

**Asymptotic Results for Penalized Quasi-Likelihood  
Estimation in Generalized Linear Mixed Models**

*The Australian National University*

**Supplementary Material**

The supplementary material is organised as follows. In Section S1, we prove consistency results for the PQL estimator separately for the conditional and unconditional regimes. Section S2 focuses on the distributional results, and Section S3 treats the remainder term in the Taylor expansion. Section S4 contains results for unpartnered fixed effects, for some special cases. Section S5 provides extra simulation results, such as for the conditional regime and different choices of  $\hat{\mathbf{G}}$ .

In the developments, we prove all results below assuming the working dispersion parameter  $\hat{\phi}$  is equal to the true dispersion parameter  $\dot{\phi}$ . Then for the general result using any  $O_p(1)$  working  $\hat{\phi}$ , we note that solving

$$\nabla Q(\hat{\boldsymbol{\theta}}) = \begin{bmatrix} \hat{\phi}^{-1} \mathbf{X}^\top \{\mathbf{y} - \boldsymbol{\mu}(\hat{\boldsymbol{\theta}})\} \\ \hat{\phi}^{-1} \mathbf{Z}^\top \{\mathbf{y} - \boldsymbol{\mu}(\hat{\boldsymbol{\theta}})\} - (\mathbf{I}_m \otimes \hat{\mathbf{G}}^{-1}) \hat{\mathbf{b}} \end{bmatrix} = \mathbf{0}_{(m+1)p}$$

for  $\hat{\boldsymbol{\theta}}$  is equivalent to solving

$$\nabla Q(\hat{\boldsymbol{\theta}}) = \begin{bmatrix} \dot{\phi}^{-1} \mathbf{X}^\top \{\mathbf{y} - \boldsymbol{\mu}(\hat{\boldsymbol{\theta}})\} \\ \dot{\phi}^{-1} \mathbf{Z}^\top \{\mathbf{y} - \boldsymbol{\mu}(\hat{\boldsymbol{\theta}})\} - (\mathbf{I}_m \otimes \hat{\mathbf{G}}_s^{-1}) \hat{\mathbf{b}} \end{bmatrix} = \mathbf{0}_{(m+1)p},$$

where  $\hat{\mathbf{G}}_s = \dot{\phi} \hat{\phi}^{-1} \hat{\mathbf{G}}$ , whose inverse is still  $O_p(1)$  and positive definite. This is equivalent to setting  $\hat{\phi}$  to  $\dot{\phi}$  and scaling  $\hat{\mathbf{G}}$  by  $\dot{\phi} \hat{\phi}^{-1}$ . The general result then follows since the results proved under  $\hat{\phi} = \dot{\phi}$  hold for any  $\hat{\mathbf{G}}$  that has an  $O_p(1)$ , positive definite inverse.

### S0.1 Bias and Identifiability in the Conditional Regime

By differentiating (2.2), we see that the PQL estimators satisfy  $\sum_{i=1}^m \hat{\phi}^{-1} \mathbf{X}_i^\top \{\mathbf{y}_i - \boldsymbol{\mu}_i(\hat{\boldsymbol{\theta}})\} = \mathbf{0}$  and  $\hat{\phi}^{-1} \mathbf{Z}_i^\top \{\mathbf{y}_i - \boldsymbol{\mu}_i(\hat{\boldsymbol{\theta}})\} - \mathbf{G}^{-1} \hat{\mathbf{b}}_i = \mathbf{0}$ ,  $i = 1, \dots, m$ . Summing both sides of the second equation across all  $i$ , since  $\mathbf{X}_i = \mathbf{Z}_i$ , it follows that  $\sum_{i=1}^m \hat{\mathbf{b}}_i = \mathbf{0}_p$ . That is, the PQL estimators of the random effects must satisfy a sum-to-zero constraint regardless of the underlying true parameter values. Under a conditional regime, this induces an asymptotic bias as captured by the term  $\mathbf{1}_m^* \otimes (m^{-1} \sum_{i=1}^m \hat{\mathbf{b}}_i)$  in Theorem 1, which can be interpreted as shifting the mean of the random effects into the corresponding fixed effects. We can deal with the bias by reparametrising the model *a priori* to satisfy a sum-to-zero constraint. That is, we can define a reparametrized vector of true values  $\hat{\boldsymbol{\theta}}^*$  which

satisfy  $\mathbf{1}_m^* \otimes (m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i^*) = \mathbf{0}_{(m+1)p}$ , and the PQL estimator will then be asymptotically normally distributed centered around  $\dot{\boldsymbol{\theta}}^*$ . Furthermore, Theorem 1 remains practically useful as, for any given sample size, we can always reparameterise the GLMM to satisfy this identifiability constraint.

The asymptotic bias discussed above is analogous to that seen in a over-parametrized one-way analysis of variance (ANOVA) model. That is, in the ANOVA model one can always reparametrise to satisfy a sum-to-zero constraint, and the corresponding estimator is consistent for this vector of the reparametrized true values. Note however that when we work unconditionally (Section 4), reparametrising in this way will lead to a different model to the original, since the clusters are no longer independent.

## S1 Proofs for Consistency

To establish our large sample distributional results, we first require the following consistency result.

**Lemma 1.** *Suppose Conditions (C1)-(C5) hold and  $mn_L^{-2} \rightarrow 0$ . Then, as  $m, n_L \rightarrow \infty$  and unconditional on the random effects  $\dot{\mathbf{b}}$ ,  $\|\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}\|_\infty = o_p(1)$ .*

These results are required to control the remainder term in the Taylor expansions we use to derive the distributional results in Section S2. To prove the result, we wish to show that for any given  $\epsilon > 0$ , there exists a

large enough constant  $C > 0$  such that, for large  $m, n_L$ , we have

$$P \left\{ \sup_{\|\mathbf{u}\|_\infty=C} Q(\dot{\boldsymbol{\theta}} + \delta_{m,n_L}^{-1} \mathbf{u}) < Q(\dot{\boldsymbol{\theta}}) \right\} \geq 1 - \epsilon,$$

for some positive, unbounded, monotonically increasing sequence  $\delta_{m,n_L}$ .

The above result implies that with probability tending to one, there exists a local maximum  $\hat{\boldsymbol{\theta}}$  in the ball  $\{\dot{\boldsymbol{\theta}} + \delta_{m,n_L}^{-1} \mathbf{u} : \|\mathbf{u}\|_\infty \leq C\}$  so that  $\|\delta_{m,n_L}(\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})\|_\infty = O_p(1)$ , and thus  $\|\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}\|_\infty = o_p(1)$ .

Consider the difference  $Q(\dot{\boldsymbol{\theta}} + \mathbf{u}) - Q(\dot{\boldsymbol{\theta}})$ . By a Taylor expansion, we obtain

$$Q(\dot{\boldsymbol{\theta}} + \mathbf{u}) - Q(\dot{\boldsymbol{\theta}}) = \mathbf{u}^\top \{\nabla Q(\dot{\boldsymbol{\theta}})\} - 0.5 \mathbf{u}^\top \{-\nabla^2 Q(\bar{\boldsymbol{\theta}})\} \mathbf{u}. \quad (\text{S1.1})$$

where  $\bar{\boldsymbol{\theta}}$  lies on the line segment joining  $\dot{\boldsymbol{\theta}}$  and  $\dot{\boldsymbol{\theta}} + \mathbf{u}$ . If we can prove that (S1.1) is negative as  $m, n_L \rightarrow \infty$  for any choice of  $C$ , then there must exist some  $\delta_{m,n_L}$  such that  $Q(\dot{\boldsymbol{\theta}} + \delta_{m,n_L}^{-1} \mathbf{u}) - Q(\dot{\boldsymbol{\theta}})$  is negative for large enough  $C$ , and the required result follows. We have

$$\begin{aligned} \nabla Q(\dot{\boldsymbol{\theta}}) &= \begin{bmatrix} \dot{\phi}^{-1} \mathbf{X}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}}) \\ \dot{\phi}^{-1} \mathbf{Z}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}}) - (\mathbf{I}_m \otimes \hat{\mathbf{G}}^{-1}) \dot{\mathbf{b}} \end{bmatrix} = \begin{bmatrix} \dot{\phi}^{-1} \mathbf{X}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}}) \\ \dot{\phi}^{-1} \mathbf{Z}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}}) \end{bmatrix} + \begin{bmatrix} \mathbf{0}_p \\ -(\mathbf{I}_m \otimes \hat{\mathbf{G}}^{-1}) \dot{\mathbf{b}} \end{bmatrix} \\ &\triangleq \boldsymbol{\lambda}_1 + \boldsymbol{\lambda}_2, \end{aligned}$$

and

$$\begin{aligned}
 -\nabla^2 Q(\bar{\boldsymbol{\theta}}) &= \begin{bmatrix} \mathbf{X}^\top \bar{\mathbf{W}} \mathbf{X} & \mathbf{X}_1^\top \bar{\mathbf{W}}_1 \mathbf{X}_1 & \cdots & \mathbf{X}_m^\top \bar{\mathbf{W}}_m \mathbf{X}_m \\ \mathbf{X}_1^\top \bar{\mathbf{W}}_1 \mathbf{X}_1 & \mathbf{X}_1^\top \bar{\mathbf{W}}_1 \mathbf{X}_1 + \hat{\mathbf{G}}^{-1} & & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{X}_m^\top \bar{\mathbf{W}}_m \mathbf{X}_m & \mathbf{0} & & \mathbf{X}_m^\top \bar{\mathbf{W}}_m \mathbf{X}_m + \hat{\mathbf{G}}^{-1} \end{bmatrix} \\
 &= \begin{bmatrix} \mathbf{X}^\top \bar{\mathbf{W}} \mathbf{X} & \mathbf{X}_1^\top \bar{\mathbf{W}}_1 \mathbf{X}_1 & \cdots & \mathbf{X}_m^\top \bar{\mathbf{W}}_m \mathbf{X}_m \\ \mathbf{X}_1^\top \bar{\mathbf{W}}_1 \mathbf{X}_1 & \mathbf{X}_1^\top \bar{\mathbf{W}}_1 \mathbf{X}_1 & & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{X}_m^\top \bar{\mathbf{W}}_m \mathbf{X}_m & \mathbf{0} & & \mathbf{X}_m^\top \bar{\mathbf{W}}_m \mathbf{X}_m \end{bmatrix} + \text{blockdiag}(\mathbf{0}_p, \mathbf{I}_m \otimes \hat{\mathbf{G}}^{-1}) \\
 &\triangleq \Gamma_1(\bar{\boldsymbol{\theta}}) + \Gamma_2,
 \end{aligned}$$

where  $\bar{\mathbf{W}}_i = \dot{\phi}^{-1} \text{diag}\{a''(\bar{\eta}_{i1}), \dots, a''(\bar{\eta}_{in_i})\}$  and  $\bar{\mathbf{W}} = \dot{\phi}^{-1} \text{diag}\{a''(\bar{\eta}_{11}), \dots, a''(\bar{\eta}_{mn_m})\}$ .

Also, let  $\Gamma_1(\dot{\boldsymbol{\theta}}) + \Gamma_2$  denote the analogous decomposition of  $-\nabla^2 Q(\dot{\boldsymbol{\theta}})$ . For both the conditional and unconditional regimes, we will prove that the second term is positive and dominates the first. However, the treatment of the terms differs between the two cases, and as such the proofs will need to be dealt with separately. In the following three sections, we will first treat the Poisson pure random intercept example, followed by the more general conditional and unconditional regimes.

Before proceeding, we demonstrate an inequality that is used in the

proofs below. Write  $\mathbf{u} = (\mathbf{u}_1^\top, \mathbf{u}_2^\top)^\top$ ,  $\mathbf{u}_2 = (\mathbf{u}_{21}^\top, \dots, \mathbf{u}_{2m}^\top)^\top$ . First, for any  $\boldsymbol{\theta}$  we have

$$\mathbf{u}^\top \boldsymbol{\Gamma}_1(\boldsymbol{\theta}) \mathbf{u} = \mathbf{u}_1^\top \mathbf{X}^\top \mathbf{W} \mathbf{X} \mathbf{u}_1 + 2\mathbf{u}_1^\top \mathbf{X}^\top \mathbf{W} \mathbf{Z} \mathbf{u}_2 + \mathbf{u}_2^\top \mathbf{Z}^\top \mathbf{W} \mathbf{Z} \mathbf{u}_2 \geq 0.$$

Next, we have

$$\begin{aligned} & \mathbf{u}^\top \boldsymbol{\Gamma}_1(\bar{\boldsymbol{\theta}}) \mathbf{u} - c_0^2 \mathbf{u}^\top \boldsymbol{\Gamma}_1(\boldsymbol{\theta}) \mathbf{u} \\ &= \mathbf{u}_1^\top \mathbf{X}^\top (\bar{\mathbf{W}} - c_0^2 \mathbf{W}) \mathbf{X} \mathbf{u}_1 + 2\mathbf{u}_1^\top \mathbf{X}^\top (\bar{\mathbf{W}} - c_0^2 \mathbf{W}) \mathbf{Z} \mathbf{u}_2 + \mathbf{u}_2^\top \mathbf{Z}^\top (\bar{\mathbf{W}} - c_0^2 \mathbf{W}) \mathbf{Z} \mathbf{u}_2. \end{aligned}$$

If we denote  $\mathbf{W}^* = \bar{\mathbf{W}} - c_0^2 \mathbf{W}$ , then by Condition (C1)  $\mathbf{W}^*$  is a diagonal matrix with non-negative entries as the entries of  $c_0^2 \mathbf{W}$  are upper bounded by the smallest component in  $\bar{\mathbf{W}}$ . Therefore

$$\mathbf{u}^\top \boldsymbol{\Gamma}_1(\bar{\boldsymbol{\theta}}) \mathbf{u} - c_0^2 \mathbf{u}^\top \boldsymbol{\Gamma}_1(\boldsymbol{\theta}) \mathbf{u} = \mathbf{u}_1^\top \mathbf{X}^\top \mathbf{W}^* \mathbf{X} \mathbf{u}_1 + 2\mathbf{u}_1^\top \mathbf{X}^\top \mathbf{W}^* \mathbf{Z} \mathbf{u}_2 + \mathbf{u}_2^\top \mathbf{Z}^\top \mathbf{W}^* \mathbf{Z} \mathbf{u}_2 \geq 0,$$

so that  $\mathbf{u}^\top \boldsymbol{\Gamma}_1(\bar{\boldsymbol{\theta}}) \mathbf{u} \geq c_0^2 \mathbf{u}^\top \boldsymbol{\Gamma}_1(\boldsymbol{\theta}) \mathbf{u}$ . Finally, note that we can choose  $\boldsymbol{\theta} = \dot{\boldsymbol{\theta}}$  or  $\boldsymbol{\theta} = E(\dot{\boldsymbol{\theta}})$  without altering the above argument.

### S1.1 Poisson pure random intercept example

We begin with the Poisson pure random intercept example, which gives insight and covers a case where  $\mathbf{X}_i \neq \mathbf{Z}_i$ . The following result is unconditional on the random effects  $\dot{\mathbf{b}}$ .

**Lemma 2.** *Assume Conditions (C1)-(C5) hold, and let  $mn^{-2} \rightarrow 0$ . Then for the Poisson pure random intercept model, as  $m, n \rightarrow \infty$  and unconditional on the random effects  $\dot{\mathbf{b}}$ , it holds that  $\|\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}\|_\infty = o_p(1)$ .*

*Proof.* Let  $\mathbf{u} = \mathbf{u}_2 = (u_{21}, \dots, u_{2m})^\top$ ,  $\boldsymbol{\theta} = \mathbf{b} = (b_1, \dots, b_m)^\top$ ,  $\hat{\mathbf{G}} = \hat{\sigma}_b^2$  (a scalar),

$$-\nabla^2 Q(\bar{\boldsymbol{\theta}}) = \text{diag}(ne^{\bar{b}_1} + \hat{\sigma}_b^{-2}, \dots, ne^{\bar{b}_m} + \hat{\sigma}_b^{-2}) \equiv \boldsymbol{\Gamma}_1(\bar{\boldsymbol{\theta}}) + \boldsymbol{\Gamma}_2, \text{ and}$$

$$\nabla Q(\dot{\boldsymbol{\theta}}) = \begin{bmatrix} \sum_{j=1}^n (y_{1j} - e^{\dot{b}_1}) - \hat{\sigma}_b^{-2} \dot{b}_1 \\ \vdots \\ \sum_{j=1}^n (y_{mj} - e^{\dot{b}_m}) - \hat{\sigma}_b^{-2} \dot{b}_m \end{bmatrix} \equiv \boldsymbol{\lambda}_1 + \boldsymbol{\lambda}_2.$$

Let  $\mathbf{M} = E\{\text{diag}(ne^{b_1}, \dots, ne^{b_m})\}$ . Then  $\mathbf{M} = \text{Var}\{\dot{\phi}^{-1} \mathbf{Z}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}})\}$ . By Condition (C1),  $c_0^2 \mathbf{u}^\top \mathbf{M} \mathbf{u} \leq \mathbf{u}^\top \boldsymbol{\Gamma}_1(\bar{\boldsymbol{\theta}}) \mathbf{u}$ . Next, let  $\lambda = \dot{\sigma}_b^2 \hat{\sigma}_b^{-2}$ . Then  $\text{Var}(\boldsymbol{\lambda}_2) = \dot{\sigma}_b^2 \hat{\sigma}_b^{-4} \mathbf{I}_m$  and

$$\lambda^{-1} \mathbf{u}_2^\top (\dot{\sigma}_b^2 \hat{\sigma}_b^{-4} \mathbf{I}_m) \mathbf{u}_2 = \mathbf{u}_2^\top (\hat{\sigma}_b^{-2} \mathbf{I}_m) \mathbf{u}_2 = \mathbf{u}^\top \boldsymbol{\Gamma}_2 \mathbf{u}.$$

Finally, by the laws of iterated expectation and variance we have

$$\begin{aligned} E\{\nabla Q(\dot{\boldsymbol{\theta}}) \nabla Q(\dot{\boldsymbol{\theta}})^\top\} &= \text{Var}\{\nabla Q(\dot{\boldsymbol{\theta}})\} \\ &= E[\text{Var}\{\nabla Q(\dot{\boldsymbol{\theta}}) | \dot{\mathbf{b}}\}] + \text{Var}[E\{\nabla Q(\dot{\boldsymbol{\theta}}) | \dot{\mathbf{b}}\}] \\ &= E\{\text{Var}(\boldsymbol{\lambda}_1 | \dot{\mathbf{b}})\} + \text{Var}(\boldsymbol{\lambda}_2) \end{aligned}$$

$$= \text{Var}(\boldsymbol{\lambda}_1) + \text{Var}(\boldsymbol{\lambda}_2).$$

Therefore, we have that

$$\begin{aligned} \mathbf{u}^\top \{-\nabla^2 Q(\bar{\boldsymbol{\theta}})\} \mathbf{u} &\geq \min(\lambda^{-1}, c_0^2) \{\mathbf{u}^\top \mathbf{M} \mathbf{u} + \mathbf{u}_2^\top (\hat{\sigma}_b^2 \hat{\sigma}_b^{-4} \mathbf{I}_m) \mathbf{u}_2\} \\ &= \min(\lambda^{-1}, c_0^2) E\{\nabla Q(\dot{\boldsymbol{\theta}}) \nabla Q(\dot{\boldsymbol{\theta}})^\top\}, \end{aligned}$$

where the latter and hence former term grows at the same rate as  $\{\mathbf{u}^\top \nabla Q(\dot{\boldsymbol{\theta}})\}^2$ .

Since at least one component of  $\mathbf{u}$  equals  $\pm C$ , for any given  $\mathbf{u}$ ,  $\mathbf{u}^\top \{-\nabla^2 Q(\bar{\boldsymbol{\theta}})\} \mathbf{u}$  is at least of order  $O_p(m)$  in probability and hence always dominates.

Since the choice of which  $|u_{2i}| = C$  is arbitrary however, we also need to make sure that the  $m$ th order statistic  $\max_{i \in \{1, \dots, m\}} [\{\sum_{j=1}^n (y_{ij} - e^{\dot{b}_i}) - \hat{\sigma}_b^{-2} \dot{b}_i\} / (ne^{\dot{b}_i} + \hat{\sigma}_b^{-2})]$ , which grows with the dimension, is of order  $o_p(1)$ . We know that the leading term in (4.3b) is  $(\mathbf{Z}^\top \dot{\mathbf{W}} \mathbf{Z})^{-1} \{\dot{\phi}^{-1} \mathbf{Z}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}})\}$  when  $mn^{-1} \rightarrow \infty$ ; for this Poisson random intercept example, up to some smaller order terms, this simplifies to the ratio  $\{\sum_{j=1}^n (y_{ij} - e^{\dot{b}_i}) - \hat{\sigma}_b^{-2} \dot{b}_i\} / (ne^{\dot{b}_i} + \hat{\sigma}_b^{-2})$ .

Intuitively then, proving a result for  $\|\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}\|_\infty$  should involve studying

$$\max_{i \in \{1, \dots, m\}} [\{\sum_{j=1}^n (y_{ij} - e^{\dot{b}_i}) - \hat{\sigma}_b^{-2} \dot{b}_i\} / (ne^{\dot{b}_i} + \hat{\sigma}_b^{-2})].$$

Put another way, consider the set of  $\mathbf{u}$  such that one component of  $\mathbf{u}$  equals  $\pm C$  and zero elsewhere. When  $C$  is the  $i$ th component of  $\mathbf{u}$ , this corresponds to deviating away from  $\dot{\boldsymbol{\theta}}$  in the  $i$ th direction. In this case, we



need  $C\{\sum_{j=1}^n(y_{ij}-e^{\dot{b}_i})-\hat{\sigma}_b^{-2}\dot{b}_i\}$  to be dominated by  $C^2ne^{\dot{b}_i}$  for any  $C$  and all  $m, n$  large enough, i.e.,  $\{\sum_{j=1}^n(y_{ij}-e^{\dot{b}_i})-\hat{\sigma}_b^{-2}\dot{b}_i\}/ne^{\dot{b}_i} = o_p(1)$ . This is indeed true as this ratio is  $O_p(n^{-1/2})$ , since  $\sum_{j=1}^n(y_{ij}-e^{\dot{b}_i})-\hat{\sigma}_b^{-2}\dot{b}_i = O_p(n^{1/2})$  due to conditional independence and Chebyshev's inequality, and  $e^{\dot{b}_i} = O_p(1)$ . However, although the ratio is of order  $O_p(n^{-1/2})$ , for any given  $m, n$  there is still a positive probability that the ratio (a random variable) is greater than one in magnitude. On the other hand, for the consistency argument to hold we need to make sure the ratio is smaller than one in magnitude for all  $m$  directions with probability tending to one, as  $m, n \rightarrow \infty$ . In particular, it is sufficient for the maximum of  $m$  of these ratios to be  $o_p(1)$ : this maximum grows with  $m$ , corresponding to the number of directions we need to bound. Intuitively, this should hold if  $m$  does not grow too fast relative to  $n$ .

Now, Downey (1990) proves that the maximum over  $m$  realisations of independently and identically distributed random variables with a finite  $q$ th moment is  $o_p(m^{1/q})$ . By Condition (C5), the ratio  $n^{1/2}\{\sum_{j=1}^n(y_{ij}-e^{\dot{b}_i})-\hat{\sigma}_b^{-2}\dot{b}_i\}/(ne^{\dot{b}_i}+\hat{\sigma}_b^{-2})$  has finite fourth moments for all  $i$  and  $n$ . Thus, the maximum of these (normalised) ratios over  $m$  clusters is of order  $o_p(m^{1/4})$ . As a result, the maximum ratio of interest is  $o_p(m^{1/4}n^{-1/2})$ . Therefore, when  $mn^{-2} \rightarrow 0$ , there exists  $\delta_{m,n}$  such that we can always choose a large

enough  $C$  for  $\delta_{m,n}^{-1} \mathbf{u}^\top \nabla Q(\dot{\boldsymbol{\theta}})$  to be dominated by  $\delta_{m,n}^{-2} \mathbf{u}^\top \{-\nabla^2 Q(\bar{\boldsymbol{\theta}})\} \mathbf{u}$ , and hence  $\|\delta_{m,n}(\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})\|_\infty = O_p(1)$  as required.

□

To conclude this section, we remark that although  $mn^{-2} \rightarrow 0$  is needed for the consistency and thus distributional result, this is a sufficient condition. Intuitively, in the Poisson pure random effects model there are no fixed parameters to estimate, and the estimate of the random effects for each cluster only depends on observations in that cluster. Thus, the relative rates of  $m$  and  $n$  should not matter for a distributional result concerning a finite subset of the random effects.

## S1.2 Conditional on the Random Effects

In this section, we prove the consistency result under the conditional regime.

In the conditional regime, we assume without loss of generality that  $\sum_{i=1}^m \dot{\mathbf{b}}_i = \mathbf{0}_p$ , recalling that we can always reparametrise the random effect coefficients so this holds.

Let  $\mathbf{M} = \Gamma_1(\dot{\boldsymbol{\theta}})$ . Then  $\mathbf{M} = \text{Var}(\boldsymbol{\lambda}_1 | \dot{\mathbf{b}}) = E(\boldsymbol{\lambda}_1 \boldsymbol{\lambda}_1^\top | \dot{\mathbf{b}})$  since  $E(\boldsymbol{\lambda}_1 | \dot{\mathbf{b}}) = \mathbf{0}_{(m+1)p}$ . By Condition (C1), we have  $c_0^2 \mathbf{u}^\top \mathbf{M} \mathbf{u} \leq \mathbf{u}^\top \Gamma_1(\bar{\boldsymbol{\theta}}) \mathbf{u}$ .

We now consider two cases: the special case when  $\mathbf{u}_1 = -\mathbf{u}_{2i}$  for all  $i$ , and when this is not the case. For the former, we have  $\mathbf{u}^\top \boldsymbol{\lambda}_1 = \mathbf{u}^\top \mathbf{M} \mathbf{u} = 0$ .

Then we must examine  $\mathbf{u}^\top \boldsymbol{\lambda}_2$  and  $\mathbf{u}^\top \boldsymbol{\Gamma}_2 \mathbf{u}$ . In this case, we have  $\mathbf{u}^\top \boldsymbol{\lambda}_2 = \sum_{i=1}^m \mathbf{u}_{2i}^\top \hat{\mathbf{G}}^{-1} \dot{\mathbf{b}}_i = -\mathbf{u}_1^\top \hat{\mathbf{G}}^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i = 0$ , and  $\mathbf{u}^\top \boldsymbol{\Gamma}_2 \mathbf{u} = m \mathbf{u}_1^\top \hat{\mathbf{G}}^{-1} \mathbf{u}_1 > 0$  since  $\hat{\mathbf{G}}$  is a positive definite matrix. Thus the difference (S1.1) is negative for large enough  $m, n_L$  and any choice of constant  $C$ .

Next, consider the case when  $\mathbf{u}_1 = -\mathbf{u}_{2i}$  for all  $i$  does not hold. Under this setting, as  $\boldsymbol{\Gamma}_2$  is a positive semi-definite matrix, we still have  $\mathbf{u}^\top \{-\nabla^2 Q(\bar{\boldsymbol{\theta}})\} \mathbf{u} \geq \mathbf{u}^\top \boldsymbol{\Gamma}_1(\bar{\boldsymbol{\theta}}) \mathbf{u} \geq c_0^2 \mathbf{u}^\top \mathbf{M} \mathbf{u}$ , where the last and hence former terms grow at the same rate as  $(\mathbf{u}^\top \boldsymbol{\lambda}_1)^2$ . Since at least one component of  $\mathbf{u}$  equals  $\pm C$ , by Conditions (C1)-(C3) we have that  $\mathbf{u}^\top \{-\nabla^2 Q(\bar{\boldsymbol{\theta}})\} \mathbf{u}$  is at least of order  $O_p(n_L)$ , and always dominates since  $\mathbf{u}^\top \boldsymbol{\lambda}_2 = O_p(m)$  at most.

Since the choice of  $\mathbf{u}$  is arbitrary, we must take into account the growth rate of the  $m$ th order statistic. That is, for any  $1 \leq k \leq p$ , we need  $\max_{i \in \{1, \dots, m\}} [(\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \{\dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i) - \hat{\mathbf{G}}^{-1} \dot{\mathbf{b}}_i\}]_{[k]} = o_p(1)$ , as per the argument for the Poisson pure random intercept model. Since the responses  $y_{ij}$  are from the exponential family and thus the moment generating function always exists, the maximum is of order  $o_p(m^{1/r} n_L^{-1/2})$  for any  $r \in \mathbb{N}$  (Downey, 1990), and hence  $o_p(1)$  since  $m n_L^{-1} \rightarrow 0$  by taking  $r = 2$ , for example. Note that the first  $p$  components of  $\nabla Q(\dot{\boldsymbol{\theta}})$ , which are associated with the fixed effects, do not need to be bounded in this way

because the dimension is fixed.

### S1.3 Unconditional on the Random Effects

In this section, we prove the consistency result under the unconditional regime. The main differences to the derivation under the conditional regime arise from the treatment of  $\boldsymbol{\lambda}_2$ , and the distribution of  $\mathbf{y}$ . In the unconditional regime it holds that  $\sum_{i=1}^m \dot{\mathbf{b}}_i = O_p(m^{1/2})$ , while in the conditional regime we impose a sum to zero constraint. Furthermore, in the unconditional regime we bound  $\mathbf{u}^\top \boldsymbol{\lambda}_2$  using its variance, while in the conditional regime this is not possible because  $\boldsymbol{\lambda}_2$  is not a random variable. Finally, in the unconditional regime we cannot use the properties of the exponential family to bound the  $m$ th order statistic, instead requiring Condition (C5).

Let  $\mathbf{M} = E\{\boldsymbol{\Gamma}_1(\dot{\boldsymbol{\theta}})\}$ . Then  $\mathbf{M} = \text{Var}(\boldsymbol{\lambda}_1) = E(\boldsymbol{\lambda}_1 \boldsymbol{\lambda}_1^\top)$  since  $E(\boldsymbol{\lambda}_1) = \mathbf{0}_{(m+1)p}$ . By Condition (C1),  $c_0^2 \mathbf{u}^\top \mathbf{M} \mathbf{u} \leq \mathbf{u}^\top \boldsymbol{\Gamma}_1(\bar{\boldsymbol{\theta}}) \mathbf{u}$ .

We consider two cases: the special case when  $\mathbf{u}_1 = -\mathbf{u}_{2i}$  for all  $i$ , and when this is not the case. In the former, we have  $\mathbf{u}^\top \boldsymbol{\lambda}_1 = \mathbf{u}^\top \mathbf{M} \mathbf{u} = 0$ . Thus we must examine  $\mathbf{u}^\top \boldsymbol{\lambda}_2$  and  $\mathbf{u}^\top \boldsymbol{\Gamma}_2 \mathbf{u}$ . In this case, we have  $\mathbf{u}^\top \boldsymbol{\lambda}_2 = \sum_{i=1}^m \mathbf{u}_{2i}^\top \hat{\mathbf{G}}^{-1} \dot{\mathbf{b}}_i = -\mathbf{u}_1^\top \hat{\mathbf{G}}^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i = O_p(m^{1/2})$ , and  $\mathbf{u}^\top \boldsymbol{\Gamma}_2 \mathbf{u} = m \mathbf{u}_1^\top \hat{\mathbf{G}}^{-1} \mathbf{u}_1 > 0$  since  $\hat{\mathbf{G}}$  is a positive definite matrix. Hence the difference (S1.1) is negative for large enough  $m, n_L$ , and any choice of constant  $C$ .

Next, consider the case when  $\mathbf{u}_1 = -\mathbf{u}_{2i}$  for all  $i$  does not hold. Then we still have  $\mathbf{u}^\top \{-\nabla^2 Q(\bar{\boldsymbol{\theta}})\} \mathbf{u} \geq c_0^2 \mathbf{u}^\top \mathbf{M} \mathbf{u}$ . Letting  $\lambda = \lambda_{\max}(\hat{\mathbf{G}}^{-1} \dot{\mathbf{G}} \hat{\mathbf{G}}^{-1}) / \lambda_{\min}(\hat{\mathbf{G}}^{-1})$ , we have

$$\text{Var}(\boldsymbol{\lambda}_2) = \mathbf{I}_m \otimes \hat{\mathbf{G}}^{-1} \dot{\mathbf{G}} \hat{\mathbf{G}}^{-1}$$

and

$$\lambda^{-1} \mathbf{u}_2^\top (\mathbf{I}_m \otimes \hat{\mathbf{G}}^{-1} \dot{\mathbf{G}} \hat{\mathbf{G}}^{-1}) \mathbf{u}_2 \leq \mathbf{u}_2^\top (\mathbf{I}_m \otimes \hat{\mathbf{G}}^{-1}) \mathbf{u}_2 = \mathbf{u}^\top \boldsymbol{\Gamma}_2 \mathbf{u}.$$

Now, by the laws of iterated expectation and variance,

$$\begin{aligned} E\{\nabla Q(\dot{\boldsymbol{\theta}}) \nabla Q(\dot{\boldsymbol{\theta}})^\top\} &= \text{Var}\{\nabla Q(\dot{\boldsymbol{\theta}})\} \\ &= E[\text{Var}\{\nabla Q(\dot{\boldsymbol{\theta}}) | \dot{\mathbf{b}}\}] + \text{Var}[E\{\nabla Q(\dot{\boldsymbol{\theta}}) | \dot{\mathbf{b}}\}] \\ &= E\{\text{Var}(\boldsymbol{\lambda}_1 | \dot{\mathbf{b}})\} + \text{Var}(\boldsymbol{\lambda}_2) \\ &= \text{Var}(\boldsymbol{\lambda}_1) + \text{Var}(\boldsymbol{\lambda}_2). \end{aligned}$$

Thus we have that

$$\begin{aligned} \mathbf{u}^\top \{-\nabla^2 Q(\bar{\boldsymbol{\theta}})\} \mathbf{u} &\geq \min(\lambda^{-1}, c_0^2) \{\mathbf{u}^\top \mathbf{M} \mathbf{u} + \mathbf{u}_2^\top (\mathbf{I}_m \otimes \hat{\mathbf{G}}^{-1} \dot{\mathbf{G}} \hat{\mathbf{G}}^{-1}) \mathbf{u}_2\} \\ &= \min(\lambda^{-1}, c_0^2) \mathbf{u}^\top E\{\nabla Q(\dot{\boldsymbol{\theta}}) \nabla Q(\dot{\boldsymbol{\theta}})^\top\} \mathbf{u}, \end{aligned}$$

where the latter and hence former term grows at the same rate as  $\{\mathbf{u}^\top \nabla Q(\dot{\boldsymbol{\theta}})\}^2$ .

Since at least one component of  $\mathbf{u}$  equals  $\pm C$ , for any given  $\mathbf{u}$  we have that  $\mathbf{u}^\top \{-\nabla^2 Q(\bar{\boldsymbol{\theta}})\} \mathbf{u}$  is at least of order  $O_p(n_L)$  and always dominates.

Since the choice of  $\mathbf{u}$  is arbitrary, we must take into account the growth rate of the  $n$ th order statistic. That is, for any  $1 \leq k \leq p$ , we require  $\max_{i \in \{1, \dots, m\}} [(\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \{\dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i) - \hat{\mathbf{G}}^{-1} \dot{\mathbf{b}}_i\}]_{[k]} = o_p(1)$ , as per the argument for the Poisson pure random intercept model. By Condition (C5), this term is of order  $o_p(m^{1/4} n_L^{-1/2})$ , and hence the result follows. Note that the first  $p$  components of  $\nabla Q(\dot{\boldsymbol{\theta}})$ , which are associated with the fixed effects, do not need to be bounded in this way because the dimension is fixed.

## S2 Proofs of Distributional Results

For both the conditional and unconditional regimes, our proof relies on examining the behaviour of the leading term in the Taylor expansion of the estimating function. Under Conditions (C1) and (C3), we take the Taylor expansion of  $\nabla Q(\hat{\boldsymbol{\theta}})$  around  $\dot{\boldsymbol{\theta}}$  and obtain, as  $m, n_L \rightarrow \infty$ ,

$$\nabla Q(\hat{\boldsymbol{\theta}}) = \mathbf{0}_{(m+1)p} = \nabla Q(\dot{\boldsymbol{\theta}}) + \nabla^2 Q(\dot{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) + \frac{1}{2} \mathbf{R}(\tilde{\boldsymbol{\theta}}), \quad (\text{S2.1})$$

where  $\tilde{\boldsymbol{\theta}}$  is a  $(m+1)p \times (m+1)p$  matrix with each row lying on the line segment between  $\dot{\boldsymbol{\theta}}$  and  $\hat{\boldsymbol{\theta}}$  and  $\mathbf{R}(\tilde{\boldsymbol{\theta}})$  is the remainder term. Rearranging,

we have

$$\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}} = -\{\nabla^2 Q(\dot{\boldsymbol{\theta}})\}^{-1} \nabla Q(\dot{\boldsymbol{\theta}}) - \frac{1}{2} \{\nabla^2 Q(\dot{\boldsymbol{\theta}})\}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}). \quad (\text{S2.2})$$

We show in Section S3 that the remainder term is of smaller order than the leading term and thus negligible in the limit, in both the conditional and unconditional regimes.

