

Accounting for Disease Model Uncertainty in Mapping Heterogeneous Traits – A Bayesian Model Averaging Approach

Swati Biswas^a Charalampos Papachristou^b

^aDepartment of Biostatistics, School of Public Health, University of North Texas Health Science Center, Fort Worth, Tex., and ^bDepartment of Mathematics, Physics, and Statistics, University of the Sciences in Philadelphia, Philadelphia, Pa., USA

Key Words

Linkage analysis • Locus heterogeneity • Markov chain Monte Carlo • Reversible jump • Alzheimer's disease

Abstract

Background: Locus heterogeneity, wherein a disease can be caused in different individuals by different genes and/or environmental factors, is a ubiquitous feature of complex traits. A Bayesian approach has been proposed to account for variable rates of heterogeneity across families in a parametric linkage analysis setup [Biswas and Lin: J Am Stat Assoc 2006;101:1341–1351]. As with any parametric approach, its application requires specification of the disease model, which limits its practical utility. **Methods:** We address this limitation by proposing a Bayesian model averaging (BMA) approach. We consider a finite number of disease models and treat the model as an unknown parameter. In practice, we use simple single-locus disease models as various categories for model. **Results:** Our simulations as well as analysis of Genetic Analysis Workshop 13 simulated data show that BMA retains at least 80% of the power that is obtained by analyzing under the true disease model. The coverage probability of interval for disease gene is maintained around the nominal level. Finally, we apply BMA to a Late-

Onset Alzheimer's Disease dataset and find evidence for linkage on chromosomes 19, 9, and 21. **Conclusion:** We conclude that the BMA approach utilizing simple single-locus models for averaging is effective for mapping heterogeneous traits.

Copyright © 2010 S. Karger AG, Basel

Introduction

Locus heterogeneity is a ubiquitous feature of complex genetic traits and is a major inhibiting factor in the success of linkage analysis for localizing susceptibility genes. It refers to the situations when a disease can be caused in different individuals by different genes and/or environmental factors. Unless the presence of heterogeneity is properly modeled in the linkage analysis, the signals for linkage can get washed out. Biswas and Lin [1] have proposed a Bayesian approach to account for heterogeneity in a parametric linkage analysis setup. A distinctive feature of this approach is that it allows for varying rates of heterogeneity across families by means of family-specific random effects, denoted by α_j , which represent the probability that the j -th family is of a linked type, i.e. its disease-segregating gene is linked to the map of the marker(s)

under study. These α_j 's are assumed to follow a common (informative or non-informative) Beta prior distribution. For a sample of k families, the overall mixture likelihood is written as

$$L(\boldsymbol{\alpha}, d | \mathbf{x}) = \prod_{j=1}^k [\alpha_j L_j(d | x_j) + (1 - \alpha_j) L_j(\infty | x_j)], \quad (1)$$

where d is the position of the disease gene on the chromosome, and x_j and $L_j(d | x_j)$ denote the observed data and the corresponding (homogeneity) likelihood, respectively, of the j -th family, $j = 1, \dots, k$. Also, $\mathbf{x} = (x_1, x_2, \dots, x_k)$ and $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_k)$.

This approach was shown to be more powerful than the commonly used heterogeneity LOD (HLOD) score approach in detecting linkage [1, 2]. More importantly, the Bayesian approach yields interval estimates with known statistical properties as opposed to the support intervals customarily reported with HLOD, whose properties are unclear in the presence of heterogeneity. Recently, Biswas and Lin [3] extended this approach to incorporate covariates that may partially aid in discerning between linked and unlinked families. They showed that incorporating covariates leads to higher powers and narrower interval estimates compared to when no covariates are utilized as in Biswas and Lin [1].

Notwithstanding the obvious advantages of the Bayesian paradigm and its potential for further useful extensions, a fundamental requirement tied to its application is the specification of the disease model (disease allele frequency and penetrances) as is the case with any parametric approach. It is well known that misspecification of the disease model can lead to a loss of power, the extent of which depends on the type and degree of misspecification. Although nonparametric methods do not require model specification, they are statistically equivalent to a LOD score analysis that assumes recessive inheritance [4]. Moreover, the nonparametric methods are less powerful than their parametric counterparts when the latter are applied using approximately correct models. The burden of specifying approximately correct models for parametric methods is, in fact, less daunting than it appears since for LOD score analysis, simple one-locus genetic models can serve as good approximations to underlying complex multi-loci models as long as the mode of inheritance is specified approximately correctly at the linked locus [5]. That is, the full disease model, which undoubtedly is complex for complex traits, does not play as critical role as the assumed model for the particular disease locus under investigation.

The need for specification of the disease model in the Bayesian approach for heterogeneity proposed by Biswas and Lin [1] limits the practical utility of this approach. Not only is there arbitrariness in specifying the disease model, the uncertainty in model specification is also ignored. To address this practical limitation, we propose to utilize a Bayesian model averaging (BMA) approach. A literature review of BMA and its applications can be found in Hoeting et al. [6]. BMA has been applied successfully in many areas of applications of statistics resulting in an improvement in performance (see www.research.att.com/~volinsky/bma.html for details of these applications, related papers, and software). Applications of BMA can also be found in some areas of genetics such as association mapping [7] and mapping of quantitative trait loci [8].

In the context of the Bayesian framework as proposed in Biswas and Lin [1], applying BMA is relatively straightforward. We can consider a finite number of disease models and treat the 'model' as an unknown parameter that can take a finite number of categories. Although, in principle, BMA allows us to consider an infinite number of models by varying one or more components of the disease model over their whole or part of the parameter space, it is not feasible from a practical point of view. The disease model contributes in the mixture likelihood equation (1) via the homogeneity likelihood $L_j(d | x_j)$, and this likelihood is not usually available in a functional form. Rather it is computed at a grid of pre-specified d values using software such as GENEHUNTER [9] and Allegro [10] for a given disease model. Calculating the values of the homogeneity likelihood at all these locations for more than a handful number of disease models is impractical. However, this should not pose a major limitation since, as discussed above, the LOD score analysis is robust to the specific values of penetrances as long as an approximately correct mode of inheritance (such as dominant or recessive) is assumed at the locus linked to the marker [5]. By the same token, it may be adequate to consider only some representative disease models without worrying about the true specific values of the model parameters. An attractive feature of model averaging is that the hypothesis test is conducted only once after averaging because of which there is no need for multiplicity adjustment unlike when separate tests are conducted with each model.

