

## Handout 16: Hierarchical Bayes modeling

Lecturer &amp; author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

### Aim

To be able to specify and analyze a Hierarchical Bayesian, as well as to extend previously introduced concepts in the Hierarchical Bayes framework.

### Basic reading list:

- Berger, J. O. (2013; Section 4.6). Statistical decision theory and Bayesian analysis. Springer.
- Robert, C. (2007, Sections 10.1-10.3). The Bayesian choice: from decision-theoretic foundations to computational implementation. Springer Science & Business Media.
- Robert, C. P., & Reber, A. (1998). Bayesian modelling of a pharmaceutical experiment with heterogeneous responses. Sankhyā: The Indian Journal of Statistics, Series B, 145-160.

### R scripts:

**Bayesian Normal Mixture model (Appendix A):** [http://htmlpreview.github.io/?https://github.com/georgios-stats/Bayesian\\_Statistics\\_Michaelmas\\_2021/blob/main/Computer\\_practical/Normal\\_Mixture\\_model/Bayesian\\_Normal\\_Mixture\\_Model.nb.html](http://htmlpreview.github.io/?https://github.com/georgios-stats/Bayesian_Statistics_Michaelmas_2021/blob/main/Computer_practical/Normal_Mixture_model/Bayesian_Normal_Mixture_Model.nb.html)

**Bayesian Variable Selection in Logistic regression model (Appendix B):** [http://htmlpreview.github.io/?https://github.com/georgios-stats/Bayesian\\_Statistics\\_Michaelmas\\_2021/blob/main/Computer\\_practical/Bernoulli\\_regression\\_model\\_variable\\_selection/Bernoulli\\_Regression\\_Model\\_VS\\_full.nb.html](http://htmlpreview.github.io/?https://github.com/georgios-stats/Bayesian_Statistics_Michaelmas_2021/blob/main/Computer_practical/Bernoulli_regression_model_variable_selection/Bernoulli_Regression_Model_VS_full.nb.html)

**Random effect model (Appendix C):** [https://github.com/georgios-stats/Bayesian\\_Statistics\\_Michaelmas\\_2021/blob/main/Lecture\\_handouts/Rscripts/Hierarchical\\_bayes/HierarchicalBayesPharmaceutical.R](https://github.com/georgios-stats/Bayesian_Statistics_Michaelmas_2021/blob/main/Lecture_handouts/Rscripts/Hierarchical_bayes/HierarchicalBayesPharmaceutical.R)

## 1 Hierarchical Bayesian Model

A hierarchical Bayesian model involves several levels of conditional distributions. It can be hierarchical due to the sampling distribution modeling or due to the decomposition of the prior information.

**Definition 1.** A hierarchical Bayes model is a Bayesian statistical model with sampling distribution  $y \sim f(y|\theta)$  and prior  $\theta \sim \pi(\theta)$ , where the prior distribution  $\pi(\theta)$  is decomposed in conditional distributions. The Bayesian model is

$$\left\{ \begin{array}{l} y|\theta \sim f(y|\theta), \text{ sampling distribution} \\ \theta \sim \pi(\theta) \text{ marginal prior specified} \end{array} \right. \xrightarrow{\text{extend space}} \left\{ \begin{array}{ll} y|z, \theta \sim f(y|z, \theta) \\ z|\theta \sim f(z|\theta) \\ \theta|\phi_1 \sim \pi_1(\theta|\phi_1) & \text{1st level prior} \\ \phi_1|\phi_2 \sim \pi_2(\phi_1|\phi_2) & \text{2nd level hyper-prior} \\ \vdots & \\ \phi_j|\phi_{j+1} \sim \pi_{j+1}(\phi_j|\phi_{j+1}) & j\text{th level hyper-prior} \\ \vdots & \\ \phi_{m-1}|\phi_m \sim \pi_m(\phi_{m-1}|\phi_m) & m\text{th level hyper-prior} \end{array} \right. \quad (1)$$

**Remark 2.** The joint distribution  $p(y, z, \theta, \phi_1, \dots, \phi_j, \dots, \phi_{m-1})$  has pdf

$$p(y, z, \theta, \phi_1, \dots, \phi_j, \dots, \phi_{m-1}) = f(y|z, \theta) f(z|\theta) \pi_1(\theta|\phi_1) \pi_2(\phi_1|\phi_2) \pi_3(\phi_2|\phi_3) \dots \pi(\phi_{m-1}|\phi_m)$$

**Remark 3.** Hierarchical Bayesian model is simply a special type of Bayesian model, where

$$\begin{cases} y|\theta & \sim f(y|\theta) \\ \theta|\phi & \sim \pi(\theta|\phi) \\ \phi|\phi_m & \sim \pi(\phi|\phi_m) \end{cases} \quad (2)$$

for  $\phi = (\phi_1, \dots, \phi_{m-1})$ , and  $\phi_m$  fixed hyper-parameter.

**Remark 4.** The Bayesian model with sampling distribution  $y \sim f(y|\theta)$  and prior  $\theta \sim \pi(\theta)$ , can be recovered from (2) by marginalizing the prior as

$$\begin{aligned} f(y|z, \theta) &= \int_{\Phi} f(y|z, \theta) f(z|\theta) d\theta \\ \pi(\theta|\phi_m) &= \int_{\Phi} \pi(\theta|\phi) \pi(\phi|\phi_m) d\phi = \int_{\Phi_1 \times \Phi_{m-1}} \pi(\theta|\phi_1) \pi(\phi_1|\phi_2) d\phi_1 \dots \pi(\phi_{m-1}|\phi_m) d\phi_{m-1}, \end{aligned} \quad (3)$$

where  $\phi_m$  is just a fixed hyper-parameter. Hence, hierarchical models are indeed included in the Bayesian paradigm.

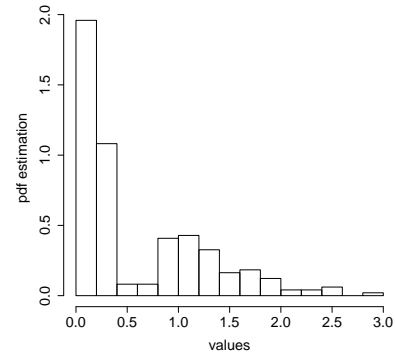
**Note 5.** A hierarchical Bayesian model can be used as a mean to specify more diverse priors. This is achieved by setting  $\phi$  to be a random hyper-parameter with  $\phi|\phi_m \sim \pi_2(\phi|\phi_m)$  instead of setting  $\phi$  to have a fixed value.

**Note 6.** A hierarchical Bayesian model can be used when the sampling distribution or the prior distributions justify a certain structure.

**Example 7.** [Normal mixture model] Consider the following application where our concern is the distribution of enzymatic activity in the blood, for an enzyme involved in the metabolism of carcinogenic substances, among a group of  $n = 245$  unrelated individuals; aka cluster analysis.

Our observables are  $y = (y_1, \dots, y_n)$  with  $n = 245$ . In the Boxplot on the right, we can clearly see that the distribution is multimodal, suggesting the existence of subpopulations/groups.

Interest lies on identifying subgroups of slow or fast metabolizers as a marker of genetic polymorphism in the general population. As we are interested in learning/identifying the sub-populations as the identity/label of the group from which each observation is drawn is unknown.



Questions of interest:

- How many groups exist?
- To which group each observation belongs?

