

Supplement for:
**Fast and Accurate Binary Response Mixed Model
Analysis via Expectation Propagation**

BY P. HALL¹, I.M. JOHNSTONE², J.T. ORMEROD³, M.P. WAND⁴ AND J.C.F. YU⁴

¹*University of Melbourne*, ²*Stanford University*, ³*University of Sydney*
and ⁴*University of Technology Sydney*

S.1 Proof of Theorem 1

For $\mathbf{x} \in \mathbb{R}^d$, define

$$\phi_{\Sigma}(\mathbf{x}) \equiv (2\pi)^{-d/2} |\Sigma|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{x}^T \Sigma^{-1} \mathbf{x}\right)$$

so that

$$\phi_{\mathbf{I}}(\mathbf{x}) = (2\pi)^{-d/2} \exp\left(-\frac{1}{2}\mathbf{x}^T \mathbf{x}\right).$$

We continue to use an unadorned ϕ to denote the Univariate Normal density function:

$$\phi(x) = (2\pi)^{-1/2} \exp(-\frac{1}{2}x^2).$$

The notation $\|\mathbf{v}\| = \sqrt{\mathbf{v}^T \mathbf{v}}$ for a column vector \mathbf{v} is also used.

Lemma 1. *For any function $g : \mathbb{R} \rightarrow \mathbb{R}$ and $d \times 1$ vectors $\boldsymbol{\alpha}_1$, $\boldsymbol{\alpha}_2$ and $\boldsymbol{\alpha}_3$ such that the integrals exist:*

$$\int_{\mathbb{R}^d} g(\boldsymbol{\alpha}_1^T \mathbf{x}) \phi_{\mathbf{I}}(\mathbf{x}) d\mathbf{x} = \int_{-\infty}^{\infty} g(\|\boldsymbol{\alpha}_1\| z) \phi(z) dz, \quad (\text{S.1})$$

$$\int_{\mathbb{R}^d} g(\boldsymbol{\alpha}_1^T \mathbf{x})(\boldsymbol{\alpha}_2^T \mathbf{x}) \phi_{\mathbf{I}}(\mathbf{x}) d\mathbf{x} = \{(\boldsymbol{\alpha}_1^T \boldsymbol{\alpha}_2)/\|\boldsymbol{\alpha}_1\|\} \int_{-\infty}^{\infty} z g(\|\boldsymbol{\alpha}_1\| z) \phi(z) dz \quad (\text{S.2})$$

$$\begin{aligned} \text{and } \int_{\mathbb{R}^d} g(\boldsymbol{\alpha}_1^T \mathbf{x})(\boldsymbol{\alpha}_2^T \mathbf{x})(\boldsymbol{\alpha}_3^T \mathbf{x}) \phi_{\mathbf{I}}(\mathbf{x}) d\mathbf{x} &= (\boldsymbol{\alpha}_2^T \boldsymbol{\alpha}_3) \int_{-\infty}^{\infty} g(\|\boldsymbol{\alpha}_1\| z) \phi(z) dz \\ &\quad + \{(\boldsymbol{\alpha}_1^T \boldsymbol{\alpha}_2)(\boldsymbol{\alpha}_1^T \boldsymbol{\alpha}_3)/\|\boldsymbol{\alpha}_1\|^2\} \int_{-\infty}^{\infty} (z^2 - 1) g(\|\boldsymbol{\alpha}_1\| z) \phi(z) dz. \end{aligned} \quad (\text{S.3})$$

Proof of Lemma 1. Note that the integrals on the left-hand side in (S.1)–(S.3) are, respectively,

$$E\{g(\boldsymbol{\alpha}_1^T \mathbf{x})\}, E\{g(\boldsymbol{\alpha}_1^T \mathbf{x})(\boldsymbol{\alpha}_2^T \mathbf{x})\} \quad \text{and} \quad E\{g(\boldsymbol{\alpha}_1^T \mathbf{x})(\boldsymbol{\alpha}_2^T \mathbf{x})(\boldsymbol{\alpha}_3^T \mathbf{x})\}$$

where

$$\mathbf{x} \sim N(\mathbf{0}_d, \mathbf{I}_d).$$

We now focus on simplification of the third integral (S.3). Simplification of the first and second integrals is similar and simpler. Make the change of variables

$$\mathbf{s} \equiv \begin{bmatrix} s_1 \\ s_2 \\ s_3 \end{bmatrix} = \mathbf{A}\mathbf{x} \quad \text{where} \quad \mathbf{A} \equiv \begin{bmatrix} \boldsymbol{\alpha}_1^T \\ \boldsymbol{\alpha}_2^T \\ \boldsymbol{\alpha}_3^T \end{bmatrix}$$

so that

$$E\{g(\boldsymbol{\alpha}_1^T \mathbf{x})(\boldsymbol{\alpha}_2^T \mathbf{x})(\boldsymbol{\alpha}_3^T \mathbf{x})\} = E\{g(s_1)s_2s_3\} \quad \text{where } \mathbf{s} \sim N(\mathbf{0}_3, \mathbf{A}\mathbf{A}^T).$$

We then note that,

$$\begin{aligned} E\{g(s_1)s_2s_3\} &= \int_{-\infty}^{\infty} g(s_1) \left\{ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} s_2 s_3 p(s_2 s_3 | s_1) ds_2 ds_3 \right\} p(s_1) ds_1 \\ &= \int_{-\infty}^{\infty} g(s_1) \left\{ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \{\text{Cov}(s_2, s_3 | s_1) + E(s_2 | s_1)E(s_3 | s_1)\} ds_2 ds_3 \right\} p(s_1) ds_1 \end{aligned}$$

and make use of the result (see e.g. Theorem 3.2.4 of Mardia, Kent & Bibby, 1979)

$$\begin{bmatrix} s_2 \\ s_3 \end{bmatrix} \mid s_1 \sim N \left((s_1 / \|\boldsymbol{\alpha}_1\|^2) \begin{bmatrix} \boldsymbol{\alpha}_1^T \boldsymbol{\alpha}_2 \\ \boldsymbol{\alpha}_1^T \boldsymbol{\alpha}_3 \end{bmatrix}, \begin{bmatrix} \|\boldsymbol{\alpha}_2\|^2 & \boldsymbol{\alpha}_2^T \boldsymbol{\alpha}_3 \\ \boldsymbol{\alpha}_2^T \boldsymbol{\alpha}_3 & \|\boldsymbol{\alpha}_3\|^2 \end{bmatrix} - (1 / \|\boldsymbol{\alpha}_1\|^2) \begin{bmatrix} (\boldsymbol{\alpha}_1^T \boldsymbol{\alpha}_2)^2 & (\boldsymbol{\alpha}_1^T \boldsymbol{\alpha}_2)(\boldsymbol{\alpha}_1^T \boldsymbol{\alpha}_3) \\ (\boldsymbol{\alpha}_1^T \boldsymbol{\alpha}_2)(\boldsymbol{\alpha}_1^T \boldsymbol{\alpha}_3) & (\boldsymbol{\alpha}_1^T \boldsymbol{\alpha}_3)^2 \end{bmatrix} \right).$$

Result (S.3) then follows via simple algebraic manipulations.