In the following, we consider the Bayesian approach of Biswas and Lin [1] and incorporate averaging over a finite number of disease models. After describing the general methodology for an arbitrary number of models, we will

describe some specific models that we will use in applications. Then we will present results from a simulation study wherein the true underlying disease models are single-locus models. Next, in order to investigate the properties of BMA for complex models with multiple interacting loci, we analyze all 100 replicates of the Genetic Analysis Workshop 13 (GAW13) simulated data that were generated to mimic the real data from the Framingham Heart Study [11]. Finally, we apply BMA to the National Institute on Aging's (NIA) Late-Onset Alzheimer's Disease (LOAD) data obtained from NIH's database of genotypes and phenotypes (dbGaP) [12].

Methods

General Methodology

We begin by considering the likelihood in (1) and α, x , and d as defined before. The homogeneity likelihood $L_j(d|x_j)$ in (1) implicitly depends on the disease model. To explicitly bring the disease model in the picture, let us denote it as m and its index by I_m . Suppose there are M models under consideration, then $I_m \in \{1, \dots, M\}$. In this setup, we have three sets of parameters – the α and d as before, and additionally the parameter m . Also, we denote the overall mixture likelihood, $L(\alpha, d)$ in (1) as $L(\alpha, d, m)$ to explicitly show the contribution of the disease model.

Next in the Bayesian framework, we define the prior distributions for the three sets of parameters. The prior for α_j is denoted as $\pi_j(\alpha_j)$, which is taken to be the non-informative $U(0, 1)$. For d , an important and crucial distinction is made between two models of different dimensionalities – linkage consisting of $k + 2$ parameters (α, d, m) and no linkage with no parameters as $d = \infty$ renders α and m meaningless. The respective prior probabilities of these models are denoted by $\pi_{d<\infty}$ and $\pi_{d=\infty}$. Further, under linkage (i.e. $d < \infty$), we have a probability distribution on all possible (discretized) values of d , the disease gene location on the chromosome (or map) under study. If we denote these locations as $1, \dots, N$, and let I_d denote the index of d values, then we have $I_d \in \{1, \dots, N\}$. The probability distribution of d under linkage is then defined on these N locations and is denoted by $\pi_d(I_d)$. For our analyses, we use $\pi_{d<\infty} = 1/22$, $\pi_{d=\infty} = 21/22$ and $\pi_d(I_d) = 1/N$, $I_d = 1, \dots, N$. This prior assigns any given chromosome (autosome) a prior probability of linkage of $1/22$ and each of the N positions on the chromosome is assigned a probability of $1/N$. Next, we specify a prior distribution for I_m , which is denoted as $\pi(I_m)$. Unless there is some a priori knowledge available about the disease model (e.g. by segregation analysis), a natural non-informative prior would be discrete uniform, i.e. $\pi(I_m) = 1/M$, $I_m = 1, \dots, M$. Now we can write down the joint posterior distribution of parameters as

$$\pi(\alpha, d, m|x) \propto \left[\pi(I_m) \prod_{j=1}^k \pi_j(\alpha_j) \right]^{I(d<\infty)} \times \left[\pi_{d<\infty} \pi_d(I_d) I(I_d \in \{1, \dots, N\}) + \pi_{d=\infty} I(d=\infty) \right] L(\alpha, d, m|x),$$

where $I(\cdot)$ is the indicator function and $L(\alpha, d, m|x)$ is the same as in (1) with the homogeneity likelihood computed using model m .

Our goal is to estimate the posterior distribution of d so that inference regarding linkage can be conducted, and if linkage is inferred, point and interval estimates for the location of the disease gene can be obtained. This is accomplished through Markov chain Monte Carlo (MCMC) methods. Since $d < \infty$ (linked: L) and $d = \infty$ (unlinked: U) are subspaces with different numbers of parameters, the sampler that we employ should allow moves between subspaces of varying dimensionalities. So we use the reversible jump MCMC algorithm [13]. At each iteration, the Markov chain can be currently in either L or U subspace and a proposal will be made to either remain in the current subspace or move to the other subspace, leading to four possible move types: $L \rightarrow L$, $L \rightarrow U$, $U \rightarrow U$, and $U \rightarrow L$. Details of these moves can be found in the Appendix.

The posterior distributions are obtained by running a large number of iterations after a burn-in period. From the estimated posterior distribution of d , we first calculate the posterior probability of linkage, \hat{p} , i.e. the proportion of times the chain stays in the L subspace. The obtained value of \hat{p} is then converted into an estimated Bayes Factor (BF) given by

$$\widehat{BF} = \frac{\hat{p}/(1-\hat{p})}{\pi_{d<\infty}/(1-\pi_{d<\infty})}.$$

If the \widehat{BF} exceeds a pre-specified threshold BF_0 , we infer linkage. In such a case, we can easily obtain a point and an interval estimate for the location of the gene d by using the mean and the 95% credible set of the posterior distribution of d under linkage, i.e. when the chain stays in the L subspace.

For all our analyses, the chain is run for 300,000 iterations after a burn-in period of 10,000 iterations. We use $BF_0 = 25$, which corresponds to a strong linkage signal following Raftery [14], and was also used by Biswas and Lin [1]. Since no multiplicity adjustment is needed for model averaging, we kept the same threshold.

Note that the BF we estimate here is the BF for testing linkage in our mixture setting. It actually amounts to testing the number of components in the mixture model – two (corresponding to linkage under heterogeneity with prior $\pi_{d<\infty}$) versus one (corresponding to no linkage with prior $\pi_{d=\infty}$). Working in this kind of mixture setting and testing for the number of components, Richardson and Green [13] comment that the BF is theoretically independent of the prior used for the number of components. They also show it empirically, and similar empirical results were obtained by Thomas et al. [15], Biswas and Lin [1], and Wang et al. [16]. In particular, Biswas and Lin [1] obtained the BF using three different values of prior, $\pi_{d<\infty}$, namely, $1/22$ (the same as in this article), $1/\text{length of chromosome}$ (another non-informative prior), and 0.1 (an informative prior), and found them to be similar. This feature of insensitivity of the BF to the prior for the number of components is attractive as we do not need to reference the prior used in the final inference using the BF.

Specific Models Used for Model Averaging

For the analyses using BMA, we used the single-locus disease models listed in table 1. The disease allele frequency in these models was chosen in such a manner that the disease prevalence is around 10%. Note that these models can be classified as either of dominant, recessive, or intermediate type. We considered various combinations of these models to investigate the perfor-

mance of BMA (described in the next section). The first four models were constructed with a non-phenocopy penetrance value of 0.5, which appears to be a natural choice in the absence of any additional information. Since for an intermediate model this value may be assigned to either genotype DD or Dd , where D (d) is the disease (wild-type) allele, we constructed two such models.

Simulations

Simulation Study Setup

To investigate the performance of the proposed BMA approach, we conducted a simulation study with data generated from the single-locus disease models listed in table 2. Specifically, we chose two models of each type – dominant, recessive, and intermediate – with values of penetrances and disease allele frequencies chosen to reflect various scenarios in which linkage analysis may be helpful in mapping genes. The prevalence corresponding to these models are not all close to 10% as in the BMA analysis models given in table 1. Thus, results of this simulation will also show the robustness of the approach to misspecification of disease prevalence reflected in penetrance and disease allele frequency values.

