

## **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 averages or proportions, possibly together with covariates of each group but without population-level data. The approximate Bayesian tool equipped with the adjustment for density maximization produces point and interval estimates for each random effect, point estimates and their standard errors for regression coefficients, and a point estimate and its standard error for a second-level variance component. For the Binomial data, the package provides an option to draw independent posterior samples of all the model parameters from their posterior distributions via the acceptance-rejection method. The main goal of **Rgb** is to produce approximate, or exact for the Binomial model, Bayesian interval estimates for the random effects that meet their nominal confidence levels. For this purpose, our models adopt unique improper hyper-prior distributions and **Rgbp** provides a quick way to check whether the resultant 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**

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 the extra-variability, called overdispersion, a two-level con-

jugate hierarchical model regards first-level mean parameters as random effects that come from a population-level conjugate prior distribution. The conjugate prior distribution is non-exchangeable if the model incorporates covariate information of each group via a linear, log-linear, or logistic regression according to the data type, and exchangeable if no covariates are available with only an intercept term for the regression.

With an assumption of homogeneity within each group, the observed data are sufficient statistics for the random effects, such as averages or proportions, possibly together with each group's covariate information. This type of data is common for a biological analysis on litter data, a meta analysis on independent studies, or small area estimation problems. For these data, the two-level model that **Rgbp** assumes can be considered as a conjugate hierarchical generalized linear model (Lee and Nelder 1996; Lee, Nelder, and Pawitan 2006) where each random effect has a conjugate prior distribution.

**Rgbp** takes a Bayesian approach with our special improper hyper-prior distributions on hyper-parameters, the parameters of the conjugate prior distribution, to produce Bayesian interval estimates for the random effects that achieve nominal confidence levels. The hyper-prior distributions lead to Stein's harmonic prior known to produce good repeated sampling coverage rates of the Bayesian interval estimates for random effects in a two-level Gaussian hierarchical model (Morris and Tang 2011; Morris and Lysy 2012; Kelly 2014). We apply an analog to Stein's harmonic prior to Poisson and Binomial hierarchical models.

When it comes to fitting the model, **Rgbp** adopts the adjustment for density maximization (Morris 1988a; Christiansen and Morris 1997; Morris and Tang 2011) (ADM), a Pearson family approximation via maximization. For example, a Delta method is a Normal case of the ADM that uses a Normal distribution to obtain an approximate distribution of a function of a parameter with its MLE and observed information plugged-in. In this article, the ADM uses Beta distributions to approximate the posterior distributions of shrinkage factors, functions of the second-level variance component, via maximization. Using the approximate Beta posterior distributions of shrinkage factors, **Rgbp** estimates the first two (three for the Gaussian case) posterior moments of the random effects. Finally, **Rgbp** approximates the posterior distribution of each random effect by a skewed Normal distribution for a Gaussian case, a Gamma distribution for a Poisson case, and a Beta distribution for a Binomial case whose two parameters (three for the skewed Normal distribution) are matched to the previously estimated posterior moments of the random effects.

For the Binomial hierarchical model, **Rgbp** provides an option to draw independent posterior samples of all the model parameters from their posterior distributions via an acceptance-rejection method instead of the approximate Bayesian tool.

In addition to fitting hierarchical models, **Rgbp** provides a quick way to evaluate 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). It is a unique procedure that distinguishes **Rgbp** from any other R packages for hierarchical modeling such as **hglm** (Rönnegård, Shen, and Alam 2010, 2011) for conjugate hierarchical generalized models and **arm** (Gelman, Su, Yajima, Hill, Pittau, Kerman, and Zheng 2014) for Bayesian hierarchical regression models. The evaluation procedure which we call "frequency method checking" adopts a parametric bootstrapping, a Monte Carlo simulation, that generates mock data sets given the values of the hyper-parameters and estimates the coverage rates based on the generated mock data sets.

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 the ADM and the acceptance-rejection method in Section 4 and 5, respectively. We introduce the frequency method checking in Section 6. We explain the usages of main functions in **Rgbp** in Section 7, and apply them to three examples in Section 8.

## 2. Conjugate hierarchical modeling structure

One of the functions in **Rgbp**, **gbp**, fits a conjugate hierarchical model in which the first-level has distributions of observed data and the second-level has conjugate distributions on the first-level mean parameters (random effects) without any observed data on this second-level. Specifically, with an assumption of the homogeneity within each group, the observed data of the first-level are sufficient statistics for random effects such as sample means or numbers of successful outcomes. These data appear often in a biological analysis on litter data, a meta analysis on independent studies, or small area estimation problems. The function **gbp** allows users to choose one of three hierarchical models according to the type of data, namely Normal-Normal, Poisson-Gamma, and Binomial-Beta models. There are possibly more hierarchical models depending on the types of data such as Exponential-Gamma (Tang 2002), but we choose the three models because these are the most commonly-used models.

In this section, we specify the modeling details. Our motivation for the parametrization is to provide an intuitive shrinkage interpretation in inference and to facilitate the estimation procedure, considering that the shrinkage factors under our parametrization are a function of only a second-level variance component.

### 2.1. Normal-Normal (“g”=Gaussian data)

The following is the general Normal-Normal hierarchical model (hereafter the Gaussian model) assumed by **Rgbp**. The subscript  $j$  below indicates the  $j$ -th group among  $k$  groups in the dataset. For  $j = 1, 2, \dots, k$ ,

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

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

where  $y_j$  is an observed unbiased estimate, sample mean, for random effect  $j$ ,  $V_j$  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 regression coefficients to be estimated. The default of the function **gbp** is to set  $x_{j,1}$  to 1 for an intercept term, though **gbp** 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 with an 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 distributions of the random effects. For  $j = 1, 2, \dots, k$ ,

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

where  $B_j \equiv V_j/(V_j + A)$ ,  $j = 1, \dots, k$ , are called shrinkage factors, and  $\mathbf{y}^T = (y_1, y_2, \dots, y_k)$ . Note that the conditional posterior mean of the random effect,  $\mu_j^* \equiv E(\mu_j | \beta, A, \mathbf{y}) = (1 - B_j)y_j + B_j\mu_j^E$ , 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 more accurate second-level conjugate prior distribution.

## 2.2. Poisson-Gamma (“p”=Poisson data)

The function **gbp** is also capable of estimating a conjugate Poisson-Gamma hierarchical model (hereafter the Poisson model), though its usage is limited to the case where the expected random effects are known ( $m = 0$ ). For instance, from previous studies we may be able to obtain some information about the mean of the population from which we collect the first-level group data. For  $j = 1, 2, \dots, k$ ,

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

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

where  $y_j$  is the number of events happening,  $n_j$  is the exposure of group  $j$ , which is not necessarily an integer,  $\lambda^E = E(\lambda_j | r)$  is the known expected random effect, and  $r$  is the unknown second-level variance component. The mean and variance of the conjugate Gamma prior distribution are  $\lambda^E$  and  $\lambda^E/r$ , respectively. We interpret  $r$  as the amount of prior information as  $n_j$  represents the amount of observed information, which makes intuitive sense because the uncertainty of the conjugate prior distribution increases as  $r$  decreases; in the limit of  $r$  going to 0, the conjugate prior distribution gets flatter. The conditional posterior distribution of the random effect  $\lambda_j$  for this Poisson-Gamma model is

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

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

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

where  $B_j \equiv r/(r + n_j)$  is the shrinkage factor, the relative amount of information in the prior compared to the data. The conditional posterior mean is a convex combination of the unbiased estimate  $\bar{y}_j = y_j/n_j$  and the known expected random effect  $\lambda^E$  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%, borrowing more information from the second-level distribution.

Note that the conditional posterior variance in Equation 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 their means. The parametrization of the Poisson model adopted in **Rgbp** makes shrinkage factors as a function of only the hyper-parameter  $r$ , while the shrinkage factors based on

the parametrization in [Christiansen and Morris \(1997\)](#) are functions of  $r$  and  $\beta$ . Thus, the estimation procedure and the interpretation of the shrinkage factors are more intuitive under the parametrization in **Rgbp**.

