

Evaluating Haplotype Effects in Case-Control Studies via Penalized-Likelihood Approaches: Prospective or Retrospective Analysis?

Megan L. Koehler,¹ Howard D. Bondell,¹ and Jung-Ying Tzeng^{1,2*}

¹Department of Statistics, North Carolina State University, Raleigh, North Carolina

²Bioinformatics Research Center, North Carolina State University, Raleigh, North Carolina

Penalized likelihood methods have become increasingly popular in recent years for evaluating haplotype-phenotype association in case-control studies. Although a retrospective likelihood is dictated by the sampling scheme, these penalized methods are typically built on prospective likelihoods due to their modeling simplicity and computational feasibility. It has been well documented that for unpenalized methods, prospective analyses of case-control data can be valid but less efficient than their retrospective counterparts when testing for association, and result in substantial bias when estimating the haplotype effects. For penalized methods, which combine effect estimation and testing in one step, the impact of using a prospective likelihood is not clear. In this work, we examine the consequences of ignoring the sampling scheme for haplotype-based penalized likelihood methods. Our results suggest that the impact of prospective analyses depends on (1) the underlying genetic mode and (2) the genetic model adopted in the analysis. When the correct genetic model is used, the difference between the two analyses is negligible for additive and slight for dominant haplotype effects. For recessive haplotype effects, the more appropriate retrospective likelihood clearly outperforms the prospective likelihood. If an additive model is incorrectly used, as the true underlying genetic mode is unknown a priori, both retrospective and prospective penalized methods suffer from a sizeable power loss and increase in bias. The impact of using the incorrect genetic model is much bigger on retrospective analyses than prospective analyses, and results in comparable performances for both methods. An application of these methods to the Genetic Analysis Workshop 15 rheumatoid arthritis data is provided. *Genet. Epidemiol.* 34:892–911, 2010. © 2010 Wiley-Liss, Inc.

Key words: haplotype-based association analysis; variable selection; regularized regression; prospective likelihood; retrospective likelihood

Contract grant sponsor: NIH; Contract grant numbers: 5R01-HL049609-14; 1R01-AG021917-01A1; T32GM081057; R01 MH084022-01; R01 MH084022-01; 1P01-CA142538-0; R01 MH084022-01; 1P01-CA142538-01; Contract grant sponsors: University of Minnesota; Minnesota Supercomputing Institute; GAW; Contract grant numbers: R01-GM031575; AR44422; Contract grant sponsor: NSF; Contract grant number: DMS-0705968.

*Correspondence to: Jung-Ying Tzeng, Department of Statistics, Campus Box 7566, North Carolina State University, Raleigh, NC 27695. E-mail: jytzeng@stat.ncsu.edu

Received 29 April 2010; Revised 3 September 2010; Accepted 15 September 2010

Published online 18 November 2010 in Wiley Online Library (www.wileyonlinelibrary.com).

DOI: 10.1002/gepi.20545

INTRODUCTION

Haplotype-based association analysis evaluates the joint effects of closely linked genetic markers on a trait of interest. When compared to its single-marker counterparts, this multi-marker approach can be more powerful to detect associations when the causal variants are not genotyped [de Bakker et al., 2005; Zaitlen et al., 2007], have low frequency [de Bakker et al., 2005; Schaid, 2004], or exhibit cis-acting effects [Clark, 2004; Schaid, 2004]. A standard approach for performing haplotype-based analysis is to regress the trait value on the haplotypes and to test the significance of the regression parameters [Balding, 2006]. In recent years, applying penalized likelihood methods to identify important haplotypic factors has become increasingly popular in the literature. For example, Li et al. [2007] use the least absolute shrinkage

and selection operator (LASSO) [Tibshirani, 1996] to perform selection among numerous possible haplotypes resulting from different haplotype window lengths. Guo and Lin [2009] use LASSO regression to evaluate the effects of rare haplotypes and high-dimension haplotype-environment interactions. Tzeng et al. [2010] use adaptive LASSO regression [Zou, 2006] to study high-dimensional gene-treatment interactions in a haplotype-based pharmacogenetic analysis. These methods introduce a penalty on the regression coefficients and shrink the coefficient estimates of non-important covariates toward zero. The motivation behind using penalized methods in haplotype-based analysis is that while the model space under consideration may be large (e.g. 6–16 haplotypic predictors with a sample size of 500–1,000 [Chen and Kao 2006; Epstein and Satten, 2003; French et al., 2006; Stram et al., 2003], which can yield 2^6 to 2^{16} possible models), many of the haplotypic predictors are not likely to be associated

with the phenotype. In this case, it is more efficient to shrink these effect estimates to zero than to estimate them purely. This shrinkage leads to a reduction in variance and can increase the power to detect important haplotypic predictors [Guo and Lin, 2009].