Lemma 2. For all $a \in \mathbb{R}$ and $d \times 1$ vectors \mathbf{b}

$$\int_{\mathbb{R}^d} \Phi(a + \mathbf{b}^T \mathbf{x}) \phi_I(\mathbf{x}) d\mathbf{x} = \Phi \left(\frac{a}{\sqrt{\mathbf{b}^T \mathbf{b} + 1}} \right), \quad (\text{S.4})$$

$$\int_{\mathbb{R}^d} \mathbf{x} \Phi(a + \mathbf{b}^T \mathbf{x}) \phi_I(\mathbf{x}) d\mathbf{x} = \frac{\mathbf{b}}{\sqrt{\mathbf{b}^T \mathbf{b} + 1}} \phi \left(\frac{a}{\sqrt{\mathbf{b}^T \mathbf{b} + 1}} \right) \quad \text{and} \quad (\text{S.5})$$

$$\int_{\mathbb{R}^d} \mathbf{x} \mathbf{x}^T \Phi(a + \mathbf{b}^T \mathbf{x}) \phi_I(\mathbf{x}) d\mathbf{x} = \Phi \left(\frac{a}{\sqrt{\mathbf{b}^T \mathbf{b} + 1}} \right) \mathbf{I} - \frac{a \mathbf{b} \mathbf{b}^T}{\sqrt{(\mathbf{b}^T \mathbf{b} + 1)^3}} \phi \left(\frac{a}{\sqrt{\mathbf{b}^T \mathbf{b} + 1}} \right). \quad (\text{S.6})$$

Proof of Lemma 2.

Suppose that Z_1 and Z_2 are independent $N(0, 1)$ random variables. As defined in Section 4.2, for a logical proposition \mathcal{P} , let $I(\mathcal{P}) = 1$ if \mathcal{P} is true and $I(\mathcal{P}) = 0$ if \mathcal{P} is false. Then

$$P(Z_1 \leq a + \|\mathbf{b}\| Z_2) = E\{I(Z_1 \leq a + \|\mathbf{b}\| Z_2)\} = E[E\{I(Z_1 \leq a + \|\mathbf{b}\| Z_2) | Z_2\}] = E\{h(Z_2)\}$$

where $h(z_2) \equiv E\{I(Z_1 \leq a + \|\mathbf{b}\| Z_2) | Z_2 = z_2\}$. But note that

$$h(z_2) = P(Z_1 \leq a + \|\mathbf{b}\| z_2) = \Phi(a + \|\mathbf{b}\| z_2)$$

which implies that

$$P(Z_1 \leq a + \|\mathbf{b}\| Z_2) = E\{\Phi(a + \|\mathbf{b}\| Z_2)\} = \int_{-\infty}^{\infty} \Phi(a + \|\mathbf{b}\| z) \phi(z) dz.$$

Then

$$P(Z_1 \leq a + \|\mathbf{b}\| Z_2) = P(Z_1 - \|\mathbf{b}\| Z_2 \leq a) = P(X_3 \leq a)$$

where $X_3 \equiv Z_1 - \|\mathbf{b}\| Z_2 \sim N(0, 1 + \mathbf{b}^T \mathbf{b})$ by independence of Z_1 and Z_2 . Then (S.4) follows immediately.

Next, let e_i denote the $d \times 1$ vector with i th entry equal to 1 and with zeroes elsewhere. Then, using Lemma 1, the i th entry of the right-hand side of (S.5) is

$$\begin{aligned} \int_{\mathbb{R}^d} (\mathbf{e}_i^T \mathbf{x}) \Phi(a + \mathbf{b}^T \mathbf{x}) \phi_{\mathbf{I}}(\mathbf{x}) d\mathbf{x} &= \{(\mathbf{e}_i^T \mathbf{b}) / \|\mathbf{b}\|\} \int_{-\infty}^{\infty} z \Phi(a + \|\mathbf{b}\|z) \phi(z) dz \\ &= -\{(\mathbf{e}_i^T \mathbf{b}) / \|\mathbf{b}\|\} \int_{-\infty}^{\infty} \Phi(a + \|\mathbf{b}\|z) \phi'(z) dz \\ &= (\mathbf{e}_i^T \mathbf{b}) \int_{-\infty}^{\infty} \phi(a + \|\mathbf{b}\|z) \phi(z) dz \end{aligned}$$

where the last result follows via integration by parts. The last integrand is

$$(2\pi)^{-1} \exp\{-\frac{1}{2}(a + \|\mathbf{b}\|z)^2 - \frac{1}{2}z^2\} = \phi\left(\frac{a}{\sqrt{\mathbf{b}^T \mathbf{b} + 1}}\right) \phi\left(\frac{z + a\|\mathbf{b}\|/(1 + \mathbf{b}^T \mathbf{b})}{1/\sqrt{\mathbf{b}^T \mathbf{b} + 1}}\right)$$

and (S.5) is an immediate consequence.

Lastly, because of (S.3), the (i, j) entry of the right-hand side of (S.6) is

$$\begin{aligned} \int_{\mathbb{R}^d} (\mathbf{e}_i^T \mathbf{x})(\mathbf{e}_j^T \mathbf{x}) \Phi(a + \mathbf{b}^T \mathbf{x}) \phi_{\mathbf{I}}(\mathbf{x}) d\mathbf{x} &= (\mathbf{e}_i^T \mathbf{e}_j) \int_{-\infty}^{\infty} \Phi(a + \|\mathbf{b}\|z) \phi(z) dz + \frac{(\mathbf{e}_i^T \mathbf{b})(\mathbf{e}_j^T \mathbf{b})}{\|\mathbf{b}\|^2} \int_{-\infty}^{\infty} \Phi(a + \|\mathbf{b}\|z) \phi''(z) dz \\ &= (\mathbf{e}_i^T \mathbf{e}_j) \Phi\left(\frac{a}{\sqrt{\mathbf{b}^T \mathbf{b} + 1}}\right) - \frac{(\mathbf{e}_i^T \mathbf{b})(\mathbf{e}_j^T \mathbf{b})}{\|\mathbf{b}\|} \int_{-\infty}^{\infty} \phi(a + \|\mathbf{b}\|z) \phi'(z) dz \end{aligned}$$

where the last result follows via integration by parts. The last integrand is

$$-(2\pi)^{-1} z \exp\{-\frac{1}{2}(a + \|\mathbf{b}\|z)^2 - \frac{1}{2}z^2\} = -z \phi\left(\frac{a}{\sqrt{\mathbf{b}^T \mathbf{b} + 1}}\right) \phi\left(\frac{z + a\|\mathbf{b}\|/(1 + \mathbf{b}^T \mathbf{b})}{1/\sqrt{\mathbf{b}^T \mathbf{b} + 1}}\right)$$

and (S.6) follows. ■

Next, we note a key connection between Kullback-Leibler projection onto the unnormalized and normalized Multivariate Normal families. For the latter, we introduce the notation

$$\text{proj}_{\mathbb{N}}[p](\mathbf{x}) = q(\mathbf{x})$$

where q is the Multivariate Normal density function minimizes $\text{KL}(p\|q)$.

