aligned}$$

$$\begin{aligned}
 \mathbf{M3} \quad & w_n | w_{\pi(n)}, \Sigma_n \sim \mathcal{N}(w_{\pi(n)}, \Sigma_n) \quad \forall n \\
 & \alpha_n \sim \Gamma(a_n, b_n) \quad \forall n \notin T
 \end{aligned}$$

1.1 Variational Inference for model M2

The posterior of the model parameters is given by,

$$\begin{aligned}
 p(\mathbf{W}, \boldsymbol{\alpha} | \mathbf{D}) & \propto p(\mathbf{D} | \mathbf{W}, \boldsymbol{\alpha}) p(\mathbf{W}, \boldsymbol{\alpha}) \\
 p(\mathbf{D} | \mathbf{W}, \boldsymbol{\alpha}) & = \prod_{(x, t) \in \mathbf{D}} \frac{\exp(w_t^\top x)}{\sum_{t' \in T} \exp(w_{t'}^\top x)}
 \end{aligned} \tag{2}$$

$$\begin{aligned}
 p(\mathbf{W}, \boldsymbol{\alpha}) & = \prod_{n \in \mathcal{N} \setminus T} \prod_{i=1}^d p(\alpha_n^{(i)} | a_n^{(i)}, b_n^{(i)}) \prod_{n \in \mathcal{N}} p(w_n | w_{\pi(n)}, \Sigma_{\pi(n)}) \\
 & = \prod_{n \in \mathcal{N} \setminus T} \prod_{i=1}^d \Gamma(\alpha_n^{(i)} | a_n^{(i)}, b_n^{(i)}) \prod_{n \in \mathcal{N}} \mathcal{N}(w_n | w_{\pi(n)}, \Sigma_{\pi(n)})
 \end{aligned} \tag{3}$$

The posterior has a logistic likelihood term with the \mathbf{W} in (2) and with a Gamma and Normal prior over α , \mathbf{W} in (3). 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 $q(\mathbf{W}, \alpha)$,

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

$$\begin{aligned} &\geq \int q(\mathbf{W}, \alpha) \log \frac{p(\mathbf{W}, \alpha, D)}{q(\mathbf{W}, \alpha)} d\mathbf{W} d\alpha \\ &= \int q(\mathbf{W}, \alpha) \log p(\mathbf{W}, \alpha, D) d\mathbf{W} d\alpha - \int q(\mathbf{W}, \alpha) \log q(\mathbf{W}, \alpha) d\mathbf{W} d\alpha \end{aligned} \quad (5)$$

$$VLB = E_q [\log p(\mathbf{W}, \alpha, D)] + H(q) \quad (6)$$

Firstly, note that equation(4) can be easily verified by combining the RHS terms. Equation (6) represents the variational lower-bound for the likelihood of the data, this is quantity that we maximize w.r.t q .

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.

$$\begin{aligned} q(\mathbf{W}, \alpha) &= \prod_{n \in \mathcal{N} \setminus T} q(\alpha_n) \prod_{n \in \mathcal{N}} q(w_n) \\ q(w_n) &\propto \mathcal{N}(\cdot | \mu_n, \Psi_n) \\ q(\alpha_n) &= \prod_{i=1}^d q(\alpha_n^{(i)}) \propto \prod_{i=1}^d \Gamma(\cdot | \tau_n^{(i)}, \nu_n^{(i)}) \end{aligned}$$

Here $q(\mathbf{W}, \alpha) \equiv q(\mathbf{W} | \boldsymbol{\mu}, \boldsymbol{\Psi}) q(\alpha | \boldsymbol{\tau}, \boldsymbol{\nu})$ where $\boldsymbol{\tau}, \boldsymbol{\nu}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$ are variational parameters which we can optimize one at a time to maximize (6).