Histogram of Enzyme dataset which is available from:  
<https://people.maths.bris.ac.uk/~mapjg/mixdata>

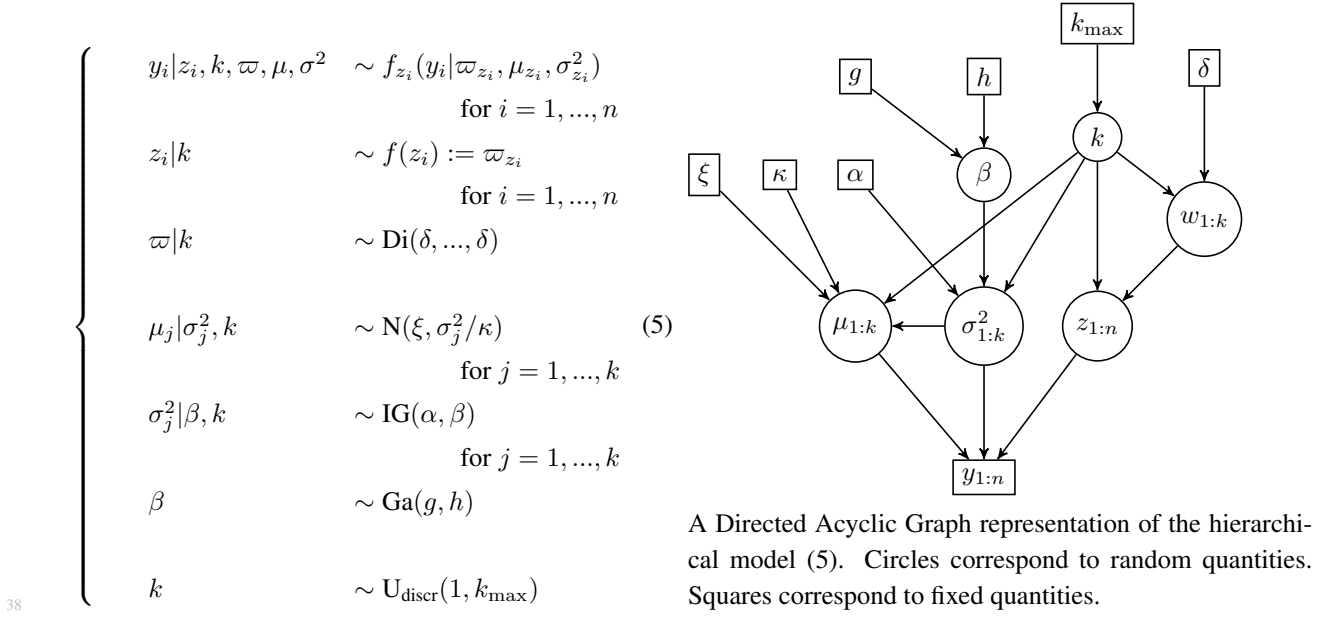
As for the sampling model, we can assume that the  $i$ -th observation  $y_i$  is randomly drawn from the  $j$ -th group which has proportion  $\varpi_j$  in the population and which is distributed according to the sampling distribution  $y_i|\theta_j \sim f_j(y_i|\theta_j)$ .

For simplicity, let's assume that all groups are Normally distributed but with different parameter values  $\{\theta_j\}$ ; hence  $j$ -th group is  $y_i|\mu_j, \sigma_j^2 \sim N(y_i|\mu_j, \sigma_j^2)$  with  $\theta_j = (\mu_j, \sigma_j^2)$ .

It is natural to regard the group label  $z_i$  for the  $i$ th observation as a latent allocation variable: then  $z_i$  is supposed to be distributed as  $z_i \sim f(z_i) = \varpi_{z_i}$  for  $z_i \in \{1, \dots, k\}$ , and  $y_i$  is supposed to be distributed as  $y_i|z_i, \theta_{z_i} \sim f_{z_i}(y_i|z_i, \theta_{z_i}) := N(y_i|\mu_{z_i}, \sigma_{z_i}^2)$ , for  $i = 1, \dots, n$ ; i.e.

$$\begin{cases} y_i|z_i, \mu_{z_i}, \sigma_{z_i}^2 & \sim f_{z_i}(y_i|\mu_{z_i}, \sigma_{z_i}^2) \\ z_i & \sim f(z_i) \end{cases} \implies \begin{cases} y_i|z_i, \mu_{z_i}, \sigma_{z_i}^2 & \sim N(y_i|\mu_{z_i}, \sigma_{z_i}^2) \\ z_i & \sim f(z_i) := \varpi_{z_i} \end{cases} \quad (4)$$

To complete the Bayesian model, we specify priors on the unknown quantities: Given there are  $k$  groups,  $\varpi_{1:k} \sim \text{Di}(\delta)$  for the group proportions,  $\mu_j \sim N(\xi_j, \sigma_j^2/\kappa)$  for the mean, and  $\sigma_j^2 \sim \text{Ga}(\alpha, \beta)$  for the variances. Assume we wish a more spread prior for  $\sigma_j^2$  (for some reason), and hence we specify a hyper-prior on  $\beta$  as  $\beta \sim \text{Ga}(g, h)$ . As the number of the groups is unknown, we assign prior  $k \sim \pi(k) \in \text{U}_{\text{discr}}(1, k_{\max})$ .



The joint distribution has pdf

$$p(y_{1:n}, z_{1:n}, k, \varpi_{1:k}, \mu_{1:k}, \sigma_{1:k}^2) = \underbrace{\prod_{i=1}^n N(y_i|\mu_{z_i}, \sigma_{z_i}^2)}_{f(y_{1:n}|z_{1:n}, \mu_{1:k}, \sigma_{1:k}^2)} \underbrace{\prod_{i=1}^n \varpi_{z_i}}_{f(z_{1:n}|k)} \underbrace{\text{Di}(\varpi_{1:k}|\delta)}_{\pi(\varpi_{1:k}|k)} \underbrace{\prod_{j=1}^k N(\mu_j|\xi, \sigma_j^2/\kappa)}_{\pi(\mu_{1:k}|\sigma_{1:k}^2, k)} \underbrace{\prod_{j=1}^k \text{Ga}(\sigma_j^2|\alpha, \beta)}_{\pi(\sigma_{1:k}^2|\beta, k)} \underbrace{\text{Ga}(\beta|g, h)}_{\pi(\beta)} \underbrace{\frac{1}{|k_{\max}|} \pi(k)}_{\pi(k)} \quad \text{Appendix A}$$

The joint posterior  $\pi(k, \varpi, \mu, \sigma^2, z|y)$  can be computed with the Bayesian theorem, and factorized as

$$\pi(k, \varpi, \mu, \sigma^2, \beta, z|y) = \frac{p(y, z, k, \varpi, \mu, \sigma^2, \beta)}{\int p(y, z, k, \varpi, \mu, \sigma^2, \beta) d(z, k, \varpi, \mu, \sigma^2, \beta)} \quad (6)$$

$$= \pi(z|y, k, \varpi, \mu, \sigma^2) \pi(\mu_{1:k}, \sigma_{1:k}^2|y, k) \pi(\varpi_{1:k}|y, k) \pi(k|y) \quad (7)$$

where to infer the number of groups from  $\pi(k|y)$ , the proportions in each group from  $\pi(\varpi_{1:k}|y, k)$ , the moments of each group from  $\pi(\mu_{1:k}, \sigma_{1:k}^2|y, k)$ , and the allocation of each observation to each group with  $\pi(z|y, k, \varpi, \mu, \sigma^2)$ .

As the required integrals are intractable, we can resolve to numerical methods, etc... Monte Carlo e.g. via JAGS...

*Note 8.* A particularly appealing aspect of hierarchical models is that they allow for conditioning on all levels, and this easy decomposition of the posterior. Consider the Bayesian hierarchical model (2) a parametric model  $f(y|\theta)$  with a

hierarchical prior  $\theta \sim \pi_1(\theta|\phi)$ , and  $\phi \sim \pi(\phi)$ . The posterior distribution of  $\theta$  is

$$\pi(\theta|y) = \int_{\Phi} \pi(\theta|y, \phi) \pi(\phi|y) d\phi \quad (8)$$

where

$$\begin{aligned} \pi(\theta|y, \phi) &= \frac{f(y|\theta) \pi_1(\theta|\phi)}{f_1(y|\phi)}; & \pi(\phi|y) &= \frac{f_1(y|\phi) \pi_2(\phi)}{f(y)}; \\ f_1(y|\phi) &= \int_{\Theta} f(y|\theta) \pi_1(\theta|\phi) d\theta; & f(y) &= \int_{\Theta} f_1(y|\phi) \pi_2(\phi) d\phi \end{aligned}$$

**Remark 9.** Note 8 has important consequences in terms of the computation of Bayes estimators, since it shows that  $\pi(\theta|y)$  can be simulated by generating, first,  $\phi$  from  $\pi(\phi|y)$  and then  $\theta$  from  $\pi(\theta|y, \phi)$ , if these two conditional distributions are easier to work with. (Snapshot from Term 2).

**Note 10.** Hierarchical decomposition (2) may facilitate the computation of intractable posterior moments. Let  $h$  be a function  $h : \Theta \rightarrow \mathbb{R}$ , then

$$E_{\pi}(h(\theta)|y) = E_{\pi}(E_{\pi}(h(\theta)|y, \phi) | y).$$

If  $E_{\pi}(h(\theta)|y) = \int h(\theta) \pi(\theta|y) d\theta$  is intractable and  $\theta$  has high dimensionality, one could possibly try to specify the prior decomposition  $\pi(\theta) = \int_{\Phi} \pi_1(\theta|\phi) \pi_2(\phi|\phi_m) d\phi$  in (3) such that  $E_{\pi}(h(\theta)|y, \phi)$  can be computed analytically, and  $\phi$  has low dimensionality. In that case one would have to compute the equivalent but lower dimensional (and hence easier) integral  $E_{\pi}(E_{\pi}(h(\theta)|y, \phi) | y) = \int E_{\pi}(h(\theta)|y, \phi) \pi(\phi|y) d\phi$ .

