

Incorporating Random effect into GP Sparse Approximation

Daeyoung Lim*
Department of Statistics
Korea University

February 13, 2016

1 GP linear mixed model

1.1 Model specifications

The model specifications remain the same except that the random effect term has been added.

$$\begin{aligned}y|\theta &\sim \mathcal{N}(Z\alpha + A\beta, \gamma^2 I_n), (n \times n) \\ \alpha|\sigma^2 &\sim \mathcal{N}\left(0, \frac{\sigma^2}{m} I_{2m}\right), (2m \times 1) \\ \beta &\sim \mathcal{N}(0, \Sigma_\beta), (s \times 1) \\ \lambda &\sim \mathcal{N}(\mu_\lambda, \Sigma_\lambda), (d \times 1) \\ \sigma &\sim \text{half-Cauchy}(A_\sigma) \\ \gamma &\sim \text{half-Cauchy}(A_\gamma)\end{aligned}$$

where A is the design matrix for the random effects and β is the parameter vector of the random effects.

1.2 Lower bound

$$p(y, \theta) = \mathcal{N}(y|Z\alpha + A\beta, \gamma^2 I_n) \mathcal{N}\left(\alpha \middle| 0, \frac{\sigma^2}{m} I_{2m}\right) \mathcal{N}(\beta|0, \Sigma_\beta) \mathcal{N}(\lambda, \mu_\lambda, \Sigma_\lambda) \text{HC}(A_\sigma) \text{HC}(A_\gamma)$$

*Prof. Taeryon Choi

$$\begin{aligned}
(1) \mathbb{E} [\log p(y|\theta)] &= -\frac{n}{2} \log(2\pi) - \frac{n}{2} \mathbb{E} [\log \gamma^2] \\
&\quad - \frac{1}{2} \mathbb{E} \left[\frac{1}{\gamma^2} \right] \left\{ y^T y - 2 (\mathbb{E} [Z] m_\alpha + A m_\beta)^T y + \text{Tr} (\mathbb{E} [Z^T Z] S_\alpha) + m_\alpha^T \mathbb{E} [Z^T Z] m_\alpha \right. \\
&\quad \left. + 2 m_\beta^T A^T \mathbb{E} [Z] m_\alpha + \text{Tr} (A^T A S_\beta) + m_\beta^T A^T A m_\beta \right\} \\
(2) \mathbb{E} [\log p(\alpha|\sigma)] &= -m \log(2\pi) - m \mathbb{E} [\log \sigma^2] + m \log m - \frac{m}{2} \mathbb{E} \left[\frac{1}{\sigma^2} \right] \text{Tr} (S_\alpha + m_\alpha m_\alpha^T) \\
(3) \mathbb{E} [\log p(\beta)] &= -\frac{s}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_\beta| - \frac{1}{2} \left\{ \text{Tr} (\Sigma_\beta^{-1} S_\beta) + m_\beta^T \Sigma_\beta^{-1} m_\beta \right\} \\
(4) \mathbb{E} [\log p(\lambda)] &= -\frac{d}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_\lambda| - \frac{1}{2} (m_\lambda - \mu_\lambda)^T \Sigma_\lambda^{-1} (m_\lambda - \mu_\lambda) - \frac{1}{2} \text{Tr} (\Sigma_\lambda^{-1} S_\lambda) \\
(5) \mathbb{E} [\log p(\sigma)] &= \log(2A_\sigma) - \log \pi - \mathbb{E} [\log (A_\sigma^2 + \sigma^2)] \\
(6) \mathbb{E} [\log p(\gamma)] &= \log(2A_\gamma) - \log \pi - \mathbb{E} [\log (A_\gamma^2 + \gamma^2)] \\
(1) \mathbb{E} [\log q(\alpha)] &= -m \log(2\pi) - \frac{1}{2} \log |S_\alpha| - m \\
(2) \mathbb{E} [\log q(\lambda)] &= -\frac{d}{2} \log(2\pi) - \frac{1}{2} \log |S_\lambda| - \frac{d}{2} \\
(3) \mathbb{E} [\log q(\beta)] &= -\frac{s}{2} \log(2\pi) - \frac{1}{2} \log |S_\beta| - \frac{s}{2} \\
(4) \mathbb{E} [\log q(\sigma)] &= -C_\sigma \frac{\mathcal{H}(2m, C_\sigma, A_\sigma^2)}{\mathcal{H}(2m-2, C_\sigma, A_\sigma^2)} - \log \mathcal{H}(2m-2, C_\sigma, A_\sigma^2) - 2m \mathbb{E} [\log \sigma] - \mathbb{E} [\log (A_\sigma^2 + \sigma^2)] \\
(5) \mathbb{E} [\log q(\gamma)] &= -C_\gamma \frac{\mathcal{H}(n, C_\gamma, A_\gamma^2)}{\mathcal{H}(n-2, C_\gamma, A_\gamma^2)} - \log \mathcal{H}(n-2, C_\gamma, A_\gamma^2) - n \mathbb{E} [\log \gamma] - \mathbb{E} [\log (A_\gamma^2 + \gamma^2)] \\
\mathcal{L} &= (1) + (2) + (3) + (4) + (5) + (6) - (1) - (2) - (3) - (4) - (5)
\end{aligned}$$