From (S2.2), to study the asymptotic behaviour of the PQL estimator we will first apply the blockwise matrix inversion formula to obtain an expression for  $-\{\nabla^2 Q(\dot{\boldsymbol{\theta}})\}^{-1}$ . Using this result, we will then obtain an expression for  $-\{\nabla^2 Q(\dot{\boldsymbol{\theta}})\}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$ , and subsequently study the asymptotic behaviour of each constituent term. Note that since  $\nabla Q(\dot{\boldsymbol{\theta}})$  is a  $(m+1)p$ -vector and  $-\{\nabla^2 Q(\dot{\boldsymbol{\theta}})\}^{-1}$  is a  $(m+1)p \times (m+1)p$  matrix, we cannot simply take their limits as per standard fixed dimension asymptotics. Instead, we must evaluate  $-\{\nabla^2 Q(\dot{\boldsymbol{\theta}})\}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$  as a whole.

We can write

$$\nabla Q(\dot{\boldsymbol{\theta}}) = \begin{bmatrix} \dot{\phi}^{-1} \mathbf{X}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}}) \\ \dot{\phi}^{-1} \mathbf{Z}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}}) - (\mathbf{I}_m \otimes \hat{\mathbf{G}}^{-1}) \dot{\mathbf{b}} \end{bmatrix} = \begin{bmatrix} \dot{\phi}^{-1} \sum_{i=1}^m \sum_{j=1}^{n_i} \mathbf{x}_{ij} (y_{ij} - \dot{\mu}_{ij}) \\ \dot{\phi}^{-1} \sum_{j=1}^{n_1} \mathbf{x}_{1j} (y_{1j} - \dot{\mu}_{1j}) - \hat{\mathbf{G}}^{-1} \dot{\mathbf{b}}_1 \\ \vdots \\ \dot{\phi}^{-1} \sum_{j=1}^{n_m} \mathbf{x}_{mj} (y_{mj} - \dot{\mu}_{mj}) - \hat{\mathbf{G}}^{-1} \dot{\mathbf{b}}_m \end{bmatrix}$$

$$\begin{aligned}
 & \triangleq \begin{bmatrix} \mathbf{S}_1 \\ \mathbf{S}_{21} + \mathbf{S}_{31} \\ \vdots \\ \mathbf{S}_{2m} + \mathbf{S}_{3m} \end{bmatrix} \triangleq \begin{bmatrix} \mathbf{S}_1 \\ \mathbf{S}_4 + \mathbf{S}_5 \end{bmatrix} \triangleq \begin{bmatrix} \mathbf{S}_1 \\ \mathbf{S}_6 \end{bmatrix}, \\
 \mathbf{B}(\dot{\boldsymbol{\theta}}) = -\nabla^2 Q(\dot{\boldsymbol{\theta}}) &= \begin{bmatrix} \mathbf{X}^\top \dot{\mathbf{W}} \mathbf{X} & \mathbf{X}^\top \dot{\mathbf{W}} \mathbf{Z} \\ \mathbf{Z}^\top \dot{\mathbf{W}} \mathbf{X} & \mathbf{Z}^\top \dot{\mathbf{W}} \mathbf{Z} + \mathbf{I}_m \otimes \hat{\mathbf{G}}^{-1} \end{bmatrix} \triangleq \begin{bmatrix} \mathbf{B}_1 & \mathbf{B}_2 \\ \mathbf{B}_2^\top & \mathbf{B}_3 + \mathbf{B}_4 \end{bmatrix}.
 \end{aligned}$$

Letting  $\mathbf{C} = \mathbf{B}_1 - \mathbf{B}_2(\mathbf{B}_3 + \mathbf{B}_4)^{-1}\mathbf{B}_2^\top$ , by the matrix block inversion formula

we have

$$\mathbf{B}^{-1} = \begin{bmatrix} \mathbf{C}^{-1} & -\mathbf{C}^{-1}\mathbf{B}_2(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \\ -(\mathbf{B}_3 + \mathbf{B}_4)^{-1}\mathbf{B}_2^\top\mathbf{C}^{-1} & (\mathbf{B}_3 + \mathbf{B}_4)^{-1} + (\mathbf{B}_3 + \mathbf{B}_4)^{-1}\mathbf{B}_2^\top\mathbf{C}^{-1}\mathbf{B}_2(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \end{bmatrix}. \tag{S2.3}$$

Next, based on the forms of  $\mathbf{B}_2$  and  $(\mathbf{B}_3 + \mathbf{B}_4)$ , we obtain

$$\mathbf{B}_2(\mathbf{B}_3 + \mathbf{B}_4)^{-1} = [\mathbf{I}_p - \hat{\mathbf{G}}^{-1}(\mathbf{X}_1^\top \dot{\mathbf{W}}_1 \mathbf{X}_1 + \hat{\mathbf{G}}^{-1})^{-1}, \dots, \mathbf{I}_p - \hat{\mathbf{G}}^{-1}(\mathbf{X}_m^\top \dot{\mathbf{W}}_m \mathbf{X}_m + \hat{\mathbf{G}}^{-1})^{-1}].$$

Then since  $\mathbf{Z}_i = \mathbf{X}_i$  for all  $i$ , we can show that

$$\begin{aligned}
 \mathbf{B}_2(\mathbf{B}_3 + \mathbf{B}_4)^{-1}\mathbf{B}_2^\top &= \sum_{i=1}^m \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i \\
 &= \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1} - \hat{\mathbf{G}}^{-1}) (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i
 \end{aligned}$$



$$\begin{aligned}
&= \sum_{i=1}^m \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i - \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i \\
&= \mathbf{B}_1 - \sum_{i=1}^m \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i.
\end{aligned}$$

It follows that

$$\mathbf{C} = \sum_{i=1}^m \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i = \sum_{i=1}^m \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \hat{\mathbf{G}}^{-1}, \tag{S2.4}$$

where the second equality arises from the fact that as a covariance matrix,

$\mathbf{C}$  must be symmetric. We may also write  $\mathbf{C}$  as

$$\begin{aligned}
\sum_{i=1}^m \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i &= \sum_{i=1}^m \{ \mathbf{I}_p - \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \} \hat{\mathbf{G}}^{-1} \\
&= \hat{\mathbf{G}}^{-1} \sum_{i=1}^m \{ \mathbf{I}_p - (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \hat{\mathbf{G}}^{-1} \}.
\end{aligned} \tag{S2.5}$$

Note that  $\mathbf{C}$  is of order  $O_p(m)$  component-wise in probability in both the conditional and unconditional regimes. Using the fact that  $\mathbf{C}^{-1}$  must also be symmetric, we obtain

$$\mathbf{C}^{-1} = \left\{ \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i \right\}^{-1} \hat{\mathbf{G}} = \hat{\mathbf{G}} \left\{ \sum_{i=1}^m \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \right\}^{-1} \tag{S2.6}$$

or equivalently

$$\begin{aligned} \mathbf{C}^{-1} &= \mathbf{C}^{-1\top} = \left[ \sum_{i=1}^m \{ \mathbf{I}_p - (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \hat{\mathbf{G}}^{-1} \} \right]^{-1} \hat{\mathbf{G}} \\ &= \left\{ m^{-1} \mathbf{I}_p + m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \hat{\mathbf{G}}^{-1} \mathbf{C}^{-1} \right\} \hat{\mathbf{G}}, \quad (\text{S2.7}) \end{aligned}$$

where the last line is derived from (a special case of) the Woodbury identity, given by  $(\mathbf{Q} - \mathbf{R})^{-1} = \mathbf{Q}^{-1} + \mathbf{Q}^{-1} \mathbf{R} (\mathbf{Q} - \mathbf{R})^{-1}$  for arbitrary matrices  $\mathbf{Q}$  and  $\mathbf{R}$  such that  $\mathbf{Q}$  and  $(\mathbf{Q} - \mathbf{R})$  are invertible. The first term in (S2.7) is the dominating term, being of order  $O(m^{-1})$ , while the second term is  $O_p(m^{-1}n_L^{-1})$  in both the conditional and unconditional regimes. We will use all the above forms of  $\mathbf{C}$  and  $\mathbf{C}^{-1}$  in subsequent developments. Similarly, we can apply the Woodbury identity to  $(\mathbf{B}_3 + \mathbf{B}_4)^{-1}$  and  $(\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1}$  to obtain  $n_L(\mathbf{B}_3 + \mathbf{B}_4)^{-1} = n_L \mathbf{B}_3^{-1} - n_L \mathbf{B}_3^{-1} \mathbf{B}_4 (\mathbf{B}_3 + \mathbf{B}_4)^{-1} = O_p(1) + O_p(n_L^{-1})$  and  $n_i(\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} = n_i(\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} - n_i(\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} = O_p(1) + O_p(n_i^{-1})$ , where the order results hold component-wise. These hold irrespective of whether we are conditioning on the random effects.

To further simplify expressions, for the rest of this article we will only use order results when representing quantities associated with these smaller order terms. Furthermore, as we want the derivations for the remainder of

this section to be applicable to both the conditional and unconditional regime, we will not distinguish between  $O(\cdot)$  and  $O_p(\cdot)$  in the following developments, and simply use  $O_p(\cdot)$  to represent both as appropriate. The terms we use “big-O notation” for will have the same order under both the conditional and unconditional regime. To simplify expressions, we will also drop the dependence on  $\boldsymbol{\theta}$ , unless stated otherwise.

Finally, it is worth emphasising that

$$[-\mathbf{I}_p, \mathbf{I}_p, \dots, \mathbf{I}_p] \begin{bmatrix} \dot{\phi}^{-1} \mathbf{X}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}}) \\ \dot{\phi}^{-1} \mathbf{Z}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}}) \end{bmatrix} = -\mathbf{S}_1 + \sum_{i=1}^m \mathbf{S}_{2i} = \mathbf{S}_1 - \sum_{i=1}^m \mathbf{S}_{2i} = \mathbf{0}_p, \quad (\text{S2.8})$$

due to the  $\mathbf{X}_i = \mathbf{Z}_i$  assumption. This is a key identity that is critical to the proofs throughout this article.

We now use the expressions above to multiply out  $-\{\nabla^2 Q(\dot{\boldsymbol{\theta}})\}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$  and obtain expressions for  $\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}}$  and  $\hat{\mathbf{b}} - \dot{\mathbf{b}}$ . From equation (S2.2), the first  $p$  components of  $\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}$  are

$$\begin{aligned} \hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}} &= \left[ \mathbf{C}^{-1} \quad -\mathbf{C}^{-1} \mathbf{B}_2 (\mathbf{B}_3 + \mathbf{B}_4)^{-1} \right] \nabla Q + \frac{1}{2} \{ \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \}_{[1:p]} \\ &= \mathbf{C}^{-1} \left[ \mathbf{I}_p \quad -[\mathbf{I}_p - \hat{\mathbf{G}}^{-1} (\mathbf{X}_1^\top \dot{\mathbf{W}}_1 \mathbf{X}_1 + \hat{\mathbf{G}}^{-1})^{-1}, \dots, \mathbf{I}_p - \hat{\mathbf{G}}^{-1} (\mathbf{X}_m^\top \dot{\mathbf{W}}_m \mathbf{X}_m + \hat{\mathbf{G}}^{-1})^{-1}] \right] \nabla Q \\ &\quad + \frac{1}{2} \{ \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \}_{[1:p]} \end{aligned}$$

$$\begin{aligned}
 &= \mathbf{C}^{-1} \left( \mathbf{S}_1 - \sum_{i=1}^m \mathbf{S}_{2i} - \sum_{i=1}^m \mathbf{S}_{3i} \right) + \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} \left\{ \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{S}_{2i} \right. \\
 &+ \left. \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{S}_{3i} \right\} + \frac{1}{2} \{ \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \}_{[1:p]} \\
 &= \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} \left\{ \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{S}_{2i} - \hat{\mathbf{G}} \sum_{i=1}^m \mathbf{S}_{3i} \right. \\
 &+ \left. \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{S}_{3i} \right\} + \frac{1}{2} \{ \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \}_{[1:p]},
 \end{aligned}$$

where the final equality uses equation (S2.8). Thus, letting  $\mathbf{V}_1 = \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{S}_{2i} - \hat{\mathbf{G}} \sum_{i=1}^m \mathbf{S}_{3i} + \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{S}_{3i}$  and applying equation (S2.7), we obtain

$$\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}} = m^{-1} \mathbf{V}_1 + \frac{1}{2} \{ \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \}_{[1:p]} + O_p(n_L^{-1}) \times m^{-1} \mathbf{V}_1.$$

Finally, using the Woodbury identity for  $(\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1}$ , we have that

$$\sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{S}_{2i} = \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} \mathbf{S}_{2i} + \sum_{i=1}^m O_p(n_L^{-2}) \mathbf{S}_{2i}.$$

Letting  $\mathbf{V}_2 = \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} \mathbf{S}_{2i} - \hat{\mathbf{G}} \sum_{i=1}^m \mathbf{S}_{3i} + \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{S}_{3i}$ , we obtain

$$\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}} = m^{-1} \mathbf{V}_2 + \frac{1}{2} \{ \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \}_{[1:p]} + O_p(n_L^{-1}) \times m^{-1} \mathbf{V}_1 + m^{-1} \sum_{i=1}^m O_p(n_L^{-2}) \mathbf{S}_{2i}.$$

Next, the last  $mp$  components of  $\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}$  are

$$\hat{\mathbf{b}} - \dot{\mathbf{b}} = [-(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{B}_2^\top \mathbf{C}^{-1} \quad (\mathbf{B}_3 + \mathbf{B}_4)^{-1} + (\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{B}_2^\top \mathbf{C}^{-1} \mathbf{B}_2 (\mathbf{B}_3 + \mathbf{B}_4)^{-1}] \nabla Q$$

$$\begin{aligned}
& + \frac{1}{2} \{ \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \}_{[p+1:(m+1)p]} \\
& = [ -(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{B}_2^\top \mathbf{C}^{-1} \quad (\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{B}_2^\top \mathbf{C}^{-1} \mathbf{B}_2 (\mathbf{B}_3 + \mathbf{B}_4)^{-1} ] \nabla Q \\
& + [ \mathbf{0}_{mp \times p} \quad (\mathbf{B}_3 + \mathbf{B}_4)^{-1} ] \nabla Q + \frac{1}{2} \{ \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \}_{[p+1:(m+1)p]} \\
& = -(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{B}_2^\top [ \mathbf{C}^{-1} \quad - \mathbf{C}^{-1} \mathbf{B}_2 (\mathbf{B}_3 + \mathbf{B}_4)^{-1} ] \nabla Q \\
& + [ \mathbf{0}_{mp \times p} \quad (\mathbf{B}_3 + \mathbf{B}_4)^{-1} ] \nabla Q + \frac{1}{2} \{ \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \}_{[p+1:(m+1)p]}.
\end{aligned}$$

Notice that we already have an expression for  $[ \mathbf{C}^{-1} \quad - \mathbf{C}^{-1} \mathbf{B}_2 (\mathbf{B}_3 + \mathbf{B}_4)^{-1} ] \nabla Q$  from the fixed effects above. Namely, it is  $m^{-1} \mathbf{V}_1 + O_p(n_L^{-1}) \times m^{-1} \mathbf{V}_1$ . Thus we have

$$\begin{aligned}
\hat{\mathbf{b}} - \dot{\mathbf{b}} & = -(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{B}_2^\top (m^{-1} \mathbf{V}_1 + O_p(n_L^{-1}) \times m^{-1} \mathbf{V}_1) \\
& + (\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{S}_6 + \frac{1}{2} \{ \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \}_{[p+1:(m+1)p]}.
\end{aligned}$$

Applying the Woodbury identity for  $(\mathbf{B}_3 + \mathbf{B}_4)^{-1}$ , we obtain

$$\begin{aligned}
\hat{\mathbf{b}} - \dot{\mathbf{b}} & = -\mathbf{1}_m \otimes (m^{-1} \mathbf{V}_1 + O_p(n_L^{-1}) \times m^{-1} \mathbf{V}_1) + O_p(n_L^{-1}) (m^{-1} \mathbf{V}_1 + O_p(n_L^{-1}) \times m^{-1} \mathbf{V}_1) \\
& + \mathbf{B}_3^{-1} \mathbf{S}_6 + O_p(n_L^{-2}) \mathbf{S}_6 + \frac{1}{2} \{ \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \}_{[p+1:(m+1)p]} \\
& = -\mathbf{1}_m \otimes m^{-1} \mathbf{V}_1 + O_p(n_L^{-1}) \times m^{-1} \mathbf{V}_1 + O_p(n_L^{-2}) \times m^{-1} \mathbf{V}_1 \\
& + \mathbf{B}_3^{-1} \mathbf{S}_4 + \mathbf{B}_3^{-1} \mathbf{S}_5 + O_p(n_L^{-2}) \mathbf{S}_6 + \frac{1}{2} \{ \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \}_{[p+1:(m+1)p]}.
\end{aligned}$$

Replacing all the  $\mathbf{V}$  and  $\mathbf{S}$  terms in the above with their definitions, we

finally obtain

$$\begin{aligned}
 \hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}} &= m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} \dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i) + m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i \\
 &\quad - m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \hat{\mathbf{G}}^{-1} \dot{\mathbf{b}}_i + \frac{1}{2} \{\mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}})\}_{[1:p]} \\
 &\quad + O_p(n_L^{-1}) \left\{ m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i) + m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i \right. \\
 &\quad \left. - m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \hat{\mathbf{G}}^{-1} \dot{\mathbf{b}}_i \right\} + m^{-1} \sum_{i=1}^m O_p(n_L^{-2}) \dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i),
 \end{aligned} \tag{S2.9}$$

and

$$\begin{aligned}
 \hat{\mathbf{b}} - \dot{\mathbf{b}} &= -\mathbf{1}_m \otimes \left\{ m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i) \right. \\
 &\quad \left. + m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i - m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \hat{\mathbf{G}}^{-1} \dot{\mathbf{b}}_i \right\} \\
 &\quad + \mathbf{B}_3^{-1} \{\dot{\phi}^{-1} \mathbf{Z}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}})\} - \mathbf{B}_3^{-1} \{(\mathbf{I}_m \otimes \hat{\mathbf{G}}^{-1}) \dot{\mathbf{b}}\} + \frac{1}{2} \{\mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}})\}_{[p+1:(m+1)p]} \\
 &\quad + O_p(n_L^{-1}) \left\{ m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i) \right. \\
 &\quad \left. + m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i - m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \hat{\mathbf{G}}^{-1} \dot{\mathbf{b}}_i \right\} \\
 &\quad + O_p(n_L^{-2}) \left\{ m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i) \right. \\
 &\quad \left. + m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i - m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \hat{\mathbf{G}}^{-1} \dot{\mathbf{b}}_i \right\}
 \end{aligned}$$

$$+ O_p(n_L^{-2})\{\dot{\phi}^{-1}\mathbf{Z}^\top(\mathbf{y} - \dot{\boldsymbol{\mu}}) - (\mathbf{I}_m \otimes \hat{\mathbf{G}}^{-1})\dot{\mathbf{b}}\}. \quad (\text{S2.10})$$

The expressions for  $\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}}$  and  $\hat{\mathbf{b}} - \dot{\mathbf{b}}$  above underlie our proofs. We use these same expressions in both the conditional and unconditional regimes, but the asymptotic behaviours of the terms on the right hand side, and the way we treat them, will differ greatly between the two cases.

As we will show later, the key leading terms for the fixed effects are  $m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} \dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i)$  and  $m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i$ . The key leading terms for the random effects are  $-\mathbf{1}_m \otimes m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i$  and  $\mathbf{B}_3^{-1} \{\dot{\phi}^{-1} \mathbf{Z}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}})\}$ . When conditioning on the random effects  $\dot{\mathbf{b}}$ , we have  $m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i = O(1)$ , while in the unconditional regime the same quantity is of order  $O_p(m^{-1/2})$  in probability. In both the conditional and unconditional regimes, we have that  $m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} \dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i)$  is of order  $O_p(N^{-1/2})$  component-wise, while the quantity  $\mathbf{B}_3^{-1} \{\dot{\phi}^{-1} \mathbf{Z}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}})\}$  is of order  $O_p(n_L^{-1/2})$  component-wise.

### S2.1 Proof of Theorem 1

The dominating terms on the right hand sides of equations (S2.9) and (S2.10) are  $m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i$  and  $\mathbf{1}_m \otimes m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i$  for the fixed and random effects, respectively. Conditional on the random effects  $\dot{\mathbf{b}}_i$ , these dominat-

ing terms are deterministic and of order  $O(1)$ . Thus we treat them as bias terms and move them to the left hand side. Next, by Conditions (C1)-(C2),  $\mathbf{B}_3^{-1}$  is a component-wise  $O(n_L^{-1})$  block-diagonal matrix, while we also have  $\mathbf{B}_2^\top = O(n_U)$ ,  $\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i^\top = O(n_i)$ , and  $\mathbf{C}^{-1} = O(m^{-1})$  component-wise. Since  $E\{\mathbf{Z}^\top(\mathbf{y} - \dot{\boldsymbol{\mu}})|\dot{\mathbf{b}}\} = \mathbf{0}_{mp}$  and  $\text{Var}\{\mathbf{Z}^\top(\mathbf{y} - \dot{\boldsymbol{\mu}})|\dot{\mathbf{b}}\} = \mathbf{Z}^\top \dot{\mathbf{W}} \mathbf{Z}$ , we obtain  $\dot{\phi}^{-1} \mathbf{D}_r^{-1} \mathbf{Z}^\top(\mathbf{y} - \dot{\boldsymbol{\mu}}) = O_p(1)$  using Chebyshev's inequality and the conditional independence.

Multiplying both sides of (S2.9) and (S2.10) by  $N^{1/2}$  and  $\mathbf{D}_r$  respectively, and applying the order results for the remainder term in Section S3.1, we obtain

$$N^{1/2} \left( \hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}} - m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i \right) = m^{-1/2} \sum_{i=1}^m n^{1/2} n_i^{-1/2} (n_i^{-1} \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} n_i^{-1/2} \dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i) + O_p(m^{1/2} n_L^{-1/2}),$$

and

$$\mathbf{D}_r \left( \hat{\mathbf{b}} - \dot{\mathbf{b}} + \mathbf{1}_m \otimes m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i \right) = \mathbf{D}_r \mathbf{B}_3^{-1} \mathbf{D}_r \mathbf{D}_r^{-1} \{ \dot{\phi}^{-1} \mathbf{Z}^\top(\mathbf{y} - \dot{\boldsymbol{\mu}}) \} + O_p(n_L^{-1/2}).$$

Recalling that  $\mathbf{X}_i = \mathbf{Z}_i$ , to prove Theorem 1 we will show a Lindeberg



condition for

$$\mathbf{A} \begin{bmatrix} m^{-1/2} \sum_{i=1}^m n^{1/2} n_i^{-1/2} (n_i^{-1} \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} n_i^{-1/2} \dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i) \\ (n_1^{-1} \mathbf{X}_1^\top \dot{\mathbf{W}}_1 \mathbf{X}_1)^{-1} \{n_1^{-1/2} \dot{\phi}^{-1} \mathbf{X}_1^\top (\mathbf{y}_1 - \dot{\boldsymbol{\mu}}_1)\} \\ \vdots \\ (n_m^{-1} \mathbf{X}_m^\top \dot{\mathbf{W}}_m \mathbf{X}_m)^{-1} \{n_m^{-1/2} \dot{\phi}^{-1} \mathbf{X}_m^\top (\mathbf{y}_m - \dot{\boldsymbol{\mu}}_m)\} \end{bmatrix} =: S,$$

and thus apply the Lindeberg-Feller central limit theorem, from which the result follows from Slutsky's theorem.

To prove the condition, first define  $\mathbf{U} = [\mathbf{Z}\mathbf{B}_3^{-1}(\mathbf{1}_m \otimes \mathbf{I}_p), \mathbf{Z}\mathbf{B}_3^{-1}]$ , and  $\mathbf{U}_k$  as the  $k$ th row of  $\mathbf{U}$ , noting it only has  $2p$  non-zero components. Then we can write  $S = \sum_{k=1}^N \mathbf{A} \mathbf{D} \mathbf{U}_k \dot{\phi}^{-1} \{y_k - \mu_k(\dot{\boldsymbol{\theta}})\} \triangleq \sum_{k=1}^N \boldsymbol{\xi}_k$ , where  $y_k$  is the  $k$ th component in  $(y_{11}, y_{12}, \dots, y_{1n_1}, y_{21}, \dots, y_{mn_m})^\top$ , and similarly for  $\mu_k(\dot{\boldsymbol{\theta}})$ .

Conditional on  $\dot{\mathbf{b}}$ , the quantities  $\{\boldsymbol{\xi}_k\}_{k=1}^N$  are independent  $q$ -vectors with expectation zero and covariance  $\text{Var}(\boldsymbol{\xi}_k | \dot{\mathbf{b}}) = \mathbf{A} \mathbf{D} \mathbf{U}_k \mathbf{W}_k \mathbf{U}_k^\top \mathbf{D} \mathbf{A}^\top$ , where  $\mathbf{W}_k$  is the  $k$ th diagonal component in  $\dot{\mathbf{W}}$ . Therefore, we have that

$$\sum_{k=1}^N \text{Var}(\boldsymbol{\xi}_k | \dot{\mathbf{b}}) = \sum_{k=1}^N \mathbf{A} \mathbf{D} \mathbf{U}_k \mathbf{W}_k \mathbf{U}_k^\top \mathbf{D} \mathbf{A}^\top$$

$$= \mathbf{A} \begin{bmatrix} \frac{1}{m} \sum_{i=1}^m \frac{n}{n_i} \left( \frac{\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i}{n_i} \right)^{-1} & \frac{1}{\sqrt{m}} \sqrt{\frac{n}{n_1}} \left( \frac{\mathbf{X}_1^\top \dot{\mathbf{W}}_1 \mathbf{X}_1}{n_1} \right)^{-1} & \cdots & \frac{1}{\sqrt{m}} \sqrt{\frac{n}{n_m}} \left( \frac{\mathbf{X}_m^\top \dot{\mathbf{W}}_m \mathbf{X}_m}{n_m} \right)^{-1} \\ \frac{1}{\sqrt{m}} \sqrt{\frac{n}{n_1}} \left( \frac{\mathbf{X}_1^\top \dot{\mathbf{W}}_1 \mathbf{X}_1}{n_1} \right)^{-1} & \left( \frac{\mathbf{X}_1^\top \dot{\mathbf{W}}_1 \mathbf{X}_1}{n_1} \right)^{-1} & \mathbf{0} & \mathbf{0} \\ \vdots & \mathbf{0} & \ddots & \mathbf{0} \\ \frac{1}{\sqrt{m}} \sqrt{\frac{n}{n_m}} \left( \frac{\mathbf{X}_m^\top \dot{\mathbf{W}}_m \mathbf{X}_m}{n_m} \right)^{-1} & \mathbf{0} & \mathbf{0} & \left( \frac{\mathbf{X}_m^\top \dot{\mathbf{W}}_m \mathbf{X}_m}{n_m} \right)^{-1} \end{bmatrix} \mathbf{A}^\top.$$

Hence using the finite selection property of  $\mathbf{A}$ , and the fact that  $m^{-1/2} n^{1/2} n_i^{-1/2} \left( n_i^{-1} \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i \right)^{-1} = o(1)$  component-wise, we obtain

$$\begin{aligned} & \lim_{m, n_L \rightarrow \infty} \sum_{k=1}^N \text{Cov}(\boldsymbol{\xi}_k | \dot{\mathbf{b}}) \\ &= \lim_{m, n_L \rightarrow \infty} \mathbf{A} \text{bdiag} \left\{ \frac{1}{m} \sum_{i=1}^m \frac{n}{n_i} \left( \frac{\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i}{n_i} \right)^{-1}, \left( \frac{\mathbf{X}_1^\top \dot{\mathbf{W}}_1 \mathbf{X}_1}{n_1} \right)^{-1}, \dots, \left( \frac{\mathbf{X}_m^\top \dot{\mathbf{W}}_m \mathbf{X}_m}{n_m} \right)^{-1} \right\} \mathbf{A}^\top \\ &= \boldsymbol{\Omega}. \end{aligned}$$

Next, by the Cauchy-Schwarz inequality, we have

$$E\{\|\boldsymbol{\xi}_k\|^2 I(\|\boldsymbol{\xi}_k\| > \epsilon) | \dot{\mathbf{b}}\} \leq E(\|\boldsymbol{\xi}_k\|^4 | \dot{\mathbf{b}})^{1/2} P(\|\boldsymbol{\xi}_k\| > \epsilon | \dot{\mathbf{b}})^{1/2}.$$

Finally, we make a note about the form of  $\text{Cov}[\mathbf{DU}_k\{y_k - \mu_k(\boldsymbol{\theta})\}]$ . Without loss of generality, suppose  $k = 1$ . Then

$$\text{Cov}[\mathbf{DU}_1\{y_1 - \mu_1(\boldsymbol{\theta})\}] =$$

$$\begin{bmatrix} n(mn_1^2)^{-1}\mathbf{H}_1\mathbf{x}_{11}W_1\mathbf{x}_{11}^\top\mathbf{H}_1^\top & n_1^{-1}m^{-1/2}(nn_1^{-1})^{1/2}\mathbf{H}_1\mathbf{x}_{11}W_1\mathbf{x}_{11}^\top\mathbf{H}_1^\top & \mathbf{0} \\ n_1^{-1}m^{-1/2}(nn_1^{-1})^{1/2}\mathbf{H}_1\mathbf{x}_{11}W_1\mathbf{x}_{11}^\top\mathbf{H}_1^\top & n_1^{-1}\mathbf{H}_1\mathbf{x}_{11}W_1\mathbf{x}_{11}^\top\mathbf{H}_1^\top & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}. \quad (\text{S2.11})$$

Again without loss of generality, consider the case  $\mathbf{A} = [\mathbf{I}_{2p}, \mathbf{0}_{(p+p) \times (m-1)p}]$ . Then by equation (S2.11) and Chebyshev's inequality, when  $k \in \{1, 2, \dots, n_1\}$  we have that  $P(\|\boldsymbol{\xi}_k\| > \epsilon|\dot{\mathbf{b}}) \leq \text{tr}\{\text{Cov}(\boldsymbol{\xi}_k|\dot{\mathbf{b}})\}/\epsilon^2 = O(n_1^{-1})$ . Thus, given  $\dot{\mathbf{b}}$ , we obtain  $\|\boldsymbol{\xi}_k\| = O_p(n_1^{-1/2})$  and  $E(\|\boldsymbol{\xi}_k\|^4|\dot{\mathbf{b}}) = O(n_1^{-2})$  by Conditions (C1)-(C3) and the properties of the exponential family. However when  $k > n_1$ , by equation (S2.11) and Chebyshev's inequality, we have that  $P(\|\boldsymbol{\xi}_k\| > \epsilon|\dot{\mathbf{b}}) \leq \text{tr}\{\text{Cov}(\boldsymbol{\xi}_k|\dot{\mathbf{b}})\}/\epsilon^2 = O(N^{-1})$  since  $n(mn_1^2)^{-1} = O(N^{-1})$ . Thus given  $\dot{\mathbf{b}}$ , it holds that  $\|\boldsymbol{\xi}_k\| = O_p(N^{-1/2})$  and  $E(\|\boldsymbol{\xi}_k\|^4|\dot{\mathbf{b}}) = O(N^{-2})$ . Therefore

$$\begin{aligned} \sum_{k=1}^N E\{\|\boldsymbol{\xi}_k\|^2 I(\|\boldsymbol{\xi}_k\| > \epsilon)|\dot{\mathbf{b}}\} &\leq \sum_{k=1}^N E(\|\boldsymbol{\xi}_k\|^4|\dot{\mathbf{b}})^{1/2} P(\|\boldsymbol{\xi}_k\| > \epsilon|\dot{\mathbf{b}})^{1/2} \\ &= \sum_{k=1}^{n_1} E(\|\boldsymbol{\xi}_k\|^4|\dot{\mathbf{b}})^{1/2} P(\|\boldsymbol{\xi}_k\| > \epsilon|\dot{\mathbf{b}})^{1/2} \\ &\quad + \sum_{k=n_1+1}^N E(\|\boldsymbol{\xi}_k\|^4|\dot{\mathbf{b}})^{1/2} P(\|\boldsymbol{\xi}_k\| > \epsilon|\dot{\mathbf{b}})^{1/2} \\ &\leq n_1 \max_{1 \leq k \leq n_1} \{E(\|\boldsymbol{\xi}_k\|^4|\dot{\mathbf{b}})^{1/2} P(\|\boldsymbol{\xi}_k\| > \epsilon|\dot{\mathbf{b}})^{1/2}\} \end{aligned}$$

$$\begin{aligned}
 & + (N - n_1) \sup_{k > n_1} \{E(\|\boldsymbol{\xi}_k\|^4 | \dot{\mathbf{b}})^{1/2} P(\|\boldsymbol{\xi}_k\| > \epsilon |\dot{\mathbf{b}})^{1/2}\} \\
 & = n_1 \times O(n_1^{-3/2}) + (N - n_1) \times O(N^{-3/2}) \\
 & = O(n_1^{-1/2}) + O(N^{-1/2}) \\
 & = o(1).
 \end{aligned}$$

The required result follows by Conditions (C1)-(C2) and the Lindeberg-Feller Central Limit Theorem. Furthermore, the general result holds straightforwardly by replacing  $n_1$  with  $O(n_L)$  in the above argument, noting that any row of  $\mathbf{A}$  can only select a fixed number of clusters.

## S2.2 Proof of Equation (4)

For the Poisson pure random intercept model, we have  $\mathbf{B} = \text{diag}(ne^{\dot{b}_1} + \hat{\sigma}_b^{-2}, \dots, ne^{\dot{b}_m} + \hat{\sigma}_b^{-2})$  and  $\mathbf{R}(\tilde{\boldsymbol{\theta}}) = \{ne^{\dot{b}_1}(\hat{b}_1 - \dot{b}_1)^2, \dots, ne^{\dot{b}_m}(\hat{b}_m - \dot{b}_m)^2\}^\top$ .

Next, suppose that  $\mathbf{A}$  picks out the first random intercept, i.e.,  $\mathbf{A} = [1, \mathbf{0}_{m-1}^\top]$ . Then we have

$$\begin{aligned}
 n^{1/2}(\hat{b}_1 - \dot{b}_1) & = n^{1/2} \mathbf{A} \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}}) + \frac{1}{2} n^{1/2} \mathbf{A} \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \\
 & = n^{-1/2} \left\{ \left( \sum_{j=1}^n y_{1j} - e^{\dot{b}_1} \right) - \dot{b}_1 / \hat{\sigma}_b^2 \right\} / \left\{ e^{\dot{b}_1} + 1 / (\hat{\sigma}_b^2 n) \right\} \\
 & \quad - \frac{1}{2} \left\{ n^{1/2} e^{\dot{b}_1} (\hat{b}_1 - \dot{b}_1)^2 \right\} / \left\{ e^{\dot{b}_1} + 1 / (\hat{\sigma}_b^2 n) \right\}
 \end{aligned}$$

$$\begin{aligned}
 &= \left[ n^{-1/2} \left\{ \left( \sum_{j=1}^n y_{1j} - e^{\hat{b}_1} \right) - \hat{b}_1 / \hat{\sigma}_b^2 \right\} / \left\{ e^{\hat{b}_1} + 1 / (\hat{\sigma}_b^2 n) \right\} \right] / \\
 &\quad \left[ 1 + \left\{ \frac{1}{2} e^{\tilde{b}_1} (\hat{b}_1 - \dot{b}_1) \right\} / \left\{ e^{\hat{b}_1} + 1 / (\hat{\sigma}_b^2 n) \right\} \right] \\
 &= n^{-1/2} \sum_{j=1}^n (y_{1j} e^{-\hat{b}_1} - 1) + o_p(1),
 \end{aligned}$$

where  $\tilde{b}_1$  lies between  $\hat{b}_1$  and  $\dot{b}_1$ , and for the last line we have used the fact that  $\hat{b}_1 - \dot{b}_1 = o_p(1)$ .

Now,  $\{y_{1j} e^{-\hat{b}_1} - 1\}_{j=1}^n$  is an exchangeable collection of uncorrelated random variables with mean zero and finite non-zero variance. Furthermore, we have for  $k \neq l$

$$\begin{aligned}
 \text{Cov}\{(y_{1k} e^{-\hat{b}_1} - 1)^2, (y_{1l} e^{-\hat{b}_1} - 1)^2\} &= E[\text{Cov}\{(y_{1k} e^{-\hat{b}_1} - 1)^2, (y_{1l} e^{-\hat{b}_1} - 1)^2 | \hat{b}_1\}] \\
 &\quad + \text{Cov}[E\{(y_{1k} e^{-\hat{b}_1} - 1)^2 | \hat{b}_1\}, E\{(y_{1l} e^{-\hat{b}_1} - 1)^2 | \hat{b}_1\}] \\
 &= 0 + \text{Cov}(e^{-\hat{b}_1}, e^{-\hat{b}_1}) \\
 &= e^{\hat{\sigma}_b^2} (e^{\hat{\sigma}_b^2} - 1) \neq 0.
 \end{aligned}$$

Thus by the Central Limit Theorem for exchangeable random variables (Blum et al., 1958), it holds that  $n^{-1/2} \sum_{j=1}^n (y_{1j} e^{-\hat{b}_1} - 1) \xrightarrow{D} N(0, e^{\hat{\sigma}_b^2})$ . Since we know  $\text{Var}\{n^{-1/2} \sum_{j=1}^n (y_{1j} e^{-\hat{b}_1} - 1)\} = e^{\hat{\sigma}_b^2}/2$  and also that  $n^{-1/2} \sum_{j=1}^n (y_{1j} e^{-\hat{b}_1} - 1) = O_p(1)$  by Chebyshev's inequality, there is no other normalization possible for an asymptotic normality result to hold.