**Example 11.** (Cont. Example 7) From another point of view, recall that the compound distribution  $f(y_i|k, \varpi_{1:k} \theta_{1:k})$  of (4) is mixture model of distribution

$$y_i|k, \theta_{1:k} \sim f(y_i|k, \varpi_{1:k}, \theta_{1:k}) = \sum_{j=1}^k \varpi_j f_j(y_i|\theta_j) = \sum_{j=1}^k \varpi_j N(y_i|\mu_j, \sigma_j^2), \text{ for } i = 1, \dots, n \quad (9)$$

is a suitable sampling distribution for modeling heterogeneous populations. Then by marginalizing, we can get the equivalent model

$$\begin{cases} y_i|k, \varpi, \mu, \sigma^2 & \sim f(y_i|k, \varpi, \mu, \sigma^2) := \sum_{j=1}^k \varpi_j N(y_i|\mu_j, \sigma_j^2), \text{ for } i = 1, \dots, n \\ \varpi|k & \sim \text{Di}(\delta) \\ \mu_j|\sigma_j^2, k & \sim N(\mu_j|\xi, \sigma_j^2), \text{ for } j = 1, \dots, k \\ \sigma_j^2|k & \sim \text{Ga}(a, \beta), \text{ for } j = 1, \dots, k \\ \beta & \sim \text{Ga}(g, h) \\ k & \sim \text{U}_{\text{discr}}(1, k_{\max}) \end{cases} \quad (10)$$

The joint distribution that admits pdf

$$p(y, k, \varpi, \mu, \sigma^2, \beta) = f(y|k, \varpi, \mu, \sigma^2) \pi(\varpi|k) \pi(\mu|\sigma^2, k) \pi(\sigma^2|k, \beta) \pi(\beta) \pi(k).$$

The posterior  $\pi(k, \varpi, \mu, \sigma^2|y)$  can be computed with the Bayesian theorem, and factorized as

$$\pi(k, \varpi, \mu, \sigma^2, \beta|y) = \frac{p(y, \varpi, \mu, \sigma^2, \beta)}{\int p(y, k, \varpi, \mu, \sigma^2, \beta) d(k, \varpi, \mu, \sigma^2, \beta)} \quad (11)$$

$$= \pi(\mu_{1:k}, \sigma_{1:k}^2|y, k) \pi(\varpi_{1:k}|y, k) \pi(k|y) \quad (12)$$

Models (5) and (10) are the same model in the sense that posterior (11) is the marginal of the posterior (6).

**Example 12.** [Variable selection in logistic regression] Consider the 'Challenger O-ring' example from the Computer practicals. Let  $y_i$  denote the presence of a defective O-ring in the  $i$ th flight (0 for absence, and 1 for presence).

Assume that  $y_i$  can be modeled as observations generated independently from a Bernoulli distribution with parameter  $p_i$ . Here,  $p_i$  denotes the relative frequency of defective O-rings at flight  $i$ . We study if 'presence of a defective O-ring' ( $y$ ) depends on the 'temperature' ( $t$ ), or the 'pressure' ( $s$ ).

Let  $t_i$  denote the temperature (in F) in the platform, and let  $s_i$  denote the Leak check pressure (in PSI) before the  $i$ th flight. Here are some possible models of interest:

$$\begin{aligned} \mathcal{M}^I : p(t; \beta_{\mathcal{M}^I}, \mathcal{M}^I) &= \frac{\exp(\beta_0)}{1 + \exp(\beta_0)} & ; \mathcal{M}^{IV} : p(t; \beta_{\mathcal{M}^{IV}}, \mathcal{M}^{IV}) &= \frac{\exp(\beta_0 + \beta_1 t + \beta_2 s)}{1 + \exp(\beta_0 + \beta_1 t + \beta_2 s)} \\ \mathcal{M}^{II} : p(t; \beta_{\mathcal{M}^{II}}, \mathcal{M}^{II}) &= \frac{\exp(\beta_0 + \beta_1 t)}{1 + \exp(\beta_0 + \beta_1 t)} & ; \mathcal{M}^V : p(t; \beta_{\mathcal{M}^V}, \mathcal{M}^V) &= \frac{\exp(\beta_0 + \beta_1 t + \beta_2 s + \beta_3 ts)}{1 + \exp(\beta_0 + \beta_1 t + \beta_2 s + \beta_3 ts)} \\ \mathcal{M}^{III} : p(t; \beta_{\mathcal{M}^{III}}, \mathcal{M}^{III}) &= \frac{\exp(\beta_0 + \beta_2 s)}{1 + \exp(\beta_0 + \beta_2 s)} & \text{etc...} \end{aligned}$$

The Bayesian hierarchical model under consideration is:

$$\left\{ \begin{array}{l} y_i | \theta \sim f(y_i | \theta) :: \left\{ y_i | \mathcal{M}, \beta_{\mathcal{M}} \sim \text{Br} \left( y_i | \frac{\exp(x_{i,\mathcal{M}}^\top \beta_{\mathcal{M}})}{1 + \exp(x_{i,\mathcal{M}}^\top \beta_{\mathcal{M}})} \right), x_{i,\mathcal{M}}^\top \beta_{\mathcal{M}} = \sum_{j \in \mathcal{M}} x_{i,j}^\top \beta_j \quad \text{for } i = 1, \dots, n \right. \\ \theta | \phi_1 \sim \pi(\theta | \phi_1) :: \left\{ \begin{array}{l} \beta_j | \mathcal{M} \sim \text{N}(\beta_j | \mu_0, \sigma_0^2) \text{ for all } j \in \mathcal{M} \\ \mathcal{M} = \{j \in \{1, \dots, d\}, \text{ s.t. } \gamma_j = 1\} \\ \gamma_j | \varpi \sim \text{Br}(\varpi), \quad j = 1, \dots, d \end{array} \right. \\ \phi_1 | \phi_2 \sim \pi(\phi_1 | \phi_2) :: \left\{ \varpi \sim \text{Be}(a_0, b_0) \right. \end{array} \right. \quad (13)$$

where  $\theta = (\mathcal{M}, \beta_{\mathcal{M}})$ ,  $\phi_1 = \varpi$ , and  $\phi_2 = (a_0, b_0)$ . Above, in the prior we considered an extra level of uncertainty by considering  $\varpi \sim \text{Be}(a_0, b_0)$ .

Here we added an additional level of uncertainty, and set  $\varpi \sim \text{Be}(a_0, b_0)$  which creates a more diverse prior model, compared to the computer practical handout example where we had set  $\varpi = 0.5$ .

Now the joint probability distribution has pdf

$$p(y, \beta_{\mathcal{M}}, \mathcal{M}, \varpi) = \underbrace{\prod_{i=1}^n \text{Br} \left( y_i | \frac{\exp(x_{i,\mathcal{M}}^\top \beta_{\mathcal{M}})}{1 + \exp(x_{i,\mathcal{M}}^\top \beta_{\mathcal{M}})} \right)}_{f(y|\theta)} \underbrace{\prod_{j \in \mathcal{M}} \text{N}(\beta_j | \mu_0, \sigma_0^2)}_{\pi(\theta|\phi_1)} \underbrace{\prod_{j=1}^d \text{Br}(\gamma_j | \varpi) \text{Be}(\varpi | a_0, b_0)}_{\pi(\phi_1|\phi_2)}$$

**Example 13.** Robert and Reber (1998) considers an experiment under which rats are intoxicated by a substance, then treated by either a placebo or a drug. (See: <https://www.jstor.org/stable/pdf/25053027.pdf>)

**Statistical model** ( $f(y|\theta)$ ): The model associated with this experiment is a linear additive model effect: given  $x_{ij}$ ,  $y_{ij}$  and  $z_{ij}$ ,  $j$ th responses of the  $i$ th rat at the control, intoxication and treatment stages, respectively. The statistical

model was specified such as that ( $1 \leq i \leq I$ )

$$\begin{aligned} x_{i,j} &\sim N(\theta_i, \sigma_c^2) & , 1 \leq j \leq J_i^c \\ y_{i,j} &\sim N(\theta_i + \delta_i, \sigma_a^2) & , 1 \leq j \leq J_i^a, \\ z_{i,j} &\sim N(\theta_i + \delta_i + \xi_i, \sigma_t^2) & , 1 \leq j \leq J_i^t, \end{aligned}$$

where  $\theta_i$  is the average control measurement,  $\delta_i$  the average intoxication effect and  $\xi_i$  the average treatment effect for the  $i$ th rat, the variances of these measurements being constant for the control, the intoxication and the treatment effects. An additional (observed) variable is  $w_i$ , which is equal to 1 if the rat is treated with the drug, and 0 otherwise.