Modifications of classic penalized methods have also been developed to perform haplotype-based analysis and attempt to address issues specific to this type of analysis. Tanck et al. [Souverain et al., 2006, 2008; Tanck et al., 2003] use a modified version of Ridge regression to stabilize inference for rare haplotypes. By constructing an L_2 -norm penalty term on the differences in coefficients of similar haplotypes, the coefficients of rare haplotypes are smoothed toward that of a similar common haplotype. Chen et al. [2009] develop an adaptive penalized likelihood framework to address the precision-efficiency tradeoff encountered in retrospective methods. Motivated by the fact that typical retrospective case-control estimates of haplotype effects are efficient but sensitive to violations of underlying assumptions (e.g. Hardy-Weinberg equilibrium and gene-environment independence), they construct a penalized estimator based on either an L_1 -norm or an L_2 -norm penalty that combines the merits of assumption-free estimators (i.e. robust) and assumption-dependent estimators (i.e. efficient). Tzeng and Bondell [2010] modify traditional adaptive LASSO regression by placing an L_1 -norm penalty on pair-wise differences of the regression coefficients. This allows for effect comparisons between all pairs of distinct haplotypes, rather than with respect to an arbitrary baseline haplotype, during the estimation process. As a result, the approach is able to sort haplotypes into different groups according to their effect sizes and eliminates the need for a post-hoc pair-wise analysis of haplotype effects. In general, the key of a penalized regression method lies in the form of the penalty—by carefully designing the form of the penalty, one can gear the penalized likelihood approach toward accomplishing various desired tasks.

Penalized regression methods rely on the underlying data likelihood. When analyzing data from case-control studies, one can implement methods based on a prospective likelihood (modeling the probability of disease status conditional on exposure) or a retrospective likelihood (modeling the probability of exposure conditional on disease status). Under a case-control design, a retrospective likelihood should be used because data are collected based on disease status. However, in practice, it is common for researchers to use a prospective likelihood, as it does not require specifying a model for the joint distribution of the genetic and environmental effects. Bypassing this step makes implementing prospective methods much easier than retrospective methods [Lin et al., 2005]. This approach seems congruent with the well-known result that optimizing the prospective likelihood yields the same inference on the disease model parameters as optimizing the retrospective likelihood [Prentice and Pyke, 1979]. This result requires that the distribution of the covariates be free of restrictions, which does not generally hold in haplotype-based analysis. Haplotypes are not directly observed from unphased genotype data. In order to reconstruct the haplotypes and estimate their effects, some assumptions must be placed on their frequency distribution (typically Hardy-Weinberg equilibrium).

Most of the penalized regression approaches mentioned above utilized a prospective likelihood. It has been well

documented that when using non-penalized regression methods in haplotype-based analysis of case-control data, ignoring the ascertainment scheme can be detrimental. A prospective analysis can lead to a loss of efficiency and severe bias when assessing the haplotype effects [Cordell, 2006; Satten and Epstein, 2004; Stram et al., 2003]. The aim of this work is to determine whether similar consequences occur when using penalized regression for case-control studies. Specifically, we consider the adaptive LASSO penalty, and use simulation studies to examine the relative performance in parameter estimation and model selection between the penalized method using a prospective likelihood and using a retrospective likelihood. Our results suggest that the impact of using a prospective likelihood in place of a retrospective likelihood depends on (1) the underlying genetic mode of the causal variants, and (2) the genetic model used in the analysis. If the correct genetic model is used, then the difference between the two analyses is negligible for additive and slight for dominant haplotype effects. For recessive haplotype effects, the more appropriate retrospective likelihood clearly outperforms the prospective likelihood. If an additive model is used regardless of the underlying genetic mode, then both retrospective and prospective penalized methods suffer from a sizeable power loss and increase in bias. The impact of using the incorrect genetic model is much bigger on retrospective analyses than prospective analyses, and results in comparable performances for both methods. In addition to extensive simulation studies, we present an application to the Genetic Analysis Workshop 15 rheumatoid arthritis data.

METHODS

PROSPECTIVE AND RETROSPECTIVE LIKELIHOODS

Let the vector (Y_i, G_i, E_i) represent the observed data for individual i in a case-control sample of size n . Let Y_i be a binary indicator of disease status, where $Y_i = 1$, if individual i is a case and 0 otherwise. Let G_i denote the unphased genotype of individual i at m biallelic SNPs and E_i denote any environmental covariates measured on individual i . Let H_i represent the vector of haplotype counts for individual i . Although researchers want to investigate the relationship between Y_i and H_i , they only have access to G_i ; therefore, the individual's haplotype set must be inferred from their unphased genotypes.

The relationship between the disease phenotype and the covariates can be characterized by the conditional density function $P(Y | H, E)$. A standard approach for binary trait values is logistic regression, which models the conditional probability as

$$P(Y = y | H, E) = \frac{\exp\{y \cdot (\beta_0 + \mathcal{Z}(H, E)^T \beta)\}}{1 + \exp\{\beta_0 + \mathcal{Z}(H, E)^T \beta\}},$$

where β_0 is an intercept, β is the vector of disease model parameters representing the log-odds ratios, and $\mathcal{Z}(H, E)$ is a specified vector-valued function of the vector of haplotype counts H and the vector of environmental covariates E . For example, one can use an identity function for $\mathcal{Z}(H, E)$ so that $\mathcal{Z}(H, E)^T = [H^{*T}, E^T]$, where H^* is the vector H with baseline haplotype element removed. Other examples of $\mathcal{Z}(H, E)$ are described in the data generation

section and the choice depends on the genetic model of the haplotype effects adopted in the analysis. The dimension of β is determined by the dimension of $[H^*, E]$; it is the sum of the number of haplotypes (excluding the baseline) and the number of environmental covariates included in the model, along with any interaction terms that may be used. The dimension of β is denoted by p . Throughout this work, we assume that the sample size is greater than the dimension of β , i.e. $n > p$. This is consistent with case-control data sets, which typically have sample sizes in the order of 10^2 to 10^3 and the number of potential predictors in the order of 10 [Chen and Kao 2006; Epstein and Satten, 2003; French et al., 2006; Stram et al., 2003].