Finally, we also have

$$\begin{aligned}
 n^{1/2}(\hat{b}_1 - \dot{b}_1) &= n^{-1/2} \sum_{j=1}^n (y_{1j} e^{-\dot{b}_1} - 1) + O_p(n^{-1/2}) \\
 \implies \hat{b}_1 &= \dot{b}_1 + n^{-1} \sum_{j=1}^n (y_{1j} e^{-\dot{b}_1} - 1) + O_p(n^{-1}) \\
 &= \dot{b}_1 + o_p(1), \quad \text{by the Weak Law of Large Numbers.}
 \end{aligned}$$

### S2.3 Proof of Theorem 2

We begin by developing two key equations, (S2.12) and (S2.13), that will be used throughout the unconditional regime. These are derived from equations (S2.9) and (S2.10) and are used in the proofs of Theorems 2-5 as well as Corollary 1. Under Conditions (C1)-(C2), the following order results are used:  $\mathbf{B}_3^{-1}$  is a component-wise  $O_p(n_L^{-1})$  block-diagonal matrix,  $\mathbf{B}_2 = O_p(n_U)$  component-wise,  $\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i^\top = O_p(n_i)$  component-wise, and  $\mathbf{C}^{-1} = O_p(m^{-1})$  component-wise. Also, by the conditional independence, we have

$$E\{\mathbf{Z}^\top(\mathbf{y} - \dot{\boldsymbol{\mu}})\} = E[E\{\mathbf{Z}^\top(\mathbf{y} - \dot{\boldsymbol{\mu}})|\dot{\mathbf{b}}\}] = \mathbf{0}_{mp},$$

$$\text{Var}\{\mathbf{Z}^\top(\mathbf{y} - \dot{\boldsymbol{\mu}})\} = E[\text{Var}\{\mathbf{Z}^\top(\mathbf{y} - \dot{\boldsymbol{\mu}})|\dot{\mathbf{b}}\}] + \text{Var}[E\{\mathbf{Z}^\top(\mathbf{y} - \dot{\boldsymbol{\mu}})|\dot{\mathbf{b}}\}] = E(\mathbf{Z}^\top \dot{\mathbf{W}} \mathbf{Z}),$$

so that  $\dot{\phi}^{-1} \mathbf{D}_r^{-1} \mathbf{Z}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}}) = O_p(1)$  using Chebyshev's inequality. Therefore we have the key equations

$$\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}} = m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i + O_p(N^{-1/2}) + O_p(n_L^{-1}) + \frac{1}{2} \{ \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \}_{[1:p]} \quad (\text{S2.12})$$

and

$$\begin{aligned} \hat{\mathbf{b}} - \dot{\mathbf{b}} &= -\mathbf{1}_m \otimes m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i + \mathbf{B}_3^{-1} \{ \dot{\phi}^{-1} \mathbf{Z}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}}) \} \\ &\quad + O_p(N^{-1/2}) + O_p(n_L^{-1}) + \frac{1}{2} \{ \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \}_{[p+1:(m+1)p]}. \end{aligned} \quad (\text{S2.13})$$

By equation (S2.12), we have

$$m^{1/2}(\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}}) = m^{-1/2} \sum_{i=1}^m \dot{\mathbf{b}}_i + O_p(n_L^{-1/2}) + O_p(m^{1/2} n_L^{-1}) + \frac{1}{2} m^{1/2} \{ \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \}_{[1:p]}.$$

Next, we consider two separate scenarios. First, suppose that  $mn_U^{-1} \rightarrow \infty$ .

Then by the order results for the remainder term in Section S3.2, the first  $p$  components of  $\mathbf{D}^* \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}})$  are of order  $O_p(m^{1/2} n_L^{-1})$ , and so the first  $p$  components of  $\mathbf{D}^*(\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})$  can be shown to be

$$m^{1/2}(\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}}) = m^{-1/2} \sum_{i=1}^m \dot{\mathbf{b}}_i + o_p(1).$$

The required result then follows from the independence of the random effects and the normal assumption on the  $\dot{\mathbf{b}}_i$ ; note the  $mn_L^{-2} \rightarrow 0$  assumption

is required for the remainder term to be smaller order than the linear term.

On the other hand, when  $mn_L^{-1} \rightarrow 0$ , the only difference from the  $mn_L^{-1} \rightarrow \infty$  case is that the first  $p$  components of  $\mathbf{D}^+\mathbf{B}^{-1}\mathbf{R}(\tilde{\boldsymbol{\theta}})$  are now of order  $O_p(m^{-1/2})$  due to the different convergence rate of the prediction gap. The result however follows along similar lines as above.

### S2.4 Proof of Theorem 3

Again we consider two different scenarios. First, suppose  $mn_L^{-1} \rightarrow \infty$ . Then from equation (S2.13) and the order results for the remainder term in Section S3.2, we have that

$$\mathbf{D}_r(\hat{\mathbf{b}} - \dot{\mathbf{b}}) = O_p(n_U^{1/2}m^{-1/2}) + O_p(1) + O_p(n_L^{-1/2}).$$

Based on the above, we obtain  $\mathbf{D}_r\hat{\mathbf{b}} = \mathbf{D}_r\dot{\mathbf{b}} + O_p(1)$ , and thus  $\hat{\mathbf{b}} = \dot{\mathbf{b}} + O_p(n_L^{-1/2})$ . The required result follows by multiplying both sides by  $\mathbf{A}_r$ .

On the other hand, suppose now  $mn_L^{-1} \rightarrow 0$ . Then a normalization by  $m^{1/2}$  is needed instead, and the third derivative term is consequently of order  $O_p(m^{-1/2})$  in probability. We thus obtain

$$m^{1/2}(\hat{\mathbf{b}} - \dot{\mathbf{b}}) = O_p(1) + O_p(m^{1/2}n_L^{-1/2}) + O_p(m^{-1/2}),$$

and the result follows.



As an side remark, note from the above proof when  $mn_L^{-1} \rightarrow 0$ , it holds that  $\|\hat{\mathbf{b}} - \dot{\mathbf{b}}\|_2 = O_p(1)$ , where  $\|\cdot\|_2$  denotes the  $l_2$ -norm. But if  $mn_U^{-1} \rightarrow \infty$  then we instead obtain  $\|\hat{\mathbf{b}} - \dot{\mathbf{b}}\|_2 = O_p(m^{1/2}n_U^{-1/2})$ . This implies that, under the unconditional regime, a consistency result based on the  $l_2$ -norm cannot hold for the entire vector of random effects when there is a partnered fixed effect. If there is no partnered fixed effect though, consistency of the entire vector is sometimes possible. For example, in the Poisson counterexample, we demonstrate that  $\|\hat{\mathbf{b}} - \dot{\mathbf{b}}\|_2 = O_p(m^{1/2}n^{-1/2}) = o_p(1)$  when  $mn^{-1} \rightarrow 0$ .

### S2.5 Proof of Theorem 4 and Corollary 1

We will prove each of the three parts of the theorem separately. The proof of part (a) also proves Corollary 1.

Part (a): When  $mn_U^{-1} \rightarrow \infty$ , we have from equation (S2.13) and the order results for the remainder term in Section S3.2 that

$$\mathbf{D}_r(\hat{\mathbf{b}} - \dot{\mathbf{b}}) = \mathbf{D}_r \mathbf{B}_3^{-1} \mathbf{D}_r \mathbf{D}_r^{-1} \dot{\phi}^{-1} \mathbf{Z}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}}) + o_p(1).$$

This is identical to the proof of Theorem 3. Next, without loss of generality, suppose  $\mathbf{A}_r$  selects the first cluster only. Then we have

$$n_1^{1/2}(\hat{\mathbf{b}}_1 - \dot{\mathbf{b}}_1) = (n_1^{-1} \mathbf{X}_1^\top \dot{\mathbf{W}}_1 \mathbf{X}_1)^{-1} n_1^{-1/2} \{\dot{\phi}^{-1} \mathbf{X}_1^\top (\mathbf{y}_1 - \dot{\boldsymbol{\mu}}_1)\} + o_p(1)$$

$$\stackrel{\Delta}{=} \mathbf{P}_{n_1} + o_p(1).$$

We wish to study the distribution of  $\mathbf{P}_{n_1}$  as  $m, n_L \rightarrow \infty$ . By definition,

$$\lim_{m, n_L \rightarrow \infty} F_{\mathbf{P}_{n_1}}(\mathbf{x}) = \lim_{m, n_L \rightarrow \infty} \int F_{\mathbf{P}_{n_1}|\dot{\mathbf{b}}_1}(\mathbf{x}) f(\dot{\mathbf{b}}_1) d\dot{\mathbf{b}}_1.$$

Since  $F_{\mathbf{P}_{n_1}|\dot{\mathbf{b}}_1}(\mathbf{x})$  is a cdf, then  $F_{\mathbf{P}_{n_1}|\dot{\mathbf{b}}_1}(\mathbf{x}) f(\dot{\mathbf{b}}_1)$  is bounded by  $f(\dot{\mathbf{b}}_1)$ . Hence applying  $\int f(\dot{\mathbf{b}}_1) d\dot{\mathbf{b}}_1 = 1$  and the dominated convergence theorem, we obtain

$$\lim_{m, n_L \rightarrow \infty} F_{\mathbf{P}_{n_1}}(\mathbf{x}) = \int \lim_{m, n_L \rightarrow \infty} F_{\mathbf{P}_{n_1}|\dot{\mathbf{b}}_1}(\mathbf{x}) f(\dot{\mathbf{b}}_1) d\dot{\mathbf{b}}_1 = \int \Psi_{\mathbf{P}_{n_1}|\dot{\mathbf{b}}_1}(\mathbf{x}) f(\dot{\mathbf{b}}_1) d\dot{\mathbf{b}}_1,$$

where  $\Psi_{\mathbf{P}_{n_1}|\dot{\mathbf{b}}_1}(\cdot)$  is the cdf associated with  $N(\mathbf{0}, \mathbf{K}_1)$ , a result which follows from conditional independence and the Lindeberg-Feller Central Limit Theorem used in Theorem 1. The general result follows by noting that the same argument can be applied to any finite subset of the random effects. Note also that the result holds regardless of the true distribution of  $\dot{\mathbf{b}}_i$ .

Part (b): When  $mn_i^{-1} \rightarrow \gamma_i \in (0, \infty)$ , we have from (S2.13) and the order results for the remainder term in Section S3.2 that

$$n_i^{1/2}(\hat{\mathbf{b}}_i - \dot{\mathbf{b}}_i) = (n_i^{-1} \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} n_i^{-1/2} \{\phi^{j-1} \mathbf{X}_i^\top (\mathbf{y}_i - \boldsymbol{\mu}_i)\} - (\gamma_i m)^{-1/2} \sum_{i=1}^m \dot{\mathbf{b}}_i + O_p(n_L^{-1/2}),$$

from the same development as in the proof of Part (a). Letting

$$\mathbf{E}_1 = (n_i^{-1} \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} n_i^{-1/2} \phi^{j-1} \mathbf{X}_i^\top (\mathbf{y}_i - \boldsymbol{\mu}_i) \text{ and } \mathbf{E}_2 = m^{-1/2} \sum_{i=1}^m \dot{\mathbf{b}}_i, \text{ then}$$

since  $\mathbf{E}_1$  and  $\mathbf{E}_2$  are independent given  $\dot{\mathbf{b}}_i$ , we obtain for any  $i$ ,

$$\begin{aligned}
 \lim_{m, n_L \rightarrow \infty} F_{\mathbf{E}_1, \mathbf{E}_2}(\mathbf{x}, \mathbf{y}) &= \lim_{m, n_L \rightarrow \infty} \int F_{\mathbf{E}_1, \mathbf{E}_2 | \dot{\mathbf{b}}_i}(\mathbf{x}, \mathbf{y}) f(\dot{\mathbf{b}}_i) d\dot{\mathbf{b}}_i \\
 &= \lim_{m, n_L \rightarrow \infty} \int F_{\mathbf{E}_1 | \dot{\mathbf{b}}_i}(\mathbf{x}) F_{\mathbf{E}_2 | \dot{\mathbf{b}}_i}(\mathbf{y}) f(\dot{\mathbf{b}}_i) d\dot{\mathbf{b}}_i \\
 &= \int \lim_{m, n_L \rightarrow \infty} F_{\mathbf{E}_1 | \dot{\mathbf{b}}_i}(\mathbf{x}) F_{\mathbf{E}_2 | \dot{\mathbf{b}}_i}(\mathbf{y}) f(\dot{\mathbf{b}}_i) d\dot{\mathbf{b}}_i \\
 &= \int \lim_{n_L \rightarrow \infty} F_{\mathbf{E}_1 | \dot{\mathbf{b}}_i}(\mathbf{x}) \lim_{m \rightarrow \infty} F_{\mathbf{E}_2 | \dot{\mathbf{b}}_i}(\mathbf{y}) f(\dot{\mathbf{b}}_i) d\dot{\mathbf{b}}_i \\
 &= \Psi_{\mathbf{E}_2}(\mathbf{y}) \int \lim_{n_L \rightarrow \infty} F_{\mathbf{E}_1 | \dot{\mathbf{b}}_i}(\mathbf{x}) f(\dot{\mathbf{b}}_i) d\dot{\mathbf{b}}_i,
 \end{aligned}$$

where  $\Psi_{\mathbf{E}_2}(\cdot)$  is the cdf of  $N(\mathbf{0}, \dot{\mathbf{G}})$ . The third line follows from the Dominated Convergence Theorem since  $F_{\mathbf{E}_1 | \dot{\mathbf{b}}_i}(\mathbf{x})$  and  $F_{\mathbf{E}_2 | \dot{\mathbf{b}}_i}(\mathbf{y})$  are cdfs and  $\int f(\dot{\mathbf{b}}_i) d\dot{\mathbf{b}}_i = 1$ . Thus  $\mathbf{E}_1$  and  $\mathbf{E}_2$  are asymptotically independent. The result follows from this asymptotic independence.

Part (c): When  $mn_L^{-1} \rightarrow 0$ , we have from (S2.13) and the order results for the remainder term in Section S3.2 that

$$m^{1/2}(\hat{\mathbf{b}} - \dot{\mathbf{b}}) = -\mathbf{1}_m \otimes \mathbf{I}_p m^{-1/2} \sum_{i=1}^m \dot{\mathbf{b}}_i + o_p(1).$$

The result then follows immediately from the normality assumption on  $\dot{\mathbf{b}}_i$ .

## S2.6 Proof of Theorem 5

Given  $mn_L^{-2} \rightarrow 0$  and  $mn_U^{-1/2} \rightarrow \infty$ , by summing equations (S2.12) and (S2.13) we see that the  $m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i$  terms cancel. Therefore, we are left with

$$\begin{aligned} n_i^{1/2}(\hat{\boldsymbol{\beta}} + \hat{\mathbf{b}}_i - \dot{\boldsymbol{\beta}} - \dot{\mathbf{b}}_i) &= n_i(\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} n_i^{-1/2} \{\dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i)\} \\ &\quad + O_p(m^{-1/2}) + O_p(n_L^{-1/2}) + O_p(m^{-1} n_U^{1/2}) \\ &= n_i(\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} n_i^{-1/2} \dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i) + o_p(1). \end{aligned}$$

The required result follows from the Dominated Convergence Theorem.

## S2.7 Result for Difference Between the Prediction Gaps of Two Clusters

Assume Conditions (C1)-(C5) are satisfied,  $mn_L^{-2} \rightarrow 0$ ,  $mn_U^{-1/2} \rightarrow \infty$ , and  $n_i n_{i'}^{-1} \rightarrow \gamma \in (0, \infty)$ . Then as  $m, n_L \rightarrow \infty$  and unconditional on the random effects  $\dot{\mathbf{b}}$ , for each  $i \neq i' \in \{1, \dots, m\}$  we have

$$n_i^{1/2} \{(\hat{\mathbf{b}}_i - \dot{\mathbf{b}}_i) - (\hat{\mathbf{b}}_{i'} - \dot{\mathbf{b}}_{i'})\} \xrightarrow{D} \text{mixN}(\mathbf{0}, \dot{\mathbf{K}}_i, F_{\dot{\mathbf{b}}_i}) * \text{mixN}(\mathbf{0}, \gamma \dot{\mathbf{K}}_{i'}, F_{\dot{\mathbf{b}}_{i'}}).$$

Proof: Theorem 4 implies that, given  $mn_L^{-2} \rightarrow 0$ ,  $mn_U^{-1/2} \rightarrow \infty$ , and

$n_i n_{i'}^{-1} \rightarrow \gamma \in (0, \infty)$ , we have

$$\begin{aligned}
n_i^{1/2}(\hat{\mathbf{b}}_i - \dot{\mathbf{b}}_i - \hat{\mathbf{b}}_{i'} + \dot{\mathbf{b}}_{i'}) &= (n_i^{-1} \mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} n_i^{-1/2} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i) \\
&\quad + \gamma^{1/2} (n_{i'}^{-1} \mathbf{X}_{i'}^\top \dot{\mathbf{W}}_{i'} \mathbf{X}_{i'})^{-1} n_{i'}^{-1/2} \mathbf{X}_{i'}^\top (\mathbf{y}_{i'} - \dot{\boldsymbol{\mu}}_{i'}) \\
&\quad + O_p(m^{-1/2}) + O_p(n_L^{-1/2}) + O_p(m^{-1} n_U^{1/2}),
\end{aligned}$$

and the result follows by the independence of  $\dot{\mathbf{b}}_i$  and  $\dot{\mathbf{b}}_{i'}$ .

### S3 Remainder Term in the Taylor Expansion

In this section, we show that in the Taylor expansion (S2.2), the remainder term  $-\frac{1}{2}\{\nabla^2 Q(\dot{\boldsymbol{\theta}})\}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}})$  is of smaller order component-wise than  $-\{\nabla^2 Q(\dot{\boldsymbol{\theta}})\}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$ . To deal with this remainder term, we have the following from equation (S2.2)

$$\begin{aligned}
\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}} &= \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}}) + \frac{1}{2} \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \\
\Rightarrow \hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}} - \frac{1}{2} \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) &= \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}}) \\
\Rightarrow (\mathbf{I}_{(m+1)p} - \Lambda)(\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) &= \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}}) \\
\Rightarrow \hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}} &= (\mathbf{I}_{(m+1)p} - \Lambda)^{-1} \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}}) \\
&= \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}}) + \left( \sum_{s=1}^{\infty} \Lambda^s \right) \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}}),
\end{aligned}$$

where the last line is derived from repeated application of the Woodbury identity, and  $\mathbf{\Lambda}$  is the appropriate  $(m+1)p \times (m+1)p$  matrix defined in detail later on. The convergence of the geometric sum and thus invertibility of  $(\mathbf{I}_{(m+1)p} - \mathbf{\Lambda})$  is shown in Lemma 4. We will show, using the consistency result  $\|\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}\|_\infty = o_p(1)$ , that  $\sum_{s=1}^{\infty} \mathbf{\Lambda}^s \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$  is of smaller order component-wise than  $\mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$ . This is equivalent to  $0.5 \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}})$  being smaller order component-wise than  $\mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$  in (S2.2).

Let  $\mathbf{T}_1$  denote the first  $p$  components of  $\mathbf{R}(\tilde{\boldsymbol{\theta}})$ ,  $\mathbf{T}_2$  its remaining  $mp$  components, and  $\mathbf{T}_{2i}$  denote the  $\{(i-1)p+1\}$ -th to  $(ip)$ -th components of  $\mathbf{T}_2$ . We first prove a result needed for later developments.

**Lemma 3.** *Assume Conditions (C1) and (C3) are satisfied. Then irrespective of whether  $\dot{\mathbf{b}}$  is conditioned on, it holds that  $\mathbf{R}(\tilde{\boldsymbol{\theta}})_{[1:p]} = \sum_{i=1}^m \mathbf{R}(\tilde{\boldsymbol{\theta}})_{[ip+1:(i+1)p]}$ .*

*Proof.* Recall the Taylor expansion  $\nabla Q(\hat{\boldsymbol{\theta}}) = \mathbf{0} = \nabla Q(\dot{\boldsymbol{\theta}}) + \nabla^2 Q(\dot{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) + \mathbf{R}(\tilde{\boldsymbol{\theta}})$ . Then

$$\begin{aligned} \mathbf{0}_{p \times 1} &= \nabla Q(\hat{\boldsymbol{\theta}})_{[1:p]} \\ &= \{\nabla Q(\dot{\boldsymbol{\theta}}) + \nabla^2 Q(\dot{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) + \mathbf{R}(\tilde{\boldsymbol{\theta}})\}_{[1:p]} \\ &= \sum_{i=1}^m \nabla Q(\hat{\boldsymbol{\theta}})_{[ip+1:(i+1)p]} \\ &= \sum_{i=1}^m \{\nabla Q(\dot{\boldsymbol{\theta}}) + \nabla^2 Q(\dot{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) + \mathbf{R}(\tilde{\boldsymbol{\theta}})\}_{[ip+1:(i+1)p]}. \end{aligned}$$

Since  $\mathbf{Z}_i = \mathbf{X}_i$  for all  $i = 1, \dots, m$  under our simplifying assumption, and

$\sum_{i=1}^m \hat{\mathbf{b}}_i = \mathbf{0}$ , then we obtain

$$\{\nabla Q(\dot{\boldsymbol{\theta}}) + \nabla^2 Q(\dot{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) + \mathbf{R}(\tilde{\boldsymbol{\theta}})\}_{[1:p]} = \sum_{i=1}^m \{\nabla Q(\dot{\boldsymbol{\theta}}) + \nabla^2 Q(\dot{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) + \mathbf{R}(\tilde{\boldsymbol{\theta}})\}_{[ip+1:(i+1)p]}.$$

Therefore, we have  $\mathbf{T}_1 = \mathbf{R}(\tilde{\boldsymbol{\theta}})_{[1:p]} = \sum_{i=1}^m \mathbf{R}(\tilde{\boldsymbol{\theta}})_{[ip+1:(i+1)p]} = \sum_{i=1}^m \mathbf{T}_{2i}$ ,

which follows from the fact that  $\sum_{i=1}^m \nabla Q(\dot{\boldsymbol{\theta}})_{[ip+1:(i+1)p]} = \nabla Q(\dot{\boldsymbol{\theta}})_{[1:p]} -$

$\sum_{i=1}^m \hat{\mathbf{G}}^{-1} \dot{\mathbf{b}}_i$  and  $\sum_{i=1}^m \{\nabla^2 Q(\dot{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})\}_{[ip+1:(i+1)p]} = \{\nabla^2 Q(\dot{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})\}_{[1:p]} +$

$\sum_{i=1}^m \hat{\mathbf{G}}^{-1} \dot{\mathbf{b}}_i - \sum_{i=1}^m \hat{\mathbf{G}}^{-1} \hat{\mathbf{b}}_i. \quad \square$

Next, let  $\mathbf{S}(\boldsymbol{\theta}) = \nabla Q(\boldsymbol{\theta})$ ,  $\tilde{\mathbf{W}}' = \dot{\phi}^{-1} \text{diag}\{a'''(\tilde{\eta}_{11}), \dots, a'''(\tilde{\eta}_{1n_1}), \dots, a'''(\tilde{\eta}_{mn_m})\}$ .

Then the remainder term can be written as

$$\mathbf{R}(\tilde{\boldsymbol{\theta}}) = \begin{bmatrix} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})^\top \frac{\partial^2 \mathbf{S}_{[1]}(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) \\ \vdots \\ (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})^\top \frac{\partial^2 \mathbf{S}_{[(m+1)p]}(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) \end{bmatrix}.$$

Now, for  $1 \leq j \leq p$ , we have  $\mathbf{S}_{[j]}(\boldsymbol{\theta}) = \dot{\phi}^{-1} \mathbf{X}_{[j]}^\top \{\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\theta})\} = \dot{\phi}^{-1} \sum_{i=1}^m \sum_{l=1}^{n_i} x_{il[j]} \{y_{il} -$

$a'(\eta_{il})\}$ , noting this is a scalar. Thus

$$\frac{\partial}{\partial \boldsymbol{\theta}} \mathbf{S}_{[j]}(\boldsymbol{\theta}) = -\dot{\phi}^{-1} \sum_{i=1}^m \sum_{l=1}^{n_i} \begin{bmatrix} \mathbf{x}_{il} \\ \frac{\partial}{\partial \mathbf{b}} \eta_{il} \end{bmatrix} a''(\eta_{il}) x_{il[j]} = - \begin{bmatrix} \mathbf{X}^\top \mathbf{W} \mathbf{X}_{[j]} \\ \mathbf{Z}^\top \mathbf{W} \mathbf{X}_{[j]} \end{bmatrix},$$

which is an  $(m+1)p$ -vector. Hence the  $(m+1)p \times (m+1)p$  matrix can be

written as

$$\begin{aligned} \frac{\partial^2 \mathbf{S}_{[j]}(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} &= -\dot{\phi}^{-1} \sum_{i=1}^m \sum_{l=1}^{n_i} \begin{bmatrix} \mathbf{x}_{il} \\ \frac{\partial}{\partial \mathbf{b}} \eta_{il} \end{bmatrix} a'''(\tilde{\eta}_{il}) x_{il[j]} \begin{bmatrix} \mathbf{x}_{il} \\ \frac{\partial}{\partial \mathbf{b}} \eta_{il} \end{bmatrix}^\top \\ &= - \begin{bmatrix} \mathbf{X}^\top \text{diag}(\mathbf{X}_{[j]}) \tilde{\mathbf{W}}' \mathbf{X} & \mathbf{X}^\top \text{diag}(\mathbf{X}_{[j]}) \tilde{\mathbf{W}}' \mathbf{Z} \\ \mathbf{Z}^\top \text{diag}(\mathbf{X}_{[j]}) \tilde{\mathbf{W}}' \mathbf{X} & \mathbf{Z}^\top \text{diag}(\mathbf{X}_{[j]}) \tilde{\mathbf{W}}' \mathbf{Z} \end{bmatrix}, \quad 1 \leq j \leq p. \end{aligned}$$

Similarly, for  $1 \leq k \leq mp$ ,  $\mathbf{S}_{[p+k]}(\boldsymbol{\theta}) = \dot{\phi}^{-1} \mathbf{Z}_{[k]}^\top \{\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\theta})\} - \{(\mathbf{I}_m \otimes \hat{\mathbf{G}}) \mathbf{b}\}_{[k]}$ ,

such that

$$\frac{\partial}{\partial \boldsymbol{\theta}} \mathbf{S}_{[p+k]}(\boldsymbol{\theta}) = - \begin{bmatrix} \mathbf{X}^\top \mathbf{W} \mathbf{Z}_{[k]} \\ \mathbf{Z}^\top \mathbf{W} \mathbf{Z}_{[k]} + \frac{\partial}{\partial \mathbf{b}} \{(\mathbf{I}_m \otimes \hat{\mathbf{G}}) \mathbf{b}\}_{[k]} \end{bmatrix},$$

where  $\partial/\partial \mathbf{b}\{(\mathbf{I}_m \otimes \hat{\mathbf{G}}) \mathbf{b}\}_{[k]}$  is not a function of  $\boldsymbol{\theta}$ . Thus

$$\frac{\partial^2 \mathbf{S}_{[p+k]}(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} = - \begin{bmatrix} \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[k]}) \tilde{\mathbf{W}}' \mathbf{X} & \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[k]}) \tilde{\mathbf{W}}' \mathbf{Z} \\ \mathbf{Z}^\top \text{diag}(\mathbf{Z}_{[k]}) \tilde{\mathbf{W}}' \mathbf{X} & \mathbf{Z}^\top \text{diag}(\mathbf{Z}_{[k]}) \tilde{\mathbf{W}}' \mathbf{Z} \end{bmatrix}, \quad 1 \leq k \leq mp.$$

Next, recall that  $\mathbf{B}_2(\mathbf{B}_3 + \mathbf{B}_4)^{-1} = [\mathbf{I}_p - \hat{\mathbf{G}}^{-1}(\mathbf{X}_1^\top \dot{\mathbf{W}}_1 \mathbf{X}_1 + \hat{\mathbf{G}}^{-1})^{-1}, \dots, \mathbf{I}_p -$

$\hat{\mathbf{G}}^{-1}(\mathbf{X}_m^\top \dot{\mathbf{W}}_m \mathbf{X}_m + \hat{\mathbf{G}}^{-1})^{-1}]$ . By Lemma 1 and the blockwise inversion for-

mula for  $\mathbf{B}^{-1}$ , the first  $p$  components of  $\mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}})$  are given by

$$\begin{aligned} & \begin{bmatrix} \mathbf{C}^{-1} & -\mathbf{C}^{-1} \mathbf{B}_2(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \end{bmatrix} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \\ &= \mathbf{C}^{-1} \begin{bmatrix} \mathbf{I}_p & -[\mathbf{I}_p - \hat{\mathbf{G}}^{-1}(\mathbf{X}_1^\top \dot{\mathbf{W}}_1 \mathbf{X}_1 + \hat{\mathbf{G}}^{-1})^{-1}, \dots, \mathbf{I}_p - \hat{\mathbf{G}}^{-1}(\mathbf{X}_m^\top \dot{\mathbf{W}}_m \mathbf{X}_m + \hat{\mathbf{G}}^{-1})^{-1}] \end{bmatrix} \mathbf{R}(\tilde{\boldsymbol{\theta}}) \end{aligned}$$



$$= \mathbf{C}^{-1} \left\{ \mathbf{T}_1 - \sum_{i=1}^m \mathbf{T}_{2i} + \sum_{i=1}^m \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{T}_{2i} \right\}. \quad (\text{S3.1})$$

Similarly, the last  $mp$  components of  $\mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}})$  are

$$\begin{aligned} & \left[ -(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{B}_2^\top \mathbf{C}^{-1} \quad (\mathbf{B}_3 + \mathbf{B}_4)^{-1} + (\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{B}_2^\top \mathbf{C}^{-1} \mathbf{B}_2 (\mathbf{B}_3 + \mathbf{B}_4)^{-1} \right] \mathbf{R}(\tilde{\boldsymbol{\theta}}) \\ &= -(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{B}_2^\top \left[ \mathbf{C}^{-1} \quad -\mathbf{C}^{-1} \mathbf{B}_2 (\mathbf{B}_3 + \mathbf{B}_4)^{-1} \right] \mathbf{R}(\tilde{\boldsymbol{\theta}}) + (\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{T}_2. \end{aligned} \quad (\text{S3.2})$$

Hence the first  $p$  components of  $\mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}})$  are given by

$$\mathbf{F}_1 = \mathbf{C}^{-1} \sum_{i=1}^m \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{T}_{2i},$$

and the last  $mp$  components of  $\mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}})$  are given by

$$\mathbf{F}_2 = -(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{B}_2^\top \mathbf{F}_1 + (\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{T}_2.$$

Next, we have

$$\mathbf{T}_2 = \begin{bmatrix} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})^\top \frac{\partial^2 \mathcal{S}_{[p+1]}(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) \\ \vdots \\ (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})^\top \frac{\partial^2 \mathcal{S}_{[(m+1)p]}(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) \end{bmatrix} = \begin{bmatrix} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})^\top \frac{\partial^2 \mathcal{S}_{[p+1]}(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \\ \vdots \\ (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})^\top \frac{\partial^2 \mathcal{S}_{[(m+1)p]}(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \end{bmatrix} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) \triangleq \mathbf{F}_3 (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})$$

and

$$\mathbf{T}_{2i} = \begin{bmatrix} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})^\top \frac{\partial^2 \mathbf{S}_{[(i-1)p+1]}(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) \\ \vdots \\ (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})^\top \frac{\partial^2 \mathbf{S}_{[ip]}(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) \end{bmatrix} = \begin{bmatrix} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})^\top \frac{\partial^2 \mathbf{S}_{[(i-1)p+1]}(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \\ \vdots \\ (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})^\top \frac{\partial^2 \mathbf{S}_{[ip]}(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \end{bmatrix} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) \triangleq \mathbf{F}_{3i} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}).$$

Here,  $\mathbf{F}_3$  is a  $mp \times (m+1)p$  matrix and  $\mathbf{F}_{3i}$  is  $p \times (m+1)p$ . Notice that

$$\mathbf{F}_3 = [\mathbf{F}_{31}^\top, \dots, \mathbf{F}_{3n}^\top]^\top. \text{ Furthermore,}$$

$$\begin{aligned} \mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}}) &= \begin{bmatrix} \mathbf{F}_1 \\ \mathbf{F}_2 \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{T}_{2i} \\ -(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{B}_2^\top \mathbf{F}_1 + (\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{T}_2 \end{bmatrix} \\ &= \begin{bmatrix} \sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{F}_{3i} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) \\ -(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{B}_2^\top \sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{F}_{3i} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) + (\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{F}_3 (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) \end{bmatrix} \\ &= \begin{bmatrix} \sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{F}_{3i} \\ -(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{B}_2^\top \sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{F}_{3i} + (\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{F}_3 \end{bmatrix} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}) \\ &= 2\boldsymbol{\Lambda} (\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}). \end{aligned}$$

The  $k$ th row of  $\mathbf{F}_{3i}$  for  $1 \leq k \leq p$  is given by

$$\begin{aligned} & -(\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}})^\top \begin{bmatrix} \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[(i-1)p+k]}) \tilde{\mathbf{W}}' \mathbf{X} & \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[(i-1)p+k]}) \tilde{\mathbf{W}}' \mathbf{Z} \\ \mathbf{Z}^\top \text{diag}(\mathbf{Z}_{[(i-1)p+k]}) \tilde{\mathbf{W}}' \mathbf{X} & \mathbf{Z}^\top \text{diag}(\mathbf{Z}_{[(i-1)p+k]}) \tilde{\mathbf{W}}' \mathbf{Z} \end{bmatrix} \\ &= -\delta_{m,n_L}^{-1} [\delta_{m,n_L} (\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}})^\top \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[(i-1)p+k]}) \tilde{\mathbf{W}}' \mathbf{X} + \delta_{m,n_L} (\hat{\mathbf{b}} - \dot{\mathbf{b}})^\top \mathbf{Z}^\top \text{diag}(\mathbf{Z}_{[(i-1)p+k]}) \tilde{\mathbf{W}}' \mathbf{X}, \end{aligned}$$

$$\delta_{m,n_L}(\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}})^\top \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[(i-1)p+k]}) \tilde{\mathbf{W}}' \mathbf{Z} + \delta_{m,n_L}(\hat{\mathbf{b}} - \dot{\mathbf{b}})^\top \mathbf{Z}^\top \text{diag}(\mathbf{Z}_{[(i-1)p+k]}) \tilde{\mathbf{W}}' \mathbf{Z}, \quad (\text{S3.3})$$

where  $\delta_{m,n_L}$  is a positive unbounded monotonically increasing sequence such that  $\delta_{m,n_L} \|\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}\|_\infty = O_p(1)$ . The consistency results proved in Section S1 ensure that such a  $\delta_{m,n_L}$  must exist; this is true for both the conditional and unconditional regimes.

Observe that only the  $\{(\sum_{l=0}^{i-1} n_l) + 1\}$ th to  $(\sum_{l=0}^i n_l)$ th components of  $\mathbf{Z}_{[(i-1)p+k]}$  are non-zero, where we define  $n_0 := 0$ . This means that, for any  $1 \leq k \leq p$ , only the  $\{(i-1)p + 1\}$ th to  $(ip)$ th columns of both  $\mathbf{X}^\top \text{diag}(\mathbf{Z}_{[(i-1)p+k]}) \tilde{\mathbf{W}}' \mathbf{Z}$  and  $\mathbf{Z}^\top \text{diag}(\mathbf{Z}_{[(i-1)p+k]}) \tilde{\mathbf{W}}' \mathbf{Z}$  will be non-zero. In other words, other than the first  $p$  columns, only the  $(ip + 1)$ th to  $\{(i+1)p\}$ th columns of  $\mathbf{F}_{3i}$  are non-zero. Thus  $\mathbf{F}_3$ , disregarding its first  $p$  columns, is an  $mp \times mp$  block-diagonal matrix.

The non-zero components of  $\delta_{m,n_L} \mathbf{F}_3$  and  $\delta_{m,n_L} \mathbf{F}_{3i}$  are all  $O_p(n_U)$  component-wise, again because at most  $n_i$  components of  $\mathbf{Z}_{[(i-1)p+k]}$  are non-zero.