**Prior model  $\pi(\theta|\phi)$ :** The different individual averages are related through a common (conjugate) prior distribution,

$$\begin{aligned} \theta_i &\sim N(\mu_\theta, \sigma_\theta^2), & \delta_i &\sim N(\mu_\delta, \sigma_\delta^2), & \xi_i|w_i &\sim \begin{cases} N(\mu_P, \sigma_P^2) & , w_i = 0 \\ N(\mu_D, \sigma_D^2) & , w_i = 1 \end{cases} \\ \sigma_c &\sim \pi(\sigma_c) \propto \frac{1}{\sigma_c}, & \sigma_a &\sim \pi(\sigma_a) \propto \frac{1}{\sigma_a}, & \sigma_t &\sim \pi(\sigma_t) \propto \frac{1}{\sigma_t}, \end{aligned} \quad (14)$$

This modeling seems to describe the natural phenomenon realistically enough, in the sense of the responses  $x_{ij}$ ,  $y_{ij}$  and  $z_{ij}$ .

**Hyper-priors  $\pi(\phi|\phi_m)$ :** For the higher levels of prior ( $\pi(\phi|\phi_m)$  in Eq 2), they considered improper (Jeffrey's) hyper-priors.

$$(\mu_\theta, \sigma_\theta) \sim \pi(\mu_\theta, \sigma_\theta) \propto \frac{1}{\sigma_\theta}, \quad (\mu_\delta, \sigma_\delta) \sim \pi(\mu_\delta, \sigma_\delta) \propto \frac{1}{\sigma_\delta}, \quad (\mu_P, \sigma_P) \sim \pi(\mu_P, \sigma_P) \propto \frac{1}{\sigma_P}, \quad (15)$$

$$(\mu_D, \sigma_D) \sim \pi(\mu_D, \sigma_D) \propto \frac{1}{\sigma_D}. \quad (16)$$

The priors in lines (14), (15) and (16) are improper non-informative priors. One could have have specify proper priors, like Normal-Inverse Gamma which are conjugate, however in that case he/she should have to specify the values for the fixed hyper-parameters.

## 2 Non-identifiability issue

**Definition 14.** A parametric model for which an element of the parametrisation is redundant is said to be non-identified.

*Remark 15.* Let Bayesian model  $(f(y|\theta), \pi(\theta))$ , where  $\theta = (\theta_1, \theta_2) \in \Theta_1 \times \Theta_2$ , and assume that the parametric model does not depend on  $\theta_1$ ; i.e.  $f(y|\theta_1, \theta_2) = f(y|\theta_2)$ . Because the likelihood does not depend on  $\theta_1$  suggests that  $y$  does not provide information about  $\theta_1$  directly.

*Remark 16.* Bayesian analysis of a non-identified model is always possible if a suitable prior  $\Pi(\theta_1, \theta_2)$  on all the parameters is specified. For instance, if one specifies a priori that learning the value of  $\theta_2$  may change his belief about  $\theta_1$ , via  $\pi(\theta_1|\theta_2) \neq \pi(\theta_1)$ .

Factorize the prior distribution as  $\pi(\theta_1, \theta_2) = \pi(\theta_1|\theta_2)\pi(\theta_2)$ . Then, we have the following PDF/PMF

$$\begin{aligned} \pi(\theta_1, \theta_2|y) &\propto f(y|\theta_1, \theta_2)\pi(\theta_1, \theta_2) = f(y|\theta_2)\pi(\theta_1|\theta_2)\pi(\theta_2) \implies \\ \pi(\theta_1, \theta_2|y) &= \pi(\theta_2|y)\pi(\theta_1|\theta_2) \implies \\ \pi(\theta_2|y) &= \frac{f(y|\theta_2)\pi(\theta_2)}{\int_{\Theta_2} f(y|\theta_2)\pi(\theta_2)d\theta_2} . \\ \pi(\theta_1|y, \theta_2) &= \pi(\theta_1|\theta_2) \end{aligned} \quad (17)$$

$$\pi(\theta_1|y) = \int_{\Theta_2} \pi(\theta_1|\theta_2)\pi(\theta_2|y)d\theta_2 \quad (18)$$

$\theta_1$  is non-identifiable parameter from the data  $y$ , because  $y$  provides no direct information about  $\theta_1$ . Inference about  $\theta_1$  based on marginal posterior  $\pi(\theta_1|y)$  depends on  $y$  but the information provided about  $\theta_1$  comes indirectly through the marginal posterior of  $\theta_2$ , see (18). Equivalently, (18) implies that  $y$  provides no information about  $\theta_1$  given  $\theta_2$ .

If we a priori specify that learning the value of  $\theta_2$  does not change our belief about  $\theta_1$   $\pi(\theta_1|\theta_2) = \pi(\theta_1)$ , then (18) becomes  $\pi(\theta_1|y) = \pi(\theta_1)$  and hence data  $y$  provide no information about  $\theta_1$  at all.

**Example 17.** (Cont Example 11) It is not difficult to understand that the Bayesian model as defined in Example 7 is non-identifiable. For simplicity we focus on the Bayesian mixture of  $k = 2$  components with

$$\begin{aligned} y|\varpi, \mu, \sigma^2 &\sim f(y|\varpi, \mu, \sigma^2) := \varpi_1 N(y|\mu_1, \sigma_1^2) + \varpi_2 N(y|\mu_2, \sigma_2^2) \\ \pi(\varpi, \mu, \sigma^2) &= \underbrace{N(\mu_1|\xi, \sigma_1^2)N(\mu_2|\xi, \sigma_2^2)}_{\pi(\mu|\sigma^2)} \underbrace{IG(\sigma_1^2|\alpha, \beta)IG(\sigma_2^2|\alpha, \beta)Di(\varpi|\delta)}_{\pi(\sigma^2)} \end{aligned} \quad (19)$$

which leads to a posterior such as

$$\pi(\varpi, \mu, \sigma^2|y) \propto [\varpi_1 N(y|\mu_1, \sigma_1^2) + \varpi_2 N(y|\mu_2, \sigma_2^2)] N(\mu_1|\xi, \sigma_1^2)N(\mu_2|\xi, \sigma_2^2)IG(\sigma_1^2|\alpha, \beta)IG(\sigma_2^2|\alpha, \beta)Di(\varpi|\delta)$$

Here the parametrization is non-identifiable because the symmetry in the sampling distribution

$$\varpi_1 N(y|\mu_1, \sigma_1^2) + \varpi_2 N(y|\mu_2, \sigma_2^2) = \varpi_2 N(y|\mu_2, \sigma_2^2) + \varpi_1 N(y|\mu_1, \sigma_1^2)$$

and the naive prior in (19) produce a posterior such that

$$\pi\left(\varpi = \begin{pmatrix} \varpi_1 \\ \varpi_2 \end{pmatrix}, \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \sigma^2 = \begin{pmatrix} \sigma_1^2 \\ \sigma_2^2 \end{pmatrix} | y\right) = \pi\left(\varpi = \begin{pmatrix} \varpi_2 \\ \varpi_1 \end{pmatrix}, \mu = \begin{pmatrix} \mu_2 \\ \mu_1 \end{pmatrix}, \sigma^2 = \begin{pmatrix} \sigma_2^2 \\ \sigma_1^2 \end{pmatrix} | y\right)$$

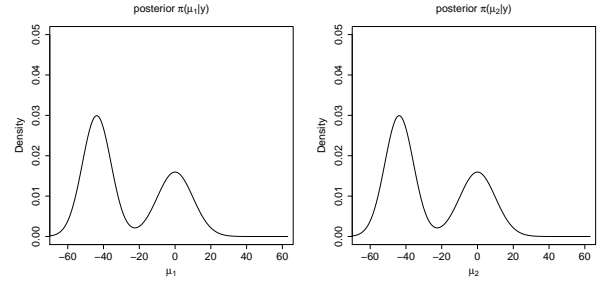
This parametrization is not meaningful for parametric inference. In Bayesian stats this can be resolved for instance by changing the prior and imposing an identifiability constrain as

$$\pi^*(\mu|\sigma^2) = \frac{N(\mu_1|\xi, \sigma_1^2)N(\mu_2|\xi, \sigma_2^2)1(\mu_1 \leq \mu_2)}{\int N(\mu_1|\xi, \sigma_1^2)N(\mu_2|\xi, \sigma_2^2)1(\mu_1 \leq \mu_2) d(\mu_1, \mu_2)} \propto \pi(\mu|\sigma^2)1(\mu_1 \leq \mu_2)$$

The non-identifiable model produces marginal posterior in Figures 1a and 1b.

