



**1351.0.55.029**

**Research Paper**

**Small Area Estimation  
Using a Multinomial Logit  
Mixed Model with Category  
Specific Random Effects**



New  
Issue

## Research Paper

# Small Area Estimation Using a Multinomial Logit Mixed Model with Category Specific Random Effects

Janice Scealy

Australian Bureau of Statistics  
and Australian National University

AUSTRALIAN BUREAU OF STATISTICS

EMBARGO: 11.30 AM (CANBERRA TIME) THURS 28 JAN 2010

ABS Catalogue no. 1351.0.55.029

© Commonwealth of Australia 2010

This work is copyright. Apart from any use as permitted under the *Copyright Act 1968*, no part may be reproduced by any process without prior written permission from the Commonwealth. Requests and inquiries concerning reproduction and rights in this publication should be addressed to The Manager, Intermediary Management, Australian Bureau of Statistics, Locked Bag 10, Belconnen ACT 2616, by telephone (02) 6252 6998, fax (02) 6252 7102, or email <intermediary.management@abs.gov.au>.

Views expressed in this paper are those of the author(s), and do not necessarily represent those of the Australian Bureau of Statistics. Where quoted, they should be attributed clearly to the author(s).

Produced by the Australian Bureau of Statistics

## INQUIRIES

The ABS welcomes comments on the research presented in this paper. For further information, please contact Ms Janice Scealy, Analytical Services Branch on Canberra (02) 6252 5764 or email <analytical.services@abs.gov.au>.

## CONTENTS

ABSTRACT .....	1
1. INTRODUCTION .....	2
2. THE MODEL .....	4
3. ESTIMATION OF $\beta$ AND $\varphi$ .....	5
4. PQL ESTIMATION OF $\beta$ AND $u$ .....	8
5. APPROXIMATE ML ESTIMATION OF $\varphi$ .....	14
6. APPROXIMATE REML ESTIMATION OF $\varphi$ .....	21
7. EMPIRICAL BEST PREDICTION, SMALL AREA ESTIMATION AND A NOTE ON MSE ESTIMATION .....	24
8. ANALYTICAL APPROXIMATION OF THE MSE .....	26
9. OUT-OF-SAMPLE SMALL AREAS .....	38
10. APPROXIMATE STANDARD ERRORS OF THE ELEMENTS OF $\hat{\beta}$ .....	44
11. PARAMETRIC BOOTSTRAP MEAN SQUARED ERRORS .....	45
12. AUXILIARY DATA .....	47
13. ESTIMATES OF MODEL PARAMETERS .....	52
14. RESIDUAL PLOTS AND GOODNESS OF FIT TESTS .....	59
15. SMALL AREA ESTIMATES AND MSE ESTIMATES .....	70
16. CONCLUSION .....	74
ACKNOWLEDGEMENTS .....	75
REFERENCES .....	76
APPENDIX .....	78



# SMALL AREA ESTIMATION USING A MULTINOMIAL LOGIT MIXED MODEL WITH CATEGORY SPECIFIC RANDOM EFFECTS

Janice Scealy  
Australian Bureau of Statistics  
and Australian National University

## ABSTRACT

This paper describes a model based approach to producing small area estimates of counts for different categories of the Australian labour force based on a multinomial logit mixed model with category specific random effects. By category specific we mean that within each small area there are two correlated random effects, one associated with the employed category and the other associated with the unemployed category. Estimates of the model parameters are produced using penalized quasi-likelihood combined with approximated restricted maximum likelihood estimation and using these, estimated counts are then produced for each small area. Mean squared error estimates of the estimated counts are approximated using two methods: 1) a parametric bootstrap and 2) analytical approximations and we compare the performance of both. Using a parametric bootstrap we also examine the properties of the combined penalized quasi-likelihood and restricted maximum likelihood estimators and discuss model goodness of fit measures and diagnostics.

Keywords: small area estimation, multinomial logit mixed model, parametric bootstrap, labour force survey.

## 1. INTRODUCTION

The Australian Bureau of Statistics (ABS) produces labour force estimates using direct survey estimators for regions with large enough sample sizes for these estimators to be reliable. In recent years there has been a growing demand for labour force estimates to be produced in smaller geographical regions. The direct estimates for these small regions are considered to be too unreliable because the standard errors are large due to small sample sizes. A way around this is to produce model based estimates which borrow strength from administrative and Census data and other types of auxiliary variables. The hope is that the model based estimators will produce estimates with mean squared error less than the direct survey estimators.

The model based approach relies on an appropriate choice of model and good auxiliary variables. The aim here is to produce estimates for each of three labour force statuses: employment, unemployment and not in the labour force for a set of small areas. Auxiliary data are available within age/sex classes for each small area and sample counts of the three labour force statuses are obtained from the Australian Labour Force Survey. The total numbers of people within each sex/age group are also assumed known and are obtained from the Estimated Resident Population (ERP) projections published by the ABS (for further details, see ABS, 2007). Molina *et al.* (2007) describe a methodology based on the application of the multinomial logit mixed model which can be used to produce estimates in this small area estimation situation. Random area effects are included in the models to account for potential correlations between the age/sex class counts in the small areas not explained by the auxiliary variables. The inclusion of random area effects in the model specifically accounts for the area level variation not explained by the auxiliary variables.

In the model described in Molina *et al.* (2007), only one random area effect is used within each small area and the random effect is therefore the same across the multinomial classes. In our situation this may not be appropriate. Some work carried out at the ABS on fitting three separate logistic mixed models to the data suggests that the variances of the random effects are not the same across each category. A more appropriate model would be to introduce category specific random effects. This allows for the variances of the random effects to differ between the categories and also allows for a potential correlation between them as well. In our case it does not make sense to assume that the category specific random effects are perfectly correlated. By allowing for a general arbitrary covariance matrix, Hartzel *et al.* (2001) make the point that the model will be structurally the same regardless of the choice of baseline category which is a good property.



In this paper we extend the model in Molina *et al.* (2007) to include category specific random effects. To estimate model parameters, we develop a similar penalized quasi-likelihood (PQL) estimation scheme with approximate maximum likelihood (ML) and/or restricted maximum likelihood (REML) for the variance components. These parameter estimates are then used to produce estimated labour force counts for each small area. Mean squared error estimates of the estimated counts are approximated using two methods:

1. a parametric bootstrap, and
2. analytical approximations,

and we compare the estimates produced using these two methods. Using a parametric bootstrap, we also examine the properties of the combined PQL and REML estimators and discuss model goodness of fit. Unlike Molina *et al.* (2007), we also consider estimation for out-of-sample small areas and briefly review some alternative estimation schemes. Note that the primary focus of this paper is to give technical details on how one might produce model based estimates for the Australian labour force. This is an experimental procedure and the ABS will not be publishing any model based estimates for the Australian labour force as part of the ABS product at this stage.

## 2. THE MODEL

Similar to Molina *et al.* (2007) let index  $i$  ( $i = 1, 2, \dots, I_d$ ) denote the sex/age groups and  $d$  ( $d = 1, 2, \dots, D$ ) denote the small areas. We have a slightly different set up because  $I_d$  is not constant ( $I_d \leq 10$ ). The reason why we have unbalanced data is because some of the sex/age groups within a small area have no sample and will be excluded from model estimation since they do not contribute to the likelihood. The labour force sample counts are denoted by  $y_{di1}$ ,  $y_{di2}$  and  $y_{di3}$  which represent respectively, employment, unemployment and not in the labour force counts in sex/age group  $i$  in the  $d$ -th small area. Let  $m_{di} = y_{di1} + y_{di2} + y_{di3}$  denote the sample size and  $p_{di1}$ ,  $p_{di2}$  and  $p_{di3}$  denote the respective probabilities of employed, unemployed and not in the labour force. Let  $u_{d1}$  and  $u_{d2}$  denote the category specific random effects. We assume that the vectors  $(y_{di1}, y_{di2}, y_{di3})^t$  given  $m_{di}$  and  $\mathbf{u}_d = (u_{d1}, u_{d2})^t$  are independent across  $d$  and  $i$  with multinomial distribution, that is with the probability density function

$$f(y_{di1}, y_{di2} | \mathbf{u}_d) = \frac{m_{di}!}{y_{di1}! y_{di2}! y_{di3}!} p_{di1}^{y_{di1}} p_{di2}^{y_{di2}} p_{di3}^{y_{di3}}. \quad (2.1)$$

It is also assumed that for  $j = 1, 2$

$$\log \frac{p_{dij}}{p_{di3}} = \mathbf{x}_{dij}^t \boldsymbol{\beta}_j + u_{dj}, \quad (2.2)$$

where  $\boldsymbol{\beta}_j$  is a vector of parameters and  $\mathbf{x}_{dij}$  is a vector of explanatory variables associated with the  $j$ -th category. Note that in (2.1) above, technically we should also be conditioning on  $\mathbf{x}_{dij}$ , that is,  $f$  should be defined as  $f(y_{di1}, y_{di2} | \mathbf{u}_d, \mathbf{x}_{dij})$ . It will be assumed throughout this paper that whenever we condition on  $\mathbf{u}_d$  we also condition on  $\mathbf{x}_{dij}$  even when it is omitted from the notation. We also assume that  $\mathbf{u}_d$  is independently and identically distributed as bivariate normal and its probability density function is

$$f(\mathbf{u}_d) = \frac{1}{2\pi |\mathbf{W}_d|^{\frac{1}{2}}} e^{-\frac{1}{2} \mathbf{u}_d^t \mathbf{W}_d^{-1} \mathbf{u}_d}$$

where

$$\mathbf{W}_d = \begin{pmatrix} \varphi_1 & \varphi_{12} \\ \varphi_{12} & \varphi_2 \end{pmatrix}.$$

Under this model we have a vector of variance components  $\boldsymbol{\varphi} = (\varphi_1, \varphi_2, \varphi_{12})^t$  that will need to be estimated along with  $\boldsymbol{\beta} = (\boldsymbol{\beta}_1^t, \boldsymbol{\beta}_2^t)^t$ . This model is a GLMM (generalised linear mixed model) and to estimate the parameters a variety of different techniques can be used. In the next section we discuss some of these.

### 3. ESTIMATION OF $\beta$ AND $\varphi$

Let  $\mathbf{y}_{di} = (y_{di1}, y_{di2})^t$  for  $i = 1, 2, \dots, I_d$  and  $d = 1, 2, \dots, D$  and let  $\mathbf{y}$  be the vector obtained by stacking the  $\mathbf{y}_{di}$ 's into a column. Using the definition of extended likelihood in Pawitan (2001) we may write the likelihood function for  $\beta$ ,  $\varphi$  and  $\mathbf{u} = (\mathbf{u}_1^t, \mathbf{u}_2^t, \dots, \mathbf{u}_D^t)^t$  as

$$\begin{aligned} L(\beta, \varphi, \mathbf{u}) &= f(\mathbf{y} | \mathbf{u})f(\mathbf{u}) \\ &= \left( \prod_{d=1}^D \prod_{i=1}^{I_d} f(y_{di1}, y_{di2} | \mathbf{u}_d) \right) \left( \prod_{d=1}^D f(\mathbf{u}_d) \right). \end{aligned} \quad (3.1)$$

Pawitan also notes that this definition of the likelihood is called hierarchical likelihood by Lee and Nelder (1996).

For the estimation of  $\beta$  and  $\varphi$ , ideally likelihood based estimation should be based on maximising  $L(\beta, \varphi)$ , where

$$\begin{aligned} L(\beta, \varphi) &= \int \left( \prod_{d=1}^D \prod_{i=1}^{I_d} f(y_{di1}, y_{di2} | \mathbf{u}_d) \right) \left( \prod_{d=1}^D f(\mathbf{u}_d) \right) d\mathbf{u} \\ &= \prod_{d=1}^D \left( \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(\mathbf{u}_d) \left( \prod_{i=1}^{I_d} f(y_{di1}, y_{di2} | \mathbf{u}_d) \right) du_{d1} du_{d2} \right). \end{aligned} \quad (3.2)$$

To estimate  $\beta$  and  $\varphi$  one could try to maximise the marginal likelihood defined at (3.2) using, for example, Monte-Carlo methods or numerical integration techniques to evaluate the integrals. A procedure like Newton-Raphson could then be used to solve the likelihood equations since the equations are non-linear. More specifically, Hartzel *et al.* (2001) describe an adaptive Gauss-Hermite, quasi-Newton algorithm which could be used in the multinomial logit mixed model case. This method is appropriate when the dimension of the integrals are small and the dataset size is not large. In our case, the dimension of the integrals are small but the dataset size is large. That is, there are a large number of double integrals to evaluate at each iteration and convergence will therefore be very slow.

Another method which could be used to maximise (3.2) is to use an automated Monte Carlo EM algorithm as described in Hartzel *et al.* (2001). Again this will be computationally intensive in our case. The method consists of implementing an EM algorithm which treats the random effects  $\mathbf{u}$  as the missing data. In the E-step a conditional expectation needs to be evaluated and this is approximated by using Monte Carlo methods (actually, an independent sample from the conditional distribution is generated). The M-Step is then undertaken by maximising this

approximate expectation. At each EM iteration, the Monte Carlo sample size is increased in an automated way until convergence. In any case, for computation reasons we do not pursue exact marginal likelihood maximisation approaches further here. Instead we will now discuss approximate methods.

The approach taken in Molina *et al.* (2007) is to use the PQL (penalized quasi-likelihood) method introduced by Breslow and Clayton (1993) combined with either ML (maximum likelihood) or REML (restricted maximum likelihood) for the variance components. The PQL method of Breslow and Clayton (1993) obtains estimates of  $\beta$  given  $\varphi$  by maximising an approximation to the marginal likelihood  $L(\beta, \varphi)$ . Part of this approximation involves a Laplace integral approximation. Estimates for  $u$  are also produced as a by-product of the approximation and hence the method produces joint estimates of  $\beta$  and  $u$  given  $\varphi$ . As Jiang (2007) states, the PQL method is equivalent to the maximum hierarchical likelihood method of Lee and Nelder (1996) in the case of normality of the random effects. The maximum hierarchical likelihood method obtains joint estimates of  $\beta$  and  $u$  by maximising the hierarchical likelihood (the log of (3.1) in our case). Interestingly in the case of normal linear mixed models the estimate of  $\beta$  obtained by maximising  $L(\beta, \varphi)$  and  $L(\beta, \varphi, u)$  given  $\varphi$  are equivalent. But this is not the case in general for GLMMs.

As mentioned by Jiang (2007), Lee and Nelder (1996) showed that in general the maximum hierarchical likelihood estimates of the fixed effects are asymptotically equivalent to the marginal maximum likelihood estimates of the fixed effects. On the surface this suggests that when the sample sizes are large, then to obtain estimates of  $\beta$ , maximising (3.2) and (3.1) are equivalent. However as Jiang (2007) states, the asymptotics here are in the sense of the cluster sample sizes approaching infinity, but the number of clusters remaining bounded. This is not satisfied in our small area estimation case since the number of clusters (small areas) is large but the cluster sample sizes are small and bounded ( $\leq 10$ , since the sampled units are the age/sex classes). Therefore estimators for  $\beta$  derived from maximising (3.1) will not be equivalent to maximising (3.2) even as we increase the number of small areas in the sample. As Jiang (2007) also states, there are a number of approximations involved in deriving the PQL and these approximations have introduced bias into the estimates and this bias does not vanish asymptotically (PQL estimators are known to be inconsistent). When the cluster sample sizes are small there is insufficient information to estimate both the random and fixed effects  $\beta$  simultaneously.

Hartzel *et al.* (2001) points out that PQL methods have been shown to be biased especially for highly non-normal cases such as Bernoulli response data and biases tend to increase as the variance components increase. Hartzel *et al.* (2001) suspect that a similar problem will exist for the multinomial logit random effects model when the multinomial sample sizes are small. In our case the variance components are

expected to be small and the average  $m_{di}$  is approximately 13. Also note that the number of multinomial observations per cluster is  $\leq 10$ . In most cases these are exactly 10. Therefore the  $m_{di}$  and cluster sizes might be sufficiently large together and the variance components sufficiently small to allow the PQL estimators to work reasonably well. This will of course need to be confirmed via simulation which is undertaken later.

In any case, two issues have been highlighted so far. The first is the computational difficulty of maximising (3.2) directly due to the presence of the integrals which have no closed form solution and the second is the inconsistency of the PQL estimates associated with maximising (3.1). These two issues give some motivation for trying to come up with alternative estimators. Jiang (1998) proposes an alternative estimation method called the method of simulated moments which is both computationally attractive and results in consistent estimators. A set of estimating equations are obtained by equating sample moments of the sufficient statistics to their expectations. The expectations are then approximated by simulating sequences of normal random variables (i.e. the integrals in the expectations are approximated by Monte Carlo simulation). The equations are then solved by a Newton-Raphson procedure. However, there is one issue associated with this method. Jiang (1998) shows that the method of simulated moment estimators can be quite inefficient. For small samples the method of simulated moment estimators seem to have substantially larger variance than estimators based on PQL. So it appears there is a bias versus variance trade-off when choosing between such methods as PQL or the method of simulated moments. There are clearly issues associated with all estimators that have been discussed so far.

#### 4. PQL ESTIMATION OF $\beta$ AND $u$

To obtain the PQL estimates of  $\beta$  and  $u$  we need to maximise the log of the joint likelihood defined at (3.1). Assume for the moment that the variance components  $\varphi$  are known. The joint log likelihood is

$$l(\beta, u) = c - \frac{1}{2} \sum_{d=1}^D u_d^t W_d^{-1} u_d + \sum_{d=1}^D \sum_{i=1}^{I_d} \sum_{j=1}^3 y_{dij} \log p_{dij}, \quad (4.1)$$

where  $c$  is a constant. The maximum likelihood estimators can be obtained by equating the first derivatives of (4.1) to zero and then solving this system of equations. Let  $q = 1, 2, \dots, Q_j$  index the components of the vectors  $\beta_j$  and  $x_{dij}$  and denote the  $q$ -th component of each by  $\beta_{j(q)}$  and  $x_{dij(q)}$  respectively. By noting that

$$p_{di1} = \frac{e^{x_{di1}^t \beta_1 + u_{d1}}}{1 + e^{x_{di1}^t \beta_1 + u_{d1}} + e^{x_{di2}^t \beta_2 + u_{d2}}}$$

$$p_{di2} = \frac{e^{x_{di2}^t \beta_2 + u_{d2}}}{1 + e^{x_{di1}^t \beta_1 + u_{d1}} + e^{x_{di2}^t \beta_2 + u_{d2}}}$$

and

$$p_{di3} = \frac{1}{1 + e^{x_{di1}^t \beta_1 + u_{d1}} + e^{x_{di2}^t \beta_2 + u_{d2}}}$$

after some algebra it can be shown that for  $j = 1, 2, 3; j' = 1, 2; q = 1, 2, \dots, Q_j; i = 1, 2, \dots, I_d$  and  $d = 1, 2, \dots, D$ ,

$$\frac{\partial \log p_{dij}}{\partial \beta_{j'(q)}} = \begin{cases} x_{dij'(q)} (1 - p_{dij'}) & \text{if } j = j' \\ -x_{dij'(q)} p_{dij'} & \text{otherwise} \end{cases}$$

and for  $j' = 1, 2$  and  $q = 1, 2, \dots, Q_{j'}$ ,

$$\frac{\partial l(\beta, u)}{\partial \beta_{j'(q)}} = \sum_{d=1}^D \sum_{i=1}^{I_d} x_{dij'(q)} (y_{dij'} - m_{di} p_{dij'}). \quad (4.2)$$

Now we need to find the first derivatives with respect to the random effects. Again after some algebra it can be shown that for  $j = 1, 2, 3; j' = 1, 2; i = 1, 2, \dots, I_d; d = 1, 2, \dots, D$  and  $d' = 1, 2, \dots, D$ ,

$$\frac{\partial \log p_{dij}}{\partial u_{d'j'}} = \begin{cases} 1 - p_{d'ij'} & \text{if } j = j' \text{ and } d = d' \\ -p_{d'ij'} & \text{if } j \neq j' \text{ and } d = d' \\ 0 & \text{otherwise,} \end{cases}$$

and for  $j' = 1, 2$  and  $d' = 1, 2, \dots, D$ ,

$$\frac{\partial \left( \sum_{d=1}^D \sum_{i=1}^{I_d} \sum_{j=1}^3 y_{dij} \log p_{dij} \right)}{\partial u_{d'j'}} = \sum_{i=1}^{I_{d'}} (y_{d'ij'} - m_{d'i} p_{d'ij'}).$$

Now we note that for  $d = 1, 2, \dots, D$ ,

$$\mathbf{W}_d^{-1} = \begin{pmatrix} \varphi_1 & \varphi_{12} \\ \varphi_{12} & \varphi_2 \end{pmatrix}^{-1} = \frac{1}{\varphi_1 \varphi_2 - \varphi_{12}^2} \begin{pmatrix} \varphi_2 & -\varphi_{12} \\ -\varphi_{12} & \varphi_1 \end{pmatrix}$$

and 
$$\mathbf{u}_d^t \mathbf{W}_d^{-1} \mathbf{u}_d = \frac{1}{\varphi_1 \varphi_2 - \varphi_{12}^2} (\varphi_2 u_{d1}^2 - 2\varphi_{12} u_{d1} u_{d2} + \varphi_1 u_{d2}^2).$$

Therefore for  $d' = 1, 2, \dots, D$  and  $j' = 1, 2$ ,

$$\frac{\partial \left( -\frac{1}{2} \sum_{d=1}^D \mathbf{u}_d^t \mathbf{W}_d^{-1} \mathbf{u}_d \right)}{\partial u_{d'j'}} = \begin{cases} \frac{-1}{\varphi_1 \varphi_2 - \varphi_{12}^2} (\varphi_2 u_{d'1} - \varphi_{12} u_{d'2}) & \text{if } j' = 1 \\ \frac{-1}{\varphi_1 \varphi_2 - \varphi_{12}^2} (\varphi_1 u_{d'2} - \varphi_{12} u_{d'1}) & \text{if } j' = 2 \end{cases}$$

and hence for  $d' = 1, 2, \dots, D$  and  $j' = 1, 2$ ,

$$\frac{\partial l(\boldsymbol{\beta}, \mathbf{u})}{\partial u_{d'j'}} = \begin{cases} \frac{-1}{\varphi_1 \varphi_2 - \varphi_{12}^2} (\varphi_2 u_{d'1} - \varphi_{12} u_{d'2}) + \sum_{i=1}^{I_{d'}} (y_{d'ij'} - m_{d'i} p_{d'ij'}) & \text{if } j' = 1 \\ \frac{-1}{\varphi_1 \varphi_2 - \varphi_{12}^2} (\varphi_1 u_{d'2} - \varphi_{12} u_{d'1}) + \sum_{i=1}^{I_{d'}} (y_{d'ij'} - m_{d'i} p_{d'ij'}) & \text{if } j' = 2 \end{cases} \quad (4.3)$$

Estimates for  $\boldsymbol{\beta}$  and  $\mathbf{u}$  are found by equating all the derivatives defined by (4.2) and (4.3) to 0 and solving the resulting system of equations. Because these equations are non-linear they cannot be solved directly. Instead they can be solved by using a Newton-Raphson algorithm. In order to use the Newton-Raphson method we will also need to work out all the second derivatives of the loglikelihood function.

For  $j' = 1, 2$ ;  $j'' = 1, 2$ ;  $q = 1, 2, \dots, Q_{j'}$  and  $q' = 1, 2, \dots, Q_{j''}$  it can be shown that

$$\frac{\partial^2 l(\boldsymbol{\beta}, \mathbf{u})}{\partial \beta_{j''(q')} \partial \beta_{j'(q)}} = \begin{cases} \sum_{d=1}^D \sum_{i=1}^{I_d} -x_{dij'(q)} x_{dij''(q')} m_{di} p_{dij'} (1 - p_{dij'}) & \text{if } j' = j'' \\ \sum_{d=1}^D \sum_{i=1}^{I_d} x_{dij'(q)} x_{dij''(q')} m_{di} p_{dij'} p_{dij''} & \text{if } j' \neq j''. \end{cases}$$

Also for  $d' = 1, 2, \dots, D$ ;  $d'' = 1, 2, \dots, D$ ;  $j' = 1, 2$  and  $j'' = 1, 2$ ,

$$\frac{\partial^2 l(\boldsymbol{\beta}, \mathbf{u})}{\partial u_{d''j''} \partial u_{d'j'}} = \begin{cases} \frac{-\varphi_2}{\varphi_1 \varphi_2 - \varphi_{12}^2} - \sum_{i=1}^{I_d} m_{d'i} p_{d'i1} (1 - p_{d'i1}) & \text{if } d'' = d' \text{ and } j'' = j' = 1 \\ \frac{-\varphi_1}{\varphi_1 \varphi_2 - \varphi_{12}^2} - \sum_{i=1}^{I_d} m_{d'i} p_{d'i2} (1 - p_{d'i2}) & \text{if } d'' = d' \text{ and } j'' = j' = 2 \\ \frac{\varphi_{12}}{\varphi_1 \varphi_2 - \varphi_{12}^2} + \sum_{i=1}^{I_d} m_{d'i} p_{d'i1} p_{d'i2} & \text{if } d'' = d' \text{ and } j'' \neq j' \\ 0 & \text{if } d'' \neq d', \end{cases}$$

and for  $d'' = 1, 2, \dots, D$ ;  $j'' = 1, 2$ ;  $j' = 1, 2$  and  $q = 1, 2, \dots, Q_{j'}$ ,

$$\frac{\partial^2 l(\boldsymbol{\beta}, \mathbf{u})}{\partial u_{d''j''} \partial \beta_{j'(q)}} = \begin{cases} -\sum_{i=1}^{I_d} x_{d''ij'(q)} m_{d''i} p_{d''ij'} (1 - p_{d''ij'}) & \text{if } j'' = j' \\ \sum_{i=1}^{I_d} x_{d''ij'(q)} m_{d''i} p_{d''ij'} p_{d''ij''} & \text{if } j'' \neq j'. \end{cases}$$

Similar to Molina *et al.* (2007), let  $\boldsymbol{\theta}_{di} = (\theta_{di1}, \theta_{di2})^t$ , where for  $j = 1, 2$ ,

$$\theta_{dij} = \log \frac{p_{dij}}{p_{di3}}.$$

We can also write (2.2) as

$$\boldsymbol{\theta}_{di} = \mathbf{X}_{di} \boldsymbol{\beta} + \mathbf{Z}_{di} \mathbf{u},$$

where

$$\mathbf{X}_{di} = \begin{pmatrix} \mathbf{x}_{di1}^t & \mathbf{0}_{1 \times Q_2} \\ \mathbf{0}_{1 \times Q_1} & \mathbf{x}_{di2}^t \end{pmatrix},$$

$$\mathbf{Z}_{di} = \begin{pmatrix} \mathbf{0}_{1 \times 2(d-1)} & 1 & 0 & \mathbf{0}_{1 \times 2(D-d)} \\ \mathbf{0}_{1 \times 2(d-1)} & 0 & 1 & \mathbf{0}_{1 \times 2(D-d)} \end{pmatrix}$$

and  $\mathbf{0}_{a^* \times b^*}$  denotes a matrix of zeros with dimension  $a^* \times b^*$ .



Denote the mean and covariance matrix of  $\mathbf{y}_{di}$  given  $\mathbf{u}$  as  $\boldsymbol{\mu}_{di}$  and  $\boldsymbol{\Sigma}_{di}$  and these are

$$\boldsymbol{\mu}_{di} = m_{di}(p_{di1}, p_{di2})^t$$

and

$$\boldsymbol{\Sigma}_{di} = m_{di} \begin{pmatrix} p_{di1}(1-p_{di1}) & -p_{di1}p_{di2} \\ -p_{di1}p_{di2} & p_{di2}(1-p_{di2}) \end{pmatrix}.$$

Again similar to Molina *et al.* (2007), let  $\mathbf{S}_\beta$  and  $\mathbf{S}_u$  be the vectors of first derivatives of the loglikelihood. That is,

$$\mathbf{S}_\beta = \left( \frac{\partial l(\boldsymbol{\beta}, \mathbf{u})}{\partial \beta_{1(1)}}, \frac{\partial l(\boldsymbol{\beta}, \mathbf{u})}{\partial \beta_{1(2)}}, \dots, \frac{\partial l(\boldsymbol{\beta}, \mathbf{u})}{\partial \beta_{1(Q_1)}}, \frac{\partial l(\boldsymbol{\beta}, \mathbf{u})}{\partial \beta_{2(1)}}, \frac{\partial l(\boldsymbol{\beta}, \mathbf{u})}{\partial \beta_{2(2)}}, \dots, \frac{\partial l(\boldsymbol{\beta}, \mathbf{u})}{\partial \beta_{2(Q_2)}} \right)^t$$

and

$$\mathbf{S}_u = \left( \frac{\partial l(\boldsymbol{\beta}, \mathbf{u})}{\partial u_{11}}, \frac{\partial l(\boldsymbol{\beta}, \mathbf{u})}{\partial u_{12}}, \frac{\partial l(\boldsymbol{\beta}, \mathbf{u})}{\partial u_{21}}, \frac{\partial l(\boldsymbol{\beta}, \mathbf{u})}{\partial u_{22}}, \dots, \frac{\partial l(\boldsymbol{\beta}, \mathbf{u})}{\partial u_{D1}}, \frac{\partial l(\boldsymbol{\beta}, \mathbf{u})}{\partial u_{D2}} \right)^t.$$

In matrix notation we have

$$\mathbf{S}_\beta = \sum_{d=1}^D \sum_{i=1}^{I_d} \mathbf{X}_{di}^t (\mathbf{y}_{di} - \boldsymbol{\mu}_{di})$$

and

$$\mathbf{S}_u = \sum_{d=1}^D \sum_{i=1}^{I_d} \mathbf{Z}_{di}^t (\mathbf{y}_{di} - \boldsymbol{\mu}_{di}) - \sum_{d=1}^D \mathbf{Z}_{d1}^t \mathbf{W}_d^{-1} \mathbf{Z}_{d1} \mathbf{u}$$

and the non-linear system of equations that we need to solve is described by

$$\mathbf{S} = (\mathbf{S}_\beta^t, \mathbf{S}_u^t)^t = \mathbf{0}_{(Q_1+Q_2+2D) \times 1}. \quad (4.4)$$

Let

$\mathbf{J}_\beta$  be a square symmetric matrix containing all the derivatives of  $\mathbf{S}_\beta$  with respect to  $\boldsymbol{\beta}$ ,

$\mathbf{J}_u$  be a square symmetric matrix containing all the derivatives of  $\mathbf{S}_u$  with respect to  $\mathbf{u}$ ,

$\mathbf{J}_{\beta u}$  be the matrix containing all the derivatives of  $\mathbf{S}_u$  with respect to  $\boldsymbol{\beta}$ .

