

Web-Supplement for:
**Accurate Logistic Variational Message Passing:
Algebraic and Numerical Details**

BY TUI H. NOLAN AND MATT P. WAND

*School of Mathematical and Physical Sciences, University of Technology Sydney,
Broadway 2007, Australia*

S.1 Proof of Theorem 1

First note that for any $a, b \in \mathbb{R}$,

$$\begin{aligned} \int_{-\infty}^{\infty} \Phi(a + bx)\phi(x) dx &= \Phi\left(\frac{a}{\sqrt{b^2 + 1}}\right) \\ \text{and } \int_{-\infty}^{\infty} x \Phi(a + bx)\phi(x) dx &= \frac{b}{\sqrt{b^2 + 1}} \phi\left(\frac{a}{\sqrt{b^2 + 1}}\right). \end{aligned} \tag{S.1}$$

From the first result in (S.1), it follows immediately that

$$\int_{-\infty}^{\infty} \text{expit}_k(\mu + \sigma x)\phi(x) dx = \sum_{i=1}^k p_{k,i} \Phi\left(\frac{\mu s_{k,i}}{\sqrt{1 + \sigma^2 s_{k,i}^2}}\right).$$

Hence, for all $\mu \in \mathbb{R}$ and $\sigma > 0$,

$$\begin{aligned} &\left| \mathcal{B}_0(\mu, \sigma^2) - \sum_{i=1}^k p_{k,i} \Phi\left(\frac{\mu s_{k,i}}{\sqrt{1 + \sigma^2 s_{k,i}^2}}\right) \right| \\ &= \left| \int_{-\infty}^{\infty} \text{expit}(\mu + \sigma x)\phi(x) dx - \int_{-\infty}^{\infty} \text{expit}_k(\mu + \sigma x)\phi(x) dx \right| \\ &\leq \int_{-\infty}^{\infty} |\text{expit}(\mu + \sigma x) - \text{expit}_k(\mu + \sigma x)| \phi(x) dx \\ &\leq \sup_{u \in \mathbb{R}} |\text{expit}(u) - \text{expit}_k(u)| \int_{-\infty}^{\infty} \phi(x) dx = \Delta_k \end{aligned}$$

where the last equality follows from (7). Part (a) of Theorem 1 follows immediately.

The second result in (S.1) implies that

$$\int_{-\infty}^{\infty} x \text{expit}_k(\mu + \sigma x)\phi(x) dx = \sigma \sum_{i=1}^k \frac{p_{k,i} s_{k,i}}{\sqrt{1 + \sigma^2 s_{k,i}^2}} \phi\left(\frac{\mu s_{k,i}}{\sqrt{1 + \sigma^2 s_{k,i}^2}}\right)$$

For all $\mu \in \mathbb{R}$ and $\sigma > 0$ we then have

$$\begin{aligned}
& \left| \mathcal{B}_1(\mu, \sigma^2) - \sigma \sum_{i=1}^k \frac{p_{k,i} s_{k,i}}{\sqrt{1 + \sigma^2 s_{k,i}^2}} \phi \left(\frac{\mu s_{k,i}}{\sqrt{1 + \sigma^2 s_{k,i}^2}} \right) \right| \\
&= \left| \int_{-\infty}^{\infty} \text{expit}(\mu + \sigma x) x \phi(x) dx - \int_{-\infty}^{\infty} \text{expit}_k(\mu + \sigma x) x \phi(x) dx \right| \\
&\leq \int_{-\infty}^{\infty} |\text{expit}(\mu + \sigma x) - \text{expit}_k(\mu + \sigma x)| |x| \phi(x) dx \\
&\leq \sup_{u \in \mathbb{R}} |\text{expit}(u) - \text{expit}_k(u)| \int_{-\infty}^{\infty} |x| \phi(x) dx \\
&= 2\Delta_k \int_0^{\infty} x \phi(x) dx = \sqrt{2/\pi} \Delta_k.
\end{aligned}$$

