

# Theory of Gaussian Variational Approximation for a Poisson Mixed Model

PETER HALL

*Department of Mathematics and Statistics, University of Melbourne, Melbourne 3000, Australia*

J.T. ORMEROD AND M.P. WAND

*School of Mathematics and Applied Statistics, University of Wollongong, Wollongong 2522,  
Australia*

26th August, 2009

---

## Abstract

Likelihood-based inference for the parameters of generalized linear mixed models is hindered by the presence of intractable integrals. Gaussian variational approximation provides a fast and effective means of approximate inference. We provide some theory for this type of approximation for a simple Poisson mixed model. In particular, we establish consistency at rate  $m^{-1/2} + n^{-1}$ , where  $m$  is the number of groups and  $n$  is the number of repeated measurements.

*Keywords:* Asymptotic theory; Generalized linear mixed models; Kullback-Liebler divergence; Longitudinal data analysis; Maximum likelihood estimation.

---

## 1 Introduction

Variational approximation has become a central component in inference in Machine Learning and other areas of Computer Science. Recent summaries of variational approximation methodology and theory are provided by Jordan (2004), Titterington (2004) and Bishop (2006). The `Infer.NET` software project (Minka *et al.*, 2008) is a manifestation of the numerous areas in which variational approximations are being applied. Almost all of this work involves Bayesian inference.

Statistical areas such as longitudinal data analysis have issues that are similar to those arising in Machine Learning. Recently, we have explored the transferral of variational approximation technology to statistical settings. One of these is likelihood-based, rather than Bayesian, inference for generalized linear mixed models. A particularly appealing approach in this context is *Gaussian variational approximation*, which involves minimum Kullback-Liebler divergence from a family of Gaussian densities. Details on Gaussian variational approximation for generalized linear mixed models are given in Ormerod & Wand (2009, 2010).

The present article is concerned with theoretical aspects of Gaussian variational approximations to maximum likelihood estimation. Almost all of the variational approximation theory of which we are aware treats Bayesian inferential settings (e.g. Humphreys & Titterington, 2000; Wang & Titterington, 2006). An exception is Hall, Humphreys & Titterington (2002), who treat likelihood-based inference for Markovian models with missingness. As we shall see, in the case of generalized linear mixed models, rigorous asymptotics for variational approximation maximum likelihood estimation is delicate and involved. For this reason attention is restricted to a simple generalized linear mixed model setting which we call the *simple Poisson mixed model* and formally define in Section 2. In Poisson mixed models, the Gaussian variational approximation admits explicit forms, which allow us to study its properties quite deeply. We show that the exact maximum likelihood estimators are well-defined. We then prove that their variational approximations are ‘root- $m$ ’ consistent, in the sense that, their discrepancy from the true parame-

ter values decreases at a rate proportional to the inverse square-root of the number of groups – denoted by  $m$ . However, this property requires the number of repeated measurements,  $n$ , to be at least as large as the square root of  $m$ . Without that condition the convergence rate is  $O_p(n^{-1})$  rather than  $O_p(m^{-1/2})$ . Hence, consistency of Gaussian variational approximation requires that both the number of groups  $m$  and the number of repeated measures  $n$  be allowed to increase. While this excludes some longitudinal data analysis settings, such as matched paired designs, there are others where it is reasonable for  $n$  to grow. Ormerod & Wand (2009, 2010) shows Gaussian variational approximation to be quite accurate for  $n \simeq 5$ . Our results also have something important to say in non-asymptotic cases, where  $n$  is small – Section 3.5 shows that Gaussian variational approximation can be inconsistent unless  $n$ , as well as  $m$ , is large.

The maximum likelihood problem and its Gaussian variational approximate solution are described in Section 2. Section 3 contains our theoretical results, and accompanying commentary. All proofs are given in Section 4. We conclude with some discussion in Section 5.

## 2 Simple Poisson Mixed Model

We now describe the simple Poisson mixed model and Gaussian variational approximate maximum likelihood estimation of its parameters. The simple Poisson mixed model is a special case of the generalized linear mixed model where the fixed effects are a simple linear relationship and the random effects correspond to a random intercept. The responses, conditional on the random effects, are assumed to be Poisson.

The observed data are  $(X_{ij}, Y_{ij})$ ,  $1 \leq i \leq m$ ,  $1 \leq j \leq n$ , where the  $Y_{ij}$ s are non-negative integers and the  $X_{ij}$ s are unrestricted real numbers. The simple Poisson mixed model is

$$Y_{ij}|X_{ij}, U_i \text{ independent Poisson with mean } \exp(\beta_0 + \beta_1 X_{ij} + U_i),$$

$$U_i \text{ independent } N(0, \sigma^2).$$

In biomedical applications the  $1 \leq i \leq m$  corresponds to  $m$  patients, and  $1 \leq j \leq n$  corresponds to  $n$  repeated measures on those patients and typically  $m \gg n$ . The random intercepts  $U_i$  invoke a within-patient correlation. See, for example, McCulloch, Searle & Neuhaus (2008) for details of this model and its longitudinal data analysis connections.

Let  $\beta \equiv (\beta_0, \beta_1)$  be the vector of fixed effects parameters. The conditional log-likelihood of  $(\beta, \sigma^2)$  is the logarithm of the joint probability mass function of the  $Y_{ij}$ s, given the  $X_{ij}$ s, as a function of the parameters. It admits the expression

$$\begin{aligned} \ell(\beta, \sigma^2) &= \sum_{i=1}^m \sum_{j=1}^n \{Y_{ij}(\beta_0 + \beta_1 X_{ij}) - \log(Y_{ij}!)\} - \frac{m}{2} \log(2\pi\sigma^2) \\ &+ \sum_{i=1}^m \log \int_{-\infty}^{\infty} \exp \left( \sum_{j=1}^n Y_{ij} u - e^{\beta_0 + \beta_1 X_{ij} + u} - \frac{u^2}{2\sigma^2} \right) du. \end{aligned} \quad (1)$$

The maximum likelihood estimates of  $\beta$  and  $\sigma^2$  are then

$$(\hat{\beta}, \hat{\sigma}^2) = \operatorname{argmax}_{\beta, \sigma^2} \ell(\beta, \sigma^2).$$

In practice, computation of  $(\hat{\beta}, \hat{\sigma}^2)$  and corresponding inference is hindered by the fact that the  $m$  integrals in (1) cannot be reduced. In this simple setting the integrals are univariate and quadrature can be entertained. However, in more elaborate grouped data

generalized linear mixed models, such as those described in Ormerod & Wand (2010), the integrals are multidimensional and quadrature is more challenging.

Gaussian variational approximation offers a remedy since it results in a closed form approximation to  $\ell(\boldsymbol{\beta}, \sigma^2)$ . So-called *variational parameters* can be chosen to optimize the quality of the approximation. Let  $\mathbf{u}$ ,  $\mathbf{x}$  and  $\mathbf{y}$  respectively denote the random vectors containing the  $U_i$ s, the  $X_{ij}$ s and the  $Y_{ij}$ s. Also, let  $p$  be the generic symbol for density or probability mass function. Then

$$\ell(\boldsymbol{\beta}, \sigma^2) = \log p(\mathbf{y}|\mathbf{x}; \boldsymbol{\beta}, \sigma^2).$$

Hence, for arbitrary density functions  $q$  on  $\mathbb{R}^m$ ,

$$\begin{aligned} \ell(\boldsymbol{\beta}, \sigma^2) &= \log p(\mathbf{y}|\mathbf{x}) \int_{\mathbb{R}^m} q(\mathbf{u}) d\mathbf{u} = \int_{\mathbb{R}^m} \log p(\mathbf{y}|\mathbf{x}) q(\mathbf{u}) d\mathbf{u} \\ &= \int_{\mathbb{R}^m} \log \left\{ \frac{p(\mathbf{y}, \mathbf{u}|\mathbf{x})/q(\mathbf{u})}{p(\mathbf{u}|\mathbf{y}, \mathbf{x})/q(\mathbf{u})} \right\} q(\mathbf{u}) d\mathbf{u} \\ &= \int_{\mathbb{R}^m} q(\mathbf{u}) \log \left\{ \frac{p(\mathbf{y}, \mathbf{u}|\mathbf{x})}{q(\mathbf{u})} \right\} d\mathbf{u} + \int_{\mathbb{R}^m} q(\mathbf{u}) \log \left\{ \frac{q(\mathbf{u})}{p(\mathbf{u}|\mathbf{y}, \mathbf{x})} \right\} d\mathbf{u}. \end{aligned}$$

The second term is the Kullback-Leibler distance between  $q(\mathbf{u})$  and  $p(\mathbf{u}|\mathbf{y}, \mathbf{x})$ . Since this is always non-negative (Kullback & Liebler, 1951) we get

$$\ell(\boldsymbol{\beta}, \sigma^2) \geq \int_{\mathbb{R}^m} q(\mathbf{u}) \log \left\{ \frac{p(\mathbf{y}, \mathbf{u}|\mathbf{x})}{q(\mathbf{u})} \right\} d\mathbf{u}.$$