For ease of notation and understanding, we now represent all terms using their orders only. Since  $\mathbf{C}^{-1} = O_p(m^{-1})$  and  $\hat{\mathbf{G}}^{-1}(\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} = O_p(n_i^{-1})$ , from the above discussion we have that

$\sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1}(\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{F}_{3i}$  is a  $p \times (m+1)p$  matrix of the form

$\delta_{m,n_L}^{-1} [O_p(1), O_p(m^{-1}), \dots, O_p(m^{-1})]$ . Next,  $(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{B}_2^\top = [\mathbf{I}_p + O_p(n_1^{-1}), \dots, \mathbf{I}_p +$

$O_p(n_m^{-1})]^\top$  and  $(\mathbf{B}_3 + \mathbf{B}_4)^{-1}$  is a block-diagonal  $O_p(n_L^{-1})$  matrix component-wise. Therefore, we find that  $\mathbf{\Lambda}$  is of the form

$$\begin{aligned}
 & 0.5 \begin{bmatrix} \sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{F}_{3i} \\ -(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{B}_2^\top \sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{F}_{3i} \end{bmatrix} + 0.5 \begin{bmatrix} \mathbf{0}_{p \times (m+1)p} \\ (\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{F}_3 \end{bmatrix} \\
 & \triangleq \mathbf{\Lambda}_1 + \mathbf{\Lambda}_2 \\
 & = \frac{1}{\delta_{m,n_L}} \begin{bmatrix} O_p(1) & O_p(m^{-1}) & \cdots & O_p(m^{-1}) \\ \vdots & \vdots & \vdots & \vdots \\ O_p(1) & O_p(m^{-1}) & \cdots & O_p(m^{-1}) \end{bmatrix} + \frac{1}{\delta_{m,n_L}} \begin{bmatrix} & & \mathbf{0}_{p \times (m+1)p} & & \\ O_p(1) & O_p(1) & \mathbf{0} & \cdots & \mathbf{0} \\ O_p(1) & \mathbf{0} & O_p(1) & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ O_p(1) & \mathbf{0} & \cdots & \mathbf{0} & O_p(1) \end{bmatrix} \\
 & = \frac{1}{\delta_{m,n_L}} \begin{bmatrix} O_p(1) & O_p(m^{-1}) & O_p(m^{-1}) & \cdots & O_p(m^{-1}) \\ O_p(1) & O_p(1) & O_p(m^{-1}) & \cdots & O_p(m^{-1}) \\ O_p(1) & O_p(m^{-1}) & O_p(1) & \cdots & O_p(m^{-1}) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ O_p(1) & O_p(m^{-1}) & \cdots & O_p(m^{-1}) & O_p(1) \end{bmatrix}. \quad (\text{S3.4})
 \end{aligned}$$

Writing  $\mathbf{\Lambda} = \delta_{m,n_L}^{-1} \mathbf{\Lambda}_\delta$ , we see that the component-wise order of  $\mathbf{\Lambda}_\delta$  remains the same no matter how many times it is multiplied by itself. Furthermore, each row of  $\mathbf{\Lambda}_\delta^s$  is  $O_p(1)$  for only a finite number of components, and  $O_p(m^{-1})$  for the others. We will use these facts to examine the behaviour

of  $\sum_{s=1}^{\infty} \mathbf{\Lambda}^s \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}}) = \sum_{s=1}^{\infty} \delta_{m, n_L}^{-s} \mathbf{\Lambda}_{\delta}^s \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$ , and we will do so separately for the conditional and unconditional regimes. Before proceeding, we first confirm the convergence of  $\sum_{s=1}^{\infty} \mathbf{\Lambda}^s$ .

**Lemma 4.** *Assume Conditions (C1)-(C5) are satisfied. Then with probability tending to one as  $m, n_L \rightarrow \infty$ , the geometric sum  $\sum_{s=1}^{\infty} \mathbf{\Lambda}^s$  converges.*

*Proof.* To prove the result we will show that, with probability tending to one as  $m, n_L \rightarrow \infty$ ,  $\|\mathbf{\Lambda}\| < 1$  for some sub-multiplicative matrix norm  $\|\cdot\|$ . In particular, we will consider the maximum absolute row sum of  $\mathbf{\Lambda}$ , denoted by  $\|\cdot\|_{\infty}$  i.e., the operator norm induced by the vector infinity norm.

From (S3.4), we have  $\|\mathbf{\Lambda}\|_{\infty} \leq \|\mathbf{\Lambda}_1\|_{\infty} + \|\mathbf{\Lambda}_2\|_{\infty}$ . We first examine  $\|\mathbf{\Lambda}_1\|_{\infty}$ . We may break up  $\mathbf{\Lambda}_1$  into

$$0.5 \begin{bmatrix} \sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^{\top} \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{F}_{3i} \\ - \sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^{\top} \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{F}_{3i} \\ - \sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^{\top} \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{F}_{3i} \\ \vdots \end{bmatrix} +$$

$$0.5 \begin{bmatrix} \mathbf{0}_p \\ (\mathbf{X}_1^\top \dot{\mathbf{W}}_1 \mathbf{X}_1 + \hat{\mathbf{G}}^{-1})^{-1} \hat{\mathbf{G}}^{-1} \sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{F}_{3i} \\ (\mathbf{X}_2^\top \dot{\mathbf{W}}_2 \mathbf{X}_2 + \hat{\mathbf{G}}^{-1})^{-1} \hat{\mathbf{G}}^{-1} \sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{F}_{3i} \\ \vdots \\ (\mathbf{X}_m^\top \dot{\mathbf{W}}_m \mathbf{X}_m + \hat{\mathbf{G}}^{-1})^{-1} \hat{\mathbf{G}}^{-1} \sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{F}_{3i} \end{bmatrix} \\ \triangleq \mathbf{\Lambda}_3 + \mathbf{\Lambda}_4 \mathbf{\Lambda}_5,$$

where  $\mathbf{\Lambda}_4 = \text{bdiag}(\mathbf{0}_{p \times p}, (\mathbf{X}_1^\top \dot{\mathbf{W}}_1 \mathbf{X}_1 + \hat{\mathbf{G}}^{-1})^{-1} \hat{\mathbf{G}}^{-1}, \dots, (\mathbf{X}_m^\top \dot{\mathbf{W}}_m \mathbf{X}_m + \hat{\mathbf{G}}^{-1})^{-1} \hat{\mathbf{G}}^{-1})$ ,

and  $\mathbf{\Lambda}_5 = (0, \mathbf{1}_m^\top)^\top \otimes \sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{F}_{3i}$ . We can also

write

$\mathbf{\Lambda}_3 = -0.5(\mathbf{1}_m^* \otimes \mathbf{I}_p) \sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{F}_{3i}$  and use the

(component-wise) order results as used in (S3.4) to see that  $\|\mathbf{\Lambda}_3\|_\infty \leq$

$\| -0.5(\mathbf{1}_m^* \otimes \mathbf{I}_p) \|_\infty \|\sum_{i=1}^m \mathbf{C}^{-1} \hat{\mathbf{G}}^{-1} (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{F}_{3i}\|_\infty = o_p(1)$ . Next,

we have  $\|\mathbf{\Lambda}_4 \mathbf{\Lambda}_5\|_\infty \leq \|\mathbf{\Lambda}_4\|_\infty \|\mathbf{\Lambda}_5\|_\infty$ . We know  $\|\mathbf{\Lambda}_5\|_\infty = o_p(1)$ , and under

conditions (C1)-(C2), we have  $\|\mathbf{\Lambda}_4\|_\infty = O_p(1)$ . Thus we obtain  $\|\mathbf{\Lambda}_1\|_\infty =$

$o_p(1)$ .

Turning to  $\mathbf{\Lambda}_2$ , we examine each row of  $(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{F}_3$ . First,  $\|(\mathbf{B}_3 + \mathbf{B}_4)^{-1} \mathbf{F}_3\|_\infty \leq \|(\mathbf{B}_3 + \mathbf{B}_4)^{-1}\|_\infty \|\mathbf{F}_3\|_\infty$  and by conditions (C1)-(C2), we have

$\|(\mathbf{B}_3 + \mathbf{B}_4)^{-1}\|_\infty = O_p(n_L^{-1})$ . Now, without loss of generality consider the

first row of  $\mathbf{F}_3$ . This is given by

$$\begin{aligned}
 & - \left[ (\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}})^\top \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{X} + (\hat{\mathbf{b}} - \dot{\mathbf{b}})^\top \mathbf{Z}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{X}, \right. \\
 & \left. (\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}})^\top \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{Z} + (\hat{\mathbf{b}} - \dot{\mathbf{b}})^\top \mathbf{Z}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{Z} \right] \\
 & = - \left[ (\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}})^\top \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{X} + (\hat{\mathbf{b}} - \dot{\mathbf{b}})^\top \mathbf{Z}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{X}, \right. \\
 & \left. (\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}})^\top \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{X} + (\hat{\mathbf{b}} - \dot{\mathbf{b}})^\top \mathbf{Z}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{X}, \mathbf{0}_{(m-1)p}^\top \right],
 \end{aligned}$$

since  $\text{diag}(\mathbf{Z}_{[1,1]})$  selects for the first cluster. Let  $\bar{\mathbf{1}}_p$  be a  $p$ -vector whose entries consist of the (component-wise) signs of  $(\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}})^\top \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{X} + (\hat{\mathbf{b}} - \dot{\mathbf{b}})^\top \mathbf{Z}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{X}$ . Then the absolute row sum of the first row of  $\mathbf{F}_3$  is given by

$$\begin{aligned}
 & 2 | \{ (\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}})^\top \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{X} + (\hat{\mathbf{b}} - \dot{\mathbf{b}})^\top \mathbf{Z}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{X} \} \bar{\mathbf{1}}_p | \\
 & = 2 | \{ (\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}})^\top \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{X} + (\hat{\mathbf{b}}_1 - \dot{\mathbf{b}}_1)^\top \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{X} \} \bar{\mathbf{1}}_p | \\
 & = 2 | \{ (\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}} + \hat{\mathbf{b}}_1 - \dot{\mathbf{b}}_1)^\top \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{X} \} \bar{\mathbf{1}}_p | \\
 & \leq 2p \| \{ \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{X} (\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}} + \hat{\mathbf{b}}_1 - \dot{\mathbf{b}}_1) \} \|_\infty \\
 & \leq 2p \| \{ \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{X} \} \|_\infty \| \hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}} + \hat{\mathbf{b}}_1 - \dot{\mathbf{b}}_1 \|_\infty \\
 & \leq 2p \| \{ \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[1,1]}) \tilde{\mathbf{W}}' \mathbf{X} \} \|_\infty (\| \hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}} \|_\infty + \| \hat{\mathbf{b}}_1 - \dot{\mathbf{b}}_1 \|_\infty) \\
 & \leq 2p \max_{k \in \{1, \dots, mp\}} \| \{ \mathbf{X}^\top \text{diag}(\mathbf{Z}_{[1,k]}) \tilde{\mathbf{W}}' \mathbf{X} \} \|_\infty (\| \hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}} \|_\infty + \| \hat{\mathbf{b}}_1 - \dot{\mathbf{b}}_1 \|_\infty)
 \end{aligned}$$

$$\begin{aligned}
 &= 2p \max_{k \in \{1, \dots, mp\}} \|\{\mathbf{X}^\top \text{diag}(\mathbf{Z}_{[k]}) \tilde{\mathbf{W}}' \mathbf{X}\}\|_\infty \|\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}}\|_\infty \\
 &\quad + 2p \max_{k \in \{1, \dots, mp\}} \|\{\mathbf{X}^\top \text{diag}(\mathbf{Z}_{[k]}) \tilde{\mathbf{W}}' \mathbf{X}\}\|_\infty \|\hat{\mathbf{b}}_1 - \dot{\mathbf{b}}_1\|_\infty \\
 &\triangleq \alpha + \alpha_1 \triangleq \omega_1,
 \end{aligned}$$

where the second equality follows from  $\text{diag}(\mathbf{Z}_{[1]})$  selecting for only the first cluster, and  $\mathbf{X}_i = \mathbf{Z}_i$ . The first inequality is due to Hölder's inequality.

Again, using Conditions (C1)-(C2) we have  $\max_{k \in \{1, \dots, mp\}} \|\{\mathbf{X}^\top \text{diag}(\mathbf{Z}_{[k]}) \tilde{\mathbf{W}}' \mathbf{X}\}\|_\infty = O_p(n_U)$ .

Now,  $p$  is a constant and the absolute row sum of any row of  $\mathbf{F}_3$  can be bounded analogously in the above way, the only difference being that for the  $k$ th row, then the quantity  $(\hat{\mathbf{b}}_1 - \dot{\mathbf{b}}_1)$  changes to the prediction gap for the cluster that  $\text{diag}(\mathbf{Z}_{[k]})$  selects for. This means that the absolute row sums for the first  $p$  rows of  $\mathbf{F}_3$  are bounded by  $\omega_1$ , the next  $p$  rows by  $\omega_2$ , and so on. Hence, to ensure  $\|\boldsymbol{\Lambda}_2\|_\infty = o_p(1)$  it suffices to ensure that  $\|\boldsymbol{\omega} \otimes \mathbf{1}_p\|_\infty = \|\boldsymbol{\omega}\|_\infty = o_p(n_L)$ , where  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_m)^\top$ .

To show this, define  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_m)^\top$ . Then  $\|\boldsymbol{\omega}\|_\infty \leq \|\alpha \mathbf{1}_m\|_\infty + \|\boldsymbol{\alpha}\|_\infty$ . By Conditions (C1)-(C2), we have  $\|\alpha \mathbf{1}_m\|_\infty = \alpha = O_p(n_U) \times o_p(1) = o_p(n_U) = o_p(n_L)$ . We also have

$$\|\boldsymbol{\alpha}\|_\infty = 2p \max_{k \in \{1, \dots, mp\}} \|\{\mathbf{X}^\top \text{diag}(\mathbf{Z}_{[k]}) \tilde{\mathbf{W}}' \mathbf{X}\}\|_\infty \max_{i \in \{1, \dots, m\}} \|\hat{\mathbf{b}}_i - \dot{\mathbf{b}}_i\|_\infty$$



$$\begin{aligned}
 &= 2p \max_{k \in \{1, \dots, mp\}} \|\{\mathbf{X}^\top \text{diag}(\mathbf{Z}_{[k]}) \tilde{\mathbf{W}}' \mathbf{X}\}\|_\infty \|\hat{\mathbf{b}} - \dot{\mathbf{b}}\|_\infty \\
 &= O_p(n_U) \times o_p(1) = o_p(n_U) = o_p(n_L),
 \end{aligned}$$

where the last line follows from conditions (C1)-(C2), and the fact that  $\|\hat{\boldsymbol{\theta}} - \dot{\boldsymbol{\theta}}\|_\infty = o_p(1)$ . The result follows since  $\|\boldsymbol{\Lambda}\|_\infty$  is therefore of order  $o_p(1)$ , and for any  $\epsilon > 0$  we have  $\|\boldsymbol{\Lambda}\|_\infty < \epsilon$  with probability tending to one as  $m, n_L \rightarrow \infty$ . The argument above holds for both the conditional and unconditional regime, and the required result follows.  $\square$

### S3.1 Conditional Regime

In the conditional regime, we assume without loss of generality that  $\sum_{i=1}^m \dot{\mathbf{b}}_i = \mathbf{0}_p$ , recalling that we can always reparametrise the random effects to satisfy this. From previous derivations, we know that when  $mn_L^{-1} \rightarrow 0$ , the quantity  $\mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$  is of order  $O_p(N^{-1/2})$  for the first  $p$  components and  $O_p(n_L^{-1/2})$  for the last  $mp$  components. By the two properties of  $\boldsymbol{\Lambda}_\delta^s$  noted above, we therefore know that  $\boldsymbol{\Lambda}_\delta^s \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$  is at most  $O_p(n_L^{-1/2})$  component-wise for any  $s$ . Hence  $\sum_{s=1}^\infty \delta_{m, n_L}^{-s} \boldsymbol{\Lambda}_\delta^s \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}}) = \delta_{m, n_L}^{-1} O_p(n_L^{-1/2}) = o_p(n_L^{-1/2})$  for sufficiently large  $m, n_L$  by the properties of a geometric sum. This is sufficient to show that the last  $mp$  components of  $\sum_{s=1}^\infty \boldsymbol{\Lambda}^s \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$  are of smaller order component-wise than  $\mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$ , so that the result

for the prediction gap holds. In particular, we thus know that  $\hat{\mathbf{b}} - \dot{\mathbf{b}} = O_p(n_L^{-1/2})$ . Furthermore, we also know that the convergence rate of  $\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}}$  is at least of order  $O_p(n_L^{-1/2})$ . As a result, we can choose  $\delta_{m,n_L} = n_L^{1/2}$  without affecting the component-wise order properties of  $\boldsymbol{\Lambda}_\delta$ . Applying  $\delta_{m,n_L} = n_L^{1/2}$ , we thus have that  $\sum_{s=1}^{\infty} \delta_{m,n_L}^{-s} \boldsymbol{\Lambda}_\delta^s \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$  is at most of order  $O_p(n_L^{-1})$  component-wise. This is smaller than  $O_p(N^{-1/2})$  when  $mn_L^{-1} \rightarrow 0$ , and the required result follows.

### S3.2 Unconditional Regime

For the unconditional regime, we consider two cases: when  $mn_L^{-1} \rightarrow 0$ , and when  $mn_L^{-1} \rightarrow \infty$  but  $mn_L^{-2} \rightarrow 0$ .

First, consider the case when  $mn_L^{-1} \rightarrow 0$ . From previous derivations, we know that when  $mn_L^{-1} \rightarrow 0$ , the quantity  $\mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$  is of order  $O_p(m^{-1/2})$  for the first  $p$  components and  $O_p(m^{-1/2})$  for the last  $mp$  components. By the two properties of  $\boldsymbol{\Lambda}_\delta^s$  noted above, we therefore know that  $\boldsymbol{\Lambda}_\delta^s \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$  is at most  $O_p(m^{-1/2})$  component-wise for any  $s$ . Hence

$\sum_{s=1}^{\infty} \delta_{m,n_L}^{-s} \boldsymbol{\Lambda}_\delta^s \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}}) = \delta_{m,n_L}^{-1} O_p(m^{-1/2}) = o_p(m^{-1/2})$  for sufficiently large  $m, n_L$ , by the properties of a geometric sum. The required result follows from this. Furthermore, this implies we may set  $\delta_{m,n_L} = m^{1/2}$  without affecting the component-wise order properties of  $\boldsymbol{\Lambda}_\delta$ . Applying

$\delta_{m,n_L} = m^{1/2}$ , we thus have that  $\sum_{s=1}^{\infty} \delta_{m,n_L}^{-s} \mathbf{\Lambda}_{\delta}^s \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$  is at most of order  $O_p(m^{-1})$  component-wise.

Next, consider the case when  $mn_U^{-1} \rightarrow \infty$  and  $mn_L^{-2} \rightarrow 0$ . From previous derivations, we know in this setting it holds that  $\mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$  is of order  $O_p(m^{-1/2})$  for the first  $p$  components and  $O_p(n_L^{-1/2})$  for the last  $mp$  components. By the two properties of  $\mathbf{\Lambda}_{\delta}^s$  noted above, we therefore obtain that  $\mathbf{\Lambda}_{\delta}^s \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$  is at most  $O_p(n_L^{-1/2})$  component-wise for any  $s$ . Hence  $\sum_{s=1}^{\infty} \delta_{m,n_L}^{-s} \mathbf{\Lambda}_{\delta}^s \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}}) = \delta_{m,n_L}^{-1} O_p(n_L^{-1/2}) = o_p(n_L^{-1/2})$  for sufficiently large  $m, n_L$ , by the properties of a geometric sum. This is sufficient to show that the last  $mp$  components of  $\sum_{s=1}^{\infty} \mathbf{\Lambda}^s \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$  are of smaller order component-wise than  $\mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$ , so that the result for the prediction gap holds. In particular, we thus know that  $\hat{\mathbf{b}} - \dot{\mathbf{b}} = O_p(n_L^{-1/2})$ . Furthermore, we also know that the convergence rate of  $\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}}$  is at least  $O_p(n_L^{-1/2})$ . As a result, we can set  $\delta_{m,n_L} = n_L^{1/2}$  without affecting the component-wise order properties of  $\mathbf{\Lambda}_{\delta}$ . Applying  $\delta_{m,n_L} = n_L^{1/2}$ , we thus have that  $\sum_{s=1}^{\infty} \delta_{m,n_L}^{-s} \mathbf{\Lambda}_{\delta}^s \mathbf{B}^{-1} \nabla Q(\dot{\boldsymbol{\theta}})$  is at most of order  $O_p(n_L^{-1})$  component-wise. This is smaller than  $O_p(m^{-1/2})$  when  $mn_U^{-1} \rightarrow \infty$ ,  $mn_L^{-2} \rightarrow 0$  and the result follows.

## S4 Unpartnered Fixed Effects

### S4.1 Generalised Linear Models

In the special case when  $\dot{\mathbf{G}} = \mathbf{0}_{p \times p}$ , i.e., all fixed effects are unpartnered in the true data generating process, the GLMM reduces to a GLM. We may then obtain a result based on a special case of our results in the conditional case, when all the true random effects are equal to zero. The result is as follows.

**Corollary A1.** Assume Conditions (C1) - (C5) are satisfied and  $mn_L^{-1} \rightarrow 0$ . Then as  $m, n_L \rightarrow \infty$  and when the true vector of random effects  $\dot{\mathbf{b}} = \mathbf{0}_{mp}$ , it holds that  $\mathbf{AD}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \xrightarrow{D} N(\mathbf{0}, \boldsymbol{\Omega})$ .

### S4.2 Linear Mixed Models

Suppose for  $i = 1, \dots, m$  and  $j = 1, \dots, n_i$  we observe data from the model  $y_{ij} = \mathbf{x}_{ij}^\top \boldsymbol{\beta} + \mathbf{z}_{ij}^\top \mathbf{b}_i + \mathbf{x}_{ij}^{(O)\top} \boldsymbol{\beta}^{(O)} + \epsilon_{ij}$ , where  $\mathbf{x}_{ij} = \mathbf{z}_{ij}$  for all  $(i, j)$ ,  $\mathbf{b}_i \stackrel{i.i.d.}{\sim} N(\mathbf{0}, \hat{\mathbf{G}})$  and  $\epsilon_{ij} \stackrel{i.i.d.}{\sim} N(0, \phi)$ . Note that this is part of the exponential family. Partition  $\boldsymbol{\beta} = (\boldsymbol{\beta}^{(P)\top}, \boldsymbol{\beta}^{(U)\top})^\top$ , corresponding to the  $p_P$  partnered and  $p_U$  unpartnered fixed effects ( $p_P + p_U = p$ ). That is, if we partition  $\mathbf{b}_i = (\mathbf{b}_i^{(P)\top}, \mathbf{b}_i^{(U)\top})^\top$ , then  $\mathbf{b}_i^{(U)} = \mathbf{0}_{p_U}$  for all  $i$ , and the corresponding elements in  $\dot{\mathbf{G}}$  are zero. Let  $\boldsymbol{\theta}^\times = (\boldsymbol{\beta}^\top, \mathbf{b}_1^{(P)\top}, \dots, \mathbf{b}_m^{(P)\top}, \mathbf{b}_1^{(U)\top}, \dots, \mathbf{b}_m^{(U)\top}, \boldsymbol{\beta}^{(O)\top})^\top$ ,

$$\boldsymbol{\theta}^- = (\boldsymbol{\beta}^\top, \mathbf{b}^\top, \boldsymbol{\beta}^{(O)\top})^\top,$$

$$\mathbf{D}^\times = \text{diag}(m^{1/2}\mathbf{1}_{p_P}, N^{1/2}\mathbf{1}_{p_U}, m^{1/2}\mathbf{1}_{mp_P}, n_1^{1/2}\mathbf{1}_{p_U}, \dots, n_m^{1/2}\mathbf{1}_{p_U}, N^{1/2}\mathbf{1}_{p_O}), \text{ and}$$

$$\mathbf{D}^- = \text{diag}(m^{1/2}\mathbf{1}_{p_P}, N^{1/2}\mathbf{1}_{p_U}, n_1^{1/2}\mathbf{1}_p, \dots, n_m^{1/2}\mathbf{1}_p, N^{1/2}\mathbf{1}_{p_O}). \text{ Also let } \mathbf{X}_i^{(O)} =$$

$$[\mathbf{x}_{i1}^{(O)}, \dots, \mathbf{x}_{ini}^{(O)}]^\top \text{ and } \mathbf{X}^{(O)} = [\mathbf{X}_1^{(O)\top}, \dots, \mathbf{X}_m^{(O)\top}]^\top. \text{ The } p_O \text{ orthogonal}$$

fixed effects  $\mathbf{x}_{ij}^{(O)}$  satisfy  $\mathbf{X}^{(O)\top} \mathbf{Z} = \mathbf{0}_{p_O \times mp}$ , for example orthogonal poly-

nomials of  $\mathbf{x}_{ij}$ . This implies  $\mathbf{X}_i^{(O)\top} \mathbf{X}_i = \mathbf{0}_{p_O \times p}$  for all  $i$ . For a  $q \times \{(m +$

$1)p + p_O\}$  matrix  $\mathbf{A}^*$  with the finite selection property, we have the following.

**Corollary A2.** Assume Conditions (C1) - (C4) are satisfied. Then as

$m, n_L \rightarrow \infty$  and unconditional on the random effects  $\dot{\mathbf{b}}$ , it holds that

1.  $\mathbf{A}^* \mathbf{D}^\times (\hat{\boldsymbol{\theta}}^\times - \dot{\boldsymbol{\theta}}^\times) \xrightarrow{D} N(\mathbf{0}, \boldsymbol{\Omega}_a)$  if  $mn_L^{-1} \rightarrow 0$ , and

2.  $\mathbf{A}^* \mathbf{D}^- (\hat{\boldsymbol{\theta}}^- - \dot{\boldsymbol{\theta}}^-) \xrightarrow{D} N(\mathbf{0}, \boldsymbol{\Omega}_b)$  if  $mn_U^{-1} \rightarrow \infty$ ,

where

$$\boldsymbol{\Omega}_a = \lim_{m, n_L \rightarrow \infty} \mathbf{A}^* \begin{bmatrix} \dot{\mathbf{G}}_{[1:p_P, 1:p_P]} & \mathbf{0}_{p_P \times p_U} & \mathbf{1}_m^\top \otimes \dot{\mathbf{G}}_{[1:p_P, 1:p_P]} & \mathbf{0}_{p_P \times mp_U} & \mathbf{0}_{p_P \times p_O} \\ \mathbf{0}_{p_U \times p_P} & \boldsymbol{\Omega}_1 & \mathbf{0}_{p_U \times mp_P} & \mathbf{0}_{p_U \times mp_U} & \mathbf{0}_{p_U \times p_O} \\ \mathbf{1}_m \otimes \dot{\mathbf{G}}_{[1:p_P, 1:p_P]} & \mathbf{0}_{mp_P \times p_U} & \mathbf{1}_{m \times m} \otimes \dot{\mathbf{G}}_{[1:p_P, 1:p_P]} & \mathbf{0}_{mp_P \times mp_U} & \mathbf{0}_{mp_P \times p_O} \\ \mathbf{0}_{mp_U \times p_P} & \mathbf{0}_{mp_U \times p_U} & \mathbf{0}_{mp_U \times mp_P} & \boldsymbol{\Omega}_2 & \mathbf{0}_{mp_U \times p_O} \\ \mathbf{0}_{p_O \times p_P} & \mathbf{0}_{p_O \times p_U} & \mathbf{0}_{p_O \times mp_P} & \mathbf{0}_{p_O \times mp_U} & \boldsymbol{\Omega}_3 \end{bmatrix} \mathbf{A}^{*\top}$$

$$\Omega_b = \lim_{m, n_L \rightarrow \infty} \mathbf{A}^* \begin{bmatrix} \dot{\mathbf{G}}_{[1:p_P, 1:p_P]} & \mathbf{0}_{p_P \times p_U} & \mathbf{0}_{p_P \times mp} & \mathbf{0}_{p_P \times p_O} \\ \mathbf{0}_{p_U \times p_P} & \Omega_1 & \mathbf{0}_{p_U \times mp} & \mathbf{0}_{p_U \times p_O} \\ \mathbf{0}_{mp \times p_P} & \mathbf{0}_{mp \times p_U} & \Omega_4 & \mathbf{0}_{mp \times p_O} \\ \mathbf{0}_{p_O \times p_P} & \mathbf{0}_{p_O \times p_U} & \mathbf{0}_{p_O \times mp} & \Omega_3 \end{bmatrix} \mathbf{A}^{*\top}$$

$$\Omega_1 = \left\{ \frac{\dot{\phi}}{m} \sum_{i=1}^m \frac{n}{n_i} \left( \frac{\mathbf{X}_i^\top \mathbf{X}_i}{n_i} \right)^{-1} \right\}_{[(p-p_U+1):p, (p-p_U+1):p]}$$

$$\Omega_2 = \text{bdiag} \left[ \left\{ \dot{\phi} \left( \frac{\mathbf{X}_1^\top \mathbf{X}_1}{n_1} \right)^{-1} \right\}_{[(p-p_U+1):p, (p-p_U+1):p]}, \dots, \left\{ \dot{\phi} \left( \frac{\mathbf{X}_m^\top \mathbf{X}_m}{n_m} \right)^{-1} \right\}_{[(p-p_U+1):p, (p-p_U+1):p]} \right]$$

$$\Omega_3 = \dot{\phi} \left( \frac{\mathbf{X}^{(O)\top} \mathbf{X}^{(O)}}{N} \right)^{-1}$$

$$\Omega_4 = \text{bdiag} \left[ \left\{ \dot{\phi} \left( \frac{\mathbf{X}_1^\top \mathbf{X}_1}{n_1} \right)^{-1} \right\}, \dots, \left\{ \dot{\phi} \left( \frac{\mathbf{X}_m^\top \mathbf{X}_m}{n_m} \right)^{-1} \right\} \right].$$

*Proof.* We use the same approach as previous proofs and examine the Taylor expansion (S2.1). In this case, we have the expressions

$$\nabla Q(\dot{\boldsymbol{\theta}}) = \begin{bmatrix} \dot{\phi}^{-1} \mathbf{X}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}}) \\ \dot{\phi}^{-1} \mathbf{Z}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}}) - (\mathbf{I}_m \otimes \hat{\mathbf{G}}^{-1}) \dot{\mathbf{b}} \\ \dot{\phi}^{-1} \mathbf{X}^{(O)\top} (\mathbf{y} - \dot{\boldsymbol{\mu}}) \end{bmatrix},$$

$$\mathbf{B}(\dot{\boldsymbol{\theta}}) = -\nabla^2 Q(\dot{\boldsymbol{\theta}}) = \begin{bmatrix} \mathbf{X}^\top \dot{\mathbf{W}} \mathbf{X} & \mathbf{X}^\top \dot{\mathbf{W}} \mathbf{Z} & \mathbf{X}^\top \dot{\mathbf{W}} \mathbf{X}^{(O)} \\ \mathbf{Z}^\top \dot{\mathbf{W}} \mathbf{X} & \mathbf{Z}^\top \dot{\mathbf{W}} \mathbf{Z} + \mathbf{I}_m \otimes \hat{\mathbf{G}}^{-1} & \mathbf{Z}^\top \dot{\mathbf{W}} \mathbf{X}^{(O)} \\ \mathbf{X}^{(O)\top} \dot{\mathbf{W}} \mathbf{X} & \mathbf{X}^{(O)\top} \dot{\mathbf{W}} \mathbf{Z} & \mathbf{X}^{(O)\top} \dot{\mathbf{W}} \mathbf{X}^{(O)} \end{bmatrix}$$

$$= \dot{\phi}^{-1} \begin{bmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \mathbf{Z} & \mathbf{0}_{p \times p_O} \\ \mathbf{Z}^\top \mathbf{X} & \mathbf{Z}^\top \mathbf{Z} + \mathbf{I}_m \otimes \hat{\mathbf{G}}^{-1} & \mathbf{0}_{mp \times p_O} \\ \mathbf{0}_{p_O \times p} & \mathbf{0}_{p_O \times mp} & \mathbf{X}^{(O)\top} \mathbf{X}^{(O)} \end{bmatrix},$$

where the last equality follows from the fact that  $\dot{\mathbf{W}} = \dot{\phi}^{-1} \mathbf{I}_N$  and  $\mathbf{X}^{(O)\top} \mathbf{Z} = \mathbf{0}_{p_O \times mp}$ . Since  $\mathbf{B}(\tilde{\boldsymbol{\theta}})$  is block diagonal, we thus know that expressions (S2.9) and (S2.10) still hold. Recall

$$\begin{aligned} \hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}} &= m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} \dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i) + m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i \\ &\quad - m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \hat{\mathbf{G}}^{-1} \dot{\mathbf{b}}_i + \frac{1}{2} \{\mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}})\}_{[1:p]} \\ &\quad + O_p(n_L^{-1}) \left\{ m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i) + m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i \right. \\ &\quad \left. - m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i + \hat{\mathbf{G}}^{-1})^{-1} \hat{\mathbf{G}}^{-1} \dot{\mathbf{b}}_i \right\} + m^{-1} \sum_{i=1}^m O_p(n_L^{-2}) \dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i) \end{aligned}$$

and

$$\begin{aligned} \hat{\mathbf{b}} - \dot{\mathbf{b}} &= -\mathbf{1}_m \otimes m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i + \mathbf{B}_3^{-1} \{\dot{\phi}^{-1} \mathbf{Z}^\top (\mathbf{y} - \dot{\boldsymbol{\mu}})\} \\ &\quad + O_p(N^{-1/2}) + O_p(n_L^{-1}) + \frac{1}{2} \{\mathbf{B}^{-1} \mathbf{R}(\tilde{\boldsymbol{\theta}})\}_{[p+1:(m+1)p]}. \end{aligned}$$

In the LMM case, the remainder term in the Taylor expansion is zero.

Thus the dominating term on the right hand side for  $\hat{\boldsymbol{\beta}}^{(U)} - \dot{\boldsymbol{\beta}}^{(U)}$  are the last  $p_U$  components of  $m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} \dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i)$ , since

the last  $p_U$  components of  $m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i$  are zero. Noting that  $\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i = (\epsilon_{i1}, \dots, \epsilon_{in_i})^\top =: \boldsymbol{\epsilon}_i$ , the result for the unpartnered fixed effects follows after normalising by  $N^{1/2}$ .

Next, again from the Taylor expansion we have from the block-diagonal structure of  $\mathbf{B}(\dot{\boldsymbol{\theta}})$  that  $\hat{\boldsymbol{\beta}}^{(O)} - \dot{\boldsymbol{\beta}}^{(O)} = (\mathbf{X}^{(O)\top} \mathbf{X}^{(O)})^{-1} \mathbf{X}^{(O)\top} (\mathbf{y} - \dot{\boldsymbol{\mu}})$  and the result follows after normalising by  $N^{1/2}$  since  $\mathbf{y} - \dot{\boldsymbol{\mu}} = (\epsilon_{11}, \dots, \epsilon_{mn_m})^\top =: \boldsymbol{\epsilon}$ .

Finally, the result for the unpartnered random effects follows from the fact that the last  $p_U$  components of  $m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i$  are zero so that the dominating term on the right hand side is  $(\mathbf{X}_i^\top \mathbf{X}_i)^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i)$ , and normalising by  $n_i$ .

The proofs for the partnered fixed and random effects are analogous to the proofs of Theorems 2 and 4, based on examining the leading term in the Taylor expansion.

For the joint behaviour of the estimator, we examine the joint behaviour of the leading terms in the Taylor Expansion. Note that  $\boldsymbol{\epsilon}$  is multivariate normal with covariance matrix  $\dot{\phi} \mathbf{I}_N$ ,  $\dot{\mathbf{b}}$  is multivariate normal with covariance matrix  $\mathbf{I}_m \otimes \dot{\mathbf{G}}$ ,  $\boldsymbol{\epsilon}$  and  $\dot{\mathbf{b}}$  are independent, and all the leading terms in the Taylor expansion are linear functions of  $\boldsymbol{\epsilon}$  and  $\dot{\mathbf{b}}$ . To determine the joint behaviour of the estimator it is thus sufficient to derive the limiting covariance between the normalised leading terms, as we see (from the lead-



ing terms) that the estimator itself is also (asymptotically) multivariate normal. For example,

$$\begin{aligned}
& \text{Cov} \left\{ N^{1/2} m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \mathbf{X}_i)^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \boldsymbol{\mu}_i), N^{1/2} (\mathbf{X}^{(O)\top} \mathbf{X}^{(O)})^{-1} \mathbf{X}^{(O)\top} (\mathbf{y} - \boldsymbol{\mu}) \right\} \\
&= n \dot{\phi} \sum_{i=1}^m (\mathbf{X}_i^\top \mathbf{X}_i)^{-1} \mathbf{X}_i^\top \mathbf{X}_i^{(O)} (\mathbf{X}^{(O)\top} \mathbf{X}^{(O)})^{-1} \\
&= \mathbf{0}_{p \times p_O}
\end{aligned}$$

due to the mutual independence of the  $\epsilon_{ij}$  and orthogonality condition of  $\mathbf{X}^{(O)}$ . The pairwise limiting covariances between the leading terms can all be derived in a similar way and the result follows. Notice here that quantities with different convergence rates are always asymptotically uncorrelated and independent in this case.

