

# Bayesian Model Search and Multilevel Inference for SNP Association Studies

Wilson, Melanie A.<sup>1</sup>, Iversen, Edwin S.<sup>1</sup>, Clyde, Merlise A.<sup>1</sup>, Schmidler, Scott C.<sup>1</sup>, Schildkraut, Joellen M.<sup>2</sup>

<sup>1</sup>Department of Statistical Science, <sup>2</sup>Department of Community and Family Medicine, Duke University, Durham, NC

## INTRODUCTION

Technological advances in genotyping have given rise to hypothesis based association studies of increasing scope. In particular, candidate pathway studies of single nucleotide polymorphisms (SNPs) have all but replaced those focusing on a small number of variants in a handful of genes. As a results, the scientific hypotheses addressed by these studies have become more complex and more difficult to address using existing analytic methodologies. Frequently, researchers assume conditional independence across markers, analyzing each SNP separately and determining pathway- or study-wide association by a secondary analysis. This approach has proven to be powerful given its simplicity, but has several significant drawbacks when used to quantify association. These include inference in the face of multiple comparisons, complications arising from correlations among the SNPs and choice of their parameterization.

We describe Multilevel Inference for SNP Association Studies (MISA), a Bayesian model search calculation applied to penalized logistic regression that searches the space of genetic markers, and over the genetic parameterizations of each, in a computationally efficient manner. This technique allows one to estimate multilevel posterior probabilities and Bayes Factors at the global, gene and SNP level.

### Multiple Markers

*How do we correct for the high false discovery rate of current marginal methods?*

*There are many frequency based methods that adjust for multiple comparisons. However, they often make use of unjustifiable distributional assumptions on p-values under  $H_0$  and  $H_a$  and the SNPs themselves.*

Conditional independence assumption of the markers is weakened in MISA methodology and all genetic markers are simultaneously included in a Bayesian model search algorithm to identify subsets of likely associated variables.

### Multiple Genetic Parameterizations

*Which genetic parameterization for each marker do we choose when calculating marginal p-values?*

*May make more sense to allow for uncertainty in the genetic models for each SNP.*

We search over multiple genetic parameterizations for each marker in the model search algorithm (Log-Additive, Dominant, and Recessive) and report posterior summaries that are averaged across the parameterizations.

### Multilevel Inference

*How do we measure the global significance of the experiment if we assume the markers are conditionally independent?*

We weaken the conditional independence assumption and compute multilevel posterior summaries at the global, gene, and SNP levels allowing us to determine the extent to which the data supports an overall association in a pathway or gene of interest and to identify the markers most likely driving the association.

## METHODS

### Model & Prior Specification

We consider case- control association studies where:  $D_i = 1$  for a disease case and  $D_i=0$  for a control.

★Use Logistic regression to relate disease status to a subset of  $p$  SNPs,  $\mathbf{x}_\gamma$ , and  $q$  confounders,  $\mathbf{z}$  of the form:

$$\text{logit}(p(D_i=1 \mid \Theta)) = \alpha_0 + \mathbf{z}_i^T \boldsymbol{\alpha} + \mathbf{x}_{\gamma i}^T \boldsymbol{\beta}_\gamma$$

- $\alpha_0$  Intercept Term
- $\boldsymbol{\alpha}$  Vector of coefficients for the fixed confounders  $\mathbf{z}_i$  included in all models
- $\mathbf{x}_{\gamma i}^T$  Parameterizations of the SNPs included in  $M_\gamma$
- $\boldsymbol{\beta}_\gamma$  Log odds ratios for SNPs included in  $M_\gamma$

★Models,  $M_\gamma$ , are specified by the  $p$  dimensional vector  $\gamma$  where  $\gamma_i$  indicates the inclusion of SNP <sub>$i$</sub>  in model  $M_\gamma$ , and if included, the genetic parameterization of the SNP. We consider log-additive ( $\gamma_i=1$ ), dominant ( $\gamma_i=2$ ), and recessive ( $\gamma_i=3$ ) genetic parameterizations of the SNPs and set  $\gamma_i=0$  to indicate SNP <sub>$i$</sub>  is not included in  $M_\gamma$

### Posterior Summaries

#### ★Model Posterior Probabilities and Bayes Factors

The degree to which the data support any model in the Bayesian framework may be assessed via model posterior probabilities or Bayes Factors.