Lemma 3. Let $f \in L_1(\mathbb{R}^d)$ be such that $f \geq 0$ and define $C_f \equiv \int_{\mathbb{R}^d} f(\mathbf{x}) d\mathbf{x}$. Then

$$\text{proj}[f](\mathbf{x}) = C_f \text{proj}_{\mathbb{N}}[f/C_f](\mathbf{x}).$$

Proof Lemma 3.

Let $g(\cdot; \boldsymbol{\eta})$ be a generic unnormalized Multivariate Normal density function with natural parameter vector $\boldsymbol{\eta}$:

$$g(\mathbf{x}; \boldsymbol{\eta}) = \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^T) \end{bmatrix}^T \boldsymbol{\eta} \right\}.$$

Then the Kullback-Leibler divergence of $g(\cdot; \boldsymbol{\eta})$ from f is

$$\text{KL}(f \| g(\cdot; \boldsymbol{\eta})) = \int_{\mathbb{R}^d} [f(\mathbf{x}) \log\{f(\mathbf{x})/g(\mathbf{x}; \boldsymbol{\eta})\} + g(\mathbf{x}; \boldsymbol{\eta}) - f(\mathbf{x})] d\mathbf{x} = \mathcal{K}(\boldsymbol{\eta}) + \text{const}$$

where ‘const’ denotes terms not depending on $\boldsymbol{\eta}$ and

$$\mathcal{K}(\boldsymbol{\eta}) \equiv (2\pi)^{d/2} \exp\{\boldsymbol{\eta}_0 + A_N(\boldsymbol{\eta}_{-0})\} - \begin{bmatrix} \int_{\mathbb{R}^d} f(\mathbf{x}) d\mathbf{x} \\ \int_{\mathbb{R}^d} \mathbf{x} f(\mathbf{x}) d\mathbf{x} \\ \int_{\mathbb{R}^d} \text{vech}(\mathbf{x}\mathbf{x}^T) f(\mathbf{x}) d\mathbf{x} \end{bmatrix}^T \boldsymbol{\eta}.$$

The derivative vector of $\mathcal{K}(\boldsymbol{\eta})$ is

$$D\mathcal{K}(\boldsymbol{\eta}) = (2\pi)^{d/2} \exp\{\boldsymbol{\eta}_0 + A_N(\boldsymbol{\eta}_{-0})\} \begin{bmatrix} 1 \\ D A(\boldsymbol{\eta}_{-0})^T \end{bmatrix}^T - \begin{bmatrix} \int_{\mathbb{R}^d} f(\mathbf{x}) d\mathbf{x} \\ \int_{\mathbb{R}^d} \mathbf{x} f(\mathbf{x}) d\mathbf{x} \\ \int_{\mathbb{R}^d} \text{vech}(\mathbf{x}\mathbf{x}^T) f(\mathbf{x}) d\mathbf{x} \end{bmatrix}^T$$

so the stationary condition, $D\mathcal{K}(\boldsymbol{\eta})^T = \mathbf{0}$, for the minimization of $\text{KL}(f \| g(\cdot; \boldsymbol{\eta}))$ is

$$(2\pi)^{d/2} \exp\{\boldsymbol{\eta}_0 + A_N(\boldsymbol{\eta}_{-0})\} \begin{bmatrix} 1 \\ \nabla A_N(\boldsymbol{\eta}_{-0}) \end{bmatrix} = \begin{bmatrix} \int_{\mathbb{R}^d} f(\mathbf{x}) d\mathbf{x} \\ \int_{\mathbb{R}^d} \mathbf{x} f(\mathbf{x}) d\mathbf{x} \\ \int_{\mathbb{R}^d} \text{vech}(\mathbf{x}\mathbf{x}^T) f(\mathbf{x}) d\mathbf{x} \end{bmatrix}. \quad (\text{S.7})$$

with $\nabla A_N(\boldsymbol{\eta}_{-0}) \equiv D A(\boldsymbol{\eta}_{-0})^T$ denoting the gradient vector of $A(\boldsymbol{\eta}_{-0})$. It is easily checked that (S.7) is satisfied by

$$(\boldsymbol{\eta}^*)_0 = \log(C_f) - A_N(\boldsymbol{\eta}_{-0}^*) - \frac{1}{2} d \log(2\pi)$$

where $\boldsymbol{\eta}_{-0}^* = (\nabla A_N)^{-1} \left(\begin{bmatrix} \int_{\mathbb{R}^d} \mathbf{x} \{f(\mathbf{x})/C_f\} d\mathbf{x} \\ \int_{\mathbb{R}^d} \text{vech}(\mathbf{x}\mathbf{x}^T) \{f(\mathbf{x})/C_f\} d\mathbf{x} \end{bmatrix} \right)$

with existence and uniqueness of $(\nabla A_N)^{-1}$ being guaranteed by Proposition 3.2 of Wainwright & Jordan (2008). The Hessian matrix of $\mathcal{K}(\boldsymbol{\eta})$ is

$$H\mathcal{K}(\boldsymbol{\eta}) = (2\pi)^{d/2} e^{\boldsymbol{\eta}_0 + A_N(\boldsymbol{\eta}_{-0})} \left\{ \begin{bmatrix} 1 \\ \nabla A_N(\boldsymbol{\eta}_{-0}) \end{bmatrix} \begin{bmatrix} 1 \\ \nabla A_N(\boldsymbol{\eta}_{-0}) \end{bmatrix}^T + \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & H A_N(\boldsymbol{\eta}_{-0}) \end{bmatrix} \right\}.$$