2 GP Logistic Model

2.1 Model specifications

For logistic models, we first postulate a link function, $g(\cdot)$ for the predictors.

$$\begin{aligned}
y &= g^{-1}(\eta) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \gamma^2 I_n) \\
g(\mathbb{E}[y]) &= \eta \\
\eta &= f(x) + A\beta \\
g^{-1}(x) &= e^x / (1 + e^x) \\
f &\sim \mathcal{GP}(m(\cdot), \kappa(\cdot, \cdot))
\end{aligned}$$

Since y is a Bernoulli random variable, $\mathbb{E}[y] = \mathbb{P}(y=1)$. Without loss of generality, we will assume the mean function to be zero and that the Gaussian process has a sparse approximation representation as in Tan & Nott. Therefore,

$$f(x) \approx \sum_{r=1}^m \left[a_r \cos \left\{ (s_r \odot x)^T \lambda \right\} + b_r \sin \left\{ (s_r \odot x)^T \lambda \right\} \right]$$

and this could further be represented in matrix form which reduces this to a linear model.

$$f(x) = Z\alpha$$

$$\begin{aligned}\eta &= Z\alpha + A\beta \\ y &= \exp \{ (Z\alpha + A\beta) - \log (\mathbf{1} + \exp \{ Z\alpha + A\beta \}) \} + \epsilon\end{aligned}$$

Every scalar function applied to a vector or a matrix is done so elementwise. We think of the following priors:

$$\begin{aligned}\alpha | \sigma &\sim \mathcal{N} \left(0, \frac{\sigma^2}{m} I_{2m} \right) \\ \beta &\sim \mathcal{N} (\mu_\beta, \Sigma_\beta) \\ \lambda &\sim \mathcal{N} (\mu_\lambda, \Sigma_\lambda) \\ \sigma &\sim \text{half-Cauchy} (A_\sigma) \\ \gamma &\sim \text{half-Cauchy} (A_\gamma) \\ \theta &= (\alpha, \beta, \lambda, \sigma, \gamma)\end{aligned}$$

$$\begin{aligned}\log p(y, \theta) &= y^T (Z\alpha + A\beta) - \mathbf{1}_n^T \log (\mathbf{1}_n + \exp \{ Z\alpha + A\beta \}) - \left(m + \frac{s+d}{2} \right) \log (2\pi) - m \log \sigma^2 + m \log m \\ &\quad - \frac{m}{2\sigma^2} \alpha^T \alpha - \frac{1}{2} \log |\Sigma_\beta| - \frac{1}{2} (\beta - \mu_\beta)^T \Sigma_\beta^{-1} (\beta - \mu_\beta) - \frac{1}{2} \log |\Sigma_\lambda| - \frac{1}{2} (\lambda - \mu_\lambda)^T \Sigma_\lambda^{-1} (\lambda - \mu_\lambda) \\ &\quad + \log (2A_\sigma) + \log (2A_\gamma) 2 \log \pi - \log (A_\sigma^2 + \sigma^2) - \log (A_\gamma^2 + \gamma^2)\end{aligned}$$