The pedigree structures used for this simulation study were chosen from the GAW13 simulated data in such a way that they can be analyzed by software GENEHUNTER [9] and Allegro [10] either as a whole or after some members are trimmed. There were a total of 274 pedigrees in a sample, whose size ranged from 7 to 19 members (mean size = 11.2). The marker map used is the chromosome 21 map provided in GAW13 consisting of six markers at locations 0, 10.02, 22.74, 36.20, 40.07, and 59.53 cM (locations adjusted so that the first marker is at 0 cM). The disease gene location is $d = 25.64$ cM. We simulated marker genotypes and disease phenotypes according to the following scheme. For each family, we first decided whether it is of a linked type, where the probability of being linked is randomly chosen from beta(3,2) distribution. Note that the mean and SD of this distribution is 0.6 and 0.2, respectively, i.e. about 60% of the families in a sample are expected to be of a linked type with a fair amount of variability across families in probabilities of a linked type. We believe this reflects a realistic heterogeneity setting. The same disease model was used for generating linked and unlinked families, with the unlinked families having the disease gene unlinked with the marker map. The ascertainment criterion was the existence of at least two affecteds in a pedigree. The affection statuses of about 45% of the people in each sample was labeled

Table 1. Single-locus disease models used in various combinations as described in the text for the analysis with the BMA approach

Model No.	f	p_{DD}	p_{Dd}	p_{dd}
1	0.08	0.5	0.5	0.03
2	0.4	0.5	0.03	0.03
3	0.14	0.5	0.3	0.03
4	0.06	0.7	0.5	0.05
5	0.04	0.7	0.7	0.05
6	0.12	0.3	0.3	0.05
7	0.3	0.7	0.05	0.05
8	0.45	0.3	0.05	0.03
9	0.35	0.3	0.1	0.05

f represents the disease allele frequency; D denotes the disease allele; p_{DD} , p_{Dd} , p_{dd} are the penetrances of the genotypes DD , Dd , and dd , respectively. All of these models correspond to a prevalence of around 0.1.

Table 2. Single-locus disease models used for simulating data

Simulation scenario	f	p_{DD}	p_{Dd}	p_{dd}	Prevalence
Dom1	0.1	0.3	0.3	0.01	0.066
Dom2	0.3	0.6	0.6	0.01	0.311
Rec1	0.2	0.7	0.05	0.05	0.076
Rec2	0.3	0.4	0.03	0.03	0.063
Inter1	0.2	0.8	0.2	0.03	0.115
Inter2	0.1	0.6	0.4	0.05	0.118

f represents the disease allele frequency; D denotes the disease allele; p_{DD} , p_{Dd} , p_{dd} are the penetrances of the genotypes DD , Dd , and dd , respectively.

missing; these are the same people whose affection statuses were treated as missing by Biswas et al. [18] in their analysis of a GAW13 replicate.

We used Allegro [10] to compute the homogeneity LOD scores at every 1 cM using the analysis models listed in table 1. These LOD scores were then inputted to our program that implements the BMA approach. We compared the performances of the BMA approach under a variety of settings that averaged over different numbers and/or types of disease models (subsets of models in table 1). These settings are: (1) $M = 2$ with models 1 and 2 (MA2); (2) $M = 3$ with models 1, 2, and 3 (MA3); (3) $M = 4$ with models 1, 2, 3, and 4 (MA4); (4) an alternative $M =$

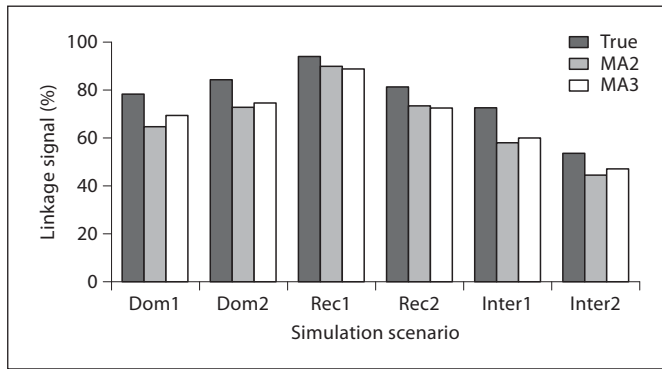


Fig. 1. Simulation results: power comparison (in %) of analyses under the true model (True) and BMA with two (MA2) and three models (MA3) based on 1,000 simulations. Power is defined as the percentage of replicates that gave linkage signal ($\widehat{BF} > 25$) when linkage is present in the data.

Table 3. Power comparison (in %) of True, MA2, MA3, MA3Alt, and MA4 based on 500 simulations

Simulation scenario	True	MA2	MA3	MA3Alt	MA4
Dom1	78.0	66.0	70.0	68.0	70.8
Dom2	80.6	73.4	75.0	75.0	75.6
Rec1	92.8	88.4	87.2	87.0	86.0
Rec2	81.2	74.0	73.2	71.4	71.2
Inter1	71.8	56.8	58.0	57.6	59.2
Inter2	54.6	45.2	47.4	49.4	50.6

Power is defined as in figure 1.

Table 4. Power comparison (in %) of True, MA2, MA3, MA6, and MA9 based on 500 simulations

Simulation scenario	True	MA2	MA3	MA6	MA9
Dom1	78.6	63.4	68.8	56.4	60.6
Dom2	88.0	72.2	74.2	66.8	70.8
Rec1	95.2	91.4	90.4	89.8	90.6
Rec2	81.4	72.8	71.8	63.4	68.2
Inter1	73.4	59.2	62.0	45.2	54.8
Inter2	52.6	43.8	46.8	38.2	43.0

Power is defined as in figure 1.

3 with models 1, 2, and 4 (MA3Alt); (5) $M = 6$ with models 4 through 9 (MA6), and (6) $M = 9$ with all models (MA9). In addition to these BMA settings, we also analyzed the data using the true generating model for comparison (referred to as ‘True’). To enable comparisons between all of the above-mentioned settings while at the same time keeping the computational load at reasonable levels, we adopted the following strategy. For each true generating model (simulation scenario) in table 2, we conducted 1,000 replications comparing MA2 and MA3 with the True model. In 500 of these 1,000 replicates, MA3Alt and MA4 were also compared while in the remaining 500, MA6 and MA9 were computed. This strategy is also guided by one of our goals to find the most parsimonious yet robust set of models for BMA. We compared the observed powers, false-positive rates, interval widths, and their coverage probabilities obtained using the BMA approach with the ones analyzed under the True model. For gauging the false-positive rates, separate sets of 1,000 simulations each were carried out under each simulation scenario in table 2 with all families taken to be of an unlinked type.

