

## **Rgbp: An R Package for Conjugate Gaussian, Poisson, and Binomial Hierarchical Modeling and Frequency Method Checking on Overdispersed Data**

**Hyungsuk Tak**  
Harvard University

**Joseph Kelly**  
Google

**Carl Morris**  
Harvard University

---

### **Abstract**

**Rgbp** is an R package that utilizes approximate Bayesian machinery to fit two-level conjugate hierarchical models on overdispersed Gaussian, Poisson, and Binomial data. The data that **Rgbp** assumes comprise of observed sufficient statistics for each random effect, such as sample means, possibly together with covariates of each group. The approximate Bayesian tool equipped with the adjustment for density maximization estimates random effects and hyper-parameters in the conjugate prior distributions. As for the Binomial model, the package has an option to draw their posterior samples via the acceptance-rejection method. The main goal of **Rgb** is to produce Bayesian interval estimates for the random effects that meet their nominal confidence levels. For this purpose, we adopt unique improper hyper-prior distributions. **Rgbp** provides a tool to check quickly whether the resulting Bayesian interval estimates for the random effects achieve the nominal confidence levels via a repeated sampling coverage evaluation, which we call “frequency method checking.”

*Keywords:* overdispersion, hierarchical model, adjustment for density maximization, acceptance-rejection method, repeated sampling coverage evaluation, frequency method checking, R.

---

## **1. Introduction (v4; 28 July 2015)**

Gaussian, Poisson, or Binomial data from several independent groups sometimes have more variation than the assumed Gaussian, Poisson, or Binomial distributions of the first-level observed data. To account for this extra-variability, called overdispersion, a two-level conjugate hierarchical model regards first-level mean parameters as random effects that come from a population-level conjugate prior distribution. The main goal of our two-level conjugate modeling is to estimate these random effects for a comparison between groups.

## 2 **Rgbp**: Hierarchical Modeling and Frequency Method Checking on Overdispersed Data

With an assumption of homogeneity within each group, the observed data are group-level aggregate data from  $k$  independent groups, composed of sufficient statistics for their  $k$  random effects (without the population-level data). Specifically, the data for the Gaussian model comprise of each group's sample mean and its standard error, those for the Poisson model use each group's outcome count and an exposure measure, and those for the Binomial model use the number of each group's successful outcomes together with the total number of trials. The data analyzed by **Rgbp** may incorporate each group's covariate information. These types of data are common, e.g., as for a biological analysis on litter data (Tamura and Young 1987), a meta analysis on independent studies (Chapter 5 in Gelman, Carlin, Stern, and Rubin (2014)), or small area estimation problems (Ghosh and Rao 1994; Rao 2003).

For such data, assuming homogeneity within each group, **Rgbp**'s two-level model may be viewed as a conjugate hierarchical generalized linear model (HGLM) (Lee and Nelder 1996; Lee, Nelder, and Pawitan 2006) where each random effect has a conjugate prior distribution. However, the HGLM focuses on estimating regression coefficients to explore associations between covariates and observed data. While **Rgbp** does that too, its emphasis concerns making valid point and interval estimates for the  $k$  random effects.

**Rgbp** combines Bayesian tools with our special improper hyper-prior distributions on the second-level parameters chosen to produce posterior interval estimates to achieve nominal confidence levels for the random effects. This hyper-prior distribution is related to Stein's harmonic prior for two-level Gaussian models where it has been seen to produce good repeated sampling coverage rates for the Bayesian interval estimates for the  $k$  random effects (Morris and Tang 2011; Morris and Lysy 2012; Kelly 2014). This prior, as extended here for **Rgbp**'s Poisson and Binomial hierarchical models, has been seen for every data set tested thus far to meet (or exceed) the pre-assigned confidence level.

For fitting the hierarchical model, **Rgbp** uses ADM, i.e. "adjustment for density maximization" (Morris 1988a; Christiansen and Morris 1997; Morris and Tang 2011). ADM approximates a posterior density or a likelihood function by fitting a selected (one dimensional) Pearson family, based on the first two derivatives of the given density function. When the Normal distribution is the chosen Pearson family, ADM reduces to maximum likelihood estimation. Because shrinkage factors lie in  $[0,1]$ , **Rgbp** approximates shrinkage factor distributions with Beta distributions because these are constrained to  $[0,1]$ , unlike the Normal.

**Rgbp** estimates the first two (actually, three in the Gaussian case) posterior moments of each random effect. Finally, **Rgbp** approximates the posterior distribution of each random effect by a skewed Normal distribution for Gaussian case, by a Gamma distribution for the Poisson case, and by a Beta distribution for the Binomial case.

For the Binomial hierarchical model, **Rgbp** provides an option to draw independent posterior samples of all the model parameters, including random effects, via an acceptance-rejection method (Robert and Casella 2013).

In addition to fitting the hierarchical models, **Rgbp** evaluates the repeated sampling coverage rates of the resulting Bayesian interval estimates for random effects (Christiansen and Morris 1997; Daniels 1999; Tang 2002; Morris and Tang 2011; Morris and Lysy 2012). This procedure distinguishes **Rgbp** from other R packages for similar hierarchical models such as **hglm** (Rönnegård, Shen, and Alam 2010, 2011) for conjugate hierarchical generalized models. The evaluation procedure which we use for this process of "frequency method checking" uses a parametric bootstrapping method that many times generates mock data sets given the fitted

values of the estimated hyper-parameters and uses the mock data to estimate the coverage rates. Using this procedure, we show that the coverage rates of the random effects obtained by our model achieve (or exceed) the nominal confidence level.

The rest of this paper is organized as follows. We specify the Bayesian hierarchical models and discuss their posterior propriety in Section 2. In Section 3, we explain the inferential models used to estimate the model parameters. We describe the estimation procedures including ADM and the acceptance-rejection method in Section 4 and 5, respectively. We introduce frequency method checking techniques in Section 6. We explain the usages of **Rgbp**'s main functions in **Rgbp** in Section 7, and apply them to three data examples in Section 8.

## 2. Conjugate hierarchical modeling structure

**Rgbp** allows users to choose one of three hierarchical models according to the type of data, namely Normal-Normal, Poisson-Gamma, and Binomial-Beta models. Although there are more hierarchical models, we choose the three models because these are based on the most common types of data we may encounter in practice. Also, their conjugacy leads to linear posterior means simplifying computations.

Our parametrization of the three hierarchical models leads to an intuitive shrinkage interpretation in inference and facilitates the estimation procedure because the shrinkage factors under our parametrization are a function of a second-level variance component (Morris 1983).

### 2.1. Normal-Normal model for Gaussian data

The following Normal-Normal hierarchical model (hereafter the Gaussian model) assumed by **Rgbp** is useful when the group-level aggregate data from  $k$  independent groups are continuous (or approximately continuous) variables with known standard errors. The subscript  $j$  below indicates the  $j$ th group among  $k$  groups in the dataset. For  $j = 1, 2, \dots, k$ ,

$$y_j | \mu_j \stackrel{\text{indep.}}{\sim} \text{Normal}(\mu_j, V_j), \quad (1)$$

$$\mu_j | \boldsymbol{\beta}, A \stackrel{\text{indep.}}{\sim} \text{Normal}(\mu_j^E, A), \quad (2)$$

where  $y_j$  is an observed unbiased estimate, e.g., sample mean, for random effect  $\mu_j$ ,  $V_j$  is a completely known standard error of  $y_j$ ,  $\mu_j^E$  is an expected random effect defined as  $E(\mu_j | \boldsymbol{\beta}, A) = \mathbf{x}_j^\top \boldsymbol{\beta} = \beta_1 x_{j,1} + \beta_2 x_{j,2} + \dots + \beta_m x_{j,m}$ , and  $m$  is the number of unknown regression coefficients. **Rgbp** sets  $x_{j,1}$  to 1 for an intercept term as a default. It also provides a usage without the intercept term. It is assumed that the second-level variance  $A$  is unknown and that the  $m \times 1$  regression coefficient vector  $\boldsymbol{\beta}$  is also unknown unless otherwise specified. If no covariates are available, but with an unknown intercept term, then  $\mathbf{x}_j^\top \boldsymbol{\beta} = \beta_1$  ( $m = 1$ ) and thus  $\mu_j^E = \mu^E = \beta_1$  for all  $j$ , resulting in an exchangeable conjugate prior distribution for the random effects. Based on these conjugate prior distributions for random effects, it is easy to derive the conditional posterior distribution of each random effect. For  $j = 1, 2, \dots, k$ ,

$$\mu_j | \boldsymbol{\beta}, A, \mathbf{y} \stackrel{\text{indep.}}{\sim} \text{Normal}((1 - B_j)y_j + B_j\mu_j^E, (1 - B_j)V_j), \quad (3)$$

where  $B_j \equiv V_j / (V_j + A)$  is a shrinkage factor of group  $j$  and  $\mathbf{y} = (y_1, y_2, \dots, y_k)^\top$ . Note that the conditional posterior mean  $E(\mu_j | \boldsymbol{\beta}, A, \mathbf{y})$ , denoted by  $\mu_j^*$ , is a convex combination of the

observed sample mean  $y_j$  and the expected random effect  $\mu_j^E$  weighted by the shrinkage factor  $B_j$ . If the variance of the conjugate prior distribution,  $A$ , is smaller than the variance of the observed distribution,  $V_j$ , then we expect the posterior mean to borrow more information from the second-level conjugate prior distribution.

## 2.2. Poisson-Gamma model for Poisson data

**Rgbp** can estimate a conjugate Poisson-Gamma hierarchical model (hereafter the Poisson model) when the group-level aggregate data from  $k$  independent groups are comprised of non-negative count data without upper limit. However, its usage is limited to the case where the expected random effect,  $\lambda_j^E = \exp(\mathbf{x}_j^\top \boldsymbol{\beta})$ , is known (or equivalently all the regression coefficients are known ( $m = 0$ )); we may be able to obtain this information from the past studies or from experts. For  $j = 1, 2, \dots, k$ ,

$$y_j | \lambda_j \stackrel{\text{indep.}}{\sim} \text{Poisson}(n_j \lambda_j), \quad (4)$$

$$\lambda_j | r \stackrel{\text{indep.}}{\sim} \text{Gamma}(r \lambda_j^E, r), \quad (5)$$

where  $y_j$  is the number of events happening,  $n_j$  is the exposure of group  $j$ , which is not necessarily an integer,  $\lambda_j^E = E(\lambda_j | r)$  is the known expected random effect ( $m = 0$ ), and  $r$  is the unknown second-level variance component. The mean and variance of this conjugate Gamma prior distribution are  $\lambda^E$  and  $\lambda^E/r$ , respectively<sup>1</sup>. [Albert \(1988\)](#) interprets  $r$  as the amount of prior information as  $n_j$  represents the amount of observed information because the uncertainty of the conjugate prior distribution increases as  $r$  decreases and vice versa. The conditional posterior distribution of the random effect  $\lambda_j$  for this Poisson model is

$$\lambda_j | r, \mathbf{y} \stackrel{\text{indep.}}{\sim} \text{Gamma}(r \lambda_j^E + n_j \bar{y}_j, r + n_j), \quad (6)$$

where  $\bar{y}_j \equiv y_j/n_j$ . The mean and variance of the conditional posterior distribution are

$$\lambda_j^* \equiv E(\lambda_j | r, \mathbf{y}) = (1 - B_j) \bar{y}_j + B_j \lambda_j^E \quad \text{and} \quad \text{VAR}(\lambda_j | r, \mathbf{y}) = \frac{\lambda_j^*}{r + n_j}. \quad (7)$$

where  $B_j \equiv r/(r + n_j)$  is the shrinkage factor of group  $j$ , the relative amount of information in the prior compared to the data. The conditional posterior mean is a convex combination of  $\bar{y}_j$  and  $\lambda_j^E$  weighted by  $B_j$ . If the conjugate prior distribution contains more information than the observed data have, *i.e.*, ensemble sample size  $r$  exceeds individual sample size  $n_j$ , then the posterior mean shrinks towards the prior mean by more than 50%.

Note that the conditional posterior variance in (7) is linear in the conditional posterior mean, whereas a slightly different Poisson-Gamma model specification has been used elsewhere ([Christiansen and Morris 1997](#)) that makes the variances quadratic functions of means.