Various likelihood models have been developed to conduct inference about the disease model parameters in haplotype-based analyses while properly accounting for phase uncertainty. The inference can be based either on a prospective likelihood or on a retrospective likelihood. In this work, we consider maximum likelihood methods developed by two groups—one focusing on a prospective approach and the other on a retrospective approach. We implement the prospective method developed in Lake et al. [2003]. Their prospective likelihood models $P(Y_i | G_i, E_i)$ and is expressed as

$$\begin{aligned} L_p &= \prod_{i=1}^n P(Y_i | G_i, E_i) = \prod_{i=1}^n \sum_{H_i \in S(G_i)} P(H_i, Y_i | G_i, E_i) \\ &= \prod_{i=1}^n \sum_{H_i \in S(G_i)} P(Y_i | H_i, E_i) P(H_i), \end{aligned}$$

where $S(G)$ is the set of all haplotype pairs consistent with G , $P(H = h) = 2 \prod_{k=1}^l \pi_k^{n_k} / n_k$ under the assumption of Hardy-Weinberg equilibrium, n_k is the number of copies of the k th haplotype in h , π_k is the population frequency of the k th haplotype, and l is the number of haplotypes included in the disease model. We implement the retrospective method developed in Lin and Zeng [2006]. Their retrospective likelihood models $P(G_i, E_i | Y_i)$ and is expressed as

$$\begin{aligned} L_r &= \prod_{i=1}^n P(G_i, E_i | Y_i) = \prod_{i=1}^n \sum_{H_i \in S(G_i)} P(H_i, G_i, E_i | Y_i) \\ &\propto \prod_{i=1}^n \sum_{H_i \in S(G_i)} P(Y_i | H_i, E_i) P(H_i) P(E_i | G_i). \end{aligned}$$

The only difference between the two likelihoods is the conditional density function $P(E_i | G_i)$ found in the retrospective likelihood. The parameters in this model are of no interest to researchers performing haplotype-based association analysis, but they must be estimated in order to make proper inference when using a retrospective design. Specifying a model for this conditional density function and the subsequent maximum likelihood estimation are computationally intensive. As a result, researchers often rely on prospective methods when analyzing case-control data even though retrospective methods are dictated by the ascertainment scheme [Lin and Zeng, 2006].

HAPLOTYPE ANALYSIS VIA PENALIZED LIKELIHOOD METHODS

While many different penalized likelihood methods can be used in haplotype-based association analysis, we

consider the adaptive LASSO (ALASSO) penalty in this work. This approach achieves simultaneous variable selection and parameter estimation and is an oracle procedure. This refers to the fact that the approach asymptotically selects the correct model, and the resulting estimator is root- n consistent and asymptotically normal with the same variance as if the true model were known beforehand [Zou, 2006].

The ALASSO effect estimates are obtained by minimizing a penalized negative log-likelihood. These estimates are expressed as

$$\hat{\beta}_\lambda = \arg \min_{\beta} -\ell_n(\beta, \phi) + \lambda \sum_{j=1}^p w_j |\beta_j|,$$

where $\ell_n(\beta, \phi)$ denotes the log-likelihood, ϕ is a (possible) set of nuisance parameters (e.g. the haplotype frequencies, π_k), λ is the non-negative regularization parameter that controls the amount of shrinkage, and w_j are data-dependent weights. By placing an L_1 -norm penalty on the regression coefficients, the ALASSO can set their estimates to exactly zero if the value of λ is large enough. It is this feature that allows the procedure to perform simultaneous variable selection and parameter estimation. Unlike its predecessor the LASSO, the ALASSO places a different penalty on each coefficient through the use of adaptive weights that are inversely proportional to their relative importance. Consequently, haplotypes with negligible effects receive larger penalties and are more readily shrunk to zero. This allows the effects of associated haplotypes to be estimated more efficiently. Zou [2006] proposed to set the weights as $w_j = |\hat{\beta}_j|^{-\gamma}$, where $\hat{\beta}_j$ is an initial root- n consistent estimator of β_j and $\gamma > 0$ is an additional tuning parameter. In our analysis, we chose $\gamma = 1$ and let $\hat{\beta}_j$ be the maximum likelihood estimate of the haplotype effect computed by `haplo.glm` in R and `HAPSTAT` in Linux for the prospective and retrospective likelihoods, respectively [Lake et al., 2003; Lin et al., 2005].