□

Note that the results hold by the Lindeberg-Feller Central Limit Theorem even if the true distribution of  $\epsilon_{ij}$  is not normal, as long as it is mean zero with finite variance. Also note that condition (C5) is no longer required, and that there is no restriction on the relative rates of  $m$  and  $n_L$ , since there is no remainder term to deal with. Our result is consistent with the results derived in Lyu and Welsh (2021a,b) who also derive a  $N^{1/2}$  convergence rate for unpartnered fixed effects that are time-varying.

In practice, we do not know if a fixed effect is truly partnered with a random effect or not, and therefore the correct asymptotic distribution and convergence rate is also unknown. In this case, an appropriate finite sample approximation, given consistent estimators  $\tilde{\mathbf{G}}$  and  $\tilde{\phi}$  of  $\dot{\mathbf{G}}$  and  $\dot{\phi}$  respectively, is

$$\hat{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}} \sim N \left\{ \mathbf{0}, m^{-1} \tilde{\mathbf{G}} + N^{-1} \frac{\tilde{\phi}}{m} \sum_{i=1}^m \frac{n}{n_i} \left( \frac{\mathbf{X}_i^\top \mathbf{X}_i}{n_i} \right)^{-1} \right\},$$

which is based on the distribution of  $m^{-1} \sum_{i=1}^m (\mathbf{X}_i^\top \dot{\mathbf{W}}_i \mathbf{X}_i)^{-1} \dot{\phi}^{-1} \mathbf{X}_i^\top (\mathbf{y}_i - \dot{\boldsymbol{\mu}}_i) + m^{-1} \sum_{i=1}^m \dot{\mathbf{b}}_i$ , noting that the two terms are independent.

## S5 Additional Simulation Results

### S5.1 Main Results for the Conditional Regime

Figures 1, 2 and 3 display the empirical coverage probabilities and results from applying the Shapiro-Wilk test, respectively, under the conditional regime and for the 25 combinations of  $(m, n)$ . Although our coverage intervals often undercovered or overcovered for small cluster sizes e.g.,  $n = 25$ , especially for the Bernoulli case, they all moved toward nominal coverage as  $n$  becomes larger than  $m$ . This is consistent with Theorem 1. The fact the empirical coverage probabilities were slow in tending towards the nominal 95% level was also not overly surprising, as the third derivative term in the corresponding Taylor expansion is  $O_p(m^{1/2}n_L^{-1/2})$ . The Shapiro-Wilk tests overall did not indicate any evidence of deviations away from normality when  $m < n$ , although there were occasionally a few  $p$ -values less than 0.05. Overall, these results strongly support the use of Theorem 1 for inference under the conditional regime.

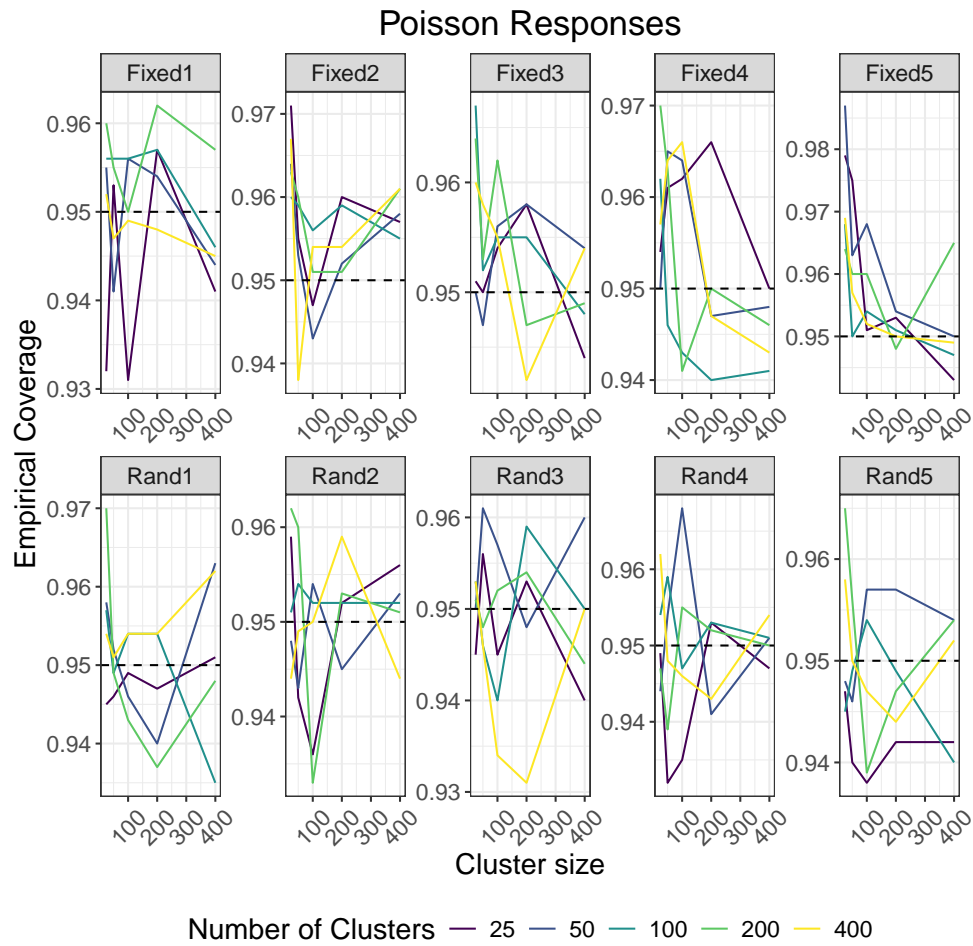


Figure 1: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the conditional regime with Poisson responses.

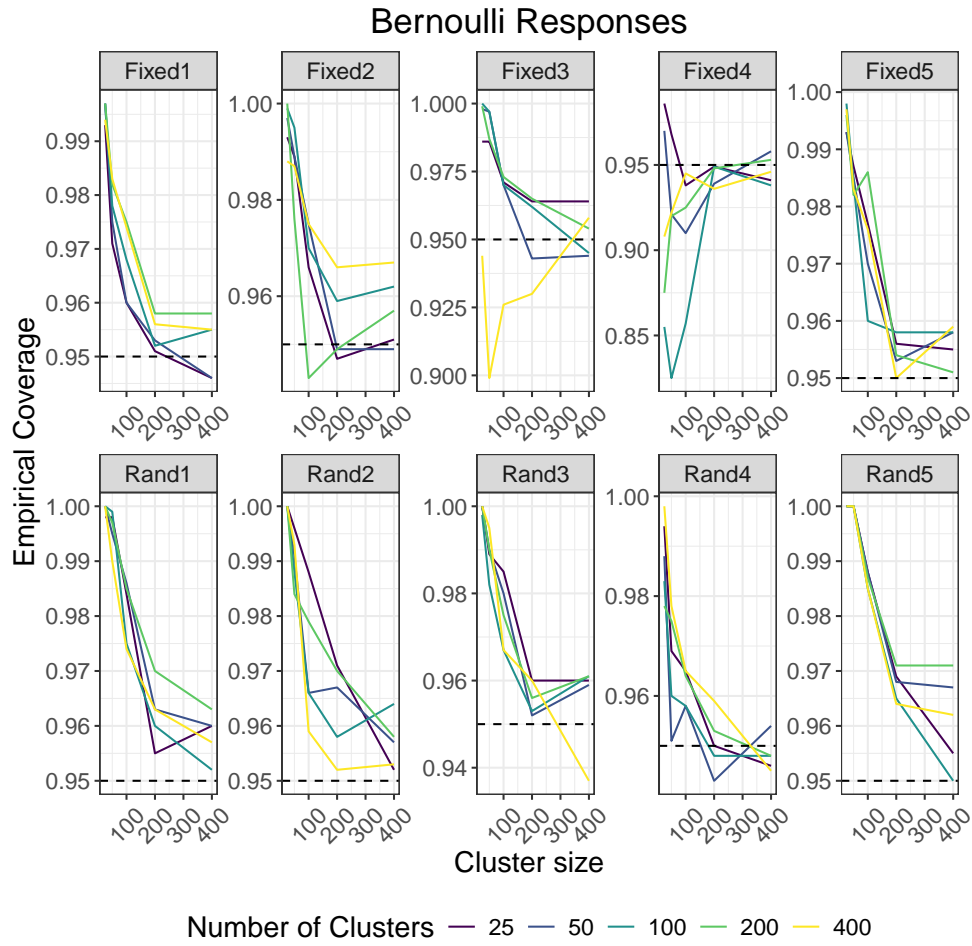


Figure 2: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the conditional regime with Bernoulli responses.

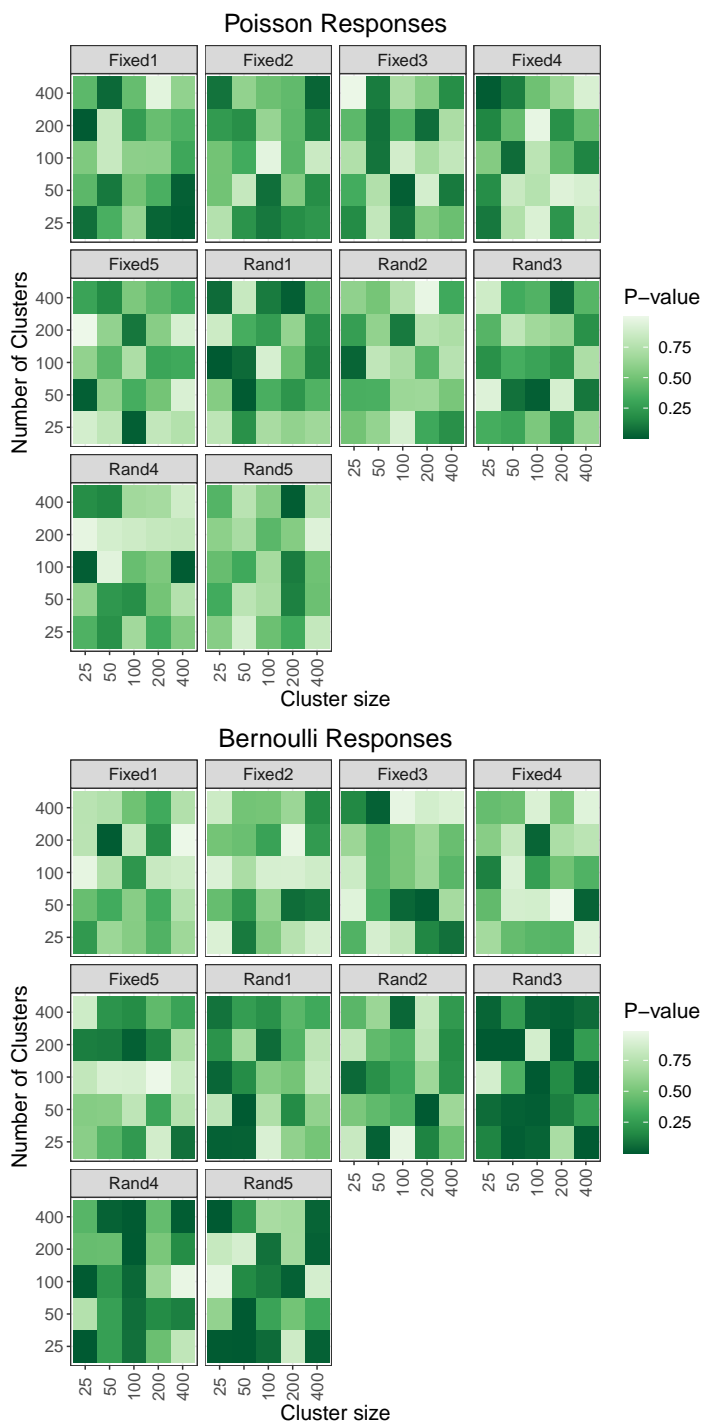


Figure 3:  $p$ -values from Shapiro-Wilk tests applied to the fixed and random effects estimates obtained using maximum PQL estimation, under the conditional regime.

S5. ADDITIONAL SIMULATION RESULTS

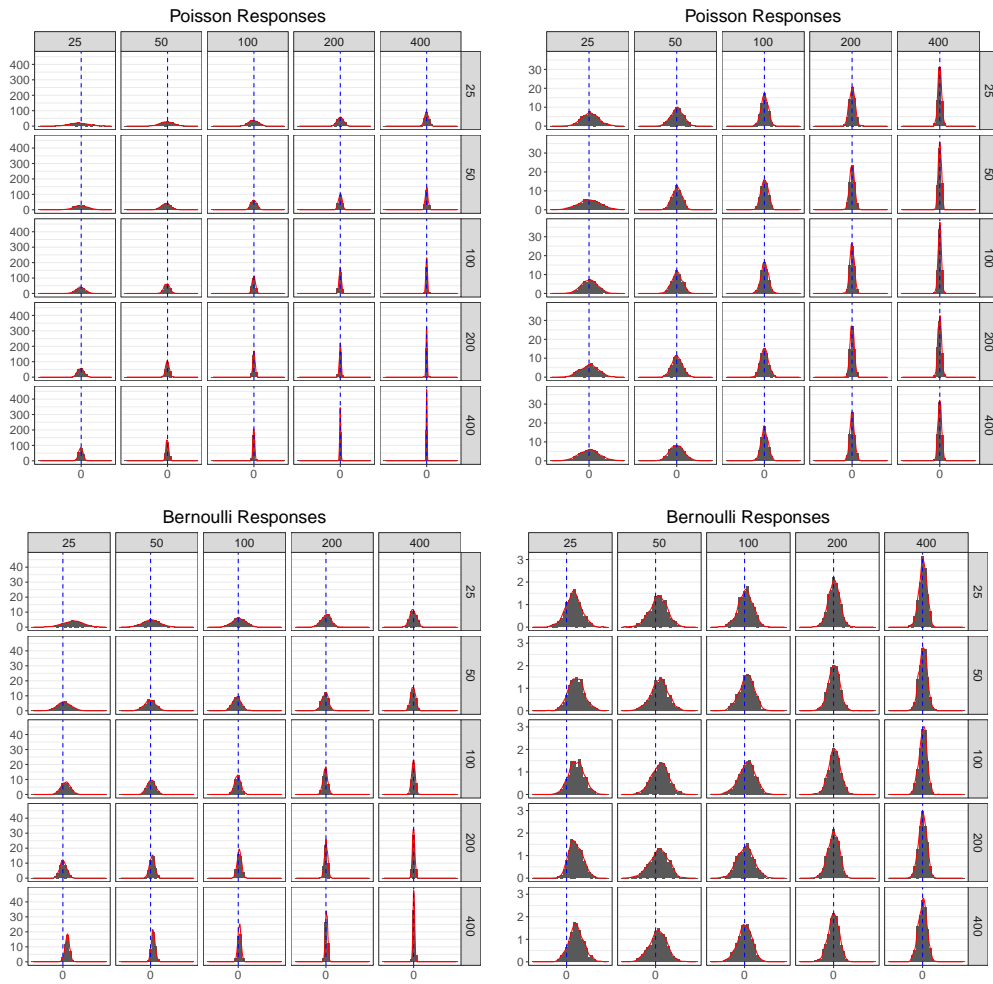


Figure 4: Histograms for the third components of  $\hat{\beta} - \beta$  (left panels) and  $\hat{b}_1 - b_1$  (right panels), under the conditional regime. Vertical facets represent the cluster sizes, while horizontal facets represent the number of clusters. The dotted blue line indicates zero, and the red curve is a kernel density smoother.

## S5.2 Frobenius Norm

Table 1: Empirical mean Frobenius norm of the difference between estimated and true random effects covariance matrix.

	$m$	Poisson					Bernoulli				
		$n = 25$	$n = 50$	$n = 100$	$n = 200$	$n = 400$	$n = 25$	$n = 50$	$n = 100$	$n = 200$	$n = 400$
$\hat{G} = m^{-1} \sum_{i=1}^m \hat{b}_i$	25	1.06	1.06	1.06	1.05	1.06	1.76	1.47	1.24	1.12	1.09
	50	0.77	0.75	0.75	0.76	0.75	1.79	1.40	1.03	0.83	0.77
	100	0.54	0.54	0.54	0.54	0.54	1.80	1.38	0.89	0.63	0.56
	200	0.39	0.38	0.38	0.38	0.38	1.80	1.35	0.82	0.51	0.41
	400	0.27	0.27	0.27	0.27	0.27	1.81	1.35	0.78	0.43	0.32
$\hat{G} = 0.25\mathbf{I}_2$	25	1.02	1.01	1.03	1.05	1.04	1.90	1.71	1.47	1.23	1.04
	50	0.73	0.73	0.74	0.74	0.76	1.90	1.70	1.44	1.15	0.91
	100	0.56	0.53	0.52	0.53	0.54	1.90	1.69	1.42	1.11	0.84
	200	0.44	0.39	0.37	0.38	0.38	1.90	1.68	1.41	1.09	0.79
	400	0.38	0.29	0.27	0.27	0.27	1.89	1.68	1.40	1.08	0.77
$\hat{G} = 0.5\mathbf{I}_2$	25	1.02	1.03	1.04	1.05	1.05	1.61	1.39	1.16	1.01	0.96
	50	0.74	0.75	0.75	0.75	0.74	1.61	1.34	1.07	0.86	0.75
	100	0.53	0.52	0.54	0.54	0.54	1.60	1.32	1.03	0.77	0.61
	200	0.39	0.38	0.38	0.38	0.38	1.59	1.31	1.01	0.73	0.52
	400	0.30	0.27	0.27	0.27	0.27	1.59	1.30	0.99	0.70	0.47
$\hat{G} = \mathbf{I}_2$	25	1.06	1.05	1.04	1.04	1.06	1.21	1.06	0.98	0.97	1.00
	50	0.74	0.75	0.75	0.75	0.76	1.17	0.93	0.78	0.73	0.71
	100	0.53	0.53	0.54	0.53	0.54	1.13	0.86	0.65	0.56	0.53
	200	0.38	0.38	0.38	0.38	0.38	1.12	0.82	0.58	0.44	0.39
	400	0.27	0.27	0.27	0.27	0.27	1.10	0.80	0.55	0.38	0.30
$\hat{G} = 2\mathbf{I}_2$	25	1.06	1.06	1.06	1.05	1.06	0.84	0.98	1.03	1.04	1.05
	50	0.75	0.74	0.75	0.76	0.75	0.71	0.71	0.74	0.74	0.75
	100	0.54	0.54	0.54	0.53	0.53	0.56	0.51	0.53	0.53	0.54
	200	0.38	0.38	0.38	0.38	0.38	0.47	0.38	0.37	0.38	0.38
	400	0.27	0.27	0.27	0.27	0.27	0.42	0.29	0.27	0.27	0.27
$\hat{G} = 4\mathbf{I}_2$	25	1.06	1.06	1.06	1.06	1.06	1.33	1.42	1.25	1.16	1.11
	50	0.77	0.76	0.76	0.75	0.76	1.18	1.07	0.93	0.83	0.80
	100	0.55	0.54	0.54	0.54	0.54	0.97	0.86	0.70	0.61	0.57
	200	0.39	0.38	0.38	0.38	0.38	0.86	0.72	0.55	0.45	0.41
	400	0.27	0.27	0.27	0.27	0.27	0.80	0.65	0.46	0.34	0.30



**S5.3**  $\hat{\mathbf{G}} = 0.25 \mathbf{I}_2$ 

Using a large  $\hat{\mathbf{G}}$  of  $4\mathbf{I}_2$  had the least impact on the results, while a small  $\hat{\mathbf{G}}$ , e.g.,  $0.25\mathbf{I}_2$  had more of a noticeable impact at small sample sizes. This is not surprising since the latter corresponds to more shrinkage, such that larger sample sizes are needed before asymptotic results apply.

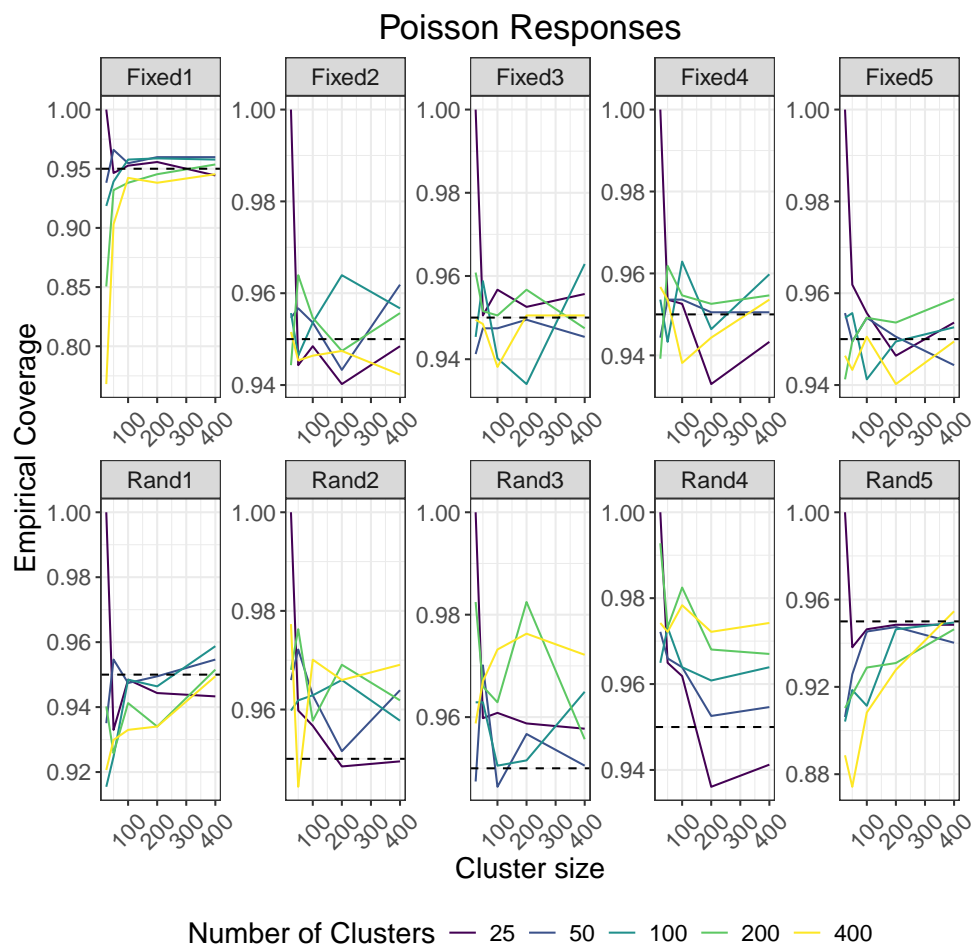


Figure 5: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the unconditional regime with Poisson responses.

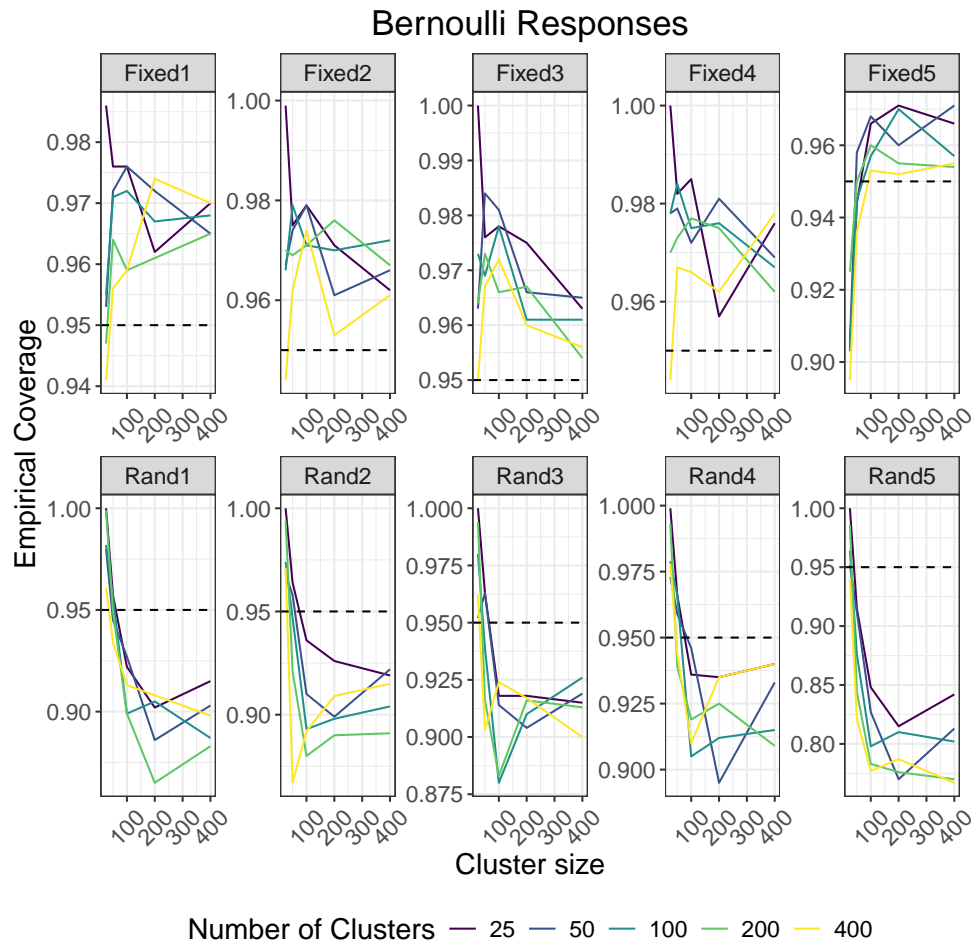


Figure 6: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the unconditional regime with Bernoulli responses.

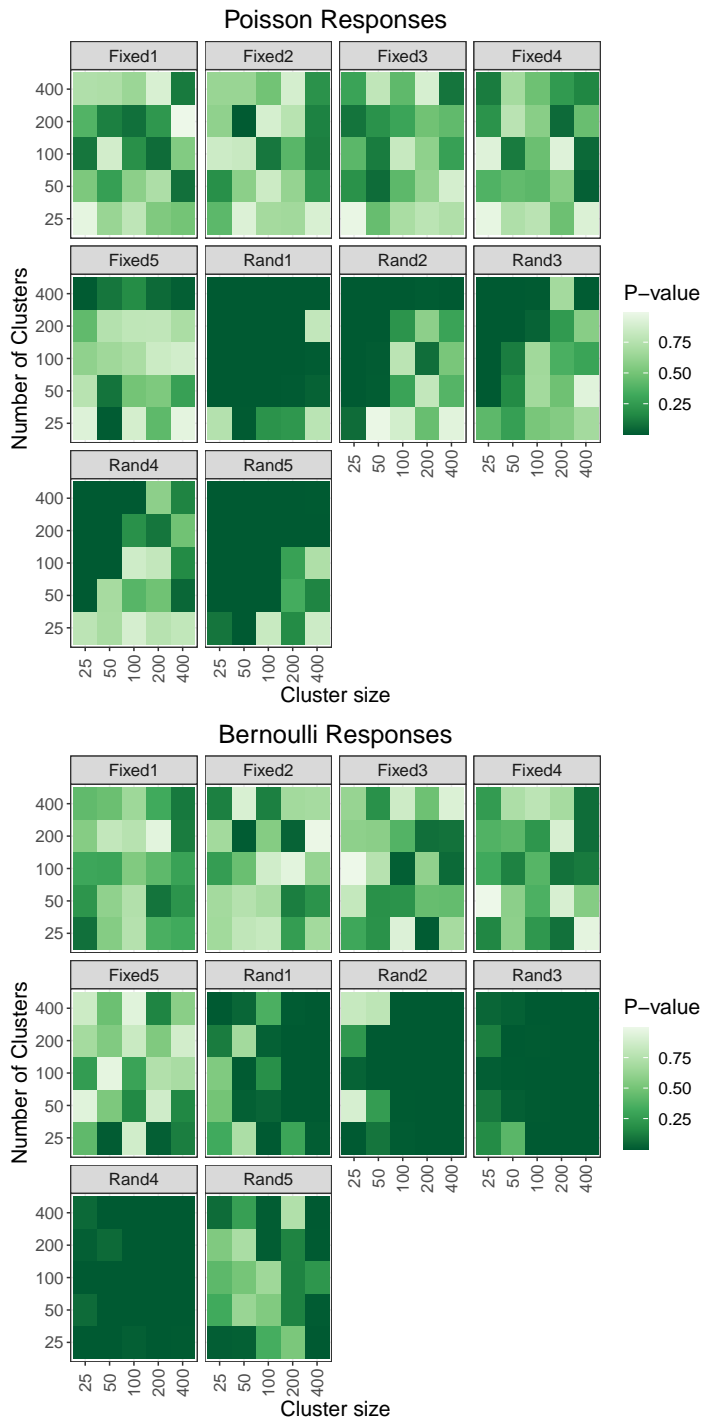


Figure 7:  $p$ -values from Shapiro-Wilk tests applied to the fixed and random effects estimates obtained using maximum PQL estimation, under the unconditional regime.

S5. ADDITIONAL SIMULATION RESULTS

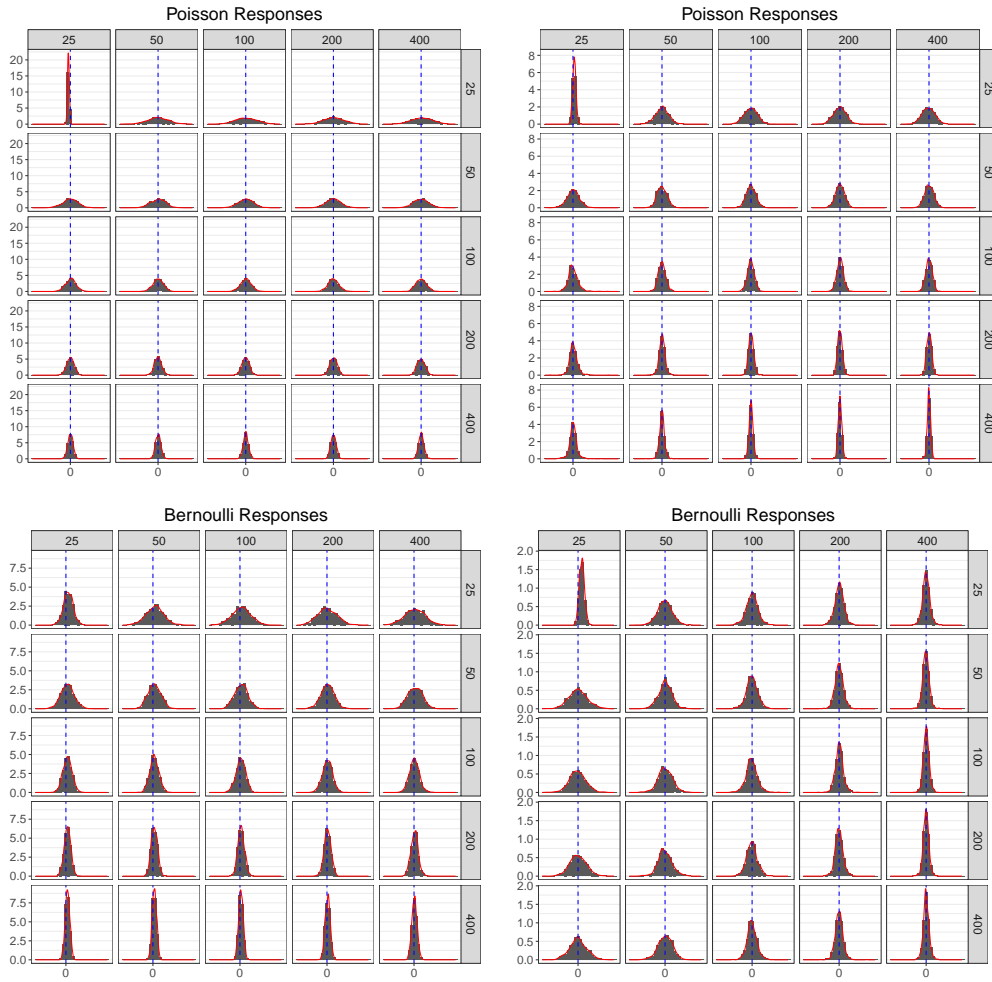


Figure 8: Histograms for the third components of  $\hat{\beta} - \beta$  (left panels) and  $\hat{b}_1 - b_1$  (right panels), under the unconditional regime. Vertical facets represent the cluster sizes, while horizontal facets represent the number of clusters. The dotted blue line indicates zero, and the red curve is a kernel density smoother.

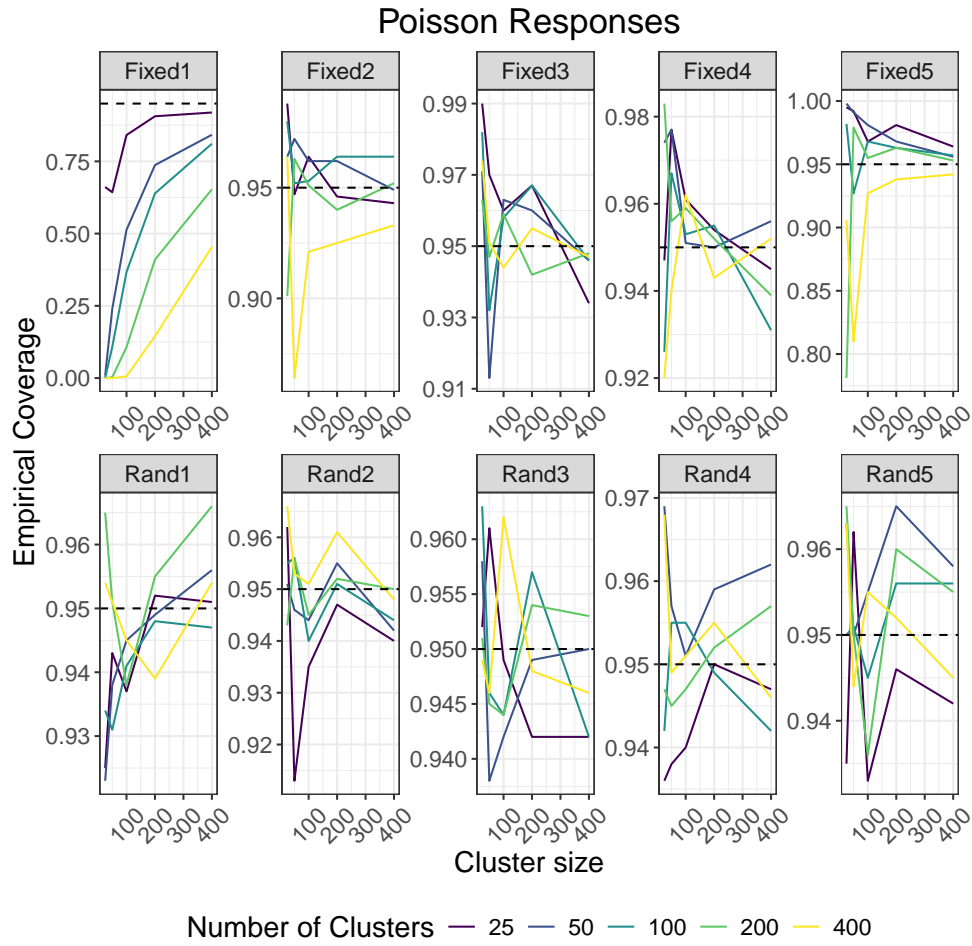


Figure 9: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the conditional regime with Poisson responses.

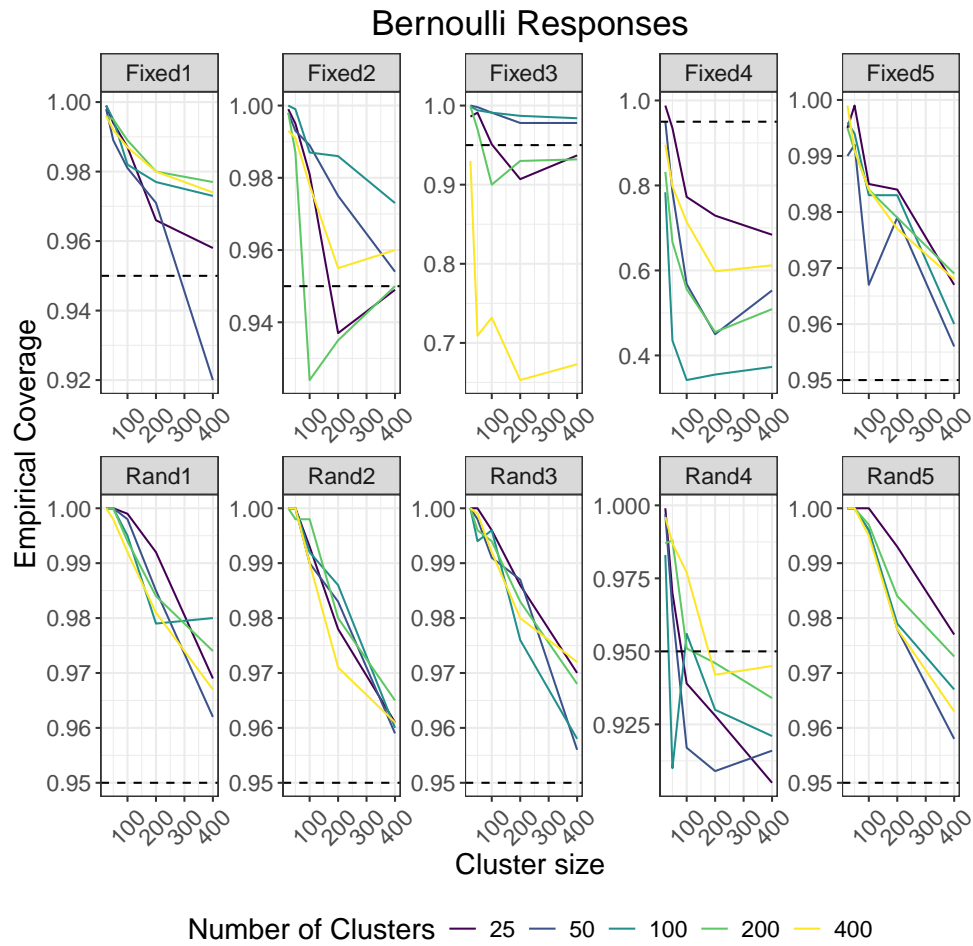


Figure 10: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the conditional regime with Bernoulli responses.