Because $-\mathbf{1}_n^T \log (\mathbf{1}_n + \exp \{ Z\alpha + A\beta \})$ is analytically intractable for expectation which is essentially integration, we come up with the following approximation:

$$\begin{aligned}-\log (1 + e^x) &= \max_{\xi \in \mathbb{R}} \left\{ B(\xi) x^2 - \frac{1}{2} x + C(\xi) \right\}, \quad \forall x \in \mathbb{R} \\ B(\xi) &= -\tanh (\xi/2) / (4\xi) \\ C(\xi) &= \xi/2 - \log (1 + e^\xi) + \xi \tanh (\xi/2) / 4\end{aligned}$$

then

$$\begin{aligned}-\mathbf{1}_n^T \log \{ \mathbf{1}_n^T + \exp (Z\alpha + A\beta) \} &\geq \mathbf{1}_n^T \left\{ B(\xi) \odot (Z\alpha + A\beta)^2 - \frac{1}{2} (Z\alpha + A\beta) + C(\xi) \right\} \\ &= (Z\alpha + A\beta)^T \text{Dg} \{ B(\xi) \} (Z\alpha + A\beta) - \frac{1}{2} \mathbf{1}_n^T (Z\alpha + A\beta) + \mathbf{1}_n^T C(\xi),\end{aligned}$$

where $\xi = (\xi_1, \dots, \xi_n)$.

$$\begin{aligned}\log \underline{p}(y, \theta; \xi) &= y^T (Z\alpha + A\beta) + (Z\alpha + A\beta)^T \text{Dg} \{ B(\xi) \} (Z\alpha + A\beta) - \frac{1}{2} \mathbf{1}_n^T (Z\alpha + A\beta) + \mathbf{1}_n^T C(\xi) \\ &\quad - \left(m + \frac{s+d}{2} \right) \log (2\pi) - m \log \sigma^2 + m \log m - \frac{m}{2\sigma^2} \alpha^T \alpha - \frac{1}{2} \log |\Sigma_\beta| - \frac{1}{2} (\beta - \mu_\beta)^T \Sigma_\beta^{-1} (\beta - \mu_\beta) \\ &\quad - \frac{1}{2} \log |\Sigma_\lambda| - \frac{1}{2} (\lambda - \mu_\lambda)^T \Sigma_\lambda^{-1} (\lambda - \mu_\lambda) + \log (2A_\sigma) + \log (2A_\gamma) 2 \log \pi - \log (A_\sigma^2 + \sigma^2) \\ &\quad - \log (A_\gamma^2 + \gamma^2)\end{aligned}$$

3 Sparse GP Probit Regression

3.1 Bayesian probit regression without random effects

The ordinary formulation goes like this:

$$\mathbb{P}(y_i = 1|x_i, \beta) = \Phi(x_i^T \beta),$$

where Φ is the cumulative distribution of standard Gaussian distribution. Since the following likelihood

$$\mathbb{P}(Y = y|X, \beta) = \prod_{i=1}^n \Phi(x_i^T \beta)^{y_i} [1 - \Phi(x_i^T \beta)]^{1-y_i}$$

hinders the tractable calculation of the posterior, we devise the following latent variable:

$$Z_i = x_i^T \beta + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, 1)$$

for $i = 1, \dots, n$ and let

$$y_i = \begin{cases} 1, & \text{if } Z_i \geq 0 \\ 0, & \text{otherwise} \end{cases}.$$

It automatically follows that

$$\begin{aligned} Z_i|y_i = 0, x_i, \beta &\sim \mathcal{N}(x_i^T \beta, 1) \mathbf{1}[Z_i < 0] \\ Z_i|y_i = 1, x_i, \beta &\sim \mathcal{N}(x_i^T \beta, 1) \mathbf{1}[Z_i \geq 0] \end{aligned}$$

suggesting a truncated Gaussian for each Z_i depending on what value y_i takes on. According to probit regression without random effects, the introduction of latent variables enables the tractable computation of the posterior with regards to the parameters.