When performing penalized likelihood methods, it is typical to center and scale the design matrix. Scaling assures that each column of the design matrix has the same variance and the resulting estimator is scale-equivariant (i.e. multiplication of any predictor by any constant will simply divide the resulting slope estimate by the identical constant; hence the linear predictor remains unchanged). This is desirable so that if, for example, the units of a predictor are changed, such as feet to inches, the resulting predicted values will remain unchanged. Often the predictors are also centered, so that in the normal linear regression setting, the intercept can be omitted and the slope parameter estimates are orthogonal to the intercept estimate. However, in the generalized linear models as considered here, this is not the case; hence, the design matrix is typically not centered. Furthermore, in the ALASSO analysis, we also do not scale the imputed haplotype design matrix because the adaptive weights we set (i.e. $|\hat{\beta}_j|^{-1}$) are scale-equivariant. The use of scale-equivariant weights automatically forces the resulting estimator to be scale-equivariant.

The ALASSO solution ($\hat{\beta}_\lambda$) also depends on the value of λ . The regularization parameter controls the tradeoff between model fit and model sparsity. By including more predictors, one can continually improve the fit on the training data at the expense of interpretability and over fitting. Many model selection criteria, such as Mallows'

C_p , Akaike information criterion (AIC), Bayesian information criterion (BIC), and cross validation [Arlot and Celisse, 2010; Hastie et al., 2009; Shao, 1997], can be used to determine the appropriate value of λ from an exhaustive grid search. Because the goal of haplotype-based association analysis is more aligned with selecting the true model than minimizing prediction error, we use the BIC for tuning, which can achieve consistent model selection [Yang, 2005]. BIC is defined as

$$\text{BIC} = -2\ell_n(\hat{\beta}_\lambda, \hat{\phi}) + df_\lambda \cdot \log(n),$$

where $\ell_n(\hat{\beta}_\lambda, \hat{\phi})$ is the log-likelihood evaluated at the estimated regression coefficients and maximized over ϕ for a given λ , and df_λ is the degrees of freedom, which equals the number of non-zero elements in $(\hat{\beta}_\lambda, \hat{\phi})$. The λ that minimizes the BIC is chosen as the regularization parameter, and its corresponding $\hat{\beta}_\lambda$ is the ALASSO estimate. For comparison, we also present some of the results using AIC as a tuning method. In the definition of AIC, the penalty on the degrees of freedom is changed from $\log(n)$ to 2. As a result, models selected using AIC incur less shrinkage, and the chosen ALASSO estimate will be closer to the unpenalized MLE estimate than those found using BIC.

For computational convenience, the objective function created via the least squares approximation (LSA) method was used to calculate the ALASSO solution. The LSA method replaces the objective function of the original ALASSO problem with a least squares objective function [Wang and Leng, 2007]. The method is motivated by a standard Taylor series expansion of $-\ell_n(\beta, \phi)$ about $(\tilde{\beta}, \tilde{\phi})$, the function's unpenalized minimizer, and shows that the ALASSO estimate has the exact same asymptotic distribution as the estimator given by

$$\hat{\beta}_\lambda = \arg \min_{\beta} (\beta - \tilde{\beta})^T \tilde{\Sigma}^{-1} (\beta - \tilde{\beta}) + \lambda \sum_{j=1}^p w_j |\beta_j|,$$

where $\tilde{\Sigma}$ is the estimated covariance matrix of $\tilde{\beta}$. Note that the minimizer of the unpenalized least squares objective function is exactly the maximum likelihood estimator. Hence, as with the penalized likelihood, varying the tuning parameter yields a continuous solution path from the MLE to the solution with all coefficients equal to zero. Because the underlying data likelihoods are not quadratic in the regression coefficients, using the alternative least squares objective function greatly reduces the computational costs for finding the ALASSO solution [Wang and Leng, 2007]. Using the LSA method eliminates the need for an iterative procedure to perform optimization; it only requires one unpenalized fit of the original objective function and then a grid search to determine λ . The final estimate is again chosen by minimizing the BIC (and AIC for some results).

SIMULATION STUDIES

We performed simulation studies to examine the performance of the ALASSO method under two competing data likelihoods when analyzing case-control data. Specifically, we wanted to determine if using a prospective likelihood in place of the more appropriate retrospective likelihood was detrimental when performing haplotype-based analyses using a penalized likelihood method.

To answer this question, we compared the parameter estimation and model selection properties of each approach. For ease of discussion, let *aPro* refer to ALASSO coupled with a prospective likelihood and *aRetro* refer to ALASSO coupled with a retrospective likelihood.

SIMULATION SETTINGS

Our simulation studies were based on two haplotype distributions (given in Table I) studied by Lin and Huang [2008]. These distributions are based on the common haplotypes formed by five SNPs on chromosome 18 in the CEU sample of the HapMap data. The SNPs used to build the first haplotype distribution were in strong linkage disequilibrium, while those used to build the second haplotype distribution were not. Distribution 1 represents a haplotype distribution with a few high frequency haplotypes, while the haplotype frequencies in Distribution 2 are more uniform. Each distribution was normalized so that the haplotype frequencies summed to 1. Because eight haplotypes define Distribution 1 and 11 haplotypes define Distribution 2, the specific dimension of β is $p = 7$ and $p = 10$, respectively.