### 2.3. Binomial-Beta (“b”=Binomial data)

The Binomial-Beta hierarchical model is the last model that **gbp** can fit. The notation  $y_j$  below is the number of successes (or failures) out of  $n_j$  trials. Unlike the Poisson model, the expected random effect is either known ( $m = 0$ ) or unknown ( $m \geq 1$ ) a priori.

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

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

where  $p_j^E \equiv E(p_j | \beta, r) = \exp(\mathbf{x}_j^\top \beta) / (1 + \exp(\mathbf{x}_j^\top \beta))$  is the expected random effect of group  $j$  ( $j = 1, 2, \dots, k$ ). The  $m \times 1$  vector of the logistic regression coefficient  $\beta$  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  $j$  is

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

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

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

The conditional posterior mean is a convex combination of the unbiased estimate  $\bar{y}_j = y_j/n_j$  and the expected random effect  $p_j^E$  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 shrinks towards the expected random effect by more than 50%.

### 2.4. Hyper-prior Distribution

Hyper-prior distributions are the distributions assigned to the second-level parameters called hyper-parameters. With the goal of objectivity in mind, our choice for the hyper-prior distributions is

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

where  $m$  is the number of the regression coefficients to be estimated. The improper flat hyper-prior distribution on  $\beta$  is a common non-informative choice. In the Gaussian case, the flat hyper-prior distribution on the second-level variance  $A$  produces good repeated sampling coverage properties of the Bayesian interval estimates for the random effects. 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 ( $\pi(B_j) \propto B_j^{-2} dB_j$ ) as does  $A$  with the  $\text{Uniform}(0, \infty)$  for the Gaussian case. The resultant full posterior distribution of random effects and hyper-parameters for the Binomial case is data-dependent. 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 the entire  $k$  groups. Then, 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  $m$  (Tak and Morris in preparation).

The Poisson model with the hyper-prior distributions in Equation 12 provides posterior propriety if and only if there are at least two groups whose observed values  $y_j$  are non-zero and the expected random effect  $\lambda^E$  is a completely known constant ( $m = 0$ ). If the expected random effect is unknown a priori, then we recommend staying with the Binomial model with the same hyper-prior distributions because the Poisson model is actually an approximation to the Binomial model.

### 3. The inferential model

The likelihood function of hyper-parameters  $A$  and  $\beta$  for the Gaussian model is derived from the independent Normal marginal distributions of the observed data with random effects,  $\mu_1, \mu_2, \dots, \mu_k$ , integrated;

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

where  $\mu_j^E = \mathbf{x}^\top \beta$ . The joint posterior density function of hyper-parameters  $f(A, \beta | \mathbf{y})$  for the Gaussian hierarchical model is proportional to their likelihood function in Equation 13 because we use flat improper hyper-prior density functions for  $A$  and  $\beta$ ;

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

The likelihood function of hyper-parameters  $r$  and  $\beta$  for the Binomial model is derived from the independent Beta-Binomial marginal distributions of the observed data with random effects,  $p_1, p_2, \dots, p_k$ , integrated out (Skellam 1948);

$$L_b(r, \beta) = \prod_{j=1}^k f(y_j | r, \beta) = \prod_{j=1}^k \binom{n_j}{y_j} \frac{B(y_j + rp_j^E, n_j - y_j + r(1 - p_j^E))}{B(rp_j^E, r(1 - p_j^E))}, \quad (15)$$

where  $p_j^E = \exp(\mathbf{x}^\top \beta) / (1 + \exp(\mathbf{x}^\top \beta))$  and 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 Equation 15 multiplied by the hyper-prior density functions of  $r$  and  $\beta$  in Equation 12 as follows;

$$f(r, \beta | \mathbf{y}) \propto L_b(r, \beta) d\beta dr / r^2. \quad (16)$$

Similarly, the likelihood function of  $r$  for the Poisson hierarchical model comes from the independent Negative-Binomial marginal distributions of the observed data with the random

effects,  $\lambda_1, \lambda_2, \dots, \lambda_k$ , integrated out;

$$L_p(r) = \prod_{j=1}^k f(y_j|r) = \prod_{j=1}^k \frac{\Gamma(r\lambda^E + y_j)}{\Gamma(r\lambda^E)(y_j!)} (1 - B_j)^{y_j} B_j^{r\lambda^E}, \quad (17)$$

where  $\Gamma(a)$  is a gamma function defined as  $\int_0^\infty x^{a-1} e^{-x} dx$  for a positive constant  $a$ . The posterior density function of  $r$ ,  $f(r|\mathbf{y})$ , for the Poisson hierarchical model is the likelihood function in Equation 17 times the hyper-prior density function of  $r$ ,  $dr/r^2$ ;

$$f(r|\mathbf{y}) \propto L_p(r) dr/r^2. \quad (18)$$

Our goal is to obtain the point and interval estimates of the random effects from their unconditional posterior distributions; for the Gaussian model,

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

for the Binomial model,

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

and lastly for the Poisson model,

$$f(\boldsymbol{\lambda}|\mathbf{y}) = \int f(\boldsymbol{\lambda}|r, \mathbf{y}) \cdot f(r|\mathbf{y}) dr. \quad (21)$$

## 4. Estimation via the adjustment for density maximization

In this section, we illustrate our estimation procedure equipped with the adjustment for density maximization (hereafter ADM) (Morris 1988a; Christiansen and Morris 1997; Morris and Tang 2011), a way to approximate a distribution of the parameter of interest by one of Pearson family distributions based on derivatives like the Delta method.

Overall, obtaining point and interval estimates of each random effect is our primary inferential interest. The approximate Bayesian tool assumes that the unconditional posterior distribution of each random effect follows a skewed-Normal distribution for the Gaussian case, a Beta distribution for the Binomial case, or a Gamma distribution for the Poisson case, whose parameters are matched to the estimated unconditional posterior moments of each random effect. We use these assumed unconditional posterior distributions to make point and interval estimates of each random effect.

Because the ADM procedure for the Gaussian model and its frequency property are well documented in Kelly (2014), we describe the ADM procedure mainly for the Binomial model.

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

Estimating the unconditional posterior moments of the shrinkage factors,  $B_1, B_2, \dots, B_k \equiv r/(r + n_k)$  ( $B_j \equiv V_j/(V_j + A)$  for the Gaussian model) and the conditional posterior moments



of the expected random effects,  $p_1^E, p_2^E, \dots, p_k^E$  ( $\mu_1^E, \mu_2^E, \dots, \mu_k^E$  for the Gaussian model or  $\lambda^E$  for the Poisson model), is the main estimation problem for the models that **gbp** assumes. This is because these moment estimates are used to estimate the unconditional posterior moments of the random effects,  $p_1, p_2, \dots, p_k$  ( $\mu_1, \mu_2, \dots, \mu_k$  for the Gaussian model or  $\lambda_1, \lambda_2, \dots, \lambda_k$  for the Poisson model). Taking the Binomial model as an example, with the assumption that hyper-parameters  $r$  and  $\beta$  are independent a posteriori, the unconditional posterior mean and variance of random effect  $j$  are

$$E(p_j|\mathbf{y}) = E(E(p_j|r, \beta, \mathbf{y})|\mathbf{y}) = (1 - E(B_j|\mathbf{y}))\bar{y}_j + E(B_j|\mathbf{y})E(p_j^E|\mathbf{y}) \quad (22)$$

$$Var(p_j|\mathbf{y}) = E(Var(p_j|r, \beta, \mathbf{y})|\mathbf{y}) + Var(E(p_j|r, \beta, \mathbf{y})|\mathbf{y}) \quad (23)$$

$$= E(p_j^*(1 - p_j^*)/(r + n_j + 1)|\mathbf{y}) + Var(B_j(\bar{y}_j - p_j^E)|\mathbf{y}) \quad (24)$$

$$\approx E(p_j^*(1 - p_j^*)(1 - B_j)/n_i|\mathbf{y}) + Var(B_j(\bar{y}_j - p_j^E)|\mathbf{y}) \quad (25)$$

$$= h_b(E(B_j|\mathbf{y}), E(B_j^2|\mathbf{y}), E(B_j^3|\mathbf{y}), E(p_j^E|\mathbf{y}), E((p_j^E)^2|\mathbf{y})). \quad (26)$$

Note that the unconditional posterior mean and approximate variance of random effect  $j$  in Equation 22 and 26 are functions of the unconditional posterior moments of shrinkage factors and expected random effects. We specify the function  $h_b$  in Appendix A.