From Proposition 3.1 of Wainwright & Jordan (2008), A_N is strictly convex on its domain and therefore $H A_N(\boldsymbol{\eta}_{-0})$ is positive definite. Hence $H\mathcal{K}(\boldsymbol{\eta})$ is positive definite for all $\boldsymbol{\eta}$ and so (S.8) is the unique minimizer of $\text{KL}(f \| g(\cdot; \boldsymbol{\eta}))$. Therefore,

$$\text{proj}[f](\mathbf{x}) = \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^T) \end{bmatrix}^T \boldsymbol{\eta}^* \right\}$$

where $\boldsymbol{\eta}^*$ is as given by (S.8). However, $\boldsymbol{\eta}_{-0}^*$ is the same natural parameter vector that arises via projection of f/C_f onto the family of Multivariate Normal density functions and so

$$\text{proj}_N[f/C_f](\mathbf{x}) = \exp \left\{ \begin{bmatrix} \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^T) \end{bmatrix}^T \boldsymbol{\eta}_{-0}^* - A_N(\boldsymbol{\eta}_{-0}^*) \right\} (2\pi)^{-d/2}$$

which immediately leads to Lemma 3. ■

The proof of Theorem 1 involves transferral between the common $N(\mu, \Sigma)$ parameters of the d -variate Normal distribution and the natural parameters corresponding to the sufficient statistics \mathbf{x} and $\text{vech}(\mathbf{x}\mathbf{x}^T)$. The transformations in each direction are

$$\begin{cases} \boldsymbol{\eta}_1 = \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} \\ \boldsymbol{\eta}_2 = -\frac{1}{2}\mathbf{D}_d^T \text{vec}(\boldsymbol{\Sigma}^{-1}) \end{cases} \quad \text{and} \quad \begin{cases} \boldsymbol{\mu} = -\frac{1}{2}\{\text{vec}^{-1}(\mathbf{D}_d^{+T}\boldsymbol{\eta}_2)\}^{-1}\boldsymbol{\eta}_1 \\ \boldsymbol{\Sigma} = -\frac{1}{2}\{\text{vec}^{-1}(\mathbf{D}_d^{+T}\boldsymbol{\eta}_2)\}^{-1} \end{cases} \quad (\text{S.9})$$

Recall the notation

$$\mathbf{v}^{\otimes k} \equiv \begin{cases} 1 & \text{for } k = 0 \\ \mathbf{v} & \text{for } k = 1 \\ \mathbf{v}\mathbf{v}^T & \text{for } k = 2 \end{cases}$$

and consider Kullback-Leibler projection of $f_{\text{input}}/C_{f_{\text{input}}}$ onto the family of d -variate Normal density functions where

$$f_{\text{input}}(\mathbf{x}) \equiv \Phi(c_0 + \mathbf{c}_1^T \mathbf{x}) \exp \left\{ \left[\begin{array}{c} \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^T) \end{array} \right]^T \begin{bmatrix} \boldsymbol{\eta}_1^{\text{input}} \\ \boldsymbol{\eta}_2^{\text{input}} \end{bmatrix} \right\},$$

and $C_{f_{\text{input}}} \equiv \int_{\mathbb{R}^d} f_{\text{input}}(\mathbf{x}) d\mathbf{x}$. Then the projection has mean and covariance matrix

$$\boldsymbol{\mu}^* = \mathcal{M}_1/\mathcal{M}_0 \quad \text{and} \quad \boldsymbol{\Sigma}^* = \mathcal{M}_2/\mathcal{M}_0 - (\mathcal{M}_1/\mathcal{M}_0)(\mathcal{M}_1/\mathcal{M}_0)^T \quad (\text{S.10})$$

where

$$\mathcal{M}_k \equiv \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \Phi(c_0 + \mathbf{c}_1^T \mathbf{x}) \exp \left\{ \left[\begin{array}{c} \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^T) \end{array} \right]^T \begin{bmatrix} \boldsymbol{\eta}_1^{\text{input}} \\ \boldsymbol{\eta}_2^{\text{input}} \end{bmatrix} \right\} d\mathbf{x}.$$

Letting

$$\boldsymbol{\Sigma}^{\text{input}} \equiv -\frac{1}{2}\{\text{vec}^{-1}(\mathbf{D}_d^{+T}\boldsymbol{\eta}_2^{\text{input}})\}^{-1} \quad \text{and} \quad \boldsymbol{\mu}^{\text{input}} \equiv \boldsymbol{\Sigma}^{\text{input}}\boldsymbol{\eta}_1^{\text{input}}$$

be the common parameters corresponding to $\boldsymbol{\eta}^{\text{input}}$ and making the change of variable $\mathbf{z} = (\boldsymbol{\Sigma}^{\text{input}})^{-1/2}(\mathbf{x} - \boldsymbol{\mu}^{\text{input}})$ we obtain

$$\begin{aligned} \mathcal{M}_k &= (2\pi)^{d/2} e^{A_N(\boldsymbol{\eta}^{\text{input}})} \int_{\mathbb{R}^d} (\boldsymbol{\mu}^{\text{input}} + (\boldsymbol{\Sigma}^{\text{input}})^{1/2}\mathbf{z})^{\otimes k} \\ &\quad \times \Phi((c_0 + \mathbf{c}_1^T \boldsymbol{\mu}^{\text{input}}) + \{(\boldsymbol{\Sigma}^{\text{input}})^{1/2}\mathbf{c}_1\}^T \mathbf{z}) \phi_I(\mathbf{z}) d\mathbf{z}. \end{aligned}$$

Lemma 2 and simple algebraic manipulations then give

$$\mathcal{M}_1/\mathcal{M}_0 = \boldsymbol{\mu}^{\text{input}} + \frac{\boldsymbol{\Sigma}^{\text{input}}\mathbf{c}_1\zeta'(r_2)}{\sqrt{\mathbf{c}_1^T \boldsymbol{\Sigma}^{\text{input}} \mathbf{c}_1 + 1}}.$$

and

$$\begin{aligned} \mathcal{M}_2/\mathcal{M}_0 &= \boldsymbol{\mu}^{\text{input}}(\boldsymbol{\mu}^{\text{input}})^T + \frac{\{\boldsymbol{\Sigma}^{\text{input}}\mathbf{c}_1(\boldsymbol{\mu}^{\text{input}})^T + \boldsymbol{\mu}^{\text{input}}\mathbf{c}_1^T \boldsymbol{\Sigma}^{\text{input}}\}\zeta'(r_2)}{\sqrt{\mathbf{c}_1^T \boldsymbol{\Sigma}^{\text{input}} \mathbf{c}_1 + 1}} \\ &\quad + \boldsymbol{\Sigma}^{\text{input}} - \frac{r\zeta'(r_2)\boldsymbol{\Sigma}^{\text{input}}\mathbf{c}_1\mathbf{c}_1^T \boldsymbol{\Sigma}^{\text{input}}}{\mathbf{c}_1^T \boldsymbol{\Sigma}^{\text{input}} \mathbf{c}_1 + 1}. \end{aligned}$$

where

$$r_2 \equiv \frac{2c_0 - \mathbf{c}_1^T \{\text{vec}^{-1}(\mathbf{D}_{d^R}^{+T} \boldsymbol{\eta}_2^{\text{input}})\}^{-1} \boldsymbol{\eta}_1^{\text{input}}}{\sqrt{2 \left[2 - \mathbf{c}_1^T \{\text{vec}^{-1}(\mathbf{D}_{d^R}^{+T} \boldsymbol{\eta}_2^{\text{input}})\}^{-1} \mathbf{c}_1 \right]}}.$$