For each haplotype distribution, we considered two simulation studies—one in which a single haplotype was associated with the disease (Simulation I) and one in which two haplotypes were associated with the disease (Simulation II). Because our focus was on identifying and estimating disease-haplotype associations, only genetic covariates were considered in our simulation studies (i.e. E is taken to be \emptyset). We took the sample size to be $n = 1,000$ with an equal number of cases and controls. In both simulation studies, we examined the effect of varying the genetic mode of the associated haplotype(s) on the performance of *aPro* and *aRetro*. We allowed the associated haplotype(s) to act additively, dominantly, or recessively with respect to disease risk. In practice, the genetic mode of a risk haplotype is unknown a priori, and researchers typically analyze the data additively regardless of the true genetic model. To mimic this scenario, we analyzed each data set using the correct genetic model and again using an additive model. We use the term “genetic mode” to refer to the true underlying architecture of the relationship between the disease and the risk haplotypes, and the term “genetic model” to refer to the assumed architecture used in the analysis.

TABLE I. Haplotype distributions used in simulations

Hap ID	Distribution 1		Distribution 2	
	Haplotype	Frequency	Haplotype	Frequency
1	00000	0.406	00010	0.131
2	00001	0.213	00001	0.105
3	01111	0.141	10010	0.103
4	10000	0.132	10101	0.100
5	10001	0.055	00100	0.088
6	01000	0.021	10100	0.088
7	01100	0.018	00101	0.086
8	01001	0.014	01101	0.084
9			10001	0.081
10			10000	0.079
11			00000	0.055

TABLE II. Settings for SIM I and SIM II (odds ratios)

Hap ID	Distribution 1					Distribution 2					
	Sim I	Sim II				Sim I	Sim II				
Freq	R	C	R/R	R/C	C/C	Freq	R	C	R/R	R/C	C/C
1	0.406					0.131					
2	0.213	1	1	1	1	0.105	1	0	1	0	0
3	0.141	1	1	1	0	0.103	1	1	1	1	0
4	0.132	1	0	1	0	0.100	1	1	1	1	1
5	0.055	0	1	0	0	1	0.088	1	1	1	1
6	0.021	1	1	0	1	1	0.088	1	1	1	1
7	0.018	1	1	1	1	1	0.086	1	1	1	1
8	0.014	1	1	1	1	1	0.084	1	1	1	1
9						0.081	1	1	1	1	1
10						0.079	1	1	0	1	1
11						0.055	0	1	0	0	1

In addition to genetic mode, we varied the frequency and effect size of the associated haplotype. In Simulation I (SIM I), a rare or a common haplotype was chosen to be the associated genetic variant for each haplotype distribution. A haplotype with a frequency less than 0.10 was considered rare; otherwise, it was considered common. We set the effect sizes of the associated haplotype (in terms of the odds ratio θ) so that the power of finding the effect fell in a reasonable range. Under additive and dominant modes, we set $\theta = \{1.0, 1.3, 1.5, 1.7, 2.0\}$, and under a recessive mode, we set $\theta = \{1.0, 2.0, 2.5, 3.0, 3.5\}$. We let $\theta = 1$ to examine the performance of the approaches under a null model. In Simulation II (SIM II), we allowed two haplotypes to be associated with the disease, where the associated haplotypes were both rare, one rare and one common, or both common. The odds ratios of both associated haplotypes were set to $\theta = 1.7$ for additive and dominant modes and $\theta = 3.0$ for a recessive mode. The settings for Simulations I and II can be found in Table II. In all, 78 different simulation settings were studied.

DATA GENERATION

We generated the haplotype pair of an individual conditional on their disease status and then dissolved the haplotype pair into its unphased genotypes. Let $P(H = h | Y = y)$ denote the probability of having a particular haplotype pair conditional on disease status. This probability can be expressed as

$$P(H = h | Y = y) = \frac{P(Y = y | H = h) \cdot P(H = h)}{\sum_h P(Y = y | H = h) \cdot P(H = h)}$$

For a case individual, $P(Y = 1 | H = h)$ was found using the logistic regression model

$$P(Y = 1 | H) = \frac{\exp\{\beta_0 + \mathcal{Z}(H)^T \beta\}}{1 + \exp\{\beta_0 + \mathcal{Z}(H)^T \beta\}}$$

For a control individual, $P(Y = 0 | H = h) = 1 - P(Y = 1 | H = h)$. The function $\mathcal{Z}(\cdot)$ depends on the genetic mode of the haplotype(s) associated with the disease. If the haplotype acts additively with respect to disease risk, then $\mathcal{Z}(H) = H^*$ where H^* is the haplotype-count vector H with the baseline haplotype element removed. If the haplotype acts

dominantly, then $\mathcal{Z}(H) = I\{H^* \geq 1\}$, where the inequality is taken component wise, and $I\{A\} = 1$ if A is true. If the haplotype acts recessively, then $\mathcal{Z}(H) = I\{H^* = 2\}$. The vector β was taken to be the log of the vectors given in Table II for each simulation setting. The value of β_0 was set to maintain a disease prevalence between 3% and 5%. Once $P(Y = y | H = h)$ was calculated for each haplotype pair formed from the haplotype distributions given in Table I, the vectors $\hat{P}_{H|Y=y} = (P(H = h_1 | Y = y) \cdots P(H = h_q | Y = y))$ were calculated for $Y = 0$ and $Y = 1$ where q is total number of haplotype pairs. The sample was generated by taking $n/2$ draws from the multinomial distribution parameterized by $P_{H|Y=0}$ to determine the haplotype pairs of the controls, and by taking $n/2$ draws from the multinomial distribution parameterized by $P_{H|Y=1}$ to determine the haplotype pairs of the cases. The haplotype pair of each individual was then dissolved into its unphased genotype.