★The posterior model probability of any model  $M_\gamma$  in the model space  $M$  is expressed as:

$$\frac{p(D \mid M_\gamma)}{p(M_\gamma)} \propto \frac{p(D \mid M_\gamma) p(M_\gamma)}{p(M_\gamma)} \quad \text{for } M_\gamma \text{ in } M$$

Proportional to the (marginal) likelihood of  $M_\gamma$  after integrating out model specific parameterization  $\boldsymbol{\theta}_\gamma = (\alpha_\gamma, \boldsymbol{\alpha}, \boldsymbol{\beta}_\gamma)$  and approximated as  $p(D \mid M_\gamma) \approx \exp(-.5 \text{AIC}(M_\gamma))$

Prior probability of  $M_\gamma$

★The Bayes Factor comparing the evidence of any two competing models  $M_{\gamma_1}$  and  $M_{\gamma_2}$  is expressed as:

$$\frac{BF(M_{\gamma_1}: M_{\gamma_2})}{\text{Poster Odds}(M_{\gamma_1}: M_{\gamma_2})} = \frac{\text{Post. Odds}(M_{\gamma_1}: M_{\gamma_2}) + \text{Prior Odds}(M_{\gamma_1}: M_{\gamma_2})}{p(M_{\gamma_1} \mid D) / p(M_{\gamma_2} \mid D)} \quad \& \quad \frac{\text{Prior Odds}(M_{\gamma_1}: M_{\gamma_2})}{p(M_{\gamma_1}) / p(M_{\gamma_2})} = \frac{\text{Post. Odds}(M_{\gamma_1}: M_{\gamma_2})}{p(M_{\gamma_1}) / p(M_{\gamma_2})}$$

#### ★Multilevel Posterior Summaries

Results are summarized by reporting the models that are sampled from the stationary distribution (from the chain where  $t_i = 1$ ) defined as  $M^s$ .

★**Global Bayes Factor:** Gives evidence that at least one SNP is associated with disease:

$$BF(H_A:H_0) = \sum_{M_\gamma \text{ in } H_A} BF(M_\gamma:M_0) p(M_\gamma \mid H_A)$$

★**SNP Inclusion Probability and BF:** Marginal measures of significance for each genetic marker:

$$\Pr(\gamma_i > 0 \mid D) = \sum_{M_\gamma \text{ in } M^s} 1(\gamma_i \neq 0) \Pr(M_\gamma \mid D, M^s) \quad \& \quad BF(\gamma_i > 0) = \text{Post. Odds}(\gamma_i > 0) + \text{Prior Odds}(\gamma_i > 0)$$

★**Gene Inclusion Probability and BF:** Measures of significance of one or more of the SNPs within a given gene being associated.

$$\Pr(\gamma_{g1} = 1 \mid D) = \sum_{M_\gamma \text{ in } M^s} 1(\gamma_{g1} = 1) \Pr(M_\gamma \mid D, M^s) \quad \& \quad BF(\gamma_{g1} = 1) = \text{Post. Odds}(\gamma_{g1} = 1) + \text{Prior Odds}(\gamma_{g1} = 1)$$

### Model Search

Our model search algorithm is based on the Evolutionary Monte Carlo algorithm of Liang and Wong (2000) where individuals in a population of model specifications compete and mate in order to produce increasingly stronger individuals via a mixture of parallel tempering (Geyer,1991) and population updates of a genetic algorithm (Holland,1975).

**Population:** ( $N \times p$ ) matrix that indicates the set of  $N$  different model specifications of the current states of all of the parallel chains.

**Individual:** Model specification vector of length  $p$  that indicates which SNPs and the genetic effect of those SNPs in the model at a current state in one of the parallel chains.

**Strength:** Determined for an individual in the population based on their value of the fitness function,  $p(M_\gamma)$ :

$$p(M_\gamma) = .5 \text{AIC}(M_\gamma) - \log(p(M_\gamma))$$

$$\text{AIC}(M_\gamma) = 2i - 2\text{DEV}(M_\gamma)$$

and  $i$  is the number of predictors in model  $M_\gamma$ .

#### ★Parallel Tempering

Interested in sampling from the model space,  $M$ , by running parallel chains where each chain is associated with a decreasing temperature  $t_i$ .

★Chains with high temp. values have flat proposal densities that allow us to make large global moves to models across regions of low probability.

★Chains with low temp. values have peaked proposal densities that allow us to make small local moves to specific models.

#### ★Genetic Algorithm Population Updates

**Mutation (Metropolis Update):** Model,  $M_\gamma$ , is chosen from the current population of models and values in the model indicator  $\gamma_i$  are mutated.

**Crossover (Partial State Swap):** One model pair, ( $M_i$ ,  $M_j$ ), is selected from the current model population and two new offspring are produced and replace the old pair in the population.

**Exchange (Full State Swap):** Given the current population and attached temperature ladder, we propose a new population by making an exchange between models  $M_i$  and  $M_j$  without changing the values of the temperature ladder.