$$\pi(Z, y, \beta) = C\pi(\beta) \prod_{i=1}^n \{1[Z_i \geq 0] \mathbf{1}[y_i = 1] + 1[Z_i < 0] \mathbf{1}[y_i = 0]\} \phi(Z_i - x_i^T \beta)$$

where C is the proportionality constant and ϕ is the density of standard Gaussian. Following the mean-field approximation and computing the optimal density for Z_i , it becomes a truncated normal about zero with mean $x_i^T \beta$ and variance 1.

3.2 Incorporating random effects and nonparametric statistics

Taking one step forward, the random effects can be taken into consideration in many situations. Furthermore, we often do not know exactly what functional form the parameters should take on. Such situations greatly call for the use of nonparametric statistics.

$$\mathbb{P}(y_i = 1|f) = \Phi(f(x_i)).$$

In the previous section where linear probit regression was considered, the functional form was assumed to be linear, i.e. $x_i^T \beta$. However, we now insist that f could be any function. Adopting Gaussian process prior for f and using sparse approximation,

$$\mathbb{P}(y_i = 1|f) = \Phi(Z_i^T \alpha).$$

Note that Z_i is different from the one we used in the previous section. Moreover, we will assume there exist additive random effects:

$$\mathbb{P}(y_i = 1|f, \alpha, \beta) = \Phi(z_i^T \alpha + a_i^T \beta).$$

Assume latent variables

$$y_i^* = z_i^T \alpha + a_i^T \beta + \epsilon, \quad \epsilon \sim \mathcal{N}(0, 1).$$

Let

$$y_i = \begin{cases} 1, & \text{if } y_i^* \geq 0 \\ 0, & \text{if } y_i^* < 0. \end{cases}$$

Let's assume priors

$$\begin{aligned} \alpha|\sigma &\sim \mathcal{N}\left(0, \frac{\sigma^2}{m} I_{2m}\right) \\ \beta &\sim \mathcal{N}(\mu_\beta, \Sigma_\beta) \\ \lambda &\sim \mathcal{N}(\mu_\lambda, \Sigma_\lambda) \\ \sigma &\sim \text{half-Cauchy}(A_\sigma) \\ \theta &= (\alpha, \beta, \lambda, \sigma) \end{aligned}$$

$$\pi(y, y^*, \theta) = C \pi(\alpha|\sigma) \pi(\beta) \pi(\lambda) \pi(\sigma) \prod_{i=1}^n \{1[y_i^* \geq 0] 1[y_i = 1] + 1[y_i^* < 0] 1[y_i = 0]\} \phi(y_i^* - z_i^T \alpha - a_i^T \beta)$$

3.2.1 $q(\alpha)$

Let $y^* = [y_1^* \dots y_n^*]^T$.