Now take  $q$  to be the  $m$ -variate Gaussian density function with mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Lambda}$ . This leads to

$$\ell(\boldsymbol{\beta}, \sigma^2) \geq \underline{\ell}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \quad (2)$$

where

$$\begin{aligned} \underline{\ell}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\Lambda}) &\equiv \sum_{i=1}^m \sum_{j=1}^n \{Y_{ij}(\beta_0 + \beta_1 X_{ij} + \mu_i) - e^{\beta_0 + \beta_1 X_{ij} + \mu_i + \frac{1}{2}\lambda_i} - \log(Y_{ij}!)\} \quad (3) \\ &\quad - \frac{m}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^m (\mu_i^2 + \lambda_i) + \frac{1}{2} \log |\boldsymbol{\Lambda}| \end{aligned}$$

is the *Gaussian variational approximation* to  $\ell(\boldsymbol{\beta}, \sigma^2)$ . Here  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)$  and  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_m)$  are the diagonal entries of  $\boldsymbol{\Lambda}$ . Since (2) holds for all choices of  $\boldsymbol{\mu}$  and  $\boldsymbol{\Lambda}$  we obtain the tightest lower bound by maximizing over these *variational* parameters. Theorem 1 in Section 3.1 implies that the off-diagonal entries of  $\boldsymbol{\Lambda}$  do not improve the lower bound so there is no loss from working with

$$\begin{aligned} \underline{\ell}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) &\equiv \sum_{i=1}^m \sum_{j=1}^n \{Y_{ij}(\beta_0 + \beta_1 X_{ij} + \mu_i) - e^{\beta_0 + \beta_1 X_{ij} + \mu_i + \frac{1}{2}\lambda_i} - \log(Y_{ij}!)\} \quad (4) \\ &\quad - \frac{m}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^m (\mu_i^2 + \lambda_i) + \frac{1}{2} \sum_{i=1}^m \log(\lambda_i). \end{aligned}$$

The Gaussian variational approximate maximum likelihood estimators are:

$$(\hat{\underline{\boldsymbol{\beta}}}, \hat{\underline{\sigma^2}}) = (\boldsymbol{\beta}, \sigma^2) \text{ component of } \underset{\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}}{\operatorname{argmax}} \underline{\ell}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}).$$

Note that maximization over  $\boldsymbol{\mu}$  and  $\boldsymbol{\lambda}$  makes the lower bound as tight as possible, and hence improves the accuracy of the variational approximation.

### 3 Theoretical Results

In this section we provide several theoretical results concerned with the maximum likelihood problem presented in Section 2 and its approximate solution via Gaussian variational approximation. All proofs are deferred to Section 4.

#### 3.1 Sufficiency of a Diagonal Covariance Matrix

**THEOREM 1.** *If  $\Sigma$  is a symmetric, positive definite matrix then, given the components down the main diagonal of  $\Sigma$ ,  $|\Sigma|$  is uniquely maximized by taking the off-diagonal components to vanish.*

Theorem 1 provides a justification for dropping the off-diagonal terms of  $\Lambda$  between (3) and (4). This means that the optimal  $q$  density factorizes into a product of  $m$  univariate normal densities. This result is in accordance with the fact that the integral over  $\mathbf{u}$  in the exact log-likelihood (1) reduces to  $m$  univariate integrals.

#### 3.2 Similarities Between the Log-likelihood and its Lower Bound

In this section we give formulae for the log-likelihood and its approximation. Assume that the  $X_{ij}$ s and  $U_i$ s are totally independent; the  $X_{ij}$ s are identically distributed as  $X$ ; and the  $U_i$ s are all normal  $N(0, \sigma^2)$ . Also, for  $1 \leq i \leq m$ , let

$$Y_{i\bullet} \equiv \sum_{j=1}^n Y_{ij} \quad \text{and} \quad B_i = B_i(\beta_0, \beta_1) \equiv \sum_{j=1}^n \exp(\beta_0 + \beta_1 X_{ij}).$$

Then the log-likelihood and its approximation are:

$$\begin{aligned} \ell(\boldsymbol{\beta}, \sigma^2) &= \ell_0(\boldsymbol{\beta}, \sigma^2) + \ell_1(\boldsymbol{\beta}, \sigma^2), \\ \text{and } \underline{\ell}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) &= \ell_0(\boldsymbol{\beta}, \sigma^2) + \ell_2(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) \end{aligned}$$

where

$$\ell_0(\boldsymbol{\beta}, \sigma^2) \equiv \sum_{i=1}^m \sum_{j=1}^n \{Y_{ij}(\beta_0 + \beta_1 X_{ij}) - \log(Y_{ij}!)\} - \frac{m}{2} \log \sigma^2, \quad (5)$$

$$\ell_1(\boldsymbol{\beta}, \sigma^2) \equiv \sum_{i=1}^m \log \left\{ \int_{-\infty}^{\infty} \exp(Y_{i\bullet} u - B_i e^u - \frac{1}{2} \sigma^{-2} u^2) du \right\}$$

$$\begin{aligned} \text{and } \ell_2(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) &\equiv \sum_{i=1}^m \left\{ \mu_i Y_{i\bullet} - B_i \exp\left(\mu_i + \frac{1}{2} \lambda_i\right) \right\} \\ &\quad - \frac{1}{2\sigma^2} \sum_{i=1}^m (\mu_i^2 + \lambda_i) + \frac{1}{2} \sum_{i=1}^m \log \lambda_i. \end{aligned} \quad (6)$$

A first step is to find  $\boldsymbol{\lambda}$  and  $\boldsymbol{\mu}$  to maximize  $\ell_2(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda})$ . It is clear from the definition of  $\ell_2$ , as a series in functions of  $(\lambda_i, \mu_i)$ , that if we keep  $\beta_0, \beta_1$  and  $\sigma^2$  fixed then the resulting  $\mu_i$  will be a function of  $\lambda_i$ , and vice versa.

#### 3.3 Properties of the Variational Parameters

Here we discuss relationships among the parameters that produce an extremum of  $\ell_2(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\lambda}, \boldsymbol{\mu})$ .

**THEOREM 2.** *If  $\sigma^2 > 0$  then: (i)  $\ell_2(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda})$  has a unique maximum in  $(\boldsymbol{\lambda}, \boldsymbol{\mu})$ ; (ii) the maximum occurs when*

$$\mu_i = \sigma^2 Y_{i\bullet} + 1 - \sigma^2 \lambda_i^{-1}, \quad \text{for } 1 \leq i \leq m; \quad (7)$$

(iii) at the maximum, the parameter values satisfy

$$0 < \lambda_i < \sigma^2 \quad \text{and} \quad \mu_i < \sigma^2 Y_{i\bullet}; \quad (8)$$

and (iv)  $\mu_i$  is defined uniquely, in terms of  $B_i$  and  $Y_{i\bullet}$ , by

$$\sigma^2 B_i \exp \left\{ \mu_i + \frac{1}{2} \sigma^2 (\sigma^2 Y_{i\bullet} + 1 - \mu_i)^{-1} \right\} = \sigma^2 Y_{i\bullet} - \mu_i. \quad (9)$$

It is worth noting that the values of the components  $(\lambda_i, \mu_i)$  at which the maximum in  $(\boldsymbol{\lambda}, \boldsymbol{\mu})$  of  $\ell_2(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda})$  occurs, are determined index-by-index and do not require a more complex maximization. Of course, this is an immediate consequence of the diagonalization noted in Theorem 1.

### 3.4 “True Values” of Parameters

In this section we derive the almost-sure limits of the values of  $\beta$  and  $\sigma^2$  that maximize  $\ell(\boldsymbol{\beta}, \sigma^2)$  and  $\underline{\ell}(\boldsymbol{\beta}, \sigma^2)$ . First, however, we derive the limits of  $m^{-1} \ell_j$  for  $j = 0, 1, 2$ . For this purpose we impose the following conditions:

- (A1) for  $1 \leq j \leq n$ , the pairs  $(X_{ij}, Y_{ij}, U_i)$  are independent and identically distributed as  $(X_i, Y_i, U_i)$ , say, which in turn is distributed as  $(X, Y, U)$ ;
- (A2) the random variables  $X$  and  $U$  are independent;
- (A3) the sets of variables  $\mathcal{S}_i = \{(X_{ij}, Y_{ij}, U_i) : 1 \leq j \leq n\}$ , for  $1 \leq i \leq m$ , are independent and identically distributed;
- (A4) each  $Y_{ij}$ , conditional on both  $X_{ij}$  and  $U_i$ , is Poisson-distributed with mean  $\exp(\beta_0^0 + \beta_1^0 X_{ij} + U_i)$ , where  $\beta_0^0$  and  $\beta_1^0$  denote the true values of  $\beta_0$  and  $\beta_1$ ;
- (A5) each  $U_i$  is normal  $N(0, (\sigma^2)^0)$ , where  $(\sigma^2)^0$  denotes the true value of  $\sigma^2$ ; and that  $(\sigma^2)^0 > 0$ ;
- (A6) the moment generating function of  $X$ ,  $\phi(t) = E\{\exp(tX)\}$ , is finite for  $|t| < 2c$ , for some  $c > 0$ ; and that  $|\beta_1^0| < c$ .