We assumed that hyper-parameters  $r$  ( $A$  for the Gaussian model) and  $\beta$  were independent a posteriori for the Binomial model, considering that they are independent a posteriori in the limit of  $k$  going to infinity because they are asymptotically Normally distributed. Also, Christiansen and Morris (1997) empirically showed that their covariance from the observed information matrix of the Poisson model, though with a different parametrization, was close to 0 in a small sample setting.

The unconditional posterior mean and variance of random effect  $j$  of the Gaussian model are

$$E(\mu_j|\mathbf{y}) = E(E(\mu_j|A, \beta, \mathbf{y})|\mathbf{y}) = (1 - E(B_j|\mathbf{y}))y_j + E(B_j|\mathbf{y})E(\mu_j^E|\mathbf{y}) \quad (27)$$

$$Var(\mu_j|\mathbf{y}) = E(Var(\mu_j|A, \beta, \mathbf{y})|\mathbf{y}) + Var(E(\mu_j|A, \beta, \mathbf{y})|\mathbf{y}) \quad (28)$$

$$= (1 - V_j)E(B_j|\mathbf{y}) + Var(B_j(y_j - \mu_j^E)|\mathbf{y}) \quad (29)$$

$$= h_g(E(B_j|\mathbf{y}), E(B_j^2|\mathbf{y}), E(\mu_j^E|\mathbf{y}), E((\mu_j^E)^2|\mathbf{y})). \quad (30)$$

The unconditional posterior mean and variance of random effect  $j$  under the Gaussian model are also functions of the unconditional posterior moments of the shrinkage factors and expected random effects. We specify the function  $h_g$  in Appendix B.

For the Poisson model, the unconditional posterior mean and variance of random effect  $j$  are

$$E(\lambda_j|\mathbf{y}) = E(E(\lambda_j|r, \mathbf{y})|\mathbf{y}) = (1 - E(B_j|\mathbf{y}))\bar{y}_j + E(B_j|\mathbf{y})\lambda^E \quad (31)$$

$$Var(\lambda_j|\mathbf{y}) = E(Var(\lambda_j|r, \mathbf{y})|\mathbf{y}) + Var(E(\lambda_j|r, \mathbf{y})|\mathbf{y}) \quad (32)$$

$$= E(\lambda_j^*/(r + n_j)|\mathbf{y}) + Var(B_j(\bar{y}_j - \lambda_j^E)|\mathbf{y}) \quad (33)$$

$$= h_p(E(B_j|\mathbf{y}), E(B_j^2|\mathbf{y})). \quad (34)$$

The unconditional posterior mean and variance of random effect  $j$  under the Poisson model are functions of the unconditional posterior moments of the shrinkage factors. We specify the function  $h_p$  in Equation 34 in Appendix C.

Next, we estimate the unconditional posterior moments of the shrinkage factors after approximating their unconditional posterior distributions by Beta distributions via the ADM.



### Unconditional posterior moments of shrinkage factors

It is noted that the shrinkage factors  $(B_1, \dots, B_k)$  are a function of  $r$ , i.e.,  $B_j = r/(r + n_j) = B_j(r)$  (or a function of  $A$  for the Gaussian model). One way to approximate the distribution of  $B_j$  is to find the maximum likelihood estimate of  $r$ ,  $\hat{r}_{MLE}$ , with its Hessian value and to use a Delta method for an asymptotic Normal distribution of  $B_j(\hat{r}_{MLE})$ . 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 the 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 as

$$B_j|\mathbf{y} \sim \text{Beta}(a_{1j}, a_{0j}), \text{ for } j = 1, 2, \dots, k, \quad (35)$$

and the ADM estimates the two parameters of the Beta distribution, i.e.,  $a_{1j}$  and  $a_{0j}$ .

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  $A(B_j|\mathbf{y})dB_j \propto B_j(1 - B_j)f(B_j|\mathbf{y})dB_j$  so that its mode is the same as the mean of the original Beta distribution. The assumed posterior mean and variance of the shrinkage factor are

$$E(B_j|\mathbf{y}) = \frac{a_{1j}}{a_{1j} + a_{0j}} = \arg \max_{B_j} A(B_j|\mathbf{y}) \equiv B_j^*, \quad (36)$$

$$\text{Var}(B_j|\mathbf{y}) = \frac{B_j^*(1 - B_j^*)}{a_{1j} + a_{0j} + 1} = \frac{B_j^*(1 - B_j^*)}{B_j^*(1 - B_j^*)[-\frac{d^2}{dB_j^2} \log(A(B_j|\mathbf{y}))|_{B_j=B_j^*}] + 1}. \quad (37)$$

The ADM estimates these mean and variance using the marginal posterior distribution of  $r$ ,  $f(r|\mathbf{y}) \propto L(r)dr/r^2$ , where the marginal likelihood  $L(r) = \int L(\beta, r)d\beta$  for the Binomial model is obtained via the Laplace approximation with a Lebesgue measure on  $\beta$  and that for the Poisson model is specified in Equation 17; see Berger, Liseo, Wolpert *et al.* (1999) for the integrated likelihood in detail. For the Gaussian model, the marginal likelihood function of  $A$  is available in a closed form with  $\beta$  integrated out, see Morris and Tang (2011).

Considering that Equation 36 and 37 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. Since  $A(B_j|\mathbf{y})$  is proportional to the marginal posterior density  $f(\alpha|\mathbf{y}) \propto e^\alpha L(\alpha)$ , the estimated posterior mean in Equation 36 is

$$\hat{B}_j^* = \frac{e^{-\hat{\alpha}}}{n_j + e^{-\hat{\alpha}}}, \quad (38)$$

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

We need the invariance information introduced in [Morris and Tang \(2011\)](#) to estimate the variance in Equation 37, which is defined as

$$\begin{aligned} \text{inv.info} &\equiv -\frac{d^2 \log(A(B_j|\mathbf{y}))}{d[\text{logit}(B_j)]^2} \Big|_{B_j=\hat{B}_j^*} = -\frac{d^2 \log(A(B_j(r)|\mathbf{y}))}{d[\log(r)]^2} \Big|_{r=\hat{r}} \\ &= -\frac{d^2 \log(A(B_j(r(\alpha))|\mathbf{y}))}{d\alpha^2} \Big|_{\alpha=\hat{\alpha}} \end{aligned} \quad (39)$$

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 posterior variance in Equation 37 as

$$\widehat{\text{Var}}(B_j|\mathbf{y}) = \frac{\hat{B}_j^{*2}(1 - \hat{B}_j^*)^2}{\text{inv.info} + \hat{B}_j^*(1 - \hat{B}_j^*)}. \quad (40)$$

After matching the estimated unconditional posterior mean and variance of shrinkage factor  $j$  in Equation 38 and 40 to the two parameters of the Beta distribution in Equation 35, i.e.,  $a_{1j}$  and  $a_{0j}$ , we get their estimates as

$$\hat{a}_{1j} = \frac{\text{inv.info}}{1 - \hat{B}_j^*} \quad \text{and} \quad \hat{a}_{0j} = \frac{\text{inv.info}}{\hat{B}_j^*}. \quad (41)$$

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

$$\hat{E}(B_j^c|\mathbf{y}) = \frac{B(\hat{a}_{1j} + c, \hat{a}_{0j})}{B(\hat{a}_{1j}, \hat{a}_{0j})}. \quad (42)$$

The distributional approximation of the ADM gets better if the true posterior distribution of the shrinkage factor is closer to the Beta distribution ([Christiansen and Morris 1997](#); [Morris and Tang 2011](#); [Morris and Lysy 2012](#)).

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

The Gaussian and Binomial models need to estimate the unconditional posterior moments of expected random effects. We estimate them using their relationship to the conditional posterior moments. For a non-negative constant  $c$ , the unconditional posterior moments are

$$E((\mu_j^E)^c|\mathbf{y}) = E(E((\mu_j^E)^c|A, \mathbf{y})|\mathbf{y}) \quad \text{and} \quad E((p_j^E)^c|\mathbf{y}) = E(E((p_j^E)^c|r, \mathbf{y})|\mathbf{y}). \quad (43)$$

We use the ADM to estimate the conditional expectations inside, and plug-in  $\hat{A} = \exp(\hat{\alpha})$  or  $\hat{r} = \exp(-\hat{\alpha})$  to estimate the outer unconditional expectations.