S.2 Derivation of Algorithm 1

The message passed from $p(\mathbf{y}|\boldsymbol{\theta})$ to $\boldsymbol{\theta}$ is

$$m_{p(\mathbf{y}|\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}}(\boldsymbol{\theta}) = \exp[\mathbf{y}^T \mathbf{A} \boldsymbol{\theta} - \mathbf{1}^T \log\{\mathbf{1} + \exp(\mathbf{A} \boldsymbol{\theta})\}]$$

is not conjugate with Multivariate Normal messages passed to $\boldsymbol{\theta}$ from other factors. A non-conjugate VMP remedy (Knowles & Minka, 2011) involves replacement of $m_{p(\mathbf{y}|\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}}(\boldsymbol{\theta})$ by

$$\tilde{m}_{p(\mathbf{y}|\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}}(\boldsymbol{\theta}) \equiv \exp \left\{ \left[\begin{array}{c} \boldsymbol{\theta} \\ \text{vec}(\boldsymbol{\theta} \boldsymbol{\theta}^T) \end{array} \right]^T \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}} \right\} \quad (\text{S.2})$$

to enforce conjugacy with Multivariate Normal messages. Under pre-specification (S.2) the current $q(\boldsymbol{\theta})$ density function satisfies

$$q(\boldsymbol{\theta}) \propto \tilde{m}_{p(\mathbf{y}|\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}}(\boldsymbol{\theta}) \times (\text{product of messages passed to } \boldsymbol{\theta} \text{ from its other neighbours})$$

From (10) of Wand (2017), we then get

$$q(\boldsymbol{\theta}) \propto \tilde{m}_{p(\mathbf{y}|\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}}(\boldsymbol{\theta}) m_{\boldsymbol{\theta} \rightarrow p(\mathbf{y}|\boldsymbol{\theta})}(\boldsymbol{\theta}) = \exp \left\{ \left[\begin{array}{c} \boldsymbol{\theta} \\ \text{vec}(\boldsymbol{\theta} \boldsymbol{\theta}^T) \end{array} \right]^T \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}} \right\}$$

where

$$\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}} \equiv \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}} + \boldsymbol{\eta}_{\boldsymbol{\theta} \rightarrow p(\mathbf{y}|\boldsymbol{\theta})}.$$

Let $\boldsymbol{\mu}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}}$ and $\boldsymbol{\Sigma}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}}$ be the corresponding common parameters. The natural parameters and common parameters are the following functions of each other:

$$\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}} = \begin{bmatrix} (\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}})_1 \\ (\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}})_2 \end{bmatrix} = \begin{bmatrix} \boldsymbol{\Sigma}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}}^{-1} \boldsymbol{\mu}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}} \\ -\frac{1}{2} \text{vec}(\boldsymbol{\Sigma}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}}^{-1}) \end{bmatrix} \quad (\text{S.3})$$

$$\text{and } \begin{cases} \boldsymbol{\mu}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}} = -\frac{1}{2} \{\text{vec}^{-1}((\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}})_2)\}^{-1} (\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}})_1 \\ \boldsymbol{\Sigma}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}} = -\frac{1}{2} \{\text{vec}^{-1}((\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}})_2)\}^{-1}. \end{cases}$$

In the upcoming arguments we will use the shorthand:

$$\boldsymbol{\eta}_{q(\boldsymbol{\theta})} \equiv \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}}, \quad \boldsymbol{\mu}_{q(\boldsymbol{\theta})} \equiv \boldsymbol{\mu}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}} \quad \text{and} \quad \boldsymbol{\Sigma}_{q(\boldsymbol{\theta})} \equiv \boldsymbol{\Sigma}_{p(\mathbf{y}|\boldsymbol{\theta}) \leftrightarrow \boldsymbol{\theta}}.$$

Using the set-up of Section 2.2 of Rohde & Wand (2016), the θ -localized approximate marginal log-likelihood is

$$\log \underline{p}(\mathbf{y}; q, \boldsymbol{\eta}_{q(\theta)})^{[\theta]} = \text{Entropy}\{q(\boldsymbol{\theta}; \boldsymbol{\eta}_{q(\theta)})\} + \text{NonEntropy}\{q(\boldsymbol{\theta}; \boldsymbol{\eta}_{q(\theta)})\}$$

where

$$\text{Entropy}\{q(\boldsymbol{\theta}; \boldsymbol{\eta}_{q(\theta)})\} = \frac{1}{2} \log \left| -\frac{1}{2} \{\text{vec}^{-1}(\boldsymbol{\eta}_{q(\theta)})_2\}^{-1} \right| + d\{1 + \log(2\pi)\}/2,$$