Appendix  
A

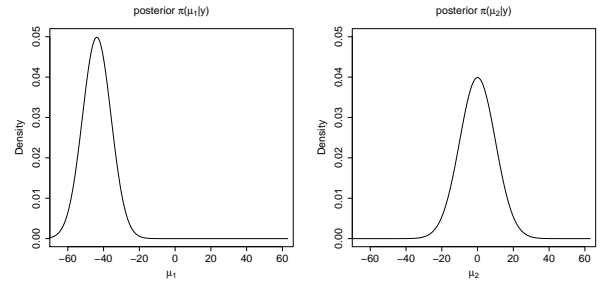
$$\begin{cases} y_i | \varpi, \mu, \sigma^2 & \sim \varpi_1 N(y | \mu_1, \sigma_1^2) + \varpi_2 N(y | \mu_2, \sigma_2^2) \\ \varpi & \sim \text{Di}(\delta) \\ \mu | \sigma^2 & \sim N(\mu_1 | \xi, \sigma_1^2) N(\mu_2 | \xi, \sigma_2^2) \\ \sigma^2 & \sim \text{Ga}(\sigma_1^2 | \alpha, \beta) \text{Ga}(\sigma_2^2 | \alpha, \beta) \end{cases} \quad (20)$$



(a) Marginal posterior  $\pi(\mu_1 | y)$  (b) Marginal posterior  $\pi(\mu_2 | y)$

After non-identifiability is resolved, the identifiable model produces marginal posterior in Figures 1c and 1d.

$$\begin{cases} y_i | \varpi, \mu, \sigma^2 & \sim \varpi_1 N(y | \mu_1, \sigma_1^2) + \varpi_2 N(y | \mu_2, \sigma_2^2) \\ \varpi & \sim \text{Di}(\delta, \dots, \delta) \\ \mu | \sigma^2 & \sim N(\mu_1 | \xi, \sigma_1^2) N(\mu_2 | \xi, \sigma_2^2) 1(\mu_1 \leq \mu_2) \\ \sigma^2 & \sim \text{Ga}(\sigma_1^2 | \alpha, \beta) \text{Ga}(\sigma_2^2 | \alpha, \beta) \end{cases} \quad (21)$$



(c) Marginal posterior  $\pi(\mu_1 | y)$  (d) Marginal posterior  $\pi(\mu_2 | y)$

### 3 Augmentation

**Definition 18.** Augmentation technique on a Bayesian model ( $y \sim F(y | \theta)$ ,  $\theta \sim \pi(\theta)$ ) is when the vector  $(y, \theta)$  is imputed by auxiliary random quantities  $x \sim Q(\cdot)$  where the distribution  $Q(\cdot)$  is such that the joint distribution of  $(y, \theta, x)$  admits the joint distribution of  $(y, \theta)$  as its marginal, e.g.

$$p(y, \theta) = \int p(y, \theta, x) dx$$

**Remark 19.** According to Definition 18, the posterior of  $\theta$  given the observables  $y$  remains the same; e.g.

$$\pi(\theta | y) = \int \pi(\theta, x | y) dx = \frac{\int p(y, \theta, x) dx}{\int p(y, \theta, x) d\theta dx}$$

**Example 20.** (Cont. Examples 7 and 11) Imputation of the hierarchical model (10) by considering the latent variables  $z_i \sim f(z_i) = \varpi_{z_i}$  for  $z_i \in \{1, \dots, k\}$  leads to the augmented hierarchical model (5).

**Remark 21.** Software JAGS, aiming at facilitating Bayesian computations, requires as inputs only Bayesian hierarchical models where the dimension of the unknown quantities is fixed (not random). Consequently many hierarchical models for model comparison such as (5), (10), and 13 cannot be used directly. This limitation can be addressed via augmentation, i.e. the imputation of the hierarchical model with suitable random variables such that the resulting joint posterior distribution admits the posterior distribution as its marginal.

**Example 22.** (Cont Example 12) Notice that the length of the join unknown parameter vector

$$(\mathcal{M}, \beta_{\mathcal{M}}, \gamma_{1:d}, \varpi)$$



is not fixed but random because  $\mathcal{M} = \{j \in \{1, \dots, d\}, \text{ s.t. } \gamma_j = 1\}$  is random/unknown. Current versions of JAGS work only when the join unknown parameter vector of the Bayesian hierarchical model is fixed. Augment the hierarchical model properly so that it can be used with RJAGS.

**Solution.** We augment model (13) with auxiliary random quantities such that  $\beta_j | \mathcal{M} \sim \text{N}(\beta_j | \mu_0, \sigma_0^2)$  for all  $j \notin \mathcal{M}$ , i.e.

$$\left\{ \begin{array}{l} \left\{ y_i | \mathcal{M}, \beta_{\mathcal{M}} \sim \text{Br} \left( y_i | \frac{\exp(x_{i,\mathcal{M}}^\top \beta_{\mathcal{M}})}{1 + \exp(x_{i,\mathcal{M}}^\top \beta_{\mathcal{M}})} \right), x_{i,\mathcal{M}}^\top \beta_{\mathcal{M}} = \sum_{j \in \mathcal{M}} x_{i,j}^\top \beta_j \quad \text{for } i = 1, \dots, n \right. \\ \left. \left\{ \beta_j | \mathcal{M} \sim \text{N}(\beta_j | \mu_0, \sigma_0^2) \quad \text{for all } j \notin \mathcal{M} \right. \right. \\ \left. \left\{ \beta_j | \mathcal{M} \sim \text{N}(\beta_j | \mu_0, \sigma_0^2) \quad \text{for all } j \in \mathcal{M} \right. \right. \\ \left. \left\{ \begin{array}{l} \mathcal{M} = \{j \in \{1, \dots, d\}, \text{ s.t. } \gamma_j = 1\} \\ \gamma_j | \varpi \sim \text{Br}(\varpi), \quad j = 1, \dots, d \end{array} \right. \right. \\ \left. \left\{ \varpi \sim \text{Be}(a_0, b_0) \right. \right. \end{array} \right.$$

then

$$p(y, \beta_{\mathcal{M}} \beta_{-\mathcal{M}}, \mathcal{M}, \varpi) = p(y, \beta_{\mathcal{M}}, \mathcal{M}, \varpi) \prod_{j \notin \mathcal{M}} \text{N}(\beta_j | \mu_0, \sigma_0^2)$$

and

$$\int p(y, \beta_{\mathcal{M}} \beta_{-\mathcal{M}}, \mathcal{M}, \varpi) d\beta_{-\mathcal{M}} = p(y, \beta_{\mathcal{M}}, \mathcal{M}, \varpi) \int \prod_{j \notin \mathcal{M}} \text{N}(\beta_j | \mu_0, \sigma_0^2) d\beta_{-\mathcal{M}} = p(y, \beta_{\mathcal{M}}, \mathcal{M}, \varpi)$$

## A Appendix: Bayesian Normal Mixture model

### A.1 Case where $k$ is a fixed known parameter

The complete R / RJAGS code is available from

- [http://htmlpreview.github.io/?https://github.com/georgios-stats/Bayesian\\_Statistics\\_Michaelmas\\_2021/blob/main/Computer\\_practical/Normal\\_Mixture\\_model/Bayesian\\_Normal\\_Mixture\\_Model.nb.html](http://htmlpreview.github.io/?https://github.com/georgios-stats/Bayesian_Statistics_Michaelmas_2021/blob/main/Computer_practical/Normal_Mixture_model/Bayesian_Normal_Mixture_Model.nb.html)
- [https://raw.githubusercontent.com/georgios-stats/Bayesian\\_Statistics\\_Michaelmas\\_2021/main/Computer\\_practical/Normal\\_Mixture\\_model/enz.dat](https://raw.githubusercontent.com/georgios-stats/Bayesian_Statistics_Michaelmas_2021/main/Computer_practical/Normal_Mixture_model/enz.dat)

For simplification, consider the hierarchical model where the number of groups  $k$  is a fixed/known quantity.

$$\left\{ \begin{array}{ll} y_i | z_i, k, \varpi, \mu, \sigma^2 & \sim f_{z_i}(y_i | \varpi_{z_i}, \mu_{z_i}, \sigma_{z_i}^2) \\ & \text{for } i = 1, \dots, n \\ z_i | k & \sim f(z_i) := \varpi_{z_i} \\ & \text{for } i = 1, \dots, n \\ \varpi | k & \sim \text{Di}(\delta, \dots, \delta) \\ \mu_j | \sigma_j^2, k & \sim N(\xi, \sigma_j^2 / \kappa) \\ & \text{for } j = 1, \dots, k \\ \sigma_j^2 | \beta, k & \sim \text{IG}(\alpha, \beta) \\ & \text{for } j = 1, \dots, k \\ \beta & \sim \text{Ga}(g, h) \end{array} \right. \quad (22)$$

*Note 23.* The RJAGS code for the analysis of the data set in the Example 7 is provided in Algorithm 1, however here the number of groups  $k$  is fixed/known quantity and not random.