To be specific, the ADM for the Binomial model assumes that the conditional posterior distribution of each expected random effect given  $\hat{\alpha}$  follows a Beta distribution because the moment calculation of the expected random effects  $(p_1^E, p_2^E, \dots, p_k^E)$  involves an intractable integration. For example, the first conditional posterior moment is

$$E(p_j^E|\hat{\alpha}, \mathbf{y}) = E\left(\frac{e^{x_j^\top \beta}}{1 + e^{x_j^\top \beta}} \Big| \hat{\alpha}, \mathbf{y}\right) = \int_{\mathbf{R}^m} \frac{e^{x_j^\top \beta}}{1 + e^{x_j^\top \beta}} f(\beta|\hat{\alpha}, \mathbf{y}) d\beta. \quad (44)$$

We assume the conditional posterior distribution of expected random effect  $j$  is a Beta distribution as follows;

$$p_j^E | \hat{\alpha}, \mathbf{y} = \frac{e^{x_j^\top \beta}}{1 + e^{x_j^\top \beta}} \Big| \hat{\alpha}, \mathbf{y} \sim \text{Beta}(b_{1j}, b_{0j}) \sim \frac{G(b_{1j})}{G(b_{1j}) + G(b_{0j})}, \quad (45)$$

where  $G(b_{1j})$  is a random variable following a  $\text{Gamma}(b_{1j}, 1)$  distribution with a unit scale and independently  $G(b_{0j})$  has a  $\text{Gamma}(b_{0j}, 1)$  distribution. Note that the representation in Equation 45 is equivalent to saying  $e^{x_j^\top \beta} | \hat{\alpha}, \mathbf{y} \sim G(b_{1j})/G(b_{0j})$ , a ratio of two independent Gamma random variables. Its mean and variance are

$$E(e^{x_j^\top \beta} | \hat{\alpha}, \mathbf{y}) = E\left(\frac{G(b_{1j})}{G(b_{0j})}\right) = \frac{b_{1j}}{b_{0j} - 1} \equiv \eta_j, \quad (46)$$

$$\text{Var}(e^{x_j^\top \beta} | \hat{\alpha}, \mathbf{y}) = \text{Var}\left(\frac{G(b_{1j})}{G(b_{0j})}\right) = \frac{\eta_j(1 + \eta_j)}{b_{0j} - 2}. \quad (47)$$

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  $\text{Normal}[\hat{\beta}, \hat{\Sigma}]$ , where  $\hat{\beta}$  is the mode of  $p(\beta | \hat{\alpha}, \mathbf{y})$  and  $\hat{\Sigma}$  is a 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 get the numerical values of the posterior mean and variance in Equation 46 and 47 as

$$\eta_j = e^{x_j^\top \hat{\beta} + x_j^\top \hat{\Sigma} x_j / 2}, \quad (48)$$

$$\widehat{\text{Var}}(e^{x_j^\top \beta} | \mathbf{y}) = \eta_j^2 (e^{x_j^\top \hat{\Sigma} x_j} - 1). \quad (49)$$

By matching the mean and variance in Equation 48 and 49 to  $b_{1j}$  and  $b_{0j}$  in Equation 46 and 47, we obtain the values of  $b_{1j}$  and  $b_{0j}$  as follows;

$$b_{1j} = \eta_j + \frac{\eta_j + 1}{e^{x_j^\top \hat{\Sigma} x_j} - 1} \quad \text{and} \quad b_{0j} = \frac{\eta_j + 1}{\eta_j (e^{x_j^\top \hat{\Sigma} x_j} - 1)} + 2. \quad (50)$$

Finally we estimate the unconditional posterior moments of the expected random effects,  $E((p_j^E)^c | \mathbf{y})$ , by  $E((p_j^E)^c | \hat{\alpha}, \mathbf{y}) = B(b_{1j} + c, b_{0j})/B(b_{1j}, b_{0j})$  for  $c \geq 0$ .

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

$$(X^\top D_{V+\hat{A}}^{-1} X)^{-1} X^\top D_{V+\hat{A}}^{-1} \mathbf{y} \quad \text{and} \quad (X^\top D_{V+\hat{A}}^{-1} X)^{-1}, \quad \text{respectively,} \quad (51)$$

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 obtain the conditional posterior moments of  $\mu_j^E = \mathbf{x}^\top \beta$  with ease and use them to estimate their unconditional posterior moments.

## 4.2. Estimation for random effects

It is intractable to derive analytically the unconditional posterior distribution of each random effect for our three models, which motivates us to assume again that each random effect has

a skewed Normal distribution ([Azzalini 1985](#)) for the Gaussian model, a Beta distribution for the Binomial model, or a Gamma distribution for the Poisson model; for  $j = 1, 2, \dots, k$ ,

$$\mu_j | \mathbf{y} \sim \text{skewed-Normal}(\phi, \omega, \delta), \quad (52)$$

$$p_j | \mathbf{y} \sim \text{Beta}(t_{1j}, t_{0j}), \quad (53)$$

$$\lambda_j | \mathbf{y} \sim \text{Gamma}(s_{1j}, s_{0j}). \quad (54)$$

[Kelly \(2014\)](#) matches the three estimated unconditional posterior moments of the random effects to the three parameters  $(\phi, \omega, \delta)$  of the skewed-Normal distribution in Equation 52, i.e., location parameter  $\phi$ , scale parameter  $\omega$ , and a skewness parameter  $\delta$ . (Hey Joey, could you add a brief description about how to match moments?) [Kelly \(2014\)](#) also shows that the skewed-Normal approximation improves a Normal approximation ([Morris and Tang 2011](#); [Morris and Lysy 2012](#)) in terms of the frequency properties of random effects.

We know that the unconditional posterior mean and variance of each random effect for the Binomial model specified in Equation 22 and 26 are functions of the unconditional posterior moments of shrinkage factors and expected random effect. By plugging-in the estimated unconditional posterior moments of shrinkage factors and those of expected random effect, 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. The estimates of two parameters  $t_{1j}$  and  $t_{0j}$  in Equation 53 come as follows;

$$\hat{t}_{1j} = \left( \frac{\hat{\mu}_{p_j}(1 - \hat{\mu}_{p_j})}{\hat{\sigma}_{p_j}^2} - 1 \right) \hat{\mu}_{p_j}, \text{ and } \hat{t}_{0j} = \left( \frac{\hat{\mu}_{p_j}(1 - \hat{\mu}_{p_j})}{\hat{\sigma}_{p_j}^2} - 1 \right) (1 - \hat{\mu}_{p_j}). \quad (55)$$

For the Poisson model, let  $\hat{\mu}_{\lambda_j}$  and  $\hat{\sigma}_{\lambda_j}^2$  denote the estimated unconditional posterior mean and variance, respectively, with the estimated unconditional posterior moments of shrinkage factors plugged-in to Equation 31 and 34. The estimates of the two parameters  $s_{1j}$  and  $s_{0j}$  in Equation 54 are

$$\hat{s}_{1j} = \frac{\hat{\mu}_{\lambda_j}^2}{\hat{\sigma}_{\lambda_j}^2}, \text{ and } \hat{s}_{0j} = \frac{\hat{\mu}_{\lambda_j}}{\hat{\sigma}_{\lambda_j}^2}. \quad (56)$$

Finally, the approximate unconditional posterior distribution of random effect  $j$  for the Binomial model is

$$p_j | \mathbf{y} \sim \text{Beta}(\hat{t}_{1j}, \hat{t}_{0j}), \quad (57)$$

and that for the Poisson model is

$$\lambda_j | \mathbf{y} \sim \text{Gamma}(\hat{s}_{1j}, \hat{s}_{0j}). \quad (58)$$

We use the mean and (2.5%, 97.5%) quantiles (if we assign 95% confidence level) of the estimated approximate posterior distribution, a skewed Normal distribution for the Gaussian model, a Beta distribution in Equation 57 for the Binomial model, or a Gamma distribution in Equation 58 for the Poisson model, as the point and interval estimates of random effects.

