

TECHNICAL REPORT 55
MARCH 2008

Posterior Computation for Hierarchical Dirichlet Process Mixture Models:
Application to Genetic Association Studies of Quantitative Traits in the the
Presence of Population Strati cation

Nicholas M. Pajewski¹ and Purushottam W. Laud
Division of Biostatistics
Department of Population Health
Medical College of Wisconsin

defined as follows.

$$\begin{aligned}
 L(Y_i | \mu_i) &= \frac{1}{\sqrt{2}} \exp \left[-\frac{1}{2} (Y_i - \mu_i)^2 \right] \\
 \mu_i &= \mu_{0i} + \sum_{l=1}^L X_{li} \mu_{li} \\
 L(W_{li}; V_{li} | \mu_{li}) &= \frac{2^{W_{li}} e^{-\mu_{li}(2V_{li} + W_{li})}}{(1 + e^{-\mu_{li}})^2} \quad i = 1; \dots; N \quad l = 1; \dots; L \\
 \mu_{0i}; \mu_{1i}; \dots; \mu_{Li} | G &\stackrel{i.i.d}{\sim} G \quad i = 1; \dots; N \\
 G | G_0 &\sim \text{DP}(G; G_0) \\
 G_0 &= \prod_{l=1}^L N(\mu_{0l}; \sigma_{0l}^2) \prod_{l=1}^L N(\mu_{li}; \sigma_{li}^2) \\
 \mu_{li} | H &\stackrel{i.i.d}{\sim} H \quad l = 1; \dots; L \\
 H | H_0 &\sim \text{DP}(H; H_0) \\
 H_0 &= \delta_{(0,0)}(\cdot) + (1 - \alpha) \text{MVN}_2(M; T) \\
 &\sim \text{Beta}(c_1; d_1) \\
 &\sim \text{Gamma}(\nu_1; \tau_1) \\
 G &\sim \text{Gamma}(\nu_2; \tau_2) \quad \text{and} \quad H \sim \text{Gamma}(\nu_3; \tau_3)
 \end{aligned}$$

Note: Throughout the document, we use the following parametrization of gamma density, $X \sim \text{Gamma}(\nu; \tau)$,

$$f(x) \propto x^{\nu-1} e^{-\tau x}$$

In the above formulation, $\mu_{li} = \text{logit}(\mu_{li})$ where μ_{li} presents the reference allele frequency for the i^{th} individual at the l^{th} SNP. $\delta_{(0,0)}(\cdot)$ represents a Dirac delta function indicating a point mass at (0;0). In addition, $N(x; \mu; \sigma^2)$ denotes a normal density with mean μ and precision σ^2 and $\text{MVN}_p(x; M; T)$ represents a p-dimensional multivariate normal with mean vector M and precision matrix T. For each of the Dirichlet Processes, we have assumed gamma priors for the scalar mass parameters G and H following ?; alternatively they could be taken as to be fixed constants. Figure 1 displays the model as a directed acyclic graph (DAG).



$0i$

Y_i

Step 1a: Perform the following proposal step for R iterations. For $i = 1; 2; \dots; N$; propose a new distinct atom membership (s_i^*) for the i^{th} observation. The approach of ? uses the conditional prior as a proposal distribution for s_i^* . Let $s_{(-i)}$ denote the set of all configuration indicators minus s_i , and let $n^{(-i)}$

Although the above log target density does not take a standard distributional form, the density is log-concave, and so a new value for μ_j^* can be sampled using Adaptive-Rejection sampling (?).

STEP 2: Update for μ_j

In order to update each μ_j , we employed the Blocked Gibbs Sampler of ?. The Blocked Gibbs Sampler is based on the stick-breaking representation of the Dirichlet Process, discussed in the work of ?. Although the stick-breaking representation of the DP involves an infinite sum of discrete points, in actual implementation, the Blocked Gibbs Sampler utilizes a finite approximation, imposing a limit F_L to the number of distinct atoms amongst the μ_j . Denote this collection of distinct points as $\mu^* = \mu_1^*; \dots; \mu_{F_L}^*$. ? show that even for large sample sizes, a limit of $F_L = 150$ provides a suitable approximation to the Dirichlet Process. Because of the point mass mixture construction in H_0 , without a loss of generality, we can include the additional distinct point μ_0^* to represent the cluster denoting no effect (i.e. $\mu_1 = 0$ and $\mu_2 = 0$) with associated model weight π_0 . Similar to the configuration representation for μ_j , define the pointers z_l where $z_l = j$ if and only if $\mu_l = \mu_j^*$ for $j = 0; 1; 2; \dots; F_L$. Then define m_j as the number of z_l currently equal to j .

Step 2a: For $j = 1; 2; \dots; F_L$; update μ_j^* . Note, because μ_0^* represents the null effect cluster, its value need not be updated. If $m_j = 0$, then $\mu_j^* \sim H_0$. Else draw $\mu_j^* \sim MVN_2(M^*; T^*)$ where

$$T^* = G_j G_j + T$$

$$M^* = (T^*)^{-1} G_j (Y - B_0 - X^{(-j)} \mu^{(-j)}) + T M$$

Y denotes a $n \times 1$ column vector of the quantitative traits Y_i . Similarly, B_0 represents a $n \times 1$ column vector where the i th element is μ_0^* . $X^{(-j)}$ is a $n \times (F_L - 1)$ matrix where the i th row is $(\mu_1^*, \dots, \mu_{j-1}^*, \mu_{j+1}^*, \dots, \mu_{F_L}^*)$.

Step 2b: For $l = 1; 2; \dots; L$; independently sample z_l where,

$$P(z_l = 0) \propto L(Y|s; \theta^*)$$

$$P(z_l = j) \propto (1 - \rho_j) L(Y|s; \theta_j^*) \quad \text{for } j = 1; 2; \dots; F_L$$

where

$$L(Y|s; \theta_j^*) \propto \exp \left\{ -\frac{1}{2} \sum_{i=1}^n \left(Y_i - \theta_{si} - \sum_{c \neq l} \rho_c (X_{ci} - z_c) \right)^2 \right\}$$

Step 2c: Update ρ_j and the stick-breaking weights (p_j) . Sample $V_j \sim \text{Beta}(c_1 + m_0; d_1 + (L - m_0))$. Then for $j = 1; 2; \dots; F_L$; set

$$\rho_1 = V_1$$

$$\rho_k = (1 - V_1)(1 - V_2)$$

STEP 3b: Update for θ_H

1. Sample $x_H | \theta_H \sim \text{Beta}(\theta_H; L)$
2. Let θ_H equal

$$\theta_H = \frac{\alpha + K_H - 1}{\alpha + K_H - 1 + L(\alpha - \log(x_H))}$$

3. Sample $G | x_G; K_G \sim$

$$\theta_H \text{ Gamma}(\alpha + K_H; \alpha - \log(x_H)) + (1 - G) \text{ Gamma}(\alpha + K_H - 1; \alpha - \log(x_H))$$

STEP 4: Update error precision

Sample $\lambda \sim \text{Gamma}(\lambda^*; \nu^*)$ where

$$\lambda^* = \frac{N}{2} + \frac{1}{2}$$

$$\nu^* = \frac{1}{2} \sum_{i=1}^N \left(Y_i - \theta_{S_i} - \sum_{l=1}^L X_{li} \theta_{z_l} \right)^2$$