COMPUTATIONAL DETAILS

For each simulation setting, 1,000 replicate data sets were generated, except for simulation under the null (i.e. $\theta = 1$) where 2,000 replicated data sets were generated. We doubled the number of simulated data sets in the null simulation to obtain more stable estimates of the false-positive counts (FP counts). For each data set, analysis began by calculating the unpenalized MLEs of the haplotype log-odds ratios. Prospective MLEs were obtained using haplo.glm in R [Lake et al., 2003] and retrospective MLEs were obtained using HAPSTAT in Linux [Lin et al., 2005]. The estimated covariance matrix of the MLEs was also obtained from each program. The final aPro and aRetro estimates were calculated by using the MLEs and their covariance matrix to compute the ALASSO solution via LSA. Based on these final estimates, estimation and model selection measures were calculated to compare the performance of the aPro and aRetro approaches. The estimation measures provided in this analysis are the bias and mean square error (MSE) of haplotype effect estimates. The model selection measures provided are the FP count (i.e. the number of non-risk haplotypes retained in the model) and the false-negative count (FN count, i.e. the number of risk haplotypes not retained in the model), for which a smaller FN count indicates better power to identify risk haplotypes. For each measure, we report the mean across the replicated data sets and the corresponding standard error of the mean. The functions used to obtain the results of this simulation study are available at the corresponding author’s website <http://www4.stat.ncsu.edu/~tzeng/publications.php>.

SIMULATION RESULTS

We present the results of Null Simulation (i.e. no-risk haplotypes) for both haplotype distributions in Table III. For the simulations involving risk haplotypes, because the pattern of results was similar across both haplotype distributions, for brevity we focus the discussion on the results for the first haplotype distribution. The results of Simulation I (single-risk haplotype) for this setting are found in Tables IV and V. The results of Simulation II (two-risk haplotypes) for this setting are found in Tables VI and VII. The discussion generalizes to the second haplotype distribution, and specific results are shown in

TABLE III. BIC-penalized results of null simulation (no-risk haplotypes)

Model ^b	Model selection results		Parameter estimation results			
	False positives ^a		Bias		MSE	
	Pro	Retro	Pro	Retro	Pro	Retro
Haplotype Distribution 1 (Correct analysis)						
Additive	0.025 ^c (0.005)	0.028 (0.006)	0.000 (0.002)	-0.001 (0.002)	0.002 (0.001)	0.002 (0.001)
Dominant	0.025 (0.006)	0.033 (0.006)	0.000 (0.002)	0.000 (0.002)	0.002 (0.001)	0.003 (0.001)
Recessive	0.026 (0.008)	0.029 (0.011)	-0.001 (0.001)	0.003 (0.002)	0.001 (0.001)	0.003 (0.002)
Haplotype Distribution 2 (Correct analysis)						
Additive	0.062 (0.020)	0.055 (0.017)	0.000 (0.002)	0.000 (0.002)	0.002 (0.001)	0.002 (0.001)
Dominant	0.059 (0.018)	0.054 (0.014)	0.001 (0.002)	0.001 (0.001)	0.001 (0.001)	0.001 (0.001)
Recessive	0.057 (0.001)	0.054 (0.017)	0.000 (0.000)	0.007 (0.005)	0.000 (0.000)	0.007 (0.005)

^aFalse positives are counts out of 7 for Distribution 1 and out of 10 for Distribution 2.

^bModel refers to genetic model adopted in the analysis.

^cMeasures in this table are found by averaging over 2,000 simulated data sets; standard errors (SE) are shown in parentheses. (Additionally, the bias and MSE were first averaged across all null haplotypes in the model)

Tables VIII–XI. For each simulation, the results are broken down into two broad categories—correct analysis versus additive analysis. Correct analysis refers to specifying the correct genetic model when analyzing the data using *haplo.glm* or HAPSTAT, while additive analysis refers to analyzing non-additive data additively.

NULL SIMULATION

For each haplotype distribution, both aPro and aRetro have desirable and similar performances under the null model (Table III). The FP count is low for both methods (at most 0.033 FPs out of 7 variables for Distribution 1, and at most 0.062 FPs out of 10 variables for Distribution 2). The effect estimations are also very similar, with the bias from both methods ranging from -0.001 to 0.007, and MSE ranging from 0.000 to 0.007. When comparing results between the two haplotype distributions, the FP count is lower for Distribution 1 than Distribution 2. This result is not unexpected, as the dimension of Distribution 2 is larger.