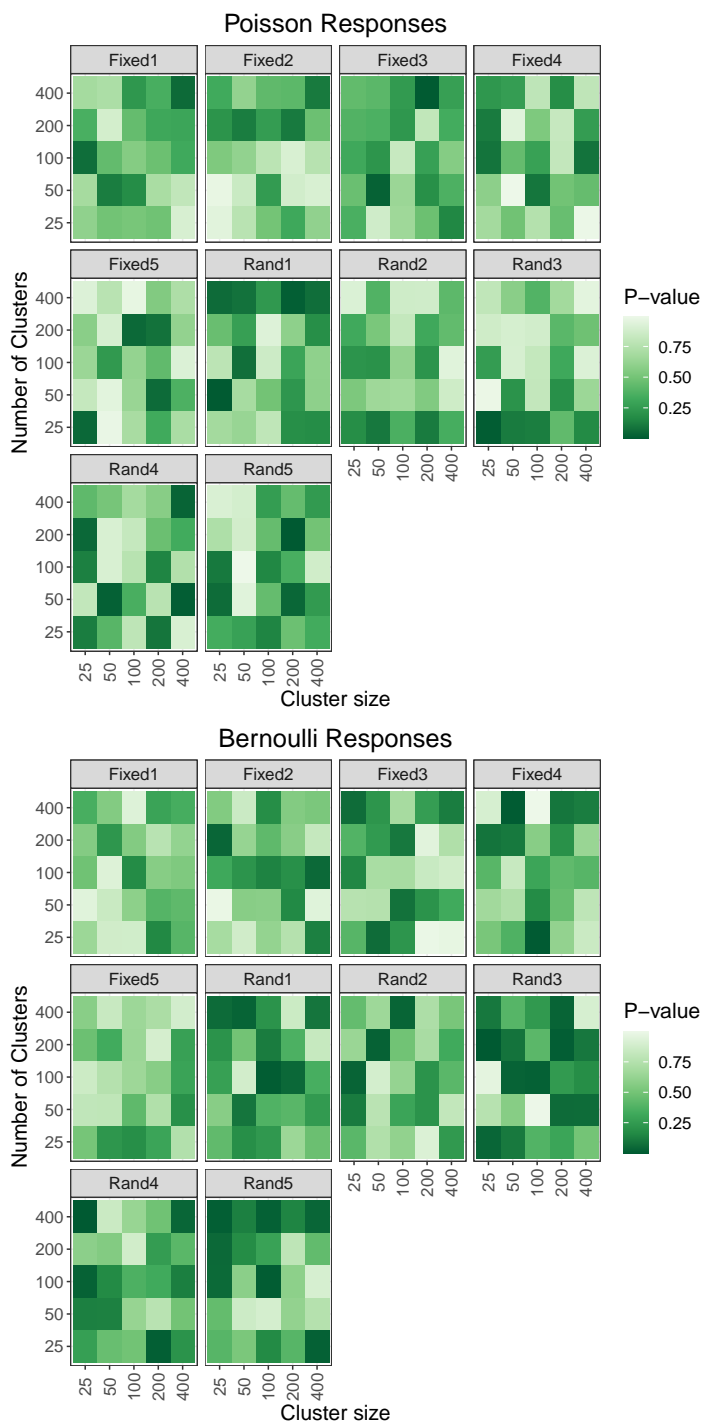


Figure 11:  $p$ -values from Shapiro-Wilk tests applied to the fixed and random effects estimates obtained using maximum PQL estimation, under the conditional regime.



S5. ADDITIONAL SIMULATION RESULTS

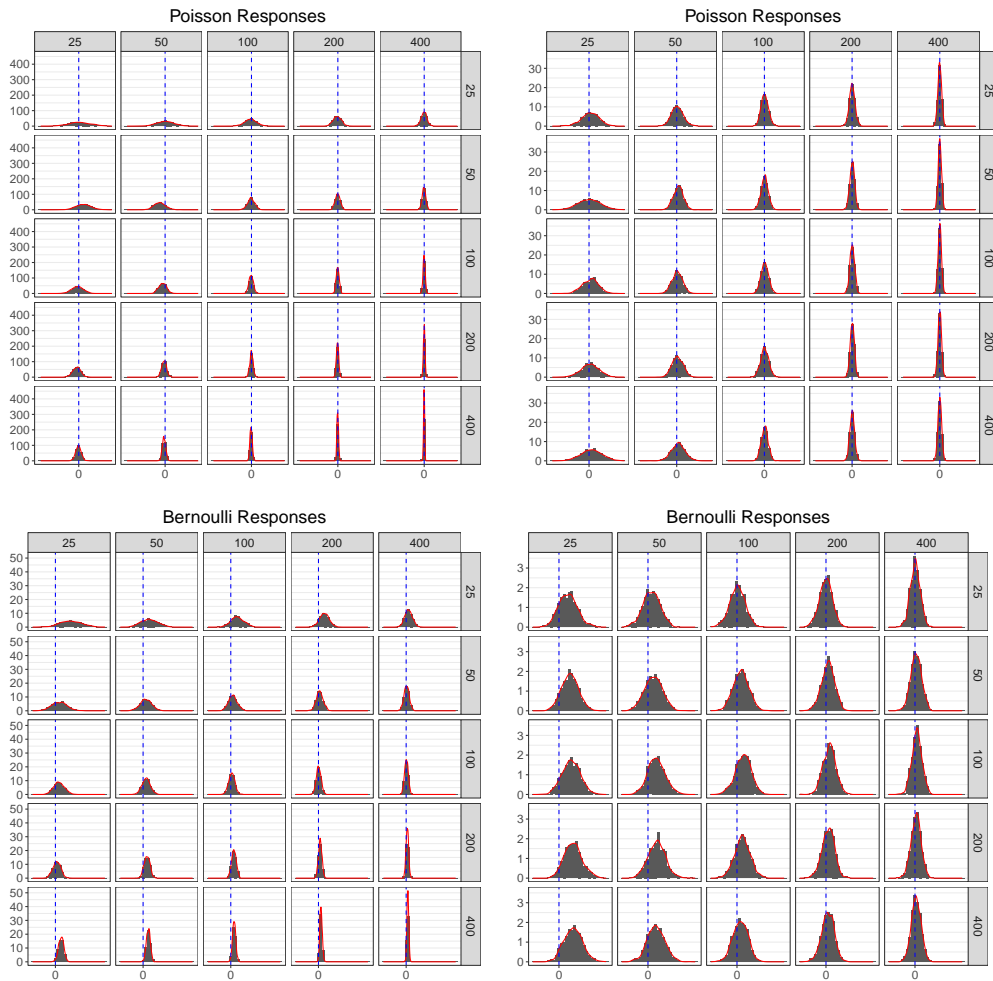


Figure 12: Histograms for the third components of  $\hat{\beta} - \beta$  (left panels) and  $\hat{b}_1 - b_1$  (right panels), under the unconditional regime. Vertical facets represent the cluster sizes, while horizontal facets represent the number of clusters. The dotted blue line indicates zero, and the red curve is a kernel density smoother.

S5.4  $\hat{G} = 0.50 I_2$

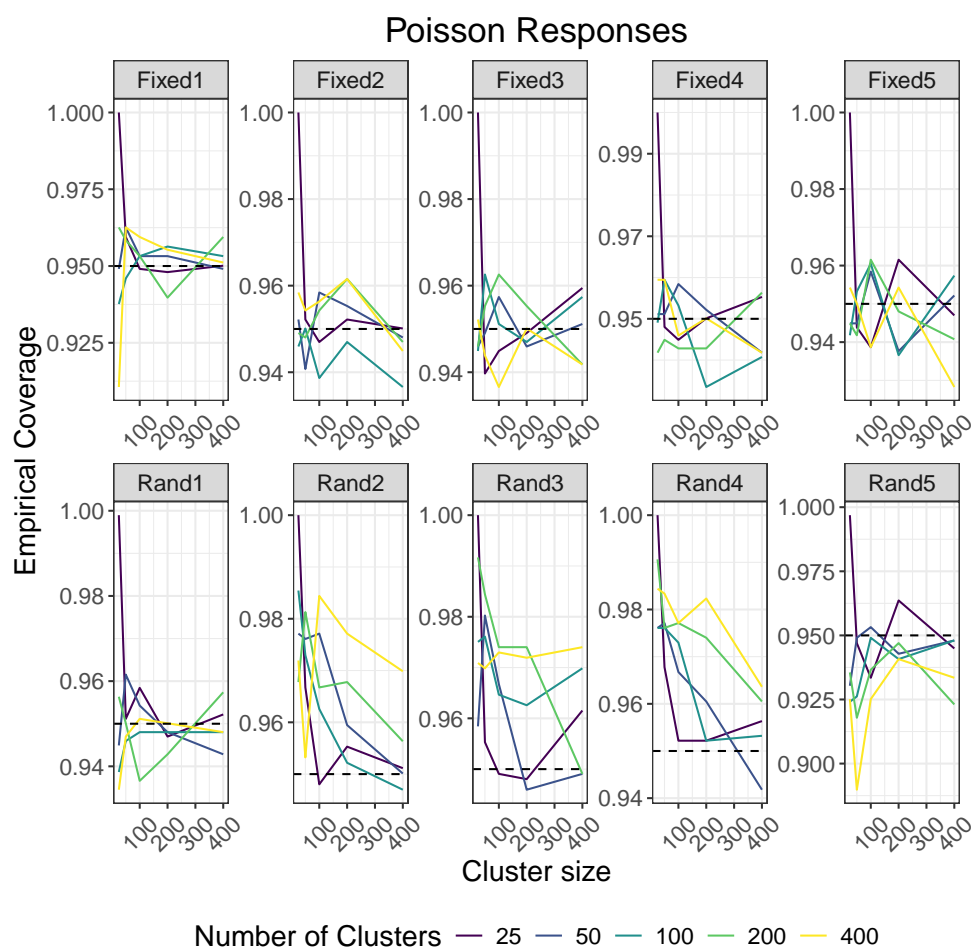


Figure 13: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the unconditional regime with Poisson responses.

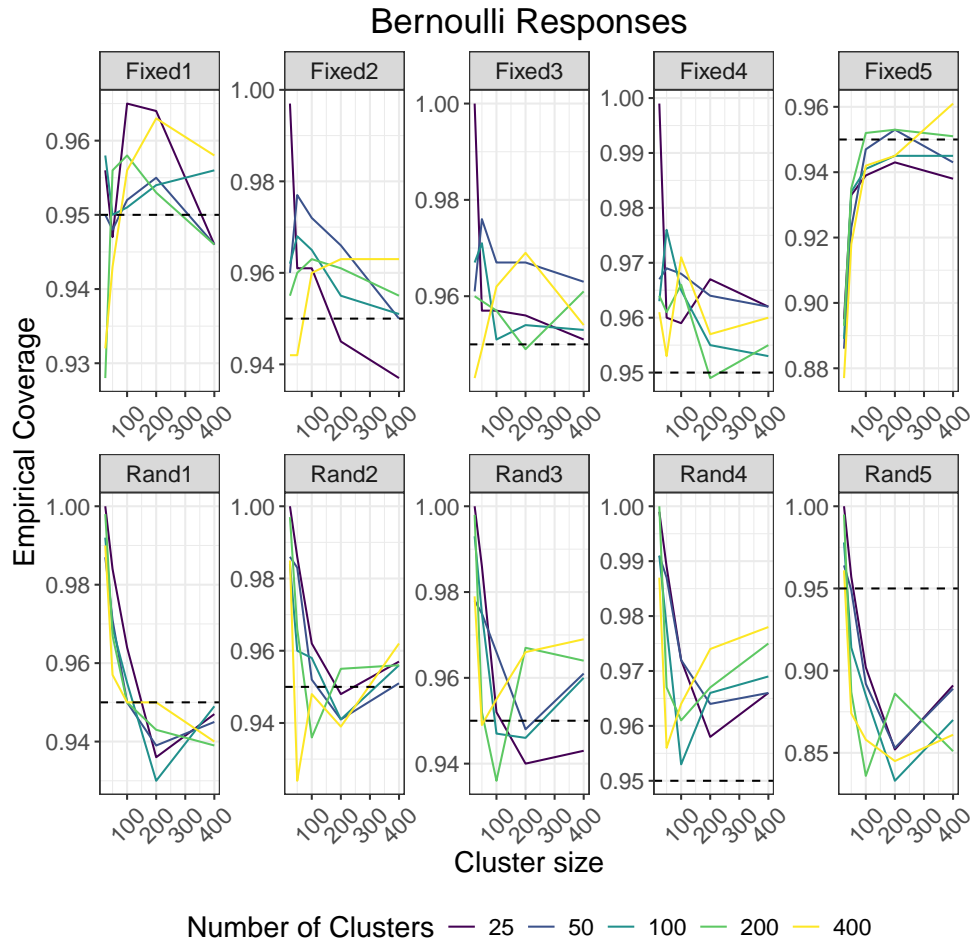


Figure 14: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the unconditional regime with Bernoulli responses.

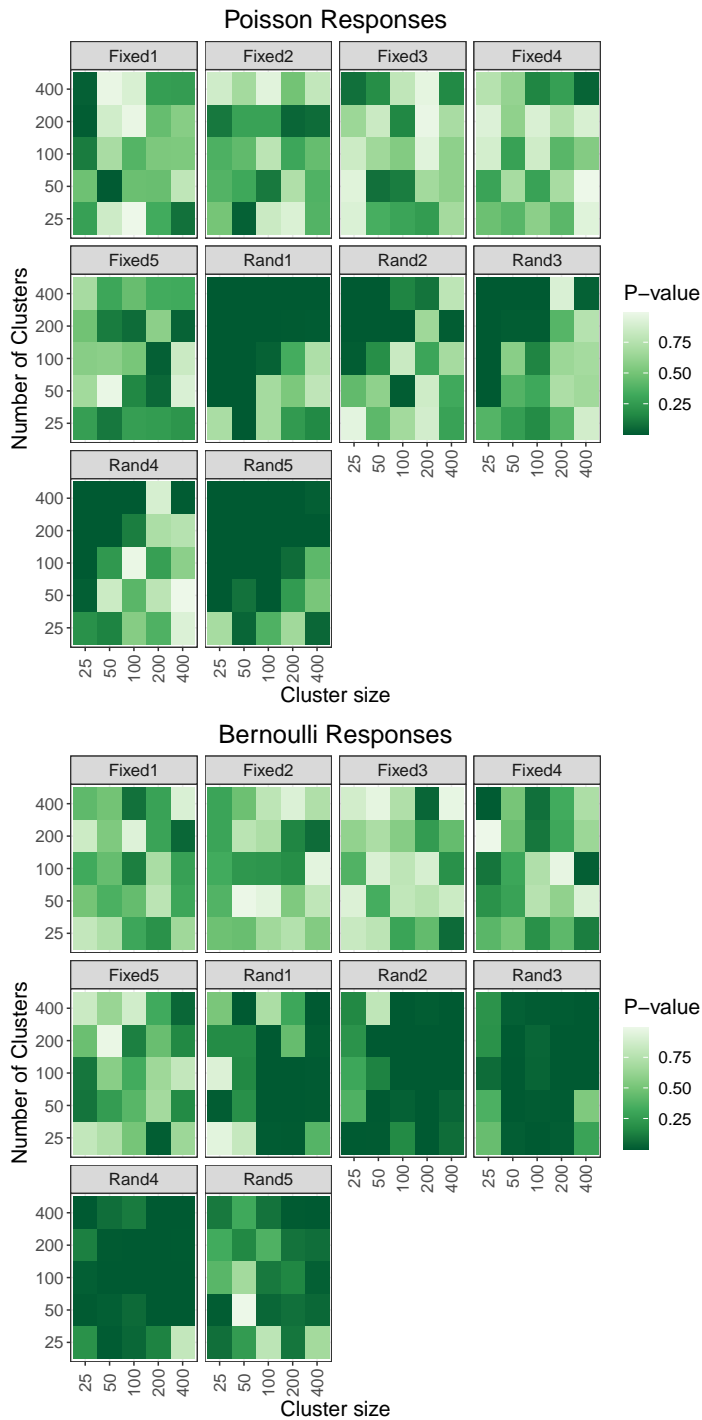


Figure 15:  $p$ -values from Shapiro-Wilk tests applied to the fixed and random effects estimates obtained using maximum PQL estimation, under the unconditional regime.

S5. ADDITIONAL SIMULATION RESULTS

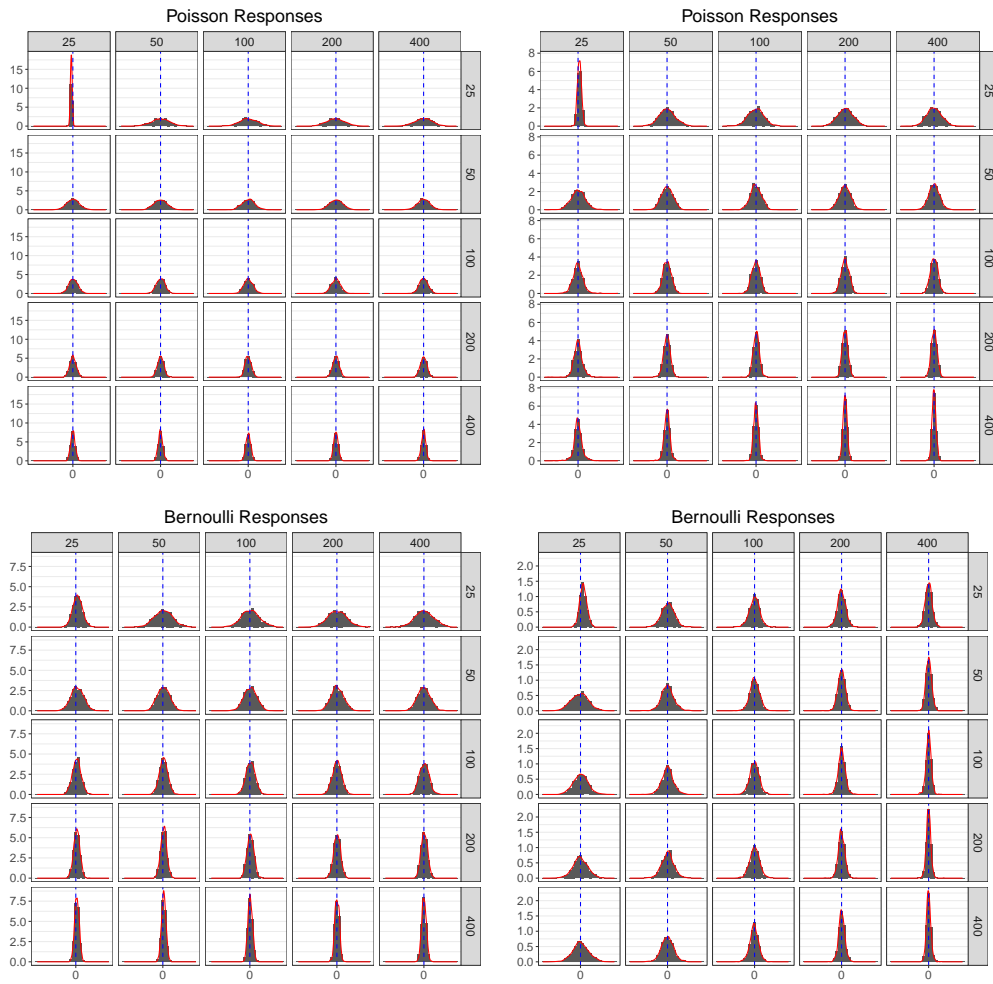


Figure 16: Histograms for the third components of  $\hat{\beta} - \beta$  (left panels) and  $\hat{b}_1 - b_1$  (right panels), under the unconditional regime. Vertical facets represent the cluster sizes, while horizontal facets represent the number of clusters. The dotted blue line indicates zero, and the red curve is a kernel density smoother.

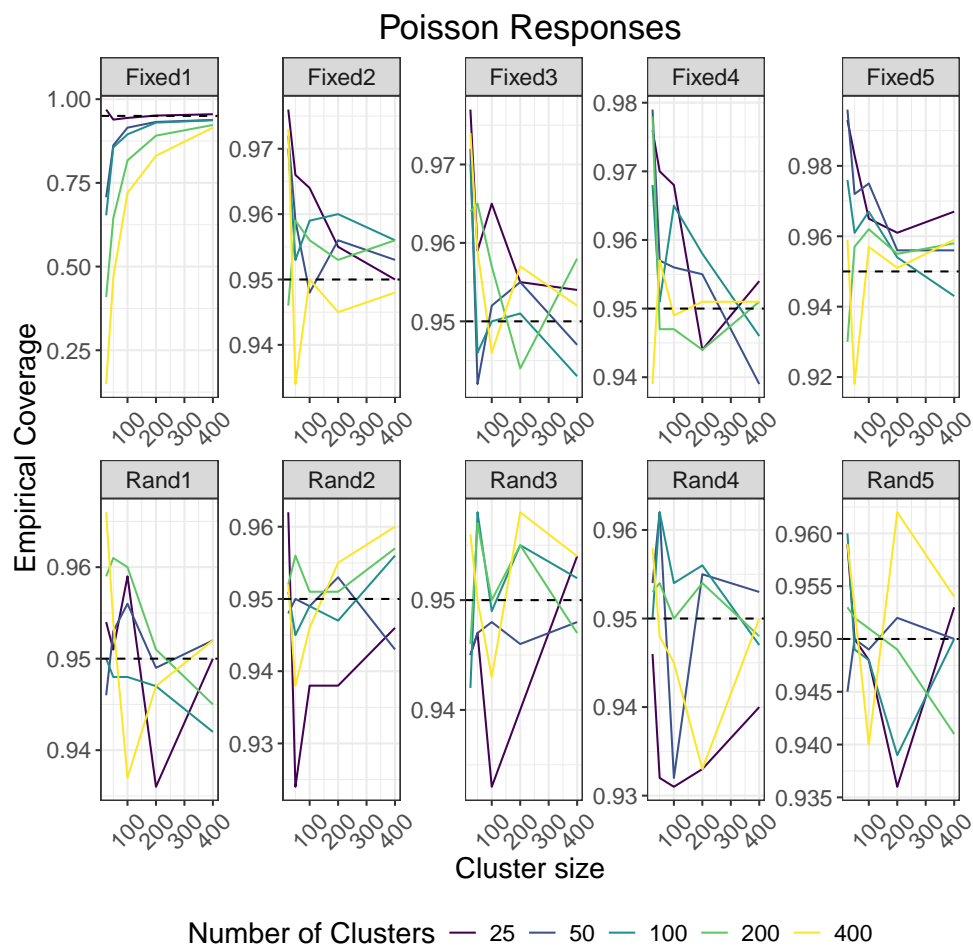


Figure 17: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the conditional regime with Poisson responses.

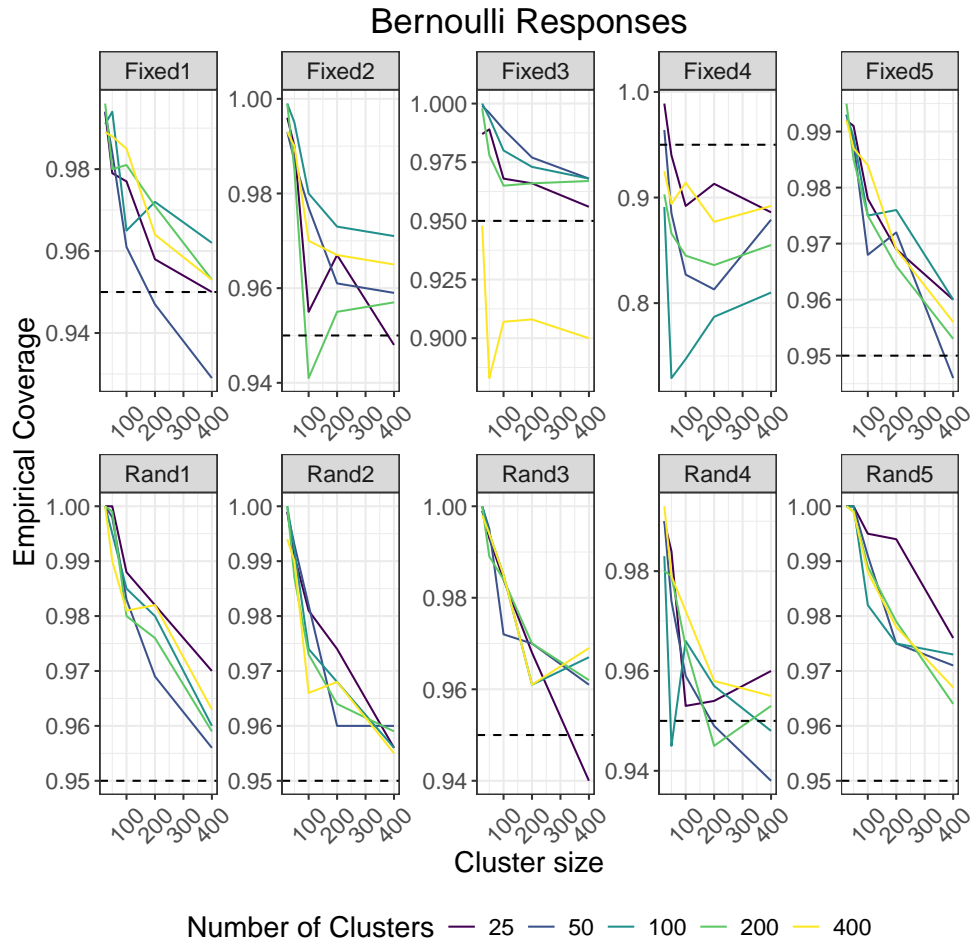


Figure 18: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the conditional regime with Bernoulli responses.

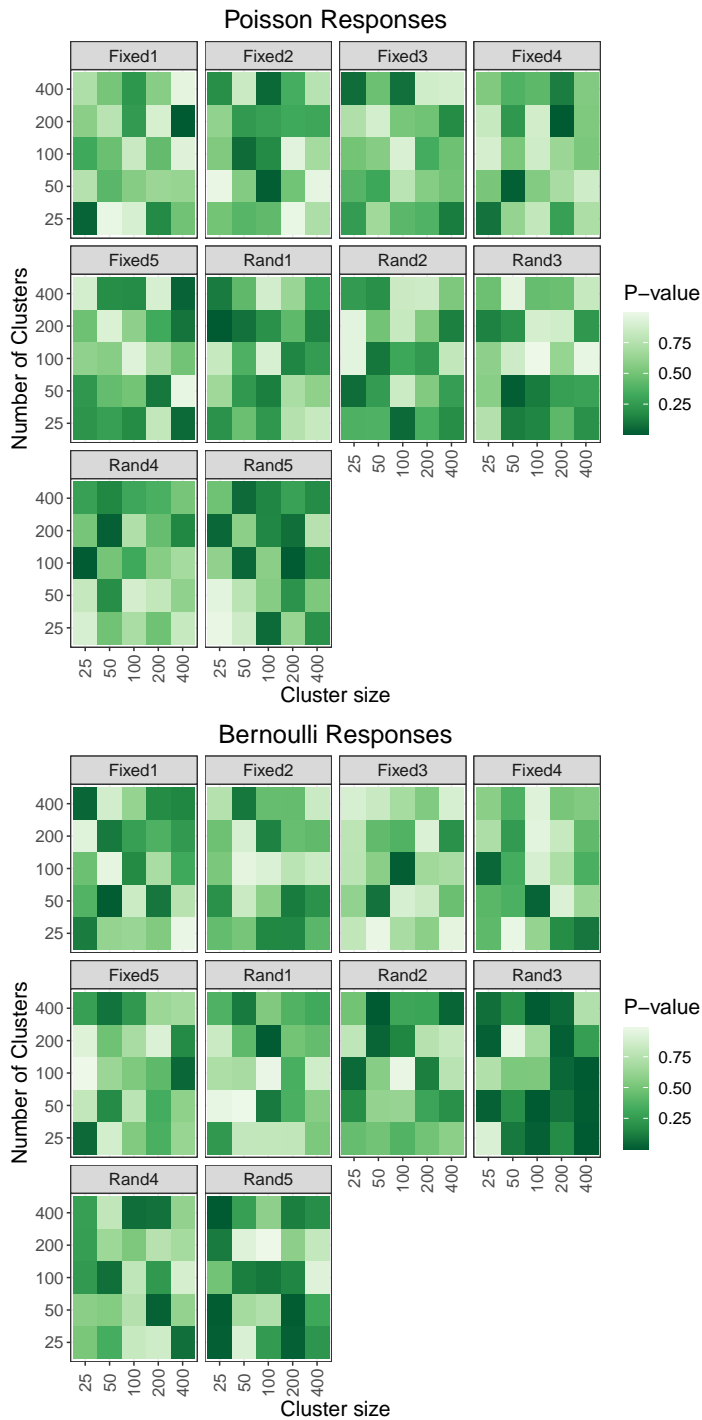


Figure 19:  $p$ -values from Shapiro-Wilk tests applied to the fixed and random effects estimates obtained using maximum PQL estimation, under the conditional regime.



S5. ADDITIONAL SIMULATION RESULTS

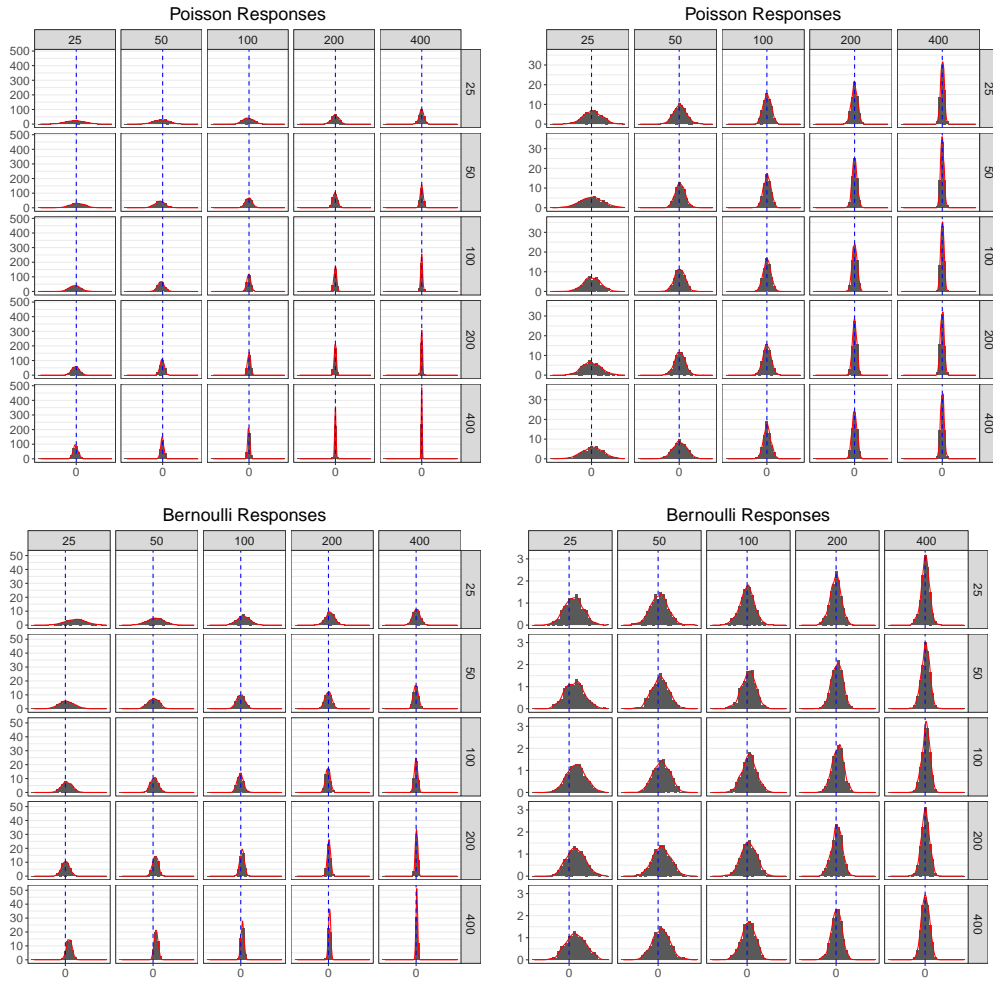


Figure 20: Histograms for the third components of  $\hat{\beta} - \beta$  (left panels) and  $\hat{b}_1 - b_1$  (right panels), under the unconditional regime. Vertical facets represent the cluster sizes, while horizontal facets represent the number of clusters. The dotted blue line indicates zero, and the red curve is a kernel density smoother.

S5.5  $\hat{G} = I_2$

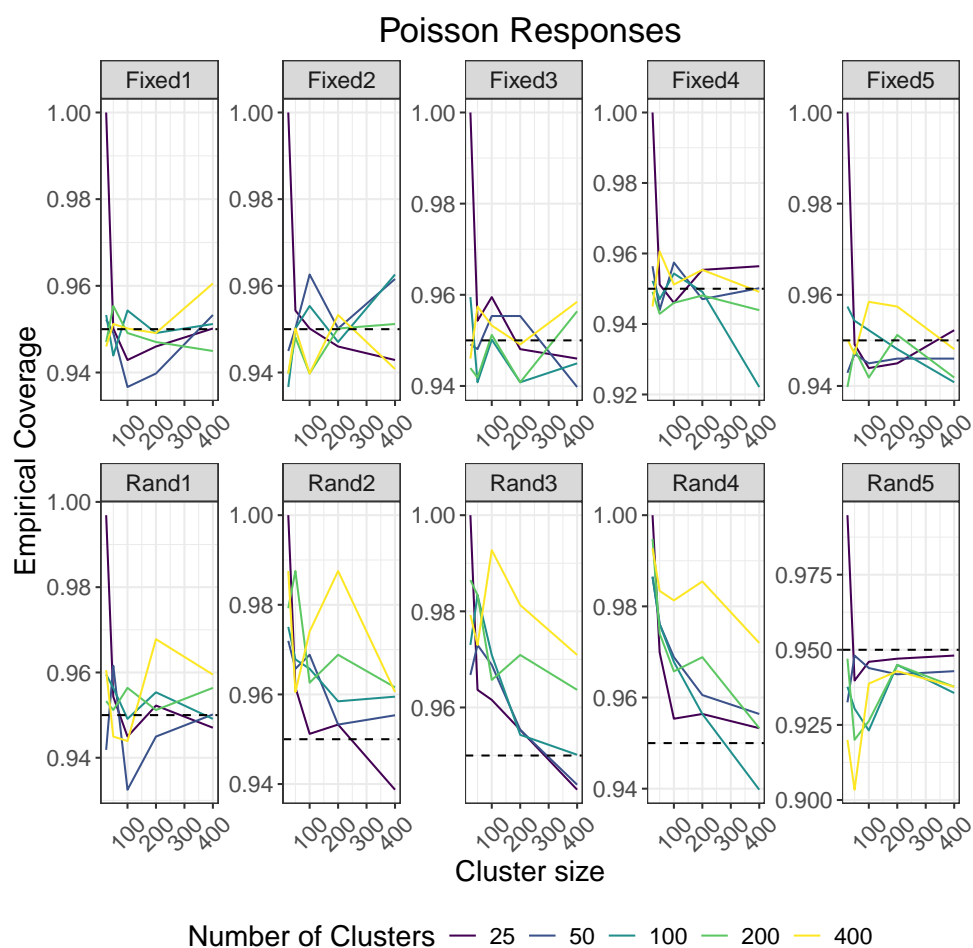


Figure 21: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the unconditional regime with Poisson responses.