Simulation Study Results

The observed powers of MA2 and MA3 are compared with that of the True model in figure 1 (based on 1,000 simulations) and in tables 3 and 4, each based on 500 simulations (first 3 columns). The BMA approach retains at least 80% of the power obtained using the True model. For example, in table 3, for Dom1, MA2 retains $66/78 = 84.6\%$ of the power from True. Comparing settings MA2 and MA3, we see that both yield similar results, with MA3 giving slightly better power when the generating disease model is dominant or intermediate. Table 3 also shows the results for MA3Alt and MA4. These analyses give similar results as MA2 and MA3, overall. MA4 gives higher powers than MA3 and especially MA2 for intermediate models. In general, BMA seems to be robust to the choice of models as well as the number of models in the range of 2 to 4 models.

However, as the number of models increases, this pattern changes. Table 4 shows the results for MA6 and MA9. Comparing these results with the ones for MA2 and MA3, we see that both MA6 and MA9 have inferior performance (except for Rec1 which is discussed in the next paragraph), with MA6 being the worst. MA6 averages over models that have nonphenocopy penetrances of either 0.7 and 0.3 (for intermediate models, this is the penetrance of the DD genotype) while MA9, in addition to these models, also averages over models with this value

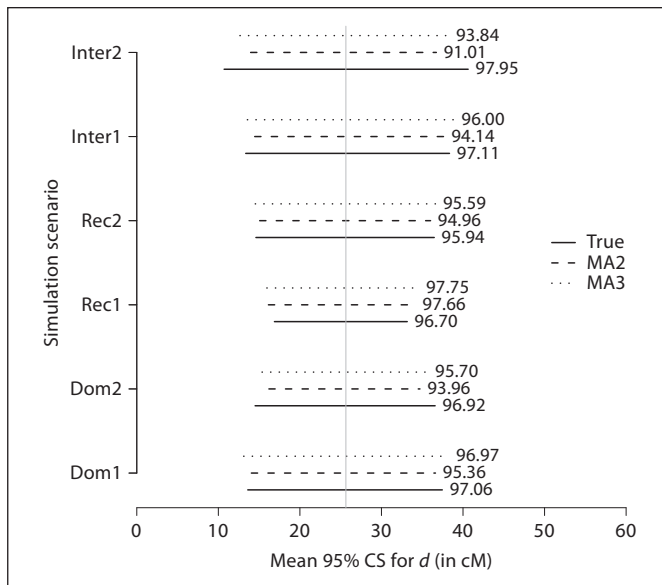


Fig. 2. Simulation results: mean 95% CS for d (averaged over replicates that showed linkage) under the True model, MA2, and MA3. Coverage probability is reported at the end of each mean CS. The vertical line shows the true location of the disease gene.

set to 0.5. Combining this observation with the results from MA2 and MA3, which are subsets of MA9 but not MA6, we may infer that one model with disease genotype frequency of 0.5 is better in extracting linkage signal compared to two models with frequencies of 0.3 and 0.7, at least in the simulation scenarios investigated here. However, surprisingly, even though MA9 is a superset of both MA2 and MA3, it performs worse than both of them, indicating that averaging over a larger number of models is not necessarily helpful. We discuss this seemingly counterintuitive result further in the Discussion section.

The setting of Rec1, however, is slightly different in the sense that MA6 seems to perform comparable with MA2 and MA3 in table 4, and perhaps better than MA4 from table 3. However, since MA4 and MA6 were calculated on different sets of 500 simulations, the minor difference in their results could be due to sampling variability. So we performed an additional set of 500 simulations for the Rec1 setting and computed MA2, MA3, MA4, and MA6. The powers obtained (in %) were 86.8, 85.8, 83.8, and 83.4, respectively. Thus MA4 and MA6 perform comparably, and slightly inferior to MA2 and MA3.

Next, in figure 2, we compare the 95% credible sets (CSs) obtained using the True, the MA2, and the MA3.

Table 5. False-positive rates of True, MA2, MA3, and MA4 based on 1,000 simulations

Simulation scenario	True	MA2	MA3	MA4
Dom1	0.003	0.004	0.004	0.004
Dom2	0.009	0.006	0.007	0.007
Rec1	0.002	0.001	0.001	0.002
Rec2	0.005	0.003	0.003	0.002
Inter1	0.004	0.004	0.005	0.005
Inter2	0.002	0.004	0.003	0.003

False-positive rate is defined as the proportion of replicates that gave linkage signal ($\widehat{BF} > 25$) when there is no linkage present in the data.

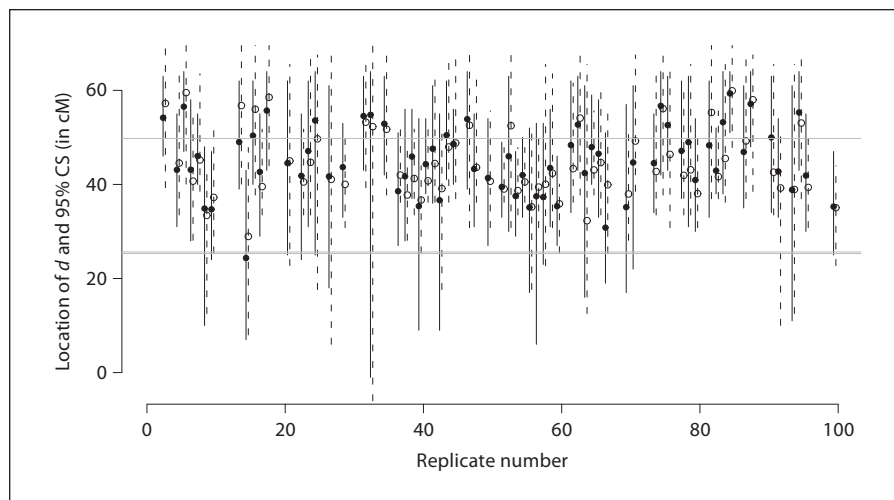
We see that, in most cases, the lengths of CSs as well as their coverage probabilities are slightly smaller with the BMA, especially with MA2 compared to True, except most notably for Rec1. However, the coverage probabilities of all BMA CSs are around the nominal value of 95%, with MA3 slightly overcovering compared to MA2. As a result, CSs with MA3 are slightly longer too.

Finally, with regard to the false-positive rates, MA2, MA3, MA4, and True are comparable (table 5). As can be seen in table 5, for some models one approach gives smaller rates while for others it is vice-versa. So there is no clear winner in this respect. Notice though, that all false-positive rates of the BMA approaches are small in absolute scale.