## RESULTS

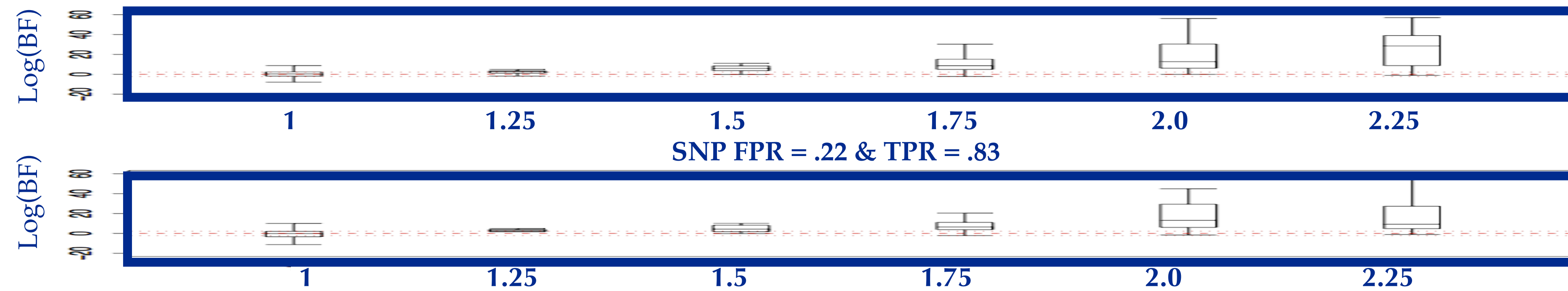
### Simulation

We explore optimal choice of the hyper-parameters,  $a$ ,  $b$ , based on a set of 30 simulated case-control data sets that are modeled after an ovarian cancer candidate pathway study of interest (NCOCS). Data sets are comprised of a binary outcome representing disease status, and genetic data on 399 ovarian cancer cases and 798 controls.

★Genetic data was simulated at the same 508 tag SNPs as genotyped in the study. The data was simulated in two stages. First, for the 53 genes represented in the data set, we phased the NCOCS control SNP genotype data and estimated recombination rates using PHASE (Stephens et al., 2001). Second, given a model of association, we generated case-control data at these tags using HAPGEN (Marchini and Su, 2006).

★Ten of the simulations are null, and in the remaining twenty we assumed that 9 random genes were associated and that within these genes, a single, randomly chosen tag SNP was the source of the association. Within ten of these associated simulations, three of the associated SNPs were accorded an odds ratio (OR) of 1.25, three an OR of 1.5, and three an OR of 1.75. The other ten associated simulations had the same structure only with OR's of 1.75, 2.0, and 2.25. In each of the twenty associated simulations, one SNP with each OR was assumed to have a dominant genetic parameterization, one of each a log-additive, and one of each a recessive.

**Simulation Results: Global FPR = .1 & TPR = .1**  
Gene FPR = .26 & TPR = .92

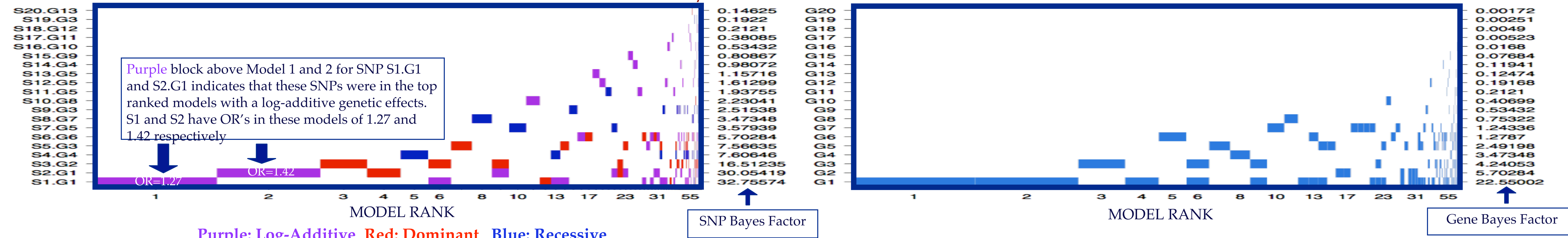


### NCOCS

The North Carolina Ovarian Cancer Study (NCOCS) is a population based case/control study. The data is comprised of cases and controls who were matched based on answers to a in-person interview. The analysis was run by looking only at the Caucasian cases/controls and a subtype of the cases. This gave us a total of 399 cases and 798 controls. For the purpose of this example we only look at the 66 SNPs in our data set that are in a candidate pathway that have passed a marginal screen. We also controlled for age of patient and previous diagnosis of breast cancer by forcing these variables into our models while performing the model selection.