Let  $(B, Y_\bullet) = (B_1, Y_{1\bullet})$ . Note that  $B$  is a function of  $\beta_0$  and  $\beta_1$ , although  $Y_\bullet$  is not. Define  $Q = Q(\beta_0, \beta_1) < 0$  to be the unique solution of the equation

$$\sigma^2 B \exp \left\{ Q + \frac{1}{2} \sigma^2 (\sigma^2 Y_\bullet + 1 - Q)^{-1} \right\} + Q - \sigma^2 Y_\bullet = 0. \quad (10)$$

Define too:

$$\ell_0^0(\boldsymbol{\beta}, \sigma^2) \equiv n \exp(\beta_0^0 + \frac{1}{2} (\sigma^2)^0) \{ \beta_0 \phi(\beta_1^0) + \beta_1 \phi'(\beta_1^0) \} - \frac{1}{2} \log(\sigma^2), \quad (11)$$

$$\ell_1^0(\boldsymbol{\beta}, \sigma^2) \equiv E \left[ \log \left\{ \int_{-\infty}^{\infty} \exp(Y_\bullet u - B e^u - \frac{1}{2} \sigma^{-2} u^2) du \right\} \right] \quad (12)$$

$$\begin{aligned} \text{and } \ell_2^0(\boldsymbol{\beta}, \sigma^2) &\equiv E(Q Y_\bullet) - \sigma^{-2} E(\sigma^2 Y_\bullet - Q) \\ &\quad - \frac{1}{2\sigma^2} E\{Q^2 + \sigma^2 (\sigma^2 Y_\bullet + 1 - Q)^{-1}\} \\ &\quad + \frac{1}{2} \log(\sigma^2) - \frac{1}{2} E\{\log(\sigma^2 Y_\bullet + 1 - Q)\}. \end{aligned} \quad (13)$$

Note that the terms in  $\frac{1}{2} \log(\sigma^2)$ , in both  $\ell_0^0$  and  $\ell_2^0$ , cancel from  $\ell_0^0 + \ell_2^0$ .

Since  $\phi(t) < \infty$  for  $|t| < 2c$  then  $E\{\exp(t|X|)\} < \infty$  for  $0 < t < 2c$ , and therefore  $|\phi'(t)| \leq E\{|X| \exp(t|X|)\} < \infty$  for  $|t| < 2c$ . Therefore  $\ell_0^0(\boldsymbol{\beta}, \sigma^2)$  is well-defined and finite provided that  $|\beta_1^0| < 2c$  and  $\sigma^2 > 0$ . The theorem below implies that  $\ell_2^0(\boldsymbol{\beta}, \sigma^2)$  is finite if  $|\beta_1^0| < c$  and  $\sigma^2 > 0$ . Clearly,  $\ell_1^0(\boldsymbol{\beta}, \sigma^2)$  is well-defined and finite whenever  $\sigma^2 > 0$ .

**THEOREM 3.** *Assume conditions (A1)–(A6). Then  $\ell_2^0(\beta, \sigma^2)$  is well-defined and finite if  $|\beta_1^0| < c$  and  $\sigma^2 > 0$ . Moreover, with probability 1, as  $m \rightarrow \infty$  and for fixed  $n$ , we have  $m^{-1} \ell_j(\beta, \sigma^2) \rightarrow \ell_j^0(\beta, \sigma^2)$  for  $j = 0, 1$ , and  $m^{-1} \sup_{\lambda, \mu} \ell_2(\beta, \sigma^2, \mu, \lambda) \rightarrow \ell_2^0(\beta, \sigma^2)$ , uniformly in*

$$\beta_0 \in [\beta_0^{(1)}, \beta_0^{(2)}], \quad \beta_1 \in [\beta_1^{(1)}, \beta_1^{(2)}], \quad \sigma^2 \in [(\sigma^2)^{(1)}, (\sigma^2)^{(2)}], \quad (14)$$

*provided that*

$$-\infty < \beta_0^{(1)} < \beta_0^{(2)} < \infty, \quad -c < \beta_1^{(1)} < \beta_1^{(2)} < c, \quad 0 < (\sigma^2)^{(1)} < (\sigma^2)^{(2)} < \infty. \quad (15)$$

Recall from Section 3.2 that the log-likelihood  $\ell$ , and its approximate form  $\underline{\ell}$ , satisfy  $\ell = \ell_0 + \ell_1$  and  $\underline{\ell} = \ell_0 + \ell_2$ . Therefore, provided the maximizations are taken over values in a range permitted by (14) and (15), the almost sure limits of the estimators of  $\beta_0$ ,  $\beta_1$  and  $\sigma^2$  that maximize  $\ell$  and  $\underline{\ell}$ , are the values of the quantities that maximize  $\ell_0^0(\beta, \sigma^2) + \ell_1^0(\beta, \sigma^2)$  and  $\ell_0^0(\beta, \sigma^2) + \ell_2^0(\beta, \sigma^2)$ , respectively.

By exploiting formulae (11)–(13) it is possible to choose the distribution of  $X$  such that (A1)–(A6) hold but the value of  $(\beta, \sigma^2)$  that maximizes  $\ell_0^0 + \ell_1^0$  is different from that which maximizes  $\ell_0^0 + \ell_2^0$ . The maximum likelihood estimators, based on maximizing  $\ell$ , are consistent and converge at rate  $m^{-1/2}$ , even if  $n$  is held fixed. Therefore, in the context described in the first sentence of this paragraph, and for fixed  $n$ , the Gaussian variational approximate estimators, based on maximizing  $\underline{\ell}$ , are inconsistent. However, as we shall show in Section 3.5, permitting  $m$  and  $n$  to diverge together leads to consistency, in fact at rate  $m^{-1/2} + n^{-1}$ .

### 3.5 Consistency at Rate $m^{-1/2} + n^{-1}$

We are now in a position to state our main results, i.e. the consistency and convergence rates of estimators of the model parameters based on Gaussian variational approximation. Recall from Section 2 that the Gaussian variational approximate maximum likelihood estimators are

$$(\hat{\underline{\beta}}_0, \hat{\underline{\beta}}_1, \hat{\underline{\sigma}}^2) = (\beta_0, \beta_1, \sigma^2) \text{ component of } \underset{\beta_0, \beta_1, \sigma^2, \mu, \lambda}{\operatorname{argmax}} \underline{\ell}(\beta, \sigma^2, \mu, \lambda).$$

We impose the following conditions:

- (A7) the moment generating function of  $X$ ,  $\phi(t) = E\{\exp(tX)\}$ , is well-defined on the whole real line;
- (A8) the mapping that takes  $\beta$  to  $\phi'(\beta)/\phi(\beta)$  is invertible;
- (A9) in some neighbourhood of  $\beta_1^0$  (the true value of  $\beta_1$ ),  $(d^2/d\beta^2) \log \phi(\beta)$  does not vanish;
- (A10) for a constant  $C > 0$ ,  $m = O(n^C)$  as  $m$  and  $n$  diverge.
- (A11) the true values  $\beta_0^0$ ,  $\beta_1^0$  and  $(\sigma^2)^0$  of  $\beta_0$ ,  $\beta_1$  and  $\sigma^2$ , respectively, lie in  $(-\infty, \infty)$ ,  $(-\infty, \infty)$  and  $(0, \infty)$ , respectively, and when choosing  $(\hat{\underline{\beta}}_0, \hat{\underline{\beta}}_1, \hat{\underline{\sigma}}^2)$  we search in the rectangular region  $[-C_1, C_1] \times [-C_1, C_1] \times [C_1^{-1}, C_1]$ , where  $C_1$  is a constant satisfying  $C_1 > \max(|\beta_0^0|, |\beta_1^0|, (\sigma^2)^0, 1/(\sigma^2)^0)$ .

**THEOREM 4.** *If (A1)–(A5) and (A7)–(A11) hold then, as  $m$  and  $n$  diverge,*

$$\hat{\underline{\beta}}_0 = \beta_0^0 + O_p(m^{-1/2} + n^{-1}), \quad \hat{\underline{\beta}}_1 = \beta_1^0 + O_p(m^{-1/2} + n^{-1}) \text{ and } \hat{\underline{\sigma}}^2 = (\sigma^2)^0 + O_p(m^{-1/2} + n^{-1}).$$

## 4 Proofs

Theorem 1, which is proved in Section 4.1, reduces the parametric complexity of the variational problem from  $O(m^2)$  to  $O(m)$  in respect of the number of groups. The proof of Theorem 2 is then relatively conventional. That theorem is then applied to prove Theorem 3, by eliminating  $\lambda$  and  $\mu$  from the variational likelihood. The proof of Theorem 4 is conducted in a sequence of three steps, each of which is essentially a lemma for the next. In particular, the first step (given in Section 4.4.1) establishes consistency of estimators of  $\beta_0$ ,  $\beta_1$  and  $\sigma^2$ ; the second step (Section 4.4.2) uses the conclusion of Step 1 to control remainder terms, so that the consistency property can be extended to a rate of convergence that is almost, but not quite, as good as the rate stated in Theorem 4. Finally, in Section 4.4.3, the conclusion of Step 2 is used to give still better control of remainders, so that the full theorem can be derived.