with d denoting the dimension of $\boldsymbol{\theta}$. Also,

$$\begin{aligned} \text{NonEntropy}\{q(\boldsymbol{\theta}; \boldsymbol{\eta}_{q(\theta)})\} &\equiv E_{q(\boldsymbol{\theta}; \boldsymbol{\eta}_{q(\theta)})} \{\log p(\mathbf{y}|\boldsymbol{\theta})\} \\ &\quad + E_{q(\boldsymbol{\theta}; \boldsymbol{\eta}_{q(\theta)})} (\text{sum of other log-factors neighboring } \boldsymbol{\theta}) \\ &= \mathbf{y}^T \mathbf{A} \boldsymbol{\mu}_{q(\theta)} - E_{q(\boldsymbol{\theta}; \boldsymbol{\eta}_{q(\theta)})} \left[\mathbf{1}^T \log\{\mathbf{1} + \exp(\mathbf{A}\boldsymbol{\theta})\} + \left[\begin{array}{c} \boldsymbol{\theta} \\ \text{vec}(\boldsymbol{\theta}\boldsymbol{\theta}^T) \end{array} \right]^T \boldsymbol{\eta}^\dagger \right] \end{aligned}$$

where $\boldsymbol{\eta}^\dagger$ is the sum of the natural parameters of the messages passed to $\boldsymbol{\theta}$ other than the message from $p(\mathbf{y}|\boldsymbol{\theta})$. Ideally we would maximize $\log \underline{p}(\mathbf{y}; q, \boldsymbol{\eta}_{q(\theta)})^{[\theta]}$ over $\boldsymbol{\eta}_{q(\theta)}$ but are thwarted by the intractability of the q -density expectation of $\log\{\mathbf{1} + \exp(\mathbf{A}\boldsymbol{\theta})\}$. To get around this we apply (1) to obtain

$$\log \underline{p}(\mathbf{y}; q, \boldsymbol{\eta}_{q(\theta)})^{[\theta]} \geq \text{Entropy}\{q(\boldsymbol{\theta}; \boldsymbol{\eta}_{q(\theta)})\} + \underline{\text{NonEntropy}}\{\boldsymbol{\eta}_{q(\theta)}, \boldsymbol{\omega}_1\}$$

where

$$\begin{aligned} \underline{\text{NonEntropy}}\{\boldsymbol{\eta}_{q(\theta)}, \boldsymbol{\omega}_1\} &\equiv \mathbf{y}^T \mathbf{A} \boldsymbol{\mu}_{q(\theta)} - \frac{1}{2} (\boldsymbol{\omega}_1^2)^T \text{diagonal}(\mathbf{A} \boldsymbol{\Sigma}_{q(\theta)} \mathbf{A}^T) \\ &\quad - \mathbf{1}^T \log\{\mathbf{1} + \exp\{\mathbf{A} \boldsymbol{\mu}_{q(\theta)} + \frac{1}{2} (\mathbf{1} - 2\boldsymbol{\omega}_1) \odot \text{diagonal}(\mathbf{A} \boldsymbol{\Sigma}_{q(\theta)} \mathbf{A}^T)\}\} \\ &\quad + \left\{ \left[\begin{array}{c} \boldsymbol{\mu}_{q(\theta)} \\ \text{vec}(\boldsymbol{\Sigma}_{q(\theta)} + \boldsymbol{\mu}_{q(\theta)} \boldsymbol{\mu}_{q(\theta)}^T) \end{array} \right] \right\}^T \boldsymbol{\eta}^\dagger \end{aligned}$$

and $\boldsymbol{\omega}_1$ is an $n \times 1$ vector of variational parameters. We now seek to maximize

$$\log \underline{p}(\mathbf{y}; q, \boldsymbol{\eta}_{q(\theta)}, \boldsymbol{\omega}_1)^{[\theta]} \equiv \text{Entropy}\{q(\boldsymbol{\theta}; \boldsymbol{\eta}_{q(\theta)})\} + \underline{\text{NonEntropy}}\{\boldsymbol{\eta}_{q(\theta)}, \boldsymbol{\omega}_1\}.$$

Using arguments analogous to those given in Section 4 of Rohde & Wand (2016), the function $\log \underline{p}(\mathbf{y}; q, \boldsymbol{\eta}_{q(\theta)}, \boldsymbol{\omega}_1)^{[\theta]}$ has a stationary point in the $[\boldsymbol{\eta}_{q(\theta)}^T, \boldsymbol{\omega}_1^T]^T$ space if and only if

$$\boldsymbol{\eta}_{q(\theta)} = \left\{ \mathbf{H}_{\boldsymbol{\eta}_{q(\theta)}} A(\boldsymbol{\eta}_{q(\theta)}) \right\}^{-1} \mathbf{D}_{\boldsymbol{\eta}_{q(\theta)}} \underline{\text{NonEntropy}}\{\boldsymbol{\eta}_{q(\theta)}, \boldsymbol{\omega}_1\}^T \text{ and} \quad (\text{S.4})$$