It can be shown that

$$J_{\beta} = -\sum_{d=1}^D \sum_{i=1}^{I_d} \mathbf{X}_{di}^t \boldsymbol{\Sigma}_{di} \mathbf{X}_{di},$$

$$J_{\beta u} = -\sum_{d=1}^D \sum_{i=1}^{I_d} \mathbf{X}_{di}^t \boldsymbol{\Sigma}_{di} \mathbf{Z}_{di}$$

and

$$J_u = -\sum_{d=1}^D \sum_{i=1}^{I_d} \mathbf{Z}_{di}^t \boldsymbol{\Sigma}_{di} \mathbf{Z}_{di} - \sum_{d=1}^D \mathbf{Z}_{d1}^t \mathbf{W}_d^{-1} \mathbf{Z}_{d1}.$$

The Newton Raphson algorithm can now be applied to find the solution to (4.4). This iterative algorithm has updating equations as follows

$$\begin{pmatrix} \boldsymbol{\beta}^{k+1} \\ \mathbf{u}^{k+1} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\beta}^k \\ \mathbf{u}^k \end{pmatrix} - \begin{pmatrix} J_{\beta}^k & J_{\beta u}^k \\ (J_{\beta u}^k)^t & J_u^k \end{pmatrix}^{-1} \mathbf{S}^k, \quad (4.5)$$

where the superscript  $k$  indicates the iteration number and the current values of all parameters are used to evaluate the functions within. Note that in its current form (4.5) contains a large matrix which needs to be inverted at each step. This term can be further simplified by noting the following partitioned matrix identity given in Henderson and Searle (1981)

$$\begin{pmatrix} J_{\beta} & J_{\beta u} \\ J_{\beta u}^t & J_u \end{pmatrix}^{-1} = \begin{pmatrix} \left( J_{\beta} - J_{\beta u} J_u^{-1} J_{\beta u}^t \right)^{-1} & -\left( J_{\beta} - J_{\beta u} J_u^{-1} J_{\beta u}^t \right)^{-1} J_{\beta u} J_u^{-1} \\ -J_u^{-1} J_{\beta u}^t \left( J_{\beta} - J_{\beta u} J_u^{-1} J_{\beta u}^t \right)^{-1} & J_u^{-1} + J_u^{-1} J_{\beta u}^t \left( J_{\beta} - J_{\beta u} J_u^{-1} J_{\beta u}^t \right)^{-1} J_{\beta u} J_u^{-1} \end{pmatrix}, \quad (4.6)$$

for any square matrices  $J_{\beta}$  and  $J_u$  with  $J_u$  nonsingular and  $J_{\beta}$  possibly singular. Molina *et al.* (2007) also make use of this identity. Note that this identity simplifies the inversion quite a lot since we no longer need to invert a square matrix of dimension  $Q_1 + Q_2 + 2D$ . Instead we need to invert  $J_{\beta} - J_{\beta u} J_u^{-1} J_{\beta u}^t$  and  $J_u$ . The matrix  $J_{\beta} - J_{\beta u} J_u^{-1} J_{\beta u}^t$  is a square matrix of dimension  $Q_1 + Q_2$  which is the total number of explanatory variables in the model and is a lot smaller than  $Q_1 + Q_2 + 2D$  since  $D$  is large. As for  $J_u$ , this is a square matrix with dimension  $2D$  which is still quite large. But note that  $J_u$  is a block diagonal matrix with block sizes of 2. So all as we need to do in this case is invert a series of 2 by 2 matrices which is trivial.

So to compute PQL estimates of  $\beta$  and  $u$  we use the iterative formula (4.5) until convergence (and using the identity given at (4.6)). An initial guess is needed for the parameters to start off the iterations. We suggest using a small value of  $u$  as its initial value and an initial value of  $\beta$  obtained by fitting a multinomial logit model without random effects.

The criterion for convergence we use is one given in Booth and Hobert (1999)

$$\max \left( \begin{array}{l} \frac{|\beta_{j(q)}^{k+1} - \beta_{j(q)}^k|}{|\beta_{j(q)}^k| + \varepsilon_1}, j = 1, 2 \text{ and } q = 1, 2, \dots, Q_j, \\ \frac{|u_{dj}^{k+1} - u_{dj}^k|}{|u_{dj}^k| + \varepsilon_1}, d = 1, 2, \dots, D \text{ and } j = 1, 2 \end{array} \right) < \varepsilon_2, \quad (4.7)$$

where  $\varepsilon_1$  and  $\varepsilon_2$  are both small positive numbers (we use  $\varepsilon_1 = 0.01$  and  $\varepsilon_2 = 0.001$ ).

In all of the above it was assumed that  $\varphi$  is known. In the next section we derive an approximate maximum likelihood estimator of  $\varphi$  given the other terms. To obtain joint estimates of  $\beta$ ,  $u$  and  $\varphi$  we will need to iterate between updating each, where  $\beta$  and  $u$  will be updated using the algorithm in this section.

## 5. APPROXIMATE ML ESTIMATION OF $\varphi$

It was mentioned in Section 3 that for the normal linear mixed model, given the variance components, estimates of  $\boldsymbol{\beta}$  obtained by maximising  $L(\boldsymbol{\beta}, \boldsymbol{\varphi}, \boldsymbol{u})$  are equivalent to maximising  $L(\boldsymbol{\beta}, \boldsymbol{\varphi})$ . However for the normal linear mixed model, given  $\boldsymbol{\beta}$ , ML estimates for the variance components based on maximising  $L(\boldsymbol{\beta}, \boldsymbol{\varphi}, \boldsymbol{u})$  and  $L(\boldsymbol{\beta}, \boldsymbol{\varphi})$  are not equivalent. This suggests that it is probably inappropriate to maximise  $L(\boldsymbol{\beta}, \boldsymbol{\varphi}, \boldsymbol{u})$  to get estimates of  $\boldsymbol{\varphi}$  in our situation too.

To obtain an approximate marginal likelihood  $L(\boldsymbol{\beta}, \boldsymbol{\varphi})$ , Molina *et al.* (2007) adapted the ideas of Schall (1991) to their bivariate setting. We will follow this approach too.

Assume that  $\boldsymbol{\beta}$  and  $\boldsymbol{u}$  are known.

$$\text{Let } g_j(\boldsymbol{y}_{di}) = \log \frac{y_{dij}}{m_{di} - y_{di1} - y_{di2}} = \log \frac{y_{dij}}{y_{di3}} \quad \text{for } j = 1, 2.$$

A first order Taylor series expansion about the point  $\boldsymbol{\mu}_{di}$  leads to

$$g_j(\boldsymbol{y}_{di}) = g_j(\boldsymbol{\mu}_{di}) + \left. \frac{\partial g_j}{\partial y_{di1}} \right|_{\boldsymbol{\mu}_{di}} (y_{di1} - \mu_{di1}) + \left. \frac{\partial g_j}{\partial y_{di2}} \right|_{\boldsymbol{\mu}_{di}} (y_{di2} - \mu_{di2}), \text{ for } j = 1, 2.$$

Let  $\boldsymbol{\xi}_{di} = (g_1(\boldsymbol{y}_{di}), g_2(\boldsymbol{y}_{di}))^t$  and  $\boldsymbol{e}_{di} = \boldsymbol{\Sigma}_{di}^{-1} (\boldsymbol{y}_{di} - \boldsymbol{\mu}_{di})$ .

Calculating the expressions of the derivatives involved and using matrix notation, the above Taylor series expansion becomes

$$\boldsymbol{\xi}_{di} = \boldsymbol{X}_{di} \boldsymbol{\beta} + \boldsymbol{Z}_{di} \boldsymbol{u} + \boldsymbol{e}_{di},$$

where  $\text{Var}(\boldsymbol{e}_{di} | \boldsymbol{u}_d) = \boldsymbol{\Sigma}_{di}^{-1}$  and  $E(\boldsymbol{e}_{di} | \boldsymbol{u}_d) = \mathbf{0}_{2 \times 2}$ .

It is also clear that  $E(\boldsymbol{\xi}_{di} | \boldsymbol{u}_d) = \boldsymbol{X}_{di} \boldsymbol{\beta} + \boldsymbol{Z}_{di} \boldsymbol{u}$  and  $\text{Var}(\boldsymbol{\xi}_{di} | \boldsymbol{u}_d) = \boldsymbol{\Sigma}_{di}^{-1}$ .

Although it is not clear in Molina *et al.* (2007), on correspondence with the author the following was established. The term  $\boldsymbol{\xi}_{di} | \boldsymbol{u}_d$  is assumed to be approximately normal. Since  $\boldsymbol{u}_d$  is normal, this also implies that the joint distribution of  $\boldsymbol{\xi}_{di}$  and  $\boldsymbol{u}_d$  is approximately normal and hence the marginal distribution of  $\boldsymbol{\xi}_{di}$  is approximately normal.

Let  $\boldsymbol{W} = \text{diag}(\boldsymbol{W}_d, d = 1, 2, \dots, D)$ . Now it is easily shown that  $E(\boldsymbol{\xi}_{di}) = \boldsymbol{X}_{di} \boldsymbol{\beta}$  and

$$\text{Var}(\boldsymbol{\xi}_{di}) = E(\text{Var}(\boldsymbol{\xi}_{di} | \boldsymbol{u}_d)) + \text{Var}(E(\boldsymbol{\xi}_{di} | \boldsymbol{u}_d)) = E(\boldsymbol{\Sigma}_{di}^{-1}) + \boldsymbol{Z}_{di} \boldsymbol{W} \boldsymbol{Z}_{di}^t.$$

The matrix  $\boldsymbol{\Sigma}_{di}$  is a function of the random effects and we now need to assume that this is approximately constant and hence  $\text{Var}(\boldsymbol{\xi}_{di}) \approx \boldsymbol{\Sigma}_{di}^{-1} + \boldsymbol{Z}_{di} \boldsymbol{W} \boldsymbol{Z}_{di}^t$ , where the random effects are replaced by their estimated values.

Let  $\xi$  denote the vector that is constructed by stacking the vectors  $\xi_{di}$  in one column and  $V = \text{Var}(\xi)$ .

Then  $V = ZWZ^t + \Sigma^{-1}$ , where  $\Sigma = \text{diag}(\Sigma_{di}, i = 1, 2, \dots, I_d, d = 1, 2, \dots, D)$ .

Now we can define the approximated normal log likelihood for  $\varphi$  as (ignoring constant terms and assuming  $\beta$  is known)

$$l(\varphi) = -\frac{1}{2} \log |V| - \frac{1}{2} (\xi - X\beta)^t V^{-1} (\xi - X\beta),$$

where  $X$  is obtained by stacking the matrices  $X_{di}$ . Also let  $Z$  be the matrix obtained by stacking  $Z_{di}$ . Given  $\beta$ , we can now obtain approximate ML estimates of  $\varphi$  by maximising  $l(\varphi)$  with respect to  $\varphi$ . However, before we do this we need to simplify  $l(\varphi)$ .

To compute the PQL estimates of  $\beta$  and  $u$ , we use the updating equations (4.5). As noted by Jiang (2007), the following alternative iterative procedure originally proposed by Breslow and Clayton (1993) can also be used to produce the same PQL estimates. That is, for fixed  $\varphi$  compute

$$\hat{\beta} = (X^t V^{-1} X)^{-1} X^t V^{-1} \xi \quad (5.1)$$

and

$$\hat{u} = WZ^t V^{-1} (\xi - X\hat{\beta}), \quad (5.2)$$

where given  $\xi$  one may first use (5.1) to update  $\hat{\beta}$ , then use (5.2) to update  $\hat{u}$  then update  $\xi$  and so on until convergence. Note that we do not use (5.1) and (5.2) to update  $\beta$  and  $u$  because (4.5) is computationally more convenient. However the form of (5.1) and (5.2) is useful here because we can use these to help simplify  $l(\varphi)$ .

From (5.2), (5.1) and the results on page 446 in Pawitan (2001),

$$(\xi - X\hat{\beta})^t V^{-1} (\xi - X\hat{\beta}) = (\xi - X\hat{\beta} - Z\hat{u})^t \Sigma (\xi - X\hat{\beta} - Z\hat{u}) + \hat{u}^t W^{-1} \hat{u}$$

and

$$|V| = |\Sigma^{-1}| |W| |Z^t \Sigma Z + W^{-1}|.$$

We can now write an approximated pseudo loglikelihood of  $\varphi$  as (ignoring terms that are not functions of  $\varphi$ )

$$l(\varphi) = -\frac{1}{2} \log |W| - \frac{1}{2} \log |Z^t \Sigma Z + W^{-1}| - \frac{1}{2} \hat{u}^t W^{-1} \hat{u}. \quad (5.3)$$

Given current estimates of  $\beta$  and  $u$ , to obtain approximate ML estimates of  $\varphi$  we need to differentiate (5.3) with respect to  $\varphi$ , set the resulting three equations to 0 and solve them.

In the derivation of  $l(\varphi)$  it was assumed that  $\Sigma$  was constant and does not depend on  $u$  or  $\beta$ . When this is the case, joint estimation of  $\beta$ ,  $\varphi$  and  $u$  is equivalent to maximising the function

$$Q(\beta, \varphi, u) = l(\beta, \varphi, u) - \frac{1}{2} \log |Z^t \Sigma Z + W^{-1}| \quad (5.4)$$

since 
$$\frac{\partial Q}{\partial \beta} = \frac{\partial l(\beta, \varphi, u)}{\partial \beta}, \quad \frac{\partial Q}{\partial u} = \frac{\partial l(\beta, \varphi, u)}{\partial u} \quad \text{and} \quad \frac{\partial Q}{\partial \varphi} = \frac{\partial l(\varphi)}{\partial \varphi}.$$

This justifies the following algorithm to obtain joint estimates of  $\beta$ ,  $\varphi$  and  $u$

1. Compute  $\hat{\beta}$  and  $\hat{u}$  given  $\varphi$  by using (4.5).
2. Fixing  $\beta$  and  $u$  at their current values  $\hat{\beta}$  and  $\hat{u}$ , update  $\varphi$  by maximising (5.3).
3. Iterate between 1 and 2 until convergence.

However this algorithm assumes that  $\Sigma$  is constant in (5.3). Pawitan (2001) notes that this algorithm is appropriate when  $\Sigma^{-1}$  is a slowly varying function of  $\mu$  ( $\mu$  is the vector obtained by stacking all the  $\mu_{di}$ ). This means we can ignore the derivative of the second term of  $Q$  with respect to  $\beta$  and  $u$ , so the first step is justified.

Pawitan (2001) notes that for certain generalised linear mixed models studied by Breslow and Clayton (1993), estimates of the variance components based on maximising (5.4) are close to the exact marginal likelihood estimates provided that the variance component is not too large. The method tends to underestimate the variance component, and the problem can be severe for large values of the variance component. However in our application, the variance components are expected to be small, so hopefully this should not be too much of an issue.

We now need to maximise (5.3) and to do this we need to differentiate this function with respect to each component of  $\varphi$ . For  $a = 1, 2$  and  $12$ ,

$$\begin{aligned} \frac{\partial l(\varphi)}{\partial \varphi_a} &= -\frac{1}{2} \frac{\partial \log |W|}{\partial \varphi_a} - \frac{1}{2} \frac{\partial \log |Z^t \Sigma Z + W^{-1}|}{\partial \varphi_a} - \frac{1}{2} \frac{\partial u^t W^{-1} u}{\partial \varphi_a} \\ &= -\frac{1}{2} Tr \left[ W^{-1} \frac{\partial W}{\partial \varphi_a} \right] - \frac{1}{2} Tr \left[ \left( Z^t \Sigma Z + W^{-1} \right)^{-1} \frac{\partial W^{-1}}{\partial \varphi_a} \right] - \frac{1}{2} u^t \frac{\partial W^{-1}}{\partial \varphi_a} u, \quad (5.5) \end{aligned}$$

where  $Tr[ \ ]$  denotes matrix trace.

After some algebra, it can be shown that

$$\text{Tr} \left[ \mathbf{W}^{-1} \frac{\partial \mathbf{W}}{\partial \varphi_1} \right] = \frac{D\varphi_2}{\varphi_1\varphi_2 - \varphi_{12}^2}, \quad (5.6)$$

$$\text{Tr} \left[ \mathbf{W}^{-1} \frac{\partial \mathbf{W}}{\partial \varphi_2} \right] = \frac{D\varphi_1}{\varphi_1\varphi_2 - \varphi_{12}^2} \quad (5.7)$$

and

$$\text{Tr} \left[ \mathbf{W}^{-1} \frac{\partial \mathbf{W}}{\partial \varphi_{12}} \right] = \frac{-2D\varphi_{12}}{\varphi_1\varphi_2 - \varphi_{12}^2}. \quad (5.8)$$

Also, 
$$\mathbf{u}^t \frac{\partial \mathbf{W}^{-1}}{\partial \varphi_1} \mathbf{u} = \frac{1}{(\varphi_1\varphi_2 - \varphi_{12}^2)^2} \sum_{d=1}^D \mathbf{u}_d^t \begin{pmatrix} -\varphi_2^2 & \varphi_2\varphi_{12} \\ \varphi_2\varphi_{12} & -\varphi_{12}^2 \end{pmatrix} \mathbf{u}_d, \quad (5.9)$$

$$\mathbf{u}^t \frac{\partial \mathbf{W}^{-1}}{\partial \varphi_2} \mathbf{u} = \frac{1}{(\varphi_1\varphi_2 - \varphi_{12}^2)^2} \sum_{d=1}^D \mathbf{u}_d^t \begin{pmatrix} -\varphi_{12}^2 & \varphi_1\varphi_{12} \\ \varphi_1\varphi_{12} & -\varphi_1^2 \end{pmatrix} \mathbf{u}_d \quad (5.10)$$

and 
$$\mathbf{u}^t \frac{\partial \mathbf{W}^{-1}}{\partial \varphi_{12}} \mathbf{u} = \frac{1}{(\varphi_1\varphi_2 - \varphi_{12}^2)^2} \sum_{d=1}^D \mathbf{u}_d^t \begin{pmatrix} 2\varphi_{12}\varphi_2 & -(\varphi_{12}^2 + \varphi_1\varphi_2) \\ -(\varphi_{12}^2 + \varphi_1\varphi_2) & 2\varphi_1\varphi_{12} \end{pmatrix} \mathbf{u}_d. \quad (5.11)$$

Now for  $d = 1, 2, \dots, D$  let

$$\begin{aligned} q_{d1} &= \sum_{i=1}^{I_d} m_{di} p_{di1} (1 - p_{di1}), \\ q_{d2} &= \sum_{i=1}^{I_d} m_{di} p_{di2} (1 - p_{di2}) \text{ and} \\ q_{d3} &= \sum_{i=1}^{I_d} -m_{di} p_{di1} p_{di2}. \end{aligned}$$

After some algebra it can be shown that

$$\begin{aligned} \text{Tr} \left[ \left( \mathbf{Z}' \boldsymbol{\Sigma} \mathbf{Z} + \mathbf{W}^{-1} \right)^{-1} \frac{\partial \mathbf{W}^{-1}}{\partial \varphi_1} \right] &= \\ - \sum_{d=1}^D \left( \frac{q_{d1}\varphi_{12}^2 + \varphi_2(1 + 2q_{d3}\varphi_{12} + q_{d2}\varphi_2)}{(-\varphi_{12}^2 + \varphi_1\varphi_2) \left( (1 + q_{d3}\varphi_{12})^2 + (q_{d2} - q_{d3}^2\varphi_1)\varphi_2 + q_{d1}(\varphi_1 - q_{d2}\varphi_{12}^2 + q_{d2}\varphi_1\varphi_2) \right)} \right), \end{aligned} \quad (5.12)$$

$$Tr \left[ \left( \mathbf{Z}' \boldsymbol{\Sigma} \mathbf{Z} + \mathbf{W}^{-1} \right)^{-1} \frac{\partial \mathbf{W}^{-1}}{\partial \varphi_2} \right] =$$

$$- \sum_{d=1}^D \left( \frac{\varphi_1 + q_{d1} \varphi_1^2 + 2q_{d3} \varphi_1 \varphi_{12} + q_{d2} \varphi_{12}^2}{\left( -\varphi_{12}^2 + \varphi_1 \varphi_2 \right) \left( \left( 1 + q_{d3} \varphi_{12} \right)^2 + \left( q_{d2} - q_{d3}^2 \varphi_1 \right) \varphi_2 + q_{d1} \left( \varphi_1 - q_{d2} \varphi_{12}^2 + q_{d2} \varphi_1 \varphi_2 \right) \right)} \right)$$

(5.13)

and

$$Tr \left[ \left( \mathbf{Z}' \boldsymbol{\Sigma} \mathbf{Z} + \mathbf{W}^{-1} \right)^{-1} \frac{\partial \mathbf{W}^{-1}}{\partial \varphi_{12}} \right] =$$

$$-2 \sum_{d=1}^D \left( \frac{\left( \varphi_{12} + q_{d1} \varphi_1 \varphi_{12} + q_{d3} \varphi_{12}^2 + q_{d3} \varphi_1 \varphi_2 + q_{d2} \varphi_{12} \varphi_2 \right)}{\left( \varphi_{12}^2 - \varphi_1 \varphi_2 \right) \left( \left( 1 + q_{d3} \varphi_{12} \right)^2 + \left( q_{d2} - q_{d3}^2 \varphi_1 \right) \varphi_2 + q_{d1} \left( \varphi_1 - q_{d2} \varphi_{12}^2 + q_{d2} \varphi_1 \varphi_2 \right) \right)} \right).$$

(5.14)

By substituting the relevant terms (5.6)–(5.14) into (5.5), we can now calculate the first derivatives of  $l(\boldsymbol{\varphi})$  with respect to  $\varphi_1, \varphi_2$  and  $\varphi_{12}$ .

Let

$$\mathbf{S}_{\boldsymbol{\varphi}} = \left( \frac{\partial l(\boldsymbol{\varphi})}{\partial \varphi_1}, \frac{\partial l(\boldsymbol{\varphi})}{\partial \varphi_2}, \frac{\partial l(\boldsymbol{\varphi})}{\partial \varphi_{12}} \right)^t.$$

To obtain an update for  $\boldsymbol{\varphi}$  given  $\boldsymbol{\beta}$  and  $\mathbf{u}$  we need to solve  $\mathbf{S}_{\boldsymbol{\varphi}} = \mathbf{0}_{3 \times 1}$ . The multinomial model described in the paper Molina *et al.* (2007) has one variance component. In this case, an explicit updating formula is available for the variance component based on rearranging the single equation

$$\frac{dl(\boldsymbol{\varphi})}{d\boldsymbol{\varphi}} = 0.$$

An experiment was undertaken to determine whether simple updating equations could be obtained by rearranging the equations  $\mathbf{S}_{\boldsymbol{\varphi}} = \mathbf{0}_{3 \times 1}$  but we could not get it to converge. Unfortunately we will need to use a single iteration of the Newton-Raphson algorithm to update  $\boldsymbol{\varphi}$  instead, which means we will also need all the second derivatives of  $l(\boldsymbol{\varphi})$ .



After some algebra it can be shown that

$$\frac{\partial^2 l(\boldsymbol{\varphi})}{\partial \varphi_1^2} = -\frac{1}{(-\varphi_{12}^2 + \varphi_1 \varphi_2)^3} \sum_{d=1}^D \mathbf{u}_d^t \begin{pmatrix} \varphi_2^3 & -\varphi_{12} \varphi_2^2 \\ -\varphi_{12} \varphi_2^2 & \varphi_{12}^2 \varphi_2 \end{pmatrix} \mathbf{u}_d \\ + \sum_{d=1}^D \frac{(q_{d1} + q_{d1} q_{d2} \varphi_2 - q_{d3}^2 \varphi_2)^2}{2 \left( (1 + q_{d3} \varphi_{12})^2 + (q_{d2} - q_{d3}^2 \varphi_1) \varphi_2 + q_{d1} (\varphi_1 - q_{d2} \varphi_{12}^2 + q_{d2} \varphi_1 \varphi_2) \right)^2},$$

$$\frac{\partial^2 l(\boldsymbol{\varphi})}{\partial \varphi_1 \partial \varphi_2} = -\frac{1}{2(-\varphi_{12}^2 + \varphi_1 \varphi_2)^3} \sum_{d=1}^D \mathbf{u}_d^t \begin{pmatrix} 2\varphi_{12}^2 \varphi_2 & -(\varphi_{12}^3 + \varphi_1 \varphi_{12} \varphi_2) \\ -(\varphi_{12}^3 + \varphi_1 \varphi_{12} \varphi_2) & 2\varphi_1 \varphi_{12}^2 \end{pmatrix} \mathbf{u}_d \\ + \sum_{d=1}^D \frac{(q_{d3} - q_{d1} q_{d2} \varphi_{12} + q_{d3}^2 \varphi_{12})^2}{2 \left( (1 + q_{d3} \varphi_{12})^2 + (q_{d2} - q_{d3}^2 \varphi_1) \varphi_2 + q_{d1} (\varphi_1 - q_{d2} \varphi_{12}^2 + q_{d2} \varphi_1 \varphi_2) \right)^2},$$

$$\frac{\partial^2 l(\boldsymbol{\varphi})}{\partial \varphi_1 \partial \varphi_{12}} = -\frac{1}{2(-\varphi_{12}^2 + \varphi_1 \varphi_2)^3} \sum_{d=1}^D \mathbf{u}_d^t \begin{pmatrix} -4\varphi_{12} \varphi_2^2 & \varphi_2 (3\varphi_{12}^2 + \varphi_1 \varphi_2) \\ \varphi_2 (3\varphi_{12}^2 + \varphi_1 \varphi_2) & -2(\varphi_{12}^3 + \varphi_1 \varphi_{12} \varphi_2) \end{pmatrix} \mathbf{u}_d \\ - \sum_{d=1}^D \frac{(q_{d3} - q_{d1} q_{d2} \varphi_{12} + q_{d3}^2 \varphi_{12})(q_{d3}^2 \varphi_2 - q_{d1} (1 + q_{d2} \varphi_2))}{\left( (1 + q_{d3} \varphi_{12})^2 + (q_{d2} - q_{d3}^2 \varphi_1) \varphi_2 + q_{d1} (\varphi_1 - q_{d2} \varphi_{12}^2 + q_{d2} \varphi_1 \varphi_2) \right)^2},$$

$$\frac{\partial^2 l(\boldsymbol{\varphi})}{\partial \varphi_2^2} = -\frac{1}{(-\varphi_{12}^2 + \varphi_1 \varphi_2)^3} \sum_{d=1}^D \mathbf{u}_d^t \begin{pmatrix} \varphi_1 \varphi_{12}^2 & -\varphi_1^2 \varphi_{12} \\ -\varphi_1^2 \varphi_{12} & \varphi_1^3 \end{pmatrix} \mathbf{u}_d \\ + \sum_{d=1}^D \frac{(q_{d2} + q_{d1} q_{d2} \varphi_1 - q_{d3}^2 \varphi_1)^2}{2 \left( (1 + q_{d3} \varphi_{12})^2 + (q_{d2} - q_{d3}^2 \varphi_1) \varphi_2 + q_{d1} (\varphi_1 - q_{d2} \varphi_{12}^2 + q_{d2} \varphi_1 \varphi_2) \right)^2},$$

$$\frac{\partial^2 l(\boldsymbol{\varphi})}{\partial \varphi_2 \partial \varphi_{12}} = -\frac{1}{2(-\varphi_{12}^2 + \varphi_1 \varphi_2)^3} \sum_{d=1}^D \mathbf{u}_d^t \begin{pmatrix} -2(\varphi_{12}^3 + \varphi_1 \varphi_{12} \varphi_2) & \varphi_1 (3\varphi_{12}^2 + \varphi_1 \varphi_2) \\ \varphi_1 (3\varphi_{12}^2 + \varphi_1 \varphi_2) & -4\varphi_1^2 \varphi_{12} \end{pmatrix} \mathbf{u}_d \\ - \sum_{d=1}^D \frac{(q_{d2} + q_{d1} q_{d2} \varphi_1 - q_{d3}^2 \varphi_1)(q_{d1} q_{d2} \varphi_{12} - q_{d3} (1 + q_{d3} \varphi_{12}))}{\left( (1 + q_{d3} \varphi_{12})^2 + (q_{d2} - q_{d3}^2 \varphi_1) \varphi_2 + q_{d1} (\varphi_1 - q_{d2} \varphi_{12}^2 + q_{d2} \varphi_1 \varphi_2) \right)^2}$$

$$\begin{aligned} \frac{\partial^2 l(\boldsymbol{\varphi})}{\partial \varphi_{12}^2} = & -\frac{1}{(-\varphi_{12}^2 + \varphi_1 \varphi_2)^3} \sum_{d=1}^D \mathbf{u}_d^t \begin{pmatrix} \varphi_2 (3\varphi_{12}^2 + \varphi_1 \varphi_2) & -(\varphi_{12}^3 + 3\varphi_1 \varphi_{12} \varphi_2) \\ -(\varphi_{12}^3 + 3\varphi_1 \varphi_{12} \varphi_2) & \varphi_1 (3\varphi_{12}^2 + \varphi_1 \varphi_2) \end{pmatrix} \mathbf{u}_d \\ & + \sum_{d=1}^D \frac{2(q_{d3} - q_{d1} q_{d2} \varphi_{12} + q_{d3}^2 \varphi_{12})^2}{\left( (1 + q_{d3} \varphi_{12})^2 + (q_{d2} - q_{d3}^2 \varphi_1) \varphi_2 + q_{d1} (\varphi_1 - q_{d2} \varphi_{12}^2 + q_{d2} \varphi_1 \varphi_2) \right)^2} \\ & - \sum_{d=1}^D \frac{-(q_{d1} q_{d2}) + q_{d3}^2}{(1 + q_{d3} \varphi_{12})^2 + (q_{d2} - q_{d3}^2 \varphi_1) \varphi_2 + q_{d1} (\varphi_1 - q_{d2} \varphi_{12}^2 + q_{d2} \varphi_1 \varphi_2)}. \end{aligned}$$

Now let

$$J_{\boldsymbol{\varphi}} = \begin{pmatrix} \frac{\partial^2 l(\boldsymbol{\varphi})}{\partial^2 \varphi_1} & \frac{\partial^2 l(\boldsymbol{\varphi})}{\partial \varphi_1 \partial \varphi_2} & \frac{\partial^2 l(\boldsymbol{\varphi})}{\partial \varphi_1 \partial \varphi_{12}} \\ \frac{\partial^2 l(\boldsymbol{\varphi})}{\partial \varphi_1 \partial \varphi_2} & \frac{\partial^2 l(\boldsymbol{\varphi})}{\partial^2 \varphi_2} & \frac{\partial^2 l(\boldsymbol{\varphi})}{\partial \varphi_2 \partial \varphi_{12}} \\ \frac{\partial^2 l(\boldsymbol{\varphi})}{\partial \varphi_1 \partial \varphi_{12}} & \frac{\partial^2 l(\boldsymbol{\varphi})}{\partial \varphi_2 \partial \varphi_{12}} & \frac{\partial^2 l(\boldsymbol{\varphi})}{\partial^2 \varphi_{12}} \end{pmatrix}.$$

To obtain an update of  $\boldsymbol{\varphi}$  given  $\boldsymbol{\beta}$  and  $\mathbf{u}$  we simply compute

$$\boldsymbol{\varphi}^{\text{update}} = \boldsymbol{\varphi}^{\text{previous}} - J_{\boldsymbol{\varphi}}^{-1} \mathbf{S}_{\boldsymbol{\varphi}},$$

where all terms on the right hand side are evaluated using the current values. Only one update is needed because after  $\boldsymbol{\varphi}$  is updated once we then have to go back to updating  $\boldsymbol{\beta}$  and  $\mathbf{u}$  by using repeats of (4.5) until convergence given the current value of  $\boldsymbol{\varphi}$ .