### 4.1 Proof of Theorem 1

Let  $\Sigma$  be  $p \times p$ , and let  $\Sigma_1$  be the  $(p-1) \times (p-1)$  matrix, let  $\mathbf{a}$  be the  $p$ -vector, and let  $b$  be the scalar such that

$$\Sigma = \begin{bmatrix} \Sigma_1 & \mathbf{a} \\ \mathbf{a}^T & b \end{bmatrix}.$$

Then,

$$|\Sigma| = |\Sigma_1| (b - \mathbf{a}^T \Sigma_1^{-1} \mathbf{a}).$$

We shall prove the theorem by induction over  $p$ , and so we may assume that  $|\Sigma_1|$  is uniquely maximized by taking the off-diagonal components of  $\Sigma_1$  to vanish. (The theorem clearly holds when  $p = 1$ .) Since  $\Sigma_1$  is positive definite then, for  $\Sigma_1$  and  $b$  fixed,  $b - \mathbf{a}^T \Sigma_1^{-1} \mathbf{a}$  is uniquely maximized by taking  $\mathbf{a} = \mathbf{0}$ , and then  $|\Sigma| = |\Sigma_1| b$ . The induction hypothesis now implies that  $|\Sigma|$  is uniquely maximized by taking the off-diagonal components of  $\Sigma$  to equal zero.

### 4.2 Proof of Theorem 2

Note that

$$\frac{\partial \ell_2(\beta, \sigma^2, \mu, \lambda)}{\partial \mu_i} = Y_{i\bullet} - B_i \exp(\mu_i + \frac{1}{2} \lambda_i) - \sigma^{-2} \mu_i, \quad (16)$$

$$\frac{\partial \ell_2(\beta, \sigma^2, \mu, \lambda)}{\partial \lambda_i} = -\frac{1}{2} B_i \exp(\mu_i + \frac{1}{2} \lambda_i) - \frac{1}{2} \sigma^{-2} + \frac{1}{2} \lambda_i^{-1}. \quad (17)$$

Equating both equations to zero to obtain a turning point in  $(\lambda_i, \mu_i)$ , and subtracting twice the second equation from the first, we see that (7) holds.

Using (7) to express  $\lambda_i$  in terms of  $\mu_i$ , the right-hand sides (16) and (17), when set equal to zero, are both equivalent to (9). The left-hand side there increases from zero to  $\sigma^2 B_i \exp\{\frac{1}{2} \sigma^2 (\sigma^2 Y_{i\bullet} + 1)^{-1}\}$ , and the right-hand side decreases from  $\infty$  to 0, as  $\mu_i$  increases from  $-\infty$  to 0. Moreover,  $\sigma^2 B_i \exp\{\frac{1}{2} \sigma^2 (\sigma^2 Y_{i\bullet} + 1)^{-1}\} > 0$ , provided that  $\sigma^2 > 0$ . Therefore, if  $\sigma^2 > 0$  then (9) has a unique solution in  $\mu_i$ , and so (i) holds.

The fact that equation formed by setting the left-hand side of (16) equal to zero has a solution means that  $\sigma^2 Y_{i\bullet} - \mu_i > 0$ , which is the second part of (8). It therefore follows from (7), and the fact that  $\lambda_i > 0$  since  $\Lambda$  must be positive definite, that  $\lambda_i \in (0, \sigma^2)$ , which is the first part of (8).

### 4.3 Proof of Theorem 3

First we establish the finiteness of  $\ell_2^0(\beta, \sigma^2)$ . Assume that  $\sigma^2 > 0$ . Since  $Q \leq 0$  and  $Y_{\bullet} \geq 0$  then  $\sigma^2 Y_{\bullet} + 1 - Q \geq 1$ , and so it suffices to ensure that  $E(Y_{\bullet}) + E(|Q| Y_{\bullet}) + E(Q^2) < \infty$ .

Now, (10) implies that  $\sigma^2 B \exp(-|Q|) \leq \sigma^2 Y_\bullet - Q \leq \sigma^2 B \exp(-|Q| + \frac{1}{2})$ . Therefore,

$$(\sigma^2 Y_\bullet)^2 + 2\sigma^2 Y_\bullet |Q| + Q^2 = (\sigma^2 Y_\bullet - Q)^2 \leq (\sigma^2 B)^2 e^{-2|Q|+1} \leq (\sigma^2 B)^2 e,$$

and so,

$$E\{(\sigma^2 Y_\bullet)^2 + 2\sigma^2 Y_\bullet |Q| + Q^2\} \leq \sigma^4 e E(B^2) = \sigma^4 \exp(2\beta_0 + 1) \phi(2\beta_1) < \infty$$

provided  $|\beta_1^0| < c$ . Hence, the latter condition implies that,  $E(Y_\bullet) + E(|Q| Y_\bullet) + E(Q^2) < \infty$ , and therefore that  $\ell_2^0(\beta, \sigma^2) < \infty$ .

Since

$$E(Y_{ij} | X_{ij}) = E\{\exp(\beta_0^0 + \beta_1^0 X_{ij} + U_i) | X_{ij}\} = \exp(\beta_0^0 + \beta_1^0 X_{ij} + \frac{1}{2}(\sigma^2)^0)$$

and  $E\{X \exp(\beta_1 X)\} = \phi'(\beta_1)$ , then

$$\begin{aligned} E\{Y_{ij}(\beta_0 + \beta_1 X_{ij})\} &= E\left\{\exp(\beta_0^0 + \beta_1^0 X_{ij} + \frac{1}{2}(\sigma^2)^0)(\beta_0 + \beta_1 X_{ij})\right\} \\ &= \exp(\beta_0^0 + \frac{1}{2}(\sigma^2)^0) \{\beta_0 \phi(\beta_1^0) + \beta_1 \phi'(\beta_1^0)\}. \end{aligned}$$

Therefore, if we take  $m$  to diverge to infinity and keep  $n$  fixed, then by the law of large numbers, and with probability 1,  $m^{-1} \ell_j(\beta, \sigma^2) \rightarrow \ell_j^0(\beta, \sigma^2)$  for  $j = 0$ , and analogously the result when  $j = 1$  also holds. This establishes pointwise convergence. Uniform convergence follows from equicontinuity of the functions  $\ell_0$  and  $\ell_1$  if they are interpreted as indexed by different values of their random arguments. For example, in the case of  $\ell_0$  we have, for different versions  $(\beta'_0, \beta'_1, (\sigma^2)')$  and  $(\beta''_0, \beta''_1, (\sigma^2)'')$  of  $(\beta_0, \beta_1, \sigma^2)$ :

$$m^{-1} |\ell_0(\beta', (\sigma^2)') - \ell_0(\beta'', (\sigma^2)'')| \leq (|\beta'_0 - \beta''_0| + |\beta'_1 - \beta''_1|) S + \frac{1}{2} |\log\{(\sigma^2)' / (\sigma^2)''\}|, \quad (18)$$

where  $S = m^{-1} \sum_i \sum_j Y_{ij} (1 + |X_{ij}|)$  and converges almost surely to a finite limit as  $m \rightarrow \infty$ . (Here we have used the fact that  $|\beta_1^0| < c$ , where  $c$  is as in (A6).) It follows from (18) that the almost sure limit as  $m \rightarrow \infty$ , of the supremum of  $|\ell_0(\beta', (\sigma^2)') - \ell_0(\beta'', (\sigma^2)'')|$  over  $|\beta'_0 - \beta''_0| + |\beta'_1 - \beta''_1| + |(\sigma^2)' - (\sigma^2)''| \leq \epsilon$ , where the parameter values are constrained to satisfy (14) and (15), converges to zero as  $\epsilon \downarrow 0$ . The case of  $\ell_1$  is similar.

To treat the convergence of  $\ell_2(\beta, \sigma^2, \mu, \lambda)$  we note that, in view of Theorem 2 and particularly (9), with probability 1,

$$\begin{aligned} m^{-1} \sup_{\lambda, \mu} \ell_2(\beta, \sigma^2, \mu, \lambda) &\rightarrow E(Q Y_\bullet) - E\{B \exp(Q + \frac{1}{2} R)\} \\ &\quad - \frac{1}{2} \sigma^{-2} E(Q^2 + R) + \frac{1}{2} E(\log R), \end{aligned}$$

where the random variables  $B, Q$  and  $R$  are jointly distributed such that  $R = \sigma^2 (\sigma^2 Y_\bullet + 1 - Q)^{-1}$  and  $Q$  solves (10). In this notation, and with probability 1,

$$\begin{aligned} m^{-1} \sup_{\lambda, \mu} \ell_2(\beta, \sigma^2, \mu, \lambda) &\rightarrow E(Q Y_\bullet) - E[B \exp\{Q + \frac{1}{2} \sigma^2 (\sigma^2 Y_\bullet + 1 - Q)^{-1}\}] \\ &\quad - \frac{1}{2} \sigma^{-2} E\{Q^2 + \sigma^2 (\sigma^2 Y_\bullet + 1 - Q)^{-1}\} \\ &\quad + \frac{1}{2} \log(\sigma^2) - \frac{1}{2} E\{\log(\sigma^2 Y_\bullet + 1 - Q)\}. \end{aligned}$$

That is equivalent to asserting that, with probability 1,  $m^{-1} \sup_{\lambda, \mu} \ell_2(\beta, \sigma^2, \mu, \lambda) \rightarrow \ell_2^0(\beta, \sigma^2)$ . Again, uniform convergence follows via an equicontinuity argument.