$$\mathbf{0} = \mathbf{D}_{\boldsymbol{\omega}_1} \underline{\text{NonEntropy}}\{\boldsymbol{\eta}_{q(\theta)}, \boldsymbol{\omega}_1\}^T \quad (\text{S.5})$$

with A denoting the Multivariate Normal log-partition function and $\mathbf{D}_{\boldsymbol{\eta}_{q(\theta)}}$ and $\mathbf{H}_{\boldsymbol{\eta}_{q(\theta)}}$ respectively denoting the derivative vector and Hessian matrix with respect to $\boldsymbol{\eta}_{q(\theta)}$ as defined in Rohde & Wand (2016). Standard vector calculus steps (e.g. Wand, 2002) lead to

$$\begin{aligned} \mathbf{D}_{\boldsymbol{\omega}_1} \underline{\text{NonEntropy}}\{\boldsymbol{\eta}_{q(\theta)}, \boldsymbol{\omega}_1\}^T &= \left[\text{expit}\{\mathbf{A} \boldsymbol{\mu}_{q(\theta)} + \frac{1}{2} (\mathbf{1} - 2\boldsymbol{\omega}_1) \odot \text{diagonal}(\mathbf{A} \boldsymbol{\Sigma}_{q(\theta)} \mathbf{A}^T)\} - \boldsymbol{\omega}_1 \right] \\ &\quad \odot \text{diagonal}(\mathbf{A} \boldsymbol{\Sigma}_{q(\theta)} \mathbf{A}^T). \end{aligned}$$

Substitution of this result into (S.5) and a reworking of the arguments that lead to Result 2 of Rohde & Wand (2016), but applied to $\underline{\text{NonEntropy}}\{\boldsymbol{\eta}_{q(\theta)}, \boldsymbol{\omega}_1\}$, rather than $\text{NonEntropy}\{q(\boldsymbol{\theta}; \boldsymbol{\eta}_{q(\theta)})\}$, lead to the fixed-point updates:

$$\left\{ \begin{array}{l} \boldsymbol{\omega}_1 \leftarrow \text{expit}\{\mathbf{A} \boldsymbol{\mu}_{q(\theta)} + \frac{1}{2} (\mathbf{1} - 2\boldsymbol{\omega}_1) \odot \text{diagonal}(\mathbf{A} \boldsymbol{\Sigma}_{q(\theta)} \mathbf{A}^T)\} \\ \mathbf{v}_{q(\theta)} \leftarrow \mathbf{D}_{\boldsymbol{\mu}_{q(\theta)}} \underline{\text{NonEntropy}}\{\boldsymbol{\eta}_{q(\theta)}, \boldsymbol{\omega}_1\}^T \\ \boldsymbol{\Sigma}_{q(\theta)} \leftarrow - \left[\mathbf{H}_{\boldsymbol{\mu}_{q(\theta)}} \underline{\text{NonEntropy}}\{\boldsymbol{\eta}_{q(\theta)}, \boldsymbol{\omega}_1\} \right]^{-1} \\ \boldsymbol{\mu}_{q(\theta)} \leftarrow \boldsymbol{\mu}_{q(\theta)} + \boldsymbol{\Sigma}_{q(\theta)} \mathbf{v}_{q(\theta)} \end{array} \right. \quad (\text{S.6})$$

Further vector calculus leads to the explicit forms:

$$\begin{aligned} D_{\boldsymbol{\mu}_{q(\boldsymbol{\theta})}} \underline{\text{NonEntropy}}\{\boldsymbol{\eta}_{q(\boldsymbol{\theta})}, \boldsymbol{\omega}_1\}^T &= \mathbf{A}^T \left(\mathbf{y} - \text{expit}[\mathbf{A}\boldsymbol{\mu}_{q(\boldsymbol{\theta})} + \frac{1}{2}(\mathbf{1} - 2\boldsymbol{\omega}_1) \odot \text{diagonal}(\mathbf{A}\boldsymbol{\Sigma}_{q(\boldsymbol{\theta})}\mathbf{A}^T)] \right) \\ &\quad + (\boldsymbol{\eta}^\dagger)_1 + 2 \text{vec}^{-1}((\boldsymbol{\eta}^\dagger)_2) \boldsymbol{\mu}_{q(\boldsymbol{\theta})} \end{aligned}$$

and

$$\begin{aligned} H_{\boldsymbol{\mu}_{q(\boldsymbol{\theta})}} \underline{\text{NonEntropy}}\{\boldsymbol{\eta}_{q(\boldsymbol{\theta})}, \boldsymbol{\omega}_1\}^T &= \\ -\mathbf{A}^T \text{diag} \left(\frac{\frac{1}{2}\mathbf{1}}{\mathbf{1} + \cosh[\mathbf{A}\boldsymbol{\mu}_{q(\boldsymbol{\theta})} + \frac{1}{2}(\mathbf{1} - 2\boldsymbol{\omega}_1) \odot \text{diagonal}(\mathbf{A}\boldsymbol{\Sigma}_{q(\boldsymbol{\theta})}\mathbf{A}^T)]} \right) \mathbf{A} &+ 2 \text{vec}^{-1}((\boldsymbol{\eta}^\dagger)_2). \end{aligned}$$