The iterative algorithm for computing joint estimates of  $\boldsymbol{\beta}$ ,  $\mathbf{u}$  and  $\boldsymbol{\varphi}$  stops when both (4.7) and the following condition is satisfied for all three variance components

$$\frac{|\varphi_a^{\text{update}} - \varphi_a^{\text{previous}}|}{|\varphi_a^{\text{previous}}| + \varepsilon_1} < \varepsilon_2,$$

where subscript  $a$  denotes the particular component.

## 6. APPROXIMATE REML ESTIMATION OF $\boldsymbol{\varphi}$

Approximate REML (Restricted Maximum Likelihood) estimation could also be used as an alternative to approximate ML estimation for estimating  $\boldsymbol{\varphi}$ . This approach was also undertaken by Molina *et al.* (2007).

For normal linear mixed models, it is well known that in some cases the ML estimator of the variance components can be downwardly biased (see page 235 of Rao and Kleffe (1988)). As mentioned in Harville (1977), one criticism of the ML approach to estimating variance components is that the ML estimator takes no account of the loss of degrees of freedom for estimating the fixed effects  $\boldsymbol{\beta}$ . For estimating the variance components,  $\boldsymbol{\beta}$  can be considered as a nuisance parameter. A way of eliminating the influence of  $\boldsymbol{\beta}$  is to construct a marginal (or approximated marginal) ML estimator for the variance components. In the normal linear mixed model case, this involves transforming the original observations such that the new data are independent of  $\boldsymbol{\beta}$  and then maximising the likelihood for the variance components of this new data. The transformed data are called error contrasts as described in Harville (1977) and have smaller dimension than the original data.

In the normal linear mixed model case, the REML loglikelihood can be derived as a modified profile likelihood (see pages 286–292 in Pawitan (2001)). In our case the approximated REML loglikelihood is proportional to

$$l(\boldsymbol{\varphi}) - \frac{1}{2} \log |\mathbf{X}'\mathbf{V}^{-1}\mathbf{X}|, \quad (6.1)$$

where  $l(\boldsymbol{\varphi})$  is given by (5.3) and as Pawitan (2001) mentions, the second term can be interpreted as a penalty term, subtracting from the profile loglikelihood the ‘undeserved’ information on the nuisance parameter (which is  $\boldsymbol{\beta}$ ). To obtain approximate REML estimates of  $\boldsymbol{\varphi}$  in our case, we simply need to maximise (6.1) which can be done using similar methods to what is described in Section 5. That is, we need to implement the same Newton-Raphson algorithm except that now  $\mathbf{S}_{\boldsymbol{\varphi}}$  and  $\mathbf{J}_{\boldsymbol{\varphi}}$  contain extra terms. These extra terms can be obtained by finding all the first and second derivatives of the second term in (6.1) with respect to all the variance components.

Although we are potentially reducing bias by using the approximated REML estimator we should keep in mind that the variance may be larger than the variance of the approximated ML estimator. Therefore in some situations the MSE of the approximated ML estimator might be smaller. See Harville (1977) for further details for a discussion in relation to normal linear mixed models.

In any case we need all the first and second derivatives of  $\log |\mathbf{X}'\mathbf{V}^{-1}\mathbf{X}|$ .

For  $a = 1, 2$  and  $12$ ,

$$\frac{\partial \log |\mathbf{X}'\mathbf{V}^{-1}\mathbf{X}|}{\partial \varphi_a} = \text{Tr} \left[ \left( \mathbf{X}'\mathbf{V}^{-1}\mathbf{X} \right)^{-1} \mathbf{X}' \frac{\partial \mathbf{V}^{-1}}{\partial \varphi_a} \mathbf{X} \right].$$

Now 
$$\frac{\partial \mathbf{V}^{-1}}{\partial \varphi_a} = -\mathbf{V}^{-1} \frac{\partial \mathbf{V}}{\partial \varphi_a} \mathbf{V}^{-1} = -\mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}' \mathbf{V}^{-1}$$

and therefore

$$\frac{\partial \log |\mathbf{X}'\mathbf{V}^{-1}\mathbf{X}|}{\partial \varphi_a} = -\text{Tr} \left[ \left( \mathbf{X}'\mathbf{V}^{-1}\mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{X} \right], \quad (6.2)$$

where 
$$\frac{\partial \mathbf{W}}{\partial \varphi_a} = \text{diag} \left( \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \right) \text{ or } \text{diag} \left( \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \right) \text{ or } \text{diag} \left( \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \right),$$

depending on which value of  $a$  is chosen.

On first glimpse (6.2) looks quite difficult to calculate, but note that the matrix within the trace is only a  $Q_1 + Q_2$  dimensional matrix in total. Because  $\mathbf{V}$  and  $\mathbf{W}$  are block diagonal ( $\mathbf{V}$  has  $d = 1, 2, \dots, D$  blocks each of size  $2I_d$ ) simplifications can also be made within the matrix multiplications in the calculations.

Now we need to calculate the second derivatives. For  $a = 1, 2, 12$  and  $b = 1, 2, 12$ ,

$$\begin{aligned} \frac{\partial^2 \log |\mathbf{X}'\mathbf{V}^{-1}\mathbf{X}|}{\partial \varphi_b \partial \varphi_a} &= -2 \text{Tr} \left[ \left( \mathbf{X}'\mathbf{V}^{-1}\mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}' \frac{\partial \mathbf{V}^{-1}}{\partial \varphi_b} \mathbf{X} \right] \\ &\quad - \text{Tr} \left[ \frac{\partial \left( \mathbf{X}'\mathbf{V}^{-1}\mathbf{X} \right)^{-1}}{\partial \varphi_b} \mathbf{X}' \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{X} \right] \\ &= 2 \text{Tr} \left[ \left( \mathbf{X}'\mathbf{V}^{-1}\mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_b} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{X} \right] \\ &\quad + \text{Tr} \left[ \left( \mathbf{X}'\mathbf{V}^{-1}\mathbf{X} \right)^{-1} \mathbf{X}' \frac{\partial \mathbf{V}^{-1}}{\partial \varphi_b} \mathbf{X} \left( \mathbf{X}'\mathbf{V}^{-1}\mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{X} \right] \\ &= 2 \text{Tr} \left[ \left( \mathbf{X}'\mathbf{V}^{-1}\mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_b} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{X} \right] \\ &\quad - \text{Tr} \left[ \left( \mathbf{X}'\mathbf{V}^{-1}\mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_b} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{X} \left( \mathbf{X}'\mathbf{V}^{-1}\mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{X} \right]. \end{aligned}$$

Now let

$$\mathbf{S}_\varphi^* = \mathbf{S}_\varphi - \frac{1}{2} \left( \frac{\partial \log |\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}|}{\partial \varphi_1}, \frac{\partial \log |\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}|}{\partial \varphi_2}, \frac{\partial \log |\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}|}{\partial \varphi_{12}} \right)^t$$

and

$$\mathbf{J}_\varphi^* = \mathbf{J}_\varphi - \frac{1}{2} \begin{pmatrix} \frac{\partial^2 \log |\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}|}{\partial^2 \varphi_1} & \frac{\partial^2 \log |\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}|}{\partial \varphi_1 \partial \varphi_2} & \frac{\partial^2 \log |\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}|}{\partial \varphi_1 \partial \varphi_{12}} \\ \frac{\partial^2 \log |\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}|}{\partial \varphi_1 \partial \varphi_2} & \frac{\partial^2 \log |\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}|}{\partial^2 \varphi_2} & \frac{\partial^2 \log |\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}|}{\partial \varphi_2 \partial \varphi_{12}} \\ \frac{\partial^2 \log |\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}|}{\partial \varphi_1 \partial \varphi_{12}} & \frac{\partial^2 \log |\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}|}{\partial \varphi_2 \partial \varphi_{12}} & \frac{\partial^2 \log |\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}|}{\partial^2 \varphi_{12}} \end{pmatrix},$$

where  $\mathbf{S}_\varphi$  and  $\mathbf{J}_\varphi$  are as defined in Section 5.

To obtain a REML update of  $\boldsymbol{\varphi}$  given  $\boldsymbol{\beta}$  and  $\mathbf{u}$  we compute

$$\boldsymbol{\varphi}^{\text{update}} = \boldsymbol{\varphi}^{\text{previous}} - (\mathbf{J}_\varphi^*)^{-1} \mathbf{S}_\varphi^*,$$

where all terms on the right hand side are evaluated using the current values.

## 7. EMPIRICAL BEST PREDICTION, SMALL AREA ESTIMATION AND A NOTE ON MSE ESTIMATION

Now that we have derived the PQL and REML estimators of the model parameters, we would like to use these to help predict the small area totals of employed, unemployed and not in the labour force. These sets of totals are, for  $j = 1, 2, 3$  and  $d = 1, 2, \dots, D$ ,

$$\delta_{dj} = \sum_{i=1}^{I_d} y_{dij} + \sum_{i=1}^I y_{dij}^r$$

where  $y_{dij}^r$  denotes the total non sample in small area  $d$ , age/sex group  $i$  and labour force class  $j$ . Note that  $I = 10$  (total number of age/sex classes within a small area). In previous sections, summations were defined using  $I_d$  and classes  $I_d + 1, \dots, I$  were excluded because they did not contribute to the likelihood functions since there was no sample in these classes. However, we now need to include these.

Obviously we cannot use  $\delta_{dj}$  directly because all the  $y_{dij}^r$  are unknown. If we knew the values of  $\boldsymbol{\beta}$  and  $\boldsymbol{u}$  then we could estimate  $\delta_{dj}$  using

$$\tilde{\delta}_{dj} = \sum_{i=1}^{I_d} y_{dij} + \sum_{i=1}^I \mu_{dij}^r$$

where

$$\mu_{dij}^r = \begin{cases} m_{di}^r \frac{e^{\mathbf{x}_{di1}^t \boldsymbol{\beta}_1 + \mathbf{u}_{d1}}}{1 + e^{\mathbf{x}_{di1}^t \boldsymbol{\beta}_1 + \mathbf{u}_{d1}} + e^{\mathbf{x}_{di2}^t \boldsymbol{\beta}_2 + \mathbf{u}_{d2}}} & \text{if } j = 1 \\ m_{di}^r \frac{e^{\mathbf{x}_{di2}^t \boldsymbol{\beta}_2 + \mathbf{u}_{d2}}}{1 + e^{\mathbf{x}_{di1}^t \boldsymbol{\beta}_1 + \mathbf{u}_{d1}} + e^{\mathbf{x}_{di2}^t \boldsymbol{\beta}_2 + \mathbf{u}_{d2}}} & \text{if } j = 2 \\ m_{di}^r \frac{1}{1 + e^{\mathbf{x}_{di1}^t \boldsymbol{\beta}_1 + \mathbf{u}_{d1}} + e^{\mathbf{x}_{di2}^t \boldsymbol{\beta}_2 + \mathbf{u}_{d2}}} & \text{if } j = 3, \end{cases}$$

and  $m_{di}^r$  is the total non-sample in sex/age group  $i$ , in small area  $d$  and these are assumed known. Again  $\tilde{\delta}_{dj}$  cannot be used here directly since  $\boldsymbol{\beta}$  and  $\boldsymbol{u}$  are unknown.

The prediction problem we now have is one where we need to predict

$$\tau_{dj} = \sum_{i=1}^I \mu_{dij}^r,$$

for  $j = 1, 2$  ( $j = 3$  is not needed since it can be obtained via subtraction since  $m_{di}^r$  is known). There are two ways to predict  $\tau_{dj}$ . The first is to simply replace  $\boldsymbol{u}$  with the PQL estimate  $\hat{\boldsymbol{u}}$  and replace  $\boldsymbol{\beta}$  with the PQL estimate  $\hat{\boldsymbol{\beta}}$ .

That is, 
$$\hat{\tau}_{dj} = \sum_{i=1}^I \hat{\mu}_{dij}^r, \quad (7.1)$$

where the hat on  $\mu_{dij}^r$  means that we are using estimates of  $\beta$  and  $u$  within this function. This is the approach taken by Molina *et al.* (2007). An alternative way of predicting  $\tau_{dj}$  is to use an empirical best predictor (EBP) as defined by Jiang (2007, pp. 143–144). Assuming that  $\beta$  and  $\phi$  are known, the best predictor in the sense of minimum MSE of  $\tau_{dj}$  would be

$$E(\tau_{dj} | \mathbf{y}) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \tau_{dj} f(\mathbf{u}_d | \beta, \phi; \mathbf{y}_d) du_{d1} du_{d2}.$$

An empirical best predictor replaces  $\beta$  and  $\phi$  in the integral above with respectively the PQL and REML estimates  $\hat{\beta}$  and  $\hat{\phi}$ . The double integral above has no closed form solution but we could approximate it using Monte Carlo methods.

Because we need to approximate the double integrals in EBP using Monte Carlo methods, this adds to the computation time and MSE estimation becomes a problem. For instance, bootstrap MSEs (and those based on other resampling methods) are not computationally feasible for EBP. Alternatively Jiang (2007) on page 144 describes a method of approximating the MSEs of EBP based on a Taylor Series expansion that gives an estimate whose bias is corrected to the second order. However in the derivation it is assumed that the estimates of  $\beta$  and  $\phi$  have certain properties which PQL estimators may not satisfy (see page 158). Therefore these MSE estimates may not be correct to second order in our context. Another problem with this MSE estimator is that one of its terms relates to the expected value of the EBP squared over the distribution of  $\mathbf{y}$  and this calculation is not computationally feasible for our model.

One advantage of using the estimator at (7.1) is that MSE estimation is relatively straight forward. Bootstrap MSEs are computationally feasible or alternatively, estimators based on Taylor series approximations can be used. In particular the Taylor series approximation method in Molina *et al.* (2007) could be extended to our category specific multinomial mixed model case. In any case it is not even clear whether the MSEs will be smaller for the EBP than those for (7.1). Therefore at this stage EBP is not recommended.

## 8. ANALYTICAL APPROXIMATION OF THE MSE

For reasons discussed in the previous section we will be using the PQL estimators of  $\beta$  and  $u$  along with (7.1) to predict the small area totals. To approximate the MSEs of these small area total estimates, we use a similar approach to the one described in Molina *et al.* (2007) based on an analytical approximation. We begin by briefly outlining this approach and then later in this section we describe the technical details of the MSE approximation.

As mentioned in Molina *et al.* (2007), under linear mixed models, Prasad and Rao (1990) obtained an analytical approximation of the MSE of an estimator of the type  $t(\hat{v}) = \lambda' \hat{\beta} + m' \hat{u}$  where  $\lambda$  and  $m$  are vectors of constants and  $v$  and  $\hat{v}$  are respectively the vector of variance components and the estimated vector of variance components. The Prasad and Rao (1990) approximation takes the form

$$MSE(t(\hat{v})) = G_1(v) + G_2(v) + G_3(v)$$

and the estimator is given by

$$\widehat{MSE}(t(\hat{v})) = G_1(\hat{v}) + G_2(\hat{v}) + 2G_3(\hat{v})$$

which corrects for the bias in the  $G_1(\hat{v})$  term. In the derivation, Prasad and Rao used a result from Kackar and Harville (1984) who showed that under certain conditions

$$MSE(t(\hat{v})) = MSE(t(v)) + E((t(\hat{v}) - t(v))^2).$$

In the Prasad and Rao context

$$MSE(t(v)) = G_1(v) + G_2(v) \quad \text{and} \quad G_3(v) = E((t(\hat{v}) - t(v))^2)$$

and Prasad and Rao (1990) proposed a new approximation to  $G_3(v)$  relevant in a small area estimation context.

Prasad and Rao's formula was adapted by Baillo and Molina (2005) to a multivariate mixed linear model and a multidimensional parameter. Both the Prasad and Rao (1990) and the Baillo and Molina (2005) MSE estimators are for linear mixed models. These estimators can be adapted to our context by noting that our model can be written as an approximate bivariate linear mixed model (see Section 5 for further details). In our context we use the predictions (7.1) for the non-sample labour force counts and in order to apply the Prasad and Rao (1990) and Baillo and Molina (2005) MSE approximations, we need to linearise (7.1) using a first order Taylor series



approximation (Molina *et al.* (2007) use a similar approach). This then results in a MSE approximation of the form  $\mathbf{G}_1(\boldsymbol{\varphi}) + \mathbf{G}_2(\boldsymbol{\varphi}) + \mathbf{G}_3(\boldsymbol{\varphi})$ . There is also an additional term needed called  $\mathbf{G}_4(\boldsymbol{\varphi})$  which is added to the above MSE approximation. This extra term comes from the fact that we are estimating the actual non-sample counts  $\sum_i y_{dij}^r$  and not simply  $\tau_{dj}^r$  which introduces additional variability (this extra term is also found in the MSE approximation in Molina *et al.* (2007)).

Note that the multivariate linear mixed model in Baillo and Molina (2005) is different from our approximate multivariate linear mixed model. In Baillo and Molina (2005) there is only one random effect associated with each small area and we have two correlated random effects in each small area. Therefore the formulas in Baillo and Molina (2005) cannot be immediately applied to our case. Nonetheless we are still able to produce a MSE estimator but note that it cannot be guaranteed to be accurate to a known order even if the multivariate linear mixed model was not an approximation (a more detailed and rigorous proof along similar lines to pages 7–17 in Baillo and Molina (2005) would be needed for that). In any case, we show in later sections that our MSE estimator performs well. We now derive this MSE estimator.

For  $d = 1, 2, \dots, D$ , let  $\boldsymbol{\delta}_d = (\delta_{d1}, \delta_{d2})^t$  be the vector of small area totals that we are interested in predicting (note that the third total within each small area can be obtained via subtraction). Now

$$\boldsymbol{\delta}_d = \sum_{i=1}^I \mathbf{y}_{di} + \sum_{i=1}^I \boldsymbol{\mu}_{di}^r + \sum_{i=1}^I (\mathbf{y}_{di}^r - \boldsymbol{\mu}_{di}^r),$$

where  $\mathbf{y}_{di}$  and  $\mathbf{y}_{di}^r$  are respectively the vectors of sample totals (known) and non-sample totals (unknown) and

$$\boldsymbol{\mu}_{di}^r = m_{di}^r (p_{di1}, p_{di2})^t, \quad (8.1)$$

where  $m_{di}^r$  is the total non-sample which is assumed known.

Now define

$$\tilde{\boldsymbol{\delta}}_d = \sum_{i=1}^I \mathbf{y}_{di} + \sum_{i=1}^I \boldsymbol{\mu}_{di}^r = \sum_{i=1}^I \mathbf{y}_{di} + \boldsymbol{\tau}_d$$

and

$$\hat{\boldsymbol{\delta}}_d = \sum_{i=1}^I \mathbf{y}_{di} + \sum_{i=1}^I \hat{\boldsymbol{\mu}}_{di}^r = \sum_{i=1}^I \mathbf{y}_{di} + \hat{\boldsymbol{\tau}}_d, \quad (8.2)$$

where  $\hat{\boldsymbol{\mu}}_{di}^r$  is (8.1) with the  $\mathbf{u}$ 's and  $\boldsymbol{\beta}$ 's in  $p_{di1}$  and  $p_{di2}$  replaced with the PQL estimates  $\hat{\mathbf{u}}$  and  $\hat{\boldsymbol{\beta}}$ .

The estimator  $\hat{\boldsymbol{\delta}}_d$  is used to estimate  $\boldsymbol{\delta}_d$ . We will define the MSE matrix of  $\hat{\boldsymbol{\delta}}_d$  as

$$\begin{aligned} \text{MSE}(\hat{\boldsymbol{\delta}}_d) &= E\left(\left(\hat{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right)\left(\hat{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right)^t\right) \\ &= E\left(\left(\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d + \tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right)\left(\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d + \tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right)^t\right) \\ &= E\left(\left(\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d\right)\left(\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d\right)^t\right) + E\left(\left(\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right)\left(\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right)^t\right). \end{aligned} \quad (8.3)$$

The MSE matrix contains  $\text{MSE}(\hat{\boldsymbol{\delta}}_{d1})$  and  $\text{MSE}(\hat{\boldsymbol{\delta}}_{d2})$  and the cross product term

$$E\left(\left(\hat{\boldsymbol{\delta}}_{d1} - \boldsymbol{\delta}_{d1}\right)\left(\hat{\boldsymbol{\delta}}_{d2} - \boldsymbol{\delta}_{d2}\right)\right).$$

Note that the terms

$$E\left(\left(\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right)\left(\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d\right)^t\right)$$

and

$$E\left(\left(\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d\right)\left(\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right)^t\right)$$

are omitted from (8.3) because these are matrices containing all zeros. This is because given  $\mathbf{u}_d$ ,  $\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d$  and  $\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d$  are independent.

$$\text{Hence } E\left(\left(\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right)\left(\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d\right)^t\right) = E\left(E\left(\left(\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right)\middle|\mathbf{u}_d\right)E\left(\left(\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d\right)^t\middle|\mathbf{u}_d\right)\right)$$

$$\text{and note that } E\left(\left(\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right)\middle|\mathbf{u}_d\right) = \mathbf{0}_{2 \times 1}.$$

By a similar argument,

$$E\left(\left(\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d\right)\left(\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right)^t\right) = \mathbf{0}_{2 \times 2}.$$

Now

$$E\left(\left(\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d\right)\left(\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d\right)^t\right) = E\left(\left(\hat{\boldsymbol{\tau}}_d - \boldsymbol{\tau}_d\right)\left(\hat{\boldsymbol{\tau}}_d - \boldsymbol{\tau}_d\right)^t\right) = \text{MSE}(\hat{\boldsymbol{\tau}}_d)$$

and

$$\begin{aligned}
E\left((\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d)(\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d)^t\right) &= E\left(\left(\sum_{i=1}^I (\mathbf{y}_{di}^r - \boldsymbol{\mu}_{di}^r)\right)\left(\sum_{i=1}^I (\mathbf{y}_{di}^r - \boldsymbol{\mu}_{di}^r)\right)^t\right) \\
&= E\left(E\left(\left(\sum_{i=1}^I (\mathbf{y}_{di}^r - \boldsymbol{\mu}_{di}^r)\right)\left(\sum_{i=1}^I (\mathbf{y}_{di}^r - \boldsymbol{\mu}_{di}^r)\right)^t \middle| \mathbf{u}_d\right)\right) \\
&= E\left(\sum_{i=1}^I \boldsymbol{\Sigma}_{di}^r\right), \tag{8.4}
\end{aligned}$$

where

$$\boldsymbol{\Sigma}_{di}^r = m_{di}^r \begin{pmatrix} p_{di1}(1-p_{di1}) & -p_{di1}p_{di2} \\ -p_{di1}p_{di2} & p_{di2}(1-p_{di2}) \end{pmatrix}.$$

The term  $E(\sum_{i=1}^I \boldsymbol{\Sigma}_{di}^r)$  cannot be further simplified because  $\boldsymbol{\Sigma}_{di}^r$  is a non-linear function of the random effects and the expectation involves an integral with no closed form solution. To estimate this term we can use (as in Molina *et al.* (2007))

$$\mathbf{G}_4(\hat{\boldsymbol{\phi}}) = \sum_{i=1}^I \hat{\boldsymbol{\Sigma}}_{di}^r,$$

where  $\hat{\boldsymbol{\Sigma}}_{di}^r$  is  $\boldsymbol{\Sigma}_{di}^r$  with the  $\mathbf{u}$ 's and  $\boldsymbol{\beta}$ 's within replaced by the PQL estimates  $\hat{\mathbf{u}}$  and  $\hat{\boldsymbol{\beta}}$ .

Now we need to further simplify and approximate the matrix  $MSE(\hat{\boldsymbol{\tau}}_d)$ .

By definition,

$$\hat{\boldsymbol{\tau}}_d = \sum_{i=1}^I \hat{\boldsymbol{\mu}}_{di}^r$$

and each of the  $\hat{\boldsymbol{\mu}}_{di}^r$  can be written as functions of

$$\hat{\boldsymbol{\theta}}_{di} = (\hat{\theta}_{di1}, \hat{\theta}_{di2})^t$$

where

$$\hat{\boldsymbol{\theta}}_{di} = \mathbf{X}_{di} \hat{\boldsymbol{\beta}} + \mathbf{Z}_{di} \hat{\mathbf{u}}.$$

We can now calculate a first order Taylor series approximation for each of the vectors  $\hat{\boldsymbol{\mu}}_{di}^r$  about the point  $\boldsymbol{\theta}_{di}$ . This approximation is

$$\hat{\boldsymbol{\mu}}_{di}^r \approx \boldsymbol{\mu}_{di}^r + \boldsymbol{\Sigma}_{di}^r (\hat{\boldsymbol{\theta}}_{di} - \boldsymbol{\theta}_{di})$$

and hence

$$\begin{aligned}
\hat{\boldsymbol{\tau}}_d - \boldsymbol{\tau}_d &\approx \sum_{i=1}^I \boldsymbol{\Sigma}_{di}^r \hat{\boldsymbol{\theta}}_{di} - \sum_{i=1}^I \boldsymbol{\Sigma}_{di}^r \boldsymbol{\theta}_{di} \\
&= \sum_{i=1}^I \boldsymbol{\Sigma}_{di}^r (\mathbf{X}_{di} \hat{\boldsymbol{\beta}} + \mathbf{Z}_{di} \hat{\boldsymbol{u}}) - \sum_{i=1}^I \boldsymbol{\Sigma}_{di}^r (\mathbf{X}_{di} \boldsymbol{\beta} + \mathbf{Z}_{di} \boldsymbol{u}) \\
&= \hat{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d,
\end{aligned}$$

where

$$\boldsymbol{\tau}'_d = \mathbf{K}_d \boldsymbol{\beta} + \mathbf{M}_d \boldsymbol{u}$$

and

$$\hat{\boldsymbol{\tau}}'_d = \mathbf{K}_d \hat{\boldsymbol{\beta}} + \mathbf{M}_d \hat{\boldsymbol{u}},$$

where

$$\mathbf{K}_d = \sum_{i=1}^I \boldsymbol{\Sigma}_{di}^r \mathbf{X}_{di}$$

and

$$\mathbf{M}_d = \sum_{i=1}^I \boldsymbol{\Sigma}_{di}^r \mathbf{Z}_{di}.$$

Technically  $\mathbf{K}_d$  and  $\mathbf{M}_d$  are random variables since  $\boldsymbol{\Sigma}_{di}$  is dependent on  $\boldsymbol{u}_d$ . However, from this point forward we need to assume that  $\mathbf{K}_d$  and  $\mathbf{M}_d$  are constant and do not depend on  $\boldsymbol{u}_d$ . Note that this will be a reasonable assumption if the variance components are small. Now, using the fact that  $\hat{\boldsymbol{\tau}}_d - \boldsymbol{\tau}_d \approx \hat{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d$ , it follows that  $MSE(\hat{\boldsymbol{\tau}}_d) \approx MSE(\hat{\boldsymbol{\tau}}'_d)$  (Molina *et al.* (2007) also make this assumption) and therefore

$$\begin{aligned}
MSE(\hat{\boldsymbol{\tau}}_d) &\approx E\left((\hat{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)(\hat{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)^t\right) \\
&= E\left((\hat{\boldsymbol{\tau}}'_d - \tilde{\boldsymbol{\tau}}'_d + \tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)(\hat{\boldsymbol{\tau}}'_d - \tilde{\boldsymbol{\tau}}'_d + \tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)^t\right), \tag{8.5}
\end{aligned}$$

where

$$\tilde{\boldsymbol{\tau}}'_d = \mathbf{K}_d \tilde{\boldsymbol{\beta}} + \mathbf{M}_d \tilde{\boldsymbol{u}}$$

and  $\tilde{\boldsymbol{\beta}}$  and  $\tilde{\boldsymbol{u}}$  are the estimates of  $\boldsymbol{\beta}$  and  $\boldsymbol{u}$  assuming that the variance components  $\boldsymbol{\varphi}_1$ ,  $\boldsymbol{\varphi}_2$  and  $\boldsymbol{\varphi}_{12}$  are known.

As noted in Section 5, the PQL estimators can be obtained via equations (5.1) and (5.2). So we can write

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= \left(\mathbf{X}'\hat{\mathbf{V}}^{-1}\mathbf{X}\right)^{-1}\mathbf{X}'\hat{\mathbf{V}}^{-1}\boldsymbol{\xi}, \\ \hat{\mathbf{u}} &= \hat{\mathbf{W}}\mathbf{Z}'\hat{\mathbf{V}}^{-1}\left(\boldsymbol{\xi}-\mathbf{X}\hat{\boldsymbol{\beta}}\right), \\ \tilde{\boldsymbol{\beta}} &= \left(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X}\right)^{-1}\mathbf{X}'\mathbf{V}^{-1}\boldsymbol{\xi}\end{aligned}\tag{8.6}$$

and

$$\tilde{\mathbf{u}} = \mathbf{W}\mathbf{Z}'\mathbf{V}^{-1}\left(\boldsymbol{\xi}-\mathbf{X}\tilde{\boldsymbol{\beta}}\right),\tag{8.7}$$

where  $\mathbf{V}=\mathbf{Z}\mathbf{W}\mathbf{Z}'+\boldsymbol{\Sigma}^{-1}$  and  $\hat{\mathbf{V}}=\mathbf{Z}\widehat{\mathbf{W}}\mathbf{Z}'+\widehat{\boldsymbol{\Sigma}}^{-1}$ . Note that  $\widehat{\mathbf{W}}$  is  $\mathbf{W}$  but with  $\boldsymbol{\varphi}_1$ ,  $\boldsymbol{\varphi}_2$  and  $\boldsymbol{\varphi}_{12}$  within replaced with either the approximated ML or REML estimates. Also,  $\widehat{\boldsymbol{\Sigma}}^{-1}$  is just  $\boldsymbol{\Sigma}^{-1}$  but with  $\boldsymbol{\beta}$  and  $\mathbf{u}$  within replaced with the PQL estimates (which are also technically functions of  $\hat{\boldsymbol{\varphi}}_1$ ,  $\hat{\boldsymbol{\varphi}}_2$  and  $\hat{\boldsymbol{\varphi}}_{12}$ ).

