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

de ned as follows.

$$
L(Y_i | i;) = \frac{1-2}{\sqrt{2}} exp \frac{-}{2} (Y_i - i)^2
$$

\n
$$
V = \frac{1}{\sqrt{2}} exp \frac{-}{2} (Y_i - i)^2
$$

\n
$$
V = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
V = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
V = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
V = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
V = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
V = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
V = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
G = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
G = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
G = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
G = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
G = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
G = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
G = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
G = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
G = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
G = \frac{1}{\sqrt{2}} Var \frac{-1}{2} (Y_i - i)^2
$$

\n
$$
G = \frac{1}{\sqrt{2}}
$$

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

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

In the above formulation, $\mathbf{u}_i = \log(t)$ (\mathbf{u}_i) where \mathbf{u}_i presents the reference allele frequency for the i^{th} individual at the l^{th} SNP. $_{(0,0)}(\cdot)$ represents a Dirac delta function indicating a point mass at (0,0). In addition, $N(x, \cdot)$ denotes a normal density with mean and precision and $MVN_p(x; M; T)$ 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 xed 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 $i = 1/2, ..., N$; 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 con guration 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 $\frac{1}{11}$ can be sampled using Adaptive-Rejection sampling (?).

STEP 2: Update for \sqrt{l}

In order to update each A g, 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 in nite sum of discrete points, in actual implementation, the Blocked Gibbs Sampler utilizes a nite approximation, imposing a limit F_L to the number of distinct atoms amongst the $\,$ $_{I.}$ Denote this collection of distinct points as $* = \frac{1}{1}, \dots, \frac{1}{F_L}$. ? 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 δ to represent the cluster denoting no e ect (i.e. $\epsilon_{11} = 0$ and l_2 = 0) with associated model weight . Similar to the con-guration representation for l_i , de ne the pointers z_l where $z_l = j$ if and only if $l = \frac{1}{j}$ for $j = 0, 1, 2, \ldots, F_L$. Then de ne m_i as the number of z_i currently equal to j.

Step 2a: For $j = 1/2$; :::; F_L ; update $\frac{*}{j}$. Note, because $\frac{*}{0}$ represents the null e ect cluster, its value need not be updated. If $m_j = 0$, then $\frac{*}{j} \sim H_0$. Else draw $\frac{*}{j} \sim MVN_2\left(M^*\,;\,T^*\right)$ where

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

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

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)ith TelemFelatolun99552 Tif 19 12289003T851[{59.8]TPJ0/FTd5[(11)]975b2F1Th7.9705(37) Step 2b: For $l = 1/2$; :::; L; independently sample z_l where,

$$
P(z_{i} = 0) \propto L(Y|S; *_{0};)
$$

\n
$$
P(z_{i} = j) \propto (1 -)p_{j}L(Y|S; *_{j}; \text{ for } j = 1; 2; ...; F_{L}
$$

\nwhere
\n
$$
2 \qquad \qquad 2^{3}
$$

\n
$$
L(Y|S; *_{j};) \propto \exp 4 \frac{1}{2} \sum_{i=1}^{N} Y_{i} - \sum_{0s_{i}} X_{ii} *_{j} - \sum_{c \neq i} (X_{ci} z_{c})^{5}
$$

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

$$
p_1 = V_1 p_k = (1 - V_1)(1 - V_2)
$$

STEP 3b: Update for H

- 1. Sample x_H | $_H$ ∼ Beta($_H$; L)
- 2. Let H equal

$$
G = \frac{3 + K_H - 1}{3 + K_H - 1 + L(\frac{3}{3} - \log(X_H))}
$$

3. Sample $_G|X_G:K_G \sim$

 H Gamma ($_3 + K_H$; $_3 - log(x_H)$) + (1 – G) Gamma ($_3 + K_H - 1$; $_3 - log(x_H)$)

STEP 4: Update error precision

Sample \sim Gamma(\rightarrow ; \rightarrow) where

* =
$$
\frac{N}{2}
$$
 + 1
* = $1 + \frac{1}{2} \times \frac{N}{i-1}$ $Y_i - 0s_i - \frac{1}{k+1} \times \frac{1}{k+2}$