Combining these last two results and noting (S.10) we obtain the common parameter solutions

$$\begin{aligned} \boldsymbol{\mu}^* &= \boldsymbol{\mu}^{\text{input}} + \frac{\boldsymbol{\Sigma}^{\text{input}} \mathbf{c}_1 \zeta'(r_2)}{\sqrt{\mathbf{c}_1^T \boldsymbol{\Sigma}^{\text{input}} \mathbf{c}_1 + 1}} \\ \boldsymbol{\Sigma}^* &= \boldsymbol{\Sigma}^{\text{input}} + \left\{ \frac{\zeta''(r_2)}{\mathbf{c}_1^T \boldsymbol{\Sigma}^{\text{input}} \mathbf{c}_1 + 1} \right\} \boldsymbol{\Sigma}^{\text{input}} \mathbf{c}_1 \mathbf{c}_1^T \boldsymbol{\Sigma}^{\text{input}}. \end{aligned}$$

Transferral to natural parameters via (S.9) and some simple manipulations then lead to

$$\begin{bmatrix} \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix} = K_{\text{probit}} \left(\begin{bmatrix} \boldsymbol{\eta}_1^{\text{input}} \\ \boldsymbol{\eta}_2^{\text{input}} \end{bmatrix}; c_0, \mathbf{c}_1 \right).$$

Finally,

$$\begin{aligned} \eta_0^* &= \log(C_{f_{\text{input}}}) - \log \int_{\mathbb{R}^{d^R}} \exp \left\{ \begin{bmatrix} \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^T) \end{bmatrix}^T \begin{bmatrix} \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix} \right\} d\mathbf{x} \\ &= \log(\mathcal{M}_0) - \frac{1}{2} d^R \log(2\pi) - A_N(\boldsymbol{\eta}^*) = \log \Phi(r_2) + A_N(\boldsymbol{\eta}^{\text{input}}) - A_N(\boldsymbol{\eta}^*) \\ &= C_{\text{probit}} \left(\begin{bmatrix} \boldsymbol{\eta}_1^{\text{input}} \\ \boldsymbol{\eta}_2^{\text{input}} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix}; c_0, \mathbf{c}_1 \right). \end{aligned}$$

S.2 Derivation of Algorithm 1

We now provide full justification of Algorithm 1, starting with a derivation of the message passing representation used in Algorithm 1.

S.2.1 Message Passing Representation Derivation

The derivation of the message passing representation is based on the infrastructure and results laid out in Minka (2005). The treatment given there is for a generalization of Kullback-Leibler divergence, known as α -divergence, and for approximation of (normalized) density functions rather than general non-negative L_1 functions. The Kullback-Leibler divergence minimization problem given by (8) corresponds to $\alpha = 1$ in the notation of Minka (2005). Following Section 4.1 of Minka (2005) we then define the messages passed from the factors neighboring \mathbf{u}_i in Figure 1 to be

$$m_{p(y_{ij}|\mathbf{u}_i;\beta)} \rightarrow \mathbf{u}_i(\mathbf{u}_i) \equiv p(y_{ij}|\mathbf{u}_i;\beta) \quad \text{and} \quad m_{p(\mathbf{u}_i;\Sigma)} \rightarrow \mathbf{u}_i(\mathbf{u}_i) \equiv p(\mathbf{u}_i;\Sigma). \quad (\text{S.11})$$

Then, (54) of Minka (2005) invokes the definition

$$m_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i;\beta)}(\mathbf{u}_i) \equiv m_{p(\mathbf{u}_i;\Sigma) \rightarrow \mathbf{u}_i(\mathbf{u}_i)} \prod_{j' \neq j} m_{p(y_{ij'}|\mathbf{u}_i;\beta)} \rightarrow \mathbf{u}_i(\mathbf{u}_i). \quad (\text{S.12})$$

Result (60) of Minka (2005) with $\alpha = 1$, $s' = 1$ (since we are working with unnormalized rather than normalized Kullback-Leibler divergence) and the simplification that

there is only one stochastic node, namely \mathbf{u}_i , provides the main factor to stochastic node message passing updates:

$$m_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) \leftarrow \frac{\text{proj} [m_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})}(\mathbf{u}_i) p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})] (\mathbf{u}_i)}{m_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})}(\mathbf{u}_i)}, \quad 1 \leq j \leq n_i. \quad (\text{S.13})$$

The other factor to stochastic node message passing update is, trivially from (S.11),

$$m_{p(\mathbf{u}_i; \boldsymbol{\Sigma}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) \leftarrow p(\mathbf{u}_i; \boldsymbol{\Sigma}).$$

The stochastic node to factor updates are, from (S.12),

$$m_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})}(\mathbf{u}_i) \leftarrow m_{p(\mathbf{u}_i; \boldsymbol{\Sigma}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) \prod_{j' \neq j} m_{p(y_{ij'}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i), \quad 1 \leq j \leq n_i.$$

Next, we simplify these message updates to a programmable form.

S.2.2 Simplification of the $m_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i)$ Updates

From (S.12) it is apparent that $m_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i)$ is an unnormalized Multivariate Normal density function and therefore

$$m_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})}(\mathbf{u}_i) = \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^T) \end{bmatrix}^T \eta_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})} \right\}$$

with natural parameter vector $\eta_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})}$. Introducing the abbreviation:

$$\boldsymbol{\eta}^\otimes \equiv \eta_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})}$$

we have

$$m_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})}(\mathbf{u}_i) = \exp(\eta_0^\otimes) \exp \left\{ \begin{bmatrix} \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^T) \end{bmatrix}^T \boldsymbol{\eta}_{-0}^\otimes \right\}$$

where η_0^\otimes denotes the first entry of $\boldsymbol{\eta}^\otimes$ and $\boldsymbol{\eta}_{-0}^\otimes$ contains the remaining entries. Substitution info (S.13) leads to

$$\begin{aligned} & m_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) \leftarrow \\ & \frac{\text{proj} \left[\exp(\eta_0^\otimes) \exp \left\{ \begin{bmatrix} \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^T) \end{bmatrix}^T \boldsymbol{\eta}_{-0}^\otimes \right\} \Phi((2y_{ij} - 1)(\boldsymbol{\beta}^T \mathbf{x}_{ij}^F + \mathbf{u}_i^T \mathbf{x}_{ij}^R)) \right] (\mathbf{u}_i)}{\exp(\eta_0^\otimes) \exp \left\{ \begin{bmatrix} \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^T) \end{bmatrix}^T \boldsymbol{\eta}_{-0}^\otimes \right\}} \\ &= \frac{\text{proj} \left[\Phi(c_{0,ij} + \mathbf{c}_{1,ij}^T \mathbf{u}_i) \exp \left\{ \begin{bmatrix} \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^T) \end{bmatrix}^T \boldsymbol{\eta}_{-0}^\otimes \right\} \right] (\mathbf{u}_i)}{\exp \left\{ \begin{bmatrix} \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^T) \end{bmatrix}^T \boldsymbol{\eta}_{-0}^\otimes \right\}} \end{aligned}$$

where

$$c_{0,ij} \equiv (2y_{ij} - 1)(\boldsymbol{\beta}^T \mathbf{x}_{ij}^F) \quad \text{and} \quad \mathbf{c}_{1,ij} \equiv (2y_{ij} - 1)\mathbf{x}_{ij}^R.$$

