

# A Semiparametric Bayesian Model for Repeated Repeated Binary Outcomes

Fernando A. Quintana\*     Peter Müller and Gary L. Rosner<sup>†</sup>

Mary V. Relling<sup>‡</sup>

May 30, 2007

## Abstract

We discuss the analysis of data from single nucleotide polymorphism (SNP) arrays comparing tumor and normal tissues. The data consist of sequences of indicators for loss of heterozygosity (LOH) and involves three nested levels of repetition: chromosomes for a given patient, regions within chromosomes, and SNPs nested within regions. We propose to analyze these data using a semiparametric model for multi-level repeated binary data. At the top level of the hierarchy we assume a sampling model for the observed binary LOH sequences that arises from a partial exchangeability argument. We argue that this implies a mixture of Markov chains model. The mixture is defined with respect to the Markov transition probabilities. We assume a nonparametric prior for the random mixing measure. The resulting model takes the form of a semiparametric random effects model with the matrix of transition probabilities being the random effects. The model includes appropriate dependence assumptions for

---

\*Departamento de Estadística, Facultad de Matemáticas, Pontificia Universidad Católica de Chile, Casilla 306, Correo 22, Santiago, CHILE.

<sup>†</sup>Department of Biostatistics, The University of Texas, M. D. Anderson Cancer Center, Box 447, 1515 Holcombe Boulevard, Houston, Texas 77030, U.S.A

<sup>‡</sup>Department of Pharmaceutical Sciences, St. Jude Children's Research Hospital, 332 N. Lauderdale, Memphis, TN 38105-2794 USA

the two remaining levels of the hierarchy, i.e., for regions within chromosomes and for chromosomes within patient.

**Keywords:** Dirichlet Process; Loss of Heterozygosity; Partial Exchangeability; Semi-parametric Random Effects; Species Sampling Models.

## 1 Introduction

In many bio-medical studies investigators collect data on a number of repeat experiments for a given set of patients or subjects. We specifically consider the case where a sequence of binary responses is collected from each experiment. The discussion is motivated by inference about regions of increased loss of heterozygosity (LOH) for single nucleotide polymorphism (SNP) arrays comparing tumor and normal tissue samples from a group of children who experience treatment-related leukemia. An interesting feature of this dataset is that we can identify three nested levels of repetition: chromosomes within a patient, SNP regions within a chromosome, and SNPs within regions. We denote the number of patients by  $n$  and the collected data as  $y_{icjk} = 1$  if LOH was recorded (and zero otherwise) in the  $k$ th SNP from region  $j$  within chromosome  $c$  for patient  $i$ , where  $i = 1, \dots, n$ ,  $c = 1, \dots, 22$ ,  $j = 1, \dots, n_{ic}$  and  $k = 1, \dots, n_{icj}$ . In short, the data consists of binary sequences nested within two levels of repeat experiments recorded for each individual. We generally describe this data structure as “repeatedly repeated binary measurements.”

Newton et al. (1998) and Newton and Lee (2000) proposed the “instability-selection” model for the analysis of LOH data. The model assumes that the observed losses (deletion of genetic material) occur in chromosomes according to a single binary Markov process. An extension of the instability-selection model to pooled analysis of LOH from several experiments is described in Miller et al. (2003). Lin et al. (2004) consider permutation-based methods with windows and hidden Markov models to assess LOH. A related model can be found in Beroukhi et al. (2006).

Our modeling approach is based on similar assumptions. However, it differs from these approaches in that we consider Markov processes that are *specific* to SNP regions within

chromosomes rather than one Markov chain for the entire chromosome, and we model the dependence of these Markov chains across regions and chromosomes. See Section 2 for a description of how regions are constructed. The use of parameters that are specific to regions allows us to define region-specific rates of LOH, and enables us to report the desired inference about regions of increased LOH. The instability-selection model was developed for a different inference goal, namely the mapping of tumor suppression genes.

In summary, we define a model structure with three nested levels of repetition: sequences of binary indicators for LOH within each region, consecutive regions within each chromosome, and chromosomes within a patient. We define a hierarchical model over the three levels of repeated measurements. For the first level of repeated measurements we assume the binary LOH sequences for each SNP region within each chromosome to be partially exchangeable of order  $\ell$ . A probability model for a binary sequence  $y_k$ ,  $k = 1, \dots, n$ , is order- $\ell$  exchangeable if the joint distribution is invariant under permutations that leave the initial  $\ell$  values and all order- $\ell$  transition counts unaltered (Quintana and Newton, 2000; Quintana and Müller, 2004). For example, let  $t(0, 0) = \sum_{k=2}^n I\{y_{k-1} = 0, y_k = 0\}$  and similarly for  $t(0, 1)$ ,  $t(1, 0)$  and  $t(1, 1)$ , denote the order-1 transition counts. An order-1 exchangeable probability model for the sequence  $(y_k, k = 1, \dots, n)$ , is a probability measure that is invariant under any permutation that leaves  $t(0, 0)$ ,  $t(0, 1)$ ,  $t(1, 0)$ ,  $t(1, 1)$  and the initial response  $y_1$  unchanged. Under some additional technical conditions (Quintana and Newton, 1998), such sequences can be represented as mixtures of order- $\ell$  Markov chains. The mixture is with respect to the Markov transition probabilities. This assumption is in agreement with the instability-selection model, but we will allow higher orders of dependence, thus extending the modeling scope.

We complete the description of the mixture of Markov chains model with a nonparametric model on the mixing measure for the corresponding transition matrix (TM). By including latent parameters, the mixture model can be written as a hierarchical model, with the latent TMs interpreted as random effects. Finally, we assume a parametric structure to link the latent parameters across consecutive regions in a chromosome, and across chromosomes within a patient. See Section 3 for details.

Fully parametric versions of such models are successfully used for Bayesian inference in multi-level repeated measurement data. Related multi-level models for discrete data are reviewed in Goldstein (2003). Heagerty and Zeger (2000) discuss maximum likelihood inference using a marginal models approach, i.e., regressing the marginal mean, rather than the conditional mean given the random effects, on covariates. Mixture models with a nonparametric mixing measure to define semiparametric random effects distributions are extensively used in nonparametric Bayesian inference, including, for example, Müller and Rosner (1997), Mukhopadhyay and Gelfand (1997) and Kleinman and Ibrahim (1998). The special case of binary outcomes has been discussed, among many others, by Basu and Mukhopadhyay (2000). An advantage of the model specification with such mixtures is that models with no random effects and fully parametric models with a parametric random effects distribution can be seen as special versions of the nonparametric case.

The nonparametric component in the proposed model is based on the Dirichlet process (DP) model introduced in Ferguson (1973). The main reasons for choosing the DP model are the intuitive nature of the prior predictive distributions and computational ease of posterior simulation. Both features are true for a wider class of probability models, known as species sampling models (SSM). We will therefore first introduce the proposed model using a SSM, although the implementation is carried out with the special case of the DP prior. SSMs are introduced and discussed in Pitman (1996), Ishwaran and James (2003), and in Quintana (2006). Recent reviews of semiparametric Bayesian inference appear in Walker et al. (1999) and Müller and Quintana (2004).