$$\begin{aligned} \log q(\alpha) &\propto \left\langle \log \pi(\alpha|\sigma) \prod_{i=1}^n \phi(y_i^* - z_i^T \alpha - a_i^T \beta) \right\rangle \\ &\propto -\frac{m}{2} \left\langle \frac{1}{\sigma^2} \right\rangle \alpha^T \alpha + \left\langle \sum_{i=1}^n \left[-\frac{1}{2} \left\{ y_i^{*2} - 2(z_i^T \alpha + a_i^T \beta) y_i^* + (z_i^T \alpha + a_i^T \beta)^2 \right\} \right] \right\rangle \\ &\propto -\frac{m}{2} \left\langle \frac{1}{\sigma^2} \right\rangle \alpha^T \alpha + \left\langle -\frac{1}{2} \sum_{i=1}^n [-2y_i^* z_i^T \alpha + \alpha^T z_i z_i^T \alpha + 2\beta^T a_i z_i^T \alpha] \right\rangle \\ &\propto -\frac{m}{2} \left\langle \frac{1}{\sigma^2} \right\rangle \alpha^T \alpha - \frac{1}{2} \left\langle -2y^{*T} Z \alpha + \alpha^T \left(\sum_{i=1}^n z_i z_i^T \right) \alpha + 2\beta^T \left(\sum_{i=1}^n a_i z_i^T \right) \alpha \right\rangle \\ &\propto -\frac{m}{2} \left\langle \frac{1}{\sigma^2} \right\rangle \alpha^T \alpha - \frac{1}{2} \left\{ -2 \langle y^{*T} \rangle \langle Z \rangle \alpha + \alpha^T \langle Z^T Z \rangle \alpha + 2 \langle \beta \rangle^T A^T \langle Z \rangle \alpha \right\} \\ &\propto -\frac{1}{2} \left\{ \alpha^T \left(m \left\langle \frac{1}{\sigma^2} \right\rangle I_{2m} + \langle Z^T Z \rangle \right) \alpha - 2 \left(\langle Z \rangle^T \langle y^* \rangle - \langle Z \rangle^T A^T \langle \beta \rangle \right)^T \alpha \right\} \\ q(\alpha) &= \mathcal{N}(\mu_{q(\alpha)}, \Sigma_{q(\alpha)}) \\ \mu_{q(\alpha)} &= \Sigma_{q(\alpha)} \left(\langle Z \rangle^T \langle y^* \rangle - \langle Z \rangle^T A^T \langle \beta \rangle \right) \\ \Sigma_{q(\alpha)} &= \left(m \left\langle \frac{1}{\sigma^2} \right\rangle I_{2m} + \langle Z^T Z \rangle \right)^{-1} \end{aligned}$$

3.2.2 $q(\beta)$

$$\begin{aligned}
\log q(\beta) &\propto \log \pi(\beta) \prod_{i=1}^n \phi(y_i^* - z_i^T \alpha - a_i^T \beta) \\
&\propto -\frac{s}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_\beta| - \frac{1}{2} (\beta - \mu_\beta)^T \Sigma_\beta^{-1} (\beta - \mu_\beta) \\
&\quad + \left\langle \sum_{i=1}^T \left[-\frac{1}{2} \left(y_i^* - 2(z_i^T \alpha + a_i^T \beta) y_i^* + (z_i^T \alpha + a_i^T \beta)^2 \right) \right] \right\rangle \\
&\propto -\frac{1}{2} \left(\beta^T \Sigma_\beta^{-1} \beta - 2\mu_\beta^T \Sigma_\beta^{-1} \beta \right) + \left\langle \sum_{i=1}^n \left[-\frac{1}{2} (-2y_i^* a_i^T \beta + 2\alpha^T z_i a_i^T \beta + \beta^T a_i a_i^T \beta) \right] \right\rangle \\
&\propto -\frac{1}{2} \left(\beta^T \Sigma_\beta^{-1} \beta - 2\mu_\beta^T \Sigma_\beta^{-1} \beta \right) - \frac{1}{2} \left(-2 \langle y^{*T} \rangle A \beta + 2 \langle \alpha^T \rangle \langle Z^T \rangle A \beta + \beta^T A^T A \beta \right) \\
&\propto -\frac{1}{2} \left\{ \beta^T \Sigma_\beta^{-1} \beta + \beta^T A^T A \beta - 2\mu_\beta^T \Sigma_\beta^{-1} \beta + 2 \langle \alpha^T \rangle \langle Z^T \rangle A \beta - 2 \langle y^{*T} \rangle A \beta \right\} \\
&\propto -\frac{1}{2} \left\{ \beta^T \left(\Sigma_\beta^{-1} + A^T A \right) \beta - 2 \left(\Sigma_\beta^{-1} \mu_\beta + A^T \langle y^* \rangle - A^T \langle Z \rangle \langle \alpha \rangle \right)^T \beta \right\} \\
q(\beta) &= \mathcal{N}(\mu_{q(\beta)}, \Sigma_{q(\beta)}) \\
\mu_{q(\beta)} &= \Sigma_{q(\beta)} \left(\Sigma_\beta^{-1} \mu_\beta + A^T \langle y^* \rangle - A^T \langle Z \rangle \langle \alpha \rangle \right) \\
\Sigma_{q(\beta)} &= \left(\Sigma_\beta^{-1} + A^T A \right)^{-1}
\end{aligned}$$