Baillo and Molina (2005) apply a result which was derived by Kackar and Harville (1984) in a linear mixed model setting. We will also apply this result and assume that the estimator of the variance components is translation invariant. This assumption means that we can further simplify (8.5) to

$$\begin{aligned}MSE(\hat{\boldsymbol{\tau}}'_d) &= E\left(\left(\hat{\boldsymbol{\tau}}'_d-\tilde{\boldsymbol{\tau}}'_d\right)\left(\hat{\boldsymbol{\tau}}'_d-\tilde{\boldsymbol{\tau}}'_d\right)^t\right)+E\left(\left(\tilde{\boldsymbol{\tau}}'_d-\boldsymbol{\tau}'_d\right)\left(\tilde{\boldsymbol{\tau}}'_d-\boldsymbol{\tau}'_d\right)^t\right) \\ &= E\left(\left(\hat{\boldsymbol{\tau}}'_d-\tilde{\boldsymbol{\tau}}'_d\right)\left(\hat{\boldsymbol{\tau}}'_d-\tilde{\boldsymbol{\tau}}'_d\right)^t\right)+MSE(\tilde{\boldsymbol{\tau}}'_d).\end{aligned}\tag{8.8}$$

To approximate the first term in (8.8) we note that  $\hat{\boldsymbol{\tau}}'_d=\tilde{\boldsymbol{\tau}}'_d(\hat{\boldsymbol{\varphi}})$  and  $\tilde{\boldsymbol{\tau}}'_d=\tilde{\boldsymbol{\tau}}'_d(\boldsymbol{\varphi})$ , where  $\boldsymbol{\varphi}=(\varphi_1,\varphi_2,\varphi_{12})^t$ . This only works if we assume that  $\widehat{\boldsymbol{\Sigma}}_{di}$  is not a function of  $\hat{\boldsymbol{\varphi}}$  (technically  $\widehat{\boldsymbol{\Sigma}}_{di}$  is a function of  $\hat{\mathbf{u}}$  and  $\hat{\boldsymbol{\beta}}$  which are themselves functions of  $\hat{\boldsymbol{\varphi}}$ ). With these assumptions we can now approximate  $\hat{\boldsymbol{\tau}}'_d-\tilde{\boldsymbol{\tau}}'_d$  using a first order Taylor series expansion like in Kackar and Harville (1984). Here though we are in a bivariate setting.

First let  $\tilde{\boldsymbol{\tau}}'_d=(\tilde{\tau}'_{d1},\tilde{\tau}'_{d2})^t$ , then

$$\hat{\boldsymbol{\tau}}'_d(\hat{\boldsymbol{\varphi}})\approx\tilde{\boldsymbol{\tau}}'_d(\boldsymbol{\varphi})+\begin{pmatrix}\frac{\partial\tilde{\tau}'_{d1}}{\partial\varphi_1}&\frac{\partial\tilde{\tau}'_{d1}}{\partial\varphi_2}&\frac{\partial\tilde{\tau}'_{d1}}{\partial\varphi_{12}} \\ \frac{\partial\tilde{\tau}'_{d2}}{\partial\varphi_1}&\frac{\partial\tilde{\tau}'_{d2}}{\partial\varphi_2}&\frac{\partial\tilde{\tau}'_{d2}}{\partial\varphi_{12}}\end{pmatrix}(\hat{\boldsymbol{\varphi}}-\boldsymbol{\varphi}).$$

This then implies that

$$E\left((\hat{\boldsymbol{\tau}}'_d - \tilde{\boldsymbol{\tau}}'_d)(\hat{\boldsymbol{\tau}}'_d - \tilde{\boldsymbol{\tau}}'_d)^t\right) \approx E\begin{pmatrix} g_{11} & g_{12} \\ g_{12} & g_{22} \end{pmatrix},$$

where for  $j = 1, 2$  and  $k = 1, 2$ ,

$$\begin{aligned} g_{jk} &= \frac{\partial \tilde{\tau}'_{dj}}{\partial \varphi_1} \frac{\partial \tilde{\tau}'_{dk}}{\partial \varphi_1} (\hat{\varphi}_1 - \varphi_1)^2 + \frac{\partial \tilde{\tau}'_{dj}}{\partial \varphi_2} \frac{\partial \tilde{\tau}'_{dk}}{\partial \varphi_2} (\hat{\varphi}_2 - \varphi_2)^2 + \frac{\partial \tilde{\tau}'_{dj}}{\partial \varphi_{12}} \frac{\partial \tilde{\tau}'_{dk}}{\partial \varphi_{12}} (\hat{\varphi}_{12} - \varphi_{12})^2 \\ &+ \left( \frac{\partial \tilde{\tau}'_{dj}}{\partial \varphi_1} \frac{\partial \tilde{\tau}'_{dk}}{\partial \varphi_2} + \frac{\partial \tilde{\tau}'_{dj}}{\partial \varphi_2} \frac{\partial \tilde{\tau}'_{dk}}{\partial \varphi_1} \right) (\hat{\varphi}_1 - \varphi_1)(\hat{\varphi}_2 - \varphi_2) \\ &+ \left( \frac{\partial \tilde{\tau}'_{dj}}{\partial \varphi_1} \frac{\partial \tilde{\tau}'_{dk}}{\partial \varphi_{12}} + \frac{\partial \tilde{\tau}'_{dj}}{\partial \varphi_{12}} \frac{\partial \tilde{\tau}'_{dk}}{\partial \varphi_1} \right) (\hat{\varphi}_1 - \varphi_1)(\hat{\varphi}_{12} - \varphi_{12}) \\ &+ \left( \frac{\partial \tilde{\tau}'_{dj}}{\partial \varphi_2} \frac{\partial \tilde{\tau}'_{dk}}{\partial \varphi_{12}} + \frac{\partial \tilde{\tau}'_{dj}}{\partial \varphi_{12}} \frac{\partial \tilde{\tau}'_{dk}}{\partial \varphi_2} \right) (\hat{\varphi}_2 - \varphi_2)(\hat{\varphi}_{12} - \varphi_{12}). \end{aligned}$$

Now let for  $j = 1, 2$ ,

$$\frac{\partial \tilde{\tau}'_{dj}}{\partial \boldsymbol{\varphi}} = \left( \frac{\partial \tilde{\tau}'_{dj}}{\partial \varphi_1}, \frac{\partial \tilde{\tau}'_{dj}}{\partial \varphi_2}, \frac{\partial \tilde{\tau}'_{dj}}{\partial \varphi_{12}} \right)^t,$$

then for  $j = 1, 2$  and  $k = 1, 2$ ,

$$g_{jk} = Tr \left[ \left( \frac{\partial \tilde{\tau}'_{dj}}{\partial \boldsymbol{\varphi}} \right) \left( \frac{\partial \tilde{\tau}'_{dk}}{\partial \boldsymbol{\varphi}} \right)^t (\hat{\boldsymbol{\varphi}} - \boldsymbol{\varphi})(\hat{\boldsymbol{\varphi}} - \boldsymbol{\varphi})^t \right].$$

Now assuming that the elements of each of  $\left( \frac{\partial \tilde{\tau}'_{dj}}{\partial \boldsymbol{\varphi}} \right) \left( \frac{\partial \tilde{\tau}'_{dk}}{\partial \boldsymbol{\varphi}} \right)^t$ , for each combination of  $j$  and  $k$ , are approximately independent of the elements of  $(\hat{\boldsymbol{\varphi}} - \boldsymbol{\varphi})(\hat{\boldsymbol{\varphi}} - \boldsymbol{\varphi})^t$ , then

$$E\left((\hat{\boldsymbol{\tau}}'_d - \tilde{\boldsymbol{\tau}}'_d)(\hat{\boldsymbol{\tau}}'_d - \tilde{\boldsymbol{\tau}}'_d)^t\right) \approx \begin{pmatrix} g_{11}^* & g_{12}^* \\ g_{12}^* & g_{22}^* \end{pmatrix}, \quad (8.9)$$

where for  $j = 1, 2$  and  $k = 1, 2$ ,

$$g_{jk}^* = Tr \left[ E \left[ \left( \frac{\partial \tilde{\tau}'_{dj}}{\partial \boldsymbol{\varphi}} \right) \left( \frac{\partial \tilde{\tau}'_{dk}}{\partial \boldsymbol{\varphi}} \right)^t \right] E \left[ (\hat{\boldsymbol{\varphi}} - \boldsymbol{\varphi})(\hat{\boldsymbol{\varphi}} - \boldsymbol{\varphi})^t \right] \right].$$

This idea is motivated from Kackar and Harville (1984) and is also applied in Prasad and Rao (1990) and Baillo and Molina (2005).

The term 
$$E\left((\hat{\boldsymbol{\varphi}} - \boldsymbol{\varphi})(\hat{\boldsymbol{\varphi}} - \boldsymbol{\varphi})^t\right)$$

can be approximated by using the inverse of the approximated observed Fisher information matrix. That is, if approximated ML is being used for the variance components then we can use the final value of  $-\mathbf{J}_{\boldsymbol{\varphi}}^{-1}$  from Section 5. Or if approximated REML is being used for the variance components then we can use the final value of  $-\mathbf{J}_{\boldsymbol{\varphi}}^{-1}$  from Section 6.

As for the other terms 
$$E\left(\left(\frac{\partial \tilde{\tau}'_{dj}}{\partial \boldsymbol{\varphi}}\right)\left(\frac{\partial \tilde{\tau}'_{dk}}{\partial \boldsymbol{\varphi}}\right)^t\right)$$

for  $j = 1, 2$  and  $k = 1, 2$ , these can be simplified further by applying some ideas from Prasad and Rao (1990) and Baillo and Molina (2005) with some adjustments. We do this now.

Note that  $\tilde{\boldsymbol{\tau}}'_d$  can be written as follows

$$\tilde{\boldsymbol{\tau}}'_d = \begin{pmatrix} \tilde{\tau}'_{d1} \\ \tilde{\tau}'_{d2} \end{pmatrix} = \begin{pmatrix} \mathbf{K}_{d1}\tilde{\boldsymbol{\beta}} + \mathbf{M}_{d1}\tilde{\boldsymbol{u}} \\ \mathbf{K}_{d2}\tilde{\boldsymbol{\beta}} + \mathbf{M}_{d2}\tilde{\boldsymbol{u}} \end{pmatrix},$$

where  $\mathbf{K}_{d1}$  and  $\mathbf{K}_{d2}$  are respectively the first and second rows of the matrix  $\mathbf{K}_d$ . Similarly,  $\mathbf{M}_{d1}$  and  $\mathbf{M}_{d2}$  are respectively the first and second rows of the matrix  $\mathbf{M}_d$ .

For  $j = 1, 2$  we have

$$\begin{aligned} \tilde{\tau}'_{dj} &= \mathbf{K}_{dj}\tilde{\boldsymbol{\beta}} + \mathbf{M}_{dj}\tilde{\boldsymbol{u}} \\ &= \mathbf{K}_{dj}\left(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X}\right)^{-1}\mathbf{X}'\mathbf{V}^{-1}\boldsymbol{\xi} + \mathbf{M}_{dj}\mathbf{WZ}'\mathbf{V}^{-1}\left(\boldsymbol{\xi} - \mathbf{X}\tilde{\boldsymbol{\beta}}\right) \\ &= \mathbf{K}_{dj}\left(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X}\right)^{-1}\mathbf{X}'\mathbf{V}^{-1}\boldsymbol{\xi} + \mathbf{M}_{dj}\mathbf{WZ}'\mathbf{V}^{-1}\left(\boldsymbol{\xi} - \mathbf{X}\left(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X}\right)^{-1}\mathbf{X}'\mathbf{V}^{-1}\boldsymbol{\xi}\right). \end{aligned}$$

Now substitute  $\boldsymbol{\xi} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e}$  into the above equation and after some algebra and some simplifications we obtain

$$\begin{aligned} \tilde{\tau}'_{dj} &= \mathbf{K}_{dj}\boldsymbol{\beta} \\ &+ \left(\mathbf{K}_{dj}\left(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X}\right)^{-1}\mathbf{X}'\mathbf{V}^{-1} + \mathbf{M}_{dj}\mathbf{WZ}'\left(\mathbf{V}^{-1} - \mathbf{V}^{-1}\mathbf{X}\left(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X}\right)^{-1}\mathbf{X}'\mathbf{V}^{-1}\right)\right)(\mathbf{Z}\mathbf{u} + \mathbf{e}). \end{aligned}$$

Now for  $a = 1, 2$  and  $12$  and  $j = 1, 2$  we have

$$\begin{aligned} \frac{\partial \tilde{\tau}'_{dj}}{\partial \varphi_a} &= \frac{\partial}{\partial \varphi_a} \left( \begin{array}{c} \mathbf{K}_{dj} (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} + \\ \mathbf{M}_{dj} \mathbf{WZ}' \left( \mathbf{V}^{-1} - \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} \right) \end{array} \right) (\mathbf{Zu} + \mathbf{e}) \\ &\approx \frac{\partial (\mathbf{M}_{dj} \mathbf{WZ}' \mathbf{V}^{-1})}{\partial \varphi_a} (\mathbf{Zu} + \mathbf{e}). \end{aligned} \quad (8.10)$$

The justification for the above approximation can be roughly explained as follows. The terms that we are ignoring contain the term  $(\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1}$ . This term is expected to be small and negligible if the number of small areas  $D$  is large (which is often the case in practise) compared to the number of parameters in  $\boldsymbol{\beta}$ . This is because  $(\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1}$  is the approximate variance of  $\hat{\boldsymbol{\beta}}$  which gets smaller as the sample size increases.

Using the approximation (8.10) we can now write for  $j = 1, 2$  and  $k = 1, 2$ ,

$$\begin{aligned} &E \left( \left( \frac{\partial \tilde{\tau}'_{dj}}{\partial \boldsymbol{\varphi}} \right) \left( \frac{\partial \tilde{\tau}'_{dj}}{\partial \boldsymbol{\varphi}} \right)' \right) \\ &\approx \left( \frac{\partial (\mathbf{M}_{dj} \mathbf{WZ}' \mathbf{V}^{-1})}{\partial \boldsymbol{\varphi}} \right) E \left( (\mathbf{Zu} + \mathbf{e})(\mathbf{Zu} + \mathbf{e})' \right) \left( \frac{\partial (\mathbf{M}_{dk} \mathbf{WZ}' \mathbf{V}^{-1})}{\partial \boldsymbol{\varphi}} \right)' \\ &= \left( \frac{\partial (\mathbf{M}_{dj} \mathbf{WZ}' \mathbf{V}^{-1})}{\partial \boldsymbol{\varphi}} \right) \mathbf{V} \left( \frac{\partial (\mathbf{M}_{dk} \mathbf{WZ}' \mathbf{V}^{-1})}{\partial \boldsymbol{\varphi}} \right)', \end{aligned} \quad (8.11)$$

where

$$\frac{\partial (\mathbf{M}_{dj} \mathbf{WZ}' \mathbf{V}^{-1})}{\partial \boldsymbol{\varphi}} = \left( \left( \frac{\partial (\mathbf{M}_{dj} \mathbf{WZ}' \mathbf{V}^{-1})}{\partial \varphi_1} \right)', \left( \frac{\partial (\mathbf{M}_{dj} \mathbf{WZ}' \mathbf{V}^{-1})}{\partial \varphi_2} \right)', \left( \frac{\partial (\mathbf{M}_{dj} \mathbf{WZ}' \mathbf{V}^{-1})}{\partial \varphi_{12}} \right)' \right)'$$

To approximate each of the matrices given by (8.11) for  $j = 1, 2$  and  $k = 1, 2$ , we simply calculate the derivatives and then replace the unknown parameters by their estimated values. To calculate the elements within the matrices we can use the following simplifications which we obtained after some algebra.



For  $j = 1, 2$  and  $k = 1, 2$  and  $a = 1, 2, 12$  and  $b = 1, 2, 12$

$$\begin{aligned}
& \left( \frac{\partial \left( \mathbf{M}_{dj} \mathbf{W} \mathbf{Z}' \mathbf{V}^{-1} \right)}{\partial \varphi_a} \right) \mathbf{V} \left( \frac{\partial \left( \mathbf{M}_{dk} \mathbf{W} \mathbf{Z}' \mathbf{V}^{-1} \right)}{\partial \varphi_b} \right)^t = \mathbf{M}_{dj} \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_b} \mathbf{M}_{dk}^t \\
& - \mathbf{M}_{dj} \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_b} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{Z} \mathbf{W} \mathbf{M}_{dk}^t - \mathbf{M}_{dj} \mathbf{W} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_b} \mathbf{M}_{dk}^t \\
& + \mathbf{M}_{dj} \mathbf{W} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_b} \mathbf{Z}' \mathbf{V}^{-1} \mathbf{Z} \mathbf{W} \mathbf{M}_{dk}^t.
\end{aligned} \tag{8.12}$$

So we can now fully approximate  $E \left( (\hat{\boldsymbol{\tau}}'_d - \tilde{\boldsymbol{\tau}}'_d) (\hat{\boldsymbol{\tau}}'_d - \tilde{\boldsymbol{\tau}}'_d)^t \right)$ .

Next we need to approximate the matrix  $MSE(\hat{\boldsymbol{\tau}}'_d)$ . To do this we follow almost directly the approach taken in Molina *et al.* (2007) and after some algebra we obtain terms very similar to those obtained in that paper.

By definition

$$\tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d = \mathbf{K}_d \tilde{\boldsymbol{\beta}} + \mathbf{M}_d \tilde{\mathbf{u}} - \mathbf{K}_d \boldsymbol{\beta} - \mathbf{M}_d \mathbf{u}. \tag{8.13}$$

Let 
$$\boldsymbol{\Pi} = \mathbf{V}^{-1} - \mathbf{V}^{-1} \mathbf{X} \left( \mathbf{X}' \mathbf{V}^{-1} \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{V}^{-1}$$

and 
$$\mathbf{P} = \left( \mathbf{X}' \mathbf{V}^{-1} \mathbf{X} \right)^{-1}$$

and substitute (8.6) and (8.7) and  $\boldsymbol{\xi} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e}$  into (8.13). We then obtain

$$\tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d = \left( \mathbf{K}_d \mathbf{P} \mathbf{X}' \mathbf{V}^{-1} + \mathbf{M}_d \mathbf{W} \mathbf{Z}' \boldsymbol{\Pi} \right) (\mathbf{Z}\mathbf{u} + \mathbf{e}) - \mathbf{M}_d \mathbf{u}$$

and hence

$$\begin{aligned}
& (\tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d) (\tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)^t \\
& = \left( \mathbf{K}_d \mathbf{P} \mathbf{X}' \mathbf{V}^{-1} + \mathbf{M}_d \mathbf{W} \mathbf{Z}' \boldsymbol{\Pi} \right) (\mathbf{Z}\mathbf{u} + \mathbf{e}) (\mathbf{Z}\mathbf{u} + \mathbf{e})^t \left( \mathbf{V}^{-1} \mathbf{X} \mathbf{P} \mathbf{K}_d^t + \boldsymbol{\Pi} \mathbf{Z} \mathbf{W} \mathbf{M}_d^t \right) \\
& + \mathbf{M}_d \mathbf{u} \mathbf{u}^t \mathbf{M}_d^t - \mathbf{M}_d \mathbf{u} (\mathbf{Z}\mathbf{u} + \mathbf{e})^t \left( \mathbf{V}^{-1} \mathbf{X} \mathbf{P} \mathbf{K}_d^t + \boldsymbol{\Pi} \mathbf{Z} \mathbf{W} \mathbf{M}_d^t \right) \\
& - \left( \mathbf{K}_d \mathbf{P} \mathbf{X}' \mathbf{V}^{-1} + \mathbf{M}_d \mathbf{W} \mathbf{Z}' \boldsymbol{\Pi} \right) (\mathbf{Z}\mathbf{u} + \mathbf{e}) \mathbf{u}^t \mathbf{M}_d^t.
\end{aligned} \tag{8.14}$$

Now 
$$E\left((\mathbf{Z}\mathbf{u} + \mathbf{e})(\mathbf{Z}\mathbf{u} + \mathbf{e})^t\right) = \mathbf{V},$$

$$E(\mathbf{u}\mathbf{u}^t) = \mathbf{W},$$

$$E(\mathbf{u}\mathbf{e}^t) = \mathbf{0}$$

and so after the above substitutions and with some further simplifications we have

$$\begin{aligned} E\left((\tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)(\tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)^t\right) &= \mathbf{K}_d \mathbf{P} \mathbf{K}_d^t - \mathbf{K}_d \mathbf{P} \mathbf{X}^t \mathbf{V}^{-1} \mathbf{Z} \mathbf{W} \mathbf{M}_d^t \\ &- \mathbf{M}_d \mathbf{W} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{X} \mathbf{P} \mathbf{K}_d^t + \mathbf{M}_d \mathbf{W} \mathbf{M}_d^t - \mathbf{M}_d \mathbf{W} \mathbf{Z}^t \boldsymbol{\Pi} \mathbf{Z} \mathbf{W} \mathbf{M}_d^t. \end{aligned} \quad (8.15)$$

Now let 
$$\mathbf{T} = \left(\mathbf{W}^{-1} + \mathbf{Z}^t \boldsymbol{\Sigma} \mathbf{Z}\right)^{-1} \quad (8.16)$$

and we will need the following identity from Henderson and Searle (1981) (for example)

$$\mathbf{V}^{-1} = \boldsymbol{\Sigma} - \boldsymbol{\Sigma} \mathbf{Z} \mathbf{T} \mathbf{Z}^t \boldsymbol{\Sigma}. \quad (8.17)$$

After some algebra it can be proved that

$$\mathbf{V}^{-1} \mathbf{Z} \mathbf{W} = \boldsymbol{\Sigma} \mathbf{Z} \mathbf{T} \quad (8.18)$$

and 
$$\mathbf{W} \mathbf{Z}^t \mathbf{V}^{-1} = \mathbf{T} \mathbf{Z}^t \boldsymbol{\Sigma}. \quad (8.19)$$

Using (8.16)–(8.19) we can further simplify (8.15) to

$$\begin{aligned} E\left((\tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)(\tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)^t\right) &= \mathbf{M}_d \mathbf{T} \mathbf{M}_d^t + \mathbf{K}_d \mathbf{P} \mathbf{K}_d^t + \mathbf{M}_d \mathbf{T} \mathbf{Z}^t \boldsymbol{\Sigma} \mathbf{X} \mathbf{P} \mathbf{X}^t \boldsymbol{\Sigma} \mathbf{T} \mathbf{M}_d^t \\ &- \mathbf{K}_d \mathbf{P} \mathbf{X}^t \boldsymbol{\Sigma} \mathbf{T} \mathbf{M}_d^t - \mathbf{M}_d \mathbf{T} \mathbf{Z}^t \boldsymbol{\Sigma} \mathbf{X} \mathbf{P} \mathbf{K}_d^t \end{aligned}$$

which on collecting terms is

$$E\left((\tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)(\tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)^t\right) = \mathbf{M}_d \mathbf{T} \mathbf{M}_d^t + \left(\mathbf{K}_d - \mathbf{M}_d \mathbf{T} \mathbf{Z}^t \boldsymbol{\Sigma} \mathbf{X}\right) \mathbf{P} \left(\mathbf{K}_d - \mathbf{M}_d \mathbf{T} \mathbf{Z}^t \boldsymbol{\Sigma} \mathbf{X}\right)^t.$$

To approximate the above equation we simply substitute in the estimates of  $\boldsymbol{\beta}$ ,  $\mathbf{u}$  and  $\boldsymbol{\varphi}$ .

We now have a formula to approximate the matrices  $MSE(\hat{\delta}_d)$ . This approximation is

$$\widehat{MSE}(\hat{\delta}_d) = \mathbf{G}_1(\hat{\boldsymbol{\varphi}}) + \mathbf{G}_2(\hat{\boldsymbol{\varphi}}) + 2\mathbf{G}_3(\hat{\boldsymbol{\varphi}}) + \mathbf{G}_4(\hat{\boldsymbol{\varphi}}), \quad (8.20)$$

where

$$\mathbf{G}_1(\hat{\boldsymbol{\varphi}}) = \hat{\mathbf{M}}_d \hat{\mathbf{T}} \hat{\mathbf{M}}_d^t,$$

$$\mathbf{G}_2(\hat{\boldsymbol{\varphi}}) = \left( \hat{\mathbf{K}}_d - \hat{\mathbf{M}}_d \hat{\mathbf{T}} \hat{\mathbf{Z}}^t \hat{\boldsymbol{\Sigma}} \mathbf{X} \right) \hat{\mathbf{P}} \left( \hat{\mathbf{K}}_d - \hat{\mathbf{M}}_d \hat{\mathbf{T}} \hat{\mathbf{Z}}^t \hat{\boldsymbol{\Sigma}} \mathbf{X} \right)^t,$$

$$\mathbf{G}_4(\hat{\boldsymbol{\varphi}}) = \sum_{i=1}^I \hat{\boldsymbol{\Sigma}}_{di}^r$$

and  $\mathbf{G}_3(\hat{\boldsymbol{\varphi}})$  is given by (8.9) with estimates substituted in and with some further simplifications such as (8.12). Note that  $\mathbf{G}_3(\hat{\boldsymbol{\varphi}})$  is multiplied by two in the above approximation. This is because as in Molina *et al.* (2007),  $E(\mathbf{G}_1(\hat{\boldsymbol{\varphi}})) \approx \mathbf{G}_1(\boldsymbol{\varphi}) - \mathbf{G}_3(\boldsymbol{\varphi})$  and we therefore multiply  $\mathbf{G}_3(\hat{\boldsymbol{\varphi}})$  by two to correct for the bias. The proof that  $E(\mathbf{G}_1(\hat{\boldsymbol{\varphi}})) \approx \mathbf{G}_1(\boldsymbol{\varphi}) - \mathbf{G}_3(\boldsymbol{\varphi})$  is given in the Appendix.

## 9. OUT-OF-SAMPLE SMALL AREAS

All previous sections implicitly assumed that for  $d = 1, 2, \dots, D$ , each small area had some sample in at least one of its age/sex classes. Now suppose that as before we have  $d = 1, 2, \dots, D$  in-sample small areas, but we now have an additional set of  $R$  small areas with no sample. We can still produce estimates and MSEs for these out-of-sample small areas, but a slightly different methodology needs to be applied.

In total we have  $d = 1, 2, \dots, D, D+1, \dots, D+R$  small areas that we wish to produce estimates and MSEs for. For the  $d = 1, 2, \dots, D$  small areas in sample we simply follow the methodology in the previous sections (discard the out-of-sample small areas for estimation of parameters, small area totals and MSEs). The rest of this section discusses estimation for the out-of-sample small areas.

Given our model, the best prediction we have of the out-of-sample small area random effects  $(u_{d1}, u_{d2})^t$  for  $d = D+1, D+2, \dots, D+R$  is  $(0, 0)^t$ . This is because in our model we are assuming that the random effects are independent between small areas and knowing what the observed in-sample data are does not give us any additional information about the out-of-sample  $\mathbf{u}$ 's. If we had a model with spatially correlated random effects, then we could adapt the method outlined in Saei and Chambers (2005) to produce estimates and MSEs. But of course, our model does not have spatially correlated random effects.

We therefore need to resort to producing synthetic estimates for the out-of-sample small areas. These are for  $d = D+1, D+2, \dots, D+R$ ,

$$\hat{\delta}_d = \sum_{i=1}^I \hat{\mu}_{di}^{*r} = \hat{\tau}_d,$$

where

$$\hat{\mu}_{di}^{*r} = m_{di}^r \left( \frac{e^{\mathbf{x}_{di1}^t \hat{\beta}_1}}{1 + e^{\mathbf{x}_{di1}^t \hat{\beta}_1} + e^{\mathbf{x}_{di2}^t \hat{\beta}_2}} \right) = m_{di}^r \begin{pmatrix} p_{di1}^{syn} \\ p_{di2}^{syn} \end{pmatrix},$$

where the estimates  $\hat{\beta} = (\hat{\beta}_1^t, \hat{\beta}_2^t)^t$  are the PQL estimates based on the in-sample data  $d = 1, 2, \dots, D$ .

Most of the derivation for an approximation to the MSE matrices for the out-of-sample small area totals proceeds in a similar way as in the previous section.

First define

$$\boldsymbol{\mu}_{di}^{*r} = m_{di}^r \begin{pmatrix} \frac{e^{\mathbf{x}_{di1}^t \boldsymbol{\beta}_1}}{1 + e^{\mathbf{x}_{di1}^t \boldsymbol{\beta}_1} + e^{\mathbf{x}_{di2}^t \boldsymbol{\beta}_2}} \\ \frac{e^{\mathbf{x}_{di2}^t \boldsymbol{\beta}_2}}{1 + e^{\mathbf{x}_{di1}^t \boldsymbol{\beta}_1} + e^{\mathbf{x}_{di2}^t \boldsymbol{\beta}_2}} \end{pmatrix} = m_{di}^r \begin{pmatrix} p_{di1}^* \\ p_{di2}^* \end{pmatrix},$$

and let

$$\boldsymbol{\delta}_d = \sum_{i=1}^I \boldsymbol{\mu}_{di}^{*r} + \sum_{i=1}^I \left( \mathbf{y}_{di}^r - \boldsymbol{\mu}_{di}^r + \boldsymbol{\mu}_{di}^r - \boldsymbol{\mu}_{di}^{*r} \right)$$

denote the real totals we are trying to estimate.

Now define

$$\tilde{\boldsymbol{\delta}}_d = \sum_{i=1}^I \boldsymbol{\mu}_{di}^{*r} = \boldsymbol{\tau}_d.$$

The MSE matrix of  $\hat{\boldsymbol{\delta}}_d$  for  $d = D+1, D+2, \dots, D+R$  can now be written as

$$\begin{aligned} \text{MSE}(\hat{\boldsymbol{\delta}}_d) &= E\left( (\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d)(\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d)^t \right) + E\left( (\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d)(\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d)^t \right) \\ &\quad + E\left( (\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d)(\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d)^t \right) + E\left( (\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d)(\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d)^t \right). \end{aligned} \quad (9.1)$$

We now go about further simplifying the terms in (9.1).