## 2.3. Binomial-Beta model for Binomial data

**Rgbp** can fit a conjugate Binomial-Beta hierarchical model (hereafter the Binomial model) when the group-level aggregate data from  $k$  independent groups are composed of each group's

---

<sup>1</sup>The density function of this Gamma prior distribution in (5) is  $f(\lambda_j | r) \propto \lambda_j^{r \lambda_j^E - 1} \exp(-r \lambda_j)$

number of successes out of total number of trials. The expected random effect in the Binomial model is either known ( $m = 0$ ) or unknown ( $m \geq 1$ ). For  $j = 1, 2, \dots, k$ ,

$$y_j | p_j \stackrel{\text{indep.}}{\sim} \text{Binomial}(n_j, p_j), \quad (8)$$

$$p_j | \boldsymbol{\beta}, r \stackrel{\text{indep.}}{\sim} \text{Beta}(rp_j^E, r(1 - p_j^E)), \quad (9)$$

where  $y_j$  is the number of successes out of  $n_j$  trials,  $p_j^E$  is the expected random effect of group  $j$  defined as  $p_j^E \equiv \mathbb{E}(p_j | \boldsymbol{\beta}, r) = \exp(\mathbf{x}_j^\top \boldsymbol{\beta}) / (1 + \exp(\mathbf{x}_j^\top \boldsymbol{\beta}))$ . The  $m \times 1$  vector of the logistic regression coefficient  $\boldsymbol{\beta}$  and the second-level variance component  $r$  are unknown. The mean and variance of the conjugate Beta prior distribution for group  $j$  are  $p_j^E$  and  $p_j^E(1 - p_j^E)/(r + 1)$ , respectively. The resultant conditional posterior distribution of random effect  $p_j$  is

$$p_j | \boldsymbol{\beta}, r, \mathbf{y} \stackrel{\text{indep.}}{\sim} \text{Beta}(n_j \bar{y}_j + rp_j^E, n_j(1 - \bar{y}_j) + r(1 - p_j^E)), \quad (10)$$

where  $\bar{y}_j = y_j/n_j$  is the observed proportion of group  $j$ . The mean and variance of the conditional posterior distribution are

$$p_j^* \equiv \mathbb{E}(p_j | \boldsymbol{\beta}, r, \mathbf{y}) = (1 - B_j) \bar{y}_j + B_j p_j^E \quad \text{and} \quad \text{VAR}(p_j | \boldsymbol{\beta}, r, \mathbf{y}) = \frac{p_j^*(1 - p_j^*)}{r + n_j + 1}. \quad (11)$$

The conditional posterior mean  $p_j^*$  is a convex combination of  $\bar{y}_j$  and  $p_j^E$  weighted by  $B_j \equiv r/(r + n_j)$  like the Poisson model. If the conjugate prior distribution contains more information than the observed distribution does ( $r > n_j$ ), then the resulting conditional posterior mean borrows more information from the conjugate Beta prior distribution.

## 2.4. Hyper-prior distribution

Hyper-prior distributions are the distributions assigned to the second-level parameters called hyper-parameters. Our choices for the hyper-prior distributions are

$$\boldsymbol{\beta} \sim \text{Uniform on } \mathbf{R}^m \quad \text{and} \quad A \sim \text{Uniform}(0, \infty) \quad (\text{or } \frac{1}{r} \sim \text{Uniform}(0, \infty)). \quad (12)$$

The improper flat hyper-prior distribution on  $\boldsymbol{\beta}$  is a common non-informative choice. In the Gaussian case, the flat hyper-prior distribution on the second-level variance  $A$  is known to produce good repeated sampling coverage properties of the Bayesian interval estimates for the random effects (Morris and Tang 2011; Morris and Lysy 2012; Kelly 2014). The resulting full posterior distribution of the random effects and hyper-parameters is proper if  $k \geq m + 3$  (Morris and Tang 2011; Kelly 2014).

In the other two cases, Poisson and Binomial, the flat prior distribution on  $1/r$  induces the same improper prior distribution on shrinkages ( $f(B_j) \propto B_j^{-2} dB_j$ ) as does  $A$  with the  $\text{Uniform}(0, \infty)$  for the Gaussian case. The Poisson model with this hyper-prior distribution on  $r$ , i.e.,  $dr/r^2$ , provides posterior propriety if there are at least two groups whose observed values  $y_j$  are non-zero and the expected random effects,  $\lambda_j^E$ , are known ( $m = 0$ ); see Appendix A for its proof. If  $\lambda_j^E$  is unknown, **Rgbp** cannot yet give reliable results because we have not verified posterior propriety. If the Poisson is being used as an approximation to the Binomial and the exposures are known integer values, then we recommend using the Binomial model with the same hyper-prior distributions.

As for posterior propriety of the Binomial model, let's define an "interior group" as the group whose number of successes  $y_j$  are neither 0 nor  $n_j$ , and  $k_y$  as the number of interior groups among  $k$  groups. The full posterior distribution of random effects and hyper-parameters is proper if and only if there are at least two interior groups in the data and the  $k_y \times m$  covariate matrix of the interior groups is of full rank  $m$  ( $k_y \geq m$ ) (Tak and Morris 2015).

### 3. The inferential model

The likelihood function of hyper-parameters,  $A$  and  $\beta$ , for the Gaussian model is derived from the independent Normal distributions of the observed data with random effects integrated out;

$$L(A, \beta) = \prod_{j=1}^k f(y_j | A, \beta) = \prod_{j=1}^k \frac{1}{\sqrt{2\pi(A + V_j)}} \exp\left(-\frac{(y_j - \mu_j^E)^2}{2(A + V_j)}\right). \quad (13)$$

The joint posterior density function of hyper-parameters for the Gaussian model is proportional to their likelihood function in (13) because we use flat improper hyper-prior density functions for  $A$  and  $\beta$ ;

$$f(A, \beta | \mathbf{y}) \propto L(A, \beta) dA d\beta. \quad (14)$$

The likelihood function of  $r$  for the Poisson model comes from the independent Negative-Binomial distributions of the observed data with the random effects integrated out;

$$L(r) = \prod_{j=1}^k f(y_j | r) = \prod_{j=1}^k \frac{\Gamma(r\lambda_j^E + y_j)}{\Gamma(r\lambda_j^E)(y_j!)} (1 - B_j)^{y_j} B_j^{r\lambda_j^E}, \quad (15)$$

where  $\Gamma(a)$  is a gamma function defined as  $\int_0^\infty x^{a-1} \exp(-x) dx$  for a positive constant  $a$ . The posterior density function of  $r$  for the Poisson model is the likelihood function in (15) times the hyper-prior density function of  $r$ , i.e.,  $dr/r^2$ ;

$$f(r | \mathbf{y}) \propto L(r) dr/r^2. \quad (16)$$

The likelihood function of hyper-parameters  $r$  and  $\beta$  for the Binomial model is derived from the independent Beta-Binomial distributions of the observed data with random effects integrated out (Skellam 1948);

$$L(r, \beta) = \prod_{j=1}^k f(y_j | r, \beta) = \prod_{j=1}^k \binom{n_j}{y_j} \frac{B(y_j + rp_j^E, n_j - y_j + r(1 - p_j^E))}{B(rp_j^E, r(1 - p_j^E))}, \quad (17)$$

where the notation  $B(a, b) (\equiv \int_0^1 v^{a-1} (1-v)^{b-1} dv)$  indicates a beta function for positive constants  $a$  and  $b$ . The joint posterior density function of hyper-parameters  $f(r, \beta | \mathbf{y})$  for the Binomial model is proportional to their likelihood function in (17) multiplied by the hyper-prior density functions of  $r$  and  $\beta$  in (12);

$$f(r, \beta | \mathbf{y}) \propto L(r, \beta) d\beta dr/r^2. \quad (18)$$

Our goal is to obtain the point and interval estimates of the random effects from their joint unconditional posterior density which can be expressed as for the Gaussian model

$$f(\mu | \mathbf{y}) = \int f(\mu | A, \beta, \mathbf{y}) \cdot f(A, \beta | \mathbf{y}) dA d\beta, \quad (19)$$

where  $\boldsymbol{\mu} = (\mu_1, \mu_2, \dots, \mu_k)^\top$  and the distributions in the integrand are given in (3) and (14). For the Poisson model, the joint unconditional posterior density for the random effects is

$$f(\boldsymbol{\lambda}|\mathbf{y}) = \int f(\boldsymbol{\lambda}|r, \mathbf{y}) \cdot f(r|\mathbf{y})dr, \quad (20)$$

where  $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \dots, \lambda_k)^\top$  and the distributions in the integrand are given in (6) and (16). Lastly, for the Binomial model, the joint unconditional posterior density for the random effects is

$$f(\mathbf{p}|\mathbf{y}) = \int f(\mathbf{p}|r, \boldsymbol{\beta}, \mathbf{y}) \cdot f(r, \boldsymbol{\beta}|\mathbf{y})drd\boldsymbol{\beta}, \quad (21)$$

where  $\mathbf{p} = (p_1, p_2, \dots, p_k)^\top$  and the distributions in the integrand are given in (10) and (18).

## 4. Estimation via the adjustment for density maximization

We illustrate our estimation procedure which utilizes adjustment for density maximization (ADM) (Morris 1988a; Christiansen and Morris 1997; Morris and Tang 2011). ADM is a method to approximate a distribution by a member of Pearson family of distributions and obtain moment estimates via maximization. The ADM procedure for the Gaussian model adopted in **Rgbp** is well documented in Kelly (2014) and so in this section we describe the ADM procedure using the Poisson and Binomial model.

### 4.1. Estimation for shrinkage factors and expected random effects

Our goal here is to estimate the unconditional posterior moments of the shrinkage factors and the expected random effects because they are used to estimate the unconditional posterior moments of the random effects.

#### *Unconditional posterior moments of shrinkage factors*

It is noted that the shrinkage factors are a function of  $r$ , i.e.,  $B_j = B_j(r) = r/(r + n_j)$  (or a function of  $A$  for the Gaussian model). A common method of estimation of  $B_j$  is to approximate the likelihood of  $r$  with two derivatives and use a Delta method for an asymptotic Normal distribution of  $\hat{B}_j(\hat{r}_{MLE})$ . This Normal approximation, however, is defined on  $(-\infty, \infty)$  whereas  $B_j$  lies on the unit interval between 0 and 1, and hence in small sample sizes the Delta method can result in point estimates lying on the boundary of the parameter space, from which the restricted MLE procedure sometimes suffers (Morris and Tang 2011; Kelly 2014).

To continue with a maximization-based estimation procedure but to steer clear of aforementioned boundary issues we make use of ADM. The ADM approximates the distribution of the function of the parameter of interest by one of the Pearson family distributions using the first two derivatives as the Delta method does; the Delta method is a special case of the ADM based on the Normal distribution.

The ADM procedure specified in Morris and Tang (2011) assumes that the unconditional posterior distribution of a shrinkage factor follows a Beta distribution; for  $j = 1, 2, \dots, k$ ,

$$B_j|\mathbf{y} \sim \text{Beta}(a_{1j}, a_{0j}). \quad (22)$$



Note that the mean of Beta distribution  $a_{1j}/(a_{1j} + a_{0j})$  is not the same as its mode  $(a_{j1} - 1)/(a_{j1} + a_{j0} - 2)$ . The ADM works on an adjusted posterior distribution  $f^A(B_j|\mathbf{y})dB_j \propto B_j(1 - B_j)f(B_j|\mathbf{y})dB_j$  so that the mode of  $f^A(B_j|\mathbf{y})$  is the same as the mean of the original Beta distribution. The assumed posterior mean and variance of the  $j$ th shrinkage factor are

$$\mathbb{E}(B_j|\mathbf{y}) = \frac{a_{1j}}{a_{1j} + a_{0j}} = \arg \max_{B_j} f^A(B_j|\mathbf{y}) \equiv B'_j, \quad (23)$$

$$\text{VAR}(B_j|\mathbf{y}) = \frac{B'_j(1 - B'_j)}{a_{1j} + a_{0j} + 1} = \frac{B'_j(1 - B'_j)}{B'_j(1 - B'_j)[- \frac{d^2}{dB_j^2} \log(f^A(B_j|\mathbf{y}))]|_{B_j=B'_j} + 1}. \quad (24)$$

The ADM estimates these mean and variance using the marginal posterior distribution of  $r$ ,  $f(r|\mathbf{y}) \propto L(r)dr/r^2$ . The marginal likelihood,  $L(r) = \int L(\boldsymbol{\beta}, r)d\boldsymbol{\beta}$ , for the Binomial model is obtained via the Laplace approximation with a Lebesgue measure on  $\boldsymbol{\beta}$  and that for the Poisson model is specified in (15).

Considering that (23) and (24) involve the maximization and Hessian calculation, we work on a logarithmic scale of  $r$ , i.e.,  $\alpha = -\log(r)$  (or  $\alpha = \log(A)$  for the Gaussian model), because the distribution of  $\alpha$  is more symmetric than that of  $r$ , and  $\alpha$  is defined on a real line without any boundary issues. Because  $f^A(B_j|\mathbf{y})$  is proportional to the marginal posterior density  $f(\alpha|\mathbf{y}) \propto \exp(\alpha)L(\alpha)d\alpha$  (Morris and Tang 2011), the posterior mean in (23) is estimated by

$$\hat{B}'_j = \frac{\exp(-\hat{\alpha})}{n_j + \exp(-\hat{\alpha})}, \quad (25)$$

where  $\hat{\alpha}$  is the mode of  $f(\alpha|\mathbf{y})$ , i.e.,  $\arg \max_{\alpha} \{\alpha + \log(L(\alpha))\}$ .