Application to the Simulated Data from GAW13

The above simulation study indicates that the proposed BMA approach performs well when the underlying disease model is single-locus, even under misspecification of prevalence. In this section, we evaluate the performance of our method in a more complex and realistic setting where several loci interact to cause a disease. For this purpose, we analyze the simulated data provided in GAW13. These data were generated to mimic the real Framingham Heart Study data. There are 50 trait loci spread over 22 autosomal chromosomes, which interact through an extremely complex model to produce a number of traits. A description of the model can be found in Daw et al. [17]. We focused on mapping genes that influence the high blood pressure (HBP) trait. We followed the

Fig. 3. Point estimates and CS for d obtained from analyses of GAW13 replicates using BMA with MA2 (solid lines and circles) and 'True' model (dashed lines and open circles). Results from only those replicates are plotted that gave linkage signal under both analyses. The true locations of the three disease genes are shown as horizontal lines (two genes are very close to each other).



criteria of Biswas et al. [18] in labeling a person as affected by HBP, which combines the longitudinal data on blood pressures and a closely related variable, hypertensive treatment. Several genes over several chromosomes affect this trait directly or indirectly through other traits. We chose to apply our approach to chromosomes 20 and 21. Chromosome 21 has three genes that directly affect HBP. The marker data for this chromosome are the same as described in the 'Simulation Study Setup' subsection. The three genes lie at 25.31, 25.64, and 49.77 cM and are named as b37, s12, and s10, respectively. Note that b37 and s12 lie between markers 3 and 4 and are indistinguishable through linkage analysis while s10 lies between markers 5 and 6. Chromosome 20 has no disease gene and it is used to gauge the false-positive rate. One hundred replicates of the simulated data were provided in GAW13, each replicate consisting of 330 pedigrees.

The homogeneity LOD scores were obtained using GENEHUNTER [9] with its maxbits option set at 18. GENEHUNTER dropped individuals automatically from large pedigrees, which, although not ideal, is a practically feasible option to analyze all 100 replicates in a timely fashion using all settings of BMA. Since the true disease model is complex involving multiple trait-contributing loci, for comparison, we also analyzed the data under an approximate marginal disease model of gene s10. Specifically, we averaged the generating models of s10 for the systolic and diastolic blood pressures to get the following model: $f = 0.3$, $p_{DD} = 0.8$, $p_{Dd} = 0.4$, and $p_{dd} = 0$. Throughout the rest, we will be referring to this model as the 'True' model although it is, in fact, close to the true marginal model only. This single-locus model corre-

sponds to a disease prevalence of 24%, which is much higher than the 9–10% used in the models for the BMA; thus, we have misspecification of prevalence in the BMA approach.

We first discuss the results for chromosome 21. The True model gives linkage signal in 82 replicates while MA2, MA3, MA3Alt, MA4, MA6, and MA9 give signal in 72, 71, 72, 71, 62, and 69 replicates. So the BMA analysis with two to four models retains about 87–88% of the power obtained with the analysis under the True model. Moreover, consistent with the results of our earlier simulation study, including more models lead to a decrease in the power.

Figure 3 shows a side-by-side plot of the 95% CSs and associated point estimates of d obtained from the analyses with MA2 and the True model for the replicates that gave linkage signal under both analyses. Notice that in most of the replicates, the gene s10 is captured in the intervals. This is consistent with the simulating model of GAW13 wherein s10 has much stronger effect on blood pressure compared to the other two genes. Also, most of the point estimates lie between s10 and s12/b37. This is expected, since the analysis assumes single disease gene on the chromosome while there are two genes that can be detected by linkage analysis. Consequently, the estimates lie in between those two gene locations. Further inspection of the figure reveals that s12 and b37 are detected more often by MA2 than by True. This is again not unexpected, since the True model is close to the s10's model, while MA2, by allowing for varied models, is able to capture more often the genes that are more difficult to map. In terms of average length and coverage probability of the

CSs, MA2 and MA3 performed as good as the True model, even marginally better. We calculated these quantities over the replicates that gave signal by the corresponding method. For coverage probability, if an interval contains any of the three genes, it is considered as a success. The average length and coverage probability of MA2 is 26.85 cM and 94.4%, respectively, while for the True analysis, the corresponding values are 30.72 cM and 95%. For MA3, these values are 28.59 cM and 97.2%.

Analysis of chromosome 20 under the True model and BMA (MA2 and MA3) give 0 false positives, consistent with results from earlier simulation study. Overall, we conclude that BMA performs satisfactorily in complex settings like this one.

Application to the LOAD Data

This dataset is described and analyzed in Lee et al. [12]. It consists of 6,000 single nucleotide polymorphisms at an average intermarker distance of 0.65 cM. The authors conducted both single-point and multipoint linkage analyses, as well as family-based and case-control association studies on these data. The linkage analyses were carried out using the Kong and Cox nonparametric method [19]. Using the broad definition of LOAD according to which persons diagnosed with definite, probable, or possible LOAD are labeled as affected, and the multipoint analysis, Lee et al. [12] got a strong linkage signal on chromosome 19 and suggestive evidence for linkage on few other chromosomes where the LOD scores exceeded 2. The signal on chromosome 19 was near the apolipoprotein E (APOE) gene, which has been consistently replicated in several LOAD data [20].

Following Lee et al. [12], we restricted our analyses to white families as they constitute more than 90% of the cohort. We first checked for Mendelian inconsistencies using Merlin [21] specifically based on double recombinations. The marker allele frequencies were estimated from the data. We had a total of 166 extended families (total number of individuals = 1,829) with mean and SD of the number of members equal to 11 and 5.1, respectively. We used the broad definition of LOAD and obtained the multipoint parametric LOD scores of all families using Merlin. Some families were too big to be analyzed by Merlin as a whole, so we manually dropped individuals and reduced the pedigree size before analyzing them in Merlin. All chromosomes were analyzed using the same combinations of models for BMA as in our simulation study.

We found evidence for linkage (i.e. $BF > 25$) on chromosomes 19, 9, and 21. The strongest signal is on chromosome 19; the BFs (posterior probability of linkage, \hat{p}) for MA2, MA3, MA3Alt, MA4, MA6, and MA9 are 35.02 (62.51%), 29.22 (58.19%), 41.23 (66.26%), 32.20 (60.53%), 29.03 (58.03%), and 29.90 (58.74%), respectively. The 95% CS for MA2, MA3Alt, and MA4 is around (67, 111) cM and contains the APOE gene. We also calculated the 95% Highest Posterior Density set and it contains two disjoint sets at (67, 91) cM and (110, 111) cM. The second narrow region could be just the effect of linkage to the APOE gene or it may indicate a second disease locus.