Firstly we note that

$$\begin{aligned} &E\left( \left( \sum_{i=1}^I \mathbf{y}_{di}^r - \sum_{i=1}^I \boldsymbol{\mu}_{di}^r \right) \left( \sum_{i=1}^I \boldsymbol{\mu}_{di}^r - \sum_{i=1}^I \boldsymbol{\mu}_{di}^{*r} \right)^t \right) \\ &= E\left( \left( E\left( \sum_{i=1}^I \mathbf{y}_{di}^r \mid \mathbf{u}_d \right) - \sum_{i=1}^I \boldsymbol{\mu}_{di}^r \right) \left( \sum_{i=1}^I \boldsymbol{\mu}_{di}^r - \sum_{i=1}^I \boldsymbol{\mu}_{di}^{*r} \right)^t \right) = \mathbf{0} \end{aligned}$$

and therefore

$$\begin{aligned} E\left( (\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d)(\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d)^t \right) &= E\left( \left( \sum_{i=1}^I \mathbf{y}_{di}^r - \sum_{i=1}^I \boldsymbol{\mu}_{di}^r \right) \left( \sum_{i=1}^I \mathbf{y}_{di}^r - \sum_{i=1}^I \boldsymbol{\mu}_{di}^r \right)^t \right) \\ &\quad + E\left( \left( \sum_{i=1}^I \boldsymbol{\mu}_{di}^r - \sum_{i=1}^I \boldsymbol{\mu}_{di}^{*r} \right) \left( \sum_{i=1}^I \boldsymbol{\mu}_{di}^r - \sum_{i=1}^I \boldsymbol{\mu}_{di}^{*r} \right)^t \right). \end{aligned} \quad (9.2)$$

The first term in the above equation is simply (8.4) and to approximate this term we can use (assuming the variance components are small)

$$\sum_{i=1}^I \hat{\Sigma}_{di}^{*r},$$

where

$$\hat{\Sigma}_{di}^{*r} = m_{di}^r \begin{pmatrix} p_{di1}^{syn} (1 - p_{di1}^{syn}) & -p_{di1}^{syn} p_{di2}^{syn} \\ -p_{di1}^{syn} p_{di2}^{syn} & p_{di2}^{syn} (1 - p_{di2}^{syn}) \end{pmatrix}. \quad (9.3)$$

The second term in (9.2) can be approximated by taking a first order Taylor series expansion. Let  $\boldsymbol{\theta} = \mathbf{X}_{di} \boldsymbol{\beta} + \mathbf{u}_d$  and  $\boldsymbol{\theta}^* = \mathbf{X}_{di} \boldsymbol{\beta}$ . Then

$$\begin{aligned} \boldsymbol{\mu}_{di}^r &\approx \boldsymbol{\mu}_{di}^{*r} + \boldsymbol{\Sigma}_{di}^{*r} (\boldsymbol{\theta} - \boldsymbol{\theta}^*), \\ &= \boldsymbol{\mu}_{di}^{*r} + \boldsymbol{\Sigma}_{di}^{*r} \mathbf{u}_d \end{aligned} \quad (9.4)$$

where  $\boldsymbol{\Sigma}_{di}^{*r}$  has the same form as (9.3) but with  $p_{di1}^{syn}$  and  $p_{di2}^{syn}$  replaced respectively with  $p_{di1}^*$  and  $p_{di2}^*$ . Hence,

$$\begin{aligned} E \left( \left( \sum_{i=1}^I \boldsymbol{\mu}_{di}^r - \sum_{i=1}^I \boldsymbol{\mu}_{di}^{*r} \right) \left( \sum_{i=1}^I \boldsymbol{\mu}_{di}^r - \sum_{i=1}^I \boldsymbol{\mu}_{di}^{*r} \right)^t \right) &\approx \mathbf{M}_d^* E(\mathbf{u}_d \mathbf{u}_d^t) \mathbf{M}_d^{*t} \\ &= \mathbf{M}_d^* \mathbf{W}_d \mathbf{M}_d^{*t}, \end{aligned} \quad (9.5)$$

where

$$\mathbf{M}_d^* = \sum_{i=1}^I \boldsymbol{\Sigma}_{di}^{*r}.$$

Note that (9.5) can be approximated with  $\hat{\mathbf{M}}_d^* \hat{\mathbf{W}}_d \hat{\mathbf{M}}_d^{*t}$ ,

where

$$\hat{\mathbf{M}}_d^* = \sum_{i=1}^I \hat{\boldsymbol{\Sigma}}_{di}^{*r}$$

and

$$\hat{\mathbf{W}}_d = \begin{pmatrix} \hat{\phi}_1 & \hat{\phi}_{12} \\ \hat{\phi}_{12} & \hat{\phi}_2 \end{pmatrix}.$$

Now we need to approximate

$$E \left( \left( \hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d \right) \left( \hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d \right)^t \right) = \text{MSE}(\hat{\boldsymbol{\tau}}_d).$$

Let  $\hat{\boldsymbol{\theta}}_{di}^* = \mathbf{X}_{di}\hat{\boldsymbol{\beta}}$  and  $\boldsymbol{\theta}_{di}^* = \mathbf{X}_{di}\boldsymbol{\beta}$ . By a first order Taylor series expansion we have

$$\hat{\boldsymbol{\mu}}_{di}^{*r} \approx \boldsymbol{\mu}_{di}^{*r} + \boldsymbol{\Sigma}_{di}^{*r} \left( \hat{\boldsymbol{\theta}}_{di}^* - \boldsymbol{\theta}_{di}^* \right) \quad (9.6)$$

and hence

$$\hat{\boldsymbol{\tau}}_d - \boldsymbol{\tau}_d \approx \sum_{i=1}^I \boldsymbol{\Sigma}_{di}^{*r} \left( \mathbf{X}_{di}\hat{\boldsymbol{\beta}} - \mathbf{X}_{di}\boldsymbol{\beta} \right). \quad (9.7)$$

Now let  $\boldsymbol{\tau}'_d = \mathbf{K}_d^* \boldsymbol{\beta}$  and  $\hat{\boldsymbol{\tau}}'_d = \mathbf{K}_d^* \hat{\boldsymbol{\beta}}$ , where  $\mathbf{K}_d^* = \sum_{i=1}^I \boldsymbol{\Sigma}_{di}^{*r} \mathbf{X}_{di}$ . Now,

$$\begin{aligned} \text{MSE}(\hat{\boldsymbol{\tau}}_d) &\approx \text{MSE}(\hat{\boldsymbol{\tau}}'_d) \\ &= E\left( (\hat{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)(\hat{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)^t \right) \\ &= E\left( (\hat{\boldsymbol{\tau}}'_d - \tilde{\boldsymbol{\tau}}'_d + \tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)(\hat{\boldsymbol{\tau}}'_d - \tilde{\boldsymbol{\tau}}'_d + \tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)^t \right) \\ &= E\left( (\tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)(\tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)^t \right) + E\left( (\hat{\boldsymbol{\tau}}'_d - \tilde{\boldsymbol{\tau}}'_d)(\hat{\boldsymbol{\tau}}'_d - \tilde{\boldsymbol{\tau}}'_d)^t \right) \end{aligned} \quad (9.8)$$

where

$$\tilde{\boldsymbol{\tau}}'_d = \mathbf{K}_d^* \tilde{\boldsymbol{\beta}}$$

and to get to the last line we assume as in the previous section that the estimator of the variance components is translation invariant.

The second term in (9.8) is

$$E\left( (\hat{\boldsymbol{\tau}}'_d - \tilde{\boldsymbol{\tau}}'_d)(\hat{\boldsymbol{\tau}}'_d - \tilde{\boldsymbol{\tau}}'_d)^t \right) = E\left( \left( \mathbf{K}_d^* \hat{\boldsymbol{\beta}} - \mathbf{K}_d^* \tilde{\boldsymbol{\beta}} \right) \left( \mathbf{K}_d^* \hat{\boldsymbol{\beta}} - \mathbf{K}_d^* \tilde{\boldsymbol{\beta}} \right)^t \right) \approx \mathbf{0}.$$

To see why the above term is approximately the  $\mathbf{0}$  matrix, see the argument just after (8.10) and note that we do not have an  $\mathbf{M}_d$  type component here.

The first term in (9.8) is

$$E\left( (\tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)(\tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)^t \right) = E\left( \left( \mathbf{K}_d^* \tilde{\boldsymbol{\beta}} - \mathbf{K}_d^* \boldsymbol{\beta} \right) \left( \mathbf{K}_d^* \tilde{\boldsymbol{\beta}} - \mathbf{K}_d^* \boldsymbol{\beta} \right)^t \right),$$

where

$$\tilde{\boldsymbol{\beta}} = \left( \mathbf{X}^t \mathbf{V}^{-1} \mathbf{X} \right)^{-1} \mathbf{X}^t \mathbf{V}^{-1} \boldsymbol{\xi}$$

and  $\mathbf{X}$  and  $\mathbf{V}$  are as defined in the previous sections and are therefore based on the in-sample areas only. Again we need to assume that  $\boldsymbol{\Sigma}^{-1}$  in  $\mathbf{V}$  does not depend on the in-sample random effects  $\boldsymbol{u}$  to do the following.

After some simplifications we have

$$\mathbf{K}_d^* \tilde{\boldsymbol{\beta}} - \mathbf{K}_d^* \boldsymbol{\beta} = \mathbf{K}_d^* \left( \mathbf{X}' \mathbf{V}^{-1} \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{V}^{-1} (\mathbf{Z} \mathbf{u} + \mathbf{e}),$$

where  $\mathbf{Z}$ ,  $\mathbf{u}$  and  $\mathbf{e}$  are defined using the in-sample areas.

Hence we have

$$\begin{aligned} & E \left( (\tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d) (\tilde{\boldsymbol{\tau}}'_d - \boldsymbol{\tau}'_d)^t \right) \\ &= \mathbf{K}_d^* \left( \mathbf{X}' \mathbf{V}^{-1} \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{V}^{-1} E \left( (\mathbf{Z} \mathbf{u} + \mathbf{e}) (\mathbf{Z} \mathbf{u} + \mathbf{e})^t \right) \mathbf{V}^{-1} \mathbf{X} \left( \mathbf{X}' \mathbf{V}^{-1} \mathbf{X} \right)^{-1} \mathbf{K}_d^{*t} \\ &= \mathbf{K}_d^* \left( \mathbf{X}' \mathbf{V}^{-1} \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{V} \mathbf{V}^{-1} \mathbf{X} \left( \mathbf{X}' \mathbf{V}^{-1} \mathbf{X} \right)^{-1} \mathbf{K}_d^{*t} \\ &= \mathbf{K}_d^* \left( \mathbf{X}' \mathbf{V}^{-1} \mathbf{X} \right)^{-1} \mathbf{K}_d^{*t} \end{aligned} \tag{9.9}$$

and we can approximate (9.9) with

$$\hat{\mathbf{K}}_d^* \left( \mathbf{X}' \hat{\mathbf{V}}^{-1} \mathbf{X} \right)^{-1} \hat{\mathbf{K}}_d^{*t}$$

where

$$\hat{\mathbf{K}}_d^* = \sum_{i=1}^I \hat{\boldsymbol{\Sigma}}_{di}^* \mathbf{X}_{di}$$

and  $\hat{\mathbf{V}}$  is just  $\mathbf{V}$  with the estimates  $\hat{\boldsymbol{\varphi}}$ ,  $\hat{\boldsymbol{\beta}}$  and  $\hat{\mathbf{u}}$  substituted in as in previous sections.

Now we need to consider the terms

$$E \left( (\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d) (\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d)^t \right)$$

and

$$E \left( (\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d) (\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d)^t \right).$$

Note that  $\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d$  is a function of the in-sample data and is independent of  $\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d$  which is a function of the out-of-sample data. Therefore

$$E \left( (\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d) (\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d)^t \right) = E \left( (\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d) \right) E \left( (\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d)^t \right)$$

and

$$E \left( (\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d) (\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d)^t \right) = E \left( (\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d) \right) E \left( (\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d)^t \right).$$



Now using (9.6) and (9.7) we have

$$E\left(\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d\right) \approx \sum_{i=1}^I \boldsymbol{\Sigma}_{di}^{*r} \left( \mathbf{X}_{di} E\left(\hat{\boldsymbol{\beta}}\right) - \mathbf{X}_{di} \boldsymbol{\beta} \right) = \mathbf{0}_{2 \times 1}$$

since  $\hat{\boldsymbol{\beta}}$  is unbiased in a linear mixed model framework.

Also using (9.4)

$$\begin{aligned} E\left(\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right) &= -E\left(\sum_{i=1}^I \left(\mathbf{y}_{di}^r - \boldsymbol{\mu}_{di}^r + \boldsymbol{\mu}_{di}^r - \boldsymbol{\mu}_{di}^{*r}\right)\right) \\ &= -E\left(E\left(\sum_{i=1}^I \left(\mathbf{y}_{di}^r - \boldsymbol{\mu}_{di}^r + \boldsymbol{\mu}_{di}^r - \boldsymbol{\mu}_{di}^{*r}\right) \middle| \mathbf{u}_d\right)\right) \\ &= -E\left(\sum_{i=1}^I \boldsymbol{\mu}_{di}^r - \boldsymbol{\mu}_{di}^{*r}\right) \\ &\approx -E\left(\sum_{i=1}^I \boldsymbol{\Sigma}_{di}^{*r} \mathbf{u}_d\right) \\ &= \mathbf{0}_{2 \times 1}. \end{aligned}$$

Hence

$$E\left(\left(\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d\right)\left(\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right)^t\right) \approx \mathbf{0}_{2 \times 2}$$

and

$$E\left(\left(\tilde{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right)\left(\hat{\boldsymbol{\delta}}_d - \tilde{\boldsymbol{\delta}}_d\right)^t\right) \approx \mathbf{0}_{2 \times 2}.$$

We can now define an MSE estimator for the out-of-sample areas,  $d = D+1, D+2, \dots, D+R$  as

$$\widehat{MSE}\left(\hat{\boldsymbol{\delta}}_d\right) = \hat{\mathbf{K}}_d^* \left(\mathbf{X}^t \hat{\mathbf{V}}^{-1} \mathbf{X}\right)^{-1} \hat{\mathbf{K}}_d^{*t} + \hat{\mathbf{M}}_d^* \hat{\mathbf{W}}_d \hat{\mathbf{M}}_d^{*t} + \sum_{i=1}^I \hat{\boldsymbol{\Sigma}}_{di}^{*r}.$$

## 10. APPROXIMATE STANDARD ERRORS OF THE ELEMENTS OF $\hat{\beta}$

In Section 8 we derived the approximate MSE matrix of an estimator of the form  $K\hat{\beta} + M\hat{u}$ , where  $M$  and  $K$  are given matrices.

If we set  $K$  to be the identity matrix and  $M$  to be a matrix containing all 0's then we can apply (8.20) directly to obtain

$$\begin{aligned}\widehat{MSE}(\hat{\beta}) &\approx K\hat{P}K^t \\ &= (X^t\hat{V}^{-1}X)^{-1}\end{aligned}\tag{10.1}$$

by noting that the component  $G_4(\hat{\varphi})$  is not relevant here. The approximation (10.1) can be used to check the significance of the covariates in the model. Note that it is not appropriate to use  $-J_{\hat{\rho}}^{-1}$  defined in Section 4 for this purpose (since it is not a function of  $W$  and does not account properly for the influence of the random effects). Note also that in this framework (of linear mixed models),  $\hat{\beta}$  is unbiased, and so the approximate standard errors of the elements of  $\hat{\beta}$  are the square root of the diagonal entries in (10.1).

## 11. PARAMETRIC BOOTSTRAP MEAN SQUARED ERRORS

The mean squared error estimators in Sections 8 and 9 were derived using various approximations. In this section we consider an alternative mean squared error estimation technique based on a parametric bootstrap. This approach is very similar to the bootstrap method described in Molina *et al.* (2007). The main difference here is that instead of simulating  $\mathbf{u}_d$  from a univariate normal distribution, we need to simulate from a bivariate normal distribution. We also extend this bootstrap technique to incorporate out-of-sample small areas.

The bootstrap method is outlined as follows:

- Model fitting: fit the model to the original data (this will be for in-sample areas  $d = 1, 2, \dots, D$ ), obtaining parameter estimates  $\hat{\boldsymbol{\beta}}_1$  and  $\hat{\boldsymbol{\beta}}_2$  and  $\hat{\boldsymbol{\phi}}$ .
- Generation of random effects: For  $d = 1, 2, \dots, D, D+1, D+2, \dots, D+R$  (includes  $R$  out-of-sample areas), independently generate  $\mathbf{u}_d^* = (u_{d1}^*, u_{d2}^*)^t$  from a bivariate normal distribution with mean  $\mathbf{0}$  and covariance matrix

$$\hat{\mathbf{W}}_d = \begin{pmatrix} \hat{\phi}_1 & \hat{\phi}_{12} \\ \hat{\phi}_{12} & \hat{\phi}_2 \end{pmatrix}.$$

- Generation of a bootstrap population: for  $d = 1, 2, \dots, D, D+1, D+2, \dots, D+R$ , calculate the probabilities

$$p_{di1}^* = \frac{e^{\mathbf{x}_{di1}^t \hat{\boldsymbol{\beta}}_1 + u_{d1}^*}}{1 + e^{\mathbf{x}_{di1}^t \hat{\boldsymbol{\beta}}_1 + u_{d1}^*} + e^{\mathbf{x}_{di2}^t \hat{\boldsymbol{\beta}}_2 + u_{d2}^*}}$$

and

$$p_{di2}^* = \frac{e^{\mathbf{x}_{di2}^t \hat{\boldsymbol{\beta}}_2 + u_{d2}^*}}{1 + e^{\mathbf{x}_{di1}^t \hat{\boldsymbol{\beta}}_1 + u_{d1}^*} + e^{\mathbf{x}_{di2}^t \hat{\boldsymbol{\beta}}_2 + u_{d2}^*}}$$

and generate the following sample and non-sample multinomial vectors

$$\mathbf{y}_{di}^* = (y_{di1}^*, y_{di2}^*)^t \sim \text{Multinomial}(m_{di}, p_{di1}^*, p_{di2}^*),$$

$$\mathbf{y}_{di}^{*r} = (y_{di1}^{*r}, y_{di2}^{*r})^t \sim \text{Multinomial}(m_{di}^r, p_{di1}^*, p_{di2}^*).$$

Calculate the true area totals

$$\boldsymbol{\delta}_d^* = (\delta_{d1}^*, \delta_{d2}^*)^t = \sum_{i=1}^{10} (\mathbf{y}_{di}^* + \mathbf{y}_{di}^{*r}).$$

- (d) Model fitting to the bootstrap sample and parameter estimation: fit model to the bootstrap sample data  $\mathbf{y}_{di}^*$ , for  $d = 1, 2, \dots, D$  only, obtaining estimates  $\hat{\boldsymbol{\beta}}_1^*$ ,  $\hat{\boldsymbol{\beta}}_2^*$ ,  $\hat{\boldsymbol{\varphi}}^*$  and predicted values  $\hat{\mathbf{u}}_d^*$ . Note that the predicted values of  $\mathbf{u}_d$  for the out-of-sample areas  $d = D + 1, D + 2, \dots, D + R$  are always  $\hat{\mathbf{u}}_d^* = \mathbf{0}$ . From these, calculate individual predicted values for  $d = 1, 2, \dots, D + R$ :

$$\hat{\mathbf{y}}_{di}^{*r} = m_{di}^r \left( \frac{e^{\mathbf{x}_{di1}' \hat{\boldsymbol{\beta}}_1^* + \hat{u}_{di1}^*}}{1 + e^{\mathbf{x}_{di1}' \hat{\boldsymbol{\beta}}_1^* + \hat{u}_{di1}^*} + e^{\mathbf{x}_{di2}' \hat{\boldsymbol{\beta}}_2^* + \hat{u}_{di2}^*}}, \frac{e^{\mathbf{x}_{di2}' \hat{\boldsymbol{\beta}}_2^* + \hat{u}_{di2}^*}}{1 + e^{\mathbf{x}_{di1}' \hat{\boldsymbol{\beta}}_1^* + \hat{u}_{di1}^*} + e^{\mathbf{x}_{di2}' \hat{\boldsymbol{\beta}}_2^* + \hat{u}_{di2}^*}} \right)^t.$$

Then calculate bootstrap estimates of totals by

$$\hat{\boldsymbol{\delta}}_d^* = \left( \hat{\delta}_{d1}^*, \hat{\delta}_{d2}^* \right)^t = \sum_{i=1}^{10} \left( \mathbf{y}_{di}^* + \hat{\mathbf{y}}_{di}^{*r} \right).$$

- (e) Bootstrap replicates: repeat steps (b)–(d)  $B$  times. Let  $\delta_{d1}^{*(b)}$  and  $\delta_{d2}^{*(b)}$  denote the true values of the parameters and  $\hat{\delta}_{d1}^{*(b)}$  and  $\hat{\delta}_{d2}^{*(b)}$  the estimators that are obtained in the  $b$ -th repetition,  $b = 1, 2, \dots, B$ . The bootstrap estimators of the mean squared error matrices  $E\left((\hat{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d)(\hat{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d)^t\right)$ , for  $d = 1, 2, \dots, D + R$  are

$$\hat{E}\left(\left(\hat{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right)\left(\hat{\boldsymbol{\delta}}_d - \boldsymbol{\delta}_d\right)^t\right) = \begin{pmatrix} B^{-1} \sum_{b=1}^B \left(\hat{\delta}_{d1}^{*(b)} - \delta_{d1}^{*(b)}\right)^2 & B^{-1} \sum_{b=1}^B \left(\hat{\delta}_{d1}^{*(b)} - \delta_{d1}^{*(b)}\right)\left(\hat{\delta}_{d2}^{*(b)} - \delta_{d2}^{*(b)}\right) \\ B^{-1} \sum_{b=1}^B \left(\hat{\delta}_{d1}^{*(b)} - \delta_{d1}^{*(b)}\right)\left(\hat{\delta}_{d2}^{*(b)} - \delta_{d2}^{*(b)}\right) & B^{-1} \sum_{b=1}^B \left(\hat{\delta}_{d2}^{*(b)} - \delta_{d2}^{*(b)}\right)^2 \end{pmatrix}.$$

- (f) Other estimates: after step (e) it is also possible to estimate mean squared errors of other parameter estimates such as  $\hat{\boldsymbol{\beta}}$  and  $\hat{\boldsymbol{\varphi}}$ . For example, after implementing the repetitions  $b = 1, 2, \dots, B$  we can also approximate root relative mean squared errors of each element in  $\hat{\boldsymbol{\beta}}_1$  and  $\hat{\boldsymbol{\beta}}_2$ . Let index  $k$  denote the  $k$ -th element of  $\boldsymbol{\beta}_1$  (or  $\boldsymbol{\beta}_2$ ). Let  $\hat{\boldsymbol{\beta}}_{1(k)}$  denote the value from step (a) and  $\hat{\boldsymbol{\beta}}_{1(k)}^{*(b)}$  denote the bootstrap estimate at iteration  $b$ . A bootstrap estimate of the  $RMSE(\hat{\boldsymbol{\beta}}_{1(k)})$  is

$$\sqrt{B^{-1} \sum_{b=1}^B \left(\hat{\boldsymbol{\beta}}_{1(k)}^{*(b)} - \hat{\boldsymbol{\beta}}_{1(k)}\right)^2}$$

and this can be compared directly with the standard error estimate from Section 10.

## 12. AUXILIARY DATA

The previous sections described the theory behind producing small area estimates of labour force counts and estimating mean squared errors. We now apply this theory to real data.

The small areas in our model represent LGAs (Local Government Areas) which have boundaries defined according to the 2001 ASGC (Australian Standard Geographical Classification). These are chosen because they provide sufficiently fine geographical areas, but still have adequate samples in the area for analysis. The total number of in-sample areas in August 2001 and August 2006 are respectively 424 and 413 and the total number of out-of-sample small areas in August 2001 and August 2006 are respectively 220 and 231.

The first step is to fit an appropriate model to the in-sample data and produce PQL estimates of  $\mathbf{u}$  and  $\boldsymbol{\beta}$  and approximated REML or ML estimates of  $\boldsymbol{\varphi}$ . One essential part of this first step is to choose an appropriate set of explanatory variables to include in the models. That is, we need to fully define  $\mathbf{x}_{di1}$  and  $\mathbf{x}_{di2}$ .

Auxiliary data are available from a variety of different sources. The two main sources are administrative Centrelink benefit payment data from DEEWR (Department of Education, Employment and Workplace Relations) and the Australian Census of Population and Housing. These data are available for the time points August 2001 and August 2006 and so it will be possible to fit two separate models with the same explanatory variables for the two different time points as a comparison.

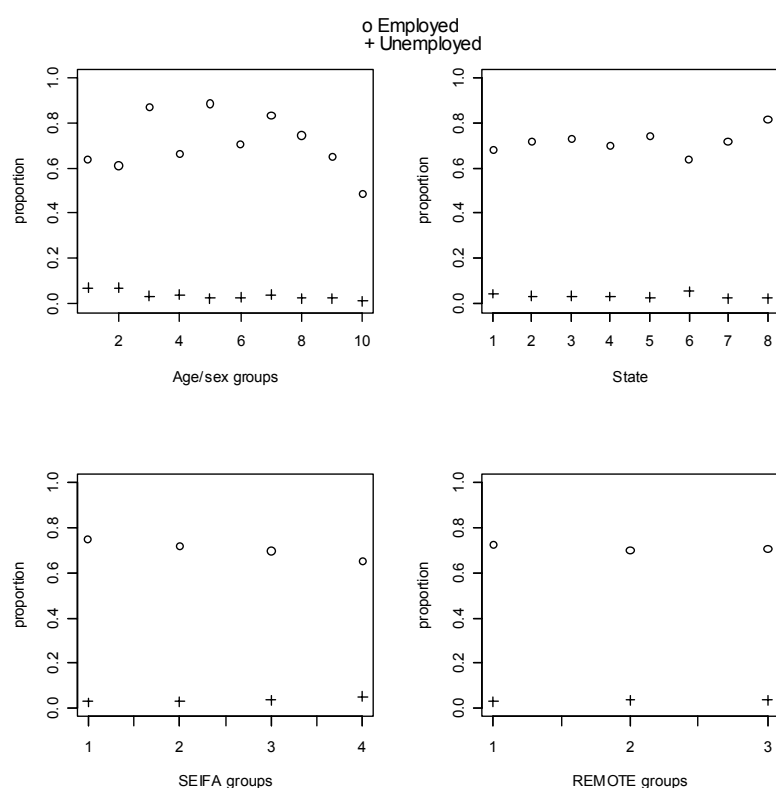
There are quite a large number of potential explanatory variables that could be used in the models. Ideally some kind of model variables selection process would need to be undertaken, however we do not consider this approach here. Instead we note that the ABS has already produced experimental estimates of small area labour force counts using three separate binomial logit mixed effects models (using the methodology outlined in Saei and Chambers, 2003). A careful model selection process was applied in this case to determine an appropriate set of explanatory variables for the binomial logit mixed models. For consistency and comparative purposes we will use the same set of explanatory variables that were used in these earlier models. Therefore each of our vectors  $\mathbf{x}_{di1}$  and  $\mathbf{x}_{di2}$  will contain the same set of 37 variables each. In summary, these variables are benefits payments variables, state indicators, age/sex indicators, remoteness indicators, socio-economic indexes for areas, and household type.

The following list outlines the explanatory variables in more detail.

- AS1–AS10: ten age/sex indicators. The five age groupings are 15–24 years, 25–34 years, 35–44 years, 45–54 years and 55–64 years. Note that we only consider the working age population for this analysis. AS1–AS10 are defined in order of increasing age and all odd AS groups are males. The base category is AS1.
- STATE1–STATE8: eight state indicators. These refer in order to New South Wales, Victoria, Queensland, South Australia, Western Australia, Tasmania, the Northern Territory and the Australian Capital Territory. The base category is New South Wales.
- NSA\_YAO: Proportion of population in the class registered to receive full payment of Newstart Allowance (unemployment benefits) or Youth Allowance (other).
- ASPAY1–ASPAY10: ten age/sex indicator by PAY interactions. PAY is the proportion of the population in the class registered to receive full other benefit payments. For example, disability support pension, parenting payments, partner allowance, wife pension, etc.. Note that this set up implies that the effect that PAY has on the probability is different for each age/sex class.
- REMOTE1–REMOTE3: three remoteness indicators. These refer to major city (REMOTE1), non-remote area (REMOTE2) and remote area (REMOTE3) as per the ASGC (Australian Standard Geographical Classification) 2001 (see ABS (2001) for further details). The base category is REMOTE1.
- SEIFA1–SEIFA4: four socio-economic index of advantage-disadvantage indicators. SEIFA1 indicates advantaged areas (whether the area is in the top 25% of SEIFA scores), SEIFA2 indicates the next 25%, SEIFA3 the next 25% and SEIFA4 indicates the most disadvantaged areas (whether the area is in the bottom 25% of SEIFA scores). The base category is SEIFA1. For further details see ABS (2003).
- HH1: Proportion of Census population in class that lives in dwelling consisting of married couple only or married couple with at least one child aged 15 or over.
- HH2: Proportion of Census population in class that lives in dwelling consisting of married couple with children all aged 0 to 14 years.
- HH3: Proportion of Census population in class that lives in dwelling consisting of one person only or one person with at least one child aged 15 or over.
- HH4: Proportion of Census population in class that lives in dwelling consisting of one person with children all aged 0 to 14 years.

Before we fit the multinomial logit mixed models, we undertake a brief exploratory data analysis using the August 2006 LFS and Census data to examine the suitability of our above chosen explanatory variables. Figure 12.1 contains average proportions of employed and unemployed within the ten AS groups, the eight STATE groups, the three REMOTE groups and the four SEIFA groups. The labour force proportions (especially employment) depend strongly on age and sex since the mean proportions vary across these categories. The AS indicators should certainly be considered as explanatory variables. Figure 12.1 also shows that the relationships between the labour force average proportions and each of STATE, SEIFA and REMOTE are in general not as strong as age/sex but there is still some variation.

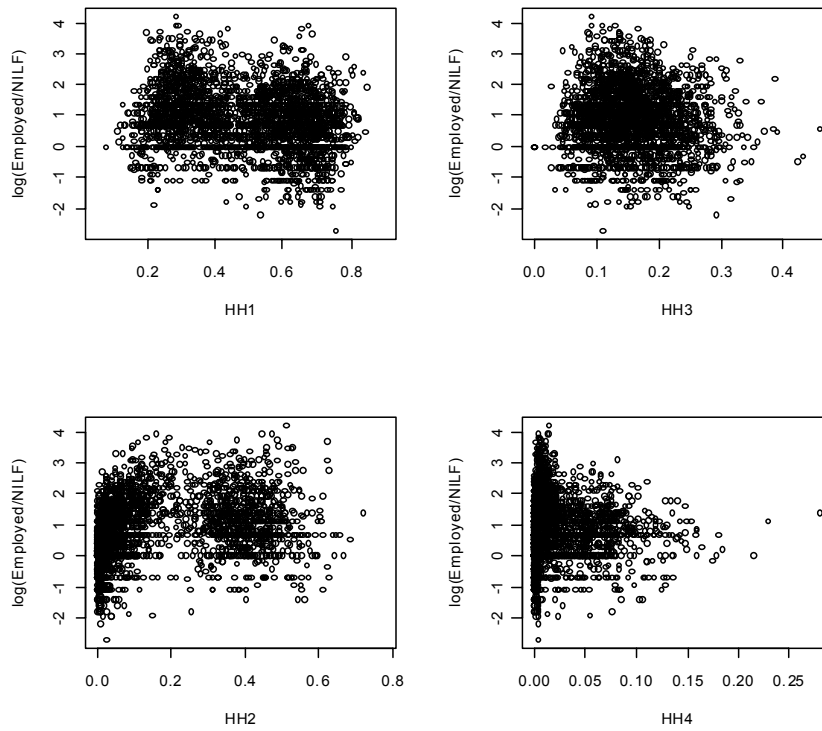
### 12.1 Plot of average labour force proportions within certain groups for 2006



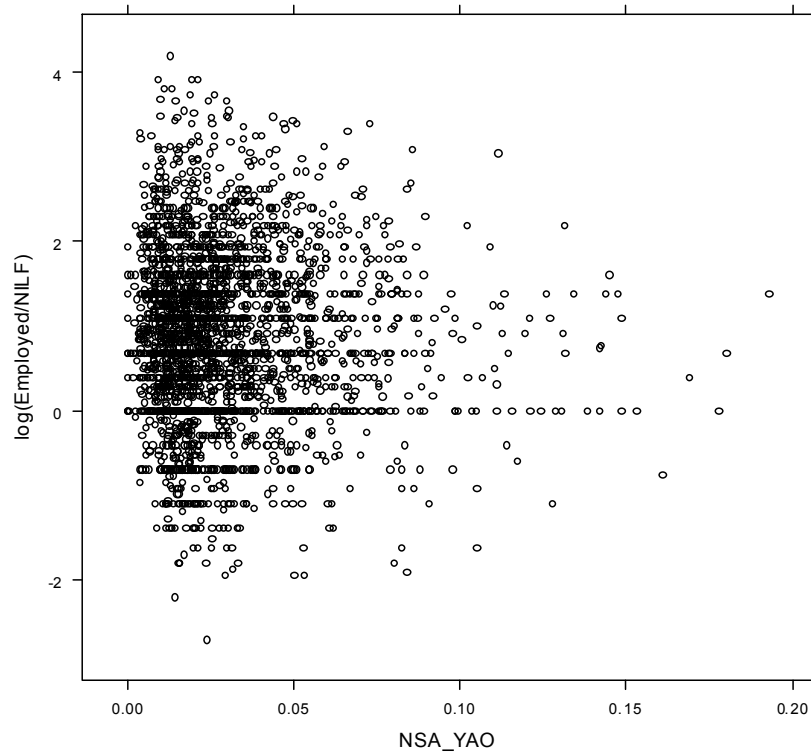
The variables we have examined so far are the categorical variables. We now examine relationships between  $\log(y_{di1}/y_{di3})$  for the cases  $y_{di3} > 0$  and  $y_{di1} > 0$  and the continuous explanatory variables for 2006. Note that we do not do this for  $y_{di2}$  since there are a large number of zeros in this case.