We need invariance information introduced in Morris and Tang (2011) to estimate the variance in (24), which is defined as

$$\begin{aligned} I_{\text{inv}} &\equiv - \frac{d^2 \log(f^A(B_j|\mathbf{y}))}{d[\text{logit}(B_j)]^2} \Big|_{B_j=\hat{B}'_j} = - \frac{d^2 \log(f^A(B_j(r)|\mathbf{y}))}{d[\log(r)]^2} \Big|_{r=\hat{r}} \\ &= - \frac{d^2 \log(f^A(B_j(r(\alpha))|\mathbf{y}))}{d\alpha^2} \Big|_{\alpha=\hat{\alpha}} \end{aligned} \quad (26)$$

Note that this invariance information is the negative Hessian value of  $\alpha + \log(L(\alpha))$  at the mode  $\hat{\alpha}$ . Using the invariance information, we estimate the unconditional posterior variance of shrinkage factor in (24) by

$$\widehat{\text{VAR}}(B_j|\mathbf{y}) = \frac{(\hat{B}'_j)^2(1 - \hat{B}'_j)^2}{I_{\text{inv}} + \hat{B}'_j(1 - \hat{B}'_j)}. \quad (27)$$

We obtain the estimates of  $a_{1j}$  and  $a_{0j}$ , the two parameters of the Beta distribution in (22), by matching them to the estimated unconditional posterior mean and variance of  $B_j$  specified in (25) and (27);

$$\hat{a}_{1j} = \frac{I_{\text{inv}}}{1 - \hat{B}'_j} \quad \text{and} \quad \hat{a}_{0j} = \frac{I_{\text{inv}}}{\hat{B}'_j}. \quad (28)$$



The moments of the Beta distribution are well defined as a function of  $a_{1j}$  and  $a_{0j}$ , i.e.,  $E(B_j^c|\mathbf{y}) = B(a_{1j} + c, a_{0j})/B(a_{1j}, a_{0j})$  for  $c \geq 0$ . Their estimates are

$$\hat{E}(B_j^c|\mathbf{y}) = \frac{B(\hat{a}_{1j} + c, \hat{a}_{0j})}{B(\hat{a}_{1j}, \hat{a}_{0j})}. \quad (29)$$

The ADM approximation to the shrinkage factors via Beta distributions is empirically proven to be more accurate than a Laplace approximation (Morris 1988a; Christiansen and Morris 1997; Morris and Tang 2011; Morris and Lysy 2012).

### *Unconditional posterior moments of expected random effects*

We estimate the unconditional posterior moments of expected random effects using their relationship to the conditional posterior moments. For a non-negative constant  $c$ , the unconditional posterior moments are

$$E((p_j^E)^c|\mathbf{y}) = E(E((p_j^E)^c|\alpha, \mathbf{y})|\mathbf{y}). \quad (30)$$

We approximate the unconditional posterior moments on the left hand side by the conditional posterior moments with  $\hat{\alpha}$  plugged-in (Kass and Steffey 1989), i.e., by  $E((p_j^E)^c|\hat{\alpha}, \mathbf{y})$ .

However, calculating conditional posterior moments of each expected random effect involves an intractable integration. For example, the first conditional posterior moment of  $p_j^E$  is

$$E(p_j^E|\hat{\alpha}, \mathbf{y}) = E\left(\frac{\exp(x_j^\top \beta)}{1 + \exp(x_j^\top \beta)} \middle| \hat{\alpha}, \mathbf{y}\right) = \int_{\mathbf{R}^m} \frac{\exp(x_j^\top \beta)}{1 + \exp(x_j^\top \beta)} f(\beta|\hat{\alpha}, \mathbf{y}) d\beta. \quad (31)$$

Thus, we use another ADM, assuming the conditional posterior distribution of each expected random effect is a Beta distribution as follows;

$$p_j^E|\hat{\alpha}, \mathbf{y} = \frac{\exp(x_j^\top \beta)}{1 + \exp(x_j^\top \beta)} \middle| \hat{\alpha}, \mathbf{y} \sim \text{Beta}(b_{1j}, b_{0j}) \sim \frac{G(b_{1j})}{G(b_{1j}) + G(b_{0j})}, \quad (32)$$

where  $G(b_{1j})$  is a random variable following a  $\text{Gamma}(b_{1j}, 1)$  distribution and independently  $G(b_{0j})$  has a  $\text{Gamma}(b_{0j}, 1)$  distribution. Note that the representation in (32) is equivalent to  $\exp(x_j^\top \beta)|\hat{\alpha}, \mathbf{y} \sim G(b_{1j})/G(b_{0j})$ , a ratio of two independent Gamma random variables. Its mean and variance are

$$E(\exp(x_j^\top \beta)|\hat{\alpha}, \mathbf{y}) = E\left(\frac{G(b_{1j})}{G(b_{0j})}\right) = \frac{b_{1j}}{b_{0j} - 1} \equiv \eta_j, \quad (33)$$

$$\text{VAR}(\exp(x_j^\top \beta)|\hat{\alpha}, \mathbf{y}) = \text{VAR}\left(\frac{G(b_{1j})}{G(b_{0j})}\right) = \frac{\eta_j(1 + \eta_j)}{b_{0j} - 2}. \quad (34)$$

In order to estimate  $b_{1j}$  and  $b_{0j}$ , we assume that the conditional posterior distribution of  $\beta$  given  $\hat{\alpha}$  and  $\mathbf{y}$  follows a Normal distribution with mean  $\hat{\beta}$  and variance-covariance matrix  $\hat{\Sigma}$ , where  $\hat{\beta}$  is the mode of  $f(\beta|\hat{\alpha}, \mathbf{y})$  and  $\hat{\Sigma}$  is an inverse of the negative Hessian matrix at the mode. Thus, the posterior distribution of  $x_j^\top \beta$  is also Normal with mean  $x_j^\top \hat{\beta}$  and variance  $x_j^\top \hat{\Sigma} x_j$ .

Using the property of the log-Normal distribution for  $\exp(x_j^\top \beta)$ , we estimate the posterior mean and variance in (33) and (34) as

$$\hat{\mathbb{E}}(\exp(x_j^\top \beta) | \hat{\alpha}, \mathbf{y}) = \exp(x_j^\top \hat{\beta} + x_j^\top \hat{\Sigma} x_j / 2) = \hat{\eta}_j, \quad (35)$$

$$\widehat{\text{VAR}}(\exp(x_j^\top \beta) | \mathbf{y}) = \hat{\eta}_j^2 (\exp(x_j^\top \hat{\Sigma} x_j) - 1). \quad (36)$$

We estimate the values of  $b_{1j}$  and  $b_{0j}$  by matching them to the estimated unconditional posterior mean and variance of  $\exp(x_j^\top \beta)$  in (35) and (36);

$$\hat{b}_{1j} = \hat{\eta}_j(\hat{b}_{0j} - 1) \quad \text{and} \quad \hat{b}_{0j} = \frac{1 + \hat{\eta}_j}{\hat{\eta}_j(\exp(x_j^\top \hat{\Sigma} x_j) - 1)} + 2. \quad (37)$$

Finally, we estimate the unconditional posterior moments of the expected random effects by

$$\hat{\mathbb{E}}((p_j^E)^c | \hat{\alpha}, \mathbf{y}) = \frac{B(\hat{b}_{1j} + c, \hat{b}_{0j})}{B(\hat{b}_{1j}, \hat{b}_{0j})} \quad \text{for } c \geq 0. \quad (38)$$

The ADM approximation to a log-Normal density via a F distribution (represented by a ratio of two independent Gamma random variables) is known to be more accurate than the Laplace approximation (Morris 1988a).

For the Gaussian model (Morris and Tang 2011), the conditional posterior distribution of  $\beta$  given  $\hat{A}$  and  $\mathbf{y}$  is Normal whose mean and variance-covariance matrix are

$$(X^\top D_{V+\hat{A}}^{-1} X)^{-1} X^\top D_{V+\hat{A}}^{-1} \mathbf{y} \quad \text{and} \quad (X^\top D_{V+\hat{A}}^{-1} X)^{-1}, \quad (39)$$

respectively, where  $X \equiv (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k)^\top$  is a  $k \times m$  covariate matrix and  $D_{V+\hat{A}}$  is a  $k \times k$  diagonal matrix with the  $j$ -th diagonal element equal to  $V_j + \hat{A}$ . Because  $\mathbf{x}^\top \beta$  given  $\hat{A}$  and  $\mathbf{y}$  is also Normally distributed, we easily obtain the conditional posterior moments of  $\mu_j^E = \mathbf{x}^\top \beta$  given  $\hat{A}$  and use them to estimate unconditional posterior moments of  $\mu_j^E$ .

## 4.2. Estimation for random effects

It is intractable to derive analytically the unconditional posterior distribution of each random effect for the three models. Thus we approximate the distributions by matching the estimated posterior moments with a skewed-Normal distribution (Azzalini 1985) for the Gaussian model, a Beta distribution for the Binomial model and a Gamma distribution for the Poisson model; for  $j = 1, 2, \dots, k$ ,

$$\mu_j | \mathbf{y} \sim \text{skewed-Normal}(\phi, \omega, \delta), \quad (40)$$

$$\lambda_j | \mathbf{y} \sim \text{Gamma}(s_{1j}, s_{0j}), \quad (41)$$

$$p_j | \mathbf{y} \sim \text{Beta}(t_{1j}, t_{0j}), \quad (42)$$

where  $(\phi, \omega, \delta)$  of the skewed-Normal distribution are location, scale, and skewness parameters, respectively.

Morris and Lysy (2012) first noted that the unconditional posterior distribution of the random effect in a two-level conjugate Gaussian model might be skewed. Kelly (2014) shows that the

skewed-Normal approximation to the unconditional posterior distribution of the random effect is better than a Normal approximation ( $\mu_j|\mathbf{y} \sim \text{Normal}$ ) in terms of the repeated sampling coverage properties of random effects. Kelly (2014) estimates the first three moments of the random effects by noting that  $\mu_j|\mathbf{y}, B_j$  is Normally distributed and thus estimates the moments by using the ADM approximation of the shrinkage factors,  $B_j$ , and the law of third cumulants (Brillinger 1969). The three estimated moments are then matched to the first three moments of the skewed-Normal distribution,  $E(\mu_j|\mathbf{y}) = \phi + \omega\delta\sqrt{2/\pi}$ ,  $\text{Var}(\mu_j|\mathbf{y}) = \omega^2(1 - 2\delta^2/\pi)$ , and  $\text{Skewness}(\mu_j|\mathbf{y}) = (4 - \pi)\delta^3/[2(\pi/2 - \delta^2)^{3/2}]$  (Azzalini 1985). The full derivation can be found in Kelly (2014).

The unconditional posterior mean and variance of random effect  $\lambda_j$  in the Poisson model are

$$E(\lambda_j|\mathbf{y}) = E(E(\lambda_j|r, \mathbf{y})|\mathbf{y}) = (1 - E(B_j|\mathbf{y}))\bar{y}_j + E(B_j|\mathbf{y})\lambda_j^E, \quad (43)$$

$$\text{VAR}(\lambda_j|\mathbf{y}) = E(\text{VAR}(\lambda_j|r, \mathbf{y})|\mathbf{y}) + \text{VAR}(E(\lambda_j|r, \mathbf{y})|\mathbf{y}) \quad (44)$$

$$= E(\lambda_j^*/(r + n_j)|\mathbf{y}) + \text{VAR}(B_j(\bar{y}_j - \lambda_j^E)|\mathbf{y}) \quad (45)$$

$$= \frac{1}{n_j} [\bar{y}_j E((1 - B_j)^2|\mathbf{y}) + \lambda_j^E E((1 - B_j)B_j|\mathbf{y})] + (\bar{y}_j - \lambda_j^E)^2 \text{VAR}(B_j|\mathbf{y}). \quad (46)$$

To estimate these, we plug the estimated unconditional posterior moments of shrinkage factors in (29) into both (43) and (46). Let  $\hat{\mu}_{\lambda_j}$  and  $\hat{\sigma}_{\lambda_j}^2$  denote the estimated unconditional posterior mean and variance, respectively. The estimates of the two parameters  $s_{1j}$  and  $s_{0j}$  in (41) are

$$\hat{s}_{1j} = \frac{\hat{\mu}_{\lambda_j}^2}{\hat{\sigma}_{\lambda_j}^2}, \text{ and } \hat{s}_{0j} = \frac{\hat{\mu}_{\lambda_j}}{\hat{\sigma}_{\lambda_j}^2}. \quad (47)$$