The rest of this article is organized as follows. Section 2 describes the LOH dataset, discussing the specific choice for modeling dependence across nonparametric random effects for chromosome regions. In section 3, we describe the main features of the proposed model, with emphasis on how we model dependence along LOH sequences from any chromosome region and among regions. We will show that the implied *marginal* model for each region remains nonparametric. Posterior simulation schemes are discussed in Section 4. Section 5 reports the resulting inference, also comparing it with the instability-selection model. Section 6 concludes with a final discussion.

## 2 The Data

The motivating dataset comes from a study conducted at the St. Jude Children’s Research Hospital in Memphis, Tennessee (SJCRH). A full description and discussion of these data can be found in Hartford et al. (2006). We briefly summarize the study. Some children develop second cancers after successful treatment for their initial cancer diagnosis. Of particular concern to the investigators are therapy-related leukemias. Several studies have identified characteristics of particular anti-cancer therapies that may affect the child’s risk of developing a later secondary leukemia (Pedersen-Bjergaard, 2005). Some of these treatments may cause genetic alterations that lead to the the patient’s subsequent secondary cancer.

Investigators at SJCRH carried out genome-wide studies of children with secondary leukemia, in order to learn about genetic factors that may contribute to a patient’s risk of a secondary leukemia. A genomic alteration of particular interest to the investigators was loss of heterozygosity (LOH). Heterozygosity refers to the presence of two different allele versions of a gene at corresponding loci of a pair of chromosomes (i.e., being heterozygous). LOH is when an allele at a particular locus is missing. LOH can occur if part or all of one of the paired chromosomes is missing or if there is a deletion or mutation of part of one of the chromosome pairs. LOH can lead to cancer if it occurs at the site of a normal tumor suppressor gene that was keeping a cancer-susceptibility gene in check.

The investigators used single nucleotide polymorphism (SNP) arrays to compare germline (normal) and tumor (secondary leukemia blasts) samples. The arrays interrogated more than 100,000 SNPs in samples from 13 patients with a diagnosis of treatment-related leukemia. These patients had enrolled in SJCRH protocols for treatment of their initial diagnosis of acute lymphoblastic leukemia (ALL) (Relling et al., 2003). Specifically, they amplified, labeled, and hybridized 500 ng of DNA from each sample to the Affymetrix GeneChip® Human Mapping100K Set. After scanning the chips, they applied the GeneChip® DNA Analysis Software to make the genotype calls for the data. With the genotype calls, the investigators declared LOH, retention, or indeterminate for each SNP investigated by the array, following the approach of Lin et al. (2004). The germline samples for these patients came from DNA

the investigators extracted from normal leukocytes (white blood cells) at the time the child achieved his or her initial remission. Bone marrow at the time of diagnosis of secondary leukemia was the source of the leukemic blast samples. The data consist of  $n = 13$  binary sequences with an outcome  $y = 1$  for a recorded LOH at a given SNP, and a zero otherwise. Each sequence is of length 116,204.

One research question of interest for these data is the identification of regions of increased LOH. Consequently, we divided the LOH sequences into regions. For chromosomes with more than 0.5% recorded LOH we used regions of length 55 SNPs each. For all other chromosomes we used regions of length 835. This resulted in a total of 874 regions of lengths either 55 (for chromosomes 5 through 9 and 15, 16 and 17) or 835 (for all other chromosomes). The two nested levels of repeated measurements are thus given by regions and the sequence of recorded indicators for LOH within each region. Besides a general notion of dependence, little is known about appropriate probability models for such data. In the process of mitosis, chromosomes cross and genetic material gets shuffled. Nucleotides closer to each other are less likely to get separated than those that are farther away. This linkage disequilibrium phenomenon suggests a Markovian dependence, in line with the sampling model to be described in Section 3.

The context of the data that we analyze here differs from that in the aforementioned papers. Our data concern patients who were previously treated for cancer and subsequently developed treatment-related leukemia. The investigators hypothesize that the chemotherapy causes chromosomal damage by mechanisms that may be different from those leading to spontaneous (new) cancers. We therefore propose a model that is developed for the goal of identifying consistent *regions* of high LOH without relying on any assumed mechanism.

# 3 A Model for Repeated Repeated Binary LOH Measurements

## 3.1 The Sampling Model

Recall that  $y_{icjk}$  represents the binary indicator of LOH for SNP  $k$  in region  $j$  of chromosome  $c$  for patient  $i$ . Let  $\mathbf{y}_{icj} = (y_{icjk}, 1 \leq k \leq n_{icj})$  be the entire LOH sequence from the  $j$ -th region for chromosome  $c$  of the  $i$ -th patient. We model correlation at the level of the observed binary outcomes by assuming the sequences  $\mathbf{y}_{icj}$  to be partially exchangeable of some order. We assume a fixed maximum order  $\ell$  that is common to all sequences. It can be shown (Quintana and Newton, 1998) that order- $\ell$  exchangeability implies that  $p(\mathbf{y}_{icj})$  can be expressed as a mixture of homogeneous order- $\ell$  Markov chains. The mixture is with respect to the order- $\ell$  transition probabilities.

Let  $\boldsymbol{\theta}_{icj}$  denote the transition matrix (TM) that defines the Markov model for subject  $i$ , chromosome  $c$  and region  $j$ . The transition probabilities  $\boldsymbol{\theta}_{icj}$  can be represented as a  $2^\ell$ -dimensional vector of transition probabilities  $\theta_{icj, m_\ell, \dots, m_1, 0}$  from state  $(m_\ell, m_{\ell-1}, \dots, m_1)$  to  $(m_{\ell-1}, \dots, m_1, 0)$ , where  $m_k \in \{0, 1\}$  for all  $k$ . Denote by  $t_{icj}(m_\ell, \dots, m_1, 0) = \sum_{k=\ell+1}^{n_{icj}} I(y_{icjk} = 0, y_{icj, k-1} = m_1, \dots, y_{icj, k-\ell} = m_\ell)$ , i.e. the count of transitions from state  $(m_\ell, \dots, m_1)$  to  $(m_{\ell-1}, \dots, m_1, 0)$ , for region  $j$  in chromosome  $c$  of patient  $i$ . The likelihood is given by

$$p(\mathbf{y}_{icj} | \boldsymbol{\theta}_{icj}) = \prod_{m_\ell, \dots, m_1 \in \{0,1\}} \left\{ \theta_{icj}(m_\ell, \dots, m_1, 0)^{t_{icj}(m_\ell, \dots, m_1, 0)} \times [1 - \theta_{icj}(m_\ell, \dots, m_1, 0)]^{t_{icj}(m_\ell, \dots, m_1, 1)} \right\}. \quad (1)$$

We adopt (1) with fixed order  $\ell = 2$ . The implied data reduction by sufficiency to a set of  $2^{\ell+1} = 8$  transition counts is critical to facilitate fast likelihood evaluation. The assumption  $\ell = 2$  implies that 4 parameters are required to represent each of the 11,362 TMs (874 per patient) involved in the likelihood model. The choice of  $\ell = 2$  is supported by exploratory analysis (not shown). Also, it generalizes the order-1 Markov models used in Newton and Lee (2000) and Lin et al. (2004).