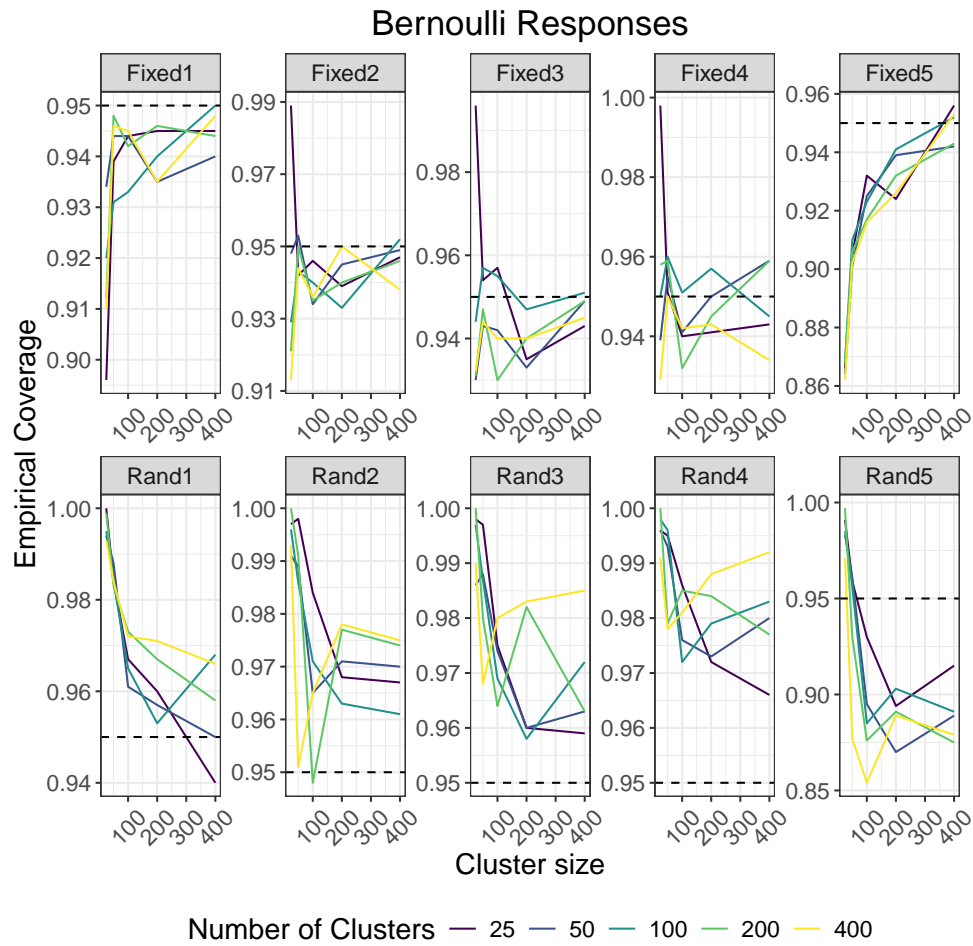


Figure 22: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the unconditional regime with Bernoulli responses.

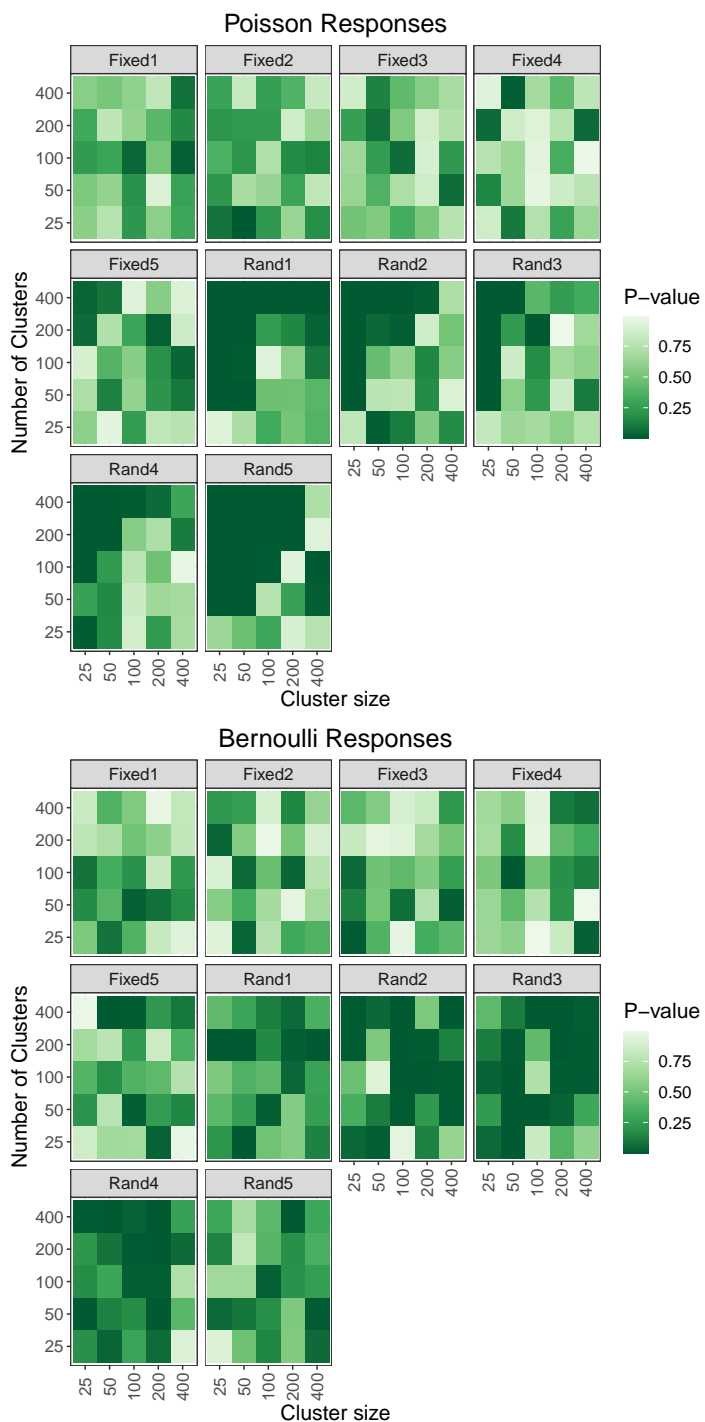


Figure 23:  $p$ -values from Shapiro-Wilk tests applied to the fixed and random effects estimates obtained using maximum PQL estimation, under the unconditional regime.

S5. ADDITIONAL SIMULATION RESULTS

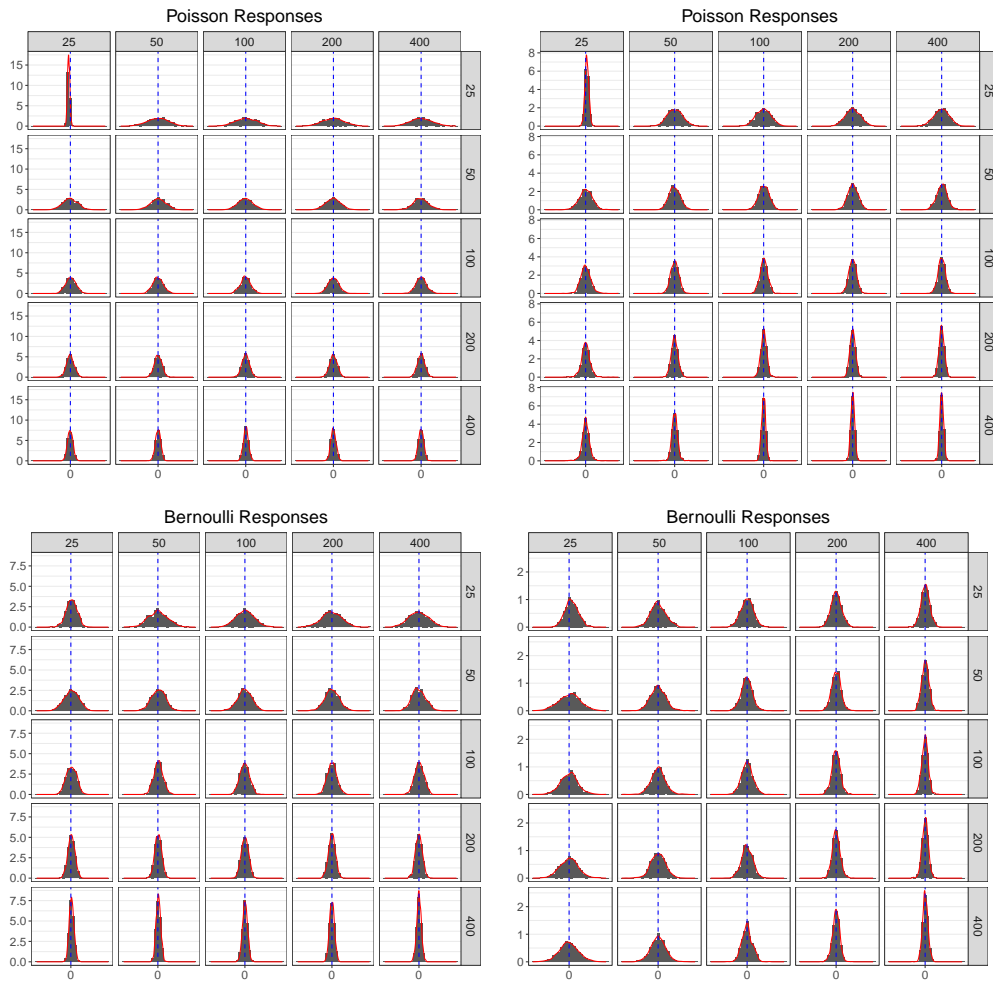


Figure 24: Histograms for the third components of  $\hat{\beta} - \beta$  (left panels) and  $\hat{b}_1 - b_1$  (right panels), under the unconditional regime. Vertical facets represent the cluster sizes, while horizontal facets represent the number of clusters. The dotted blue line indicates zero, and the red curve is a kernel density smoother.

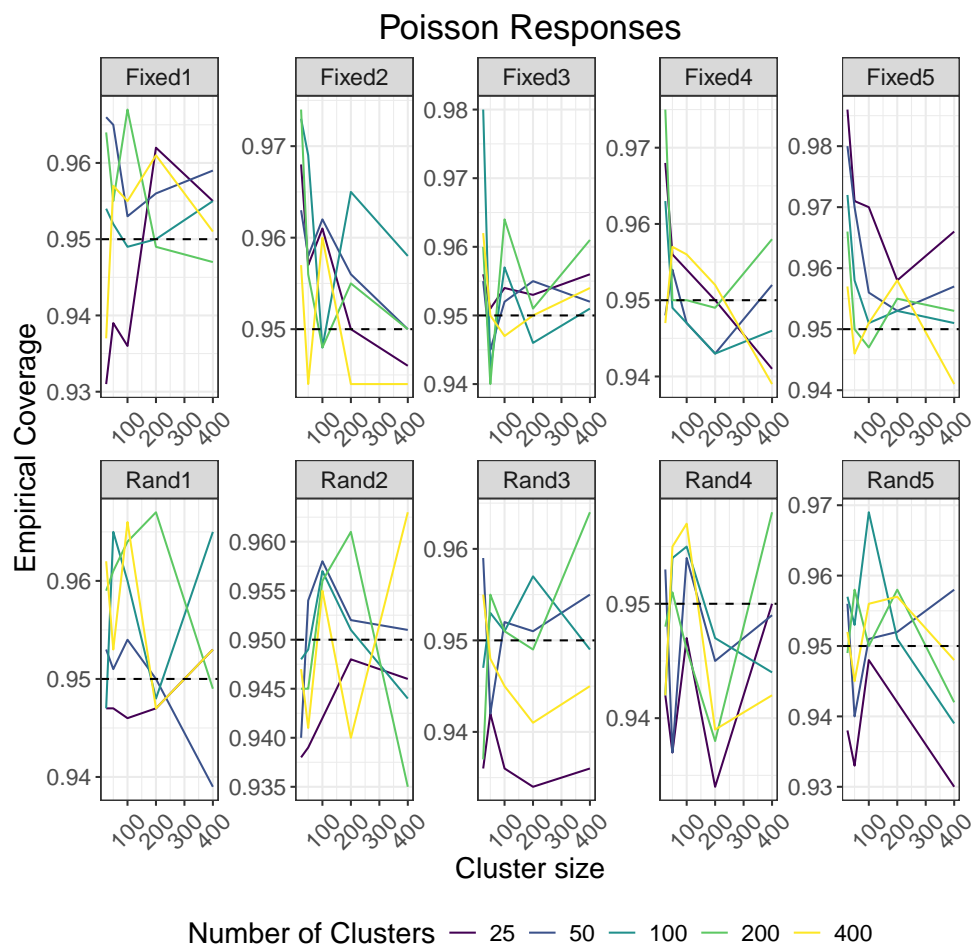


Figure 25: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the conditional regime with Poisson responses.

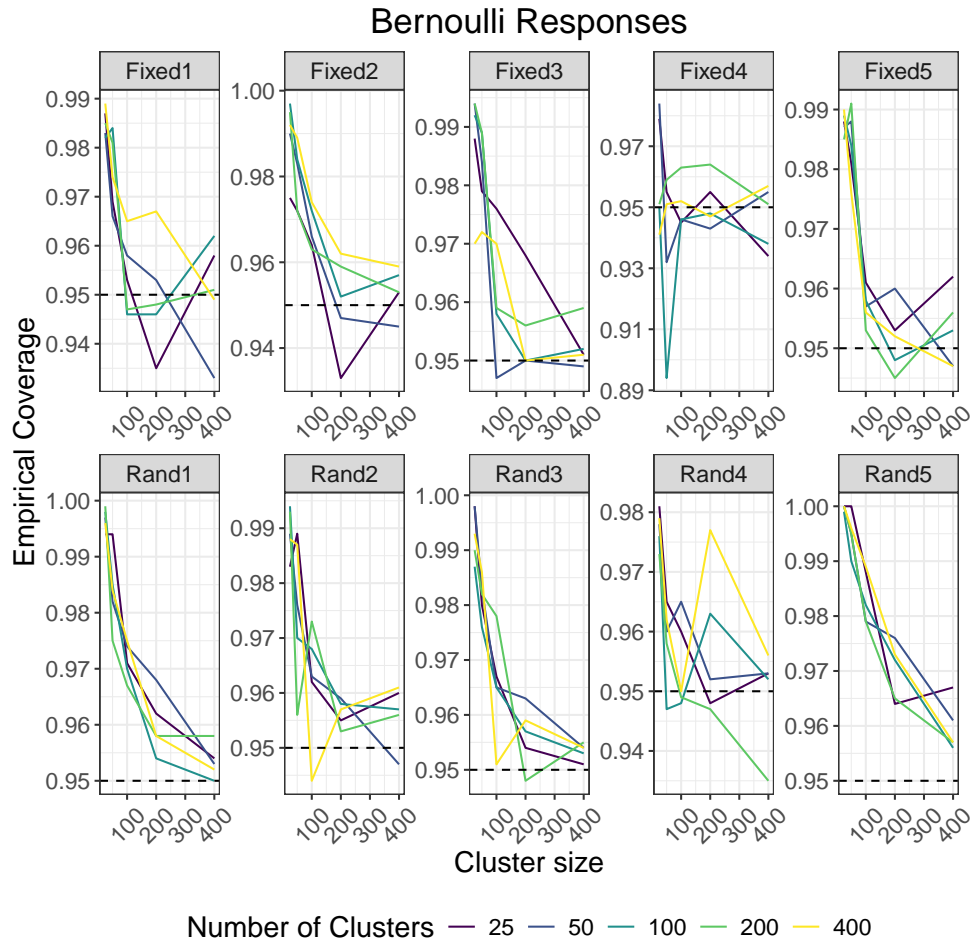


Figure 26: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the conditional regime with Bernoulli responses.

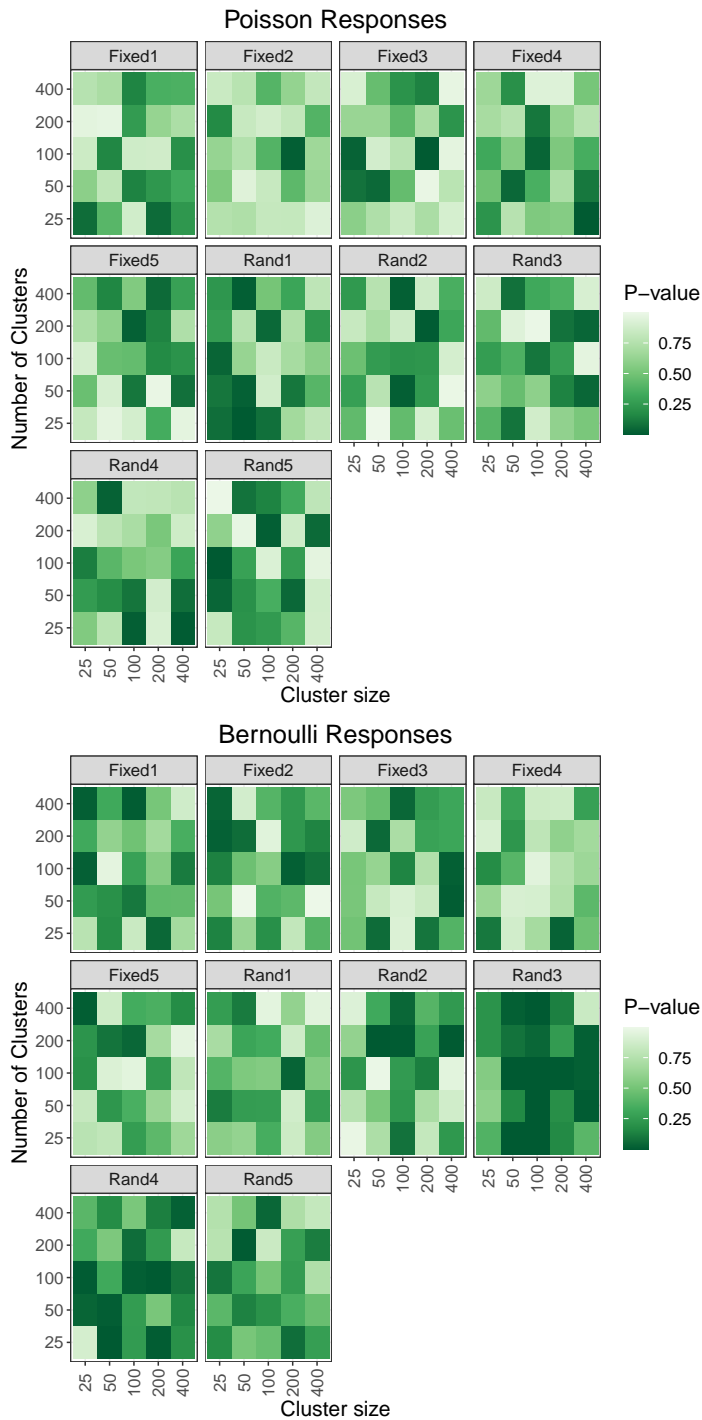


Figure 27:  $p$ -values from Shapiro-Wilk tests applied to the fixed and random effects estimates obtained using maximum PQL estimation, under the conditional regime.



S5. ADDITIONAL SIMULATION RESULTS

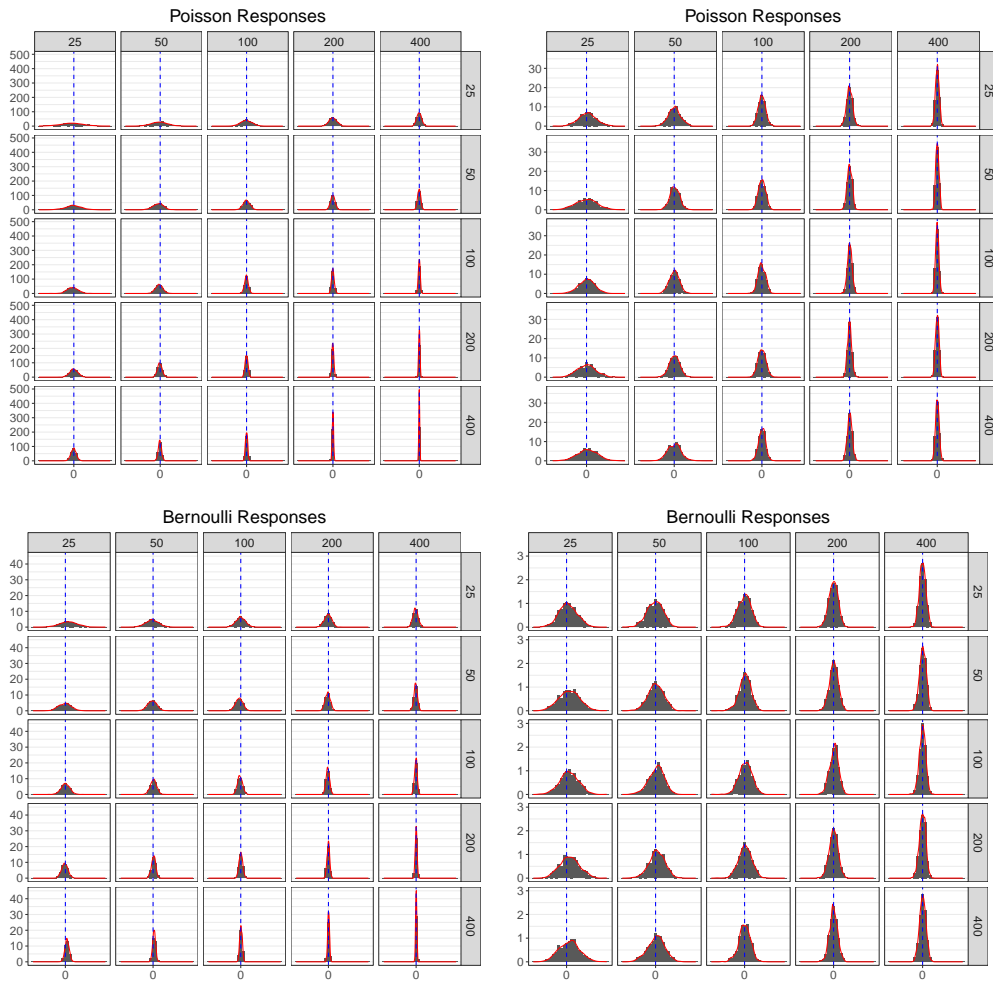


Figure 28: Histograms for the third components of  $\hat{\beta} - \beta$  (left panels) and  $\hat{b}_1 - b_1$  (right panels), under the unconditional regime. Vertical facets represent the cluster sizes, while horizontal facets represent the number of clusters. The dotted blue line indicates zero, and the red curve is a kernel density smoother.

S5.6  $\hat{G} = 2 I_2$

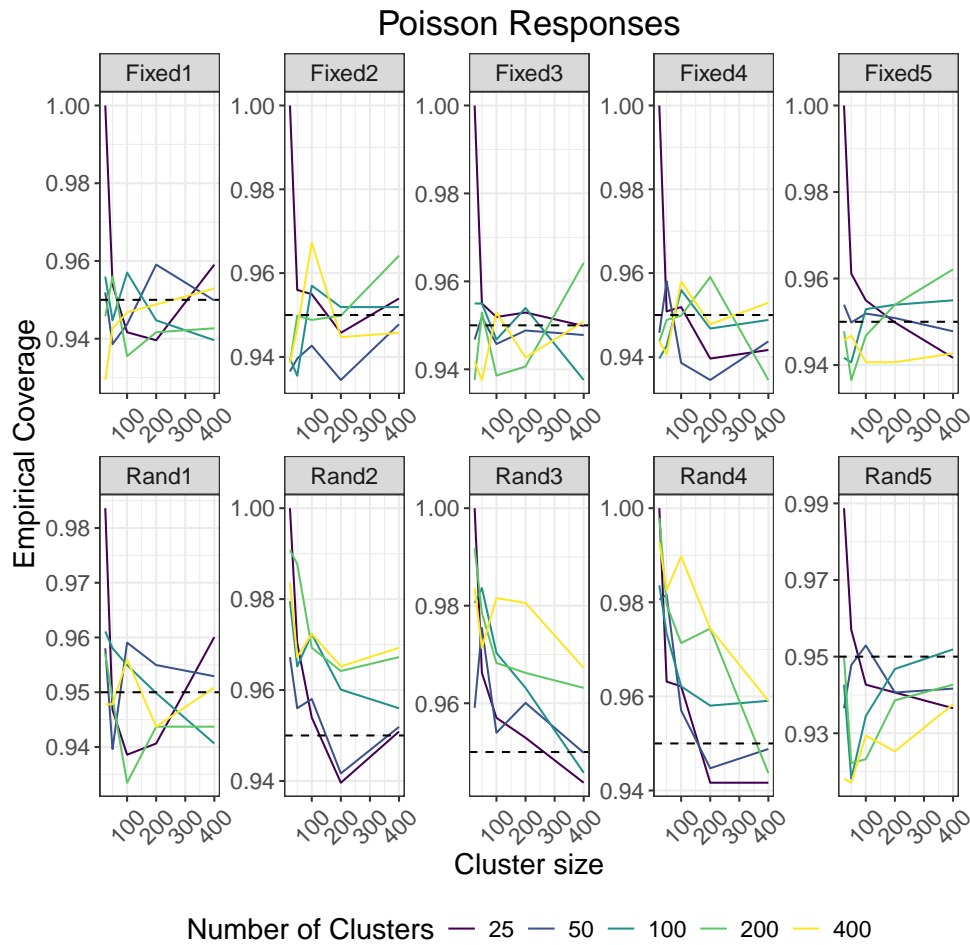


Figure 29: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the unconditional regime with Poisson responses.

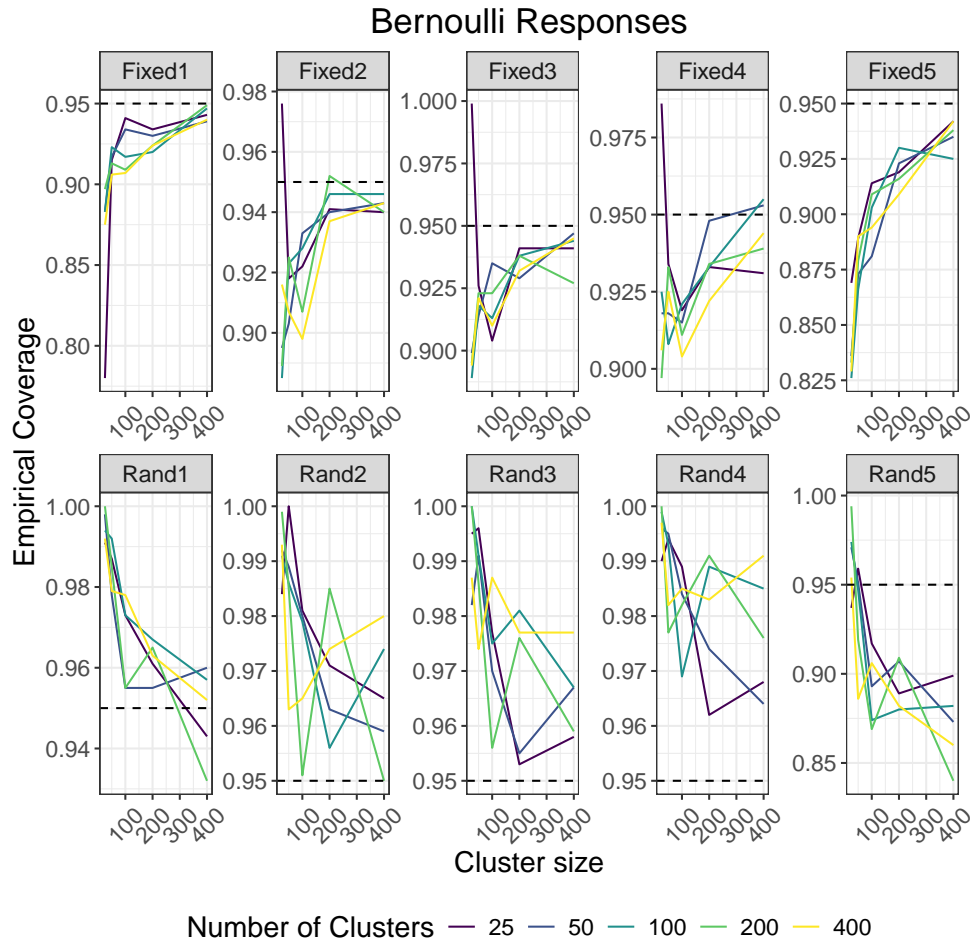


Figure 30: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the unconditional regime with Bernoulli responses.

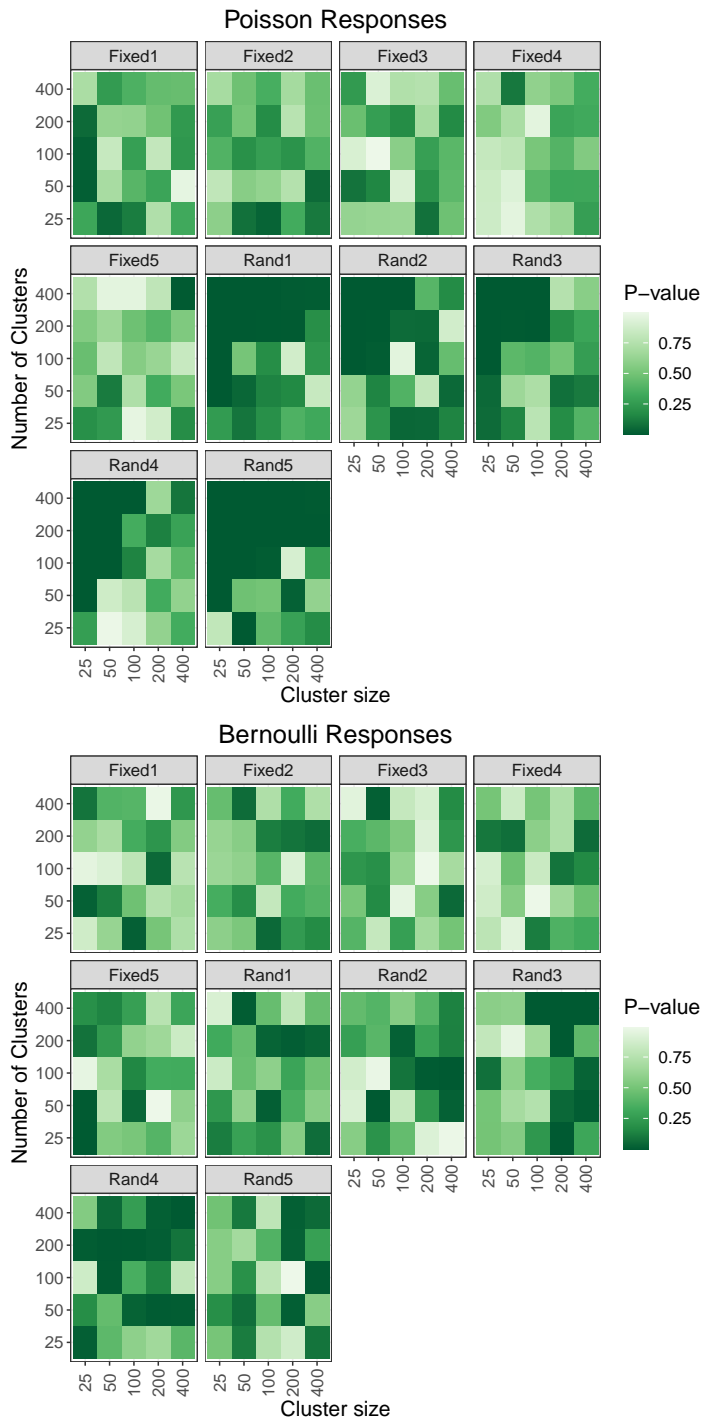


Figure 31:  $p$ -values from Shapiro-Wilk tests applied to the fixed and random effects estimates obtained using maximum PQL estimation, under the unconditional regime.

S5. ADDITIONAL SIMULATION RESULTS

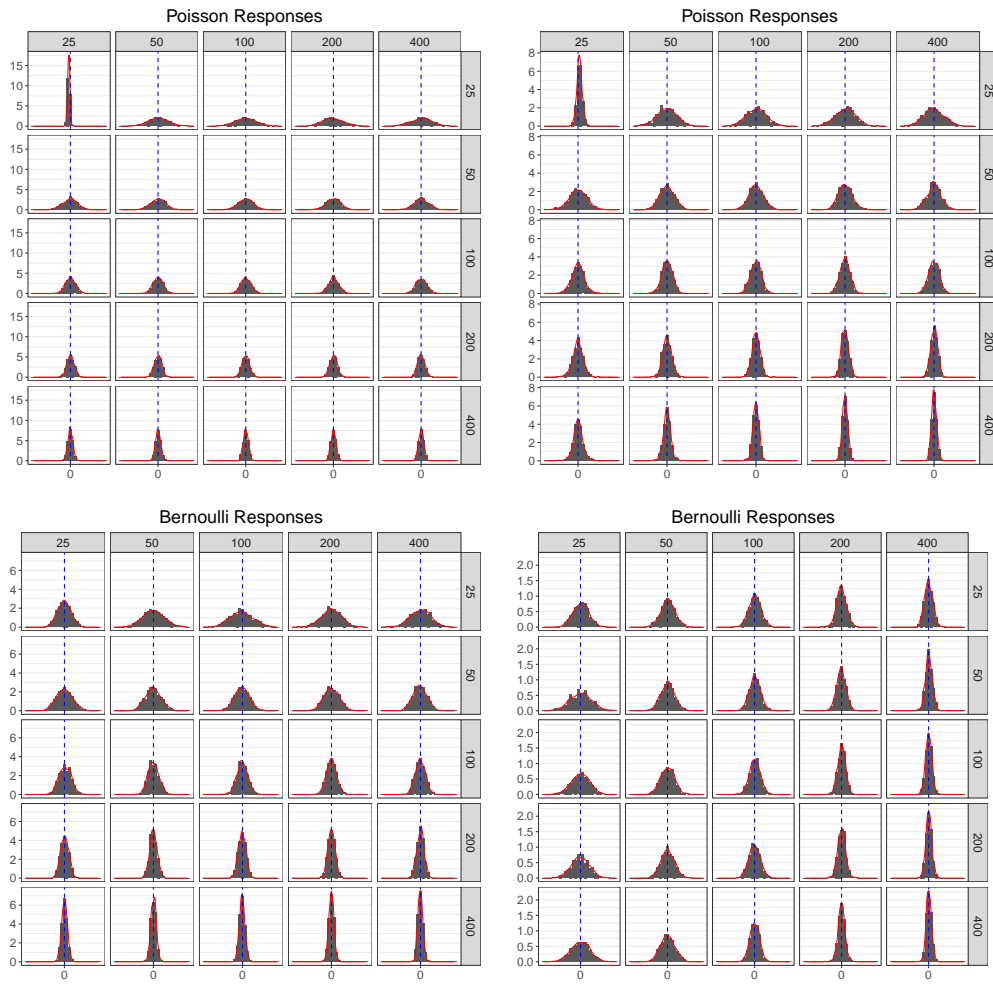


Figure 32: Histograms for the third components of  $\hat{\beta} - \beta$  (left panels) and  $\hat{b}_1 - b_1$  (right panels), under the unconditional regime. Vertical facets represent the cluster sizes, while horizontal facets represent the number of clusters. The dotted blue line indicates zero, and the red curve is a kernel density smoother.

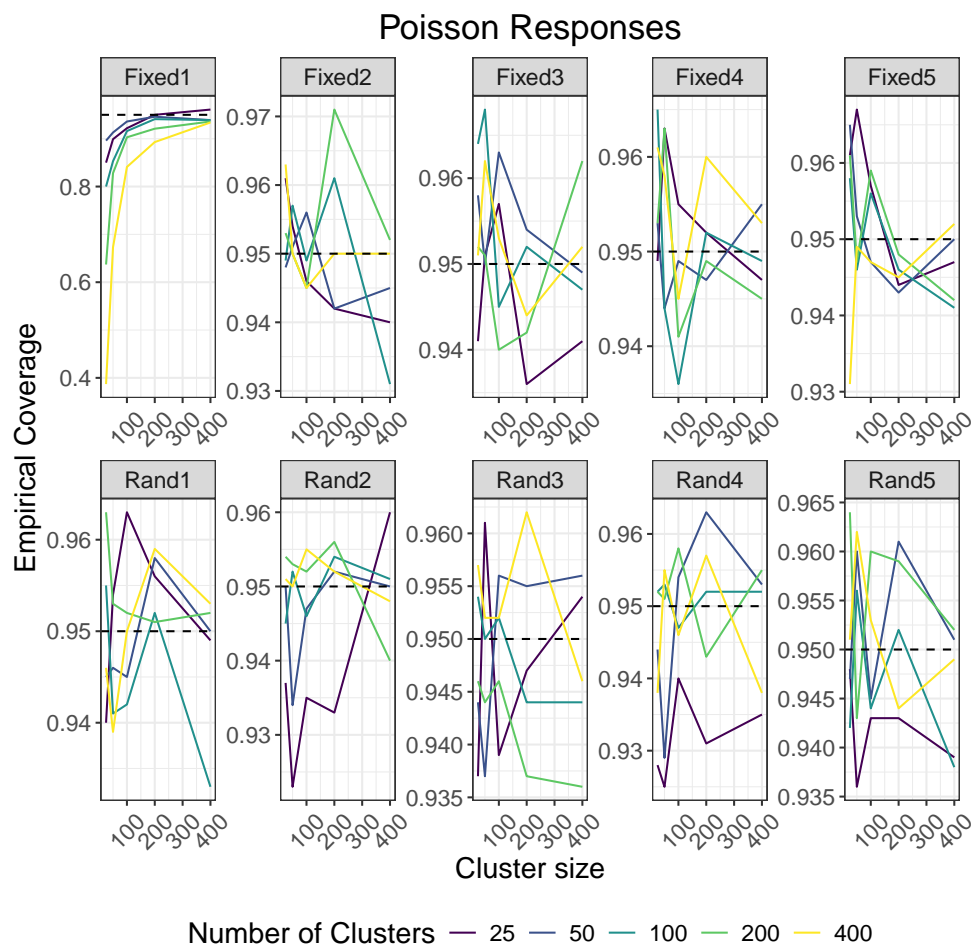


Figure 33: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the conditional regime with Poisson responses.