## 4.4 Proof of Theorem 4

### 4.4.1 Consistency

For  $1 \leq i \leq m$ , let  $\hat{\lambda}_i$  and  $\hat{\mu}_i$  denote the entries of  $\boldsymbol{\mu}$  and  $\boldsymbol{\lambda}$  that maximize  $\underline{\ell}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda})$ . Equating the right-hand sides of (16) and (17) to zero, and dividing by  $n$ , we see that  $\hat{\beta}_0$ ,  $\hat{\beta}_1$ ,  $\hat{\sigma}^2$ ,  $\hat{\lambda}_i$  and  $\hat{\mu}_i$  satisfy:

$$\hat{\mu}_i / (n \hat{\sigma}^2) + \exp(\hat{\mu}_i + \frac{1}{2} \hat{\lambda}_i + \hat{\beta}_0) \frac{1}{n} \sum_{j=1}^n \exp(\hat{\beta}_1 X_{ij}) - \frac{1}{n} Y_{i\bullet} = 0 \quad (19)$$

$$\text{and } \frac{1}{n \hat{\sigma}^2} + \exp(\hat{\mu}_i + \frac{1}{2} \hat{\lambda}_i + \hat{\beta}_0) \frac{1}{n} \sum_{j=1}^n \exp(\hat{\beta}_1 X_{ij}) - \frac{1}{n \hat{\lambda}_i} = 0. \quad (20)$$

Since  $X$  has finite moment generating function,  $\phi$ , on the real line then Markov's inequality can be used to prove that for all  $C_1, C_2 > 0$  and  $\rho \in (0, \frac{1}{2})$ ,

$$\sup_{\beta_1: |\beta_1| \leq C_1} P \left\{ \left| \frac{1}{n} \sum_{j=1}^n \exp(\beta_1 X_{ij}) - \phi(\beta_1) \right| > n^{-\rho} \right\} = O(n^{-C_2}).$$

Therefore, if  $\mathcal{G}(n)$  is a grid of points in the interval  $[-C_1, C_1]$  containing no more than  $O(n^C)$  points, for some  $C > 0$ , then

$$P \left\{ \sup_{\beta_1 \in \mathcal{G}(n)} \left| \frac{1}{n} \sum_{j=1}^n \exp(\beta_1 X_{ij}) - \phi(\beta_1) \right| > n^{-\rho} \right\} = O(n^{-C_2})$$

for all  $C_2 > 0$ . Choosing  $C$  sufficiently large, and using Hölder continuity of the exponential function and of  $\phi$ , we can extend this bound from  $\mathcal{G}(n)$  to the whole of the interval:

$$P \left\{ \sup_{\beta_1: |\beta_1| \leq C_1} \left| \frac{1}{n} \sum_{j=1}^n \exp(\beta_1 X_{ij}) - \phi(\beta_1) \right| > n^{-\rho} \right\} = O(n^{-C_2}).$$

Therefore, provided  $m = O(n^C)$  for some  $C > 0$ ,

$$P \left\{ \max_{1 \leq i \leq m} \sup_{\beta_1: |\beta_1| \leq C_1} \left| \frac{1}{n} \sum_{j=1}^n \exp(\beta_1 X_{ij}) - \phi(\beta_1) \right| > n^{-\rho} \right\} = O(n^{-C_2}) \quad (21)$$

for all  $C_1, C_2 > 0$ . The probability statement in each of the three displays immediately above (21) could have been preceded by  $\max_{1 \leq i \leq m}$ , although we chose not to do so because the distribution of  $X_{ij}$  does not depend on  $i$ . Nevertheless, the passage to (21) can be interpreted as moving the operator  $\max_{1 \leq i \leq m}$  from outside to inside the probability statement.

Recall that, conditional on  $X_{ij}$  and  $U_i$ ,  $Y_{ij}$  is Poisson-distributed with mean  $\exp(\beta_0^0 + \beta_1^0 X_{ij} + U_i)$ . It therefore can be proved using Markov's inequality that for all  $C_1, C_2 > 0$  and  $\rho \in (0, \frac{1}{2})$ ,

$$\max_{1 \leq i \leq m} P \left[ \left| \frac{1}{n} \sum_{j=1}^n \{Y_{ij} - \exp(\beta_0^0 + \beta_1^0 X_{ij} + U_i)\} \right| > n^{-\rho} \right] = O(n^{-C_2}).$$

Hence, since  $m = O(n^C)$  for some  $C > 0$ ,

$$P \left[ \max_{1 \leq i \leq m} \left| \frac{1}{n} \sum_{j=1}^n \{Y_{ij} - \exp(\beta_0^0 + \beta_1^0 X_{ij} + U_i)\} \right| > n^{-\rho} \right] = O(n^{-C_2}). \quad (22)$$

The properties  $m = O(n^C)$  and

$$P\left\{\max_{1 \leq i \leq m} U_i \geq \sigma (2D \log m)^{1/2}\right\} \leq m^{1-D} \quad (23)$$

imply that there exists  $C' > 0$  such that

$$P\left[\max_{1 \leq i \leq m} \exp(U_i) > \exp\{C' (\log n)^{1/2}\}\right] \rightarrow 0. \quad (24)$$

Combining (21), (22) and (24) we see that for all  $C_1, C_2 > 0$  and  $\rho \in (0, \frac{1}{2})$ ,

$$P\left\{\max_{1 \leq i \leq m} \left|\frac{1}{n} Y_{i\bullet} - \exp(\beta_0^0 + U_i) \phi(\beta_1^0)\right| > n^{-\rho}\right\} = O(n^{-C_2}). \quad (25)$$

From (21) and (25) we deduce that, for each  $\rho \in (0, \frac{1}{2})$ ,  $\hat{\lambda}_i$  and  $\hat{\mu}_i$  satisfy

$$\frac{\hat{\mu}_i}{\sigma^2 n} + \exp\left(\hat{\mu}_i + \frac{1}{2} \hat{\lambda}_i + \beta_0\right) \{\phi(\beta_1) + O_p(n^{-\rho})\} = \exp(\beta_0^0 + U_i) \{\phi(\beta_1^0) + O_p(n^{-\rho})\}, \quad (26)$$

$$\frac{1}{\sigma^2 n} + \exp\left(\hat{\mu}_i + \frac{1}{2} \hat{\lambda}_i + \beta_0\right) \{\phi(\beta_1) + O_p(n^{-\rho})\} = \frac{1}{n \hat{\lambda}_i}, \quad (27)$$

where the  $O_p(n^{-\rho})$  remainders are of the stated orders uniformly in  $1 \leq i \leq m$  and in  $\sigma^2, |\beta_0|, |\beta_1| \leq C_1$ , and  $\sigma^2 \geq C_1^{-1}$  with  $C_1 > 0$  arbitrary but fixed.

By (8),  $0 < \lambda_i < \sigma^2$ . Therefore, the left-hand side of (26) equals  $\{\mu_i/(\sigma^2 n)\} + \exp(\mu_i + \omega_i) \{\phi(\beta_1) + O_p(n^{-\rho})\}$ , where  $\beta_0 - \frac{1}{2} \sigma^2 \leq \omega_i \leq \beta_0 + \frac{1}{2} \sigma^2$ ; call this result (R<sub>1</sub>). We confine  $\sigma^2$  and  $\beta_0$  to compact sets, in particular to values satisfying  $\sigma^2, |\beta_0| \leq C_1$  and  $\sigma^2 \geq C_1^{-1}$ , and so  $|\omega_i|$  is bounded uniformly in  $i$ . This property, (R<sub>1</sub>) and (26) imply that

$$\exp\left(\hat{\mu}_i + \frac{1}{2} \hat{\lambda}_i + \beta_0\right) \{\phi(\beta_1) + o_p(1)\} = \exp(\beta_0^0 + U_i) \{\phi(\beta_1^0) + o_p(1)\}, \quad (28)$$

uniformly in  $1 \leq i \leq m$  and in  $\sigma^2, \beta_0$  and  $\beta_1$  satisfying  $\sigma^2, |\beta_0|, |\beta_1| \leq C_1$  and  $\sigma^2 \geq C_1^{-1}$ , for any fixed  $C_1 > 1$ . To establish (28) in detail, note first that for each  $C_1 > 0$ ,

$$\phi(\beta_1) \text{ is bounded away from zero and infinity uniformly in } |\beta_1| \leq C_1. \quad (29)$$