### 3.2 Random Effects Model

We complete the definition of the order- $\ell$  exchangeable model (with  $\ell = 2$ ) as a mixture of Markov chains by adding a probability model for the TMs  $\boldsymbol{\theta}_{icj}$  in (1). In words, we use a hierarchical normal model to define dependence across chromosomes, a non-parametric prior to define the joint distribution of subject-specific effects, and a normal autoregression model to define dependence across regions. The main features of the proposed model are the use of flexible nonparametric priors for the implied marginal distribution of the random effects at all three levels, i.e., regions, chromosomes and subjects, and the use of parsimonious parametric models to define the dependence structure across regions and chromosomes.

We start with a model for the TM  $\boldsymbol{\theta}_{ic1}$  corresponding to the first region, i.e.  $j = 1$  of chromosome  $c$  and patient  $i$ . Focusing on only one region ( $n_{ic} = 1$ ) for the moment, we reduce notation to  $\boldsymbol{\theta}_{ic} \equiv \boldsymbol{\theta}_{ic1}$ . We assume a normal hierarchical model across chromosomes,  $p(\text{logit}(\boldsymbol{\theta}_{ic}) \mid \boldsymbol{\mu}_{ic}) = N(\boldsymbol{\mu}_{ic}, \mathbf{S})$  and  $p(\boldsymbol{\mu}_{ic} \mid \boldsymbol{\mu}_i) = N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_1)$ , independently across  $c = 1, \dots, 22$ . In other words, we define the dependence across chromosomes by assuming an exchangeable normal model for the TMs on a logit scale. We complete the model, still restricting to the first region only, by assuming an exchangeable prior on  $\boldsymbol{\mu}_i$  across patients,  $\boldsymbol{\mu}_i \sim F$ . Instead of specific parametric assumptions for  $F$  we use a non-parametric prior. Formally, we define  $F$  to be a random probability measure and write  $F \sim RPM$ . We define specific choices of RPMs below. In summary, we assume

$$\text{logit}(\boldsymbol{\theta}_{ic}) \sim N(\boldsymbol{\mu}_{ic}, \mathbf{S}), \quad \boldsymbol{\mu}_{ic} \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_1), \quad \boldsymbol{\mu}_i \sim F \quad \text{and} \quad F \sim RPM. \quad (2)$$

Let  $\delta_x(\cdot)$  denote a point mass at  $x$ . To define the nonparametric component  $F$ , we consider RPMs that can be represented as  $F(\cdot) = \sum_{h \geq 1} w_h \delta_{\mathbf{z}_h}(\cdot)$  where  $\{w_h\}$  is a collection of random weights constrained by  $0 \leq w_h \leq 1$  and  $P(\sum_{h \geq 1} w_h = 1) = 1$ , and  $\{\mathbf{z}_h\}$  is a random sample from a *baseline* distribution  $F_0$ , independently of  $w$ . The baseline  $F_0$  may itself depend on additional hyper-parameters. The discrete nature of  $F$  implies a positive probability for ties among the  $\boldsymbol{\mu}_i$  values. The groups of patients sharing a common value can be interpreted as *clusters*. Let  $\boldsymbol{\mu}_1^*, \dots, \boldsymbol{\mu}_L^*$  be the unique values among  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_n$ . We define latent indicators  $s_i$  for *cluster memberships*, such that  $s_i = h$  if  $\boldsymbol{\mu}_i = \boldsymbol{\mu}_h^*$ . Let also

$m_h$  represent the *size* of the  $h$ th cluster, i.e. the number of  $\boldsymbol{\mu}_i$ s with value equal to  $\boldsymbol{\mu}_h^*$ . An important class of RPMs are the *species sampling models* (SSM) discussed in Pitman (1996) and Ishwaran and James (2003). Such models can be best understood by considering the prior predictive distribution for a sample  $\boldsymbol{\mu}_i \sim F$ ,  $i = 1, \dots, n$ , generated from a random probability model with a SSM prior. Marginalizing with respect to  $F$ , the following prior predictive probabilities apply:

$$p(\boldsymbol{\mu}_{n+1} \mid \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_n) = \begin{cases} \delta_{\boldsymbol{\mu}_h^*}(\boldsymbol{\mu}_{n+1}), & \text{with probability } \rho_h(m_1, \dots, m_L), \quad h = 1, \dots, L \\ F_0(\boldsymbol{\mu}_{n+1}), & \text{with probability } \rho_{L+1}(m_1, \dots, m_L), \end{cases} \quad (3)$$

where  $\rho_h(m_1, \dots, m_L)$  are predictive configuration probabilities depending only on the cluster sizes formed by the first  $n$  patients (Pitman, 1996). The predictive rule (3) implies that with some probability the TMs for a new patient mimic some of the previous ones (up to residual variation); with the remaining probability, the latent parameters controlling the TMs are drawn from the baseline distribution  $F_0$ . Letting  $\rho$  represent the entire collection of predictive probabilities  $\{\rho_h(\cdot)\}$ , we use the notation  $F \sim \text{SSM}(\rho, F_0)$ . The choice of  $\rho$  is subject to constraints. See, for example Pitman (1996). A computationally convenient choice for the baseline measure  $F_0$  is a conjugate prior to the kernel in (2). In the case of the normal kernel we use  $F_0(\boldsymbol{\mu}) = N(\boldsymbol{\mu}; \mathbf{m}, \mathbf{V})$ .

A popular special case of the SSM is the DP (Ferguson, 1973) for which

$$\rho_h(m_1, \dots, m_L) = m_h / (M + n), \quad h = 1, \dots, L \quad \text{and} \quad \rho_{L+1}(m_1, \dots, m_L) = M / (M + n), \quad (4)$$

where  $M$  is the *total mass parameter*. As a default choice we recommend using a DP prior unless specific prior information suggests a different set of predictive probabilities  $\rho$ .

In summary, we have defined dependence across chromosomes  $c$  by the hierarchical model (2) and dependence across subjects  $i$  by (3). A critical feature of the proposed model is that the subject specific random effects  $\boldsymbol{\mu}_i$  are of dimension  $2^\ell$ . A non-parametric prior for the joint vector of all  $\boldsymbol{\mu}_{ic}$  would be of prohibitive dimension. In other words, we use a parametric model to specify dependence across  $c$ , but the marginal model for each  $\boldsymbol{\mu}_{ic}$ , marginalizing

w.r.t.  $\boldsymbol{\mu}_i$ , is a semiparametric mixture of normals

$$p(\boldsymbol{\mu}_{ic} | F) \sim \int N(\boldsymbol{\mu}, \mathbf{S} + \boldsymbol{\Sigma}_1) dF(\boldsymbol{\mu}).$$