Using Theorem 1:

$$m_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) \leftarrow \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^T) \end{bmatrix}^T \eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i} \right\}$$

where the linear and quadratic coefficient updates are

$$\begin{aligned} (\eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i})_{-0} &\longleftarrow K_{\text{probit}} \left((\eta_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})})_{-0}; (2y_{ij} - 1)(\boldsymbol{\beta}^T \mathbf{x}_{ij}^F), (2y_{ij} - 1)\mathbf{x}_{ij}^R \right) \\ &\quad - (\eta_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})})_{-0} \end{aligned}$$

and the constant coefficient update is

$$\begin{aligned} (\eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i})_0 &\longleftarrow C_{\text{probit}} \left((\eta_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})})_{-0}, (\eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i})_{-0} \right. \\ &\quad \left. + (\eta_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})})_{-0}; (2y_{ij} - 1)(\boldsymbol{\beta}^T \mathbf{x}_{ij}^F), (2y_{ij} - 1)\mathbf{x}_{ij}^R \right). \end{aligned}$$

S.2.3 Simplification of the $m_{p(\mathbf{u}_i; \Sigma) \rightarrow \mathbf{u}_i}(\mathbf{u}_i)$ Update

The second definition in (S.11) gives

$$m_{p(\mathbf{u}_i; \Sigma) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) \longleftarrow p(\mathbf{u}_i; \Sigma) = \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^T) \end{bmatrix}^T \eta_{\Sigma} \right\}.$$

Therefore, if $\eta_{p(\mathbf{u}_i; \Sigma) \rightarrow \mathbf{u}_i}$ denotes the natural parameter vector of $m_{p(\mathbf{u}_i; \Sigma) \rightarrow \mathbf{u}_i}(\mathbf{u}_i)$ then it has the trivial update

$$\eta_{p(\mathbf{u}_i; \Sigma) \rightarrow \mathbf{u}_i} \longleftarrow \eta_{\Sigma}.$$

S.2.4 Simplification of the $m_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})}(\mathbf{u}_i)$ Updates

Given the simplified forms of the messages in the two previous subsections we have from (S.12):

$$\begin{aligned} m_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})}(\mathbf{u}_i) &\longleftarrow \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^T) \end{bmatrix}^T \eta_{p(\mathbf{u}_i; \Sigma) \rightarrow \mathbf{u}_i} \right\} \\ &\quad \times \prod_{j' \neq j} \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^T) \end{bmatrix}^T \eta_{p(y_{ij'}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i} \right\} \end{aligned}$$

which leads to

$$\begin{aligned} \eta_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})} &\longleftarrow \eta_{p(\mathbf{u}_i; \Sigma) \rightarrow \mathbf{u}_i} + \sum_{j' \neq j} \eta_{p(y_{ij'}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i} \\ &= \eta_{p(\mathbf{u}_i; \Sigma) \rightarrow \mathbf{u}_i} + \text{SUM}\{\eta_{p(y_i|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}\} - \eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}. \end{aligned}$$

S.2.5 Assembly of All Natural Parameter Updates

We now return to the message passing protocol given in Section 3.2:

Initialize all factor to stochastic node messages.

Cycle until all factor to stochastic node messages converge:

For each factor:

Compute the messages passed to the factor using (11) or (12).

Compute the messages passed from the factor using (9) or (10).

For the factors $p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})$:

computing the messages passed to each of these factors reduces to

$$\eta_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})} \leftarrow \eta_{p(\mathbf{u}_i; \Sigma) \rightarrow \mathbf{u}_i} + \text{SUM}\{\eta_{p(y_i|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}\} - \eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}$$

and computing the messages passed from these factors reduces to

$$\begin{aligned} \left(\eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i} \right)_{-0} &\leftarrow K_{\text{probit}} \left((\eta_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})})_{-0}; c_{0,ij}, \mathbf{c}_{1,ij} \right) \\ &\quad - (\eta_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})})_{-0} \end{aligned}$$

and

$$\begin{aligned} \left(\eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i} \right)_0 &\leftarrow C_{\text{probit}} \left((\eta_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})})_{-0}, (\eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i})_{-0} \right. \\ &\quad \left. + (\eta_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})})_{-0}; c_{0,ij}, \mathbf{c}_{1,ij} \right). \end{aligned}$$

For the factors $p(\mathbf{u}_i; \Sigma)$:

computing the messages passed from these factors reduces to

$$\eta_{\mathbf{u}_i \rightarrow p(\mathbf{u}_i; \Sigma)} \leftarrow \sum_{j=1}^{n_i} \eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}$$

and computing the messages passed to these factors reduces to

$$\eta_{p(\mathbf{u}_i; \Sigma) \rightarrow \mathbf{u}_i} \rightarrow \eta_{\Sigma}.$$

Algorithm 1 is essentially these natural parameter updates being cycled until convergence. The update for $\left(\eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i} \right)_0$ can be moved outside of the cycle loop without affecting convergence. Also, the $\eta_{\mathbf{u}_i \rightarrow p(\mathbf{u}_i; \Sigma)}$ updates are redundant and are omitted from Algorithm 1.

S.3 Derivation of Starting Values Recommendation

We now derive useful starting values for the $\eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}$ that have to be initialized in Algorithm 1. Note that

$$\log p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) = \zeta(a_{ij}) - \log(2) \quad \text{where} \quad a_{ij} \equiv (2y_{ij} - 1)(\boldsymbol{\beta}^T \mathbf{x}_{ij}^F + \mathbf{u}_i^T \mathbf{x}_{ij}^R)$$

and ζ is as defined in Section 3.1. Let $\hat{\mathbf{u}}_i$ be a prediction of \mathbf{u}_i and consider the following expansion of the data-dependent component of $\ell(\boldsymbol{\beta}, \Sigma)$:

$$\begin{aligned}\zeta(a_{ij}) &= \zeta(\hat{a}_{ij} + (\mathbf{u}_i - \hat{\mathbf{u}}_i)^T \mathbf{x}_{ij}^R (2y_{ij} - 1)) \\ &= \zeta(\hat{a}_{ij}) + (\mathbf{u}_i - \hat{\mathbf{u}}_i)^T \mathbf{x}_{ij}^R (2y_{ij} - 1) \zeta'(\hat{a}_{ij}) + \frac{1}{2} \{(\mathbf{u}_i - \hat{\mathbf{u}}_i)^T \mathbf{x}_{ij}^R (2y_{ij} - 1)\}^2 \zeta''(\hat{a}_{ij}) + \dots \\ &= \begin{bmatrix} 1 \\ \mathbf{u}_i - \hat{\mathbf{u}}_i \\ \text{vech}((\mathbf{u}_i - \hat{\mathbf{u}}_i)(\mathbf{u}_i - \hat{\mathbf{u}}_i)^T) \end{bmatrix}^T \check{\boldsymbol{\eta}}_{ij} + \dots\end{aligned}$$