Figure 12.2 contains plots of  $\log(y_{di1}/y_{di3})$  versus each of HH1, HH2, HH3 and HH4. From these it appears there might be a weak association, especially for HH1 and HH3. A similar plot for the variable NSA\_YAO is given in figure 12.3, although we note that the linear relationship in this case does not look very strong. However, note that irrespective of this observation, a relationship may still exist between  $\log(y_{di2}/y_{di3})$  and NSA\_YAO and this is the reason why this variable is still included.

## 12.2 Plot of $\log(y_{di1}/y_{di3})$ versus each of HH1, HH2, HH3 and HH4 for 2006



## 12.3 Plot of $\log(y_{di1}/y_{di3})$ versus NSA\_YAO for 2006





#### 12.4 Plot of $\log(y_{di1}/y_{di3})$ versus PAY within AS groups for 2006

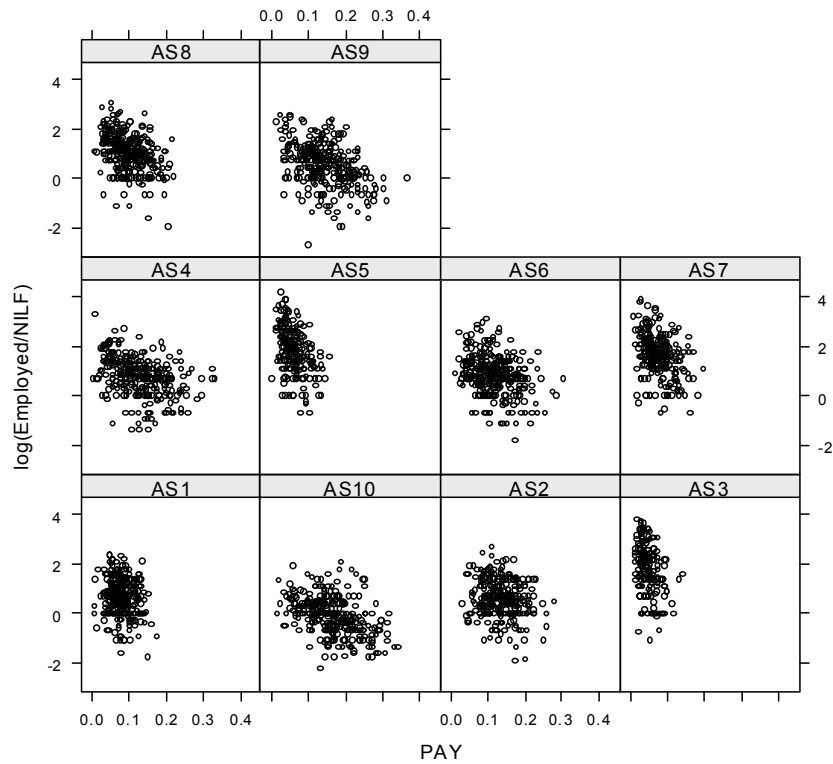


Figure 12.4 contains a plot of  $\log(y_{di1}/y_{di3})$  versus PAY within AS groups. There appears to be a negative relationship between  $\log(y_{di1}/y_{di3})$  and PAY. However, note that the slopes appear to differ between AS groups. This is the reason why we include the PAY variable in the model as an interaction between each of the AS indicators.

### 13. ESTIMATES OF MODEL PARAMETERS

Tables 13.1 and 13.2 contain the PQL estimates of  $\beta_1$  and  $\beta_2$  for the years 2006 and 2001. Estimates of the standard errors of the parameter estimates and p-values are also given in these tables. The p-values are calculated assuming that the distribution of the parameter estimate is approximately normal. Estimates which are significant at the 0.05 level are marked with a \*. Not all parameter estimates appear to be significant, but as we stated in the previous section, all variables are still included for consistency with the previously fitted binomial logit mixed models. As expected, for employment (and this holds for both years), the AS group indicators and ASPAY variables are highly significant. For both years, the variable NSA\_YAO is highly significant for unemployment. This is what we would expect as NSA\_YAO is related to unemployment.

The standard errors in tables 13.1 and 13.2 use the analytical approximation (10.1) which has many assumptions behind it. As a comparison, we also calculate parametric bootstrap RMSEs using the algorithm described in Section 11 (in particular see step (f)). We implement the parametric bootstrap using the 2006 data, with the simulation size set to  $B = 1000$  and use the combined PQL-REML algorithm for parameter estimation. A comparison between the 74 estimated analytical standard errors and the parametric bootstrap estimated root mean squared errors is given in figure 13.3 (figure 13.4 contains the smaller standard errors in the range 0–0.5 only). From these plots we can see that the differences between the estimated SEs and RMSEs are very small. Also in all cases we find that the bias component of the bootstrap RMSEs is small ( $<1.5\%$ ). Therefore the analytical SE estimates of  $\hat{\beta}$  appear to be good approximations.

The PQL parameter estimates in tables 13.1 and 13.2 use the approximated REML method to estimate the variance components  $\phi$ . As an alternative to this we could also have used the approximated ML method. Table 13.5 contains estimates of the variance components for both August 2001 and August 2006, using the two different estimation methods. To estimate the RRMSEs of the approximated ML and REML estimators, the parametric bootstrap of Section 11 (with  $B = 1000$ ) can be used with some slight modifications that we now briefly discuss:

- In step (a) of the algorithm, compute the PQL-REML estimates. These will be conditioned on in the rest of the simulation
- In step (d) we need to compute both the REML and ML estimates of the variance components.
- At the conclusion of the simulation we have 1000 REML and 1000 ML estimates of  $\phi$ . These can then be used to calculate estimates of bias and RMSEs.

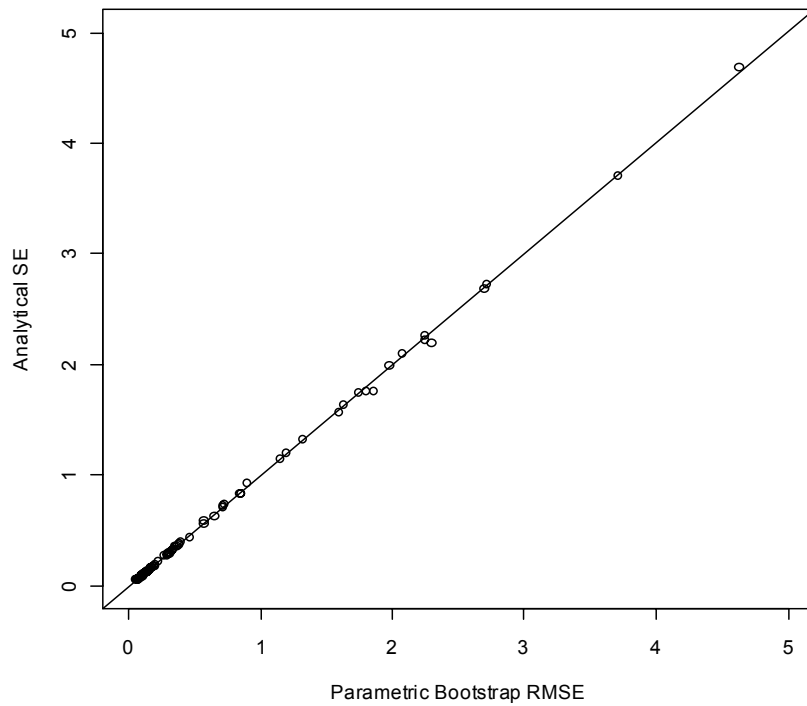
### 13.1 PQL Estimates of $\beta_1$ and $\beta_2$ for 2006 (REML used for variance components)

Variable	Employed				Unemployed			
	Estimate	SE	Z	p-value	Estimate	SE	Z	p-value
Intercept	0.428	0.192	2.23	0.0257 *	-1.890	0.358	-5.27	< 0.0001 *
STATE2	0.092	0.059	1.55	0.1210	-0.016	0.096	-0.17	0.8690
STATE3	0.175	0.062	2.81	0.0050 *	0.057	0.101	0.56	0.5740
STATE4	0.096	0.069	1.39	0.1650	-0.095	0.115	-0.83	0.4090
STATE5	0.096	0.065	1.47	0.1420	-0.162	0.114	-1.42	0.1560
STATE6	0.002	0.091	0.02	0.9870	-0.013	0.144	-0.09	0.9260
STATE7	-0.137	0.158	-0.87	0.3870	-0.730	0.276	-2.65	0.0081 *
STATE8	0.444	0.131	3.39	0.0007 *	0.155	0.199	0.78	0.4350
AS2	0.094	0.138	0.68	0.4950	-0.154	0.269	-0.57	0.5660
AS3	2.130	0.161	13.20	< 0.0001 *	0.986	0.297	3.33	0.0009 *
AS4	0.627	0.138	4.55	< 0.0001 *	-0.819	0.311	-2.63	0.0085 *
AS5	2.320	0.184	12.60	< 0.0001 *	1.210	0.382	3.18	0.0015 *
AS6	0.637	0.170	3.75	0.0002 *	-0.374	0.377	-0.99	0.3210
AS7	1.600	0.145	11.10	< 0.0001 *	-0.048	0.297	-0.16	0.8720
AS8	0.476	0.129	3.69	0.0002 *	-1.220	0.314	-3.88	0.0001 *
AS9	0.360	0.134	2.68	0.0074 *	-0.930	0.328	-2.84	0.0045 *
AS10	-0.647	0.132	-4.91	< 0.0001 *	-2.540	0.443	-5.73	< 0.0001 *
REMOTE2	0.101	0.078	1.30	0.1940	0.028	0.115	0.24	0.8090
REMOTE3	0.215	0.055	3.92	< 0.0001 *	0.063	0.089	0.71	0.4790
SEIFA2	-0.145	0.059	-2.48	0.0131 *	-0.041	0.094	-0.44	0.6630
SEIFA3	-0.157	0.073	-2.14	0.0324 *	0.057	0.122	0.46	0.6430
SEIFA4	-0.098	0.082	-1.19	0.2340	0.106	0.140	0.76	0.4500
ASPAY1	-5.430	1.210	-4.50	< 0.0001 *	-4.700	2.260	-2.08	0.0375 *
ASPAY2	-4.060	0.843	-4.81	< 0.0001 *	-1.970	1.570	-1.25	0.2110
ASPAY3	-17.100	2.720	-6.30	< 0.0001 *	-12.900	4.680	-2.76	0.0058 *
ASPAY4	-4.510	0.739	-6.10	< 0.0001 *	0.652	1.630	0.40	0.6890
ASPAY5	-16.300	1.740	-9.39	< 0.0001 *	-17.400	3.700	-4.70	< 0.0001 *
ASPAY6	-4.790	0.735	-6.52	< 0.0001 *	-2.640	1.770	-1.49	0.1360
ASPAY7	-10.500	1.330	-7.94	< 0.0001 *	-2.160	2.690	-0.80	0.4220
ASPAY8	-5.410	0.844	-6.41	< 0.0001 *	0.325	2.220	0.15	0.8840
ASPAY9	-6.390	0.638	-10.00	< 0.0001 *	-3.730	1.760	-2.12	0.0340 *
ASPAY10	-4.720	0.566	-8.34	< 0.0001 *	0.034	2.200	0.02	0.9880
NSA_YAO	-0.030	1.160	-0.03	0.9800	8.900	2.100	4.23	< 0.0001 *
HH1	0.637	0.219	2.91	0.0036 *	0.444	0.400	1.11	0.2670
HH2	-0.244	0.292	-0.84	0.4030	-0.582	0.591	-0.98	0.3250
HH3	1.510	0.362	4.18	< 0.0001 *	0.814	0.725	1.12	0.2630
HH4	1.090	0.924	1.18	0.2380	3.630	1.990	1.82	0.0688

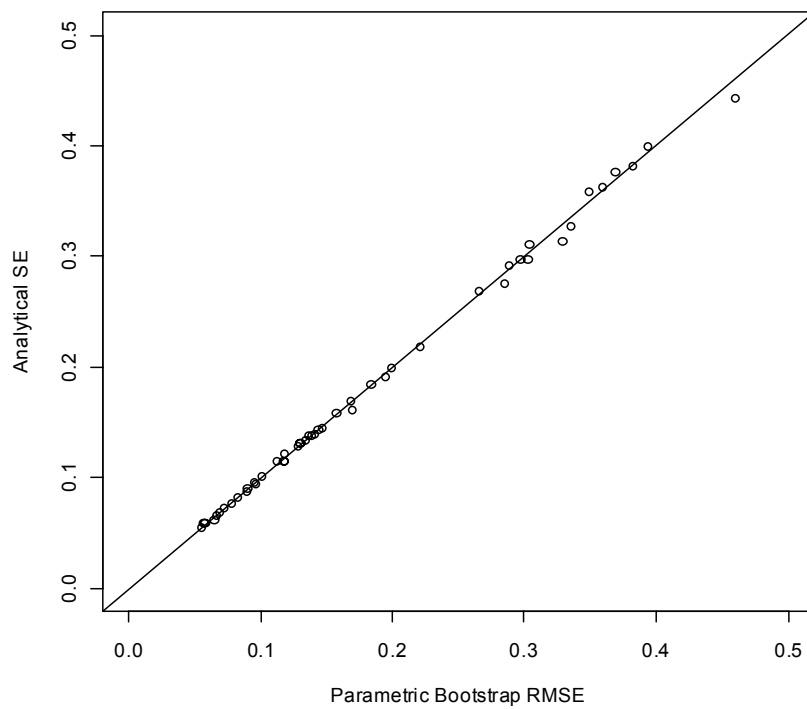
### 13.2 PQL Estimates of $\beta_1$ and $\beta_2$ for 2001 (REML used for variance components)

Variable	Employed				Unemployed			
	Estimate	SE	Z	p-value	Estimate	SE	Z	p-value
Intercept	0.779	0.204	3.82	0.0001 *	-2.610	0.370	-7.05	< 0.0001 *
STATE2	0.024	0.045	0.53	0.5980	0.122	0.084	1.45	0.1470
STATE3	0.131	0.048	2.74	0.0061 *	0.330	0.086	3.83	0.0001 *
STATE4	0.045	0.056	0.81	0.4170	0.119	0.104	1.15	0.2500
STATE5	0.039	0.052	0.75	0.4520	0.146	0.095	1.53	0.1260
STATE6	-0.018	0.071	-0.25	0.8020	0.220	0.124	1.77	0.0767
STATE7	0.075	0.124	0.61	0.5450	0.135	0.219	0.62	0.5370
STATE8	0.145	0.096	1.51	0.1310	0.101	0.177	0.57	0.5670
AS2	-0.089	0.151	-0.59	0.5540	-0.148	0.267	-0.55	0.5800
AS3	1.960	0.182	10.80	< 0.0001 *	1.650	0.287	5.77	< 0.0001 *
AS4	0.781	0.147	5.31	< 0.0001 *	0.200	0.282	0.71	0.4770
AS5	2.460	0.200	12.30	< 0.0001 *	1.800	0.344	5.24	< 0.0001 *
AS6	0.715	0.180	3.98	< 0.0001 *	-0.133	0.358	-0.37	0.7110
AS7	1.720	0.151	11.40	< 0.0001 *	0.966	0.265	3.64	0.0003 *
AS8	0.415	0.135	3.08	0.0021 *	-0.767	0.296	-2.59	0.0096 *
AS9	0.080	0.143	0.56	0.5770	-0.795	0.315	-2.52	0.0117 *
AS10	-0.766	0.151	-5.09	< 0.0001 *	-2.680	0.506	-5.30	< 0.0001 *
REMOTE2	0.036	0.058	0.62	0.5330	0.017	0.101	0.17	0.8690
REMOTE3	0.138	0.042	3.27	0.0011 *	0.077	0.075	1.02	0.3080
SEIFA2	-0.053	0.046	-1.15	0.2500	0.039	0.083	0.48	0.6350
SEIFA3	-0.069	0.060	-1.16	0.2460	-0.083	0.108	-0.77	0.4400
SEIFA4	-0.084	0.068	-1.24	0.2150	-0.141	0.122	-1.15	0.2500
ASPAY1	-3.830	1.010	-3.79	0.0002 *	0.201	1.650	0.12	0.9030
ASPAY2	-2.370	0.746	-3.18	0.0015 *	0.658	1.340	0.49	0.6230
ASPAY3	-9.230	2.630	-3.51	0.0004 *	-7.110	3.780	-1.88	0.0601
ASPAY4	-4.730	0.819	-5.78	< 0.0001 *	-1.030	1.680	-0.61	0.5410
ASPAY5	-13.800	1.770	-7.78	< 0.0001 *	-9.390	2.850	-3.29	0.0010 *
ASPAY6	-4.400	0.746	-5.90	< 0.0001 *	-0.604	1.680	-0.36	0.7200
ASPAY7	-11.200	1.300	-8.62	< 0.0001 *	-8.210	2.380	-3.45	0.0006 *
ASPAY8	-5.590	0.665	-8.41	< 0.0001 *	-1.780	1.760	-1.01	0.3120
ASPAY9	-4.460	0.507	-8.81	< 0.0001 *	-1.430	1.260	-1.14	0.2540
ASPAY10	-6.020	0.614	-9.82	< 0.0001 *	0.843	2.370	0.36	0.7230
NSA_YAO	-0.697	0.897	-0.78	0.4370	7.510	1.450	5.16	< 0.0001 *
HH1	0.164	0.217	0.75	0.4510	0.741	0.394	1.88	0.0601
HH2	-0.819	0.303	-2.70	0.0069 *	-0.911	0.557	-1.63	0.1030
HH3	1.580	0.457	3.45	0.0006 *	2.580	0.938	2.75	0.0060 *
HH4	0.271	1.310	0.21	0.8360	3.610	2.660	1.36	0.1740

### 13.3 Plot of analytical SEs versus bootstrap RMSEs of $\hat{\beta}$ for 2006



### 13.4 Plot of the small analytical SEs versus bootstrap RMSEs of $\hat{\beta}$ for 2006



### 13.5 Estimates of $\varphi$

Parameter	2001		2006	
	ML	REML	ML	REML
$\varphi_1$	0.0238	0.0280	0.0683	0.0758
$\varphi_2$	0.0468	0.0595	0.0691	0.0853
$\varphi_{12}$	0.0105	0.0125	0.0424	0.0466
$\rho$	0.316	0.307	0.617	0.580

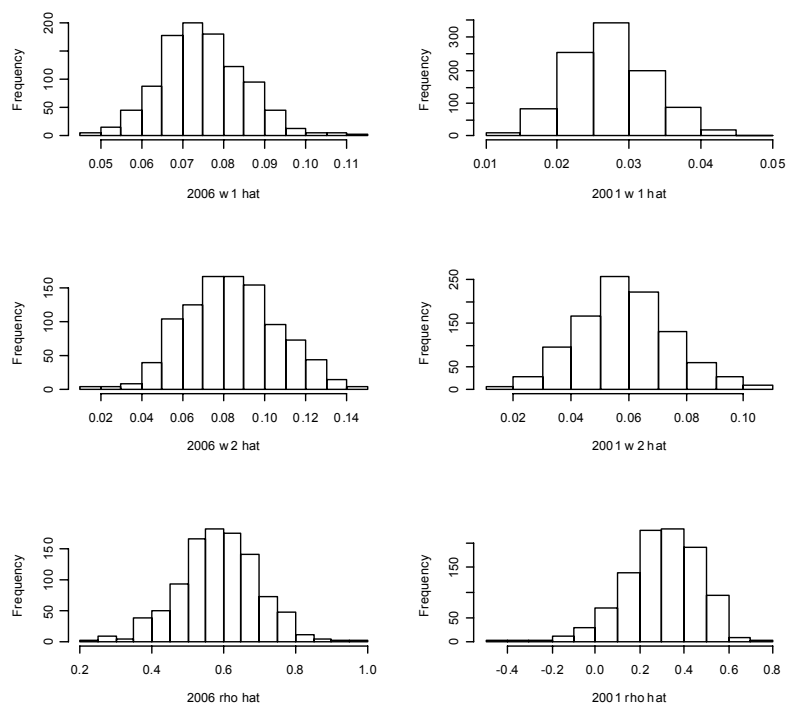
Table 13.6 contains the bootstrap results for  $\varphi$  for August 2006. In terms of bias, the REML estimator is performing reasonably well since all relative bias estimates are  $<2\%$  in absolute value. From these results it is clear that the REML estimator is to be preferred over the ML estimator since not only are the estimated biases smaller (as one might expect), but the RRMSE estimates are also smaller for the REML estimator.

### 13.6 Comparison of the REML and ML estimators of $\varphi$ for 2006

Parameter	Original estimate	Relative bias (%)		RRMSE (%)	
		ML	REML	ML	REML
$\varphi_1$	0.0758	-9.15	-1.56	15.56	13.22
$\varphi_2$	0.0853	-19.54	-1.85	31.87	27.08
$\rho$	0.580	6.71	0.80	22.68	20.15

It is also of interest to test the significance of the variance components  $\varphi$  (since if these are not significantly different from zero, then there is no point using a mixed effects model). In Molina *et al.* (2007), the authors test the significance of their single variance component using a likelihood ratio test (which is based on the approximate marginal likelihood). In theory we could apply a similar test here, but unfortunately the distribution of the likelihood ratio test statistic in our case will not be easy to derive under the null hypothesis. This is because under the null hypothesis  $\varphi = (\varphi_1, \varphi_2, \varphi_{12})^t = \mathbf{0}$  and  $\varphi$  is on the boundary of the parameter space which is a non standard condition. Molina *et al.* (2007) were able to apply this test because in the case of one variance component, the null distribution is known to be a mixture of two  $\chi^2$  distributions.

### 13.7 Histograms of $\hat{\varphi}_1$ , $\hat{\varphi}_2$ and $\hat{\rho}$ for 2006 and 2001



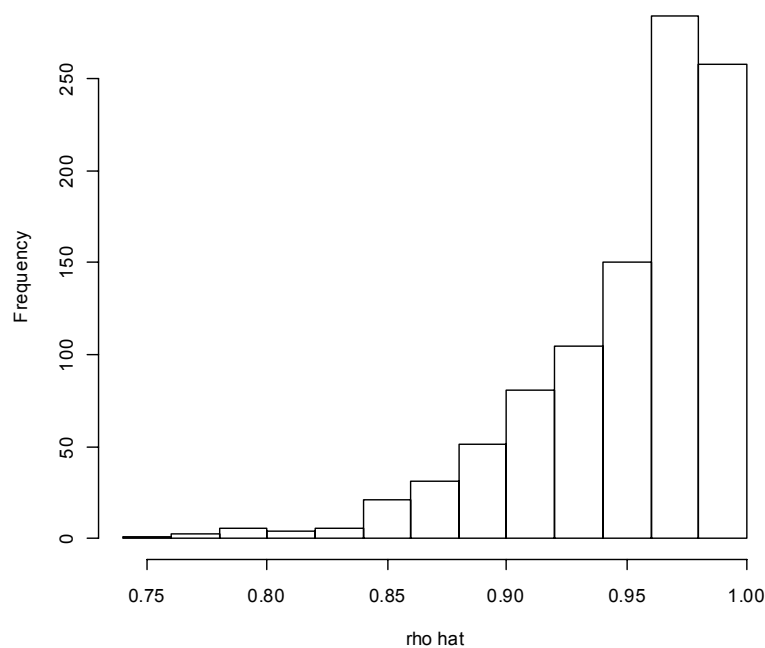
In any case, as an alternative, the previous parametric bootstrap used to produce table 13.6 will give some insight into the significance of the variance components. Using the 1000 REML estimates of each parameter, we build up empirical distributions under the fitted model of the REML estimators. Histograms of the 1000 sets of REML variance component estimates are given in figure 13.7. From these we can calculate approximate 95% parametric bootstrap confidence intervals using the percentile method. These confidence intervals are given in table 13.8. Notice that the confidence intervals for  $\varphi_1$  and  $\varphi_2$  for both years are not close to zero and the distributions look roughly symmetric, giving some evidence that the variance components are significant and a random effects model is appropriate in both years.

### 13.8 Approximate 95% confidence intervals for $\hat{\varphi}_1$ , $\hat{\varphi}_2$ and $\hat{\rho}$ for 2006 and 2001

Parameter	2006		2001	
	Lower	Upper	Lower	Upper
$\varphi_1$	0.0556	0.0948	0.0169	0.0395
$\varphi_2$	0.0439	0.1280	0.0272	0.0937
$\rho$	0.371	0.788	-0.048	0.581

If we were to apply the model in Molina *et al.* (2007) to our LFS data, we would be imposing that  $\boldsymbol{\varphi} = (\varphi_1, \varphi_1, \varphi_1)^t$ . Tables 13.5 and 13.8 give some evidence that this more restrictive model may not be appropriate especially for the 2001 data. To test the null hypothesis of  $\boldsymbol{\varphi} = (\varphi_1, \varphi_1, \varphi_1)^t$  for the 2006 data, a parametric bootstrap is used to generate empirical distributions under the null. The algorithm in Section 11 is used for this purpose, but note that we need to make one small change. Since we need to simulate under the null, in step (c) we replace all instances of  $u_{d2}^*$  with  $u_{d1}^*$  to ensure there is only one variance component. A histogram of the empirical distribution of  $\hat{\rho}$  under the null hypothesis is given in figure 13.9 for the 2006 data. There is strong evidence to suggest that  $\rho = 1$  does not hold since the observed  $\hat{\rho}$  is 0.580 and is nowhere near any of the simulated  $\hat{\rho}$  values under the null. Therefore our category specific random effects model appears to be more appropriate for our data than the more restrictive model used in Molina *et al.* (2007).

**13.9 Histogram of  $\hat{\rho}$  under the null for 2006**





## 14. RESIDUAL PLOTS AND GOODNESS OF FIT TESTS

In the case of the multinomial logit model with no random effects or variance components, it is usual to assess model fit by computing various different summary statistics, measuring differences between the observed and fitted values. For example, as stated in Dobson (2001), one could use the Pearson  $\chi^2$  statistic, Deviance or the Likelihood ratio chi-squared statistic.

Suppose we have no random effects or variance components  $\boldsymbol{\varphi}$  (the model is a GLM (generalised linear model)), then the standard Pearson  $\chi^2$  goodness of fit statistic is

$$\chi_p^2 = \sum_{j=1}^3 \sum_{d=1}^D \sum_{i=1}^{I_d} \frac{(y_{dij} - E_{dij})^2}{E_{dij}},$$

where  $E_{dij}$  is the estimated expected value of  $y_{dij}$ . The Pearson  $\chi^2$  goodness of fit test is used to determine whether or not the model appears to hold. Given that the model holds and under appropriate asymptotic conditions (expected counts  $E_{dij}$  are large), the Pearson  $\chi^2$  statistic will follow an approximate  $\chi_{2(N-p)}^2$  distribution, where  $N = \sum_{d=1}^D I_d = 7,820$  is the total number of multinomial observations for the August 2006 data (similarly we could also compute this for August 2001) and  $p = 37$  is the number of parameters in  $\boldsymbol{\beta}_1$  (or  $\boldsymbol{\beta}_2$ ). A large or small value of  $\chi_p^2$  indicates that overall the model does not appear to hold (the null hypothesis is that the model holds).

Unfortunately we cannot apply the above test to our data since our model is not a GLM. The application of goodness of fit tests in a GLMM (generalised linear mixed model) situation is not straight forward theoretically. For instance, the observations are no longer independent and the Pearson statistic and other statistics are not guaranteed to have a  $\chi^2$  distribution under the null hypothesis even when the expected counts are large. Also, another problem is that the estimates  $E_{dij}$  are not easy to calculate as they involve integrals with no closed form solution.

There are two approaches that we consider to get around the above issues. The first is to use as  $E_{dij}$ , the conditional expectation estimates  $m_{di}\hat{p}_{dij}$  predicted from the model and the second is to use an estimate of  $E(y_{dij})$  from a Taylor series expansion about  $\boldsymbol{u}_d = \mathbf{0}$ , since in our case the variance components are small. The distribution of these  $\chi^2$  statistics can then be estimated by applying the parametric bootstrap of Section 11 using  $B = 1000$ .