We now complete the model by defining dependence across the second level of repeat experiments in the data structure, i.e., dependence across regions, using a similar modeling strategy. The only difference is that an appropriate model for dependence across regions is based on spatial dependence rather than exchangeability. We return to the general case with multiple regions,  $j = 1, \dots, n_{ic}$ , for each chromosome and patient, as described in Section 2. Let  $\boldsymbol{\theta}_{ic} = (\boldsymbol{\theta}_{ic1}, \dots, \boldsymbol{\theta}_{icn_{ic}})$ ,  $\boldsymbol{\theta}_i = (\boldsymbol{\theta}_{i1}, \dots, \boldsymbol{\theta}_{i,22})$  and  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_n)$ . Let also  $\boldsymbol{\mu}_{ic}$ ,  $\boldsymbol{\mu}_i$  and  $\boldsymbol{\mu}$  be analogously defined. The likelihood remains as in (1). As in (2) we introduce latent variables  $\boldsymbol{\mu}_{icj}$  with  $\text{logit}(\boldsymbol{\theta}_{icj}) | \boldsymbol{\mu}_{icj} \sim N(\boldsymbol{\mu}_{icj}, \mathbf{S})$ , and replace (2) by an extended model to include dependence across regions. We introduce a vector of autoregressive coefficients  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_{2^\ell})$ . The coefficient  $\alpha_h$  is used to characterize the change of the  $h$ -th coefficient of the transition matrix  $\boldsymbol{\mu}_{icj}$  across regions. Let  $D(\boldsymbol{\alpha}) = \text{diag}(\boldsymbol{\alpha})$  denote the  $2^\ell \times 2^\ell$  diagonal matrix with  $\boldsymbol{\alpha}$  on the diagonal. We replace (2) by

$$\boldsymbol{\mu}_{i0} \sim F, \quad \boldsymbol{\mu}_{ic0} = \boldsymbol{\mu}_{i0} + \boldsymbol{\epsilon}_{ic0}, \quad \boldsymbol{\mu}_{ic1} = \boldsymbol{\mu}_{ic0} + \boldsymbol{\epsilon}_{ic1} \quad \text{and} \quad \boldsymbol{\mu}_{icj} = \boldsymbol{\mu}_{ic0} + D(\boldsymbol{\alpha})(\boldsymbol{\mu}_{ic,j-1} - \boldsymbol{\mu}_{ic0}) + \boldsymbol{\epsilon}_{icj}, \quad (5)$$

where  $\boldsymbol{\epsilon}_{icj}$  and  $\boldsymbol{\epsilon}_{ic0}$  are independent normal residuals with  $\boldsymbol{\epsilon}_{ic0} \sim N(\mathbf{0}, \boldsymbol{\Sigma}_1)$ , and  $\boldsymbol{\epsilon}_{icj} \sim N(\mathbf{0}, \boldsymbol{\Sigma}_2)$  for  $j \geq 1$ . The model assumes first-order stationarity, and includes a strictly stationary model as a special case by appropriate choices of  $\boldsymbol{\Sigma}_1$ ,  $\boldsymbol{\Sigma}_2$  and  $\boldsymbol{\alpha}$ . Marginally for each region  $j$ , the model still implies the nonparametric mixture of normals model for  $\text{logit}(\boldsymbol{\theta}_{icj})$ , as before, now with the kernel  $N(\text{logit}(\boldsymbol{\theta}_{ic}); \boldsymbol{\mu}_{ic}, \mathbf{S})$  replaced by a normal kernel  $p(\text{logit}(\boldsymbol{\theta}_{icj}) | \boldsymbol{\mu}_{ic0}, \mathbf{S}, \boldsymbol{\alpha}) = N(\boldsymbol{\mu}_{ic0}, V(\mathbf{S}, \boldsymbol{\alpha}))$  with the variance-covariance matrix  $V(\mathbf{S}, \boldsymbol{\alpha})$  implied by marginalizing (5) with respect to  $\boldsymbol{\mu}_{ic1}, \dots, \boldsymbol{\mu}_{ic,j-1}$ . The cluster structure on patients remains determined by the SSM assumption for  $F$ . The desired learning about regions of increased LOH is then accomplished by examining the posterior distributions of an appropriate function of the  $\boldsymbol{\theta}_{icj}$  parameters. See Section 5 below for details on how we do this.

The structure in (5) highlights how the proposed semi-parametric relates to a simpler

parametric model. If we were to assume a parametric model for  $\boldsymbol{\mu}_{i0}$ , for example the base measure of the DP,  $\boldsymbol{\mu}_{i0} \sim N(\boldsymbol{m}, \mathbf{V})$ , the model would reduce to a fully parametric hierarchical model. Besides increased flexibility the advantage of the semi-parametric extension is that it remains more faithful to the prior judgement about the binary sequences by building on only the order  $\ell$  exchangeability assumption.

The model specification is completed by defining hyperpriors on all remaining parameters. Let  $\boldsymbol{\eta}$  denote the set of all other hyper-parameters. These include the regression coefficients  $\boldsymbol{\alpha}$ , the covariance matrices  $\mathbf{S}$ ,  $\boldsymbol{\Sigma}_1$  and  $\boldsymbol{\Sigma}_2$ , and hyper-parameters from the baseline distribution  $F_0$ ,  $\boldsymbol{m}$  and  $\mathbf{V}$ . For  $\boldsymbol{\alpha}$  we use a normal prior,  $p(\boldsymbol{\alpha}) = N(\boldsymbol{\alpha}; \boldsymbol{a}_0, \mathbf{A}_0)$ . Next, for  $\mathbf{S}$  we choose an inverse-Wishart prior,  $p(\mathbf{S}) = IW(\mathbf{S}; \mathbf{S}_0, \nu_S)$ . We also assume independent conjugate inverse-Wishart priors for the residual variances:  $p(\boldsymbol{\Sigma}_1) = IW(\boldsymbol{\Sigma}_1; \boldsymbol{\Sigma}_{01}, \nu_1)$  and  $p(\boldsymbol{\Sigma}_2) = IW(\boldsymbol{\Sigma}_2; \boldsymbol{\Sigma}_{02}, \nu_2)$ . Finally, for the hyper-parameters in  $F_0$  we use  $p(\boldsymbol{m}, \mathbf{V}) = N(\boldsymbol{m}; \boldsymbol{m}_0, \mathbf{V}/\lambda_0) \times IW(\mathbf{V}; \mathbf{V}_0, \nu_V)$ . In the earlier definitions we assume  $\boldsymbol{a}_0$ ,  $\mathbf{A}_0$ ,  $\nu_A$ ,  $\mathbf{S}_0$ ,  $\nu_S$ ,  $\nu_1$ ,  $\boldsymbol{\Sigma}_{01}$ ,  $\nu_2$ ,  $\boldsymbol{\Sigma}_{02}$ ,  $\boldsymbol{m}_0$ ,  $\mathbf{V}_0$ ,  $\lambda_0$  and  $\nu_V$  to be known.

## 4 Posterior Simulation

Model (2) has the great advantage of conditional independence at various levels. This conditional independence facilitates implementation of the Gibbs sampling algorithm. In particular, the transition probabilities  $\boldsymbol{\theta}_{icj}$  are conditionally independent across  $i$ ,  $c$  and  $j$ , given all other parameters. As a consequence, the  $\boldsymbol{\theta}_{icj}$  can be updated one at a time. Sampling from the corresponding full conditionals can be accomplished using standard methods for logistic regression, as discussed, e.g., in Carlin and Louis (1996). Another important feature of the model is that it can without problem accommodate unbalanced data with variable numbers of SNP regions,  $n_{ic}$ , per chromosome within patient and varying length sequences of binary LOH measurements,  $n_{icj}$ , per region.