where, as in Section 3.3, $\hat{a}_{ij} \equiv (2y_{ij} - 1)(\boldsymbol{\beta}^T \mathbf{x}_{ij}^F + \hat{\mathbf{u}}_i^T \mathbf{x}_{ij}^R)$, and

$$\check{\boldsymbol{\eta}}_{ij} \equiv \begin{bmatrix} \zeta(\hat{a}_{ij}) \\ \mathbf{x}_{ij}^R (2y_{ij} - 1) \zeta'(\hat{a}_{ij}) \\ \frac{1}{2} \zeta''(\hat{a}_{ij}) \mathbf{D}_{d^R}^T \text{vec}(\mathbf{x}_{ij}^R (\mathbf{x}_{ij}^R)^T) \end{bmatrix}.$$

It follows that the quadratic approximation to $\log p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})$ based on Taylor expansion about $\hat{\mathbf{u}}_i$ is $\log \check{p}(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})$ where

$$\check{p}(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \equiv \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i - \hat{\mathbf{u}}_i \\ \text{vech}((\mathbf{u}_i - \hat{\mathbf{u}}_i)(\mathbf{u}_i - \hat{\mathbf{u}}_i)^T) \end{bmatrix}^T \check{\boldsymbol{\eta}}_{ij} \right\}.$$

The starting value recommendation for $\eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}$ is based on replacement of the conditional probability mass function $p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})$ by $\check{p}(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})$ in (S.13):

$$m_{\check{p}(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) \leftarrow \frac{\text{proj}[m_{\mathbf{u}_i \rightarrow \check{p}(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})}(\mathbf{u}_i) \check{p}(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})](\mathbf{u}_i)}{m_{\mathbf{u}_i \rightarrow \check{p}(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})}(\mathbf{u}_i)} = \check{p}(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})$$

with the $\text{proj}[\cdot]$ being superfluous in this case due to $\check{p}(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})$ being already in the Multivariate Normal family. The starting value for $\eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}$ that arises from this substitution is then given by

$$\exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^T) \end{bmatrix}^T \eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}^{\text{start}} \right\} = \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i - \hat{\mathbf{u}}_i \\ \text{vech}((\mathbf{u}_i - \hat{\mathbf{u}}_i)(\mathbf{u}_i - \hat{\mathbf{u}}_i)^T) \end{bmatrix}^T \check{\boldsymbol{\eta}}_{ij} \right\}.$$

By matching coefficients of like terms we arrive at

$$\eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}^{\text{start}} = \begin{bmatrix} \eta_0^{\text{start}} \\ (2y_{ij} - 1) \zeta'(\hat{a}_{ij}) \mathbf{x}_{ij}^R - \zeta''(\hat{a}_{ij}) \mathbf{x}_{ij}^R (\mathbf{x}_{ij}^R)^T \hat{\mathbf{u}}_i \\ \frac{1}{2} \zeta''(\hat{a}_{ij}) \mathbf{D}_{d^R}^T \text{vec}(\mathbf{x}_{ij}^R (\mathbf{x}_{ij}^R)^T) \end{bmatrix}$$

where

$$\eta_0^{\text{start}} = \zeta(\hat{a}_{ij}) - (2y_{ij} - 1) \zeta'(\hat{a}_{ij}) (\mathbf{x}_{ij}^R)^T \hat{\mathbf{u}}_i + \frac{1}{2} \zeta''(\hat{a}_{ij}) \{(\mathbf{x}_{ij}^R)^T \hat{\mathbf{u}}_i\}^2.$$

In Algorithm 1 the cycle loop corresponds to determination of the natural parameter vector

$$\left(\eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i} \right)_{-0}$$

implying that the first entry of $\eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}^{\text{start}}$ is not needed for these iterations. Hence, we can instead set $\eta_0^{\text{start}} = 0$ without affecting Algorithm 1. We now have (14).

S.4 Details of Confidence Interval Calculations

Here we provide full details of approximate confidence intervals calculations based on quasi-Newton maximization of $\ell(\beta, \Sigma)$. The calculations depend on the following ingredients:

- some additional convenient matrix notation.
- formulae for transformation from the parameter vector $\theta \equiv \text{vech}(\frac{1}{2} \log(\Sigma))$ to a parameter vector ω that is more appropriate for confidence interval construction.
- formulae for the reverse transformation: from ω to θ .
- a quasi-Newton optimization-based strategy for calculating confidence intervals for the entries of ω , which are then easily transformed to confidence intervals for interpretable covariance matrix parameters, as illustrated in Figures 2 and 4.

S.4.1 Additional Matrix Notation

For a $d \times d$ matrix A define $\text{diagonal}(A)$ to be the $d \times 1$ vector consisting of the diagonal entries of A and, provided $d \geq 2$, define $\text{vecbd}(A)$ to be the $\frac{1}{2}d(d-1)$ vector containing the entries of A that are below the diagonal of A in order from left to right and top to bottom. For example,

$$\text{diagonal} \begin{pmatrix} 1 & 5 & 9 & 13 \\ 2 & 6 & 10 & 14 \\ 3 & 7 & 11 & 15 \\ 4 & 8 & 12 & 16 \end{pmatrix} = \begin{bmatrix} 1 \\ 6 \\ 11 \\ 16 \end{bmatrix} \quad \text{and} \quad \text{vecbd} \begin{pmatrix} 1 & 5 & 9 & 13 \\ 2 & 6 & 10 & 14 \\ 3 & 7 & 11 & 15 \\ 4 & 8 & 12 & 16 \end{pmatrix} = \begin{bmatrix} 2 \\ 3 \\ 4 \\ 7 \\ 8 \\ 12 \end{bmatrix}.$$

In addition, if each of a and b are $d \times 1$ vectors then $a \odot b$ is $d \times 1$ vector of element-wise products and a/b is $d \times 1$ vector of element-wise quotients. Similarly, $\log(a)$ and $\tanh(a)$ are obtained in an element-wise fashion.