Introduction of $\boldsymbol{\omega}_0 \equiv \text{logit}(\boldsymbol{\omega}_1)$ and substitution into (S.6) then gives

$$\left\{ \begin{array}{l} \boldsymbol{\omega}_0 \longleftarrow \mathbf{A}\boldsymbol{\mu}_{q(\boldsymbol{\theta})} + \frac{1}{2}(\mathbf{1} - 2\boldsymbol{\omega}_1) \odot \text{diagonal}(\mathbf{A}\boldsymbol{\Sigma}_{q(\boldsymbol{\theta})}\mathbf{A}^T) \\ \boldsymbol{\omega}_1 \longleftarrow \text{expit}(\boldsymbol{\omega}_0) \ ; \ \boldsymbol{\omega}_2 \longleftarrow \mathbf{1}/[2\{\mathbf{1} + \cosh(\boldsymbol{\omega}_0)\}] \\ \boldsymbol{\Sigma}_{q(\boldsymbol{\theta})} \longleftarrow \left\{ \mathbf{A}^T \text{diag}(\boldsymbol{\omega}_2)\mathbf{A} - 2\text{vec}^{-1}((\boldsymbol{\eta}^\dagger)_2) \right\}^{-1} \\ \boldsymbol{\mu}_{q(\boldsymbol{\theta})} \longleftarrow \boldsymbol{\mu}_{q(\boldsymbol{\theta})} + \boldsymbol{\Sigma}_{q(\boldsymbol{\theta})} \left\{ \mathbf{A}^T(\mathbf{y} - \boldsymbol{\omega}_1) + (\boldsymbol{\eta}^\dagger)_1 + 2 \text{vec}^{-1}((\boldsymbol{\eta}^\dagger)_2) \boldsymbol{\mu}_{q(\boldsymbol{\theta})} \right\}. \end{array} \right. \quad (\text{S.7})$$

The remainder of the derivation of Algorithm 1 involves expressing (S.7) in terms of the input and output natural parameter vectors

$$\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}} \quad \text{and} \quad \boldsymbol{\eta}_{\boldsymbol{\theta} \rightarrow p(\mathbf{y}|\boldsymbol{\theta})}.$$

Using (S.3), the $\boldsymbol{\omega}_0$ update is equivalent to

$$\left\{ \begin{array}{l} \boldsymbol{\mu} \longleftarrow -\frac{1}{2}\mathbf{A} \left\{ \text{vec}^{-1}((\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}})_2) \right\}^{-1} (\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}})_1 \\ \boldsymbol{\sigma}^2 \longleftarrow -\frac{1}{2} \text{diagonal} \left[\mathbf{A} \left\{ \text{vec}^{-1}((\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}})_2) \right\}^{-1} \mathbf{A}^T \right] \\ \boldsymbol{\omega}_0 \longleftarrow \boldsymbol{\mu} + \frac{1}{2}(\mathbf{1} - 2\boldsymbol{\omega}_1) \odot \boldsymbol{\sigma}^2. \end{array} \right.$$

The $\boldsymbol{\Sigma}_{q(\boldsymbol{\theta})}$ update can be written as

$$-\frac{1}{2} \left\{ \text{vec}^{-1}((\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}})_2) \right\}^{-1} \longleftarrow \left\{ \mathbf{A}^T \text{diag}(\boldsymbol{\omega}_2)\mathbf{A} - 2 \text{vec}^{-1}((\boldsymbol{\eta}^\dagger)_2) \right\}^{-1}$$

which is equivalent to

$$(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}})_2 + (\boldsymbol{\eta}^\dagger)_2 \longleftarrow -\frac{1}{2} \text{vec}(\mathbf{A}^T \text{diag}(\boldsymbol{\omega}_2)\mathbf{A}) + (\boldsymbol{\eta}^\dagger)_2$$

which, in turn, is equivalent to the second component of $\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}}$ being updated according to

$$(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}})_2 \longleftarrow -\frac{1}{2} \text{vec}(\mathbf{A}^T \text{diag}(\boldsymbol{\omega}_2)\mathbf{A}). \quad (\text{S.8})$$