The evidence on chromosomes 9 and 21 is less strong than that on chromosome 19. For chromosome 9, the BF from MA2 and MA3 exceeded the cutoff with MA2's $BF = 32.87$ and $\hat{p} = 61\%$ while for chromosome 21, only MA2 showed evidence with $BF = 29.84$ and $\hat{p} = 59\%$. For both chromosomes, averaging over more models gave worse results except that MA9 gave better results than MA6, consistent with our simulation results. Also in both cases, the recessive model was visited more often by the MCMC than the dominant and intermediate, and thus, in agreement with our simulation results, MA3 performed worse than MA2. These chromosomes have been implicated earlier [20, 23–25]. However, Lee et al. [12] did not find any linkage evidence (even suggestive) on these chromosomes using the same data. So our study replicates these earlier findings using a different and independent dataset. The 95% CS for chromosome 9 is (40.22, 52.28) cM, which contains the linkage peak obtained by Scott et al. [24] and overlaps with their 1-LOD support interval. The 95% CS for chromosome 21 is very narrow (44.82, 44.95 cM). Two regions in the vicinity of this location on 21q have been implicated earlier – the APP gene [20, 23] and a region around marker D21S1440 [20, 25]. So the linkage signal could be attributed to one or both of these loci. However, as Bacanu et al. [25] noted, the results on the location may not be robust due to discrepancies in the maps for the 21q region. Finally, we also found some suggestive evidence on chromosome 18 ($BF = 15.5$, $\hat{p} = 42.46\%$) consistent with Lee et al. [12].

Discussion

We proposed a BMA approach to account for uncertainty in the disease model in mapping heterogeneous traits. This approach retains at least 80% of the power that is obtained by analyzing under the True model. From all settings considered, MA2 and MA3 seem to

give better and comparable results, with MA3 having a slight advantage when the true model is dominant or intermediate. Thus, if there is some a priori belief that the true model for the disease locus under consideration may not be recessive, MA3, MA3Alt, or MA4 can be used. MA4 gives better results than MA2 when the true model is intermediate but is more computationally intensive. In general, if there is no a priori information about the disease model, we recommend using MA2 or MA3. The penetrances and allele frequencies may be varied in these models to reflect the prevalence of the disease; however, a penetrance of 0.5 for the disease genotype seems to be a robust choice (see further discussion little later). A notable advantage of averaging over models rather than analyzing separately under different models is that averaging gets around the issue of multiplicity adjustment since testing is conducted only once after averaging over all models. Moreover, BMA, by explicitly including a parameter for the disease model and its distribution, accounts for the uncertainty in the model choice leading to more robust estimates and inferences. We note that our approach of updating the disease model is similar in essence to the updating of QTL states in Hayashi and Awatha [22], wherein the authors considered four possible patterns of combinations of QTL states in each pair of F0 grandparents.

The analyses of the GAW13 simulated data revealed a notable advantage of model averaging when the True model is complex. The CS from the BMA approach may be able to capture those genes that may be missed by using one model only, which usually is constructed to approximate the true disease model for a particular disease gene. This is an advantage of using a spectrum of disease models rather than fixing a model.

The analysis of the LOAD data confirmed the role of the APOE locus on chromosome 19, which has been replicated using the same data previously [12] as well as elsewhere [20]. On the other hand, the linkage findings on chromosomes 9 and 21, which Lee et al. [12] did not find in these data, is worth revisiting. The major difference between the linkage mapping approaches of these authors and ours is that our approach explicitly accounts for heterogeneity and that might explain the discrepancy of the results on these chromosomes. Indeed, for chromosome 21, Olson et al. [23] obtained evidence for linkage only after they accounted for heterogeneity by utilizing covariates, most notably age of onset and disease duration. Similarly for chromosome 9, accounting for heterogeneity lead to an increase in LOD scores by Scott et al. [24]. So our results illustrate the well-known fact that ac-

counting for heterogeneity can help in uncovering disease genes that may be otherwise missed for complex traits like Alzheimer's disease. Finally, our study presents an independent replication of the involvement of these chromosomes in LOAD and, to the best of our knowledge, the first one using these NIA-LOAD data for chromosomes 9 and 21.

An interesting finding is that averaging over a larger number of models does not necessarily give better power. Both MA6 and MA9 performed inferior compared to MA2 and MA3. Moreover, given that MA9 is a superset of MA2 and MA3, it shows that just including additional models may not result in an increase in the power. A possible explanation could be that with a larger number of models, the chain does not cover the posterior distribution well because of the large sample space. Although running the chain for a longer time may alleviate this problem, the increase in the computation time can be hardly justified. As it is now, MA9 takes substantially longer time than MA3, mainly because LOD scores need to be computed at nine disease models, and running the chain longer to get similar power itself is not justified, let alone running even longer to potentially get higher powers. A related finding is the lowest power obtained with MA6, which shows that averaging over a small number of 'average' kind of models (penetrance of 0.5 for disease genotypes) is more effective than averaging over larger number of models with more extreme penetrance values such as 0.7 and 0.3. This conclusion is reinforced by the fact that MA9, which includes 'average' models in addition to 'extreme' models, gives higher powers than MA6. This last point of better performance of MA9 compared to MA6 also adds another aspect to the overall conclusion – just using larger number of models in averaging is not the only explanation for lower powers. The choice of models also plays a role in this regard. Overall, it appears that depending upon the true underlying model, BMA may be somewhat sensitive to the number and combination of models used for averaging. Although MA2 and MA3 seem to perform satisfactory, at least for the settings considered here, care should be taken in the choice of models in general.

Since model averaging gets around the need for fixing a disease model, it may be compared with nonparametric methods. To this end, we applied the Kong and Cox nonparametric method [19] on all 100 replicates of the GAW13 simulated data that we analyzed earlier with BMA. As before, chromosomes 20 and 21 were analyzed. The Kong and Cox method with a cutoff of 3.09 gave 35% power with 2 false positives. While recall that BMA gave around 71%

power with 0 false positives. Increasing the cutoff for the Kong and Cox method to get 0 false positives will make its power go down further. Although this does show the superiority of our parametric approach, it should not be overlooked that the Kong and Cox approach does not account for heterogeneity, which may be part of the reason for its lower power and higher false positives compared to BMA. Nevertheless, the results are consistent with the well-known fact that nonparametric methods are less powerful than parametric methods applied with approximately correct models, which is accounted for to some extent through averaging over models.