Next, consider updating  $\boldsymbol{\mu}_{i0}$  in model (5). Details on updating the configurations of ties among the  $\{\boldsymbol{\mu}_{i0}, i = 1, \dots, n\}$ , and the unique values  $\boldsymbol{\mu}_h^*$  are described, among others, in MacEachern and Müller (2000) and Neal (2000) for the special DP case, and in Ishwaran and

James (2003) and Quintana (2006) for the more general SSMS. Any posterior simulation for DP mixture models can easily be adapted for general SSMS by replacing the special DP prior probabilities in (4) by the more general  $\rho_h(m_1, \dots, m_L)$  in (3). The complete conditional posterior for  $\boldsymbol{\mu}_{icj}$ ,  $j \geq 1$ , including conditioning on  $\text{logit}(\boldsymbol{\theta}_{icj})$  and  $\{\boldsymbol{\mu}_{ics}\}$ ,  $s \neq j$ , is a straightforward normal linear regression.

Updating the autoregression parameters in (5) proceeds by draws from the complete conditional posterior distribution. The distributions follow easily from the fact that model (5) is linear in  $\boldsymbol{\alpha}$ . The conjugate normal prior assumption for  $\boldsymbol{\alpha}$  allows for straightforward posterior simulation for  $\boldsymbol{\alpha}$  conditional on imputed values for all other parameters. Finally, the remaining parameters are easily updated from the corresponding conjugate-style conditionals. See further details in Quintana and Müller (2004) and in Müller et al. (2007).

## 5 Identifying Regions of Increased LOH

We assume model (1) with  $\ell = 2$ , that is,  $\boldsymbol{\theta}_{icj}$  is of dimension  $2^\ell = 4$ , and represents the full order-2 TM for the  $j$ th SNP region for chromosome  $c$  of the  $i$ th patient. For the random effects distribution  $p(\boldsymbol{\theta}_{icj}, j = 1, \dots, n_{ic})$ , we use model (5) with the SSM prior specifying default predictive weights for a  $DP(M, F_0)$  prior with  $M = 1$ . The model treats the responses from different regions as conditionally independent given region-specific parameters.

Figure 1 shows the estimated marginal posterior means plus and minus one posterior standard deviation for the components of the  $\boldsymbol{\alpha}$  and  $\boldsymbol{m}_0$  coefficients. They are well away from zero, suggesting a significant autoregression effect, except possibly for  $\alpha_3$ , which corresponds to the logit of transition probabilities from (1, 0) to (0, 1). Also, they are significantly away from 1, except possibly for  $\alpha_1$ , i.e the coefficient for (0, 0) to (0, 1) transitions. In other words, the data suggests an autoregressive effect that is not a random walk-type process. On the other hand, the  $\boldsymbol{m}_0$  coefficients, which control the center of the baseline measure (in the logit scale) are all well negative. This reflects the fact that over 99% of all the responses are zero, and so the baseline values for transition probabilities from any previous two values to an LOH response are very low for any given region. This is further reflected in Figure 2a

which shows the estimated posterior means for the (0,0) to (0,0) transition probabilities (on the logit scale). The estimated transition probabilities are very high for almost all regions and patients.

If one wishes to evaluate LOH in any given region, we recommend using the long-run proportion of LOH in that region. This is based on the fact that chromosome regions are large enough to justify approximating (ergodic) LOH averages by their corresponding limits. From elementary Markov chain theory (Ross, 2002), we find that these are given by

$$\begin{aligned} \lim_{k \rightarrow \infty} P(y_{icjk} = 1) = \\ \lim_{k \rightarrow \infty} \sum_{j_1, j_2 \in \{0,1\}} P(y_{icjk} = 1 | y_{icj, k-1} = j_1, y_{icj, k-2} = j_2) P(y_{icj, k-1} = j_1, y_{icj, k-2} = j_2) = \\ \sum_{j_1, j_2 \in \{0,1\}} P(y_{icj2} = 1 | y_{icj1} = j_1, y_{icj0} = j_2) \lim_{k \rightarrow \infty} P(y_{icj, k-1} = j_1, y_{icj, k-2} = j_2). \end{aligned}$$

The first term in the last summation is the transition probability from  $(j_2 j_1)$  to  $(j_1 1)$ , while the limit probabilities in the second part are given by the stationary distribution corresponding to the appropriate TM, and therefore, easily identified as functions of  $\theta_{icj}$ .

Figure 2b shows the estimated marginal posterior long-run probabilities of LOH in the logit scale, computed as indicated above. Patients 2, 6, 8, 9, and 10 show some regions with higher probability of LOH than the other patients. In general we distinguish an overall low percentage of observed LOH sites.

The estimated marginal probabilities do not yet address the original inference goal of identifying regions of increased LOH. We address this goal by defining an indicator of “increased LOH” for patient  $i$  and region  $j$  as

$$I_{icj} = \begin{cases} 1 & \text{if } \lim_{k \rightarrow \infty} P(y_{icjk} = 1) > p_0 \\ 0 & \text{otherwise,} \end{cases} \quad (6)$$

where  $0 < p_0 < 1$  is a fixed threshold. In other words, we say that a given region has increased LOH if the marginal long-run probability of LOH is greater than  $p_0$ . As noted earlier,  $\lim_{k \rightarrow \infty} P(y_{icjk} = 1)$  is a function of the transition probabilities  $\theta_{icj}$  only. Figure 3a reports the posterior expectations of  $I_{icj}$  in the logit scale, using  $p_0 = 1\%$ . Patients 2, 6, 8,

9, and 10 stand out again. We interpret this observation to mean that these five patients have longer stretches of neighboring loci that have relatively higher posterior probability of increased LOH relative to the threshold of 1%. We also notice in Figure 3a that some of the end regions on the right tend to exhibit higher probability of increased LOH than do the other chromosomal regions in our data set.

We still need to relate the reported probabilities to the desired decision of identifying regions of increased LOH. In carrying out this decision we face a massive multiplicity problem. In the context of high throughput gene expression data, several procedures have been proposed to address such decision problems based on the notion of false discovery rates (Benjamini and Hochberg, 1995; Storey, 2002). Most discussions are for the stylized setup of a two-group comparison microarray experiment. For each of a large number of genes recorded on the microarrays, we wish to make a decision about differential expression. Under the same setup and using a Bayesian decision-theoretic perspective, Müller et al. (2004) show that under a variety of loss functions the optimal decision is characterized by flagging all those comparisons with marginal probability of differential expression beyond a certain threshold. The conclusion is valid for any probability model. In particular, the probability model can include dependence, as in the proposed model for LOH indicators. Thus the same solution applies. The optimal inference about regions of increased LOH is achieved by marking all regions with marginal probability of increased LOH beyond a threshold. Let  $\delta_{icj} \in \{0, 1\}$  denote an indicator for reporting increased LOH for region  $(icj)$ . Let  $D = \sum_{icj} \delta_{icj}$  denote the total number of reported regions. Let  $I_{icj} = I_{icj}(\boldsymbol{\theta}_{icj})$  denote the unknown truth and define the false discovery proportion (FDP) as