---

**Algorithm 1** [Bayesian Normal Mixture model] RJAGS script for the analysis of the data set in the Example 7 .

---

```
rm(list=ls())
# Load rjags
library("rjags")
# define the Bayesian hierarchical model in JAGS syntax
jags_model <- "
  model{
    #sampling distribution
    for (i in 1:n) {
      y[i] ~ dnorm(mu[zeta[i]], tau[zeta[i]])
      zeta[i] ~ dcat(varpi[])
    }
    #prior distribution
    varpi ~ ddirich(delta)
    for (i in 1:k) {
      mu[i] ~ dnorm(xi, kappa)
      tau[i] ~ dgamma(alpha, beta)
      sigma[i] <- 1/sqrt(tau[i])
    }
    beta ~ dgamma(g, h)
  }
"
```

---

**Example 24.** (Cont. Example 7) As we will see in term 2, although (5) and 10 are equivalent in the sense that they produce the same inference, the expended hierarchical model (6) is computational convenient compared to hierarchical model (11) in the sense that its priors are conditional conjugate.

For simplicity assume that the number of groups is known and fixed to  $k$ . The full conditional posteriors in model (6) are:

$$\begin{aligned} w|y... &\sim \text{Di}(\delta + n_1, \dots, \delta + n_k); \text{ where } n_j = \sum_{i=1}^n 1(z_i = j) \\ \mu_j|y... &\sim \text{N}\left(\frac{\sum_{i:z_i=j} y_i - \xi\kappa}{n_j + \kappa}, \frac{\sigma_j^2}{n_j + \kappa}\right), \text{ for } j = 1, \dots, k \\ \sigma_j^2|y... &\sim \text{IG}\left(a + \frac{n_j}{2}, \beta + \frac{1}{2} \sum_{i:z_i=j} (y_i - \mu_j)^2\right), \text{ for } j = 1, \dots, k \\ z_i|y... &\sim \pi(z_i = j|y...) = \frac{\frac{w_j}{\sigma_j} \exp\left(-\frac{1}{2} \frac{(y_i - \mu_j)^2}{\sigma_j^2}\right)}{\sum_{j'=1}^k \frac{w_{j'}}{\sigma_{j'}} \exp\left(-\frac{1}{2} \frac{(y_i - \mu_{j'})^2}{\sigma_{j'}^2}\right)}; \text{ for } i = 1, \dots, n \\ \beta|y... &\sim \text{Ga}\left(g + k\alpha, h + \sum_{j=1}^k \sigma_j^2\right) \end{aligned}$$

which can be used in Monte Carlo integration.

Model (11) does not produce full conditional posteriors of standard form, due to the summation in the likelihood.

## A.2 Case where $k$ is a random/unknown parameter

The complete R / RJAGS code is available from

- [http://htmlpreview.github.io/?https://github.com/georgios-stats/Bayesian\\_Statistics\\_Michaelmas\\_2021/blob/main/Computer\\_practical/Normal\\_Mixture\\_Model\\_unknown\\_number\\_of\\_components/Bayesian\\_Normal\\_Mixture\\_Model\\_unknown\\_number\\_of\\_components.nb.html](http://htmlpreview.github.io/?https://github.com/georgios-stats/Bayesian_Statistics_Michaelmas_2021/blob/main/Computer_practical/Normal_Mixture_Model_unknown_number_of_components/Bayesian_Normal_Mixture_Model_unknown_number_of_components.nb.html)
- [https://raw.githubusercontent.com/georgios-stats/Bayesian\\_Statistics\\_Michaelmas\\_2021/main/Computer\\_practical/Normal\\_Mixture\\_Model\\_unknown\\_number\\_of\\_components/enz.dat](https://raw.githubusercontent.com/georgios-stats/Bayesian_Statistics_Michaelmas_2021/main/Computer_practical/Normal_Mixture_Model_unknown_number_of_components/enz.dat)

*Note 25.* Consider the hierarchical model 5.

*Note 26.* Notice that the length of the join unknown parameter vector

$$(z_{1:n}, k, \varpi_{1:k}, \mu_{1:k}, \sigma_{1:k}^2)$$

is not fixed but random because  $k$  is random/unknown. Current versions of JAGS work only when the join unknown parameter vector of the Bayesian hierarchical model is fixed.

*Note 27.* It is given that if  $c_j \sim \text{Ga}(\delta_j, 1)$  for  $j = 1, \dots, k$  then

$$\varpi_{1:k} = \left( \frac{c_1}{\sum_{j=1}^k c_j}, \dots, \frac{c_k}{\sum_{j=1}^k c_j} \right) \sim \text{Di}(\delta_1, \dots, \delta_k).$$

220 *Note 28.* To make JAGS work on our problem we augment  $(z_{1:n}, k, \varpi_{1:k}, \mu_{1:k}, \sigma_{1:k}^2)$  with additional random variables

$$\begin{aligned}
 221 \quad & c_j | k \sim \text{Ga}(\delta, 1) && \text{for } j = k+1, \dots, k_{\max} \\
 222 \quad & \mu_j | \sigma_j^2, k \sim \text{N}(\xi, \sigma_j^2 / \kappa) && \text{for } j = k+1, \dots, k_{\max} \\
 223 \quad & \sigma_j^2 | \beta, k \sim \text{IG}(\alpha, \beta) && \text{for } j = k+1, \dots, k_{\max}
 \end{aligned}$$

224 Precisely, we consider the augmented Bayesian hierarchical model

225 *Note 29.* Consider augmentation where (5) is augmented by  $\mu_j, \sigma_j^2 | k \sim \text{N}(\xi, \sigma_j^2 / \kappa) \text{IG}(\alpha, \beta)$  for  $j = k+1, \dots, k_{\max}$ ,  
 226 and  $c_j \sim \text{Ga}(\delta, 1)$  for  $j = k+1, \dots, k_{\max}$

$$\left\{ \begin{array}{ll}
 y_i | z_i, k, \varpi, \mu, \sigma^2 & \sim f_{z_i}(y_i | \varpi_{z_i}, \mu_{z_i}, \sigma_{z_i}^2) & \text{for } i = 1, \dots, n \\
 z_i | k & \sim f(z_i) := \varpi_{z_i} & \text{for } i = 1, \dots, n \\
 \varpi_{1:k} & = \left( \frac{c_1}{\sum_{j=1}^k c_j}, \dots, \frac{c_k}{\sum_{j=1}^k c_j} \right) \\
 \\ 
 c_j | k & \sim \text{Ga}(\delta, 1) & \text{for } j = 1, \dots, k_{\max} \\
 \mu_j | \sigma_j^2, k & \sim \text{N}(\xi, \sigma_j^2 / \kappa) & \text{for } j = k+1, \dots, k_{\max} \\
 \sigma_j^2 | \beta, k & \sim \text{IG}(\alpha, \beta) & \text{for } j = k+1, \dots, k_{\max} \\
 \\ 
 \mu_j | \sigma_j^2, k & \sim \text{N}(\xi, \sigma_j^2 / \kappa) & \text{for } j = 1, \dots, k \\
 \sigma_j^2 | \beta, k & \sim \text{IG}(\alpha, \beta) & \text{for } j = 1, \dots, k \\
 \beta & \sim \text{Ga}(g, h) \\
 k & \sim \text{U}_{\text{discr}}(1, k_{\max})
 \end{array} \right. \quad (23)$$