For example, to optimize w_n , we differentiate VLB w.r.t $q(w_n)$,

$$VLB = \int q(w_n | \mu_n, \Psi_n) [q(\mathbf{W}^{-w_n}, \alpha) \log p(\mathbf{W}, \alpha, D) d\mathbf{W} d\alpha] - \int q(w_n | \mu_n, \Psi_n) \log q(w_n | \mu_n, \Psi_n) dw_n$$

$$\text{Setting the gradient w.r.t } q(w_n) \text{ to zero yields} \quad (7)$$

$$\log q^*(w_n | \mu_n, \Psi_n) = E_{q^{-w_n}} [\log p(\mathbf{W}, \alpha, D)] + \text{constant} \quad (8)$$

where $-w_n$ denotes all parameters other than w_n . Similarly,

$$\log q^*(\alpha_n | \tau_n, \nu_n) = E_{q^{-\alpha_n}} [\log p(\mathbf{W}, \alpha, D)] + \text{constant} \quad (9)$$

In general, to update the parameters we need to calculate the expected log-likelihood with 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,

$$\begin{aligned} E_q[\log P(\mathbf{W}, \alpha)] &= \sum_{n \in \mathcal{N}} -\frac{1}{2} E_q \left[(w_n - w_{\pi(n)})^\top \Sigma_{\pi(n)}^{-1} (w_n - w_{\pi(n)}) \right] \\ &\quad - \sum_{n \in \mathcal{N}} \frac{1}{2} E_q [\log |\Sigma_n|] - \sum_{n \in \mathcal{N} \setminus T} \sum_{i=1}^d E_q [\alpha_n^{(i)}] b_n + E_q \left[(a_n^{(i)} - 1) \log(a_n^{(i)}) \right] \end{aligned} \quad (10)$$

$$E_q [\log P(\mathbf{D}|\mathbf{W}, \boldsymbol{\alpha})] = \sum_{(x,t) \in D} \mu_n^\top x - \sum_{(x,t) \in D} E_q \left[\log \left(\sum_{n \in T} \exp(w_n^\top x) \right) \right] \quad (11)$$

The expectation of the log-sum-exp function is not computable in closed form. Therefore we replace the log-sum-exp with suitable lower-bound proposed in [1] whose expectation can be computed in closed form.

$$\begin{aligned} \log \left(\sum_{n \in T} \exp(w_n^\top x) \right) &\leq \beta_x + \sum_{n \in T} \frac{(w_n^\top x - \beta_x - \xi_{xn})}{2} + \sum_{n \in T} \lambda(\xi_{xn}) \left((w_n^\top x - \beta_x)^2 - \xi_{xn}^2 \right) \\ &\quad + \sum_{n \in T} \log(1 + e^{\xi_{xn}}) \\ &\leq - \left(\frac{|T|}{2} - 1 \right) \beta_x + \sum_{n \in T} \frac{w_n^\top x}{2} - \sum_{n \in T} \frac{\xi_{xn}}{2} \\ &\quad + \sum_{n \in T} \lambda(\xi_{xn}) (w_n^\top (x x^\top) w_n - 2\beta_x (w_n^\top x) + \beta_x^2 - \xi_{xn}^2) + \sum_{n \in T} \log(1 + e^{\xi_{xn}}) \end{aligned}$$

Here we have introduced variational parameter β_x and ξ_{xn} for every $x \in D, n \in \mathcal{N}$. In order to get the tightest possible bound, we can optimize over these variational parameters. The expected log-sum-exp function can be computed as followed (the function $\lambda(x)$ is defined in the paper)

$$\begin{aligned} E_q \left[\log \left(\sum_{n \in T} \exp(w_n^\top x) \right) \right] &\leq \sum_{n \in \mathcal{N}} \left(\frac{1}{2} - 2\lambda(\xi_{xn})\beta_x \right) \mu_n^\top x + \sum_{n \in \mathcal{N}} \mu_n^\top (2\lambda(\xi_{xn})x x^\top) \mu_n \\ &\quad - \left(\frac{|T|}{2} - 1 \right) \beta_x - \sum_{n \in \mathcal{N}} \left(\lambda(\xi_{xn})\beta_x^2 - \frac{\xi_{xn}}{2} - \lambda(\xi_{xn})\xi_{xn}^2 + \log(1 + e^{\xi_{xn}}) \right) \end{aligned} \quad (12)$$