3.2.3 $q(y^*)$

$$\begin{aligned}
\log q(y^*) &\propto \left\langle \log \prod_{i=1}^n \{1[y_i^* \geq 0] 1[y_i = 1] + 1[y_i^* < 0] 1[y_i = 0]\} \phi(y_i^* - z_i^T \alpha - a_i^T \beta) \right\rangle \\
&\propto \sum_{i=1}^n \log \{1[y_i^* \geq 0] 1[y_i = 1] + 1[y_i^* < 0] 1[y_i = 0]\} \\
&\quad + \left\langle \sum_{i=1}^n \left[-\frac{1}{2} \left(y_i^{*2} - 2(z_i^T \alpha + a_i^T \beta) y_i^* + (z_i^T \alpha + a_i^T \beta)^2 \right) \right] \right\rangle \\
&\propto \sum_{i=1}^n \log \{1[y_i^* \geq 0] 1[y_i = 1] + 1[y_i^* < 0] 1[y_i = 0]\} \\
&\quad - \frac{1}{2} \left(y^{*T} y^* - 2(\langle Z \rangle \langle \alpha \rangle + A \langle \beta \rangle)^T y^* \right) \\
q(y^*) &= \mathcal{TN}(\langle Z \rangle \langle \alpha \rangle + A \langle \beta \rangle, I_n)
\end{aligned}$$

where \mathcal{TN} indicates truncated normal distribution, in this case multivariate. Each element of y^* is truncated at zero. The direction of truncation remains the same. Assuming $y_i = 1$,

$$\mu_{q(y_i^*)} = \langle z_i^T \rangle \mu_{q(\alpha)} + a_i^T \mu_{q(\beta)} + \frac{\phi(\langle z_i^T \rangle \mu_{q(\alpha)} + a_i^T \mu_{q(\beta)})}{\Phi(\langle z_i^T \rangle \mu_{q(\alpha)} + a_i^T \mu_{q(\beta)})}$$

3.2.4 $q(\sigma)$

$$\begin{aligned}\log q(\sigma) &\propto \langle \log \pi(\alpha|\sigma) \pi(\sigma) \rangle \\ &\propto -m \log \sigma^2 - \frac{m}{2} \langle \alpha^T \alpha \rangle / \sigma^2 - \log(A_\sigma^2 + \sigma^2) \\ q(\sigma) &\propto \frac{\exp(-C_\sigma / \sigma^2)}{\sigma^{2m} (A_\sigma^2 + \sigma^2)} \\ C_\sigma &= \frac{m}{2} \left(\text{Tr}(\Sigma_{q(\alpha)}) + \mu_{q(\alpha)}^T \mu_{q(\alpha)} \right)\end{aligned}$$

3.2.5 LB: $\mathbb{E}[\log p(y, y^*|\theta)]$

Recall that

$$\begin{aligned}\pi(y, y^*, \theta) &= C \pi(\alpha|\sigma) \pi(\beta) \pi(\lambda) \pi(\sigma) \prod_{i=1}^n \{1[y_i^* \geq 0] 1[y_i = 1] + 1[y_i^* < 0] 1[y_i = 0]\} \phi(y_i^* - z_i^T \alpha - a_i^T \beta) \\ \mathcal{L} &= \mathbb{E}[\log \pi(y, y^*, \theta)] - \mathbb{E}[\log q(y^*) q(\theta)] \\ &= \mathbb{E}[\log p(y|y^*, \theta)] + \mathbb{E}[\log p(y^*|\theta)] + \mathbb{E}[\log p(\theta)] - \mathbb{E}[\log q(\theta)]\end{aligned}$$

The above equation includes all the nodes in the graphical model. Therefore, we compute the expectation of the logarithm of each node.