To estimate the unconditional posterior moments of random effects in the Binomial model, we assume that hyper-parameters  $r$  and  $\beta$  are independent *a posteriori*. With this assumption, the unconditional posterior mean and variance of random effect  $p_j$  are

$$E(p_j|\mathbf{y}) = E(E(p_j|r, \beta, \mathbf{y})|\mathbf{y}) = (1 - E(B_j|\mathbf{y}))\bar{y}_j + E(B_j|\mathbf{y})E(p_j^E|\mathbf{y}), \quad (48)$$

$$\text{VAR}(p_j|\mathbf{y}) = E(\text{VAR}(p_j|r, \beta, \mathbf{y})|\mathbf{y}) + \text{VAR}(E(p_j|r, \beta, \mathbf{y})|\mathbf{y}) \quad (49)$$

$$= E(p_j^*(1 - p_j^*)/(r + n_j + 1)|\mathbf{y}) + \text{VAR}(B_j(\bar{y}_j - p_j^E)|\mathbf{y}) \quad (50)$$

$$\approx E(p_j^*(1 - p_j^*)(1 - B_j)/n_j|\mathbf{y}) + \text{VAR}(B_j(\bar{y}_j - p_j^E)|\mathbf{y}) \quad (51)$$

$$= \{(1 - \bar{y}_j)\bar{y}_j[1 - E(B_j|\mathbf{y})] + (2\bar{y}_j - 1)E(B_j(1 - B_j)|\mathbf{y})(\bar{y}_j - E(p_j^E|\mathbf{y})) \\ + E(B_j^2(1 - B_j)|\mathbf{y})E((\bar{y}_j - p_j^E)^2|\mathbf{y})\}/n_j + \text{VAR}(B_j(\bar{y}_j - p_j^E)|\mathbf{y}), \quad (52)$$

where the approximation in (51) is a first-order Taylor approximation. By plugging the estimated unconditional posterior moments of shrinkage factors in (29) and those of expected random effect in (38) into both (48) and (52), we obtain the estimates of the unconditional posterior mean and variance of each random effect, denoted by  $\hat{\mu}_{p_j}$  and  $\hat{\sigma}_{p_j}^2$ , respectively. Finally, we obtain the estimates of two parameters  $t_{1j}$  and  $t_{0j}$  in (42) as follows;

$$\hat{t}_{1j} = \left( \frac{\hat{\mu}_{p_j}(1 - \hat{\mu}_{p_j})}{\hat{\sigma}_{p_j}^2} - 1 \right) \hat{\mu}_{p_j}, \text{ and } \hat{t}_{0j} = \left( \frac{\hat{\mu}_{p_j}(1 - \hat{\mu}_{p_j})}{\hat{\sigma}_{p_j}^2} - 1 \right) (1 - \hat{\mu}_{p_j}). \quad (53)$$

Finally, the assumed unconditional posterior distribution of random effect for the Gaussian model is

$$\mu_j|\mathbf{y} \sim \text{skewed-Normal}(\hat{\phi}, \hat{\omega}, \hat{\delta}), \quad (54)$$

that for the Poisson model is

$$\lambda_j | \mathbf{y} \sim \text{Gamma}(\hat{s}_{1j}, \hat{s}_{0j}). \quad (55)$$

and that for the Binomial model is

$$p_j | \mathbf{y} \sim \text{Beta}(\hat{t}_{1j}, \hat{t}_{0j}), \quad (56)$$

Our point and interval estimates of each random effect are the mean and (2.5%, 97.5%) quantiles (if we assign 95% confidence level) of the assumed unconditional posterior distribution in (54), (56), or (55).

## 5. The acceptance-rejection method for the Binomial model

As for the Binomial model, the package **Rgbp** also provides a way to draw posterior samples of random effects and hyper-parameters via the acceptance-rejection (A-R) method (Robert and Casella 2013). Unlike the approximate Bayesian machinery specified in the previous section, this method does not assume that hyper-parameters are independent *a posteriori*. The joint posterior density function of  $\alpha = -\log(r)$  and  $\beta$  based on their joint hyper-prior density function in (12) is

$$f(\alpha, \beta | \mathbf{y}) \propto f(\alpha, \beta) L(\alpha, \beta) \propto \exp(\alpha) L(\alpha, \beta) d\alpha d\beta. \quad (57)$$

The A-R method is useful when it is difficult to sample a parameter of interest  $\theta$  directly from its target probability density  $f(\theta)$ , which is known up to a normalizing constant, but an easy-to-sample envelope function  $g(\theta)$  is available. The A-R method samples  $\theta$  from the envelope  $g(\theta)$  and accepts it with a probability  $\frac{f(\theta)}{Mg(\theta)}$ , where  $M$  is a constant making  $f(\theta)/g(\theta) \leq M$  for all  $\theta$ . The distribution of the accepted  $\theta$  exactly follows  $f(\theta)$ . The A-R method is stable as long as the tails of the envelop function are thicker than those of the target density function. The goal of the A-R method for the Binomial model is to draw posterior samples of hyper-parameters from (57), using an easy-to-sample envelop function  $g(\alpha, \beta)$  that has thicker tails than the target density function.

We factor the envelope function into two parts,  $g(\alpha, \beta) = g_1(\alpha)g_2(\beta)$  to model the tails of each function separately. We consider the tail behavior of the conditional posterior density function  $f(\alpha | \beta, \mathbf{y})$  to come up with  $g_1(\alpha)$ ;  $f(\alpha | \beta, \mathbf{y})$  behaves as  $\exp(-\alpha(k-1))$  when  $\alpha$  goes to  $\infty$  and as  $\exp(\alpha)$  when  $\alpha$  goes to  $-\infty$ . It indicates that  $f(\alpha | \beta, \mathbf{y})$  is skewed to the left because the right tail touches the  $x$ -axis faster than the left tail does as long as  $k > 1$ . A skewed  $t$ -distribution is a good candidate for  $g_1(\alpha)$  because it behaves as a power law on both tails, leading to thicker tails than those of  $f(\alpha | \beta, \mathbf{y})$ .

It is too complicated to figure out the tail behaviors of  $f(\beta | \alpha, \mathbf{y})$ . However, because  $f(\beta | \alpha, \mathbf{y})$  in the Gaussian model (as an approximation) has a multivariate Gaussian density function (Morris and Tang 2011; Kelly 2014), we consider a multivariate  $t$ -distribution with 4 degrees of freedom as a good candidate for  $g_2(\beta)$ .

Specifically, we assume

$$g_1(\alpha) = g_1(\alpha; \mu, \sigma, a, b) \equiv \text{Skewed-}t(\alpha | \mu, \sigma, a, b), \quad (58)$$

$$g_2(\beta) = g_2(\beta; \xi, S_{(m \times m)}) \equiv t_4(\beta | \xi, S), \quad (59)$$

where Skewed- $t(\alpha|\mu, \sigma, a, b)$  represents a density function of a skewed  $t$ -distribution at  $\alpha$  with location  $\mu$ , scale  $\sigma$ , degree of freedom  $a + b$ , and skewness  $a - b$  for any positive constants  $a$  and  $b$  (Jones and Faddy 2003). Jones and Faddy (2003) derive the mode of  $g_1(\alpha)$  as

$$\mu + \frac{(a - b)\sqrt{a + b}}{\sqrt{(2a + 1)(2b + 1)}}, \quad (60)$$

and provide a representation to generate random variables that follows Skewed- $t(\alpha|\mu, \sigma, a, b)$ ;

$$\alpha \sim \mu + \sigma \frac{\sqrt{a + b}(2T - 1)}{2\sqrt{T(1 - T)}}, \text{ where } T \sim \text{Beta}(a, b). \quad (61)$$

They also show that the tails of the skewed- $t$  density function follow a power law with  $\alpha^{-(2a+1)}$  on the left and  $\alpha^{-(2b+1)}$  on the right when  $b > a$ .

The notation  $t_4(\beta|\xi, S)$  in (59) indicates a density function of a multivariate  $t$ -distribution at  $\beta$  with 4 degrees of freedom, a location vector  $\xi$ , and a  $m \times m$  scale matrix  $S$  that leads to the variance-covariance matrix  $2S$ .

**Rgbp** determines the parameters of  $g_1(\alpha)$  and  $g_2(\beta)$ , i.e.,  $\mu$ ,  $\sigma$ ,  $a$ ,  $b$ ,  $\xi$ , and  $S$ , to make the product of  $g_1(\alpha)$  and  $g_2(\beta)$  similar to the target joint posterior density  $f(\alpha, \beta|\mathbf{y})$ . First, **Rgbp** obtains the mode of  $f(\alpha, \beta|\mathbf{y})$  and the inverse of the negative Hessian matrix at the modes,  $-H^{-1}$ . Let  $(\hat{\alpha}, \hat{\beta})$  denote the modes of  $f(\alpha, \beta|\mathbf{y})$ ,  $-H_{\hat{\alpha}}^{-1}$  indicate (1, 1) element of  $-H^{-1}$ , and  $-H_{\hat{\beta}}^{-1}$  represent  $-H^{-1}$  without the first row and column.

For  $g_1(\alpha)$ , **Rgbp** sets  $(a, b)$  to  $(k, 2k)$  if  $k$  is less than 10 (or otherwise to  $(\log(k), 2\log(k))$ ) for a left-skewness and for small values of  $a$  and  $b$  (thick tails). **Rgbp** matches the mode of  $g_1(\alpha)$  specified in (60) to  $\hat{\alpha}$  by setting the location parameter  $\mu$  to  $\hat{\alpha} - (a - b)\sqrt{a + b} / \sqrt{(2a + 1)(2b + 1)}$ . **Rgbp** sets the scale parameter  $\sigma$  to  $(-H_{\hat{\alpha}}^{-1})^{0.5}\psi$ , where  $\psi$  is a tuning parameter; when the A-R method produces extreme weights defined in (62) below, we need enlarge the value of  $\psi$ . For  $g_2(\beta)$ , **Rgbp** sets the location vector  $\xi$  to the mode  $\hat{\beta}$  and the scale matrix  $S$  to  $-H_{\hat{\beta}}^{-1}/2$  so that the variance-covariance matrix becomes  $-H_{\hat{\beta}}^{-1}$ .

For the implementation of the acceptance-rejection method, **Rgbp** draws four times more trial samples than the desired number of samples, denoted by  $N$ , independently from  $g_1(\alpha)$  and  $g_2(\beta)$ . **Rgbp** calculates  $4N$  weights, each of which is defined as

$$w_i \equiv w(\alpha^{(i)}, \beta^{(i)}) = \frac{f(\alpha^{(i)}, \beta^{(i)}|\mathbf{y})}{g_1(\alpha^{(i)})g_2(\beta^{(i)})}, \text{ for } i = 1, 2, \dots, 4N. \quad (62)$$

**Rgbp** accepts each pair of  $(\alpha^{(i)}, \beta^{(i)})$  with a probability  $w_i/M$  where  $M$  is set to the maximum of all the  $4N$  weights. When **Rgbp** accepts more than  $N$  pairs, it discards the redundant. If **Rgbp** accepts less than  $N$  pairs, then it additionally draws  $N'$  (six times the shortage) pairs and calculates a new maximum  $M'$  from all the previous and new weights; **Rgbp** accepts or rejects the entire pairs again with new probabilities  $w_j/M'$ ,  $j = 1, 2, \dots, 4N + N'$ .

After obtaining posterior samples of hyper-parameters, **Rgbp** draws posterior samples of random effects from  $f(\mathbf{p}|\mathbf{y})$  in (21). The integration on the right hand side of (21) can be done by sampling  $\mathbf{p}$  from  $f(\mathbf{p}|\beta, r, \mathbf{y})$  for  $j = 1, 2, \dots, k$  in (42), given  $r = \exp(-\alpha)$  and  $\beta$  that are already sampled from  $f(\alpha, \beta|\mathbf{y})$  via the A-R method.

## 6. Frequency method checking

The question as to whether the interval estimates of random effects for given confidence level obtained by a specific model achieve the nominal coverage rate for any true parameter values is one of the key model evaluation criteria. Unlike standard model checking methods that test whether a two-level model is appropriate for data (Dean 1992; Christiansen and Morris 1996), *frequency method checking* is a procedure to evaluate the coverage properties of the model. Conditioning that the two-level model is appropriate, the frequency method checking generates pseudo-data sets given specific values of hyper-parameters and estimates unknown coverage probabilities based on these mock data sets (a parametric bootstrapping). We describe the frequency method checking based on the Gaussian model because the idea can be easily applied to the other two models.

### 6.1. Pseudo-data generation

Figure 1 displays the process of generating pseudo-data sets. It is noted that the conjugate prior distribution of each random effect in (2) is completely determined by two hyper-parameters,  $A$  and  $\beta$ . Fixing these hyper-parameters at specific values, we generate  $N_{\text{sim}}$  sets of random effects from the conjugate prior distribution, i.e.,  $\{\mu^{(i)}, i = 1, \dots, N_{\text{sim}}\}$ , where the superscript  $(i)$  indicates the  $i$ -th simulation. Next, using the distribution of observed data in (1), we generate  $N_{\text{sim}}$  sets of observed data sets  $\{y^{(i)}, i = 1, \dots, N_{\text{sim}}\}$  given each  $\mu^{(i)}$ . Note that we generate one observed data set per one set of random effects.

### 6.2. Coverage probability estimation

After fitting the Gaussian model on each simulated data set, we obtain interval estimates of random effects  $\mu$ . Let  $(\hat{\mu}_{j, \text{low}}^{(i)}, \hat{\mu}_{j, \text{upp}}^{(i)})$  represent the lower and upper bounds of the interval estimate of random effect  $j$  based on the  $i$ -th simulation given a specific confidence level. Let's define a coverage indicator of random effect  $j$  on the  $i$ -th mock data set as

$$I_{A, \beta}(\mu_j^{(i)}) = \begin{cases} 1, & \text{if } \mu_j^{(i)} \in (\hat{\mu}_{j, \text{low}}^{(i)}, \hat{\mu}_{j, \text{upp}}^{(i)}) \\ 0, & \text{otherwise} \end{cases} \quad (63)$$

The subscript,  $A$  and  $\beta$ , indicates that an outcome of the coverage indicator depends on the simulated random effects and mock data generated by  $A$  and  $\beta$ .