1.1.1 Optimizing $q^*(w_n)$

Combining the bound for the logistic soft-max function with (11) and (10), the update for parameter $w_n, n \in T$ can be written as

$$\begin{aligned} \log q^*(w_n | \mu_n, \Psi_n) &= \left(\sum_{(x,t) \in D} (I(t=n) - \frac{1}{2} + 2\lambda(\xi_{xn})\beta_x)x + E_{q^{-w_n}} \left[\Sigma_{\pi(n)}^{-1} \right] \mu_{\pi(n)} \right) - w_n \\ &\quad \left(\sum_{(x,t) \in D} 2\lambda(\xi_{xn})x x^\top + E_{q^{-w_n}} \left[\Sigma_{\pi(n)}^{-1} \right] \right) \\ &\quad \left(\sum_{(x,t) \in D} (I(t=n) - \frac{1}{2} + 2\lambda(\xi_{xn})\beta_x)x + E_{q^{-w_n}} \left[\Sigma_{\pi(n)}^{-1} \right] \mu_{\pi(n)} \right) - w_n + constant \end{aligned}$$

From the above, we can directly match the sufficient statistics i.e. mean and the covariance matrix of $q^*(w_n)$ and set μ_n and Ψ_n as in the main paper. Note that we have used the fact that

$$\begin{aligned} E_{q^{-w_n}} \left[\Sigma_n^{-1} \right] &= \text{diag} \left(E_{q^{-w_n}} \left[\alpha_n^{(1)} \right], E_{q^{-w_n}} \left[\alpha_n^{(2)} \right], \dots, E_{q^{-w_n}} \left[\alpha_n^{(d)} \right] \right) \\ &= \text{diag} \left(\frac{\tau_n^{(1)}}{\nu_n^{(1)}}, \frac{\tau_n^{(2)}}{\nu_n^{(2)}}, \dots, \frac{\tau_n^{(d)}}{\nu_n^{(d)}} \right) \end{aligned}$$

For $n \notin T$, fortunately there is not convolution with a logistic function, therefore the posterior can be derived similar to convolutions of Normal distributions.

1.1.2 Optimizing $q^*(\alpha_n)$

For optimizing $q^*(\alpha_n)$, we essentially follow the same strategy as above.

$$\begin{aligned} \log q^*(\alpha_n^{(i)} | \tau_n^{(i)}, v_n^{(i)}) &= - \sum_{c \in C_n} \frac{\alpha_n^{(i)}}{2} E_q \left[(w_n^{(i)} - w_c^{(i)})^2 \right] + \left(\sum_{c \in C_n} \frac{\log(\alpha_n^{(i)})}{2} \right) - \alpha_n^{(i)} b_n^{(i)} + (a_n^{(i)} - 1) \log(\alpha_n^{(i)}) \\ &= -\alpha_n^{(i)} \left(b_n^{(i)} + \sum_{c \in C_n} \frac{1}{2} \left((\mu_c^{(i)} - \mu_n^{(i)})^2 + \Psi_n^{(i,i)} + \Psi_c^{(i,i)} \right) \right) \\ &\quad + (a_n^{(i)} + \frac{|C_n|}{2} - 1) \log(\alpha_n^{(i)}) + \text{constant} \end{aligned}$$

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

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

1.2 Variational Inference for model M1

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

If $n \in T$

$$\begin{aligned} \Psi_n^{-1} &= \sum_{(x,t) \in D} 2\lambda(\xi_{xn}) x x^\top + E_{q^{-w_n}} \left[\Sigma_{\pi(n)}^{-1} \right] \\ \mu_n &= \Psi_n \left(\sum_{(x,t) \in D} (I(t=n) - \frac{1}{2} + 2\lambda(\xi_{xn}) \beta_x) x + E_{q^{-w_n}} \left[\Sigma_{\pi(n)}^{-1} \right] \mu_{\pi(n)} \right) \end{aligned}$$