One important issue is the choice of models for averaging. Here we used models that were fixed in advance independent of the data. Although our results show that, in general, this strategy works well, a better strategy may be to use ‘data-driven’ models, i.e. models derived from and consistent with the data to be analyzed. Ideally, to avoid overfitting, one should split the dataset into two parts, and use one part for inferring models, which are then used to analyze the other part. However, this may not be feasible without a large dataset. Wan and Lin [26] derive models from the identical-by-descent (IBD) distribution estimate at the trait locus. To obtain sensible models corresponding to the estimated IBD distribution, they place constraints of $p_{DD} \geq p_{Dd} \geq p_{dd}$ and $f < 0.15$. They also considered perturbations of estimated IBD sharing distribution to account to some extent for uncertainty in the estimates. This method was applied to the rheumatoid arthritis data of GAW15 which consists of affected sib pairs. For data with general pedigrees, this approach can be used if there is a reasonable number of affected sib pairs to get some plausible estimates. Using such models ‘informed’ by the data rather than one-size-fits-all models can increase the power. Even when it is impractical to apply this strategy, one should try to use models tailored to specific application. For example, the models that we used for averaging correspond to a disease prevalence of 10%. If a disease is much rarer, one should definitely adjust the model parameters to reflect this information. Similarly, if there is a priori information about the mode of inheritance, that can be incorporated through a proper choice of models and/or through prior distribution of models used for averaging. Undoubtedly, any assumed model will be off from the true underlying model but by using a tailored model, one ensures that the models are not grossly different from and/or are infeasible with the available epidemiological information on the disease and thus increase power. Since all of our models used in averaging are off from the true

generating models and still BMA gives good power, it suggests we do not have to worry about less gross deviations. Note that the general methodology of the BMA approach is, however, completely independent of the choice of models. Model selection should be considered as an important step preceding the application of the BMA approach.

Here we viewed heterogeneity as a mixture of two groups of families – one group carrying the disease gene under consideration (having signal) while the other group consisting of families of all other types (not having signal). Alternatively, Sillanpaa and Bhattacharjee [27], in the context of an association study under regression set-up, view heterogeneity as the sample containing individuals with two competing disease mechanisms (two competing signals), each with its own set of trait-associated loci (see also [28]). Thus, it may be of interest to investigate how such competing mechanisms can be modeled in our mixture model framework for linkage analysis.

Finally, we close with a note on a possible extension of our approach to mapping multiple loci. Here we considered the mapping of a single disease gene at a time. However, the flexible framework of reversible jump MCMC can readily accommodate the mapping of several loci simultaneously with the number of loci unknown in advance. In fact, several approaches in these lines have been proposed in the linkage analysis literature earlier [15, 29–31], although they were not in the setting of heterogeneity. Some preliminary exploration of these ideas in the context of heterogeneity can be found in Biswas et al. [18], where simultaneous mapping of two loci was considered, and in Biswas [32], where an outline for mapping several loci with the number of loci unknown has been laid out. The basic idea is that corresponding to each disease locus an α parameter is defined for each family denoting the probability that the family is linked to that locus. Thus, for L loci, say, we will have an $(L + 1)$ component mixture likelihood. A critical issue in this extension is the specification of the disease model for each of the multiple loci. In this regard, rather than assigning a specific disease model, specifying a range of models for each locus and averaging over them using BMA can potentially address this issue.

Acknowledgments

The LOAD data was obtained from NIH GWAS Data Repository dbGaP with study accession number phs000160.v1.p1. Funding support for the ‘Genetic Consortium for Late Onset Alzheimer’s Disease’ was provided through the Division of Neurosci-

ence, NIA. The Genetic Consortium for Late Onset Alzheimer's Disease includes a genome-wide association study funded as part of the Division of Neuroscience, NIA. Assistance with phenotype harmonization and genotype cleaning, as well as with general study coordination, was provided by the Genetic Consortium for Late Onset Alzheimer's Disease.

We thank Dr. Shili Lin for reading a draft of this manuscript and giving valuable comments and suggestions. We are also thankful to the two anonymous reviewers for constructive comments and suggestions, which led to improvement and clearer presentation of the manuscript.

Appendix

Details of MCMC Sampling Scheme

Let λ_1 and λ_2 be the probabilities of proposing moves of $L \rightarrow L$ and $U \rightarrow U$, respectively. So when the chain is in the L subspace, a proposal is made to either remain in L with the probability λ_1 or move to U with the probability $1 - \lambda_1$. In an $L \rightarrow L$ move, α , d , and m parameters are updated. To update α_j , we note that

$$\pi(\alpha_j | \alpha_{(-j)}, d, m, \mathbf{x}) \propto \pi_j(\alpha_j) [\alpha_j L_j(d | x_j, m) + (1 - \alpha_j) L_j(\infty | x_j, m)], 0 \leq \alpha_j \leq 1,$$

where $\alpha_{(-j)}$ is the set of all α 's except α_j . The normalizing constant can be obtained easily when $\pi_j(\alpha_j)$ is a uniform distribution. So we use Gibbs sampler to update α_j . Next, to update d , its posterior distribution is given by

$$\pi(d | \alpha, \mathbf{x}, m) \propto \pi_{d < \infty} \pi_d(I_d) L(\alpha, d, m | \mathbf{x}) = g(d), d < \infty.$$

We use a Metropolis Hastings algorithm to sample d with the proposed value of $d = d^*$ obtained from within ± 5 locations of the current location $d_{(t)}$ with equal probability. The acceptance probability is calculated as

$$\min \left\{ 1, \frac{g(d^*) f(d_{(t)} | d^*)}{g(d_{(t)}) f(d^* | d_{(t)})} \right\},$$

where $f(d^* | d_{(t)})$ is the proposal distribution of d^* given $d_{(t)}$, which is discrete uniform in this case. If accepted, $d_{(t+1)}$ is set to be d^* or else it is kept same as $d_{(t)}$. For updating m , note that

$$\pi(m | \alpha, d, \mathbf{x}) \propto \pi(I_m) L(\alpha, d, m | \mathbf{x}) = g(I_m), I_m = 1, \dots, M.$$

We can calculate $g(I_m)$ for each value of I_m and hence obtain the normalizing constant easily. So, Gibbs sampling is used to sample I_m directly from its conditional distribution.

In a $U \rightarrow L$ move, α , d , and m are generated from proposal distributions which we set to be the same as their respective prior distributions $\pi_j(\alpha_j)$, $\pi_d(I_d)$, and $\pi(I_m)$. The acceptance probability of this move is $\min\{1, A(\alpha, d, m)\}$, where

$$A(\alpha, d, m) = \frac{\pi_{d < \infty} \pi_d(I_d) \prod_{j=1}^k \pi_j(\alpha_j) \pi(I_m)}{\pi_{d = \infty}} \times \frac{L(\alpha, d, m | \mathbf{x})}{\prod_{j=1}^k L_j(\infty | x_j)} \times \frac{1 - \lambda_1}{(1 - \lambda_2) \cdot \pi_d(I_d) \cdot \prod_{j=1}^k \pi_j(\alpha_j) \pi(I_m)} \times |1|.$$

Note that $A(\alpha, d, m)$ is the product of prior ratio, likelihood ratio, proposal ratio, and Jacobian of the transformation. Since we have used the same proposal distributions as the prior distributions, and we set $\lambda_1 = \lambda_2 = 0.5$ following Biswas and Lin [1], the above expression simplifies with cancellation of terms in numerator and denominator. If this move is accepted, the generated α , d and m are taken to be the updated values in the L state, otherwise the chain stays in the U state with $d_{(t+1)} = \infty$.