For the update of the first component of $\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}}$ we note that the last update of (S.7) is equivalent to

$$\boldsymbol{\Sigma}_{q(\boldsymbol{\theta})}^{-1} \boldsymbol{\mu}_{q(\boldsymbol{\theta})} \longleftarrow \boldsymbol{\Sigma}_{q(\boldsymbol{\theta})}^{-1} \boldsymbol{\mu}_{q(\boldsymbol{\theta})} + \mathbf{A}^T(\mathbf{y} - \boldsymbol{\omega}_1) + (\boldsymbol{\eta}^\dagger)_1 + 2 \text{vec}^{-1}((\boldsymbol{\eta}^\dagger)_2) \boldsymbol{\mu}_{q(\boldsymbol{\theta})} \quad (\text{S.9})$$

where, on the right-hand side,

$$\boldsymbol{\Sigma}_{q(\boldsymbol{\theta})}^{-1} = \mathbf{A}^T \text{diag}(\boldsymbol{\omega}_2)\mathbf{A} - 2 \text{vec}^{-1}((\boldsymbol{\eta}^\dagger)_2) \quad (\text{S.10})$$

according to its updated value and

$$\mu_{q(\theta)} = -\frac{1}{2}\{\text{vec}^{-1}((\boldsymbol{\eta}_{p(\mathbf{y}|\theta)} \leftrightarrow \boldsymbol{\theta})_2)\}^{-1}(\boldsymbol{\eta}_{p(\mathbf{y}|\theta)} \leftrightarrow \boldsymbol{\theta})_1 \quad (\text{S.11})$$

is the terms of the natural parameter from the previous iteration before (S.8) has taken place. Substitution of (S.10) and (S.11) into (S.9) we get

$$\begin{aligned} & (\boldsymbol{\eta}_{p(\mathbf{y}|\theta)} \rightarrow \boldsymbol{\theta})_1 + (\boldsymbol{\eta}^\dagger)_1 \longleftarrow \\ & \left\{ -\frac{1}{2}\mathbf{A}^T \text{diag}(\boldsymbol{\omega}_2)\mathbf{A} + \text{vec}^{-1}((\boldsymbol{\eta}^\dagger)_2) \right\} \left\{ \text{vec}^{-1}((\boldsymbol{\eta}_{p(\mathbf{y}|\theta)} \leftrightarrow \boldsymbol{\theta})_2) \right\}^{-1} (\boldsymbol{\eta}_{p(\mathbf{y}|\theta)} \leftrightarrow \boldsymbol{\theta})_1 \\ & + \mathbf{A}^T(\mathbf{y} - \boldsymbol{\omega}_1) + (\boldsymbol{\eta}^\dagger)_1 - \text{vec}^{-1}((\boldsymbol{\eta}^\dagger)_2) \left\{ \text{vec}^{-1}((\boldsymbol{\eta}_{p(\mathbf{y}|\theta)} \leftrightarrow \boldsymbol{\theta})_2) \right\}^{-1} (\boldsymbol{\eta}_{p(\mathbf{y}|\theta)} \leftrightarrow \boldsymbol{\theta})_1 \end{aligned}$$

which is equivalent to

$$\begin{aligned} (\boldsymbol{\eta}_{p(\mathbf{y}|\theta)} \rightarrow \boldsymbol{\theta})_1 & \longleftarrow -\frac{1}{2}\mathbf{A}^T \text{diag}(\boldsymbol{\omega}_2)\mathbf{A} \left\{ \text{vec}^{-1}((\boldsymbol{\eta}_{p(\mathbf{y}|\theta)} \leftrightarrow \boldsymbol{\theta})_2) \right\}^{-1} (\boldsymbol{\eta}_{p(\mathbf{y}|\theta)} \leftrightarrow \boldsymbol{\theta})_1 + \mathbf{A}^T(\mathbf{y} - \boldsymbol{\omega}_1) \\ & = \mathbf{A}^T(\boldsymbol{\omega}_2 \odot \boldsymbol{\mu}) + \mathbf{A}^T(\mathbf{y} - \boldsymbol{\omega}_1) = \mathbf{A}^T(\mathbf{y} - \boldsymbol{\omega}_1 + \boldsymbol{\omega}_2 \odot \boldsymbol{\mu}). \end{aligned}$$

Combining this update with that given in (S.8) we get the following update for the full natural parameter vector:

$$\boldsymbol{\eta}_{p(\mathbf{y}|\theta)} \rightarrow \boldsymbol{\theta} \longleftarrow \begin{bmatrix} \mathbf{A}^T(\mathbf{y} - \boldsymbol{\omega}_1 + \boldsymbol{\omega}_2 \odot \boldsymbol{\mu}) \\ -\frac{1}{2}\text{vec}(\mathbf{A}^T \text{diag}(\boldsymbol{\omega}_2)\mathbf{A}) \end{bmatrix}$$

which matches Algorithm 1.