*Simple unbiased coverage estimator.*

When the confidence level is 95%, the proportion of 95% interval estimates that contain random effect  $j$  is an intuitive choice for the coverage rate estimator for random effect  $j$ . This

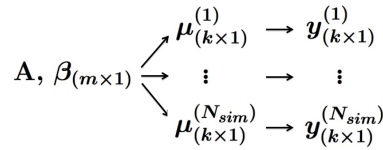


Figure 1: Pseudo-data generating process.

estimator implicitly assumes that there exist  $k$  unknown coverage probabilities of random effects, denoted by  $C_{A,\beta}(\mu_j)$  for  $j = 1, 2, \dots, k$ , depending on the values of the hyper-parameters that generate random effects and mock data sets. The coverage indicators for random effect  $j$  in (63) is assumed to follow an independent and identically distributed Bernoulli distribution given the unknown coverage rate  $C_{A,\beta}(\mu_j)$ . The sample mean of these coverage indicators is a simple unbiased coverage estimator for  $C_{A,\beta}(\mu_j)$ ;

$$\bar{I}_{A,\beta}(\mu_j) = \frac{1}{N_{\text{sim}}} \sum_{i=1}^{N_{\text{sim}}} I_{A,\beta}(\mu_j^{(i)}), \text{ for } j = 1, 2, \dots, k. \quad (64)$$

Note that  $\bar{I}_{A,\beta}(\mu_j)$  averages over possible values of  $\mu_j$  and  $y_j$  generated by specific values of  $A$  and  $\beta$ . The unbiased variance estimator of  $\text{VAR}(\bar{I}_{A,\beta}(\mu_j))$  is

$$\widehat{\text{VAR}}(\bar{I}_{A,\beta}(\mu_j)) = \frac{1}{N_{\text{sim}}(N_{\text{sim}} - 1)} \sum_{i=1}^{N_{\text{sim}}} (I_{A,\beta}(\mu_j^{(i)}) - \bar{I}_{A,\beta}(\mu_j))^2, \text{ for } j = 1, 2, \dots, k. \quad (65)$$

*Rao-Blackwellized unbiased coverage estimator.*

The frequency method checking is computationally expensive in nature because it fits a model on every mock data set. The situation deteriorates if the number of simulations or the size of data is large, or the estimation method is computationally demanding. Christiansen and Morris (1997) and Tang (2002) used a Rao-Blackwellized (RB) unbiased coverage estimator for the unknown coverage rate of each random effect, which is more efficient than the simple unbiased coverage estimator. For  $j = 1, 2, \dots, k$ ,

$$C_{A,\beta}(\mu_j) = \text{E}(\bar{I}_{A,\beta}(\mu_j)|A, \beta) = E\left[\frac{1}{N_{\text{sim}}} \sum_{i=1}^{N_{\text{sim}}} \text{E}(I_{A,\beta}(\mu_j^{(i)})|A, \beta, \mathbf{y}^{(i)}) \middle| A, \beta\right], \quad (66)$$

where the sample mean of the interior conditional expectations in (66) is the RB unbiased coverage estimator. Specifically,

$$\bar{I}_{A,\beta}^{RB}(\mu_j) = \frac{1}{N_{\text{sim}}} \sum_{i=1}^{N_{\text{sim}}} \text{E}(I_{A,\beta}(\mu_j^{(i)})|A, \beta, \mathbf{y}^{(i)}) \quad (67)$$

$$= \frac{1}{N_{\text{sim}}} \sum_{i=1}^{N_{\text{sim}}} \text{P}(\mu_j^{(i)} \in (\hat{\mu}_{j, \text{low}}^{(i)}, \hat{\mu}_{j, \text{upp}}^{(i)})|A, \beta, \mathbf{y}^{(i)}). \quad (68)$$

We can easily compute the above conditional posterior probabilities in (68) using the cumulative density function of the Gaussian conditional posterior distribution of each random effect in (3). The variance of  $\bar{I}_{A,\beta}^{RB}(\mu_j)$  does not exceed the variance of a simple unbiased coverage estimator,  $\bar{I}_{A,\beta}(\mu_j)$  (Rao 1945; Blackwell 1947).

If one dataset  $\mathbf{y}^{(i)}$  is simulated per one set of random effects  $\boldsymbol{\mu}^{(i)}$ , the variance estimator below is an unbiased estimator of  $\text{VAR}(\bar{I}_{A,\beta}^{RB}(\mu_j))$ . For  $j = 1, 2, \dots, k$ ,

$$\widehat{\text{VAR}}(\bar{I}_{A,\beta}^{RB}(\mu_j)) \equiv \frac{1}{N_{\text{sim}}(N_{\text{sim}} - 1)} \sum_{i=1}^{N_{\text{sim}}} \left( \text{E}(I_{A,\beta}(\mu_j^{(i)})|A, \beta, \mathbf{y}^{(i)}) - \bar{I}_{A,\beta}^{RB}(\mu_j) \right)^2. \quad (69)$$



*Rao-Blackwellized overall unbiased coverage estimator*

Assuming that the unknown coverage probabilities are the same for all random effects, we use the Rao-Blackwellized overall unbiased coverage estimator and its variance estimator;

$$\bar{I}_{r,\beta}^{RB} = \frac{1}{k} \sum_{j=1}^k \bar{I}_{r,\beta}^{RB}(p_j) \quad \text{and} \quad \widehat{\text{VAR}}(\bar{I}_{RB}) = \frac{1}{k^2} \sum_{j=1}^k \widehat{\text{VAR}}(\bar{I}_{r,\beta}^{RB}(p_j)). \quad (70)$$

## 7. Usage of functions in Rgbp

In this section, we describe the usage of the two main functions of **Rgbp**, i.e., **gbp** for model fitting and **coverage** for frequency method checking.

### 7.1. Model fitting

The function **gbp** creates an S3 object “gbp” on which three generic functions **plot**, **print**, and **summary** are defined.

There are two cases according to whether covariates are available or not. When no covariates are available, the function **gbp** requires fitting an intercept term or designating known values of the expected random effects, i.e., the intercept term must be either estimated or known. The default of **gbp** is to fit an intercept term. The value(s) of the known expected random effect(s) can be assigned through an optional argument **mean.PriorDist**. Note that **gbp** can fit the Poisson model only when the values of expected random effects,  $\lambda_j^E$ , are known. The usage of **gbp** to fit each model without any covariates is

```
R> g.output <- gbp(y, se.or.n, model = "gaussian")
R> b.output <- gbp(y, se.or.n, model = "binomial")
R> p.output <- gbp(y, se.or.n, mean.PriorDist, model = "poisson")
```

The argument **y** is a vector of  $k$  observed sample means for the Gaussian model,  $k$  observed numbers of successful outcomes for the Binomial model, and  $k$  observed outcome counts for the Poisson model. The argument **se.or.n** is a vector of  $k$  standard errors of each sample mean for the Gaussian model,  $k$  numbers of trials for the Binomial model, and  $k$  exposures for the Poisson model. The argument **mean.PriorDist** is either a constant (if all the known expected random effects are the same) or a  $k \times 1$  vector of known expected random effects.

If covariate information for each group is available, users can fit the Gaussian and Binomial models, using the following codes.

```
R> g.output <- gbp(y, se.or.n, X, model = "gaussian")
R> b.output <- gbp(y, se.or.n, X, model = "binomial")
```

The argument **X** is a matrix of covariate(s) each column of which corresponds to one covariate for  $k$  groups. For example, if users have two covariates for each group, the argument **X** must be  $k \times 2$  matrix to estimate three regression coefficients  $\beta = (\beta_1, \beta_2, \beta_3)$  including an intercept term,  $\beta_1$ , as a default. If users do not want to include an intercept term ( $\beta_1 = 0$ ), estimating two regression coefficients for the two covariates, users can add an optional argument **intercept** as follows.

```
R> g.output <- gbp(y, se.or.n, X, model = "gaussian", intercept = FALSE)
```

The function `gbp` contains several optional arguments. The argument `Alpha`, whose default value is 0.95, sets the confidence level, producing  $100 \times \text{Alpha}\%$  interval estimates for the random effects. For the Gaussian model, setting the argument `normal.CI` to `TRUE` lets `gbp` use a Normal approximation to the unconditional posterior distribution of the random effect (Morris and Tang 2011). The default value of `normal.CI` is `FALSE` for the skewed-Normal approximation (Kelly 2014).

The function `gbp` uses the A-R method to fit the Binomial model if users assign the desired number of posterior samples ( $N$  in (62)) through the argument `n.AR`; its default value is 0. There are several arguments related to the A-R method. The argument `n.AR.factor` determines how many trial samples the method draws; its default value is 4, meaning that the function `gbp` draws  $n.AR \times 4$  trial samples and accepts or rejects them. The argument `trial.scale` is  $\psi$  determining the scale parameter of the skewed- $t$  distribution; its default value is 1.3. The argument `save.result` indicates whether `gbp` saves the whole posterior samples of the random effects and hyper-parameters; its default value is `TRUE`. The two arguments `t` and `u`, taking on positive values, allow users to choose the joint hyper-prior density function,  $f(r, \beta) \propto d\beta dr / (t + r)^{u+1}$ ; the default values for `t` and `u` are 0 and 1, respectively, for the joint hyper-prior density function specified in (12).

For example, when there are two covariates, the following code produces 2,000 posterior samples of random effects and hyper-parameters,  $r$  and  $\beta_{(3 \times 1)}$  including an intercept term, via the A-R method with 8,000 trial samples.

```
R> b.output <- gbp(y, se.or.n, X, model = "binomial", n.AR = 2000)
```

The object `b.output` above contains 8,000 weights (`b.output$weight`), 2,000 posterior samples of  $\alpha$  (`b.output$alpha`), a  $2,000 \times 3$  matrix of  $\beta$  (`b.output$beta`) each column of which corresponds to 2,000 posterior samples of each regression coefficient, and a  $k \times 2,000$  matrix of random effects (`b.output$p`) each row of which has posterior samples of each random effect.

The S3 object “gbp” benefits from three generic functions, `print`, `summary`, and `plot`. The estimation result for all the random effects appears if users type the “gbp” object in the R console, which plays the same role of the function `print` with its default argument “`sort = TRUE`”. When the argument `sort` is set to `TRUE`, the function `print` prints out the estimation result for all the groups in the ascending order of  $n$  for the Binomial and Poisson model and the descending order of standard errors for the Gaussian model. When the argument `sort` is `FALSE`, the estimation result is returned in the order of data input.

```
R> b.output
R> print(b.output, sort = FALSE)
```

The function `summary` prints a detailed estimation result, including the estimation result for the hyper-parameters,  $A$  (or  $r$ ) and  $\beta$ .

```
R> summary(b.output)
```

The function `plot` draws a shrinkage plot and  $100 \times \text{Alpha}\%$  interval plot for random effects. Its default argument “`sort = TRUE`” displays the  $100 \times \text{Alpha}\%$  interval plot in the ascending

order of  $n$  for the Binomial and Poisson model and the descending order of standard errors for the Gaussian model. When the argument `sort` is set to `FALSE` the  $100 \times \text{Alpha}\%$  interval plot is displayed in the order of data input.

```
R> plot(b.output)
R> plot(b.output, sort = FALSE)
```

## 7.2. Frequency method checking

The function `coverage` conducts frequency method checking. It estimates the coverage properties for our estimators of the random effects at a particular value of the hyperparameters by averaging the coverage over many simulated datasets. The basic usage of `coverage` needs a “gbp” object, such as `b.output` above, as the first argument;

```
R> cov <- coverage(b.output, nsim = 1000)
```

The argument `nsim` sets the number of simulations,  $N_{\text{sim}}$ , defined in Section 6.1. If users do not assign values of the hyper-parameters through the arguments `A.or.r` and `reg.coef`, then the function `coverage` automatically sets the estimated posterior modes of hyper-parameters saved in the “gbp” object (or their posterior medians if the acceptance-rejection method for the Binomial model is used) as the generative values of hyper-parameters. If users want to conduct the frequency method checking with different generated values of hyper-parameters, for example,  $r = 100$  and  $\beta = (2, 5)^\top$  when one covariate was used with an intercept term, then users can specify them via the arguments `A.or.r` and `reg.coef`;

```
R> cov <- coverage(b.output, A.or.r = 100, reg.coef = c(2, 5), nsim = 1000)
```

When users fit a model with a known value of the expected random effect, `coverage` conducts the frequency method checking based on the known value. However, users might want to conduct the frequency method checking with different values of the expected random effects. For example, if users want to try a different value of the expected random effect, e.g., 30 (or can be a vector of different values), we add the argument `mean.PriorDist` as follows.

```
R> cov <- coverage(p.output, mean.PriorDist = 30, nsim = 1000)
```

The resulting frequency method checking is based on the estimated posterior mode of  $r$  (because it is not specified through `A.or.r`) and the newly specified value of the expected random effect, 30.

Though the function `coverage` does not produce an S3 object, the result of `coverage` contains various numerical details;  $k$  RB coverage estimates (`cov$coverageRB`) and their standard errors (`cov$se.coverageRB`), RB overall coverage estimate (`cov$overall.coverageRB`) and its standard error (`cov$se.overall.coverageRB`), etc.