S.4.2 Transformation from θ to ω

Given a $\frac{1}{2}d(d+1) \times 1$ vector θ , its corresponding ω vector of the same length is found via the steps:

1. Obtain the spectral decomposition $\text{vech}^{-1}(\theta) = U_\theta \text{diag}(\lambda_\theta) U_\theta^T$.
2. Set $\Sigma = U_\theta \text{diag}\{\exp(2\lambda_\theta)\} U_\theta^T$.
3. (a) If $d = 1$ then $\omega = \log(\sqrt{\Sigma})$.
(b) If $d > 1$ then

$$\omega = \begin{bmatrix} \log(\sqrt{\text{diagonal}(\Sigma)}) \\ \tanh^{-1} \left(\text{vecbd}(\Sigma) / \sqrt{\text{vecbd}(\text{diagonal}(\Sigma) \text{diagonal}(\Sigma)^T)} \right) \end{bmatrix}.$$

S.4.3 Transformation from ω to θ

Given a $\frac{1}{2}d(d+1) \times 1$ vector ω , its corresponding θ vector of the same length is found via the steps:

1. Form the $d \times d$ symmetric matrix Σ as follows:
 - (a) If $d = 1$ then $\Sigma = \exp(2\omega)$.
 - (b) If $d > 1$ then let ω_1 denote the first d entries of ω and ω_2 denote the remaining $\frac{1}{2}d(d-1)$ entries of ω .
 - i. Set $\text{diagonal}(\Sigma) = \exp(2\omega_1)$.
 - ii. Obtain the below-diagonal entries of Σ so that

$$\text{vecbd}(\Sigma) = \tanh(\omega_2) \odot \text{vecbd}(\exp(\omega_1) \exp(\omega_1)^T)$$

holds. Obtain the above-diagonal entries of Σ such that symmetry of Σ is enforced.

2. Obtain the spectral decomposition: $\Sigma = U_\Sigma \text{diag}(\lambda_\Sigma) U_\Sigma^T$.

3. Obtain $\theta = \text{vech}\left(\frac{1}{2}U_\Sigma \text{diag}\{\log(\lambda_\Sigma)\} U_\Sigma^T\right)$.

S.4.4 Quasi-Newton Optimization-Based Confidence Interval Calculations

The steps for obtaining confidence intervals for each of the interpretable parameters are:

1. Obtain $(\hat{\beta}, \hat{\theta})$ using a quasi-Newton optimization routine applied the expectation propagation-approximate log-likelihood ℓ with unconstrained input parameters (β, θ) .
2. Obtain $\hat{\omega}$ corresponding to $\hat{\theta}$ using the steps given in Section S.4.2.
3. Call the quasi-Newton optimization routine with input parameters (β, ω) instead of (β, θ) , and initial value $(\hat{\beta}, \hat{\omega})$. In this call, request that the Hessian matrix $H\ell(\beta, \omega)$ at the maximum $(\hat{\beta}, \hat{\omega})$ be computed. The steps given in Section S.4.3 are used to obtain the corresponding (β, θ) vector for evaluation of ℓ via the version of ℓ used in 1. for the optimization.
4. Form $100(1 - \alpha)\%$ confidence intervals for the entries of (β, ω) using

$$\begin{bmatrix} \hat{\beta} \\ \hat{\omega} \end{bmatrix} \pm \Phi^{-1}(1 - \frac{1}{2}\alpha) \sqrt{-\text{diagonal}(\{H\ell(\hat{\beta}, \hat{\omega})\}^{-1})}.$$

5. Transform the confidence intervals limits for the ω component, using the functions \exp and \tanh , to instead correspond to the standard deviation and correlation parameters:

$$\begin{bmatrix} \sqrt{\text{diagonal}(\Sigma)} \\ \text{vecbd}(\Sigma) / \sqrt{\text{vecbd}(\text{diagonal}(\Sigma) \text{diagonal}(\Sigma)^T)} \end{bmatrix}.$$

S.5 Details of Approximate Best Prediction

For the binary mixed model (1), the best prediction of \mathbf{u}_i is

$$\begin{aligned}\text{BP}(\mathbf{u}_i) &= E(\mathbf{u}_i|\mathbf{y}) = E(\mathbf{u}_i|\mathbf{y}_i) = \int_{\mathbb{R}^{d^R}} \mathbf{u}_i p(\mathbf{u}_i|\mathbf{y}_i; \boldsymbol{\beta}, \boldsymbol{\Sigma}) d\mathbf{u}_i \\ &= \int_{\mathbb{R}^{d^R}} \mathbf{u}_i \left\{ \frac{p(\mathbf{y}_i|\mathbf{u}_i; \boldsymbol{\beta})p(\mathbf{u}_i; \boldsymbol{\Sigma})}{\int_{\mathbb{R}^{d^R}} p(\mathbf{y}_i|\mathbf{u}_i; \boldsymbol{\beta})p(\mathbf{u}_i; \boldsymbol{\Sigma})} \right\} d\mathbf{u}_i\end{aligned}$$

where $\mathbf{y}_i \equiv (y_{i1}, \dots, y_{in_i})$. Now note that Algorithm 1 involves replacement of

$$p(\mathbf{y}_i|\mathbf{u}_i; \boldsymbol{\beta})p(\mathbf{u}_i; \boldsymbol{\Sigma}) \quad \text{by} \quad \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^T) \end{bmatrix}^T \widehat{\boldsymbol{\eta}}_i \right\}$$

where $\widehat{\boldsymbol{\eta}}_i$ is defined by (16). This leads to the approximation

$$\begin{aligned}\text{BP}(\mathbf{u}_i) &= E(\widehat{\mathbf{u}}_i) \text{ where } \widehat{\mathbf{u}}_i \text{ is Multivariate Normal with natural parameter } \widehat{\boldsymbol{\eta}}_i \\ &= -\frac{1}{2} \left\{ \text{vec}^{-1} \left(\mathbf{D}_d^{+T} \widehat{\boldsymbol{\eta}}_{i2} \right) \right\}^{-1} \widehat{\boldsymbol{\eta}}_{i1}.\end{aligned}$$

Using (13.7) of McCulloch, Searle & Neuhaus (2008), the covariance matrix of $\text{BP}(\mathbf{u}_i) - \mathbf{u}_i$ is

$$\text{Cov}\{\text{BP}(\mathbf{u}_i) - \mathbf{u}_i\} = E_{\mathbf{y}_i}\{\text{Cov}(\mathbf{u}_i|\mathbf{y}_i)\}.$$

The expectation propagation approximation of $\text{Cov}(\mathbf{u}_i|\mathbf{y}_i)$ is

$$\text{Cov}(\mathbf{u}_i|\mathbf{y}) = -\frac{1}{2} \left\{ \text{vec}^{-1} \left(\mathbf{D}_d^{+T} \widehat{\boldsymbol{\eta}}_{i2} \right) \right\}^{-1}.$$

However, approximation of $\text{Cov}\{\text{BP}(\mathbf{u}_i) - \mathbf{u}_i\}$ is hindered by the expectation over the \mathbf{y}_i vector.

Additional References

Mardia, K.V., Kent, J.T. and Bibby, J.M. (1979). *Multivariate Analysis*. London: Academic Press.

Wainwright, M.J. and Jordan, M.I. (2008). Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, **1**, 1–305.