If (28) fails for some  $\eta > 0$  then it can be shown from (26) that “for some  $i$  in the range  $1 \leq i \leq m$ , both  $\mu_i < 0$  and  $|\mu_i|/(\sigma^2 n) > \eta \exp(\mu_i + \omega_i)$ , for all  $\sigma^2, \beta_0$  and  $\beta_1$  satisfying  $\sigma^2, |\beta_0|, |\beta_1| \leq C_1$  and  $\sigma^2 \geq C_1^{-1}$ ,” where the property within quotation marks holds with probability bounded away from zero along an infinite subsequence of values of  $n$ . In this case, since  $|\omega_i|$  is bounded then  $\mu_i < -\log n + O(\log \log n)$ , and so, in view of (29), the left-hand side of (26), which can be nonnegative, is, for this  $i = i(n)$ , of order  $n^{-c}$  for some  $c > 0$ . Hence for this  $i$ ,  $\exp(\beta_0^0 + U_i) = O_p(n^{-c})$ , and so the probability that the latter bound holds for some  $1 \leq i \leq m$  must itself be bounded away from zero along an infinite subsequence of values of  $n$  diverging to infinity; call this result (R<sub>2</sub>). Since the distribution of  $U_i$  is symmetric then (24) holds if, on the left-hand side, we replace  $U_i$  by  $-U_i$ , and so

$$P\left[\min_{1 \leq i \leq m} \exp(U_i) \geq \exp\{-C' (\log n)^{1/2}\}\right] \rightarrow 1.$$

Since  $\exp\{-C' (\log n)^{1/2}\}$  is of strictly larger order than  $n^{-c}$  for any  $c > 0$  then property (R<sub>2</sub>) is violated, and so the original assumption that (28) fails must have been false.

Results (28) and (29) imply that  $\hat{\lambda}_i$  and  $\hat{\mu}_i$  satisfy

$$\hat{\mu}_i + \frac{1}{2} \hat{\lambda}_i + \beta_0 + \log \phi(\beta_1) = \beta_0^0 + U_i + \log \phi(\beta_1^0) + o_p(1), \quad (30)$$

$$\exp\left(\hat{\mu}_i + \frac{1}{2} \hat{\lambda}_i + \beta_0\right) \phi(\beta_1) = \exp(\beta_0^0 + U_i) \phi(\beta_1^0) \{1 + o_p(1)\}, \quad (31)$$

uniformly in  $1 \leq i \leq m$  and in  $\sigma^2, \beta_0$  and  $\beta_1$  satisfying  $\sigma^2, |\beta_0|, |\beta_1| \leq C_1$  and  $\sigma^2 \geq C_1^{-1}$ , for each fixed  $C_1 > 0$ . Substituting (31) into (27), and noting (29), we deduce that the  $\widehat{\lambda}_i$  satisfy

$$\frac{1}{\sigma^2 n} + \exp(\beta_0^0 + U_i) \phi(\beta_1^0) \{1 + o_p(1)\} = \frac{1}{n \widehat{\lambda}_i},$$

uniformly in  $1 \leq i \leq m$  and in  $\sigma^2, \beta_0$  and  $\beta_1$  satisfying  $\sigma^2, |\beta_0|, |\beta_1| \leq C_1$  and  $\sigma^2 \geq C_1^{-1}$ . This result, (23) and the version of (23) with  $U_i$  replaced by  $-U_i$  imply that

$$-\log(n \widehat{\lambda}_i) = \beta_0^0 + U_i + o_p(1), \quad (32)$$

uniformly in the same sense. In particular,  $\sup_{1 \leq i \leq m} \widehat{\lambda}_i \rightarrow 0$  in probability. Hence, by (30),

$$\widehat{\mu}_i = U_i + \beta_0^0 - \beta_0 + \log\{\phi(\beta_1^0)/\phi(\beta_1)\} + o_p(1), \quad (33)$$

uniformly in  $1 \leq i \leq m$  and in  $\sigma^2, \beta_0$  and  $\beta_1$  satisfying  $\sigma^2, |\beta_0|, |\beta_1| \leq C_1$  and  $\sigma^2 \geq C_1^{-1}$ .

The following property follows using the argument leading to (21):

$$P\left\{\max_{1 \leq i \leq m} \sup_{\beta_1: |\beta_1| \leq C_1} \left| \frac{1}{n} \sum_{j=1}^n X_{ij} \exp(\beta_1 X_{ij}) - \phi'(\beta_1) \right| > n^{-\rho}\right\} = O(n^{-C_2}); \quad (34)$$

the property below is a consequence of (25):

$$P\left\{\left| \frac{1}{mn} \sum_{i=1}^m Y_{i\bullet} - \exp(\beta_0^0 + \frac{1}{2}(\sigma^2)^0) \phi(\beta_1^0) \right| > n^{-\rho}\right\} = O(n^{-C_2}); \quad (35)$$

and the following property can be derived analogously:

$$P\left\{\left| \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n X_{ij} Y_{ij} - \exp(\beta_0^0 + \frac{1}{2}(\sigma^2)^0) \phi'(\beta_1^0) \right| > n^{-\rho}\right\} = O(n^{-C_2}). \quad (36)$$

Each of (34)–(36) holds for all  $C_1, C_2 > 0$  and all  $\rho \in (0, \frac{1}{2})$ . Formula (6) for  $\ell_2(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda})$  implies that

$$-\frac{1}{mn} \frac{\partial \ell_2(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \beta_0} = \frac{1}{mn} \sum_{i=1}^m \exp(\beta_0 + \mu_i + \frac{1}{2} \lambda_i) \sum_{j=1}^n \exp(\beta_1 X_{ij}), \quad (37)$$

$$\begin{aligned} -\frac{1}{mn} \frac{\partial \ell_2(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \beta_1} &= \frac{1}{mn} \sum_{i=1}^m \exp(\beta_0 + \mu_i + \frac{1}{2} \lambda_i) \sum_{j=1}^n X_{ij} \exp(\beta_1 X_{ij}) \\ &= \exp(\beta_0 + \frac{1}{2}(\sigma^2)^0) \phi'(\beta_1) + o_p(1), \end{aligned} \quad (38)$$

where the second identity in (38) follows from (33) and (34). The second identity in (38) holds uniformly in values of  $\lambda_i$  and  $\mu_i$  that solve (19) and (20), and for  $\sigma^2, \beta_0$  and  $\beta_1$  satisfying  $\sigma^2, |\beta_0|, |\beta_1| \leq C_1$  and  $\sigma^2 \geq C_1^{-1}$ . Also, by (5), (35) and (36),

$$\frac{1}{mn} \frac{\partial \ell_0(\boldsymbol{\beta}, \sigma^2)}{\partial \beta_0} = \frac{1}{mn} \sum_{i=1}^m Y_{i\bullet} = \exp(\beta_0^0 + \frac{1}{2}(\sigma^2)^0) \phi(\beta_1^0) + o_p(1), \quad (39)$$

$$\frac{1}{mn} \frac{\partial \ell_0(\boldsymbol{\beta}, \sigma^2)}{\partial \beta_1} = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n X_{ij} Y_{ij} = \exp(\beta_0^0 + \frac{1}{2}(\sigma^2)^0) \phi'(\beta_1^0) + o_p(1), \quad (40)$$

uniformly in the same sense. Combining (38) and (40) we deduce that:

$$\frac{1}{mn} \frac{\partial \underline{\ell}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \beta_1} = \exp((\sigma^2)^0/2) \{ \exp(\beta_0^0) \phi'(\beta_1^0) - \exp(\beta_0) \phi'(\beta_1) \} + o_p(1), \quad (41)$$

uniformly in the following sense:

$$\begin{aligned} & \text{uniformly in } \sigma^2, \beta_0, \beta_1, \boldsymbol{\lambda} \text{ and } \boldsymbol{\mu} \text{ that solve } \partial \underline{\ell} / \partial \beta_0 = \partial \underline{\ell} / \partial \beta_1 = 0, \partial \underline{\ell} / \partial \boldsymbol{\lambda} = \\ & \partial \underline{\ell} / \partial \boldsymbol{\mu} = \mathbf{0}, \text{ and which satisfy } \sigma^2, |\beta_0|, |\beta_1| \leq C_1 \text{ and } \sigma^2 \geq C_1^{-1}, \text{ provided that} \\ & C_1 \text{ is so large that } |\beta_0^0|, |\beta_1^0| < C_1 \text{ and } C_1^{-1} < (\sigma^2)^0 < C_1. \end{aligned} \quad (42)$$

Recall from (16), (37) and (39) that

$$\frac{\partial \underline{\ell}_2(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \mu_i} = Y_{i\bullet} - B_i \exp(\mu_i + \frac{1}{2} \lambda_i) - \sigma^{-2} \mu_i, \quad (43)$$

$$\begin{aligned} \frac{\partial \underline{\ell}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \beta_0} &= \sum_{i=1}^m \sum_{j=1}^n Y_{ij} - \sum_{i=1}^m \exp(\beta_0 + \mu_i + \frac{1}{2} \lambda_i) \sum_{j=1}^n \exp(\beta_1 X_{ij}) \\ &= \sum_{i=1}^m \{Y_{i\bullet} - B_i \exp(\mu_i + \frac{1}{2} \lambda_i)\}. \end{aligned} \quad (44)$$

Adding (43) over  $i$ , and subtracting from (44), we deduce that

$$\sum_{i=1}^m \widehat{\underline{\mu}}_i = 0. \quad (45)$$

Hence, by (33),

$$\beta_0^0 - \widehat{\underline{\beta}}_0 + \log\{\phi(\beta_1^0)/\phi(\widehat{\underline{\beta}}_1)\} = o_p(1), \quad (46)$$

uniformly in the sense of (42).