SIMULATION I FOR HAPLOTYPE DISTRIBUTION 1

Additive genetic mode. When the risk haplotype acts additively, the FN and FP counts are very close for aPro and aRetro (Table IV). The FN count decreases when the effect size of the risk haplotype increases or if the frequency of the risk haplotype increases, but the relative performance of the two methods stays the same. When comparing the bias and MSE of aPro and aRetro, the measures from the two procedures are also very close, which indicates that aPro and aRetro perform similarly with respect to effect estimation. For both procedures, the bias on the effect estimates is negative and the magnitudes

are larger than what have been reported for an unpenalized likelihood analysis (e.g. Lin and Zeng [2006] and Tables XIV and XV). These results are not unexpected when using a penalized likelihood approach. When a haplotype is not included in the model, its effect estimate is shrunk toward zero or set to exactly zero. Shrinkage can cause a large bias on the effect estimates. The impact of using a penalized method on the bias is greatest when the effect size is large and the power to detect the risk haplotype is low. As a result, a decrease in effect size or an increase in power does not necessarily guarantee a reduction in bias; the magnitude of the bias is a compromise between these two factors. This phenomenon is seen when examining the biases in Table IV. First, the bias for rare-risk haplotypes is larger than the corresponding common-risk haplotypes with the same effect size because the power to detect rare-risk haplotype is smaller. Second, for a given risk haplotype, the bias on the effect estimate increases as the effect size increases until the power to detect the risk haplotype becomes large enough to overcome the shrinkage, and the bias on the effect estimate begins to decrease.

A similar pattern is observed when examining the MSE of the two procedures. Again, the MSE of aPro and aRetro is larger than what has been found in an unpenalized likelihood analysis. MSE is an estimation measure that incorporates both the variance of an estimator and its bias. Because the effect estimates obtained from penalized methods are typically more efficient than those obtained from the corresponding unpenalized methods, it appears that the MSE of the effect estimates from aPro or aRetro could be dominated by their biases.

Dominant genetic mode. Under a dominant genetic mode, aRetro performs slightly better than aPro when the data are analyzed under the correct genetic model

TABLE IV. BIC-penalized results of Sim I (one-risk haplotype) for Distribution 1—correct analysis^a

Mode ^c	Freq	OR	Model selection results				Parameter estimation results			
			False negatives ^b		False positives		Bias		MSE	
			Pro	Retro	Pro	Retro	Pro	Retro	Pro	Retro
Additive	Rare	1.3	0.943 ^d (0.007)	0.945 (0.007)	0.019 (0.004)	0.024 (0.005)	-0.230 (0.004)	-0.231 (0.004)	0.071 (0.001)	0.071 (0.001)
		1.5	0.820 (0.012)	0.805 (0.013)	0.010 (0.007)	0.015 (0.006)	-0.304 (0.005)	-0.298 (0.006)	0.142 (0.003)	0.140 (0.004)
		1.7	0.623 (0.015)	0.614 (0.015)	0.051 (0.008)	0.050 (0.008)	-0.305 (0.01)	-0.304 (0.01)	0.186 (0.004)	0.183 (0.004)
		2.0	0.244 (0.014)	0.239 (0.013)	0.059 (0.008)	0.064 (0.008)	-0.201 (0.010)	-0.200 (0.010)	0.144 (0.006)	0.142 (0.006)
	Common	1.3	0.801 (0.013)	0.792 (0.013)	0.030 (0.005)	0.027 (0.005)	-0.190 (0.005)	-0.188 (0.005)	0.059 (0.001)	0.058 (0.001)
		1.5	0.465 (0.016)	0.460 (0.016)	0.075 (0.006)	0.075 (0.007)	-0.189 (0.005)	-0.187 (0.007)	0.082 (0.002)	0.081 (0.001)
		1.7	0.120 (0.010)	0.112 (0.010)	0.066 (0.009)	0.073 (0.009)	-0.103 (0.006)	-0.101 (0.006)	0.050 (0.003)	0.048 (0.003)
		2.0	0.003 (0.002)	0.003 (0.002)	0.054 (0.008)	0.056 (0.008)	-0.071 (0.005)	-0.072 (0.005)	0.027 (0.001)	0.027 (0.001)
Dominant	Rare	1.3	0.946 (0.007)	0.941 (0.007)	0.018 (0.004)	0.028 (0.005)	-0.230 (0.004)	-0.228 (0.004)	0.072 (0.001)	0.071 (0.001)
		1.5	0.870 (0.011)	0.840 (0.012)	0.025 (0.008)	0.015 (0.007)	-0.324 (0.006)	-0.311 (0.007)	0.151 (0.002)	0.146 (0.002)
		1.7	0.659 (0.015)	0.632 (0.015)	0.034 (0.006)	0.031 (0.006)	-0.317 (0.010)	-0.307 (0.010)	0.196 (0.004)	0.188 (0.004)
		2.0	0.332 (0.015)	0.303 (0.015)	0.064 (0.009)	0.069 (0.009)	-0.251 (0.011)	-0.239 (0.011)	0.184 (0.007)	0.170 (0.007)
	Common	1.3	0.864 (0.011)	0.849 (0.011)	0.024 (0.006)	0.025 (0.005)	-0.207 (0.004)	-0.203 (0.005)	0.063 (0.001)	0.062 (0.001)
		1.5	0.590 (0.016)	0.555 (0.016)	0.060 (0.007)	0.055 (0.007)	-0.231 (0.006)	-0.224 (0.007)	0.101 (0.002)	0.096 (0.002)
		1.7	0.263^e (0.014)	0.206 (0.013)	0.064 (0.008)	0.067 (0.009)	-0.161 (0.008)	-0.145 (0.007)	0.088 (0.004)	0.073 (0.003)
		2.0	0.035 (0.006)	0.025 (0.005)	0.062 (0.008)	0.065 (0.008)	-0.090 (0.006)	-0.091 (0.006)	0.046 (0.003)	0.041 (0.003)
Recessive	Rare	2.0	0.990 (0.003)	0.910 (0.011)	0.010 (0.007)	0.035 (0.013)	-0.693 (0.000)	-0.583 (0.026)	0.480 (0.000)	0.474 (0.005)
		2.5	0.990 (0.003)	0.870 (0.013)	0.010 (0.007)	0.010 (0.007)	-0.916 (0.000)	-0.753 (0.031)	0.840 (0.000)	0.753 (0.016)
		3.0	0.990 (0.002)	0.770 (0.014)	0.010 (0.007)	0.015 (0.009)	-1.099 (0.000)	-0.764 (0.045)	1.207 (0.000)	0.987 (0.030)
		3.5	0.995 (0.001)	0.720 (0.009)	0.005 (0.005)	0.005 (0.005)	-1.253 (0.000)	-0.865 (0.046)	1.569 (0.000)	1.163 (0.047)
	Common	2.0	0.925 (0.008)	0.695 (0.015)	0.000 (0.000)	0.010 (0.007)	-0.603 (0.023)	-0.441 (0.028)	0.467 (0.006)	0.347 (0.014)
		2.5	0.735 (0.014)	0.350 (0.015)	0.010 (0.007)	0.030 (0.012)	-0.549 (0.045)	-0.329 (0.033)	0.711 (0.023)	0.321 (0.027)
		3.0	0.480 (0.016)	0.110 (0.010)	0.015 (0.009)	0.045 (0.016)	-0.440 (0.047)	-0.226 (0.026)	0.631 (0.040)	0.188 (0.026)
		3.5	0.355 (0.015)	0.060 (0.008)	0.030 (0.012)	0.020 (0.010)	-0.367 (0.052)	-0.242 (0.024)	0.671 (0.055)	0.176 (0.026)