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

As for the Binomial model, the package **Rgbp** also provides a way to independently draw exact posterior samples of random effects and hyper-parameters via the acceptance-rejection method. We continue working on a logarithmic scale of  $r$ ,  $\alpha = \log(1/r) = -\log(r)$ . The joint posterior density function of  $\alpha$  and  $\beta$  based on their joint hyper-prior density function in Equation 12 is

$$f(\alpha, \beta | \mathbf{y}) \propto f(\alpha, \beta) L(\alpha, \beta) \propto e^\alpha L(\alpha, \beta) d\alpha d\beta. \quad (59)$$

The Acceptance-Rejection (A-R) method (Everson and Morris 2000; Tang 2002) 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 independently draw posterior samples of hyper-parameters from  $f(\alpha, \beta | \mathbf{y})$ , 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  $e^{-\alpha(k-1)}$  when  $\alpha$  goes to  $\infty$  and as  $e^\alpha$  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 it 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, since  $f(\beta | \alpha, \mathbf{y})$  of the approximate Gaussian counterpart 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

$$g_1(\alpha) = g_1(\alpha; \mu, \sigma, a, b) \equiv \text{Skewed-}t(\alpha | \mu, \sigma, a, b), \quad (60)$$

$$g_2(\beta) = g_2(\beta; \boldsymbol{\mu}^*, S_{(m \times m)}) \equiv t_4(\beta | \boldsymbol{\mu}^*, S), \quad (61)$$

where the notation  $\text{Skewed-}t(\alpha | \mu, \sigma, a, b)$  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). The article of Jones and Faddy (2003) derives the mode of  $g_1(\alpha)$  as

$$\mu + \frac{(a - b)\sqrt{a + b}}{\sqrt{(2a + 1)(2b + 1)}}. \quad (62)$$

The article also show that the tails of the skewed- $t$  density function follow a power law as  $\alpha^{-(2a+1)}$  on the left and  $\alpha^{-(2b+1)}$  on the right when  $b > a$ . It also provides a representation to generate the random variable following their skewed- $t$  distribution as

$$\alpha \sim \mu + \sigma \frac{\sqrt{a + b}(2T - 1)}{2\sqrt{T(1 - T)}}, \text{ where } T \sim \text{Beta}(a, b). \quad (63)$$

The notation  $t_4(\beta|\mu^*, S)$  in Equation 61 indicates a density function of a multivariate  $t$ -distribution at  $\beta$  with 4 degrees of freedom, a location vector  $\mu^*$ , and a  $m \times m$  scale matrix  $S$  that leads to the variance-covariance matrix  $2S$ .

We set the parameters of  $g_1(\alpha)$  and  $g_2(\beta)$ , i.e.,  $\mu$ ,  $\sigma$ ,  $a$ ,  $b$ ,  $\mu^*$ , 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, we obtain 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.

Next, we set  $(a, b)$  to  $(k, 2k)$  if  $k < 10$  (or otherwise  $(\log(k), 2\log(k))$ ) to maintain a left-skewness of  $g_1(\alpha)$  and to keep  $a$  and  $b$  small enough for thick tails. We match the mode of  $g_1(\alpha)$  specified in Equation 62 to  $\hat{\alpha}$  by fixing the location parameter  $\mu$  at  $\hat{\alpha} - (a - b)\sqrt{a + b}/\sqrt{(2a + 1)(2b + 1)}$ . We set the scale parameter  $\sigma$  to  $(-H_{\hat{\alpha}}^{-1})^{0.5}\psi$ , where  $\psi$  is a tuning parameter;  $\psi = 1.3$  is the default. When the A-R method produces extreme values of weights defined in Equation 65 below, we increase the value of  $\psi$ .

As for  $g_2(\beta)$ , we matches the location vector  $\mu^*$  to the mode  $\hat{\beta}$  and the scale matrix  $S$  to  $-H_{\hat{\beta}}^{-1}/2$  so that the variance-covariance matrix becomes  $-H_{\hat{\beta}}^{-1}$ ;

$$g_2(\beta) \equiv t_4(\beta|\mu^* = \hat{\beta}, S = -H_{\hat{\beta}}^{-1}/2). \quad (64)$$

For the implementation of the acceptance-rejection method, we obtain four times more trial samples than the desired number of samples  $N$  independently from  $g_1(\alpha)$  and  $g_2(\beta)$ . We calculate  $4N$  weights, each of which is defined as

$$w_i \equiv w(\alpha^{(i)}, \beta^{(i)}) = \frac{f(\alpha^{(i)}, \beta^{(i)}|\mathbf{y})}{g_1(\alpha^{(i)})g_2(\beta^{(i)})}, \text{ for } i = 1, 2, \dots, 4N. \quad (65)$$

We accept 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. The usual acceptance rates from our data examples are around 25%. In a case where we accept more than the desired number of samples  $N$ , we discard the redundant. If the number of accepted samples is smaller than  $N$ , then we sample additional pairs (6 times more than the shortage) and calculate a new maximum  $M'$  from all the previous and new weights, accepting or rejecting the entire pairs again with new probabilities  $w_i/M'$ .

Once we have posterior samples of hyper-parameters, it is easy to obtain posterior samples of random effects via a Monte Carlo integration below.

$$f(\mathbf{p}|\mathbf{y}) = \int f(\mathbf{p}|\alpha, \beta, \mathbf{y}) \cdot f(\alpha, \beta|\mathbf{y}) d\alpha d\beta. \quad (66)$$

The integration can be done by sampling  $(p_1, p_2, \dots, p_k)$  from the independent Beta conditional posterior distributions  $f(p_j|\beta, r, \mathbf{y})$  in Equation 53 given  $r (= e^{-\alpha})$  and  $\beta$  already sampled from  $f(\alpha, \beta|\mathbf{y})$  via the A-R method.

## 6. Frequency method checking

Whether the 95% interval estimates of random effects obtained by a specific model achieve the nominal 95% confidence level for any true parameter values is one of the key model evaluation

criteria. A frequency method checking is a procedure to evaluate it, which is different from a model checking that tests whether a two-level model is appropriate for data (overdispersion exists in data) (Dean 1992; Christiansen and Morris 1996). Conditioning that the two-level model is appropriate, the frequency method checking generates pseudo-data sets given specific values of hyper-parameters (a parametric bootstrapping) and estimates unknown coverage probabilities based on these mock data sets.

From now on, the explanation will be 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 Equation 2 is completely determined by two hyper-parameters,  $A$  and  $\beta$ . Fixing these hyper-parameters at specific values, we generate  $N_{sim}$  sets of random effects from the conjugate prior distribution, i.e.,  $\{\mu^{(i)}, i = 1, \dots, N_{sim}\}$ , where the superscript  $(i)$  indicates the  $i$ -th simulation. Next, using the distribution of observed data in Equation 1, we generate  $N_{sim}$  sets of observed data sets  $\{y^{(i)}, i = 1, \dots, N_{sim}\}$  given each  $\mu^{(i)}$ . Note that we generate one observed data set per one set of random effects.

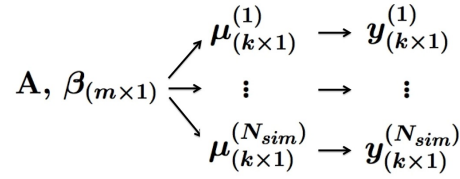


Figure 1: Pseudo-data generating process

### 6.2. Coverage probability estimation

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

$$I_{A, \beta}(\mu_j^{(i)}) = \begin{cases} 1, & \text{if } \mu_j^{(i)} \in (\hat{\mu}_{j, low}^{(i)}, \hat{\mu}_{j, upp}^{(i)}) \\ 0, & \text{otherwise} \end{cases} \quad (67)$$

We consider the coverage indicators as functions of  $A$  and  $\beta$  because outcomes of indicators depend on the simulated random effects and mock data generated by these hyper-parameters.