Formula (41) [respectively, (46)] equates the value of  $\exp(\beta_0) \phi'(\beta_1)$  [respectively,  $\exp(\beta_0) \phi(\beta_1)$ ] to its true value, plus a negligible remainder. Since, by assumption, there is a one-to-one relation between values of  $\beta$  and values of  $\phi'(\beta)/\phi(\beta)$ , then the value of  $\phi'(\beta_1)/\phi(\beta_1)$  (which equals  $\exp(\beta_0) \phi'(\beta_1)/\{\exp(\beta_0) \phi(\beta_1)\}$ ), together with  $\exp(\beta_0)$  (which equals  $\exp(\beta_0) \phi(\beta_1)/\phi(\beta_1)$ ), uniquely determine the values of  $\beta_0$  and  $\beta_1$ . Therefore consistency of estimation of  $\beta_0$  and  $\beta_1$  follows from (41) and (46).

Note too that

$$\begin{aligned} \frac{2}{m} \frac{\partial \underline{\ell}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \sigma^2} &= \sigma^{-4} \frac{1}{m} \sum_{i=1}^m (\lambda_i + \mu_i^2) - \sigma^{-2} \\ &= \sigma^{-4} \frac{1}{m} \sum_{i=1}^m [U_i + \beta_0^0 - \beta_0 + \log\{\phi(\beta_1^0)/\phi(\beta_1)\}]^2 - \sigma^{-2} + o_p(1), \end{aligned} \quad (47)$$

where the first identity follows from (5) and (6), and the second comes from (33). Therefore if  $\beta_0$  and  $\beta_1$  solve the variational approximate likelihood equations then

$$\widehat{\underline{\sigma}}^2 = \frac{1}{m} \sum_{i=1}^m U_i^2 + o_p(1) = (\sigma^2)^0 + o_p(1).$$

and  $(\sigma^2)^0$  is estimated consistently.

#### 4.4.2 Convergence rate for variational approximate estimators equals $O_p(m^{-1/2} + n^{-\rho})$

By (32),

$$\widehat{\underline{\lambda}}_i = n^{-1} \exp\{- (\beta_0^0 + U_i) + o_p(1)\}, \quad (48)$$

where the remainder is of the stated order uniformly in  $1 \leq i \leq m$ . Also, by (23), the version of (23) for  $-U_i$  rather than  $U_i$ , and (33),

$$\max_{1 \leq i \leq m} |\widehat{\underline{\mu}}_i| = O_p(n^\eta), \quad (49)$$

for each  $\eta > 0$ . Hence, by (26), (48) and the consistency of the estimator of  $(\sigma^2)^0$ , proved in Section 4.4.1,

$$\exp(\widehat{\underline{\mu}}_i + \widehat{\underline{\beta}}_0) \{ \phi(\widehat{\underline{\beta}}_1) + O_p(n^{-\rho}) \} = \exp(\beta_0^0 + U_i) \{ \phi(\beta_1^0) + O_p(n^{-\rho}) \}, \quad (50)$$

uniformly in  $1 \leq i \leq m$  and for each  $\rho \in (0, \frac{1}{2})$ .

Equation (50), together with (29), imply the following stronger form of (33):

$$\widehat{\underline{\mu}}_i = U_i + \beta_0^0 - \widehat{\underline{\beta}}_0 + \log\{ \phi(\beta_1^0) / \phi(\widehat{\underline{\beta}}_1) \} + O_p(n^{-\rho}), \quad (51)$$

uniformly in  $1 \leq i \leq m$  and for each  $\rho \in (0, \frac{1}{2})$ .

The argument from (34) down, but with (34) replaced by (51), can now be used to prove that for each  $\rho \in (0, \frac{1}{2})$ ,

$$\widehat{\underline{\sigma}}^2 = (\sigma^2)^0 + O_p(m^{-1/2} + n^{-\rho}), \quad \widehat{\underline{\beta}}_0 = \beta_0^0 + O_p(m^{-1/2} + n^{-\rho}), \quad \widehat{\underline{\beta}}_1 = \beta_1^0 + O_p(m^{-1/2} + n^{-\rho}). \quad (52)$$

The term in  $m^{-1/2}$  derives from the standard deviations of means of functions of the  $U_i$ s.

#### 4.4.3 Concise convergence rate for variational approximate estimators

Without loss of generality,  $\phi'(t) \neq 0$  in some compact neighbourhood of  $\beta_1^0$ . (If not, add a constant to the random variable  $X$  to ensure that  $\phi'(\beta_1^0) \neq 0$ .) We shall take  $\beta_1^0$  to lie in that neighbourhood, and assume too that  $\widehat{\underline{\beta}}_1$  is there. In view of the already-proved consistency of  $\widehat{\underline{\beta}}_1$  for  $\beta_1^0$ , this assumption too can be made without loss of generality.

Using (21), (22), (29), (48), (49) and the consistency of  $\widehat{\underline{\sigma}}^2$  for  $(\sigma^2)^0$ , it can be proved that

$$\frac{1}{n} \sum_{j=1}^n \exp(\beta_1 X_{ij}) = \phi(\beta_1) \exp\{ \Delta_{i1}(\beta_1) \}, \quad (53)$$

$$\frac{1}{n} \sum_{j=1}^n X_{ij} \exp(\beta_1 X_{ij}) = \phi'(\beta_1) \exp\{ \Delta_{i2}(\beta_1) \}, \quad (54)$$

$$\frac{1}{n} Y_{i\bullet} - \frac{\mu_i}{\sigma^2 n} = \exp\{ U_i + \beta_0^0 + \Delta_{i3}(\beta_0, \beta_1) \} \phi(\beta_1^0), \quad (55)$$

uniformly in  $i$ , and where  $\Delta_{i1}$ ,  $\Delta_{i2}$  and  $\Delta_{i3}$  satisfy, for all  $C_1 > 0$  and each  $\rho \in (0, \frac{1}{2})$

$$\max_{1 \leq i \leq m} \sup_{|\beta_0|, |\beta_1| \leq C_1} |\Delta_{ik}(\beta_0, \beta_1)| = O_p(n^{-\rho}). \quad (56)$$

(When  $k = 1$  or  $2$  the dependence of  $\Delta_{ik}(\beta_0, \beta_1)$  on  $\beta_0$  is degenerate. Note that, by (8), the left-hand side of (55) is strictly positive. ) It can also be proved that

$$\max_{k=1,2,3} \max_{r=1,2} \sup_{|\beta_0|, |\beta_1| \leq C_1} \frac{1}{m} \left| \sum_{i=1}^m \exp(U_i) \Delta_{ik}(\beta_0, \beta_1)^r \right| = O_p(m^{-1/2} + n^{-1}) \quad (57)$$

$$\max_{k=1,2,3} \max_{r_1=0,1} \max_{r_2=1,2} \sup_{|\beta_0|, |\beta_1| \leq C_1} \frac{1}{m} \left| \sum_{i=1}^m U_i^{r_1} \Delta_{ik}(\beta_0, \beta_1)^{r_2} \right| = O_p(m^{-1/2} + n^{-1}). \quad (58)$$

In the notation of (53)–(55), (19) is equivalent to:

$$\exp\{ \beta_0 + \mu_i + \frac{1}{2} \lambda_i + \Delta_{i1}(\beta_1) \} \phi(\beta_1) = \exp\{ \beta_0^0 + U_i + \Delta_{i3}(\beta_0, \beta_1) \} \phi(\beta_1^0). \quad (59)$$

Taking logarithms of both sides of (59), and adding over  $i$ , we deduce that if  $\beta_0$  and  $\beta_1$  satisfy that equation,

$$\beta_0 - \beta_0^0 + \log\{ \phi(\beta_1) / \phi(\beta_1^0) \} = \frac{1}{m} \sum_{i=1}^m \left\{ U_i - \mu_i - \frac{1}{2} \lambda_i + \Delta_{i3}(\beta_0, \beta_1) - \Delta_{i1}(\beta_1) \right\}. \quad (60)$$