A coverage plot summarizing the result of `coverage` automatically appears. If the result is saved in a variable such as `cov` above, then users can recall the coverage plot, using the function `coverage.plot`.

```
R> coverage.plot(cov)
```

## 8. Examples

Applications to three data sets: Medical profiling of 31 hospitals with Poisson distributed fatality counts; Educational assessment of 8 schools with Normally distributed data; and evaluation of 18 baseball hitters with Binomial success rates and one covariate. For each example, we use 95% confidence level.

### 8.1. Poisson data with 31 hospitals: Known expected random effect

We analyze a data set of 31 hospitals in New York state comprising of the outcomes of the coronary artery bypass graft (CABG) surgery ([Morris and Lysy 2012](#)). The data set contains the number of deaths,  $\mathbf{y}$ , for a specified period after CABG surgeries out of the total number of patients,  $\mathbf{n}$ , receiving CABG surgeries in each hospital. A goal would be to obtain the point and interval estimates for the unknown true fatality rates (random effects) of 31 hospitals to evaluate each hospital's reliability on the CABG surgery ([Morris and Christiansen \(1995\)](#) use a similar Poisson model to handle these hospital profile data). We interpret the caseloads,  $\mathbf{n}$ , as exposures and assume that the state-level fatality rate per exposure of this surgery is known,  $\lambda_j^E = 0.03$  ( $m = 0$ ).

The following code can be used to load these data into R.

```
R> library("Rgbp")
R> data("hospital")
R> y <- hospital$d
R> n <- hospital$n
```

The function `gbp` can then be used to fit a Poisson-Gamma to the fatality rates in New York states with the expected random effect,  $\lambda_j^E$ , equal to 0.03.

```
R> p.output <- gbp(z, n, mean.PriorDist = 0.03, model = "poisson")
R> p.output
```

Summary for each unit (sorted by n):

	obs.mean	n	prior.mean	shrinkage	low.intv	post.mean	upp.intv	post.sd
1	0.0448	67	0.03	0.911	0.0199	0.0313	0.0454	0.00653
2	0.0294	68	0.03	0.910	0.0189	0.0299	0.0435	0.00631
3	0.0238	210	0.03	0.765	0.0185	0.0285	0.0407	0.00566
4	0.0430	256	0.03	0.728	0.0225	0.0335	0.0467	0.00619
5	0.0335	269	0.03	0.718	0.0208	0.0310	0.0432	0.00573
6	0.0438	274	0.03	0.714	0.0229	0.0339	0.0472	0.00621
7	0.0432	278	0.03	0.711	0.0228	0.0338	0.0469	0.00617
8	0.0136	295	0.03	0.699	0.0157	0.0250	0.0366	0.00534
9	0.0288	347	0.03	0.663	0.0200	0.0296	0.0410	0.00536
10	0.0372	349	0.03	0.662	0.0222	0.0325	0.0446	0.00571
11	0.0391	358	0.03	0.656	0.0228	0.0331	0.0454	0.00579
12	0.0177	396	0.03	0.633	0.0165	0.0255	0.0363	0.00506
13	0.0278	431	0.03	0.613	0.0200	0.0292	0.0400	0.00511

14	0.0249	441	0.03	0.608	0.0191	0.0280	0.0387	0.00502
15	0.0273	477	0.03	0.589	0.0199	0.0289	0.0394	0.00499
16	0.0455	484	0.03	0.585	0.0256	0.0364	0.0491	0.00601
17	0.0304	494	0.03	0.580	0.0211	0.0302	0.0409	0.00506
18	0.0220	501	0.03	0.577	0.0180	0.0266	0.0369	0.00483
19	0.0277	505	0.03	0.575	0.0202	0.0290	0.0395	0.00494
20	0.0204	540	0.03	0.559	0.0173	0.0258	0.0358	0.00474
21	0.0284	563	0.03	0.548	0.0206	0.0293	0.0395	0.00485
22	0.0236	593	0.03	0.535	0.0187	0.0270	0.0369	0.00466
23	0.0150	602	0.03	0.532	0.0147	0.0230	0.0329	0.00466
24	0.0238	629	0.03	0.521	0.0188	0.0271	0.0368	0.00460
25	0.0204	636	0.03	0.518	0.0173	0.0254	0.0351	0.00455
26	0.0480	729	0.03	0.484	0.0286	0.0393	0.0516	0.00587
27	0.0306	849	0.03	0.446	0.0223	0.0303	0.0397	0.00445
28	0.0274	914	0.03	0.428	0.0208	0.0285	0.0374	0.00423
29	0.0213	940	0.03	0.421	0.0176	0.0249	0.0335	0.00407
30	0.0293	1193	0.03	0.364	0.0223	0.0296	0.0379	0.00397
31	0.0201	1340	0.03	0.338	0.0170	0.0235	0.0310	0.00360
Mean		517	0.03	0.600	0.0201	0.0293	0.0403	0.00517

The output contains information about (from the left) the observed fatality rates  $\bar{y}_j$ , caseloads  $n_j$ , known expected random effect  $\lambda_j^E$ , shrinkage estimates  $\hat{B}'_j$ , lower bounds (2.5%) of posterior interval estimates  $\hat{\lambda}_{j,\text{low}}$ , posterior means  $\hat{E}(\lambda_j|\mathbf{y})$ , upper bounds (97.5%) of posterior interval estimates  $\hat{\lambda}_{j,\text{upp}}$ , and posterior standard deviations  $\widehat{\text{SD}}(\lambda_j|\mathbf{y})$  for random effects based on the assumed unconditional Gamma posterior distributions in (55).

A function `summary` shows selective information about hospitals with minimum, median, and maximum exposures and the estimation result of the hyper-parameter  $\alpha = -\log(r)$ .

```
R> summary(p.output)
```

Main summary:

	obs.mean	n	prior.mean	shrinkage	low.intv	post.mean
Unit with min(n)	0.0448	67	0.03	0.911	0.0199	0.0313
Unit with median(n)	0.0455	484	0.03	0.585	0.0256	0.0364
Unit with max(n)	0.0201	1340	0.03	0.338	0.0170	0.0235
Overall Mean		517	0.03	0.600	0.0201	0.0293
	upp.intv	post.sd				
	0.0454	0.00653				
	0.0491	0.00601				
	0.0310	0.00360				
	0.0403	0.00517				

Second-level Variance Component Estimation Summary:

alpha=log(A) for Gaussian or alpha=log(1/r) for Binomial and Poisson data:

```
post.mode.alpha post.sd.alpha post.mode.r
      -6.53          0.576          684
```

The output of `summary` shows that  $\hat{r} = \exp(6.53) = 684$ , which is an indicator of how valuable and informative the second-level hierarchy is. It means that the 25 hospitals with caseload less than 684 patients shrink their sample means towards the prior mean (0.03) more than 50%. For example, the shrinkage estimate of the first hospital ( $\hat{B}_1 = 0.911$ ) was calculated by  $684 / (684 + 67)$ , where 67 is its caseload ( $n_1$ ). As for this hospital, using more information from the conjugate prior distribution is an appropriate choice because the amount of observed information (67) is much less than the amount of state-level information (684).

To obtain a graphical summary, we use the function `plot`.

```
R> plot(p.output)
```

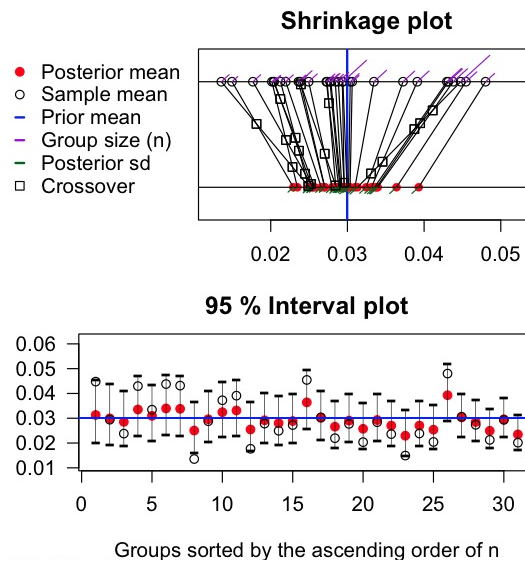


Figure 2: Shrinkage plot and 95% interval plot for fatality rates at 31 hospitals sorted by their caseloads.

In Figure 2 the regression towards the mean (RTTM) is obvious in the first plot; the observed fatality rates, denoted by empty dots on the upper horizontal line, are shrinking towards the known expected random effect, denoted by a blue vertical line at 0.03, to the different extents. Note that some hospitals' ranks have changed by shrinking much harder towards 0.03 than the others. For example, an empty square at the crossing point of the two left-most lines (8th and 23rd hospitals on the list above) indicates that the seemingly safest hospital in terms of the observed mortality rate is probably not the safest in terms of the estimated posterior mean accounting for the different caseloads of these two hospitals.

To be specific, their observed fatality rates ( $y_j$ ,  $j = 8, 23$ ) are 0.0136 and 0.0150 and caseloads ( $n_j$ ,  $j = 8, 23$ ) are 295 and 602, respectively. Considering solely the observed fatality rates may lead to an unfair comparison because the latter hospital handled twice the caseload.

**Rgbp** accounts for this caseload difference, making the posterior mean for the random effect of the former hospital shrink toward the state-level mean ( $\lambda_j^E=0.03$ ) much harder than that for the latter hospital.

Note that the point estimates are not enough to evaluate hospital reliability because one hospital may have a lower point estimate but larger uncertainty (variance) than the other. The second plot of Figure 2 displays the 95% interval estimates. Each posterior mean (red dot) is between the sample mean (empty dot) and the known expected random effect (a blue horizontal line).

This 95% interval plot reveals that the 31st hospital has the lowest upper bound even though its point estimate ( $\hat{\lambda}_{31} = 0.0235$ ) is slightly larger than that of the 23rd hospital ( $\hat{\lambda}_{23} = 0.0230$ ). The observed mortality rates for these two hospitals ( $y_j, j = 23, 31$ ) are 0.0150 and 0.0201 and the caseloads ( $n_j, j = 23, 31$ ) are 602 and 1340 each. The 31st hospital has twice the caseload, which leads to borrowing less information from the New York state-level hierarchy (or shrinking less toward the state-level mean, 0.03) with smaller variance. Based on the point and interval estimates, the 31st hospital seems the most reliable one.

Next, we do frequency method checking to question how reliable the estimation procedure is, assuming  $r$  equals its estimated value,  $\hat{r} = 683.53$ . The function `coverage` generates pseudo-datasets starting with the estimated value of  $r$  as a generative value. For reference, we could designate other generative values of  $r$  and  $\lambda_j^E$  by adding two arguments, `A.or.r` and `mean.PriorDist`, into the code below.

```
R> p.coverage <- coverage(p.output, nsim = 1000)
```

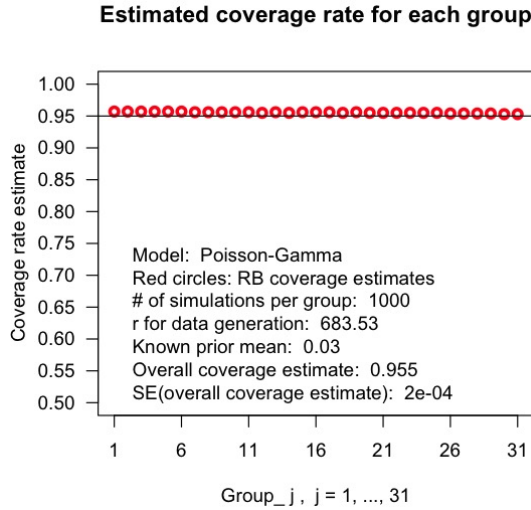


Figure 3: Coverage plot via frequency method checking for 31 hospitals.

In Figure 3, the black horizontal line at 0.95 represents the nominal confidence level and the red circles indicate RB unbiased coverage estimates,  $\bar{I}_{r, \lambda_j^E}^{RB}(\lambda_j)$  for  $j = 1, 2, \dots, 31$ . The RB overall unbiased coverage estimate across all the hospitals ( $\bar{\bar{I}}_{r, \lambda_j^E}^{RB}$ ) is 0.955. None of RB unbiased coverage estimates for the 31 hospitals are less than 0.95 regardless of their caseloads,



which range from 67 for hospital 1 to 1,340 for hospital 31. This result shows that the interval estimates for this particular dataset adequately achieves a 95% confidence level if  $r = \hat{r}$ .

The following code provides 31 RB unbiased coverage estimates and their standard errors (the output is omitted for space reasons).

```
R> p.coverage$coverageRB
R> p.coverage$se.coverageRB
```

The code below produces 31 simple unbiased coverage estimates and their standard errors.

```
R> p.coverage$coverageS
R> p.coverage$se.coverageS
```

It turns out that the variance estimate of the RB unbiased coverage estimate for the first hospital ( $0.0016^2$ ) is about 19 times smaller than that of the simple one ( $0.0070^2$ ). It means that the RB unbiased coverage estimates based on 1,000 simulations ( $N_{\text{sim}}$ ) are as precise as the simple unbiased coverage estimates based on 19,000 simulations in terms of estimating the coverage probability for the first hospital,  $C_{r,\lambda^E}(\lambda_1)$ .