228 then

$$229 \quad p(y_{1:n}, z_{1:n}, k, \varpi_{1:k}, \mu_{1:k_{\max}}, \sigma_{1:k_{\max}}^2, c_{k+1:1:k_{\max}}) = p(y_{1:n}, z_{1:n}, k, \varpi_{1:k}, \mu_{1:k}, \sigma_{1:k}^2) \prod_{j=k+1}^{k_{\max}} \text{N}(\mu_j | \xi, \sigma_j^2 / \kappa) \text{Ga}(\sigma_k^2 | \alpha, \beta) \text{Ga}(c_j | \delta, 1)$$

230 and

$$231 \quad \int p(y_{1:n}, z_{1:n}, k, \varpi_{1:k}, \mu_{1:k_{\max}}, \sigma_{1:k_{\max}}^2, c_{k+1:1:k_{\max}}) d(\mu_{k+1:k_{\max}}, \sigma_{k+1:k_{\max}}^2, c_{k+1:1:k_{\max}}) = p(y_{1:n}, z_{1:n}, k, \varpi_{1:k}, \mu_{1:k}, \sigma_{1:k}^2)$$

232 The RJAGS code for the analysis of the data set in the Example 12 is provided in Algorithm 2.

---

**Algorithm 2** [Bayesian Normal Mixture model] RJAGS script for the analysis of the data set in the Example 7 .

---

```
rm(list=ls())
# Load rjags
library("rjags")
# define the Bayesian hierarchical model in JAGS syntax
jags_model <- "
model{
  # sampling distribution
  for (i in 1:n) {
    y[i] ~ dnorm(mu[zeta[i]], tau[zeta[i]])
    zeta[i] ~ dcat( varpi ) }
  # within model parameter priors
  for (i in 1:length(pk) ) {
    x[i] ~ dgamma(delta[i],1)
    ind[i] <- (i<=k)
  }
  varpi <- (x*ind)/sum(x*ind)
  for (i in 1:length(pk) ) {
    mu[i] ~ dnorm(xi,kappa)
    tau[i] ~ dgamma(alpha,beta)
    sigma[i] <- 1/sqrt(tau[i])
  }
  beta ~ dgamma(g,h)
  k~dcat(pk)
}
"
```

---

## B Appendix: Bayesian Variable Selection in Logistic regression model

The complete R / RJAGS code is available from

- [http://htmlpreview.github.io/?https://github.com/georgios-stats/Bayesian\\_Statistics\\_Michaelmas\\_2021/blob/main/Computer\\_practical/Bernoulli\\_regression\\_model\\_variable\\_selection/Bernoulli\\_Regression\\_Model\\_VS\\_full.nb.html](http://htmlpreview.github.io/?https://github.com/georgios-stats/Bayesian_Statistics_Michaelmas_2021/blob/main/Computer_practical/Bernoulli_regression_model_variable_selection/Bernoulli_Regression_Model_VS_full.nb.html)

*Note 30.* Notice that the length of the join unknown parameter vector

$$(\mathcal{M}, \beta_{\mathcal{M}}, \gamma_{1:d}, \varpi)$$

is not fixed but random because  $\mathcal{M} = \{j \in \{1, \dots, d\}, \text{ s.t. } \gamma_j = 1\}$  is random/unknown. Current versions of JAGS work only when the join unknown parameter vector of the Bayesian hierarchical model is fixed.

*Note 31.* To make JAGS work on our problem we 'cheat' by augmenting  $(\mathcal{M}, \beta_{\mathcal{M}}, \gamma_{1:d}, \varpi)$  with additional random variables  $\beta_j | j \notin \mathcal{M} \sim N(\beta_j | \mu_0, \sigma_0^2)$ . Precisely, we consider the augmented Bayesian hierarchical model

$$\left\{ \begin{array}{ll} y_i | \mathcal{M}, \beta_{\mathcal{M}} & \sim \text{Bernoulli}(p(x_i; \mathcal{M}, \beta_{\mathcal{M}})), \quad \text{for } i = 1, \dots, n \\ p(x_i; \mathcal{M}, \beta_{\mathcal{M}}) & = \frac{\exp(x_{i,\mathcal{M}}^\top \beta_{\mathcal{M}})}{1 + \exp(x_{i,\mathcal{M}}^\top \beta_{\mathcal{M}})}; \quad \text{where } x_{i,\mathcal{M}}^\top \beta_{\mathcal{M}} = \sum_{j \in \mathcal{M}} x_{i,j} \beta_j \\ \beta_j | j \notin \mathcal{M} & \sim N(\beta_j | \mu_0, \sigma_0^2), \quad j \notin \mathcal{M} \\ \beta_j | j \in \mathcal{M} & \sim N(\beta_j | \mu_0, \sigma_0^2), \quad j \in \mathcal{M} \\ \mathcal{M} & = \{j \in \{1, \dots, d\}, \text{ s.t. } \gamma_j = 1\} \\ \gamma_j & \sim \text{Bernoulli}(\varpi), \quad j = 1, \dots, d \\ \varpi & \sim \text{Be}(a_0, b_0) \end{array} \right.$$

Now we work on the augmented vector  $(\mathcal{M}, \beta_{\mathcal{M}}, \beta_{\mathcal{M}^c}, \gamma_{1:d}, \varpi)$ . The augmented hierarchical model admits the hierarchical model as its marginal because

$$\pi(\mathcal{M}, \beta_{\mathcal{M}}, \gamma_{1:d}, \varpi | y) = \int \pi(\mathcal{M}, \beta_{\mathcal{M}}, \beta_{\mathcal{M}^c}, \gamma_{1:d}, \varpi | y) d\beta_{\mathcal{M}^c}$$

*Note 32.* The RJAGS code for the analysis of the data set in the Example 12 is provided in Algorithm 3.

---

**Algorithm 3** [Bayesian Mixed Effect Model] RJAGS script for the analysis of the data set in the Example 7 .

---

```
rm(list=ls())
# Load rjags
library("rjags")
# define the Bayesian hierarchical model in JAGS syntax
hierarhicalmodel<-"
model {
  # sampling distribution
  for (i in 1:n) {
    eta[i] <- inprod(X[i,],beta*ind)
    mean[i] <- exp( eta[i] ) / ( 1 + exp( eta[i] ) )
    y[i] ~ dbern(mean[i])
  }
  # within model prior + augmentation
  for ( j in 1:dmax ) {
    beta[j] ~ dnorm( 0 , 0.1 )
  }
  # marginal model prior
  ind[1] <- 1
  for (j in 2:dmax) {
    ind[j] ~ dbern( pp )
  }
  # hyper-prior
  pp ~ dbeta(1.0,1.0)
}
"
```

---

## C Appendix: Random effect model

The complete R / RJAGS code is available from

- [https://github.com/georgios-stats/Bayesian\\_Statistics\\_Michaelmas\\_2021/blob/main/Lecture\\_handouts/Rscripts/Hierarchical\\_bayes/HierarchicalBayesPharmaceutical.R](https://github.com/georgios-stats/Bayesian_Statistics_Michaelmas_2021/blob/main/Lecture_handouts/Rscripts/Hierarchical_bayes/HierarchicalBayesPharmaceutical.R)

*Note 33.* The RJAGS code for the analysis of the data set in the Example 13 is provided in Algorithm 4.