Additionally,  $m^{-1} \sum_i U_i = O_p(m^{-1/2})$ ; by (45),  $\sum_i \hat{\underline{\mu}}_i = 0$ ; using (48),  $m^{-1} \sum_i \hat{\underline{\lambda}}_i = O_p(n^{-1})$ ; and by (58),

$$\sup_{|\beta_0|, |\beta_1| \leq C_1} \frac{1}{m} \left| \sum_{i=1}^m \{\Delta_{i3}(\beta_0, \beta_1) - \Delta_{i1}(\beta_1)\} \right| = O_p(m^{-1/2} + n^{-1}).$$

Combining the results from (60) down we deduce that

$$\hat{\underline{\beta}}_0 - \beta_0^0 + \log\{\phi(\hat{\underline{\beta}}_1)/\phi(\beta_1^0)\} = O_p(m^{-1/2} + n^{-1}). \quad (61)$$

Observe that

$$\Delta \equiv \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n X_{ij} \{Y_{ij} - \exp(\beta_0^0 + \beta_1^0 X_{ij} + U_i)\} = O_p(m^{-1/2}). \quad (62)$$

Using (38) and (40) we deduce that the equation  $\partial \ell / \partial \beta_1 = 0$  is equivalent to:

$$\frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n X_{ij} \left\{ Y_{ij} - \exp\left(\beta_0 + \mu_i + \frac{1}{2} \lambda_i + \beta_1 X_{ij}\right) \right\} = 0,$$

which in turn is equivalent to:

$$\begin{aligned} \Delta + \exp(\beta_0^0) \phi'(\beta_1^0) \frac{1}{m} \sum_{i=1}^m \exp\{U_i + \Delta_{i2}(\beta_1^0)\} \\ = \exp(\beta_0) \phi'(\beta_1) \frac{1}{m} \sum_{i=1}^m \exp\{\mu_i + \frac{1}{2} \lambda_i + \Delta_{i2}(\beta_1)\}. \end{aligned} \quad (63)$$

But by (59),

$$\hat{\underline{\mu}}_i + \frac{1}{2} \hat{\underline{\lambda}}_i = \beta_0^0 - \hat{\underline{\beta}}_0 + \log\{\phi(\beta_1^0)/\phi(\hat{\underline{\beta}}_1)\} + U_i + \Delta_{i3}(\hat{\underline{\beta}}_0, \hat{\underline{\beta}}_1) - \Delta_{i1}(\hat{\underline{\beta}}_1), \quad (64)$$

and therefore,

$$\begin{aligned} \exp\left\{\mu_i + \frac{1}{2} \lambda_i + \Delta_{i2}(\hat{\underline{\beta}}_0)\right\} &= \exp\left\{\beta_0^0 - \hat{\underline{\beta}}_0 + U_i + \Delta_{i3}(\hat{\underline{\beta}}_0, \hat{\underline{\beta}}_1) - \Delta_{i1}(\hat{\underline{\beta}}_1) \right. \\ &\quad \left. + \Delta_{i2}(\hat{\underline{\beta}}_1)\right\} \phi(\beta_1^0)/\phi(\hat{\underline{\beta}}_1). \end{aligned} \quad (65)$$

Combining (63) and (65) we deduce that

$$\begin{aligned} \Delta \exp(-\beta_0^0) \phi(\beta_1^0)^{-1} + \phi'(\beta_1^0) \phi(\beta_1^0)^{-1} \frac{1}{m} \sum_{i=1}^m \exp\{U_i + \Delta_{i2}(\beta_1^0)\} \\ = \phi'(\hat{\underline{\beta}}_1) \phi(\hat{\underline{\beta}}_1)^{-1} \frac{1}{m} \sum_{i=1}^m \exp\left\{U_i + \Delta_{i3}(\hat{\underline{\beta}}_0, \hat{\underline{\beta}}_1) - \Delta_{i1}(\hat{\underline{\beta}}_1) + \Delta_{i2}(\hat{\underline{\beta}}_1)\right\}. \end{aligned} \quad (66)$$

Together, (52), (57), (58), (62) and (66) imply that,

$$\begin{aligned} \phi'(\beta_1^0) \phi(\beta_1^0)^{-1} \frac{1}{m} \sum_{i=1}^m \exp(U_i) \\ = \phi'(\hat{\underline{\beta}}_1) \phi(\hat{\underline{\beta}}_1)^{-1} \frac{1}{m} \sum_{i=1}^m \exp(U_i) + O_p(m^{-1/2} + n^{-1}). \end{aligned} \quad (67)$$

(To derive (67) we Taylor-expanded quantities  $\exp(U_i + \Delta_i)$  as  $\exp(U_i) (1 + \Delta_i + \frac{1}{2} \Delta_i^2)$ , plus a remainder dominated by  $\frac{1}{6} |\Delta_i|^3 \exp(U_i + |\Delta_i|) = O_p(n^{\eta-(3/2)})$ , uniformly in  $1 \leq i \leq m$  for all  $\eta > 0$ ; here we used (56).) Result (67) implies that

$$\phi'(\beta_1^0) \phi(\beta_1^0)^{-1} = \phi'(\hat{\underline{\beta}}_1) \phi(\hat{\underline{\beta}}_1)^{-1} + O_p(m^{-1/2} + n^{-1}). \quad (68)$$

Together, (61) and (68) imply that

$$\widehat{\underline{\beta}}_0 = \beta_0^0 + O_p(m^{-1/2} + n^{-1}), \quad \widehat{\underline{\beta}}_1 = \beta_1^0 + O_p(m^{-1/2} + n^{-1}). \quad (69)$$

Results (48), (64) and (69) imply that

$$\widehat{\underline{\mu}}_i = U_i + \Delta_{i3}(\widehat{\underline{\beta}}_0, \widehat{\underline{\beta}}_1) - \Delta_{i1}(\widehat{\underline{\beta}}_1) + O_p(m^{-1/2} + n^{-1}), \quad (70)$$

uniformly in  $1 \leq i \leq m$ . Combining (58) and (70) we obtain:

$$\frac{1}{m} \sum_{i=1}^m \mu_i^2 = \frac{1}{m} \sum_{i=1}^m U_i^2 + O_p(m^{-1/2} + n^{-1}) = (\sigma^2)^0 + O_p(m^{-1/2} + n^{-1}). \quad (71)$$

From (47), (48) and (71) we deduce that  $\widehat{\underline{\sigma}}^2 = (\sigma^2)^0 + O_p(m^{-1/2} + n^{-1})$ . The theorem follows from this property and (69).

## 5 Discussion

The preceding two sections represent an important first step in understanding the theoretical properties of variational approximations in likelihood-based inference. The simple Poisson mixed model lends itself to a deep understanding of such properties since it is the one of the simplest generalized linear mixed model that is complicated enough to benefit from approximation methods. Of course, there are several extensions that could be entertained: non-equal sample sizes within groups, multiple fixed effects, multiple variance components, more general covariance structures and other generalized response families. The case of unequal sample sizes is straightforward to address provided we take the sizes to lie between two constant multiples of  $n$ ; more disparate sample sizes lead to more complex results. General covariance structures and response families are more challenging to address, but a problem of arguably greater interest is that of finding practical approximations to the distributions of estimators. Methods for solving that problem are currently being developed. Asymptotic distribution theory is another extension of interest. The current article is likely to provide a basis for such future work.

## Acknowledgements

This research was partially supported by Australian Research Council grants to University of Melbourne and University of Wollongong.

## References

- Bishop, C.M. (2006). *Pattern Recognition and Machine Learning*. New York: Springer.
- Hall, P., Humphreys, K. & Titterton, D.M. (2002). On the adequacy of variational lower bound functions for likelihood-based inference in Markovian models with missing values. *Journal of the Royal Statistical Society, Series B*, **64**, 549–564.
- Humphreys, K. & Titterton, D.M. (2000). Approximate Bayesian inference for simple mixtures. In *Proceedings of Computational Statistics 2000* (eds J.G. Bethlehem and P.G.M. van der Heijden), pp. 331–336. Heidelberg: Physica.
- Jordan, M.I. (2004). Graphical models. *Statistical Science*, **19**, 140–155.

- Kullback, S. & Leibler, R.A. (1951). On information and sufficiency. *The Annals of Mathematical Statistics*, **22**, 79–86.
- McCulloch, C.E., Searle, S.R. & Neuhaus, J.M. (2008). *Generalized, Linear, and Mixed Models, Second Edition*. New York: John Wiley & Sons.
- Minka, T., Winn, J., Guiver, G. & Kannan, A. (2008). Infer.Net. Microsoft Research Cambridge, Cambridge, UK.
- Ormerod, J.T. and Wand, M.P. (2009). Comment on paper by Rue, Martino and Chopin. *Journal of the Royal Statistical Society, Series B*, **71**, 377-378.
- Ormerod, J.T. and Wand, M.P. (2010). Gaussian variational approximate inference for generalized linear mixed models. Unpublished manuscript.
- Titterton, D.M. (2004). Bayesian methods for neural networks and related models. *Statistical Science* **19**, 128–139.
- Wang, B. & Titterton, D.M. (2006). Convergence properties of a general algorithm for calculating variational Bayesian estimates for a normal mixture model. *Bayesian Analysis* **1**, 625–650.