## 8.2. Gaussian data with 8 schools: Unknown expected random effect and no covariates

The Education Testing Service (ETS) conducted randomized experiments in eight separate schools (groups) to test whether students (units) SAT scores are effected by coaching. The dataset contains the estimated coaching effects on SAT scores ( $y_j, j = 1, \dots, 8$ ) and standard errors ( $se_j, j = 1, \dots, 8$ ) of the eight schools (Rubin 1981). These data are contained in the package and can be loaded into R

```
R> library("Rgbbp")
R> data("schools")
R> y <- schools$y
R> se <- schools$se
```

Due to the nature of the test each school's coaching effect has an approximately Normal sampling distribution with approximately known sampling variances, based on large sample consideration. At the second hierarchy, the mean for each school is assumed to be drawn from a common Normal distribution ( $m = 1$ ).

```
R> g.output <- gbp(y, se, model = "gaussian")
R> g.output
```

Summary for each group (sorted by the descending order of se):

	obs.mean	se	prior.mean	shrinkage	low.intv	post.mean	upp.intv	post.sd
8	12.00	18.0	8.17	0.734	-10.21	9.19	29.9	10.23
3	-3.00	16.0	8.17	0.685	-17.13	4.65	22.5	10.10
1	28.00	15.0	8.17	0.657	-2.32	14.98	38.8	10.56

4	7.00	11.0	8.17	0.507	-8.78	7.59	23.6	8.26
6	1.00	11.0	8.17	0.507	-13.03	4.63	20.1	8.44
2	8.00	10.0	8.17	0.459	-7.25	8.08	23.4	7.81
7	18.00	10.0	8.17	0.459	-1.29	13.48	30.8	8.18
5	-1.00	9.0	8.17	0.408	-13.30	2.74	16.7	7.63
Mean		12.5	8.17	0.552	-9.16	8.17	25.7	8.90

This output from **gbp** summarizes the results. In this Gaussian model the amount of shrinkage for each unit is governed by the shrinkage factor,  $B_j = V_j/(V_j + A)$ . As such, schools whose variation within the school ( $V_j$ ) is less than the between school variation ( $A$ ) will shrink greater than 50%. The results provided by **gbp** suggests that there is little evidence that the training provided much added benefit due to the fact that every school's 95% posterior interval contains 0. In the case where the number of groups is large **Rgbp** provides a summary feature:

```
R> summary(g.output)
```

Main summary:

	obs.mean	se	prior.mean	shrinkage	low.intv	post.mean
Unit with min(se)	-1.00	9.0	8.17	0.408	-13.30	2.74
Unit with median(se)1	1.00	11.0	8.17	0.507	-13.03	4.63
Unit with median(se)2	7.00	11.0	8.17	0.507	-8.78	7.59
Unit with max(se)	12.00	18.0	8.17	0.734	-10.21	9.19
Overall Mean		12.5	8.17	0.552	-9.16	8.17

	upp.intv	post.sd
	16.7	7.63
	20.1	8.44
	23.6	8.26
	29.9	10.23
	25.7	8.90

Second-level Variance Component Estimation Summary:

alpha=log(A) for Gaussian or alpha=log(1/r) for Binomial and Poisson data:

post.mode.alpha	post.sd.alpha	post.mode.A
4.77	1.14	118

Regression Summary:

	estimate	se	z.val	p.val
beta1	8.168	5.73	1.425	0.154

The summary provides results regarding the second level hierarchy parameters. It can be seen that the estimate of the expected random effect,  $\mu^E = \beta_1$  (**beta1**), is not significantly

different from 0 suggesting that there was no effect of the coaching program on SAT math scores.

**Rgbp** also provides functionality to plot the results of the analysis as seen in Figure 4. Plotting the results provides a visual aid to understanding but is only largely beneficial when the number of groups ( $k$ ) is small.

```
R> plot(g.output)
```

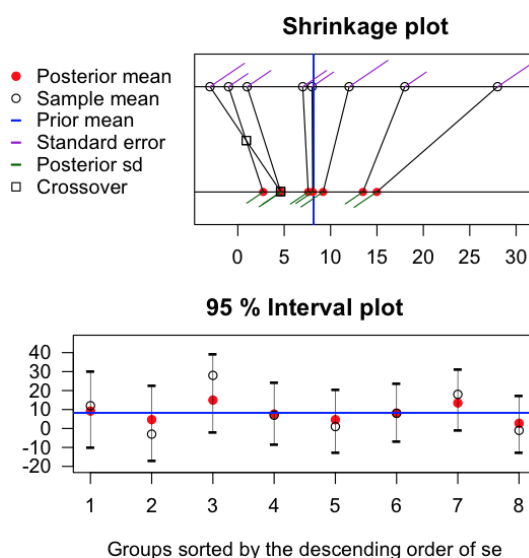


Figure 4: Shrinkage plot and 95% interval plot for 8 schools.

The frequency method checking generates new pseudo-data from our assumed model. Unless otherwise specified, the procedure fixes the hyper-parameter values at their estimates ( $\hat{A}$  and  $\hat{\beta}_1$  in this example) and then simulates random effects  $\theta_j$  for each group  $j$ . The model is then estimated and this is repeated an  $N_{\text{sim}}$  (`nsim`) number of times to estimate the coverage probabilities of the procedure.

```
R> g.coverage <- coverage(g.output, nsim = 1000)
```

As seen in Figure 5 the desired 95% confidence level, denoted by a black horizontal line at 0.95, is achieved for each school in this example. Note that all the coverage estimates depend on the chosen generative values of  $A$  and  $\beta_1$ , and the assumption that the model is valid.

In addition, RB unbiased coverage estimate and its standard error for each school can be gotten with the command below.

```
R> g.coverage$coverageRB
```

```
[1] 0.966 0.959 0.967 0.960 0.959 0.962 0.960 0.966
```

```
R> g.coverage$se.coverageRB
```

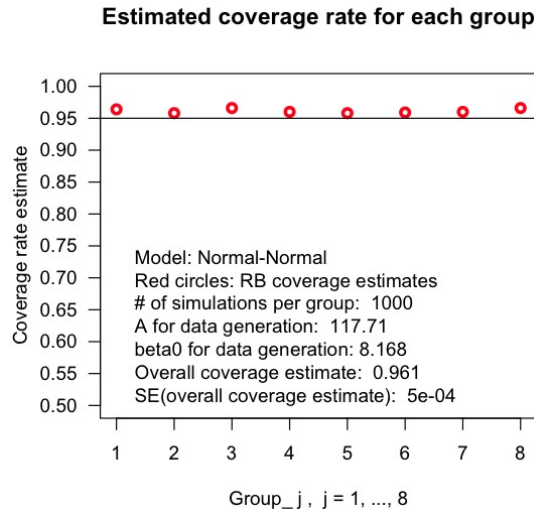


Figure 5: Coverage plot via frequency method checking for 8 schools.

```
[1] 0.0013 0.0012 0.0013 0.0013 0.0011 0.0011 0.0010 0.0017
```

### 8.3. Binomial data with 18 baseball players: Unknown expected random effect and one covariate

The data of 18 major league baseball players contain the batting averages through their first 45 official at-bats of the 1970 season (Efron and Morris 1975). A binary covariate is created that takes on the value one if a player is an outfielder and zero otherwise. The data can be loaded into R with the following code

```
R> library("Rgbp")
R> data("baseball")
R> y <- baseball$Hits
R> n <- baseball$At.Bats
R> x <- ifelse(baseball$Position == "fielder", 1, 0)
```

Conditional on the unknown true batting average (random effect) of each player it is assumed that the at-bats are independent and therefore,  $y_j|p_j \stackrel{\text{indep.}}{\sim} \text{Binomial}(45, p_j)$ ,  $j = 1, \dots, 18$ . Our goal is to obtain point and interval estimates of each random effect whilst considering the additional information on whether the player is an outfielder or not. The function `gbp` provides a way to incorporate such covariate information seamlessly into the model so that the regression towards the mean (RTTM) occurs within outfielders and non-outfielders separately.

```
R> b.output <- gbp(z, n, x, model = "binomial")
R> b.output
```

Summary for each unit (sorted by n):

	obs.mean	n	X1	prior.mean	shrinkage	low.intv	post.mean	upp.intv	post.sd
1	0.400	45	1.0	0.310	0.715	0.248	0.335	0.429	0.0462
2	0.378	45	1.0	0.310	0.715	0.244	0.329	0.420	0.0448
3	0.356	45	1.0	0.310	0.715	0.240	0.323	0.411	0.0437
4	0.333	45	1.0	0.310	0.715	0.236	0.316	0.403	0.0429
5	0.311	45	1.0	0.310	0.715	0.230	0.310	0.396	0.0424
6	0.311	45	0.0	0.233	0.715	0.179	0.256	0.341	0.0415
7	0.289	45	0.0	0.233	0.715	0.175	0.249	0.331	0.0400
8	0.267	45	0.0	0.233	0.715	0.171	0.243	0.323	0.0388
9	0.244	45	0.0	0.233	0.715	0.166	0.237	0.315	0.0380
10	0.244	45	1.0	0.310	0.715	0.210	0.291	0.379	0.0432
11	0.222	45	0.0	0.233	0.715	0.161	0.230	0.308	0.0377
12	0.222	45	0.0	0.233	0.715	0.161	0.230	0.308	0.0377
13	0.222	45	0.0	0.233	0.715	0.161	0.230	0.308	0.0377
14	0.222	45	1.0	0.310	0.715	0.202	0.285	0.375	0.0441
15	0.222	45	1.0	0.310	0.715	0.202	0.285	0.375	0.0441
16	0.200	45	0.0	0.233	0.715	0.155	0.224	0.302	0.0377
17	0.178	45	0.0	0.233	0.715	0.148	0.218	0.297	0.0381
18	0.156	45	0.0	0.233	0.715	0.140	0.211	0.292	0.0389
Mean		45	0.4	0.267	0.715	0.191	0.267	0.351	0.0410

Note that the shrinkage estimates are the same for all players because all players have the same 45 at-bats.

```
R> summary(b.output)
```

Main summary:

	obs.mean	n	X1	prior.mean	shrinkage	low.intv
Unit with min(obs.mean)	0.156	45	0.000	0.233	0.715	0.140
Unit with median(obs.mean)1	0.244	45	0.000	0.233	0.715	0.166
Unit with median(obs.mean)2	0.244	45	1.000	0.310	0.715	0.210
Unit with max(obs.mean)	0.400	45	1.000	0.310	0.715	0.248
Overall Mean		45	0.444	0.267	0.715	0.191

post.mean	upp.intv	post.sd
0.211	0.292	0.0389
0.237	0.315	0.0380
0.291	0.379	0.0432
0.335	0.429	0.0462
0.267	0.351	0.0410

Second-level Variance Component Estimation Summary:

alpha=log(A) for Gaussian or alpha=log(1/r) for Binomial and Poisson data:

```
post.mode.alpha post.sd.alpha post.mode.r
```

-4.73                      0.957                      113

Regression Summary:

	estimate	se	z.val	p.val
beta1	-1.194	0.131	-9.129	0.000
beta2	0.389	0.187	2.074	0.038

The regression coefficient for the outfielder indicator is significant, considering that  $p$  value for  $\hat{\beta}_2$  is 0.038. It means that the two estimates for the expected random effects for the outfielders and infielders are significantly different. Also, the positive sign of  $\hat{\beta}_2$  indicates that the population batting average for outfielders tends to be higher than that for infielders. The estimated odds ratio is  $\exp(0.389) = 1.48$ .

R> `plot(b.output)`

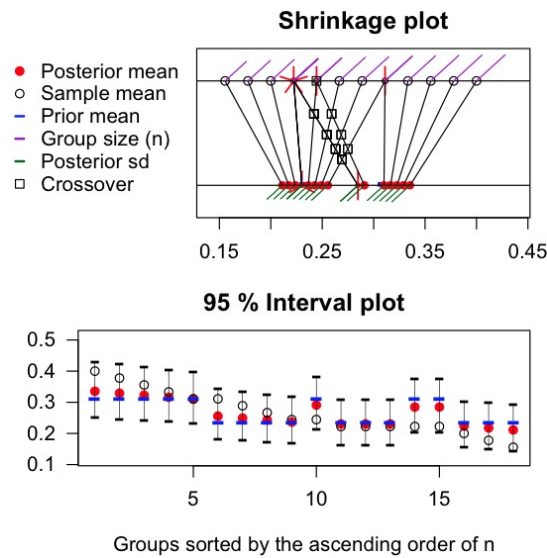


Figure 6: Shrinkage plot and 95% interval plot for 18 baseball players.

The shrinkage plot in Figure 6 shows that the observed batting averages (empty dots) on the upper horizontal line shrink towards the two expected random effects, 0.233 and 0.310. The short red line symbols near some empty dots are for when two or more points have the same mean and are plotted over each other. For example, five players (from the 11th player to the 15th) have the same batting average, 0.222, and at this point on the upper horizontal line, there are short red lines toward five directions.