To account for missing data, we use fastPHASE (Stephens et al., 2001) to sample haplotypes and missing genotypes given the observed unphased genotypes and create 100 imputed data sets. Results were found by using the average value of our fitness function,  $p(M_\gamma)$ , across the imputed data sets.

**NCOCS Pathway MISA Results: Global BF = 2.3721**



★Plots summarize the marginal associations of the 20 most highly associated of the 66 SNPs and genes they represent. SNPs and genes in the pathway are represented by a two level name where the number represents the rank of the SNP or gene by means of the marginal Bayes Factor. Summary of model inclusion for the top 100 sampled models on the basis of their posterior model probability are plotted. The models are ordered on the x-axis in descending probability and the width of the column associated with a model is proportional to that probability. A SNP's or gene's presence in a model is indicated by a block at the intersection of the model's column and a SNP's or gene's row. The color of the block in the SNP plot indicates the genetic parameterization of the SNP in that model.

★SNPs 1,2, and 3 have Bayes Factors greater than 10, providing strong evidence that marginally these SNPs are associated with ovarian cancer. Both S1 and S2 are in G1 which has a gene Bayes Factor of 22.5502, giving strong evidence that at least one of the SNPs within gene G1 is associated with ovarian cancer.

## CONCLUSION

In this paper, we describe an analytic strategy of pathway association studies that allows one to quantify evidence of associations at 3 levels: global, gene, and SNP while allowing for uncertainty in the genetic parameterization of the markers. Our methodology for SNP association studies accounts for several problems encountered by commonly used marginal analyses. These include the hard to justify conditional independence assumption made on p-values under  $H_0$ ,  $H_A$  and/or on the markers themselves. While our algorithm is not assumption free, it does relax the typical conditional independence assumption, does not require an *a priori* choice of SNP-specific genetic parameterization, facilitates a multilevel assessment of statistical significance, and allows optimal prior parameters to be chosen via simulation. Given the multilevel assessment for a pathway or gene of interest, we are able to ascertain if any genes warrant further investigation and which regions of these genes to tag more densely.

Our analytic strategy can readily accommodate different modeling strategies. The simulation study plays a critical role in this strategy by aiding in the investigation of the optimal hyper-parameters and operating characteristics of these parameters within our method. Similar simulations can be produced based on the study of interest and optimal hyper-parameters can be chosen to achieve a desired balance between the false and true positive rates. Our strategy can also accommodate models associated with studies other than case-control associations studies in that other penalized likelihood functions may be used instead of AIC. Likelihoods such as those arising from survival models may be used in place of the logistic regression model. Also, missing data on SNPs is a common phenomenon; we address this problem by generating multiple imputations for the missing data. This allows us to use all cases while reporting valid uncertainty summaries. Other methods may be used to account for missing data in the study of interest. Finally, when computing the posterior probabilities of each of the sampled models, we assume that the prior of each model comes from the beta-binomial distribution where the SNPs are assumed to be independent of each other with a common inclusion probability. It would be interesting to explore other prior distributions that build in more genetic information either on pathway structure or known genetic function of the SNPs.

### Literature Cited

- [1] C. Geyer. *Markov chain monte carlo maximum likelihood*. Computing Science and Statistics, page 156-163, 2001.
- [2] J.H. Holland. *Adaptation in natural and artificial systems*. University of Michigan Press, 1975.
- [3] F. Liang & W.H. Wong. *Evolutionary monte carlo: Applications to Cp model sampling and change point problem*. Statistica Sinica, 10:317, 2000.
- [4] M. Stephens, N. Smith, P. Donnelly. *A New Statistical Method for Haplotype Reconstruction for Population Data*. The American Journal of Human Genetics, 68:978-989, 2001.
- [5] J. Marchini & Z. Su. *A C++ Program for Simulating Case and Control SNP Haplotypes*. 2006
- [6] R. Kass & A. Raftery. *Bayes Factors*. J. Amer. Statist. Assoc. 90:733-795, 1995.

### Acknowledgements

This work was supported in part by:

- Duke SPORC in Breast Cancer: 5P50CA068438-07
- North Carolina Ovarian Cancer Study: 1R01CA076016
- Bayesian Modeling and Optimal Design for Studies of Gene-Environment Association: 1R01HL090559-01
- NSF SCREMS award DMS-0422400.

### Web Resources



★Email:

maw27@stat.duke.edu

★R Package & Pdf copy of poster:

www.stat.duke.edu/gbye/MISA.html