$$\text{FDP} = \sum_{i,j} (1 - I_{icj}) \delta_{icj} / D.$$

The FDP is a function of the unknown parameters  $I_{icj}$  and the data, implicitly through  $\delta_{icj}$ . Let  $v_{icj} = E(I_{icj} \mid \text{data})$  denote the marginal posterior probability of increased LOH in region  $(icj)$ . These probabilities are reported in Figure 3a. The posterior expected FDP,  $\overline{\text{FDR}} = E(\text{FDP} \mid \text{data})$ , is evaluated as  $\overline{\text{FDR}} = \sum_{i,j} (1 - v_{icj}) \delta_{icj} / D$ . Using a threshold of 0.29, i.e.,  $\delta_{icj} = I(v_{icj} > 0.29)$ , we find  $\overline{\text{FDR}} = 10\%$ . This inference is reported in Figure 3b,

with black bars indicating the decision to flag a region as exhibiting “increased LOH.” Note the double thresholding that is implicit in the definition of  $\delta$  by thresholding the posterior expectation of  $I_{icj}$ , which in turn is defined by a threshold on the limiting probabilities. This arises because we are interested in regions of increased LOH, rather than regions of LOH (the latter would be very few for a comparable FDP). Figure 3b confirms the initial inference about the rightmost end regions having relatively high posterior probability of increased LOH. For two patients (2 and 6) we see uniformly increased LOH across all regions.

Finally, for comparison we applied three alternative approaches to the dataset at hand. We used the methods proposed in Newton et al. (1998) and Lin et al. (2004) and a fully parametric version of the proposed model. For Newton et al. (1998)’s method, we used a set of 100 equally spaced loci. We found all the log of odds (LOD) scores to be zero, i.e. no region of high LOH was detected. This somewhat surprising conclusion is explained by the fact that our binary sequences are very long and the proportion of recorded LOH is very low. Thus a model that assumes a single Markovian process across all regions, may not be flexible enough to capture local behavior as our approach does. In this case, the MLEs required to implement the model in Newton et al. (1998) are essentially driven by the overwhelming proportion of observed “no loss” and thus the result. Results for the approach proposed in Lin et al. (2004) are shown in Figure 4. We used the implementation in dChip, the public domain software that is introduced in Lin et al. (2004). Figure 4a plots the probability of LOH using the hidden Markov model score defined in dChip. Compare the inference with Figure 2b. While the general patterns are similar under both methods, inference under the proposed model includes more extensive borrowing strength within the hierarchical model. Also, the reported inference are region-specific probabilities under a coherent joint probability model across regions, chromosomes and samples. This is important when the investigator is interested in summaries like the reported inference about regions of increased LOH that depend on joint probability models. On the other hand, an important advantage of the approach in Lin et al. (2004) is the highly reduced reliance on a specific model. In fact the approach includes an option to compute simple entirely model-independent LOH scores. Significance is judged by a permutation test with appropriate multiplicity control. Finally, we

implemented a fully parametric model by replacing the nonparametric model  $\boldsymbol{\mu}_{i0} \sim F$  in (5) with a parametric model  $\boldsymbol{\mu}_{i0} \sim N(\boldsymbol{m}, \boldsymbol{V})$ . All other model choices are left unchanged. Figure 4b shows the resulting probability of LOH. Compared with Figure 2b we see less smoothing across regions in the reported probabilities, but otherwise similar results. The additional smoothing in the semi-parametric model arises from the clustering that is implied by the DP prior. In contrast the parametric model can be described as assuming all singleton clusters, i.e., all random effects  $\boldsymbol{\mu}_{i0}$  in (5) are distinct. The main argument for the semi-parametric model is that it naturally follows from the prior judgement about partial exchangeability of the binary sequences.

## 6 Conclusion

Motivated by the analysis of LOH data, we have introduced a semiparametric Bayesian model for binary measurements with two nested levels of repetition. The top-level repetition is modeled by means of mixtures of Markov chains on the TM for a given order of dependence  $\ell$  (SNPs within chromosome segments). Marginally, for each second level repeated measurement unit (chromosome region), a nonparametric model is postulated for the random effects related to that unit. The proposed approach completes the model by defining dependence across regions and across chromosomes, using a parametric hierarchical model.

The data analysis carried out can be extended and complemented in several ways. The goal of the inference in our motivating application was to detect regions of increased LOH. In other applications with LOH data one might be interested in modeling for loss of tumor suppression genes that are hypothesized to cause the observed LOH. Newton and Lee (2000) propose an instability-selection model that facilitates such inference. The instability-selection model could be used to replace the partially exchangeable sampling model in our approach, while still keeping dependence across SNPs as in (5).

For other applications it might be useful to extend the proposed model to include occasion-specific covariates in (1). This could be achieved by using a log-linear model for the transition probabilities. The log-linear model would include the regression on the lagged

outcomes as well as additional occasion-specific covariates. The regression coefficients in the log-linear model would then replace  $\text{logit}(\theta_i)$  in (2).

## Acknowledgments

This research was supported, in part, by grants CA075981 and GM061393 from the U.S. National Cancer Institute, and by grants FONDECYT 1060729 and Laboratorio de Análisis Estocástico PBCT-ACT13. We thank Dr. Wenjian Yang for help with the data files and advice about the analysis in dChip. We also thank the Associate Editor and Referees for their valuable comments and suggestions.

## References

- Basu, S. and Mukhopadhyay, S. (2000), “Bayesian analysis of binary regression using symmetric and asymmetric links,” *Sankhyā*, 62, 372–387.
- Benjamini, Y. and Hochberg, Y. (1995), “Controlling the false discovery rate: a practical and powerful approach to multiple testing,” *Journal of the Royal Statistical Society. Series B. Methodological*, 57, 289–300.
- Beroukhi, R., et al. (2006), “Inferring Loss-of-Heterozygosity from Unpaired Tumors Using High-Density Oligonucleotide SNP Arrays,” *PLoS Computational Biology*, 2.
- Carlin, B. P. and Louis, T. A. (1996), *Bayes and Empirical Bayes Methods for Data Analysis*, London: Chapman & Hall.
- Ferguson, T. S. (1973), “A Bayesian analysis of some nonparametric problems,” *The Annals of Statistics*, 1, 209–230.
- Goldstein, H. (2003), *Multilevel Statistical Models, Third Edition*, Kendall’s Library of Statistics, 3, London: Arnold Publishers.