S.3 Approximation of $\text{corr}(\beta_0, \beta_1|\mathbf{y})$

Consider the Bayesian linear regression model (8). Then, given the approximate noninformativity of prior distribution of $\boldsymbol{\beta} = [\beta_0 \ \beta_1]$, the posterior covariance matrix of $\boldsymbol{\beta}$, $\text{Cov}(\boldsymbol{\beta}|\mathbf{y})$ is such that

$$\begin{aligned} \text{Cov}(\boldsymbol{\beta}|\mathbf{y}) & \approx \text{the inverse Fisher information matrix of } \boldsymbol{\beta} = [\mathbf{X}^T \text{diag}\{b''(\mathbf{X}\boldsymbol{\beta})\}\mathbf{X}]^{-1} \\ & = \left(\mathbf{X}^T \text{diag} \left[\frac{\mathbf{1}}{2\{1 + \cosh(\mathbf{X}\boldsymbol{\beta})\}} \right] \mathbf{X} \right)^{-1} \end{aligned}$$

where

$$\mathbf{X} \equiv \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}.$$

Straightforward algebra then leads to the approximate posterior correlation between β_0 and β_1 being

$$\text{corr}(\beta_0, \beta_1|\mathbf{y}) \approx \frac{-\frac{1}{n} \sum_{i=1}^n [x_i / \{1 + \cosh(\beta_0 + \beta_1 x_i)\}]}{\sqrt{\frac{1}{n} \sum_{i=1}^n [1 / \{1 + \cosh(\beta_0 + \beta_1 x_i)\}] \frac{1}{n} \sum_{i=1}^n [x_i^2 / \{1 + \cosh(\beta_0 + \beta_1 x_i)\}]}}.$$

However, the x_i s are uniformly distributed on $(0, 1)$ so replacement of sample means by population means leads to the final approximation

$$\text{corr}(\beta_0, \beta_1|\mathbf{y}) \approx \frac{-\int_0^1 [x / \{1 + \cosh(\beta_0 + \beta_1 x)\}] dx}{\sqrt{\int_0^1 [1 / \{1 + \cosh(\beta_0 + \beta_1 x)\}] dx \int_0^1 [x^2 / \{1 + \cosh(\beta_0 + \beta_1 x)\}] dx}}.$$

S.4 Approximate Marginal Log-Likelihood Expressions

The simulation study described in Section 4 concerning variational inference for the Bayesian logistic regression model (2) used approximate marginal log-likelihood expressions, appropriate for the particular approach, as a means to assess convergence. Each of the expressions are given in this section. The last one uses the definition

$$B(\mu, \sigma^2) \equiv \int_{-\infty}^{\infty} b(\mu + \sigma x) \phi(x) dx$$

where, as defined in Section 3.2,

$$b(x) \equiv \log(1 + e^x).$$

As with the \mathcal{B}_r notation given there, evaluations of $B(\boldsymbol{\mu}, \boldsymbol{\sigma}^2)$ when $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}^2$ are equal-sized column vectors are defined in an element-wise fashion, as illustrated by

$$B\left(\begin{bmatrix} 19 \\ 11 \end{bmatrix}, \begin{bmatrix} 36 \\ 28 \end{bmatrix}\right) \equiv \begin{bmatrix} B(19, 36) \\ B(11, 28) \end{bmatrix}.$$