---

**Algorithm 4** [Random effect model] RJAGS script for the analysis of the data set in the Example 13.

---

```
rm(list=ls())
# Load rjags
library("rjags")
# define the Bayesian hierarchical model in JAGS syntax
hierarhicalmodel <- "
  model {
    for ( i in 1 : I ) {
      for ( j in 1 : J ) {
        x[i,j] ~ dnorm( theta[i] , tau_c )
        y[i,j] ~ dnorm( theta[i] + delta[i] , tau_a )
        z[i,j] ~ dnorm( theta[i] + delta[i] + xi[i] , tau_t )
      }
      theta[i] ~ dnorm( mu_theta , tau_theta )
      delta[i] ~ dnorm( mu_delta , tau_delta )
      w_ind[i] <- ifelse(w[i] == 0, 0, 1)
      xi[i] <- w_ind[i]*x_d +(1-w_ind[i])*x_p
    }
    sig2_c <- 1/tau_c tau_c ~ dgamma(0.01 , 0.01)
    sig2_a = 1/tau_a tau_a ~ dgamma(0.01 , 0.01)
    sig2_t <- 1/tau_t tau_t ~ dgamma(0.01 , 0.01)
    mu_theta ~ dnorm( 0 , 0.001 )
    sig2_theta = 1/tau_theta
    tau_theta ~ dgamma(0.01 , 0.01)
    mu_delta ~ dnorm( 0 , 0.001 )
    sig2_delta <- 1/tau_delta
    tau_delta ~ dgamma(0.01 , 0.01)
    x_p ~ dnorm( mu_p , tau_p )
    x_d ~ dnorm( mu_d , tau_d )
    mu_p ~ dnorm( 0 , 0.001 )
    sig2_p <- 1/tau_p
    tau_p ~ dgamma(0.01 , 0.01)
    mu_d ~ dnorm( 0 , 0.001 )
    sig2_d <- 1/tau_d
    tau_d ~ dgamma(0.01 , 0.01)
  }
"
```

---

**Example 34.** (Cont...) You may use

$$-\frac{1}{2} \sum_{i=1}^n \frac{(x - \mu_i)^2}{\sigma_i^2} = -\frac{1}{2} \frac{(x - \hat{\mu})^2}{\hat{\sigma}^2} + C; \quad \hat{\sigma}^2 = \left( \sum_{i=1}^n \frac{1}{\sigma_i^2} \right)^{-1}; \quad \hat{\mu} = \hat{\sigma}^2 \left( \sum_{i=1}^n \frac{\mu_i}{\sigma_i^2} \right); \quad C = \frac{1}{2} \frac{\left( \sum_{i=1}^n \frac{\mu_i}{\sigma_i^2} \right)^2}{\sum_{i=1}^n \frac{1}{\sigma_i^2}} - \frac{1}{2} \sum_{i=1}^n \frac{\mu_i^2}{\sigma_i^2}$$



256 The joint posterior pdf of  $\vartheta = (\theta_{1:I}, \delta_{1:I}, \xi_{1:I}, \sigma_c^2, \sigma_a^2, \sigma_t^2, \sigma_\theta^2, \sigma_\delta^2, \sigma_P^2, \sigma_D^2, \mu_\theta, \mu_\delta, \mu_P, \mu_D)$  given obs.  $x, y, z$  is

$$\begin{aligned}
 257 \quad \pi(\vartheta|x, y, z) &\propto \prod_{i=1}^I \left[ \exp \left( -\frac{(\theta_i - \mu_\theta)^2}{2\sigma_\theta^2} - \frac{(\delta_i - \mu_\delta)^2}{2\sigma_\delta^2} \right) \prod_{j=1}^{J_i^c} \exp \left( -\frac{(x_{i,j} - \theta_i)^2}{2\sigma_c^2} \right) \times \prod_{j=1}^{J_i^a} \exp \left( -\frac{(y_{i,j} - \theta_i - \delta_i)^2}{2\sigma_a^2} \right) \right. \\
 258 \quad &\times \prod_{j=1}^{J_i^t} \exp \left( -\frac{(z_{i,j} - \theta_i - \delta_i - \xi_i)^2}{2\sigma_t^2} \right) \times \prod_{w_i=0} \exp \left( -\frac{(\xi_i - \mu_P)^2}{2\sigma_P^2} \right) \prod_{w_i=0} \exp \left( -\frac{(\xi_i - \mu_D)^2}{2\sigma_D^2} \right) \Big] \\
 259 \quad &\times \sigma_c^{-\sum_i J_i^c - 1} \sigma_a^{-\sum_i J_i^a - 1} \sigma_t^{-\sum_i J_i^t - 1} \sigma_\theta^{I-1} \sigma_\delta^{I-1} \sigma_P^{I_D-1} \sigma_D^{I_P-1}.
 \end{aligned}$$

260 The joint posterior distributions is not of standard form, and its pdf is intractable. However the full conditionals are of  
 261 standard form. For instance, the full conditional posterior distribution density

$$\begin{aligned}
 262 \quad \pi(\delta_{1:I}|x_{\text{all}}, y_{\text{all}}, z_{\text{all}}, \theta_{1:I}, \xi_{1:I}, \sigma_c^2, \sigma_a^2, \sigma_t^2, \sigma_\theta^2, \sigma_\delta^2, \sigma_P^2, \sigma_D^2, \mu_\theta, \mu_\delta, \mu_P, \mu_D) \\
 263 \quad &\propto \prod_{i=1}^I \left[ \exp \left( -\frac{(\delta_i - \mu_\delta)^2}{2\sigma_\delta^2} \right) \times \prod_{j=1}^{J_i^a} \exp \left( -\frac{(y_{i,j} - \theta_i - \delta_i)^2}{2\sigma_a^2} \right) \times \prod_{j=1}^{J_i^t} \exp \left( -\frac{(z_{i,j} - \theta_i - \delta_i - \xi_i)^2}{2\sigma_t^2} \right) \right] \\
 264 \quad &\propto \prod_{i=1}^I \left[ \exp \left( -\frac{(\delta_i - \mu_\delta)^2}{2\sigma_\delta^2} - \sum_{j=1}^{J_i^a} \frac{(\delta_i - (y_{i,j} - \theta_i))^2}{2\sigma_a^2} - \sum_{j=1}^{J_i^t} \frac{(\delta_i - (z_{i,j} - \theta_i - \xi_i))^2}{2\sigma_t^2} \right) \right] \\
 265 \quad &\propto \prod_{i=1}^I \left[ \exp \left( -\frac{(\delta_i - \mu_{\delta,i}^*)^2}{2(\sigma_{\delta,i}^*)^2} + \text{const...} \right) \right] \propto \prod_{i=1}^I \left[ \exp \left( -\frac{(\delta_i - \mu_{\delta,i}^*)^2}{2(\sigma_{\delta,i}^*)^2} + \text{const...} \right) \right] \\
 266 \quad &\propto \prod_{i=1}^I \text{N}(\delta_i | \mu_{\delta,i}^*, (\sigma_{\delta,i}^*)^2)
 \end{aligned}$$

267 with

$$268 \quad \delta_i | \text{rest}, \dots \stackrel{\text{ind}}{\sim} \text{N}(\mu_{\delta,i}^*, (\sigma_{\delta,i}^*)^2), \forall i = 1, \dots, n$$

269 where

$$270 \quad (\sigma_{\delta,i}^*)^2 = \left( \frac{1}{\sigma_\delta^2} + \frac{1}{\sigma_a^2} J_i^a + \frac{1}{\sigma_t^2} J_i^t \right)^{-1}; \quad \mu_{\delta,i}^* = (\sigma_{\delta,i}^*)^2 \left( \frac{\mu_\delta}{\sigma_\delta^2} + \frac{\sum_{j=1}^{J_i^a} y_{i,j} - J_i^a \theta_i}{\sigma_a^2} + \frac{\sum_{j=1}^{J_i^t} y_{i,j} - J_i^t \theta_i - J_i^t \xi_i}{\sigma_t^2} \right)$$

271 Notice that  $\delta_i$  are a postriori independent given all the resp unknown parameters  
 272  $(\theta_{1:I}, \xi_{1:I}, \sigma_c^2, \sigma_a^2, \sigma_t^2, \sigma_\theta^2, \sigma_\delta^2, \sigma_P^2, \sigma_D^2, \mu_\theta, \mu_\delta, \mu_P, \mu_D)$ . Notice that the prior  $\delta_i \sim \text{N}(\mu_\delta, \sigma_\delta^2)$  in Example 13 is  
 273 conditional conjugate prior of  $\delta_i$ .

274 Try to compute the rest

$$\begin{aligned}
 275 \quad &\pi(\theta_{1:I} | \text{rest}, \dots) \sim?; \quad \pi(\sigma_t^2 | \text{rest}, \dots) \sim?; \quad \pi(\sigma_c^2 | \text{rest}, \dots) \sim?; \quad \pi(\sigma_a^2 | \text{rest}, \dots) \sim?; \\
 276 \quad &\pi(\xi_{1:I} | \text{rest}, \dots) \sim?; \quad \pi(\sigma_\theta^2 | \text{rest}, \dots) \sim?; \quad \pi(\sigma_\delta^2 | \text{rest}, \dots) \sim?; \quad \pi(\sigma_P^2 | \text{rest}, \dots) \sim?, \text{ etc...}
 \end{aligned}$$

277 See the solutions in: Robert, C. P., & Reber, A. (1998) from the link ([https://www.jstor.org/stable/pdf/](https://www.jstor.org/stable/pdf/25053027.pdf)  
 278 25053027.pdf).