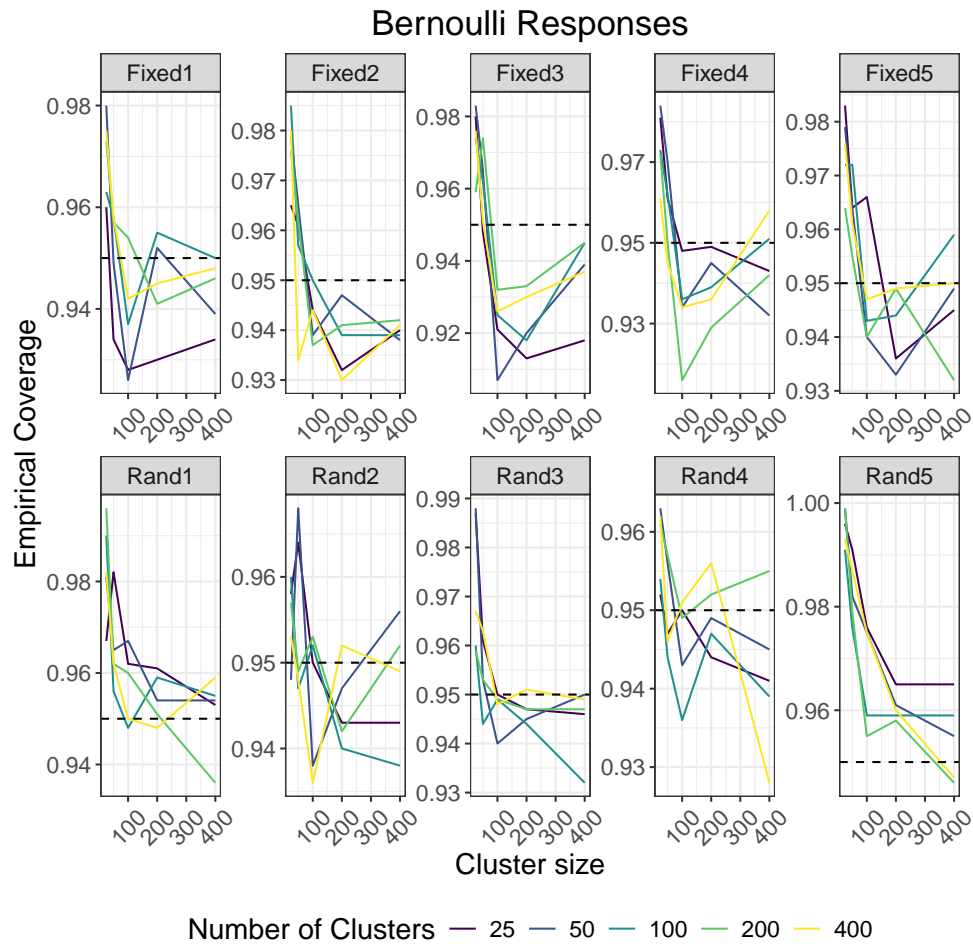


Figure 34: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the conditional regime with Bernoulli responses.

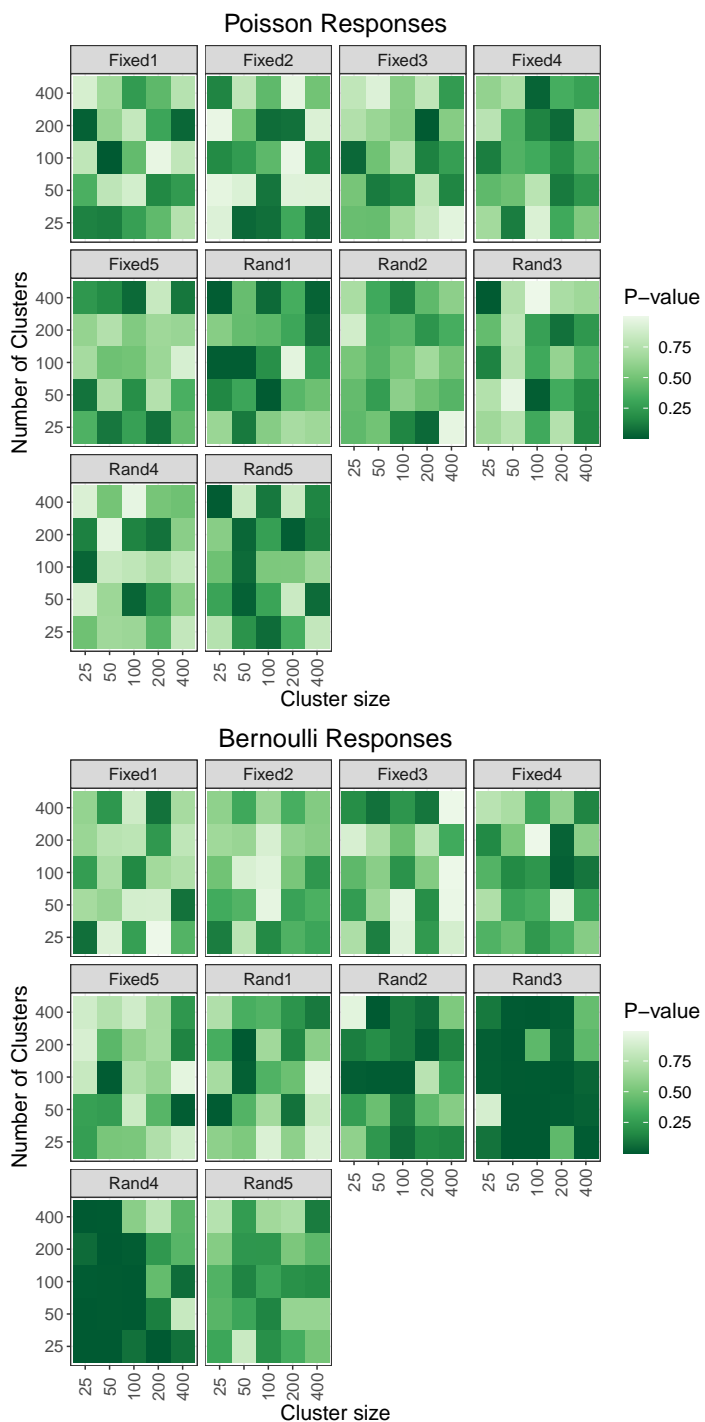


Figure 35:  $p$ -values from Shapiro-Wilk tests applied to the fixed and random effects estimates obtained using maximum PQL estimation, under the conditional regime.



S5. ADDITIONAL SIMULATION RESULTS

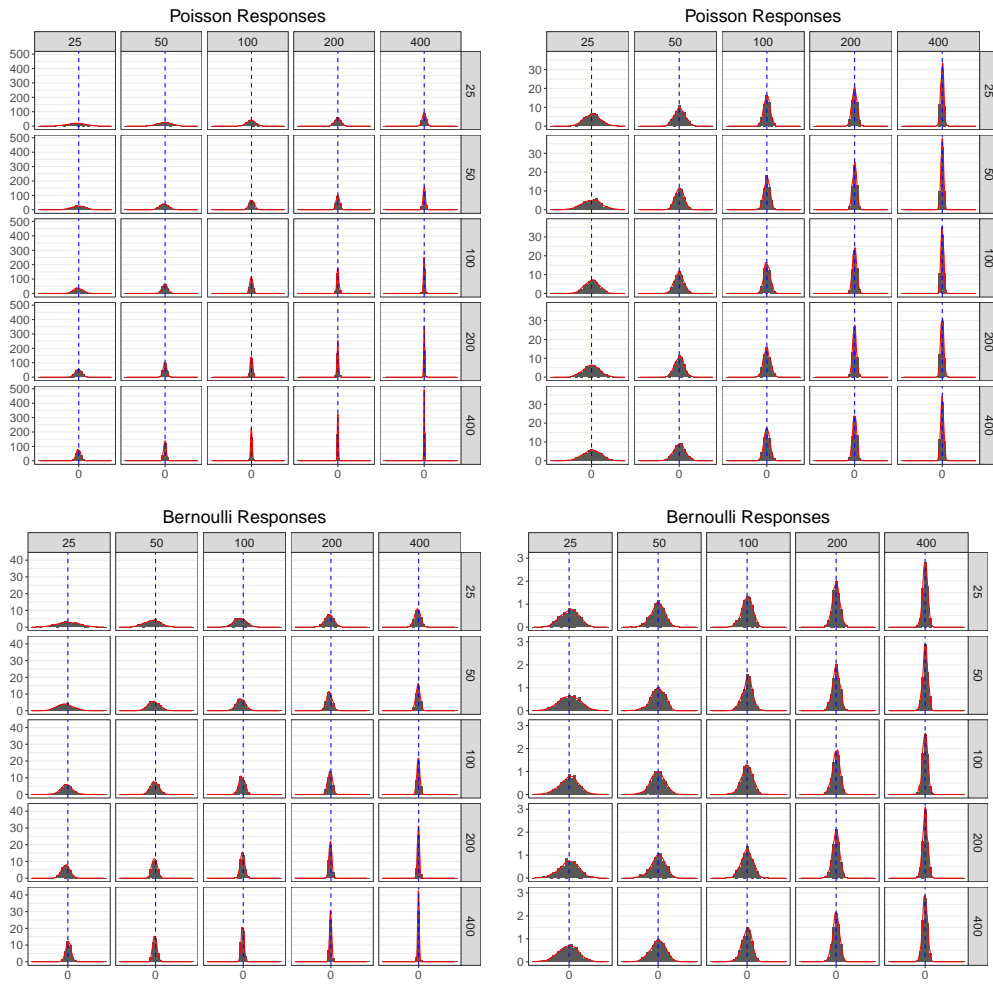


Figure 36: Histograms for the third components of  $\hat{\beta} - \beta$  (left panels) and  $\hat{b}_1 - b_1$  (right panels), under the unconditional regime. Vertical facets represent the cluster sizes, while horizontal facets represent the number of clusters. The dotted blue line indicates zero, and the red curve is a kernel density smoother.

S5.7  $\hat{G} = 4 I_2$

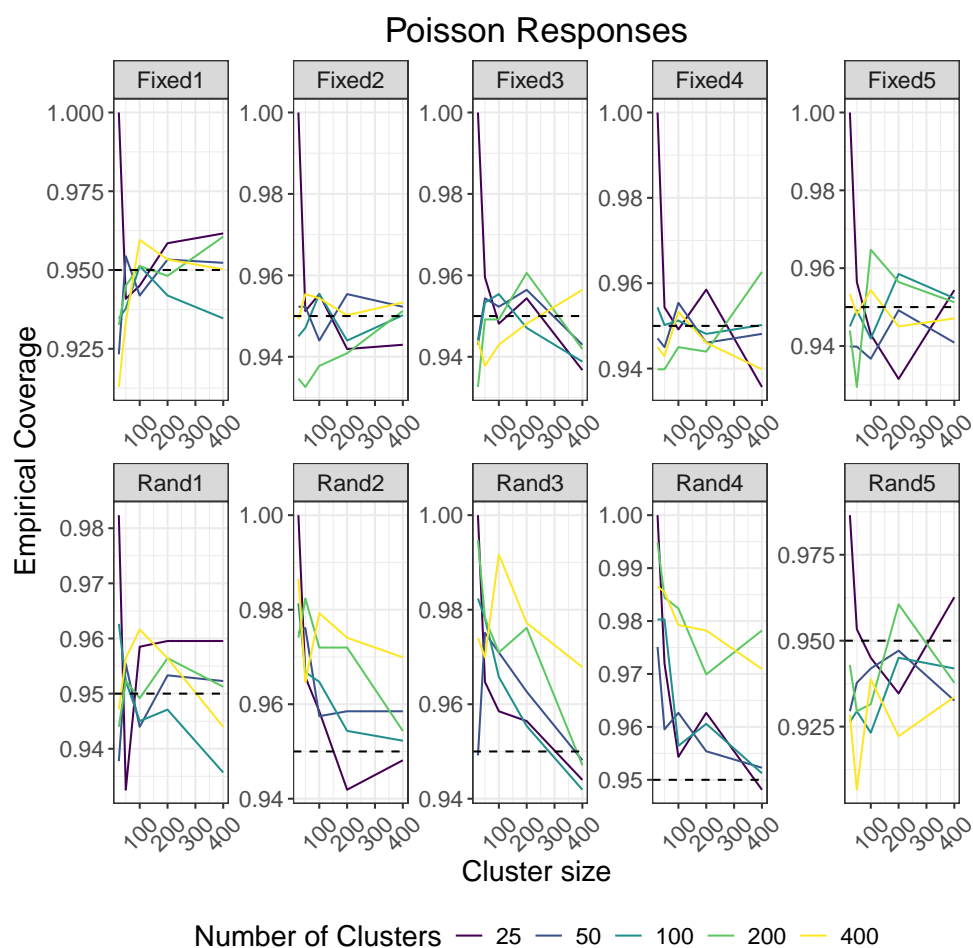


Figure 37: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the unconditional regime with Poisson responses.

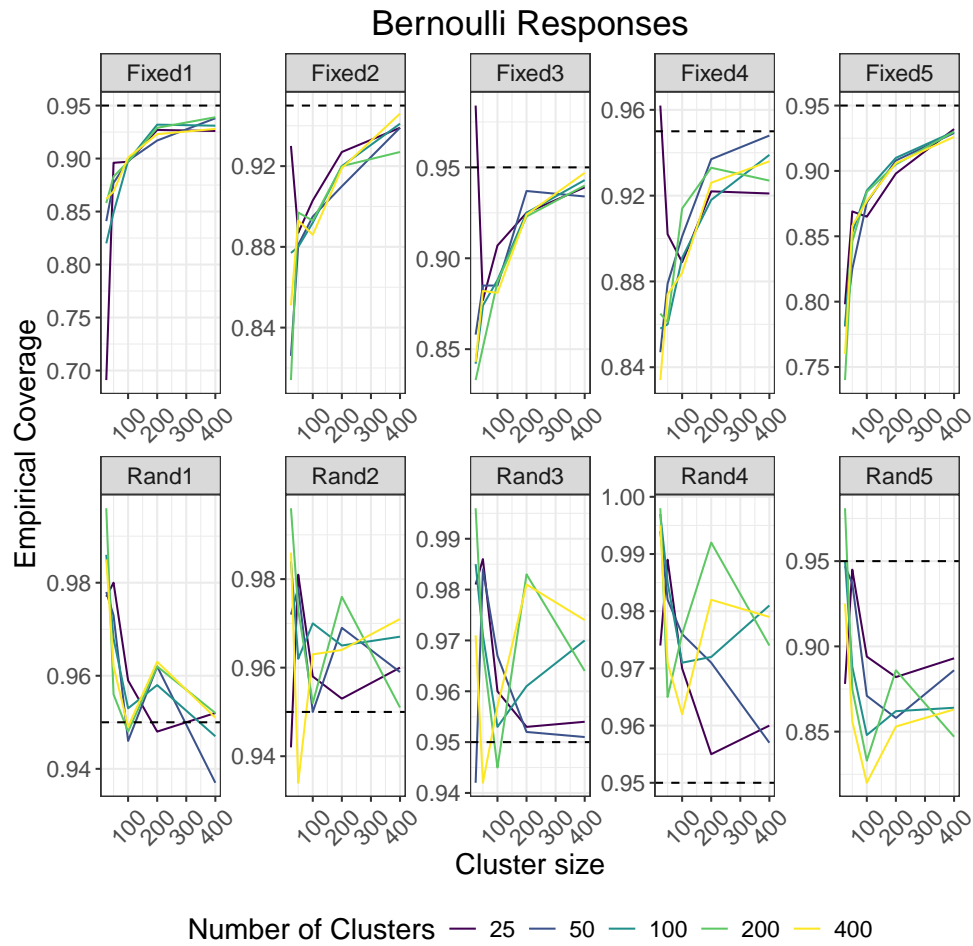


Figure 38: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the unconditional regime with Bernoulli responses.

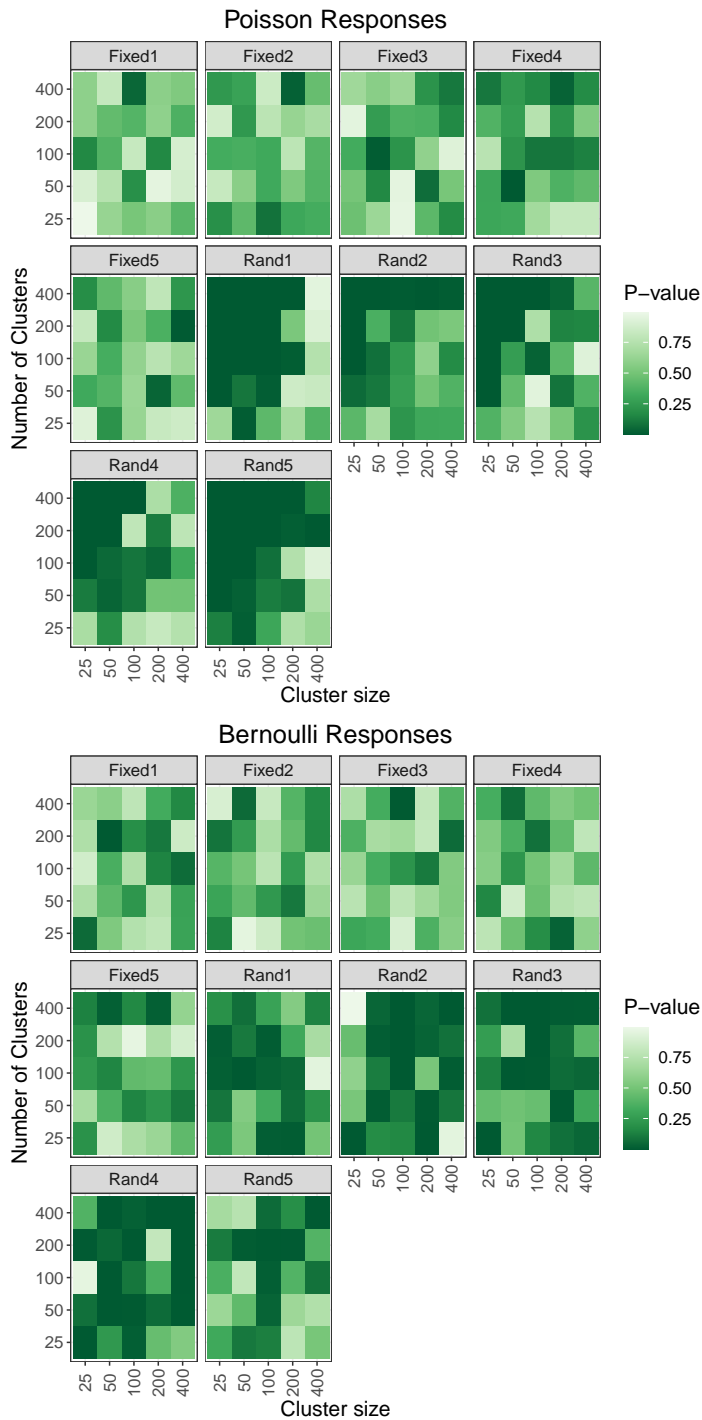


Figure 39:  $p$ -values from Shapiro-Wilk tests applied to the fixed and random effects estimates obtained using maximum PQL estimation, under the unconditional regime.

S5. ADDITIONAL SIMULATION RESULTS

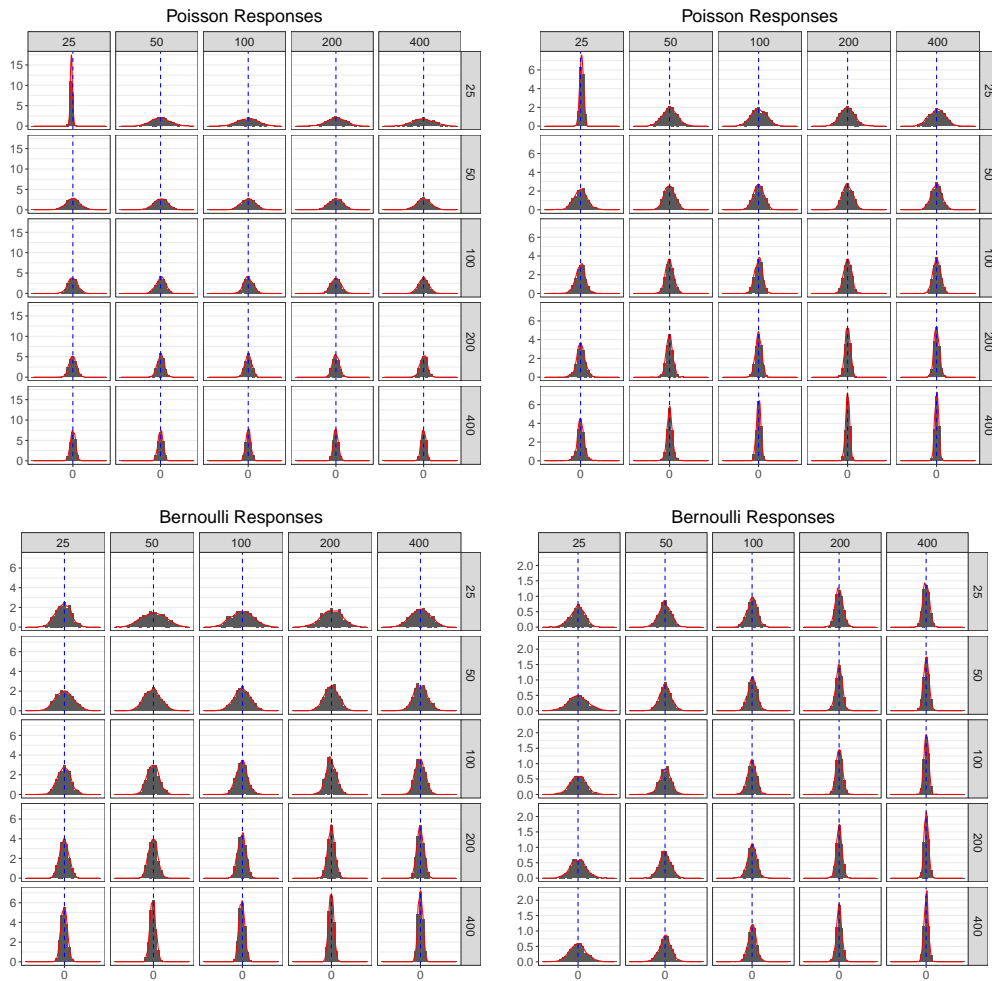


Figure 40: Histograms for the third components of  $\hat{\beta} - \beta$  (left panels) and  $\hat{b}_1 - b_1$  (right panels), under the unconditional regime. Vertical facets represent the cluster sizes, while horizontal facets represent the number of clusters. The dotted blue line indicates zero, and the red curve is a kernel density smoother.

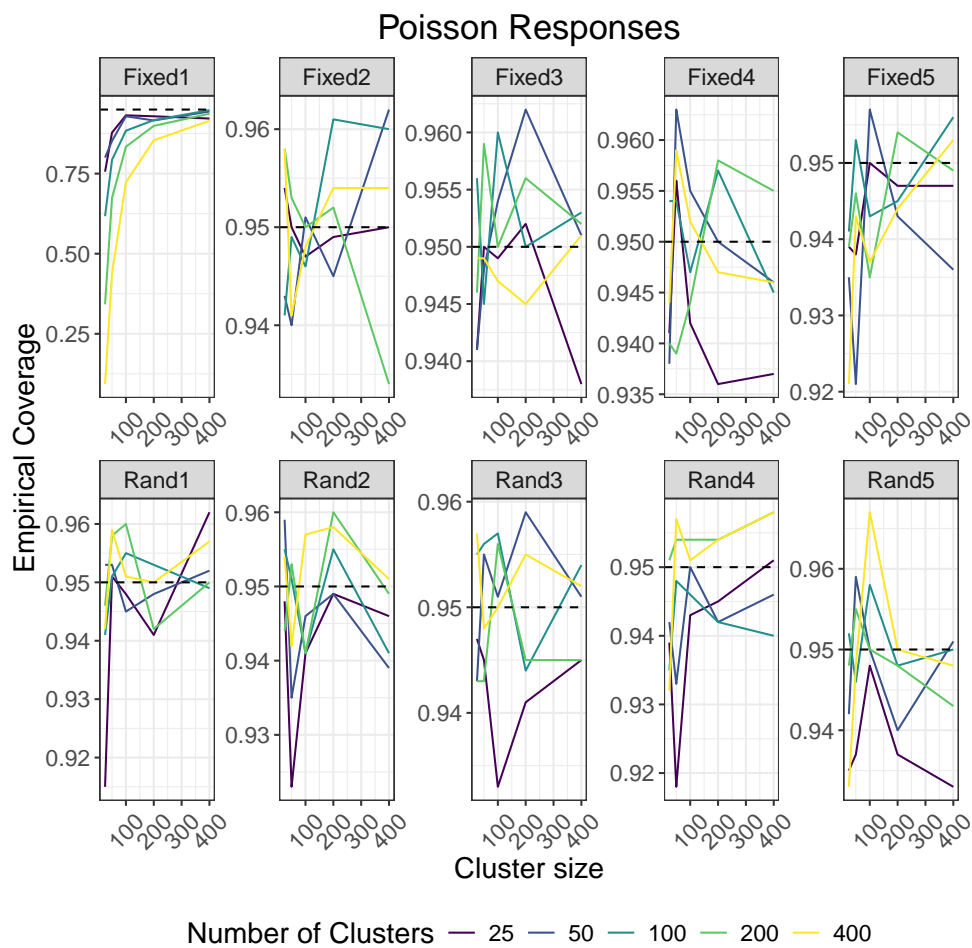


Figure 41: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the conditional regime with Poisson responses.

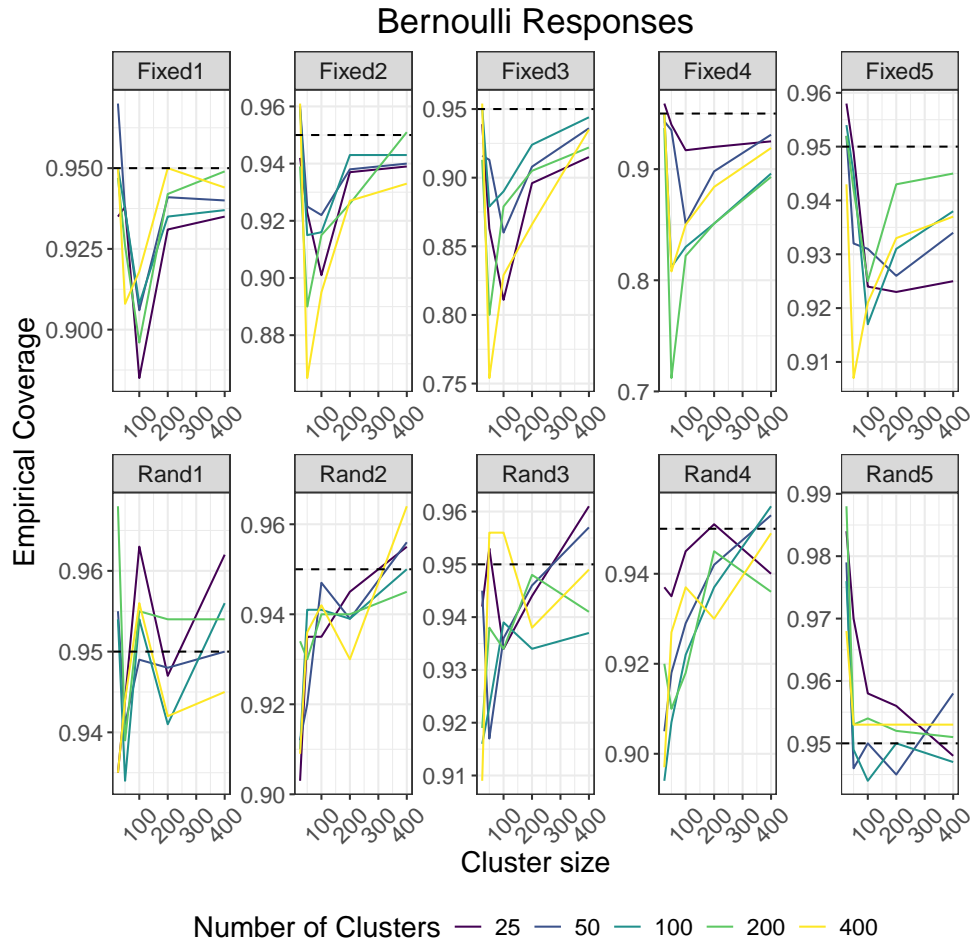


Figure 42: Empirical coverage probability of 95% coverage intervals for the five fixed and random effects estimates, obtained under the conditional regime with Bernoulli responses.

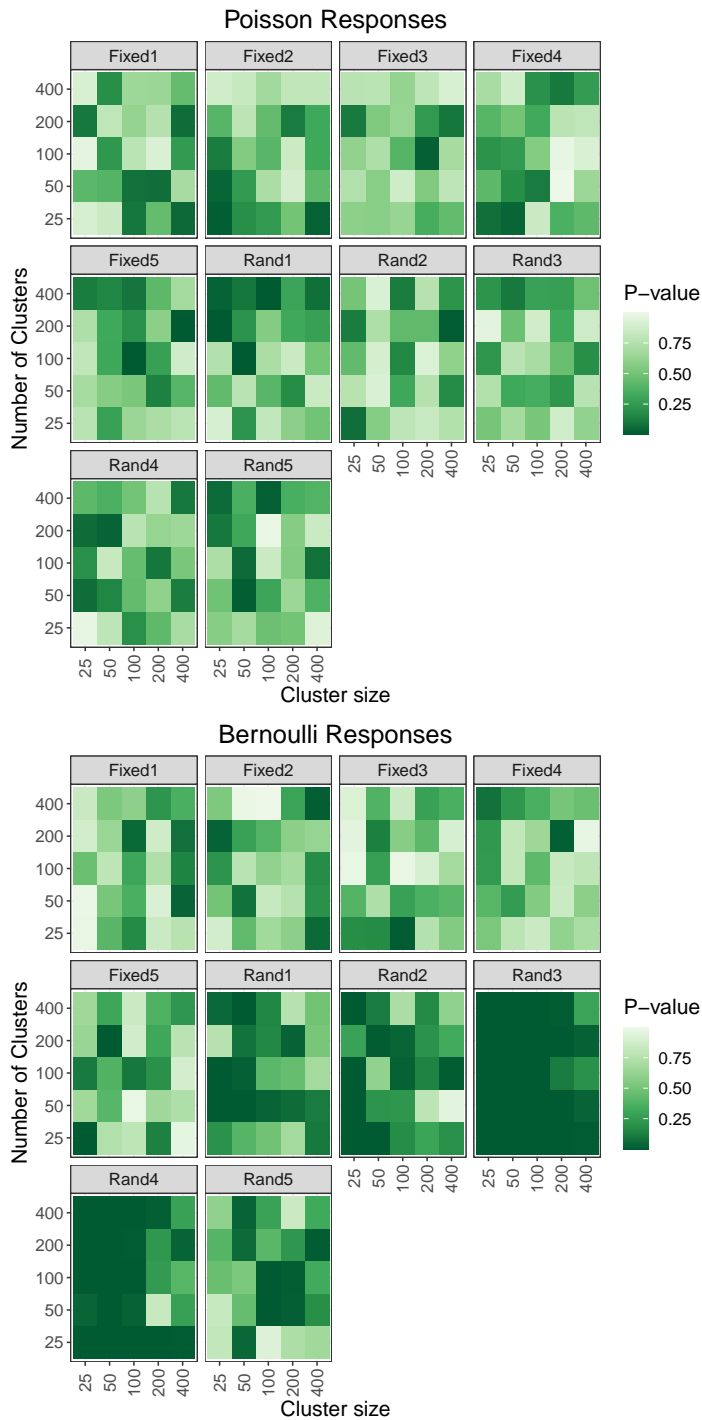


Figure 43:  $p$ -values from Shapiro-Wilk tests applied to the fixed and random effects estimates obtained using maximum PQL estimation, under the conditional regime.



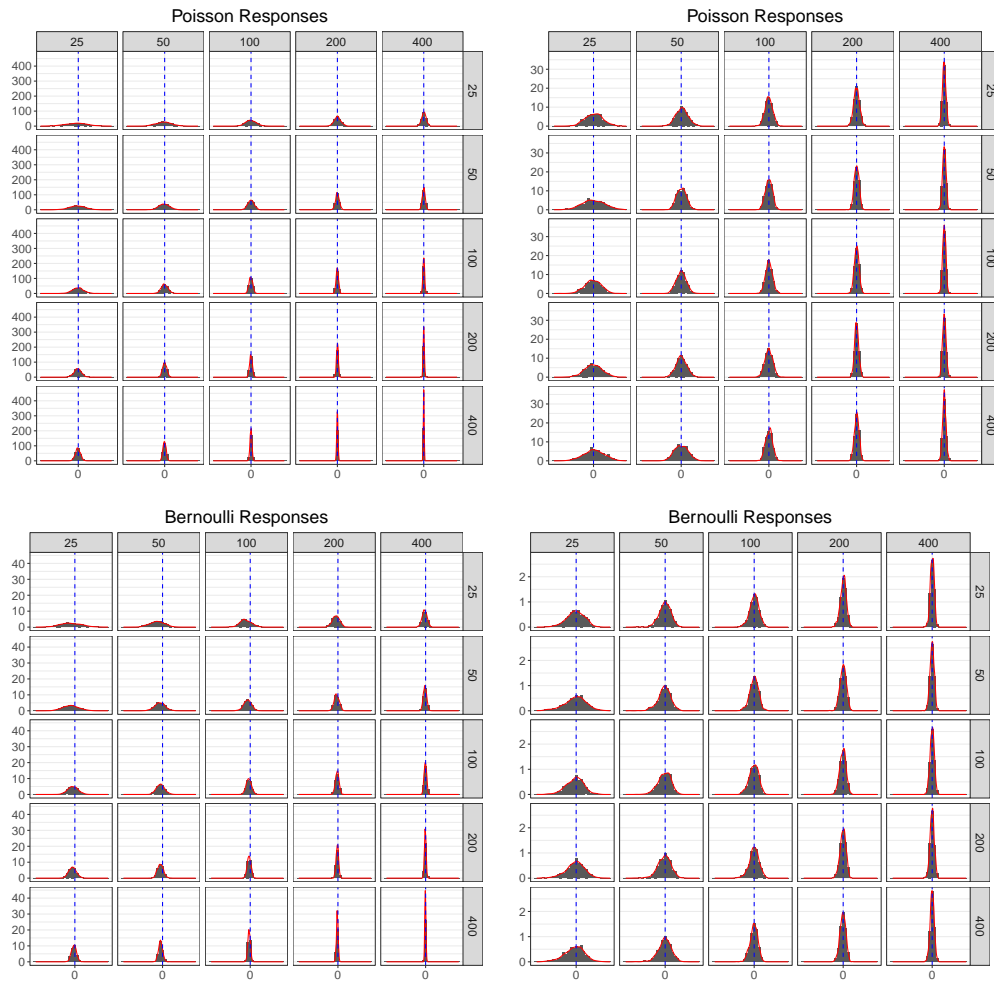


Figure 44: Histograms for the third components of  $\hat{\beta} - \beta$  (left panels) and  $\hat{b}_1 - b_1$  (right panels), under the unconditional regime. Vertical facets represent the cluster sizes, while horizontal facets represent the number of clusters. The dotted blue line indicates zero, and the red curve is a kernel density smoother.

## Bibliography

Blum, J., H. Chernoff, M. Rosenblatt, and H. Teicher (1958). Central limit theorems for interchangeable processes. *Canadian Journal of Mathemat-*

*ics* 10, 222–229.

Downey, P. J. (1990). Distribution-free bounds on the expectation of the maximum with scheduling applications. *Operations Research Letters* 9, 189–201.

Lyu, Z. and A. H. Welsh (2021a). Asymptotics for EBLUPs: Nested error regression models. *Journal of the American Statistical Association* 117, 1–15.

Lyu, Z. and A. H. Welsh (2021b). Increasing cluster size asymptotics for nested error regression models. *Journal of Statistical Planning and Inference* 217, 52–68.