$$\begin{aligned}\log p(y, y^*|\theta) &= \log \prod_{i=1}^n \{1[y_i^* \geq 0] 1[y_i = 1] + 1[y_i^* < 0] 1[y_i = 0]\} \phi(y_i^* - z_i^T \alpha - a_i^T \beta) \\ &\implies \mathbb{E} \left[\sum_{i=1}^n \log \{1[y_i^* \geq 0] 1[y_i = 1] + 1[y_i^* < 0] 1[y_i = 0]\} \right] + \sum_{i=1}^n \mathbb{E} [\log \phi(y_i^* - z_i^T \alpha - a_i^T \beta)]\end{aligned}$$

If we set $U = \{1[y_i^* \geq 0] 1[y_i = 1] + 1[y_i^* < 0] 1[y_i = 0]\}$ and let it be a new random variable, it is always 1 with probability 1. Therefore, the logarithm of U , i.e. $\log U$ is always 0 which terminates the need for calculating the expected value since it's 0 regardless. We proceed to the remaining terms.

$$\sum_{i=1}^n \mathbb{E} [\log \phi(y_i^* - z_i^T \alpha - a_i^T \beta)] \quad \dots (1)$$

The entropy of the variational distribution of $q(y_i^*)$ should be coupled with eqn (1) to simplify the calculation.

$$\begin{aligned}
& \sum_{i=1}^n \mathbb{E} [\log \phi(y_i^* - z_i^T \alpha - a_i^T \beta) - \log \phi(y_i^* - \langle z_i^T \rangle \mu_{q(\alpha)} - a_i^T \mu_{q(\beta)})] \\
& + \sum_{i=1}^n \log \left(\left\{ \Phi \left(\langle z_i \rangle^T \mu_{q(\alpha)} + a_i^T \mu_{q(\beta)} \right) \right\}^{y_i} \left\{ 1 - \Phi \left(\langle z_i \rangle^T \mu_{q(\alpha)} + a_i^T \mu_{q(\beta)} \right) \right\}^{1-y_i} \right) \\
& = \sum_{i=1}^n \mathbb{E} \left[-\frac{1}{2} \left(-2 (z_i^T \alpha + a_i^T \beta) y_i^* + (z_i^T \alpha + a_i^T \beta)^2 + 2 \left(\langle z_i \rangle^T \mu_{q(\alpha)} + a_i^T \mu_{q(\beta)} \right) y_i^* \right) \right. \\
& \quad \left. + \frac{1}{2} \sum_{i=1}^n \left(\left(\langle z_i \rangle^T \mu_{q(\alpha)} + a_i^T \mu_{q(\beta)} \right)^2 \right) \right. \\
& \quad \left. + \sum_{i=1}^n \log \left(\left\{ \Phi \left(\langle z_i \rangle^T \mu_{q(\alpha)} + a_i^T \mu_{q(\beta)} \right) \right\}^{y_i} \left\{ 1 - \Phi \left(\langle z_i \rangle^T \mu_{q(\alpha)} + a_i^T \mu_{q(\beta)} \right) \right\}^{1-y_i} \right) \right]
\end{aligned}$$

$$\begin{aligned}
& = -\frac{1}{2} \left(\text{Tr} \left(\langle Z^T Z \rangle \Sigma_{q(\alpha)} \right) + \mu_{q(\alpha)}^T \left(\langle Z^T Z \rangle - \langle Z \rangle^T \langle Z \rangle \right) \mu_{q(\alpha)} + \text{Tr} \left(A^T A \Sigma_{q(\beta)} \right) \right) \\
& \quad + \sum_{i=1}^n \log \left(\left\{ \Phi \left(\langle z_i \rangle^T \mu_{q(\alpha)} + a_i^T \mu_{q(\beta)} \right) \right\}^{y_i} \left\{ 1 - \Phi \left(\langle z_i \rangle^T \mu_{q(\alpha)} + a_i^T \mu_{q(\beta)} \right) \right\}^{1-y_i} \right)
\end{aligned}$$