The 95% interval plot in Figure 6 shows the range of true batting average for each player, which clarifies the regression toward the mean (RTTM) within two groups. The 10th, 14th, and 15th players, for example, are outfielders but their observed batting averages are far lower than the first five outfielders. This can be attributed to their bad luck because their observed

batting averages are close to the lower bounds of their interval estimates. The RTTM indicates that their batting averages shrink towards the expected random effect of outfielders (0.310) in the long run.

To check the level of trust in these interval estimates, we proceed to frequency method checking by assuming the estimates, 112.95 for  $\hat{r}$  and (-1.194, 0.389) for  $\hat{\beta}$ , are the generative values.

```
R> b.coverage <- coverage(b.output, nsim = 1000)
```

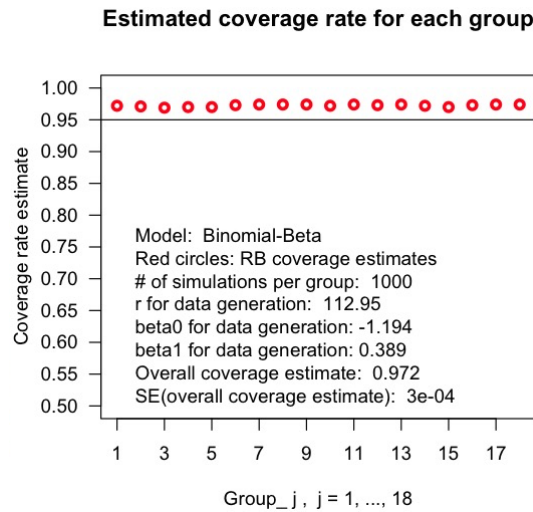


Figure 7: Coverage plot via frequency method checking for 18 players.

In Figure 7, the estimated coverage probabilities for random effects are beyond 0.95, conservatively satisfying the 95% confidence level if  $r = \hat{r}$  and  $\beta = \hat{\beta}$ . The RB overall unbiased coverage estimate across all the players is 0.972.

We can check the RB unbiased coverage estimates and their standard errors for each player.

```
R> bcv$coverageRB
```

```
[1] 0.971 0.973 0.972 0.972 0.970 0.973 0.973 0.974 0.973 0.973 0.971 0.973
[13] 0.973 0.972 0.972 0.971 0.973 0.971
```

```
R> bcv$se.coverageRB
```

```
[1] 0.0015 0.0012 0.0013 0.0014 0.0016 0.0010 0.0012 0.0010 0.0010 0.0013
[11] 0.0015 0.0013 0.0019 0.0013 0.0014 0.0015 0.0011 0.0014
```

If we want to draw 2,000 posterior samples of random effects and hyper-parameters from their full posterior distribution via the A-R method, we use the following R code.

```
R> b.output <- gbp(y, n, x, model = "binomial", n.AR = 2000)
```



The “gbp” object `b.output` contains 8,000 weights (`b.output$weight`), 2,000 posterior samples of  $\alpha$  (`b.output$alpha`), a  $2,000 \times 2$  matrix of  $\beta$  (`b.output$beta`) each column of which corresponds to 2,000 posterior samples of each regression coefficient, and a  $k \times 2,000$  matrix of random effects (`b.output$p`) each row of which has posterior samples of each random effect. If we run the frequency method checking using this “gbp” object obtained via the A-R method, the  $N_{\text{sim}}$  simulations also run the A-R method each time.

## 9. Discussion

**Rgbp** is an R package for estimating and validating two-level Gaussian, Poisson, and Binomial hierarchical models. The package aims to provide a procedure that is computationally efficient with good frequency properties and includes “frequency method checking” functionality to examine repeated sampling properties and to test that the method is valid at specified hyperparameter values.

As an alternative to other maximization based estimation methods such as MLE and REML, **Rgbp** provides approximate point and interval estimates of parameters via ADM. Using the ADM approach, with our specified choice of priors, protects from cases of overshrinkage and undercoverage from which the aforementioned methods suffer (Morris 1988b).

A benefit of **Rgbp** is that it produces non-random output (except the A-R method for the Binomial model) and so results are easily reproduced and compared across studies. In addition to being a standalone analysis tool the package can be used as an aid in a broader estimation procedure. For example, by checking the similarity of output of **Rgbp** and that of another estimation procedure such as MCMC (Markov Chain Monte Carlo), the package can be used as a confirmatory tool to check whether the alternative procedure has been programmed correctly. In addition, the parameter estimates obtained via **Rgbp** can be used to initialize a MCMC thus decreasing time to convergence. Lastly, due to its speed and ease of use, **Rgbp** can be used as a method of preliminary data analysis. Such results may tell statisticians and practitioners alike whether a more intensive method in terms of implementation and computational time, such as MCMC, is needed.

### A. Posterior propriety of the Poisson model

If the posterior distribution of  $r$  is proper, then the full posterior distribution of random effects and  $r$  is also proper because

$$f(\boldsymbol{\lambda}, r | \mathbf{y}) = f(\boldsymbol{\lambda} | \mathbf{y}) \cdot f(r | \mathbf{y}), \quad (71)$$

where  $f(\boldsymbol{\lambda} | \mathbf{y})$  is a product of  $k$  proper conditional posterior density function in (41). Thus, our goal is to show that  $\int_0^\infty f(r | \mathbf{y}) dr < \infty$ ;

$$f(r | \mathbf{y}) \propto \frac{1}{r^2} L(r) \propto \frac{1}{r^2} \prod_{j=1}^k \frac{\Gamma(r\lambda_j^E + y_j)}{\Gamma(r\lambda_j^E)} (1 - B_j)^{y_j} B_j^{r\lambda_j^E} \quad (72)$$

$$= \frac{1}{r^2} \left[ r^{\sum_{j=1}^k y_j} + \dots + a_k r^k \right] \exp \left( -r \sum_{j=1}^k \lambda_j^E \log(1 + n_j/r) \right) \prod_{j=1}^k \left( \frac{n_j}{n_j + r} \right)^{y_j}, \quad (73)$$

where the polynomial function of  $r$  in the bracket has constant coefficients.

If there are at least two groups whose observed values  $y_j$  are non-zero, then  $f(r|\mathbf{y})$  goes to zero as  $r$  goes to 0 due to the polynomial function of  $r$  in (73); the following two factors in (73) approach one. As  $r$  becomes infinite,  $f(r|\mathbf{y})$  touches zero exponentially fast due to the exponential term in the middle of (73). Thus, the integration of  $f(r|\mathbf{y})$  must be finite.

## Acknowledgments

The authors thank Professor Cindy Christiansen, Professor Phil Everson and the 2012 class of Harvard's Stat 324r: Parametric Statistical Inference and Modeling for their valuable inputs.

## References

- Albert JH (1988). "Computational methods using a Bayesian hierarchical generalized linear model." *Journal of the American Statistical Association*, **83**(404), 1037–1044.
- Azzalini A (1985). "A Class of Distributions which Includes the Normal Ones." *Scandinavian journal of statistics*, pp. 171–178.
- Blackwell D (1947). "Conditional Expectation and Unbiased Sequential Estimation." *The Annals of Mathematical Statistics*, pp. 105–110.
- Brillinger D (1969). "The calculation of cumulants via conditioning." *Annals of the Institute of Statistical Mathematics*, **21**(1), 215–218. ISSN 0020-3157. doi:10.1007/BF02532246. URL <http://dx.doi.org/10.1007/BF02532246>.
- Christiansen C, Morris C (1996). "Fitting and Checking a Two-Level Poisson Model: Modeling Patient Mortality Rates in Heart Transplant Patients." In D Berry, D Stangl (eds.), *Bayesian Biostatistics*, pp. 467–501. CRC press.
- Christiansen C, Morris C (1997). "Hierarchical Poisson Regression Modeling." *Journal of the American Statistical Association*, **92**(438), pp. 618–632. ISSN 01621459. URL <http://www.jstor.org/stable/2965709>.
- Daniels MJ (1999). "A Prior for the Variance in Hierarchical Models." *Canadian Journal of Statistics*, **27**(3), 567–578.
- Dean CB (1992). "Testing for Overdispersion in Poisson and Binomial Regression Models." *Journal of the American Statistical Association*, **87**(418), 451–457.
- Efron B, Morris C (1975). "Data Analysis Using Stein's Estimator and its Generalizations." *Journal of the American Statistical Association*, **70**(350), pp. 311–319. ISSN 01621459. URL <http://www.jstor.org/stable/2285814>.
- Everson PJ, Morris CN (2000). "Inference for Multivariate Normal Hierarchical Models." *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **62**(2), 399–412.

- Gelman A, Carlin JB, Stern HS, Rubin DB (2014). *Bayesian data analysis*, volume 2. Taylor & Francis.
- Ghosh M, Rao J (1994). “Small area estimation: an appraisal.” *Statistical science*, pp. 55–76.
- Jones M, Faddy M (2003). “A Skew Extension of the t-Distribution, with Applications.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **65**(1), 159–174.
- Kass RE, Steffey D (1989). “Approximate Bayesian Inference in Conditionally Independent Hierarchical Models (Parametric Empirical Bayes Models).” *Journal of the American Statistical Association*, **84**(407), pp. 717–726. ISSN 01621459. URL <http://www.jstor.org/stable/2289653>.
- Kelly J (2014). *Advances in the Normal-Normal Hierarchical Model*. Ph.D. thesis, Harvard University.
- Lee Y, Nelder JA (1996). “Hierarchical Generalized Linear Models.” *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 619–678.
- Lee Y, Nelder JA, Pawitan Y (2006). *Generalized Linear Models with Random Effects: A Unified Analysis via h-Likelihood*. Chapman & Hall/ CRC, New York.
- Morris C (1988a). “Approximating Posterior Distributions and Posterior Moments.” In J Bernardo, MH DeGroot, DV Lindley, AFM Smith (eds.), *Bayesian Statistics 3*, pp. 327–344. Oxford University Press.
- Morris C (1988b). “Determining the Accuracy of Bayesian Empirical Bayes Estimates in the Familiar Exponential Families.” In S Gupta, J Berger (eds.), *Statistical Decision Theory and Related Topics IV*, pp. 251–263. Springer-Verlag.
- Morris C, Christiansen C (1995). “Hierarchical Models for Ranking and for Identifying Extremes, with Application.” In J Bernardo, J Berger, A Dawid, A Smith (eds.), *Bayesian Statistics 5*, pp. 227–296. New York: Oxford University Press.
- Morris C, Lysy M (2012). “Shrinkage Estimation in Multilevel Normal Models.” *Statistical Science*, **27**(1), 115–134.
- Morris C, Tang R (2011). “Estimating Random Effects via Adjustment for Density Maximization.” *Statistical Science*, **26**(2), pp. 271–287. ISSN 08834237. URL <http://www.jstor.org/stable/23059992>.
- Morris CN (1983). “Natural exponential families with quadratic variance functions: statistical theory.” *The Annals of Statistics*, pp. 515–529.
- Rao CR (1945). “Information and Accuracy Attainable in the Estimation of Statistical Parameters.” *Bulletin of the Calcutta Mathematical Society*, **37**(3), 81–91.
- Rao JN (2003). *Small area estimation*. Wiley Online Library.
- Robert C, Casella G (2013). *Monte Carlo statistical methods*. Springer Science & Business Media.

- Rönnegård L, Shen X, Alam M (2010). “**hglm**: A Package for Fitting Hierarchical Generalized Linear Models.” *The R Journal*, **2**(2), 20–28. ISSN 20734859.
- Rönnegård L, Shen X, Alam M (2011). *The **hglm** Package*. URL <http://CRAN.R-project.org/package=hglm>.
- Rubin DB (1981). “Estimation in Parallel Randomized Experiments.” *Journal of Educational Statistics*, **6**(4), pp. 377–401. ISSN 03629791. URL <http://www.jstor.org/stable/1164617>.
- Skellam J (1948). “A Probability Distribution Derived from the Binomial Distribution by Regarding the Probability of Success as Variable between the Sets of Trials.” *Journal of the Royal Statistical Society. Series B (Methodological)*, **10**(2), 257–261.
- Tak H, Morris C (2015). “Posterior Propriety and Frequency Coverage Evaluation of Bayesian Beta-Binomial Logistic Regression Model.” *in preparation*.
- Tamura RN, Young SS (1987). “A stabilized moment estimator for the beta-binomial distribution.” *Biometrics*, pp. 813–824.
- Tang R (2002). *Fitting and Evaluating Certain Two-Level Hierarchical Models*. Ph.D. thesis, Harvard University.

**Affiliation:**

Hyungsuk Tak  
Department of Statistics  
Harvard University  
1 Oxford Street, Cambridge, MA  
E-mail: [hyungsuk.tak@gmail.com](mailto:hyungsuk.tak@gmail.com)

Joseph Kelly  
Google  
76 Ninth Avenue, New York, NY  
E-mail: [josephkelly@google.com](mailto:josephkelly@google.com)

Carl Morris  
Department of Statistics  
Harvard University  
1 Oxford Street, Cambridge, MA  
E-mail: [morris@fas.harvard.edu](mailto:morris@fas.harvard.edu)