For our multinomial logit mixed model sample data we have

$$E(y_{dij}) = E(E(y_{dij} | \boldsymbol{u}_d)) = m_{di}E(p_{dij}) \quad (14.1)$$

and

$$\begin{aligned} \text{Var}(y_{dij}) &= E(\text{Var}(y_{dij} | \mathbf{u}_d)) + \text{Var}(E(y_{dij} | \mathbf{u}_d)) \\ &= m_{di} E(p_{dij}) (1 - E(p_{dij})) + m_{di} (m_{di} - 1) \text{Var}(p_{dij}). \end{aligned} \quad (14.2)$$

It can be shown using Taylor series expansions that when the variance components are small,

$$E(p_{di1}) \approx p_{di1}^* + \frac{1}{2} \begin{pmatrix} \varphi_1 p_{di1}^* (1 - 2p_{di1}^*) (1 - p_{di1}^*) \\ + \varphi_2 (p_{di2}^{*2} p_{di1}^* - p_{di1}^* p_{di2}^* (1 - p_{di2}^*)) \\ + 2\varphi_{12} (2p_{di1}^{*2} p_{di2}^* - p_{di1}^* p_{di2}^*) \end{pmatrix}, \quad (14.3)$$

$$E(p_{di2}) \approx p_{di2}^* + \frac{1}{2} \begin{pmatrix} \varphi_1 (p_{di1}^{*2} p_{di2}^* - p_{di2}^* p_{di1}^* (1 - p_{di1}^*)) \\ + \varphi_2 p_{di2}^* (1 - 2p_{di2}^*) (1 - p_{di2}^*) \\ + 2\varphi_{12} (2p_{di2}^{*2} p_{di1}^* - p_{di1}^* p_{di2}^*) \end{pmatrix} \quad (14.4)$$

and

$$E(p_{di3}) \approx p_{di3}^* + \frac{1}{2} \begin{pmatrix} \varphi_1 (p_{di1}^{*2} p_{di3}^* - p_{di3}^* p_{di1}^* (1 - p_{di1}^*)) \\ + \varphi_2 (p_{di2}^{*2} p_{di3}^* - p_{di3}^* p_{di2}^* (1 - p_{di2}^*)) \\ + 4\varphi_{12} p_{di1}^* p_{di2}^* p_{di3}^* \end{pmatrix}, \quad (14.5)$$

where  $p_{di1}^*$ ,  $p_{di2}^*$  and  $p_{di3}^*$  are respectively  $p_{di1}$ ,  $p_{di2}$  and  $p_{di3}$  with  $\mathbf{u}_d$  replaced with the zero vector (the Taylor series expansions were taken about the point  $\mathbf{u}_d = \mathbf{0}$ ). Also, we can show that

$$\text{Var}(p_{di1}) \approx \varphi_1 p_{di1}^{*2} (1 - p_{di1}^*)^2 + \varphi_2 p_{di1}^{*2} p_{di2}^{*2} - 2\varphi_{12} p_{di1}^{*2} (1 - p_{di1}^*) p_{di2}^*, \quad (14.6)$$

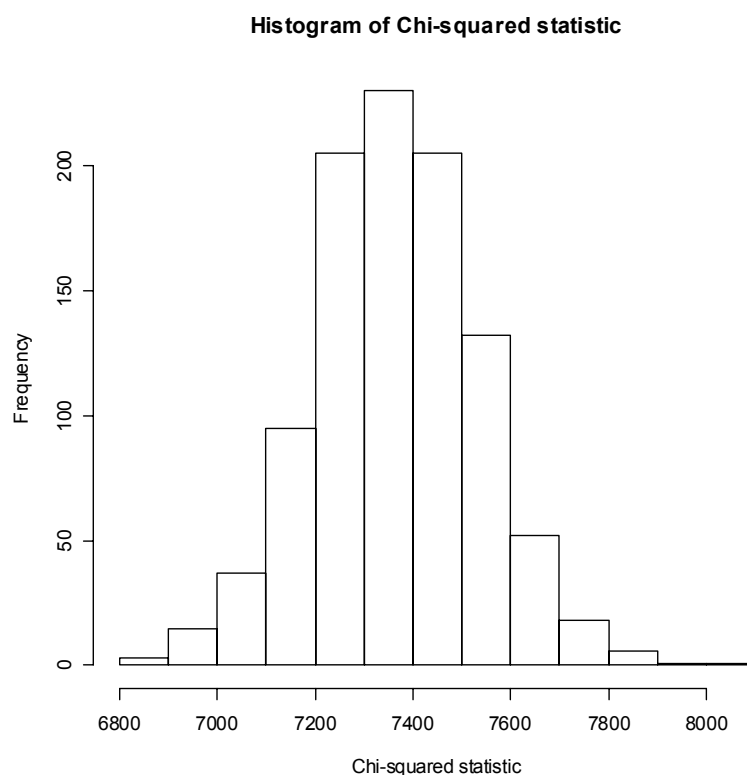
$$\text{Var}(p_{di2}) \approx \varphi_1 p_{di1}^{*2} p_{di2}^{*2} + \varphi_2 p_{di2}^{*2} (1 - p_{di2}^*)^2 - 2\varphi_{12} p_{di2}^{*2} p_{di1}^* (1 - p_{di2}^*) \quad (14.7)$$

and

$$\text{Var}(p_{di3}) \approx \varphi_1 p_{di1}^{*2} p_{di3}^{*2} + \varphi_2 p_{di2}^{*2} p_{di3}^{*2} + 2\varphi_{12} p_{di1}^* p_{di2}^* p_{di3}^{*2}. \quad (14.8)$$

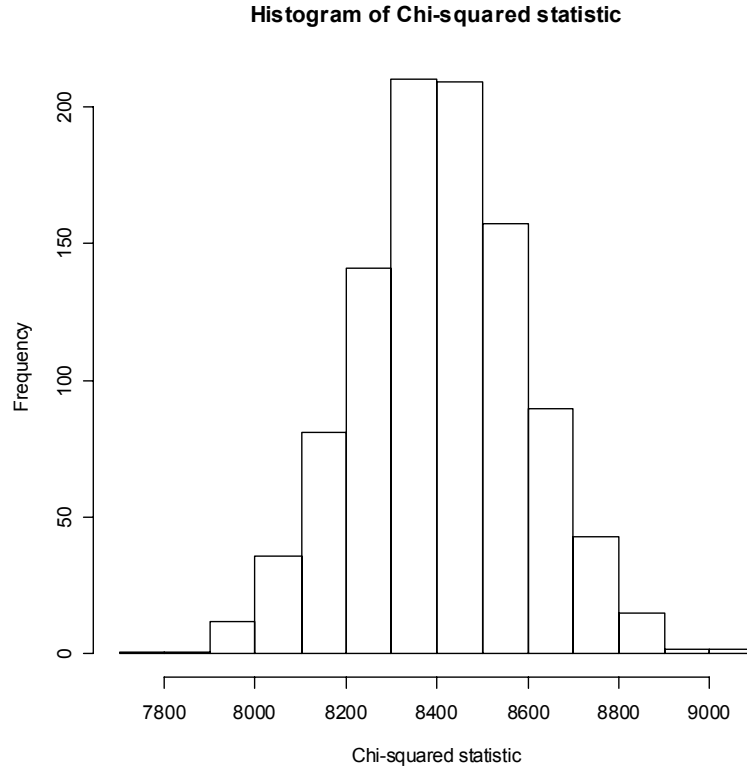
We are now at a point where we can approximate  $E_{dij}$  in the  $\chi^2$  statistic using estimates of  $E(y_{dij})$  based on the Taylor series approximations (14.3)–(14.5) and noting the relationship (14.1). Note that within these approximations, the parameters  $\beta$  and  $\phi$  are replaced by their estimates. Call this the unconditional approximation method. The other method as we mentioned earlier is to use the estimates of the conditional expectation  $E(y_{dij} | \mathbf{u}_d)$  for  $E_{dij}$  in the  $\chi^2$  statistic. Call this method the conditional method.

#### 14.1 Histogram of the $\chi^2$ statistic using the conditional method for 2006



Using the parametric bootstrap of Section 11, we generate empirical distributions of the two  $\chi^2$  statistics under the fitted model based on a simulation of size  $B = 1000$ . The two distributions are summarised in figures 14.1 and 14.2 for the August 2006 data. An approximate 95% confidence interval for the conditional method is (7035, 7708) and an approximate 95% confidence interval for the unconditional approximation method is (8045, 8775). The values of the  $\chi^2$  statistic for the sample data are 7669 (conditional method) and 8703 (unconditional approximation method). Clearly these values are within the appropriate approximate confidence intervals, suggesting that overall both the models for  $y_{dij}$  and  $y_{dij} | \mathbf{u}_d$  are adequate for the August 2006 sample data (similar conclusions can also be drawn for the August 2001 data, but details are not given here).

## 14.2 Histogram of the $\chi^2$ statistic using the unconditional approximation method for 2006



One issue associated with goodness of fit tests is that they only produce one overall summary measure. It is of interest to also examine individual deviations and look for outliers/influential points. Residual plots are useful for this purpose. For each labour force status  $j = 1, 2, 3$ , define the following conditional standardised residuals for the in-sample data (Molina *et al.* (2007) call these Pearson residuals)

$$r_{dij}^c = \frac{y_{dij} - m_{di} \hat{p}_{dij}}{\sqrt{m_{di} \hat{p}_{dij} (1 - \hat{p}_{dij})}}.$$

We can also calculate approximate unconditional standardised residuals as follows

$$r_{dij}^{uc} = \frac{y_{dij} - \hat{E}(y_{dij})}{\sqrt{\widehat{Var}(y_{dij})}}$$

where  $\hat{E}(y_{dij})$  and  $\widehat{Var}(y_{dij})$  are estimates based on the earlier Taylor series approximations (14.3)–(14.8) and the relationships (14.1) and (14.2) with parameters  $\boldsymbol{\beta}$  and  $\boldsymbol{\varphi}$  replaced with estimates.

Interestingly, we can also calculate the following summary statistics for each labour force category  $j = 1, 2, 3$

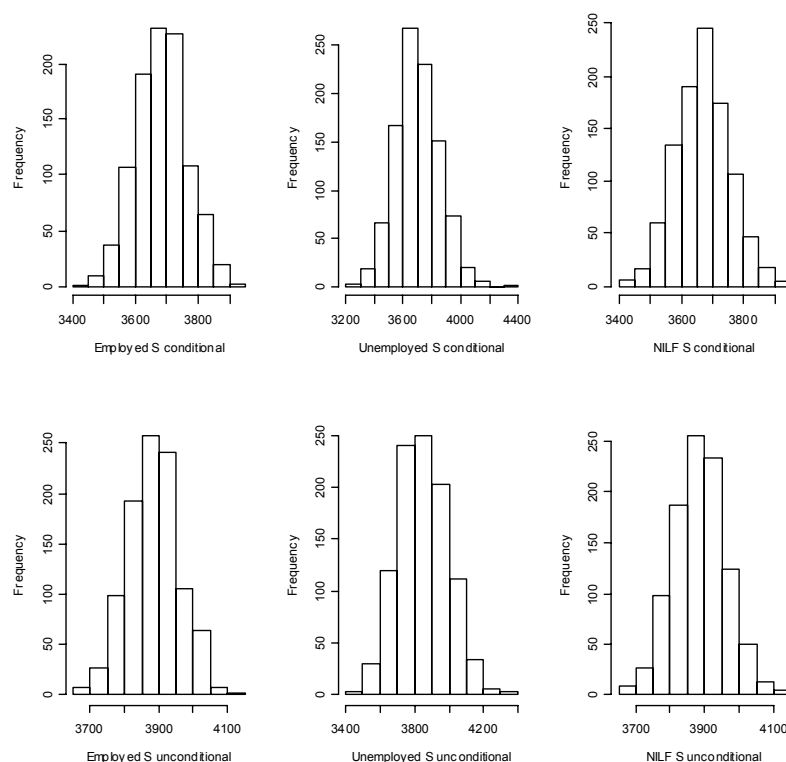
$$S_j^c = \sum_{d=1}^D \sum_{i=1}^{I_d} (r_{dij}^c)^2$$

and

$$S_j^{uc} = \sum_{d=1}^D \sum_{i=1}^{I_d} (r_{dij}^{uc})^2$$

and compare these with appropriate empirical distributions generated from a parametric bootstrap under the fitted model with  $B = 1000$  (again we use the bootstrap approach in Section 11). The empirical distributions are given in figure 14.3. Table 14.4 contains the  $S^c$  and  $S^{uc}$  values for the August 2006 sample data as well as approximate 95% parametric bootstrap confidence intervals using the percentile method. None of the  $S^c$  or  $S^{uc}$  values are within their corresponding confidence intervals. These statistics indicate that there may be underdispersion present for the Unemployed counts and overdispersion for both the Employed and NILF counts. Previously when we calculated the overall  $\chi^2$  statistics we did not obtain significant values. This is because the over and underdispersion in a way averaged themselves out overall.

### 14.3 Histogram of the $S^c$ and $S^{uc}$ statistics for 2006 from a parametric bootstrap



#### 14.4 $S^C$ and $S^{UC}$ values and associated parametric bootstrap 95% confidence intervals

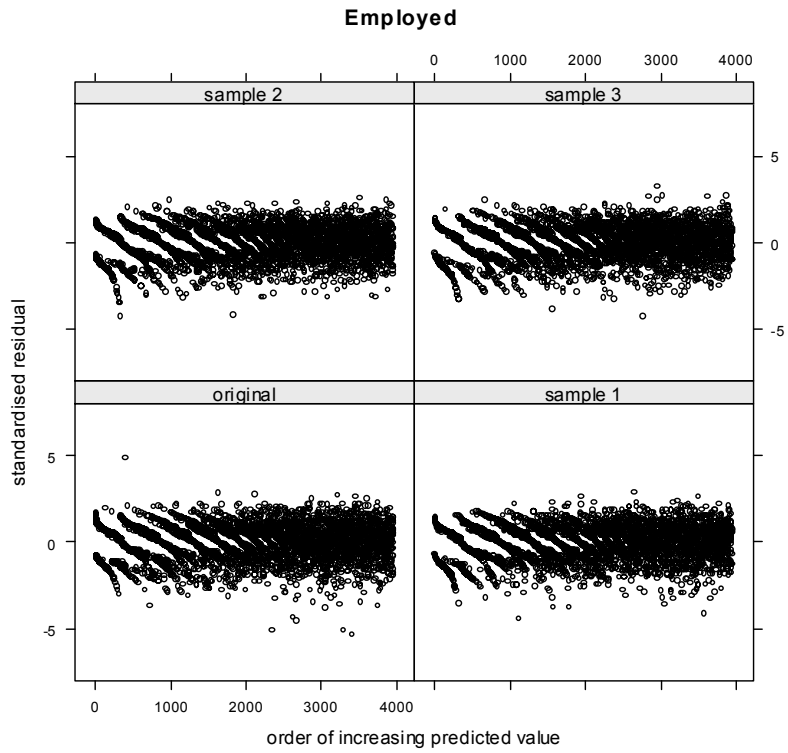
LFS status	Conditional			Unconditional		
	$S^C$	Lower	Upper	$S^{UC}$	Lower	Upper
Employed	4,224	3,522	3,848	4,306	3,741	4,030
Unemployed	3,393	3,407	4,000	3,565	3,578	4,137
Not in labour force	4,253	3,503	3,842	4,407	3,736	4,043

We now examine some residual plots to look more closely at this potential under and overdispersion issue. Figure 14.5 contains the standardised residuals  $r_{di1}^c$  versus the order of increasing predicted values  $m_{di}\hat{P}_{di1}$  for the 2006 employed sample data. The panel labelled original contains the original sample and the other three panels contain values calculated from three parametric bootstrap samples. Similarly, figure 14.6 contains the unconditional standardised residuals  $r_{di1}^{uc}$  versus the order of increasing predicted values  $\hat{E}(y_{di1})$ . Similar plots for NILF and unemployed are given in figures 14.7–14.10.

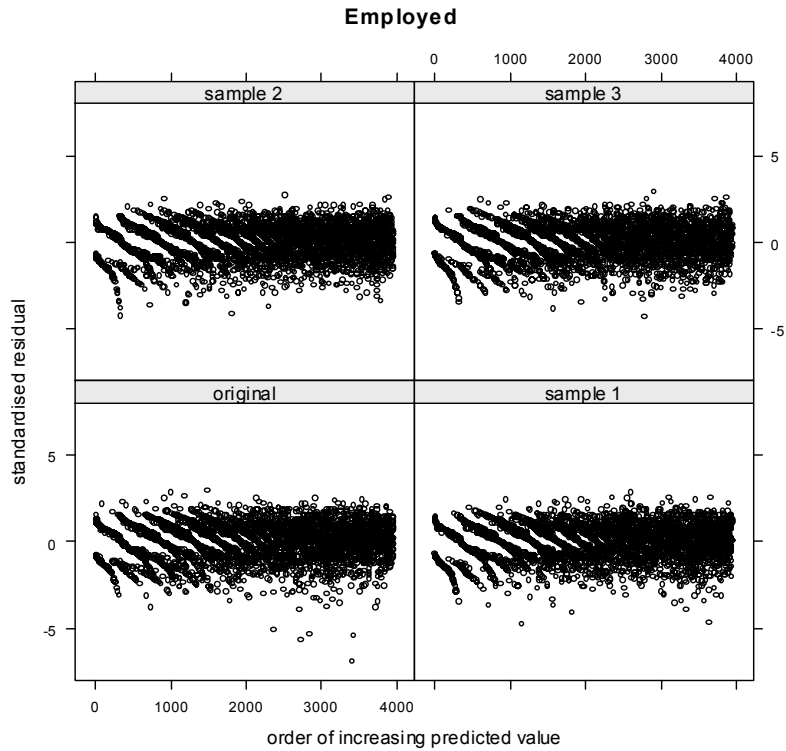
When comparing the bootstrap generated employed and NILF samples with the original data, it appears that the original data contain a small number of larger absolute residual values for both Employed and NILF. We actually recalculated the  $S^C$  and  $S^{UC}$  statistics by setting these few larger absolute residuals to zero for employed and NILF, but the resulting  $S^C$  and  $S^{UC}$  statistics were still significantly large. Hence these couple of larger absolute residuals are not solely responsible for the apparent overdispersion. In any case, when ignoring the small number of larger absolute residuals, the sample and original data distributions look roughly similar and hence we argue that the apparent under and overdispersion is not large enough for us to be overly concerned. In the significance tests we are clearly only picking up small significant differences and we suspect this is because our sample sizes are large.

We mentioned above that there were a small number of larger residuals in absolute value than one might expect when compared to the bootstrap samples. Most of these do not appear to be overly large and the NILF and employed ones mostly correspond to the same units. Upon further investigation, the only thing these outlier units have in common is that they are mostly all from remote areas (REMOTE3=1). Therefore the model may not be doing as well in the remote areas. Note that there is not much we can do about this issue because the sample sizes are small in the remote areas and we do not have any further covariates available.

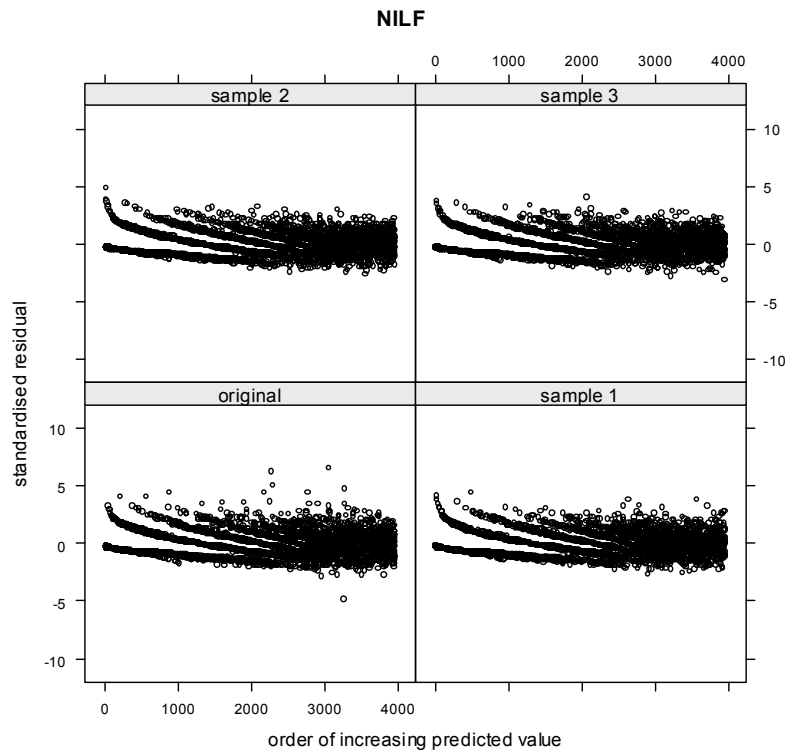
14.5 Plots of the conditional employed standardised residuals versus order of predicted values for original 2006 data and for three parametric bootstrap simulations



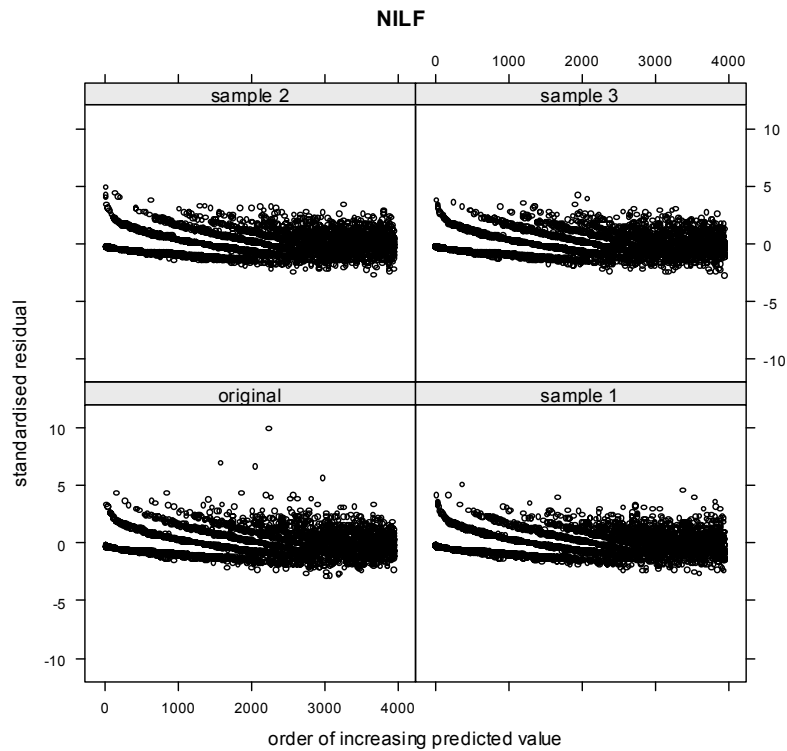
14.6 Plots of the unconditional employed standardised residuals versus predicted values for original 2006 data and for three parametric bootstrap simulations



**14.7 Plots of the conditional NILF standardised residuals versus order of predicted values for original 2006 data and for three parametric bootstrap simulations**

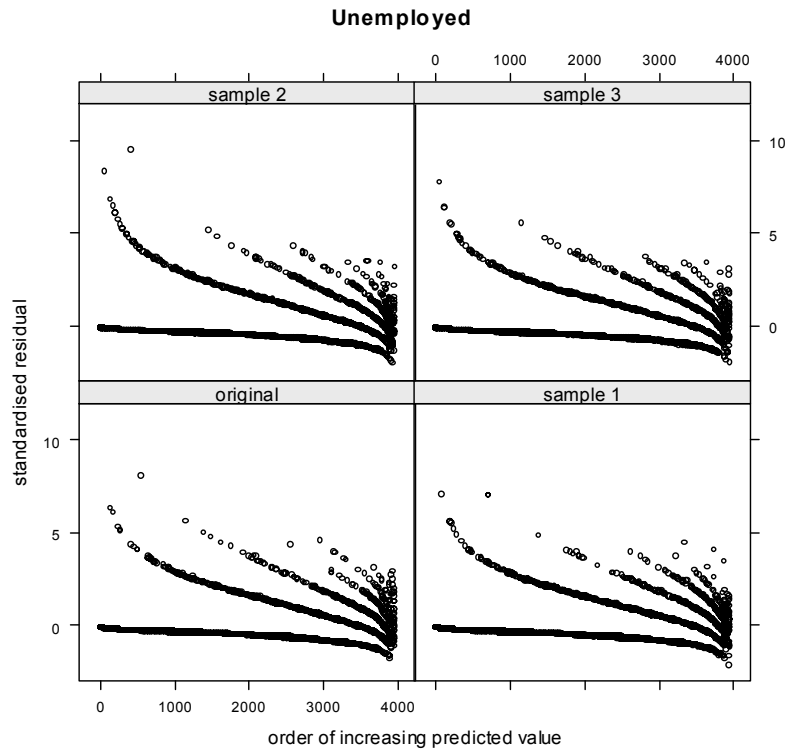


**14.8 Plots of the unconditional NILF standardised residuals versus predicted values for original 2006 data and for three parametric bootstrap simulations**

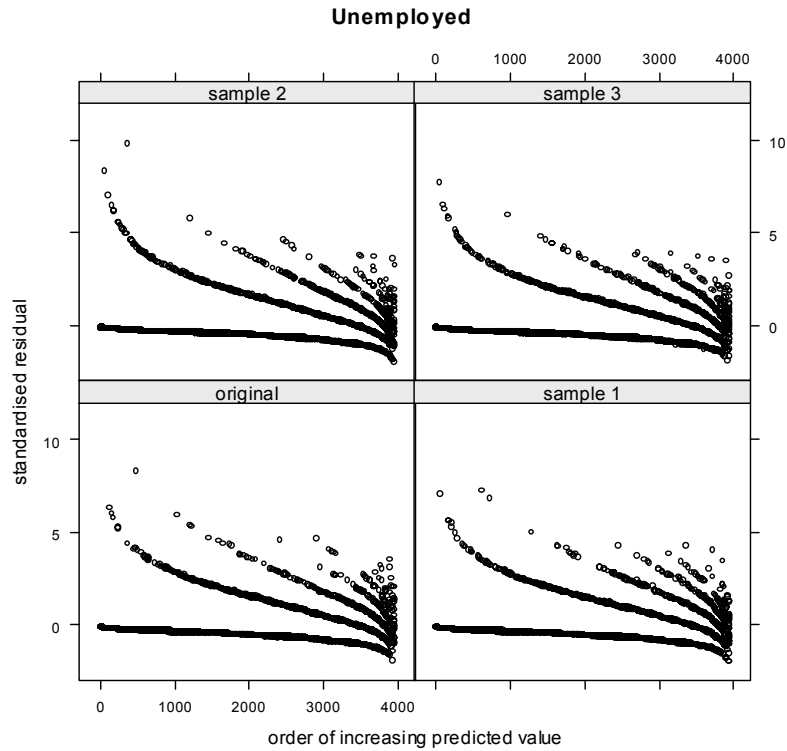




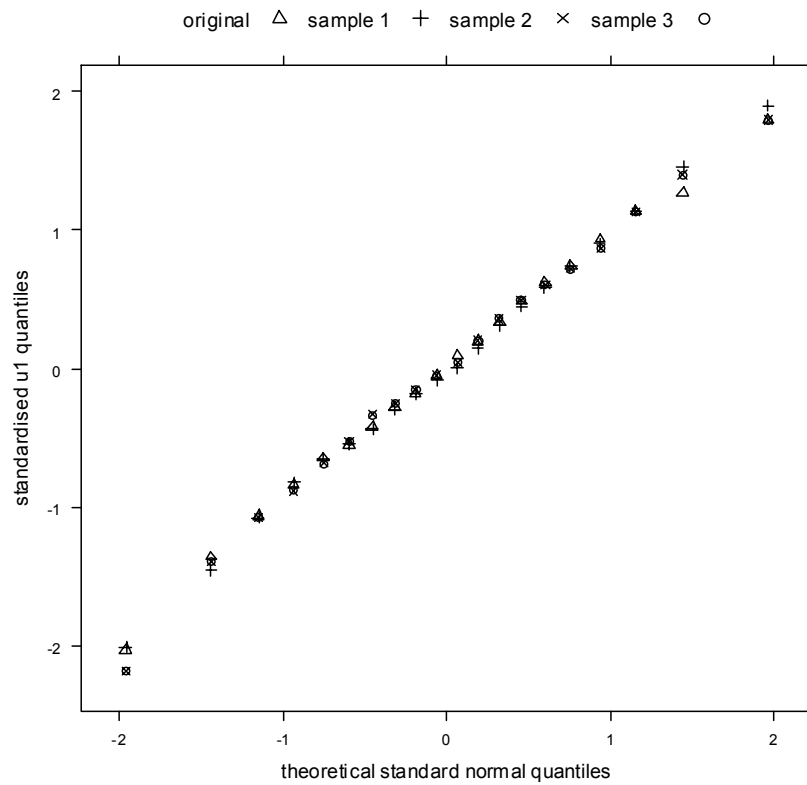
**14.9 Plots of the conditional unemployed standardised residuals versus order of predicted values for original 2006 data and for three parametric bootstrap simulations**



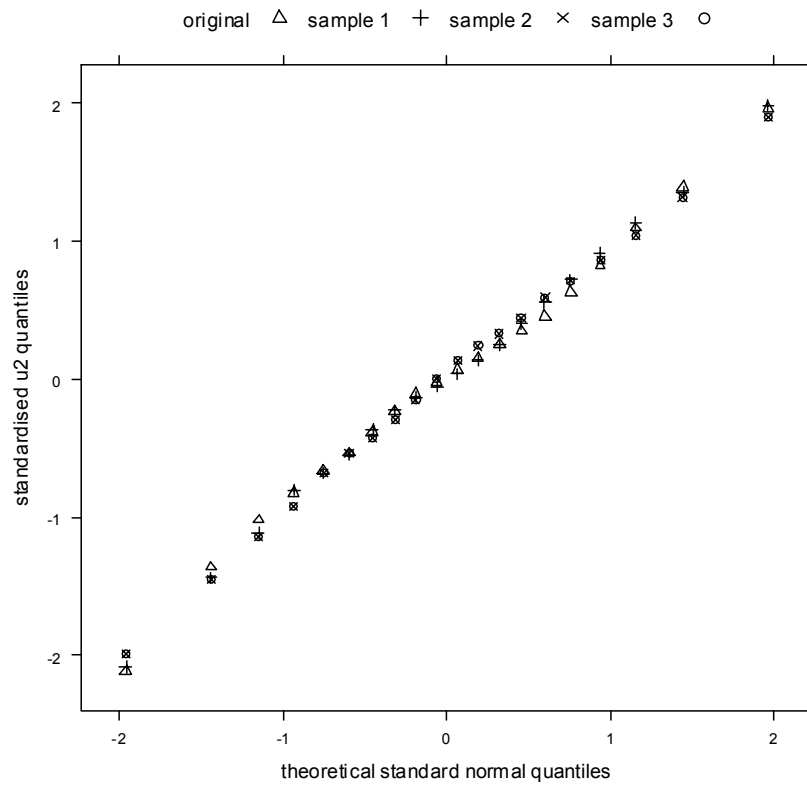
**14.10 Plots of the unconditional unemployed standardised residuals versus predicted values for original 2006 data and for three parametric bootstrap simulations**



14.11 Q-Q plots of  $\hat{u}_{d1}$  for the original 2006 sample and three bootstrap samples

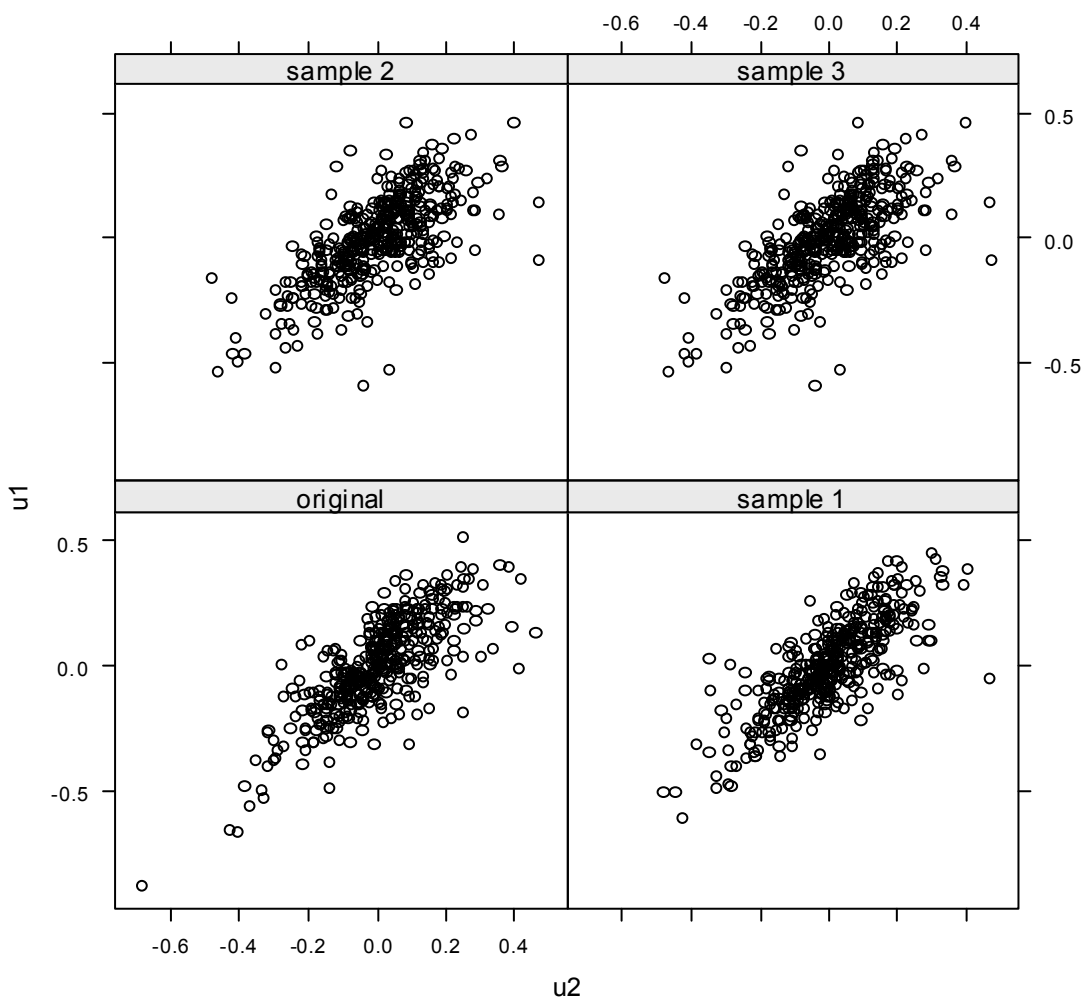


14.12 Q-Q plots of  $\hat{u}_{d2}$  for the original 2006 sample and three bootstrap samples



So far we have said nothing about the predicted random effects  $\hat{\mathbf{u}}$ . Figure 14.11 contains a plot of the quantiles of standardised  $\hat{u}_{d1}$  values versus theoretical standard normal quantiles. This is done for the original 2006 sample and three samples generated using the parametric bootstrap. Similarly, figure 14.12 contains a plot of the quantiles of standardised  $\hat{u}_{d2}$  values versus theoretical standard normal quantiles. Figure 14.13 contains a plot of  $\hat{u}_{d1}$  versus  $\hat{u}_{d2}$ . From these figures it appears that the original 2006 sample  $\hat{\mathbf{u}}_d$  values behave similar to those obtained from ‘typical’ samples. The estimated  $\hat{\mathbf{u}}_d$  values from the original sample therefore do not appear to give any indication of model departure.