If $n \notin T$

$$\begin{aligned} \Psi_n^{-1} &= E_{q^{-w_n}} \left[\Sigma_{\pi(n)}^{-1} \right] + |C_n| E_{q^{-w_n}} \left[\Sigma_n^{-1} \right] \\ \mu_n &= \Psi_n \left(E_{q^{-w_n}} \left[\Sigma_{\pi(n)}^{-1} \right] \mu_{\pi(n)} + E_{q^{-w_n}} \left[\Sigma_n^{-1} \right] \sum_{c \in C_n} \mu_c \right) \\ v_n &= b_n + \frac{1}{2} \sum_{c \in C_n} \text{trace}(\Psi_n) + \text{trace}(\Psi_c) + (\mu_n - \mu_c)^\top (\mu_n - \mu_c) \\ \tau_n &= a_n + \frac{|C_n|d}{2} \end{aligned}$$

where

$$E_{q^{-w_n}} \left[\Sigma_n^{-1} \right] = \text{diag} \left(\frac{\tau_n}{v_n}, \frac{\tau_n}{v_n}, \dots, \frac{\tau_n}{v_n} \right)$$

1.3 Variational Inference for model M3

The extension to HBLR-M3 follows on similar lines.

$$\begin{aligned} \text{If } n \in T \quad \Psi_n^{-1} &= \sum_{(x,t) \in D} 2\lambda(\xi_{xn})xx^\top + E_{q^{-w_n}} [\Sigma_n^{-1}] \\ \mu_n &= \Psi_n \left(\sum_{(x,t) \in D} (I(t=n) - \frac{1}{2} + 2\lambda(\xi_{xn})\beta_x)x + E_{q^{-w_n}} [\Sigma_n^{-1}] \mu_{\pi(n)} \right) \end{aligned}$$

If $n \notin T$

$$\begin{aligned} \Psi_n^{-1} &= E_{q^{-w_n}} [\Sigma_n^{-1}] + \sum_{c \in C_n} E_{q^{-w_n}} [\Sigma_c^{-1}] \\ \mu_n &= \Psi_n \left(E_{q^{-w_n}} [\Sigma_n^{-1}] \mu_{\pi(n)} + \sum_{c \in C_n} E_{q^{-w_n}} [\Sigma_c^{-1}] \mu_c \right) \\ v_n &= b_n + \frac{1}{2} (\text{trace}(\Psi_{\pi(n)}) + \text{trace}(\Psi_n) + (\mu_n - \mu_{\pi(n)})^\top (\mu_n - \mu_{\pi(n)})) \\ \tau_n &= a_n + \frac{d}{2} \end{aligned}$$

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.

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

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

$$(a_n^*, b_n^*) = \arg \max a_n \log(b_n) - \log \Gamma(a_n) + (a_n - 1)E_q[\log(\alpha_n)] - b_n E_q[\alpha_n]$$

Note that $E_q[\log(\alpha_n)] = \Psi(\tau_n) - \log(v_n)$ and $E_q[\alpha_n] = \frac{\tau_n}{v_n}$; both of which can be considered as constants in the maximization. But since there is exactly only one sample, one cannot learn the a_n, b_n effectively [2].

One way to overcome this, is to assume that all the α_n 's are commonly drawn from a single a, b . This enables a larger number of samples to successfully estimate the the single a, b . The downside is

that, this commonly shrinks all the α_n 's to the expected value of the distribution $\frac{a}{b}$, which might be not a good thing.

We conducted several preliminary experiments where we tried sharing all α_y 's under single a, b as well sharing the α_y 's of sibling 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 successfully applied.

References

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