S.4.1 Jaakkola-Jordan Updates

$$\begin{aligned} \log \underline{p}(\mathbf{y}; q, \boldsymbol{\xi}) &= \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}; \boldsymbol{\xi})}| + \frac{1}{2} \boldsymbol{\mu}_{q(\boldsymbol{\beta}; \boldsymbol{\xi})}^T \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}; \boldsymbol{\xi})}^{-1} \boldsymbol{\mu}_{q(\boldsymbol{\beta}; \boldsymbol{\xi})} - \frac{1}{2} \boldsymbol{\mu}_{\boldsymbol{\beta}}^T \boldsymbol{\Sigma}_{\boldsymbol{\beta}}^{-1} \boldsymbol{\mu}_{\boldsymbol{\beta}} \\ &\quad + \sum_{i=1}^n \{ \xi/2 - \log(1 + e^{\xi_i}) + (\xi/4) \tanh(\xi_i/2) \} - \frac{1}{2} \log |\boldsymbol{\Sigma}_{\boldsymbol{\beta}}| \end{aligned}$$

where

$$\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}; \boldsymbol{\xi})} \equiv \left[\mathbf{X}^T \text{diag} \left\{ \frac{\tanh(\boldsymbol{\xi}/2)}{2\xi} \right\} \mathbf{X} + \boldsymbol{\Sigma}_{\boldsymbol{\beta}}^{-1} \right]^{-1} \quad \text{and} \quad \boldsymbol{\mu}_{q(\boldsymbol{\beta}; \boldsymbol{\xi})} \equiv \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}; \boldsymbol{\xi})} \{ \mathbf{X}^T (\mathbf{y} - \frac{1}{2} \mathbf{1}) + \boldsymbol{\Sigma}_{\boldsymbol{\beta}}^{-1} \boldsymbol{\mu}_{\boldsymbol{\beta}} \}.$$

and $\boldsymbol{\xi}$ is the current value of the variational parameter vector that arises in the Jaakkola-Jordan device. See, for example, Section 5.1 of Wand (2017).

S.4.2 Saul-Jordan Updates

$$\begin{aligned} \log \underline{p}(\mathbf{y}; q, \boldsymbol{\omega}) &= \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}| - \frac{1}{2} \text{tr} \left[\{ \boldsymbol{\Sigma}_{\boldsymbol{\beta}}^{-1} \{ \boldsymbol{\Sigma}_{q(\boldsymbol{\beta})} + (\boldsymbol{\mu}_{q(\boldsymbol{\beta})} - \boldsymbol{\mu}_{\boldsymbol{\beta}})(\boldsymbol{\mu}_{q(\boldsymbol{\beta})} - \boldsymbol{\mu}_{\boldsymbol{\beta}})^T \} \} \right] \\ &\quad + \mathbf{y}^T \mathbf{X} \boldsymbol{\mu}_{q(\boldsymbol{\beta})} - \frac{1}{2} (\boldsymbol{\omega}^2)^T \text{diagonal}(\mathbf{X} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta})} \mathbf{X}^T) \\ &\quad - \mathbf{1}^T \log [\mathbf{1} + \exp \{ \mathbf{X} \boldsymbol{\mu}_{q(\boldsymbol{\beta})} + \frac{1}{2} (\mathbf{1} - 2\boldsymbol{\omega}) \odot \text{diagonal}(\mathbf{X} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta})} \mathbf{X}^T) \}] \\ &\quad + \frac{d}{2} - \frac{1}{2} \log |\boldsymbol{\Sigma}_{\boldsymbol{\beta}}| \end{aligned}$$

where $\boldsymbol{\omega}$ is the current value of the variational parameter vector that arises in the Saul-Jordan device.

S.4.3 Knowles-Minka-Wand Updates

$$\begin{aligned} \log \underline{p}(\mathbf{y}; q) &= \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}| - \frac{1}{2} \text{tr} \left[\{ \boldsymbol{\Sigma}_{\boldsymbol{\beta}}^{-1} \{ \boldsymbol{\Sigma}_{q(\boldsymbol{\beta})} + (\boldsymbol{\mu}_{q(\boldsymbol{\beta})} - \boldsymbol{\mu}_{\boldsymbol{\beta}})(\boldsymbol{\mu}_{q(\boldsymbol{\beta})} - \boldsymbol{\mu}_{\boldsymbol{\beta}})^T \} \} \right] \\ &\quad + \mathbf{y}^T \mathbf{X} \boldsymbol{\mu}_{q(\boldsymbol{\beta})} - \mathbf{1}^T B(\mathbf{X} \boldsymbol{\mu}_{q(\boldsymbol{\beta})}, \text{diagonal}(\mathbf{X} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta})} \mathbf{X}^T)) \\ &\quad + \frac{d}{2} - \frac{1}{2} \log |\boldsymbol{\Sigma}_{\boldsymbol{\beta}}| \end{aligned}$$

Additional References

Rohde, D. and Wand, M.P. (2016). Semiparametric mean field variational Bayes: General principles and numerical issues. *Journal of Machine Learning Research*, **17(172)**, 1–47.

Wand, M.P. (2002). Vector differential calculus in statistics. *The American Statistician*, **56**, 55–62.