- Hartford, C., et al. (2006), “Genome Scan for Therapy-Related Myeloid Leukemia,” Technical report, Department of Pharmaceutical Sciences, St. Jude Children’s Research Hospital.
- Heagerty, P. J. and Zeger, S. L. (2000), “Marginalized multilevel models and likelihood inference,” *Statistical Science*, 15, 1–19.
- Ishwaran, H. and James, L. J. (2003), “Generalized weighted Chinese restaurant processes for species sampling mixture models,” *Statistica Sinica*, 13, 1211–1235.
- Kleinman, K. and Ibrahim, J. (1998), “A Semi-parametric Bayesian Approach to the Random Effects Model,” *Biometrics*, 54, 921–938.
- Lin, M., et al. (2004), “dChIPSNP: significance curve and clustering of SNP-array-based loss-of-heterozygosity data,” *Bioinformatics*, 20, 1233–1240.
- MacEachern, S. N. and Müller, P. (2000), “Efficient MCMC Schemes for Robust Model Extensions using Encompassing Dirichlet Process Mixture Models,” in *Robust Bayesian Analysis*, eds. F. Ruggeri and D. R. Insua, New York.
- Miller, B. J., et al. (2003), “Pooled Analysis of Loss of Heterozygosity in Breast Cancer: a Genome Scan Provides Comparative Evidence for Multiple Tumor Suppressors and Identifies Novel Candidate Regions,” *American Journal of Human Genetics*, 73, 748–767.
- Mukhopadhyay, S. and Gelfand, A. E. (1997), “Dirichlet Process Mixed Generalized Linear Models,” *Journal of the American Statistical Association*, 92, 633–639.
- Müller, P. and Quintana, F. (2004), “Nonparametric Bayesian Data Analysis,” *Statistical Science*, 19, 95–110.
- Müller, P. and Rosner, G. (1997), “A Bayesian population model with hierarchical mixture priors applied to blood count data,” *Journal of the American Statistical Association*, 92, 1279–1292.
- Müller, P., Rosner, G., and Quintana, F. A. (2007), “Semiparametric Bayesian Inference for Multilevel Repeated Measurement Data,” *Biometrics*, 63, 280–289.

- Müller, P., et al. (2004), “Optimal Sample Size for Multiple Testing: the Case of Gene Expression Microarrays,” *Journal of the American Statistical Association*, 99.
- Neal, R. M. (2000), “Markov chain sampling methods for Dirichlet process mixture models,” *Journal of Computational and Graphical Statistics*, 9, 249–265.
- Newton, M. A. and Lee, Y. (2000), “Inferring the Location and Effect of Tumor Suppressor Genes by Instability-Selection Modeling of Allelic-Loss Data,” *Biometrics*, 56, 1088–1097.
- Newton, M. A., et al. (1998), “On the statistical analysis of allelic-loss data,” *Stat Med*, 17, 1425–45.
- Pedersen-Bjergaard, J. (2005), “Insights into leukemogenesis from therapy-related leukemia,” *New England Journal of Medicine*, 352, 1591–1594.
- Pitman, J. (1996), “Some Developments of the Blackwell-MacQueen Urn Scheme,” in *Statistics, Probability and Game Theory. Papers in Honor of David Blackwell*, eds. T. S. Ferguson, L. S. Shapeley, and J. B. MacQueen, IMS Lecture Notes, 245–268, IMS.
- Quintana, F. and Müller, P. (2004), “Nonparametric Bayesian Assessment of the Order of Dependence for Binary Sequences,” *Journal of Computational and Graphical Statistics*, 13, 213–231.
- Quintana, F. A. (2006), “A predictive view of Bayesian clustering,” *Journal of Statistical Planning and Inference*, 136, 2407–2429.
- Quintana, F. A. and Newton, M. A. (1998), “Assessing the Order of Dependence for Partially Exchangeable Binary Data,” *Journal of the American Statistical Association*, 93, 194–202.
- (2000), “Computational aspects of Nonparametric Bayesian analysis with applications to the modeling of multiple binary sequences,” *Journal of Computational and Graphical Statistics*, 9, 711–737.
- Relling, M.V. and Boyett, J., et al. (2003), “Granulocyte colony-stimulating factor and the risk of secondary myeloid malignancy after etoposide treatment,” *Blood*, 101, 3862–3867.

Ross, S. M. (2002), *Introduction to Probability Models, Eighth Edition*, Academic Press.

Storey, J. (2002), “A direct approach to false discovery rates,” *Journal of the Royal Statistical Society, Series B (Methodological)*, 64, 479–498.

Walker, S. G., et al. (1999), “Bayesian Nonparametric Inference for Random Distributions and Related Functions (Disc: P510-527),” *Journal of the Royal Statistical Society, Series B: Statistical Methodology*, 61, 485–509.

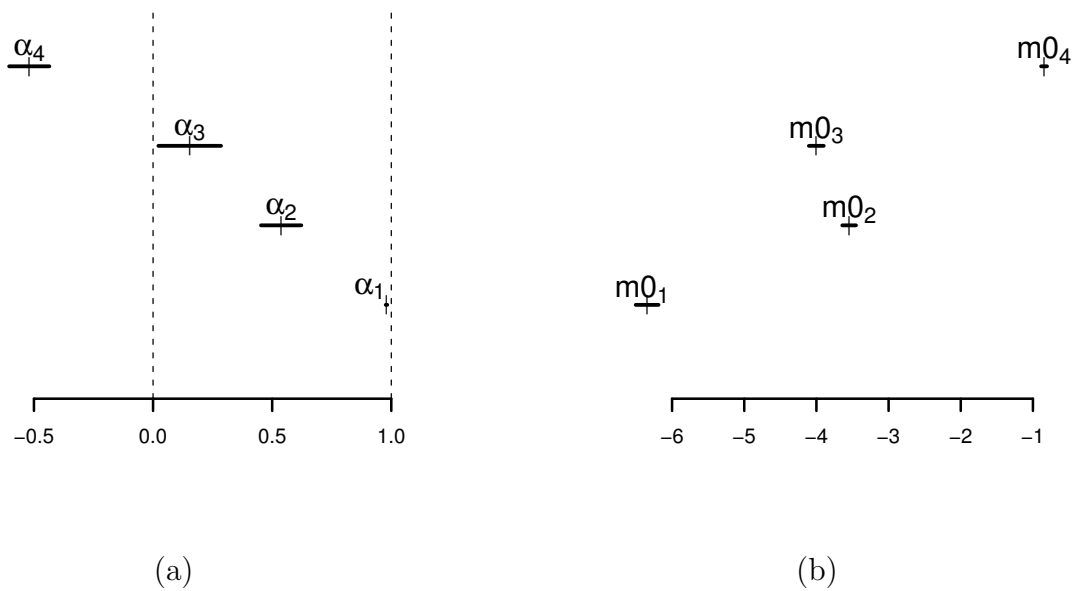
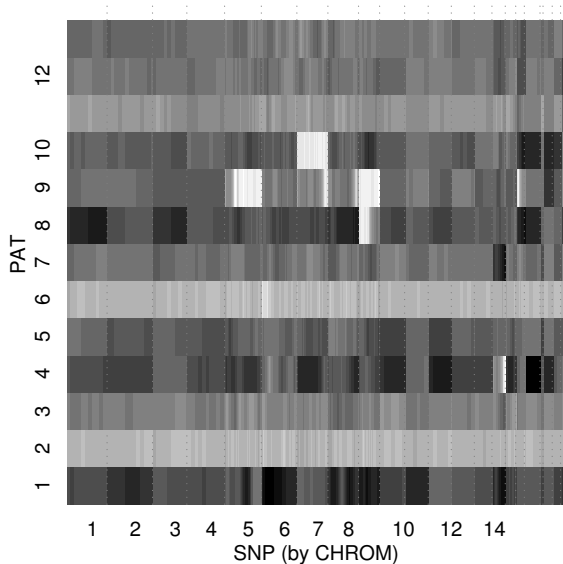
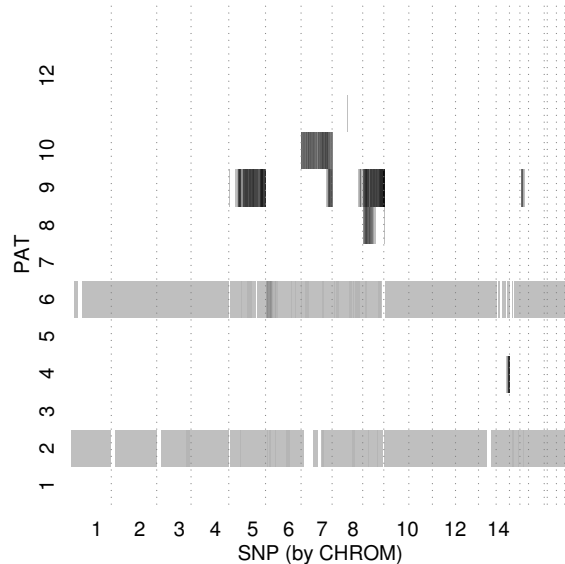