*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 of random effect  $j$ . This 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 Equation 67 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  $C_{A,\beta}(\mu_j)$ .

$$\bar{I}_{A,\beta}(\mu_j) = \frac{1}{N_{sim}} \sum_{i=1}^{N_{sim}} I_{A,\beta}(\mu_j^{(i)}), \quad j = 1, 2, \dots, k. \quad (68)$$

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  $Var(\bar{I}_{A,\beta}(\mu_j))$  is

$$\widehat{Var}(\bar{I}_{A,\beta}(\mu_j)) = \frac{1}{N_{sim}(N_{sim} - 1)} \sum_{i=1}^{N_{sim}} (I_{A,\beta}(\mu_j^{(i)}) - \bar{I}_{A,\beta}(\mu_j))^2, \quad j = 1, 2, \dots, k. \quad (69)$$

*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 effects, which is more efficient than the simple indicator-based coverage estimator. For  $j = 1, 2, \dots, k$ ,

$$C_{A,\beta}(\mu_j) = E(\bar{I}_{A,\beta}(\mu_j)|A, \beta) = E\left[\frac{1}{N_{sim}} \sum_{i=1}^{N_{sim}} E(I_{A,\beta}(\mu_j^{(i)})|A, \beta, \mathbf{y}^{(i)}) \middle| A, \beta\right], \quad (70)$$

where the sample mean of conditional expectations inside the outer expectation is the RB unbiased coverage estimator. To be specific,

$$\begin{aligned} \bar{I}_{A,\beta}^{RB}(\mu_j) &= \frac{1}{N_{sim}} \sum_{i=1}^{N_{sim}} E(I_{A,\beta}(\mu_j^{(i)})|A, \beta, \mathbf{y}^{(i)}) \\ &= \frac{1}{N_{sim}} \sum_{i=1}^{N_{sim}} Pr(\mu_j^{(i)} \in (\hat{\mu}_{j, low}^{(i)}, \hat{\mu}_{j, upp}^{(i)})|A, \beta, \mathbf{y}^{(i)}). \end{aligned} \quad (71)$$

We can easily compute the above conditional posterior probabilities using the cumulative density function of the Normal conditional posterior distribution of each random effect in Equation 3. The variance of  $\bar{I}_{A,\beta}^{RB}(\mu_j)$  is smaller than or equal to the variance of a simple coverage estimator  $\bar{I}_{A,\beta}(\mu_j)$  (Rao 1945; Blackwell 1947).

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

$$\widehat{Var}(\bar{I}_{A,\beta}^{RB}(\mu_j)) \equiv \frac{1}{N_{sim}(N_{sim} - 1)} \sum_{i=1}^{N_{sim}} \left( E(I_{A,\beta}(\mu_j^{(i)})|A, \beta, \mathbf{y}^{(i)}) - \bar{I}_{A,\beta}^{RB}(\mu_j) \right)^2. \quad (72)$$

*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 coverage estimator and its variance estimator as follows.

$$\bar{I}_{r,\beta}^{RB} = \frac{1}{k} \sum_{j=1}^k \bar{I}_{r,\beta}^{RB}(p_j) \quad \text{and} \quad \widehat{Var}(\bar{I}_{RB}) = \frac{1}{k^2} \sum_{j=1}^k \widehat{Var}(\bar{I}_{r,\beta}^{RB}(p_j)). \quad (73)$$

## 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. To use these relevant functions, we need to save the outcome of **gbp** into an object such as **g.output** in the code below).

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 knowing the value of the the expected random effect, meaning that the intercept term must be either estimated or known. The default of **gbp** in this case is to fit an intercept term. We can assign the value of the known expected random effect through an optional argument **mean.PriorDist**. Note that **gbp** can fit the Poisson model only when the value of expected random effect,  $\lambda^E$ , is 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 numbers of events happening 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. If we want to designate a known value of the expected random effect  $\mu^E = \beta_1$  for the Gaussian model or  $p^E = \exp(\beta_1)/(1+\exp(\beta_1))$  for the Binomial model, then we put that value into **gbp** using an argument **mean.PriorDist**. For example, if the  $\mu^E$  is known as 10, then we use the following code.

```
R> g.output <- gbp(y, se.or.n, mean.PriorDist = 10, model = "gaussian")
```

If covariate information for each group is available, we 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 we 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 as a default. If we do not want to include an intercept term, estimating two regression coefficients for the two covariates, we 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 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 the function **gbp** use a Normal approximation to the unconditional posterior distribution of the random effect (Morris and Tang 2011). The default of **normal.CI** is **FALSE** for the skewed-Normal approximation (Kelly 2014).

The function **gbp** uses the acceptance-rejection method to fit the Binomial model if we assign the desired number of posterior samples ( $N$  in Equation 65) through the argument **n.AR**; its default value is 0. There are several arguments related to the acceptance-rejection method. The argument **n.AR.factor** determines how many trial samples the method draws; its default value is 4, meaning that the method draws  $\text{n.AR} \times 4$  trial samples and accepts or rejects them. The argument **trial.scale** is  $\psi$  determining the scale parameter of the skewed- $t$  distribution (the envelope function); its default value is 1.3. The argument **save.result** indicates whether we save 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 chose the joint hyper-prior distribution  $f(r, \beta) \propto d\beta dr / (t + r)^{u+1}$ ; the default values for **t** and **u** are 0 and 1, respectively, for the joint hyper-prior distribution specified in Equation 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 acceptance-rejection 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 detailed outcomes. The function **names** will show the list of outcomes (too many to be listed in this article); see the document of **gbp** to check the details of the items saved.

```
R> names(b.output)
R> ?gbp
```

The S3 object “gbp” benefits from three generic functions, **print**, **summary**, and **plot**. The estimation result for all the random effects pops up if we type the “gbp” object at 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 ascending order of standard errors for the Gaussian model. When the argument **sort** is **FALSE**, the estimation result comes out as 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.

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

The function `plot` draws a shrinkage plot and  $100 \times \text{Alpha}\%$  interval plot for random effects with its default argument “`sort = TRUE`” that draws the  $100 \times \text{Alpha}\%$  interval plot in the ascending order of  $n$  for the Binomial and Poisson model and the ascending order of standard errors for the Gaussian model. When the argument `sort` is set to `FALSE` the  $100 \times \text{Alpha}\%$  interval plot will be plotted 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 a frequency method checking that estimates the coverage probabilities of random effects, conditioning on the values of the hyper-parameters that generate random effects and mock data sets. 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 `sim` sets the number of simulations  $N_{sim}$  defined in Section 6. If we do not assign values of the hyper-parameters by the arguments `A.or.r` and `reg.coef` as above, 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) to the generative values of hyper-parameters. If we want to conduct the frequency method checking with different generated values of hyper-parameters, for example,  $r = 100$  and  $\beta^T = (2, 5)$  when one covariate was used (including an intercept term) to create the “gbp” object, then we specify the code as

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

When we fit a model with a known value of the expected random effect, we may want to conduct the frequency method checking with a different value for the expected random effect. In this case, the argument `mean.PriorDist` enables it. For example, when we obtain the “gbp” object `p.output` after fitting a Poisson model with a specific value of  $\lambda^E$ , if we want to conduct the frequency method checking with the same specific value of  $\lambda^E$ , then we do not specify the argument `mean.PriorDist` in `coverage`. Once we specify the argument `mean.PriorDist` in `coverage`, for example 30 as below, the frequency method checking will be based on the estimated posterior mode of  $r$  and the newly specified value of  $\lambda^E$ , 30.

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

Although the function `coverage` does not produce an S3 object, a detailed summary appears in a plot that automatically pops up after running the function `coverage`. In addition, the object `cov` that contains the result of `coverage` contains numerical details, see `names(cov)` or type “`?coverage`” for the list of detailed outcomes. If we save the result into an object such as `cov` above, then we can always recall the plot, using the function `coverage.plot`.

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

## 8. Examples

### 8.1. Data of 31 hospitals with a 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)<sup>1</sup>. The data set contains the number of deaths,  $y$ , for a specified period after CABG surgeries out of the total number of patients,  $n$ , receiving CABG surgeries in each hospital. Health care providers may use these data to improve their care qualities, and patients may refer to the data to select a better hospital. Our goal is to obtain the point and interval estimates for the unknown true mortality rates of 31 hospitals (random effects) to evaluate each hospital's reliability on the CABG surgery. We interpret the caseloads,  $n$ , as exposures and assume that the state-level death rate per exposure of this surgery ( $\lambda^E$ ) is known as 0.030 to fit the Poisson model for an illustrative purpose. If covariate information is available or the information about the expected random effect is unavailable a priori, we recommend using the Binomial model.