An $L \rightarrow U$ move has an acceptance probability of A^{-1} and if accepted, the current values of all parameters are discarded and d is set at ∞ . In a $U \rightarrow U$ move, we set the updated value of $d = \infty$. More details about these moves can be found in Biswas and Lin [1], which also include sensitivity analyses on the choice of various parameter values and proposal distributions.

References

- Biswas S, Lin S: A Bayesian approach for incorporating variable rates of heterogeneity in linkage analysis. *J Am Stat Assoc* 2006; 101:1341–1351.
- Biswas S, Lin S, Berry DA: A new Bayesian approach incorporating covariate information for heterogeneity and its comparison with HLOD. *BMC Genet* 2005;6:S138.
- Biswas S, Lin S: Incorporating covariates in mapping heterogeneous traits – a hierarchical model using empirical Bayes estimation. *Genet Epidemiol* 2007;31:684–696.
- Knapp M, Seuchter SA, Baur MP: Linkage analysis in nuclear families. II. Relationship between affected sib-pair tests and lod score analysis. *Hum Hered* 1994;44:44–51.
- Greenberg DA, Abreu P, Hodge SE: The power to detect linkage in complex disease by means of simple LOD-score analyses. *Am J Hum Genet* 1998;63:870–879.
- Hoeting J, Madigan D, Raftery A, Volinsky C: Bayesian model averaging. *Statist Sci* 1999;14:382–401.
- Zhang L, Mukherjee B, Ghosh M, Wu R: A Bayesian framework for genetic association in case-control studies: accounting for unknown population substructure. *Statist Model* 2006;6:352–372.
- Sillanpaa MJ, Corander J: Model choice in gene mapping: what and why. *Trends Genet* 2002;18:301–307.
- Kruglyak L, Daly MJ, Reeve-Daly MP, Lander ES: Parametric and nonparametric linkage analysis: a unified multipoint approach. *Am J Hum Genet* 1996;58:1347–1363.
- Gudbjartsson JE, Jonasson HH, Frigge ML, Kong A: Allegro, a new computer program for multipoint linkage analysis. *Nat Genet* 2000;25:12–13.
- Almasy L, Amos CI, Bailey-Wilson JE, et al: Genetic analysis workshop 13: analysis of longitudinal family data for complex diseases and related risk factors. *BMC Genet* 2003; 4:S1.
- Lee JH, Cheng R, Graff-Radford N, Foroud T, Mayeux R: Analyses of the National Institute on Aging Late-Onset Alzheimer's Disease Family Study: implication of additional loci. *Arch Neurol* 2008;65:1518–1526.

- 13 Richardson S, Green PJ: On Bayesian analysis of mixtures with an unknown number of components. *J Royal Stat Soc B* 1997;59:731–792.
- 14 Raftery AE: Hypothesis testing and model selection; in Gilks WR, Richardson S, Spiegelhalter DJ (eds): *Markov Chain Monte Carlo in Practice*. London, UK, Chapman & Hall, 1996, pp 163–187.
- 15 Thomas DC, Richardson S, Gauderman J, Pitkniemi J: A Bayesian approach to multi-point mapping in nuclear families. *Genet Epidemiol* 1998;14:903–908.
- 16 Wang Z, Lin S, Popesco M, Rotter A: Modeling and analysis of multi-library, multi-group SAGE data with application to a study of mouse cerebellum. *Biometrics* 2007;63:777–786.
- 17 Daw EW, Morrison J, Zhou X, Thomas DC: Genetic analysis workshop 13: simulated longitudinal data on families for a system of oligogenic traits. *BMC Genetics* 2003;4 (suppl 1):S3.
- 18 Biswas S, Papachristou C, Irwin ME, Lin S: Linkage analysis of the simulated data evaluations and comparison of methods. *BMC Genet* 2003;4:S70.
- 19 Kong A, Cox N: Allele-sharing models: LOD scores and accurate linkage tests. *Am J Hum Genet* 1997;61:1179–1188.
- 20 Kamboh MI: Molecular genetics of late-onset Alzheimer's disease. *Ann Hum Genet* 2004;68:381–404.
- 21 Abecasis GR, Cherny SS, Cookson WO, Cardon LR: Merlin – rapid analysis of dense genetic maps using sparse gene flow trees. *Nat Genet* 2002;30:97–101.
- 22 Hayashi T, Awata T: Bayesian mapping of QTL in outbred F2 families allowing inference about whether F0 grandparents are homozygous or heterozygous at QTL. *Heredity* 2005;94:326–337.
- 23 Olson JM, Goddard KAB, Dudek DM: The amyloid precursor protein locus and very-late-onset Alzheimer disease. *Am J Hum Genet* 2001;69:895–899.
- 24 Scott WK, Hauser ER, Schmechel DE, Welsh-Bohmer KA, Small GW, Roses AD, Saunders AM, Gilbert JR, Vance JM, Haines JL, Pericak-Vance MA: Ordered-subsets linkage analysis detects novel Alzheimer disease loci on chromosomes 2q34 and 15q22. *Am J Hum Genet* 2003;73:1041–1051.
- 25 Bacanu SA, Devlin B, Chowdari KV, DeKosky ST, Nimgaonkar VL, Sweet RA: Linkage analysis of Alzheimer disease with psychosis. *Neurology* 2002;59:118–120.
- 26 Wan S, Lin S: A likelihood-based procedure for obtaining confidence intervals of disease loci with general pedigree data. *BMC Proc* 2007;1(suppl 1):S106.
- 27 Sillanpaa MJ, Bhattacharjee M: Association mapping of complex trait loci with context-dependent effects and unknown context variable. *Genetics* 2006;174:1597–1611.
- 28 Province MA, Shannon WD, Rao DC: Classification methods for confronting heterogeneity. *Adv Genet* 2001;42:273–286.
- 29 Heath SC: Markov chain Monte Carlo segregation and linkage analysis for oligogenic models. *Am J Hum Genet* 1997;61:748–760.
- 30 Lee JK, Thomas DC: Performance of Markov chain-Monte Carlo approaches for mapping genes in oligogenetic models with an unknown number of loci. *Am J Hum Genet* 2000;67:1232–1250.
- 31 Uimari P, Sillanpaa MJ: Bayesian oligogenic analysis of quantitative and qualitative traits in general pedigrees. *Genet Epidemiol* 2001;21:224–242.
- 32 Biswas S: On Incorporating Heterogeneity in Linkage Analysis. PhD dissertation, The Ohio State University, Department of Statistics 2003. Available at <http://www.ohiolink.edu/etd/view.cgi?osu1070468056>.