Figure 1: Marginal posterior means and standard deviations for (a)  $\alpha$  and (b)  $\mathbf{m}_0$ . The horizontal bars show the marginal posterior mean (a)  $E(\alpha_\ell | Y)$  and (b)  $E(m_{0\ell} | Y)$  (marked by “|”) plus/minus one posterior standard deviation for  $\ell = 1, \dots, 4$

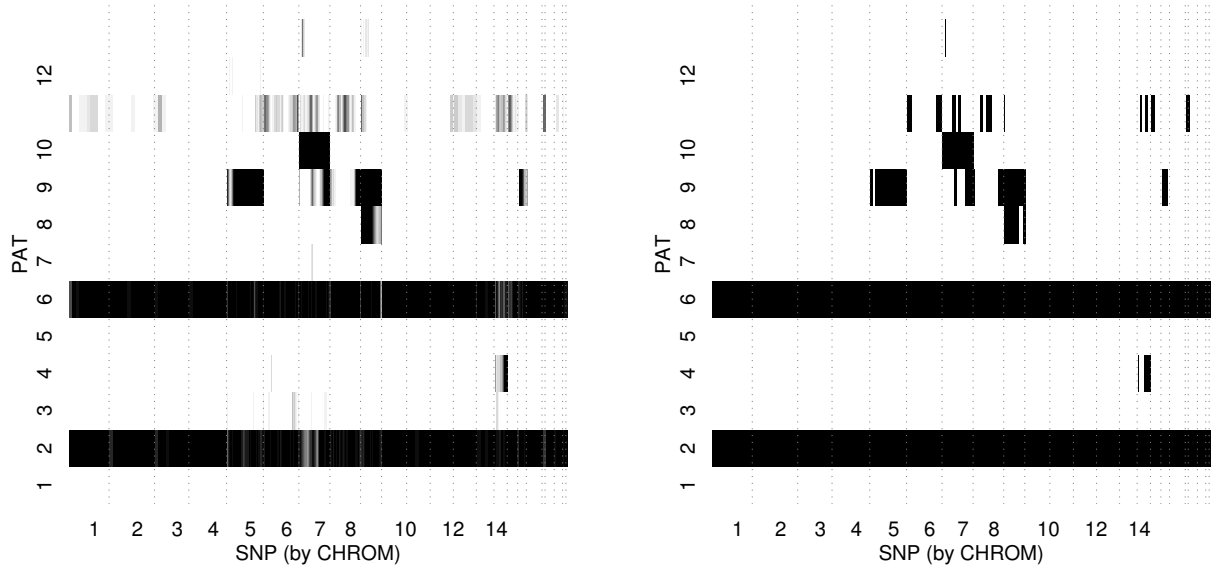


(a)  $(0, 0) \rightarrow (0, 0)$  transition probability



(b) probability of LOH

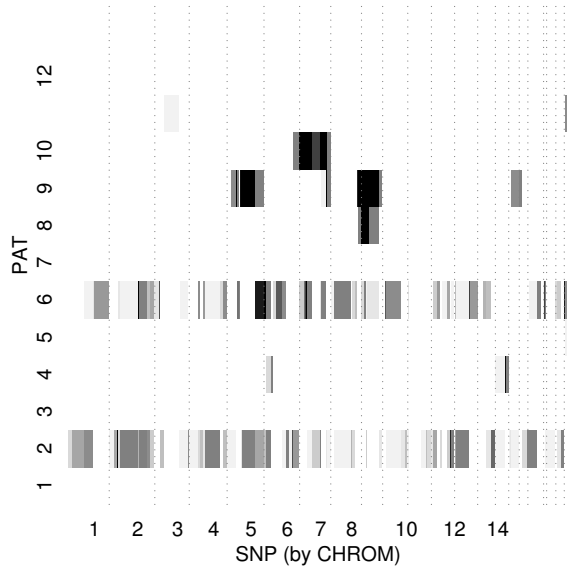
Figure 2: The left panel reports the marginal posterior means of the transition probability from state “00” to “00”, on a logit scale. The right panel shows the marginal posterior probability of LOH. Values are indicated by gray shades with black corresponding to 1.0. From left to right, columns represent regions 1 through 874. Chromosome boundaries are indicated by dotted vertical lines, and labeled for chromosomes 1 through 14. The rows correspond to patients 1 through 13. The right panel shows the probability of LOH, again arranged by chromosomes.



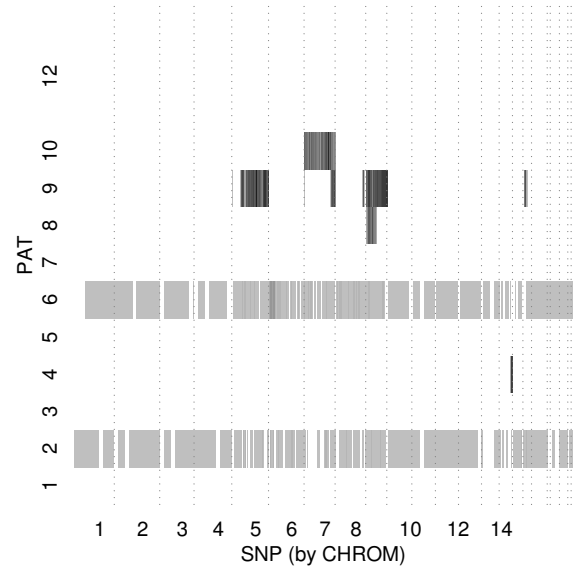
(a)  $Pr(I_{icj} | data)$

(b)  $\delta_{icj}$

Figure 3: The left panel reports the marginal posterior probability of increased LOH,  $P(I_{icj} | data)$ , with darker gray shades for higher probabilities. On the horizontal axis are regions arranged by chromosomes, as in Figure 2. The right panel shows the decisions  $\delta_{icj}$ , with a black bar indicating that region  $(ij)$  is reported as increased LOH.



(a) dChip



(b) parametric model

Figure 4: *Probability of LOH under two alternative methods. Panel (a) shows the inference reported by dChip using the hidden Markov model scores. Panel (b) shows inference under a fully parametric model. In both panels the horizontal axis shows regions arranged by chromosomes. Compare with Figure 2b.*