These data can be loaded into R using the following code.

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

The function `gbp` fits the Poisson hierarchical model with the Gamma conjugate prior distribution as a population distribution of the death rates in New York states whose mean ( $\lambda^E$ ) is 0.030. The number of regression coefficients ( $m$ ) is 0 because we do not need to estimate them via a log-linear regression.

```
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

<sup>1</sup>Visit <http://www.health.ny.gov/statistics/diseases/cardiovascular/> for more profile data of the CABG surgery in New York state.

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
colMeans		517	0.03	0.600	0.0201	0.0293	0.0403	0.00517

The output contains information about the observed death rates  $\bar{y}_j$ , caseloads  $n_j$ , known prior mean  $\lambda^E$ , shrinkage estimates  $\hat{B}_j$ , lower bounds of interval estimates  $\hat{\lambda}_{j,low}$ , posterior means  $\hat{\lambda}_j \equiv E(\lambda_j|\mathbf{y})$ , upper bounds of interval estimates  $\hat{\lambda}_{j,upp}$ , and posterior standard deviations for each random effect based on the assumed unconditional Gamma posterior distributions.

A function `summary` shows selective information about hospitals with minimum, median, and maximum exposures and more detailed estimation results about 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 observed sample means of hospitals whose caseloads are less than 684 will shrink toward the prior mean (0.030) 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 ( $n_1$ ), and its posterior mean is  $(1 - 0.911) * 0.0448 + 0.911 * 0.030 = 0.0313$ . As for this hospital, using more information from the conjugate prior distribution is an appropriate choice because the amount of observed information (67) is far less than the amount of state-level information (684).

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

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

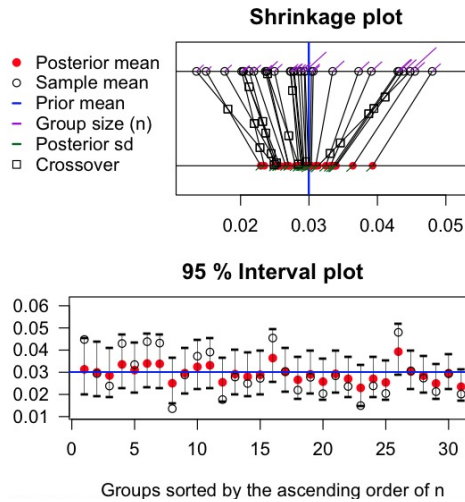


Figure 2: Shrinkage plot and 95% interval plot for 31 hospitals

In Figure 2 the regression towards the mean (RTTM) is obvious in the first graph; the observed death 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 0.030, to the different extents. Note that some hospitals' ranks have changed by shrinking much harder towards 0.030 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 among 31 hospitals in terms of the observed death rate is probably not the safest in terms of the estimated posterior mean accounting for the different caseloads of these two hospitals.



Intuitively, switching ranks for these two hospitals is reasonable. To be specific, their observed death 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 death 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 ( $\lambda^E=0.030$ ) much harder than that for of 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 estimated 95% intervals. We see that 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 death 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, 0.030) with smaller variance. Based on the point and interval estimates, the 31st hospital seems the most reliable one among all candidates.

When fitting a model it is always a good idea to question how reliable the estimation procedure is. The function `coverage` generates pseudo-datasets given the estimated value of  $r$ , 683.53, as a generative value. For reference, we can designate any generative values of  $r$  and  $\lambda^E$  by adding two arguments into the code below, for example, `A.or.r = 600` and `mean.PriorDist = 0.05`.

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

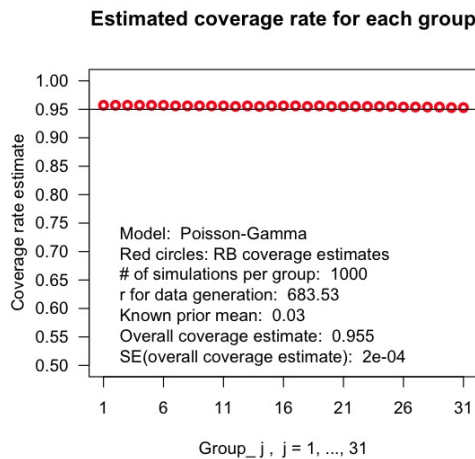


Figure 3: Coverage plot via frequency method checking for 31 hospitals

In Figure 3, the black horizontal line at 0.95 represents the nominal confidence level ( $\text{Alpha} = 0.95$  as a default) and the red circles indicate Rao-Blackwellized (RB) unbiased coverage

estimates,  $\bar{I}_{r,\lambda^E}^{RB}(\lambda_j)$  for  $j = 1, 2, \dots, 31$ . The overall RB unbiased coverage estimate across all the hospitals ( $\bar{I}_{r,\lambda^E}^{RB}$ ) is 0.955. None of RB unbiased coverage estimates for 31 hospitals are less than 0.95 regardless of their caseloads ( $n_j$ ). This result shows that the interval estimates for this particular dataset accurately achieves 95% confidence level.

The following code provides specific values of the 31 RB unbiased coverage estimates and their standard errors for each hospital (the output is omitted for a space concern).

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

And the code below shows 31 simple unbiased coverage estimates and their standard errors for each hospital.

```
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 ( $N_{sim}$ ) 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,  $C_{r,\lambda^E}(\lambda_1)$ .

See [Morris and Christiansen \(1995\)](#) for a similar ranking problem using a Poisson hierarchical modeling with a different parametrization.

## 8.2. Data of 8 schools with 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](#)). We can load this data set into R by the following codes.

```
R> library(Rgbp)
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 known sampling variance, i.e., standard error of each school is completely known. At the second hierarchy, the mean for each school is assumed to be drawn from a common Normal distribution and hence, we can use the Gaussian component of **gbp** to fit this Gaussian hierarchical model.

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

Summary for each unit (sorted by se):

	obs.mean	se	prior.mean	shrinkage	low.intv	post.mean	upp.intv	post.sd
5	-1.00	9.0	8.168	0.408	-13.297	2.737	16.692	7.634
2	8.00	10.0	8.168	0.459	-7.255	8.077	23.361	7.810
7	18.00	10.0	8.168	0.459	-1.289	13.484	30.821	8.176
4	7.00	11.0	8.168	0.507	-8.780	7.592	23.602	8.257
6	1.00	11.0	8.168	0.507	-13.027	4.633	20.131	8.441
1	28.00	15.0	8.168	0.657	-2.315	14.979	38.763	10.560
3	-3.00	16.0	8.168	0.685	-17.130	4.650	22.477	10.096
8	12.00	18.0	8.168	0.734	-10.208	9.189	29.939	10.227
colMeans		12.5	8.168	0.552	-9.163	8.168	25.723	8.900

This output from `gpb` summarizes the results. In this Gaussian hierarchical 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 `gpb` 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 **Rgpb** 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,  $\mu^E$  (**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)
```

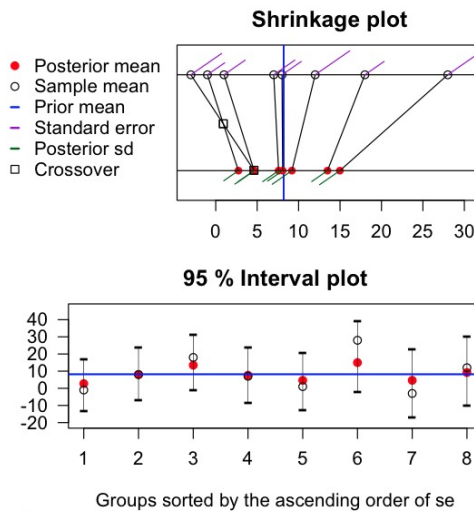


Figure 4: Shrinkage plot and 95% interval plot for 8 schools

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  $N_{sim}$  (**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 (black horizontal line at 0.95) is achieved (actually, exceeded) 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, Rao-Blackwellized (RB) unbiased coverage estimate and its standard error for each school can be gotten with the command below.