**14.13** Plots of  $\hat{u}_{d1}$  versus  $\hat{u}_{d2}$  for the original 2006 sample and three bootstrap samples



## 15. SMALL AREA ESTIMATES AND MSE ESTIMATES

Similar to figure 8 in Molina *et al.* (2007), figure 15.1 contains plots of the ratios of direct RSE estimates to model analytical RRMSE estimates versus sample sizes for 2006. The line  $y = 1$  is also plotted. A ratio greater than 1 indicates we get gains by using the model based approach, whereas a ratio less than 1 indicates we get gains by using the direct survey estimation approach. Since all ratios are greater than 1 we are always getting gains by using the model based approach. The gains are quite large when the sample sizes are small and are small when the sample sizes are larger. Therefore when the sample size is small, the model based estimates have much smaller estimated MSEs than the direct survey estimates. Hence we have successfully reduced the MSEs by using a model based approach.

**15.1 Plots of the ratios of direct RSE estimates to model analytical RRMSE estimates versus sample sizes for 2006**

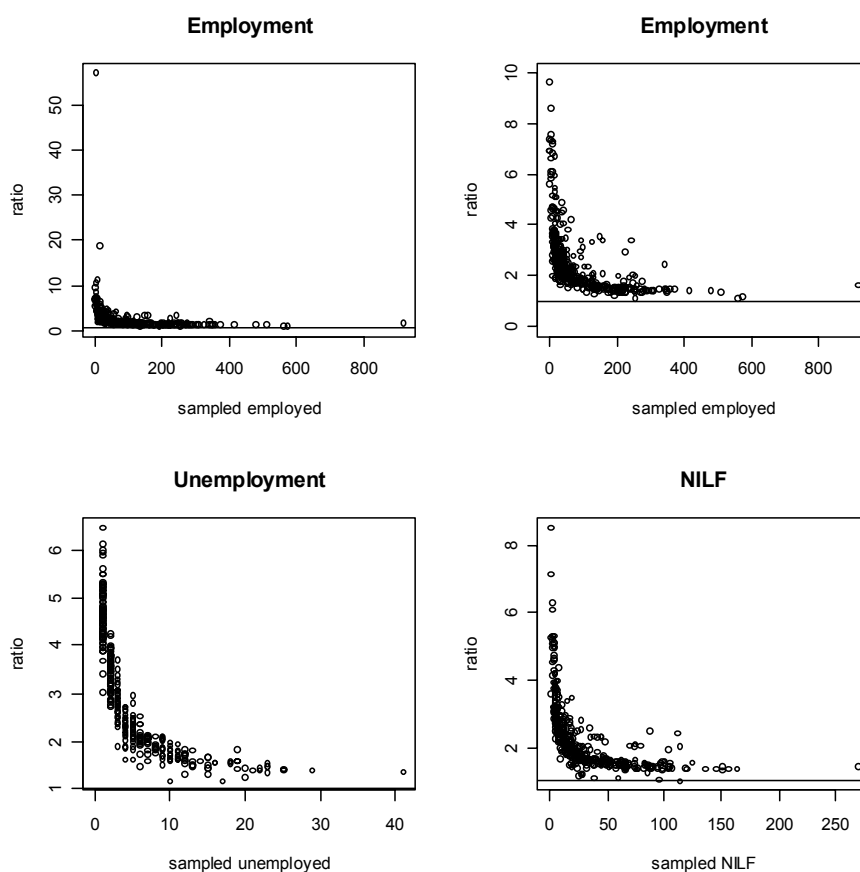
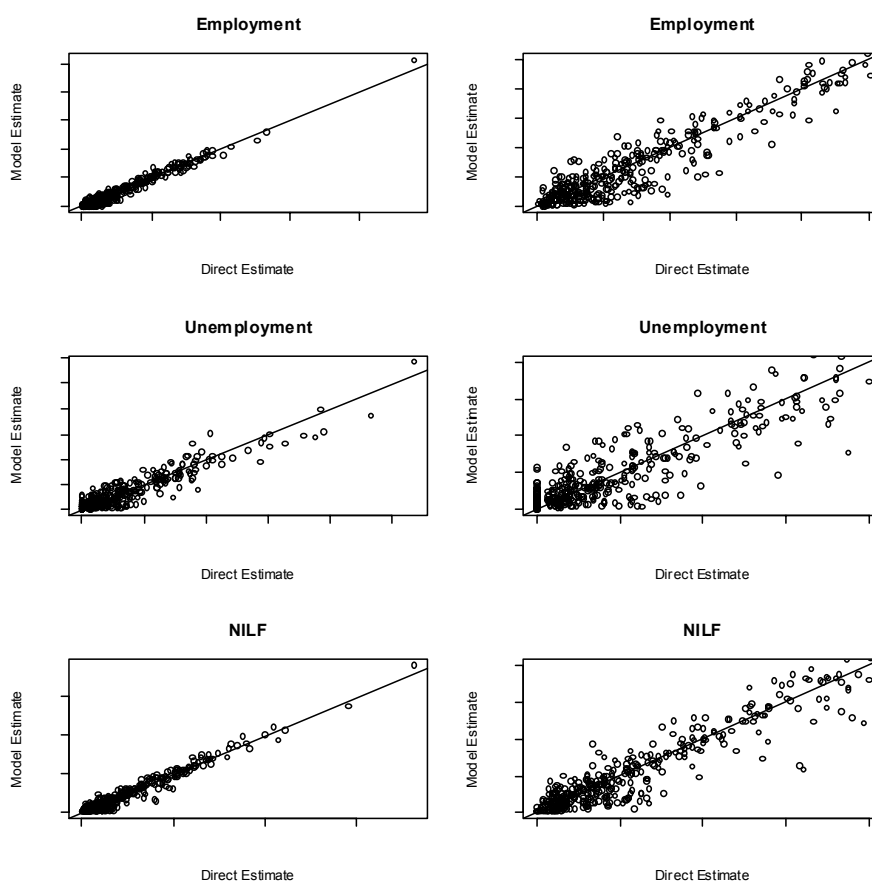


Figure 15.2 contains a comparison between the model based estimates (8.2) and the direct survey estimates for the August 2006 in-sample small areas. Note that the estimates of the NILF small area totals are obtained via subtraction. These plots are useful as a check for bias. For confidentiality reasons we have removed the actual numbers on the plots. The three plots on the left contain all of the estimates, whereas the plots on the right contain only the estimates closer to zero. The line  $y = x$  is also drawn on these plots. Note that the direct survey estimates should be approximately design unbiased but with large standard errors. Figure 15.2 suggests that the model based estimates for Employed and NILF are roughly unbiased since although there is some variation, the estimates are distributed roughly about the line  $y = x$ . The unemployment model based estimates appear to be a little worse in some cases. For instance, when the direct estimates are large, the model based estimate tends to be smaller. However, for the most part, the model based estimators appear roughly unbiased or have a small bias.

### 15.2 Comparison between model based small area estimates and direct survey estimates for 2006



We now consider mean squared error estimation for the estimated small area totals. The mean squared error matrices can be estimated using two methods, either by an analytical approximation or a parametric bootstrap. For further details on the analytical approximation and the parametric bootstrap see Sections 8, 9 and 11. A comparison of the average percentage RRMSEs for 2006 derived from the analytical approximation and the parametric bootstrap with  $B = 1000$  is given in table 15.3. On average the differences between the two methods are very small. The largest average absolute difference is for unemployment and this is only 0.83% and 0.62% for respectively the in-sample and out-of-sample small areas.

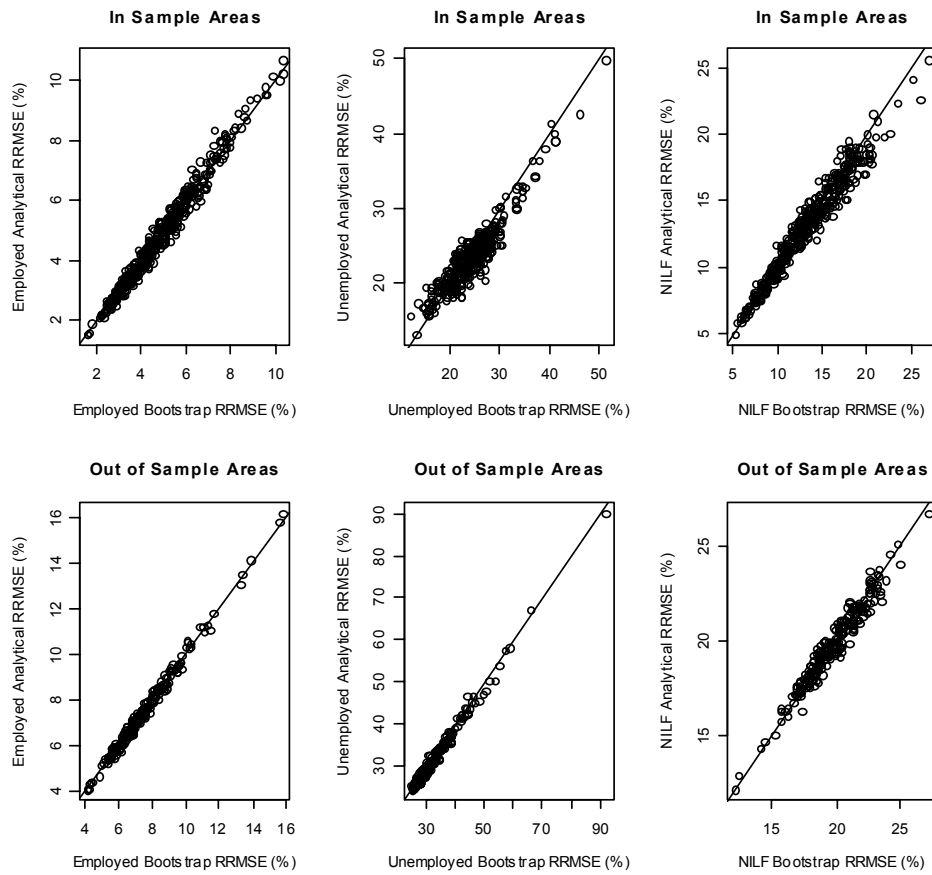
### 15.3 2006 Average RRMSE (%) estimates

<i>LFS status</i>	<i>In-sample areas</i>		<i>Out-of-sample areas</i>	
	<i>Bootstrap</i>	<i>Analytical</i>	<i>Bootstrap</i>	<i>Analytical</i>
Employed	4.84	4.86	7.44	7.47
Unemployed	24.52	23.69	32.30	31.68
Not in labour force	13.34	13.26	19.66	19.72

Figure 15.4 contains an overall comparison between the analytical and parametric bootstrap RRMSEs for August 2006. The line  $y = x$  is also plotted. Figure 15.4 also confirms that on average the RRMSEs from both methods compare well since the values are roughly distributed about the line  $y = x$ . The in-sample unemployment analytical RRMSE estimates appear a little worse than the others on average since the analytical approximation looks to be on average slightly overestimating the smaller RRMSEs and underestimating the larger RRMSEs.

Table 15.5 contains the 2.5th and 97.5th percentiles of the distribution of the differences between the 2006 parametric bootstrap and the analytical RRMSE percentages. Figure 15.4 and table 15.5 shows that there is some variability in the differences between the parametric bootstrap and analytical RRMSEs. However, this variability is not overly large with the worse case being for the in-sample unemployed and the majority of these differences are  $< 4\%$  in absolute value. Note also that some of these larger differences could also be due to extra variation resulting from the parametric bootstrap since  $B = 1000$  is only a moderately large value.

## 15.4 Comparison between analytical and bootstrap RRMSEs



## 15.5 2.5th and 97.5th percentiles of the distribution of the differences between the 2006 parametric bootstrap and analytical RRMSE percentages

	<i>In-sample areas</i>		<i>Out-of-sample areas</i>	
<i>LFS status</i>	<i>2.5th percentile</i>	<i>97.5th percentile</i>	<i>2.5th percentile</i>	<i>97.5th percentile</i>
Employed	-0.50	0.50	-0.37	0.32
Unemployed	-2.54	4.24	-1.03	2.45
Not in labour force	-1.50	2.12	-0.95	1.00

## 16. CONCLUSION

This paper successfully adapted the model and estimators described in Molina *et al.* (2007) to include category specific random effects. We showed that the category specific multinomial logit mixed model is more appropriate for our dataset than the more restrictive one given by Molina *et al.* (2007). The PQL-REML estimation procedure worked very well in our context and we showed via a parametric bootstrap that the PQL-REML estimators had good statistical properties. For instance, the bias in the REML variance component estimates was found to be small.

Similar to Molina *et al.* (2007), we described and derived two different estimators of the mean squared errors of the small area estimated totals. These two different approaches are based on using an analytical approximation and a parametric bootstrap. In the paper by Molina *et al.* (2007), the authors undertook a simulation study and concluded that the bootstrap estimator performed better than the analytical approximation and recommended the bootstrap be used. However they noted that the differences were smaller for the actual UK unemployment data. We showed that for the Australian labour force data that the analytical approximation RRMSEs compared very well with the parametric bootstrap RRMSEs and the differences were all reasonably small. In our context we recommend that the analytical RRMSEs be used because our parametric bootstrap is much more computationally intensive than the one given in Molina *et al.* (2007). We believe the small gains in accuracy will not be worth the extra computational effort involved for the parametric bootstrap in our case.

In this paper we used residual plots and  $\chi^2$  goodness of fit tests to check model assumptions. We used a parametric bootstrap to generate the empirical distributions of the  $\chi^2$  statistic. These tests and plots indicated that the model assumptions appeared to hold approximately for the sample data. There was very slight under and overdispersion present and a couple of small outliers for remote areas. In a future study we might try to improve the model for remote areas. In any case, for the most part the multinomial logit mixed model appears to work reasonably well for modelling the Australian Labour Force count data. Interestingly, the multinomial model has quite a restrictive variance and correlation structure and the fact that the multinomial model works so well here is very convenient. This is because extending the model to account for under and overdispersion in a small area context would not be straight forward. This would certainly be an interesting topic for future research.

Another future research topic could be to try account better for the sample design and any design informativeness. Our estimators like those in the Molina *et al.* (2007) paper essentially assume that the Labour Force sample has been collected using SRSWOR.



## ACKNOWLEDGEMENTS

This work was initially intended to form a small introductory part of my PhD thesis and has since evolved into a separate more detailed joint ABS and ANU research paper. As such, I would like to acknowledge my PhD supervisor Professor Alan Welsh for his technical advice and continual support throughout this work and for putting me onto this topic in the first place.

I would also like to thank my small area estimation colleagues, in particular Daniel Elazar for their support and help and for allowing me to work on this topic part time at the ABS. Also, I would especially like to thank those ABS staff who manipulated the data and got it into a useable form which I could then run the multinomial logit mixed models on.

I would also like to acknowledge DEEWR, the Department of Education, Employment and Workplace Relations, for making their administrative data available to use, and I very much appreciate their effort.

Finally, I would like to acknowledge two other ABS colleagues. Thanks especially to Peter Rossiter for his substantial help with formatting this document. Also, thanks very much to Frank Yu for proof reading this paper and providing some very useful and relevant comments.

## REFERENCES

- Australian Bureau of Statistics (2001) *Australian Standard Geographical Classification (ASGC), 2001*, cat. no. 1216.0, ABS, Canberra.
- (2003) *Information Paper: Census of Population and Housing – Socio-Economic Indexes for Areas, Australia, 2001*, cat. no. 2039.0, ABS, Canberra.
- (2007) *Regional Population Growth*, cat. no. 3218.0, ABS, Canberra.
- Baillo, A. and Molina, I. (2005) *Mean Squared Errors of Small Area Estimators under a Unit-Level Multivariate Model*, Statistics and Econometrics Working Paper ws054007, Universidad Carlos III de Madrid, Madrid. (last viewed 13 January 2010)  
< <http://ideas.repec.org/p/cte/wsrepe/ws054007.html> >
- Booth, J.G. and Hobert, J.P. (1999) “Maximising Generalised Linear Mixed Model Likelihoods with an Automated Monte Carlo EM Algorithm”, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 61(1), pp. 265–285.
- Breslow, N.E. and Clayton, D.G. (1993) “Approximate Inference in Generalized Linear Mixed Models”, *Journal of the American Statistical Association*, 88(421), pp. 9–25.
- Dobson, A.J. (2001) *An Introduction to Generalized Linear Models*, Second Edition, Chapman and Hall/CRC Press, London.
- Geweke, J. (1996) “Monte Carlo Simulation and Numerical Integration”, in H.M. Amman, D.A. Kendrick and J. Rust (eds), *Handbook of Computational Economics, Volume I*, Elsevier Science Publishers B.V., Amsterdam.
- Hartzel, J.; Agresti, A. and Caffo, B. (2001) “Multinomial Logit Random Effects Models”, *Statistical Modelling*, 1, pp. 81–102.
- Harville, D.A. (1977) “Maximum Likelihood Approaches to Variance Component Estimation and to Related Problems”, *Journal of the American Statistical Association*, 72(358), pp. 320–340.
- Henderson, H.V. and Searle, S.R. (1981) “On Deriving the Inverse of a Sum of Matrices”, *SIAM Review*, 23(1), pp. 53–60.
- Jiang, J. (1998) “Consistent Estimators in Generalised Linear Mixed Models”, *Journal of the American Statistical Association*, 93(442), pp. 720–729.
- Jiang, J. (2007) *Linear and Generalized Linear Mixed Models and Their Applications*, Springer Science and Business Media LLC, New York.

- Kackar, R.N. and Harville, D.A. (1984) “Approximations for Standard Errors of Estimators of Fixed and Random Effects in Mixed Linear Models”, *Journal of the American Statistical Association*, 79(388), pp. 853–862.
- Lee, Y. and Nelder, J.A. (1996) “Hierarchical Generalized Linear Models (with Discussion)”, *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(4), pp. 619–678.
- Molina, I.; Saei, A. and Lombardía, M.J. (2007) “Small Area Estimates of Labour Force Participation under a Multinomial Logit Mixed Model”, *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 170(4), pp. 975–1000.
- Pawitan, Y. (2001) *In All Likelihood: Statistical Modelling and Inference Using Likelihood*, Oxford University Press.
- Prasad, N.G.N. and Rao, J.N.K. (1990) “The Estimation of the Mean Squared Error of Small-Area Estimators”, *Journal of the American Statistical Association*, 85(409), pp. 163–171.
- Rao, C.R. and Kleffe, J. (1988) *Estimation of Variance Components and Applications*, Elsevier Science Publishers B.V., Amsterdam.
- Saei, A. and Chambers, R. (2003) *Small Area Estimation Under Linear and Generalized Linear Mixed Models With Time and Area Effects*, S3RI Methodology Working Papers, M03/15, Southampton Statistical Sciences Research Institute. (last viewed 13 January 2010)  
< <http://eprints.soton.ac.uk/8165/> >
- Saei, A. and Chambers, R. (2005) *Empirical Best Linear Unbiased Prediction for Out of Sample Areas*, S3RI Methodology Working Papers, M05/03, Southampton Statistical Sciences Research Institute. (last viewed 13 January 2010)  
< <http://eprints.soton.ac.uk/14073/> >
- Schall, R. (1991) “Estimation in Generalized Linear Models with Random Effects”, *Biometrika*, 78(4), pp. 719–727.

## APPENDIX

This appendix contains the proof of  $E(\mathbf{G}_1(\hat{\boldsymbol{\varphi}})) \approx \mathbf{G}_1(\boldsymbol{\varphi}) - \mathbf{G}_3(\boldsymbol{\varphi})$  (see Section 8).

First we take a second order Taylor series expansion (and assume that  $\hat{\boldsymbol{\Sigma}}$  in  $\hat{\mathbf{T}}$  does not depend on  $\hat{\boldsymbol{\varphi}}$ )

$$\begin{aligned} \hat{\mathbf{T}} &= \mathbf{T}(\hat{\boldsymbol{\varphi}}) \approx \mathbf{T}(\boldsymbol{\varphi}) + \frac{\partial \mathbf{T}}{\partial \varphi_1}(\hat{\varphi}_1 - \varphi_1) + \frac{\partial \mathbf{T}}{\partial \varphi_2}(\hat{\varphi}_2 - \varphi_2) \\ &+ \frac{\partial \mathbf{T}}{\partial \varphi_{12}}(\hat{\varphi}_{12} - \varphi_{12}) + \frac{1}{2} \frac{\partial^2 \mathbf{T}}{\partial \varphi_1^2}(\hat{\varphi}_1 - \varphi_1)^2 + \frac{1}{2} \frac{\partial^2 \mathbf{T}}{\partial \varphi_2^2}(\hat{\varphi}_2 - \varphi_2)^2 \\ &+ \frac{1}{2} \frac{\partial^2 \mathbf{T}}{\partial \varphi_{12}^2}(\hat{\varphi}_{12} - \varphi_{12})^2 + \frac{\partial^2 \mathbf{T}}{\partial \varphi_1 \partial \varphi_2}(\hat{\varphi}_1 - \varphi_1)(\hat{\varphi}_2 - \varphi_2) \\ &+ \frac{\partial^2 \mathbf{T}}{\partial \varphi_1 \partial \varphi_{12}}(\hat{\varphi}_1 - \varphi_1)(\hat{\varphi}_{12} - \varphi_{12}) + \frac{\partial^2 \mathbf{T}}{\partial \varphi_2 \partial \varphi_{12}}(\hat{\varphi}_2 - \varphi_2)(\hat{\varphi}_{12} - \varphi_{12}). \end{aligned}$$

Assuming that  $\mathbf{T}$  is approximately constant and does not depend on  $\mathbf{u}$  (technically  $\mathbf{T}$  depends on  $\mathbf{u}$ ) and  $E(\hat{\boldsymbol{\varphi}}) \approx \boldsymbol{\varphi}$ , then

$$\begin{aligned} E(\mathbf{T}(\hat{\boldsymbol{\varphi}})) &\approx \mathbf{T}(\boldsymbol{\varphi}) + \frac{1}{2} \frac{\partial^2 \mathbf{T}}{\partial \varphi_1^2} E\left((\hat{\varphi}_1 - \varphi_1)^2\right) + \frac{1}{2} \frac{\partial^2 \mathbf{T}}{\partial \varphi_2^2} E\left((\hat{\varphi}_2 - \varphi_2)^2\right) \\ &+ \frac{1}{2} \frac{\partial^2 \mathbf{T}}{\partial \varphi_{12}^2} E\left((\hat{\varphi}_{12} - \varphi_{12})^2\right) + \frac{\partial^2 \mathbf{T}}{\partial \varphi_1 \partial \varphi_2} E\left((\hat{\varphi}_1 - \varphi_1)(\hat{\varphi}_2 - \varphi_2)\right) \\ &+ \frac{\partial^2 \mathbf{T}}{\partial \varphi_1 \partial \varphi_{12}} E\left((\hat{\varphi}_1 - \varphi_1)(\hat{\varphi}_{12} - \varphi_{12})\right) + \frac{\partial^2 \mathbf{T}}{\partial \varphi_2 \partial \varphi_{12}} E\left((\hat{\varphi}_2 - \varphi_2)(\hat{\varphi}_{12} - \varphi_{12})\right). \end{aligned} \quad (\text{A.1})$$

After some algebra and making use of (8.16), (8.18) and (8.19) it can be proved that for  $a = 1, 2, 12$  and  $b = 1, 2, 12$ ,

$$\begin{aligned} \frac{\partial^2 \mathbf{T}}{\partial \varphi_a \partial \varphi_b} &= -\frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_b} - \frac{\partial \mathbf{W}}{\partial \varphi_b} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_a} \\ &+ \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_b} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{Z} \mathbf{W} + \frac{\partial \mathbf{W}}{\partial \varphi_b} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{Z} \mathbf{W} \\ &+ \mathbf{W} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_b} + \mathbf{W} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_b} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_a} \\ &- \mathbf{W} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_b} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{Z} \mathbf{W} \\ &- \mathbf{W} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_b} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{Z} \frac{\partial \mathbf{W}}{\partial \varphi_a} \mathbf{Z}^t \mathbf{V}^{-1} \mathbf{Z} \mathbf{W}. \end{aligned} \quad (\text{A.2})$$

Now, assuming that  $\widehat{\mathbf{M}}_d \approx \mathbf{M}_d$  (i.e.  $\widehat{\mathbf{M}}_d$  is approximately constant), then

$$E(\mathbf{G}_1(\hat{\boldsymbol{\varphi}})) \approx \mathbf{M}_d E(\hat{\mathbf{T}}) \mathbf{M}_d^t = \begin{pmatrix} \mathbf{M}_{d1} E(\hat{\mathbf{T}}) \mathbf{M}_{d1}^t & \mathbf{M}_{d1} E(\hat{\mathbf{T}}) \mathbf{M}_{d2}^t \\ \mathbf{M}_{d2} E(\hat{\mathbf{T}}) \mathbf{M}_{d1}^t & \mathbf{M}_{d2} E(\hat{\mathbf{T}}) \mathbf{M}_{d2}^t \end{pmatrix}. \quad (\text{A.3})$$

Now substitute (A.1) into (A.3) and using (8.12) and (A.2) we obtain

$$\begin{aligned} E(\mathbf{G}_1(\hat{\boldsymbol{\varphi}})) \approx \mathbf{G}_1(\boldsymbol{\varphi}) & - \begin{pmatrix} b_{1,1}^{(1,1)} & b_{1,1}^{(1,2)} \\ b_{1,1}^{(1,2)} & b_{1,1}^{(2,2)} \end{pmatrix} - \begin{pmatrix} b_{2,2}^{(1,1)} & b_{2,2}^{(1,2)} \\ b_{2,2}^{(1,2)} & b_{2,2}^{(2,2)} \end{pmatrix} - \begin{pmatrix} b_{12,12}^{(1,1)} & b_{12,12}^{(1,2)} \\ b_{12,12}^{(1,2)} & b_{12,12}^{(2,2)} \end{pmatrix} \\ & - \begin{pmatrix} b_{1,2}^{(1,1)} & b_{1,2}^{(1,2)} \\ b_{1,2}^{(1,2)} & b_{1,2}^{(2,2)} \end{pmatrix} - \begin{pmatrix} b_{2,1}^{(1,1)} & b_{2,1}^{(1,2)} \\ b_{2,1}^{(1,2)} & b_{2,1}^{(2,2)} \end{pmatrix} - \begin{pmatrix} b_{1,12}^{(1,1)} & b_{1,12}^{(1,2)} \\ b_{1,12}^{(1,2)} & b_{1,12}^{(2,2)} \end{pmatrix} \\ & - \begin{pmatrix} b_{12,1}^{(1,1)} & b_{12,1}^{(1,2)} \\ b_{12,1}^{(1,2)} & b_{12,1}^{(2,2)} \end{pmatrix} - \begin{pmatrix} b_{2,12}^{(1,1)} & b_{2,12}^{(1,2)} \\ b_{2,12}^{(1,2)} & b_{2,12}^{(2,2)} \end{pmatrix} - \begin{pmatrix} b_{12,2}^{(1,1)} & b_{12,2}^{(1,2)} \\ b_{12,2}^{(1,2)} & b_{12,2}^{(2,2)} \end{pmatrix} \end{aligned}$$

where for  $a = 1, 2, 12$ ,  $b = 1, 2, 12$ ,  $j = 1, 2$  and  $k = 1, 2$ ,

$$b_{a,b}^{(j,k)} = \left( \frac{\partial (\mathbf{M}_{dj} \mathbf{WZ}^t \mathbf{V}^{-1})}{\partial \varphi_a} \right) \mathbf{V} \left( \frac{\partial (\mathbf{M}_{dk} \mathbf{WZ}^t \mathbf{V}^{-1})}{\partial \varphi_b} \right)^t E((\hat{\varphi}_a - \varphi_a)(\hat{\varphi}_b - \varphi_b))$$

and therefore

$$\begin{aligned} E(\mathbf{G}_1(\hat{\boldsymbol{\varphi}})) & \approx \mathbf{G}_1(\boldsymbol{\varphi}) - \begin{pmatrix} g_{11}^* & g_{12}^* \\ g_{12}^* & g_{22}^* \end{pmatrix} \\ & = \mathbf{G}_1(\boldsymbol{\varphi}) - \mathbf{G}_3(\boldsymbol{\varphi}). \end{aligned}$$





## FOR MORE INFORMATION . . .

*INTERNET*      **www.abs.gov.au** the ABS website is the best place for data from our publications and information about the ABS.

### INFORMATION AND REFERRAL SERVICE

Our consultants can help you access the full range of information published by the ABS that is available free of charge from our website. Information tailored to your needs can also be requested as a 'user pays' service. Specialists are on hand to help you with analytical or methodological advice.

*PHONE*                      1300 135 070

*EMAIL*                      [client.services@abs.gov.au](mailto:client.services@abs.gov.au)

*FAX*                              1300 135 211

*POST*                            Client Services, ABS, GPO Box 796, Sydney NSW 2001

## FREE ACCESS TO STATISTICS

All statistics on the ABS website can be downloaded free of charge.

*WEB ADDRESS*      [www.abs.gov.au](http://www.abs.gov.au)