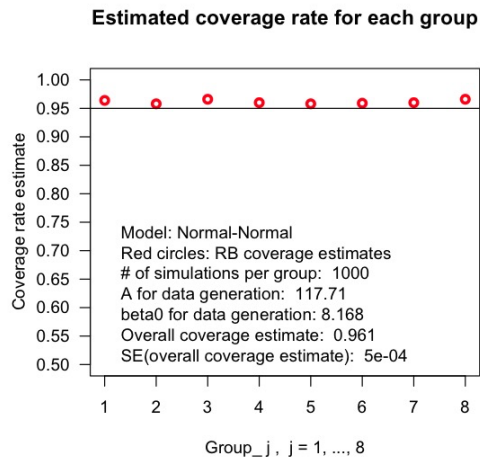


Figure 5: Coverage plot via frequency method checking for 8 schools

```
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
```

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

All the individual RB coverage estimates are saved in the  $8 \times 1,000$  matrix, each row of which is about each school. We can load this matrix by the following code.

```
R> g.coverage$raw.resultRB
```

### 8.3. Data of 18 baseball players with 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). We add a covariate about their outfielder indicators taking on 1 if a player is an outfielder and 0 otherwise. We can load the data into R with the following codes.

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

Conditioning on the unknown true batting average (random effect) for each player we assume 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,  $p_j$ , 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 conjugate Beta prior distribution 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	1.00	0.310	0.715	0.248	0.335	0.429	0.0462
2	0.378	45	1.00	0.310	0.715	0.244	0.329	0.420	0.0448
3	0.356	45	1.00	0.310	0.715	0.240	0.323	0.411	0.0437
4	0.333	45	1.00	0.310	0.715	0.236	0.316	0.403	0.0429
5	0.311	45	1.00	0.310	0.715	0.230	0.310	0.396	0.0424
6	0.311	45	0.00	0.233	0.715	0.179	0.256	0.341	0.0415
7	0.289	45	0.00	0.233	0.715	0.175	0.249	0.331	0.0400
8	0.267	45	0.00	0.233	0.715	0.171	0.243	0.323	0.0388
9	0.244	45	0.00	0.233	0.715	0.166	0.237	0.315	0.0380
10	0.244	45	1.00	0.310	0.715	0.210	0.291	0.379	0.0432
11	0.222	45	0.00	0.233	0.715	0.161	0.230	0.308	0.0377
12	0.222	45	0.00	0.233	0.715	0.161	0.230	0.308	0.0377
13	0.222	45	0.00	0.233	0.715	0.161	0.230	0.308	0.0377
14	0.222	45	1.00	0.310	0.715	0.202	0.285	0.375	0.0441
15	0.222	45	1.00	0.310	0.715	0.202	0.285	0.375	0.0441
16	0.200	45	0.00	0.233	0.715	0.155	0.224	0.302	0.0377
17	0.178	45	0.00	0.233	0.715	0.148	0.218	0.297	0.0381
18	0.156	45	0.00	0.233	0.715	0.140	0.211	0.292	0.0389
colMeans		45	0.44	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 all outfielders tends to be higher than that for infielders. The estimated odds ratio is  $\exp(0.389) = 1.48$ .

R> plot(b.output)

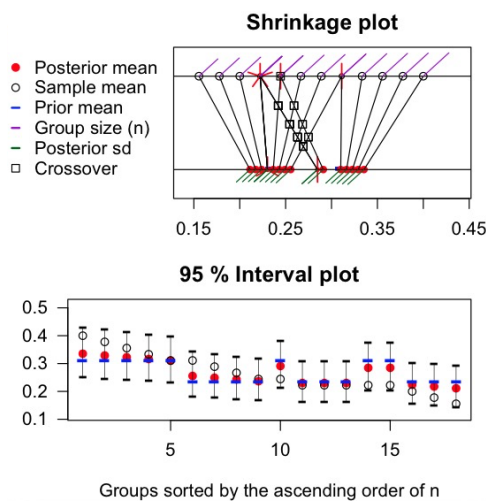


Figure 6: Shrinkage plot and 95% interval plot for 18 baseball players



It is evident in the shrinkage plot in Figure 6 that shrinkage occurs from the observed batting averages (empty dots) on the upper horizontal line towards the two prior means, 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 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 suggests that their batting averages will 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 given generative values.

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

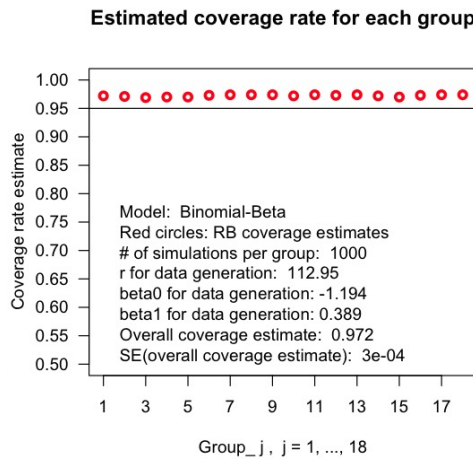


Figure 7: Coverage plot via frequency method checking for 18 players

In Figure 7, the estimated coverage probabilities for random effects are beyond 0.95, conservatively satisfying the definition of the 95% confidence interval. The Rao-Blackwellized (RB) overall unbiased coverage estimate is 0.972 (across all the players). Note that each coverage estimate depends on given true values of  $r$  and  $\beta_{(2 \times 1)}$ , and the assumption that the model is valid (overdispersion exists).

The RB unbiased coverage estimates and their standard errors for each player follow.

```
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 acceptance-rejection 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 2,000 posterior samples of  $\alpha$  (`b.output$alpha`), a  $2,000 \times 2$  matrix of  $\beta$  (`b.output$beta`) with each column corresponding to posterior samples of each regression coefficient, and a  $k \times 2,000$  matrix of random effects with each row contains posterior samples of each random effect.

If we run the frequency method checking using this “gbp” object obtained via the acceptance-rejection method, the  $N_{sim}$  simulations also run the acceptance-rejection method.

## 9. Discussion and summary

**Rgbp** is an R package for estimating and validating two-level Gaussian, Binomial and Poisson 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 hyper-parameter values.

As an alternative to other maximization based estimation methods such as MLE and REML, **Rgbp** provides 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 acceptance-rejection 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) 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.

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.

In addition to the built-in “frequency method checking” procedure, we may use the package to undergo “model checking”. For example, in the Gaussian hierarchical model, the assumed marginal distribution of the data is given in Equation 13 as part of the likelihood function. By substituting the point estimates of  $A$  and  $\beta$  from the package into this marginal distribution a test can be constructed to see whether the data follow the marginal distribution suggested by the hierarchical model.

## 10. 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.

### A. Unconditional posterior variance of the Binomial model

### B. Unconditional posterior variance of the Gaussian model

### C. Unconditional posterior variance of the Poisson model

## References

- Azzalini A (1985). "A class of distributions which includes the normal ones." *Scandinavian journal of statistics*, pp. 171–178.
- Berger JO, Liseo B, Wolpert RL, *et al.* (1999). "Integrated likelihood methods for eliminating nuisance parameters." *Statistical Science*, **14**(1), 1–28.
- Blackwell D (1947). "Conditional expectation and unbiased sequential estimation." *The Annals of Mathematical Statistics*, pp. 105–110.
- 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, Su YS, Yajima M, Hill J, Pittau MG, Kerman J, Zheng T (2014). "arm: data analysis using regression and multilevel/hierarchical models, 2010." URL <http://CRAN.R-project.org/package=arm>. *R package version*, pp. 1–3.

- 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.
- 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>.
- Patterson HD, Thompson R (1971). “Recovery of inter-block information when block sizes are unequal.” *Biometrika*, **58**(3), 545–554.
- Rao CR (1945). “Information and accuracy attainable in the estimation of statistical parameters.” *Bulletin of the Calcutta Mathematical Society*, **37**(3), 81–91.
- 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.” *R package version*, **1**.
- 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 (in preparation). “Posterior Propriety and Frequency Coverage Evaluation of Bayesian Beta-Binomial Logistic Regression Model.” *in preparation*.

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)