^aCorrect analysis means the genetic model adopted in the analysis is the same as the underlying genetic mode.

^bFalse negatives are counts out of 1 and false positives are counts out of 6.

^cMode refers to underlying genetic mode of data; Freq refers to frequency of risk haplotype; OR refers to the odds ratio of the risk haplotype.

^dMeasures in this table are found by averaging over 1,000 simulated data sets; standard errors (SE) are shown in parentheses.

^eBolded measures (pro vs. retro) are statistically significantly different at 0.05 level (i.e. the ±2SE intervals do NOT overlap).

TABLE V. BIC-penalized results of Sim I (one-risk haplotype) for Distribution 1—additive analysis^a

Mode ^c	Freq	OR	Model selection results				Parameter estimation results			
			False negatives ^b		False positives		Bias		MSE	
			Pro	Retro	Pro	Retro	Pro	Retro	Pro	Retro
Dominant	Rare	1.3	0.949 ^d (0.007)	0.956 (0.006)	0.023 (0.005)	0.020 (0.004)	-0.234 (0.004)	-0.237 (0.004)	0.070 (0.001)	0.070 (0.001)
		1.5	0.865 (0.011)	0.860 (0.011)	0.025 (0.007)	0.025 (0.006)	-0.326 (0.006)	-0.324 (0.006)	0.149 (0.003)	0.148 (0.004)
		1.7	0.665 (0.015)	0.677 (0.015)	0.040 (0.006)	0.040 (0.006)	-0.331 (0.009)	-0.342 (0.009)	0.195 (0.004)	0.197 (0.004)
		2.0	0.356 (0.015)	0.366 (0.015)	0.074 (0.009)	0.083 (0.010)	-0.286 (0.011)	-0.304 (0.010)	0.194 (0.007)	0.198 (0.007)
	Common	1.3	0.879 (0.010)	0.883 (0.010)	0.024 (0.005)	0.025 (0.006)	-0.219 (0.004)	-0.221 (0.004)	0.062 (0.001)	0.062 (0.001)
		1.5	0.650 (0.015)	0.660 (0.015)	0.060 (0.007)	0.050 (0.007)	-0.273 (0.004)	-0.281 (0.005)	0.110 (0.004)	0.111 (0.003)
		1.7	0.291 (0.014)	0.301 (0.015)	0.064 (0.008)	0.068 (-0.009)	-0.225 (0.007)	-0.242 (0.007)	0.099 (0.004)	0.103 (0.004)
		2.0	0.049 (0.007)	0.051 (0.007)	0.081 (0.009)	0.090 (0.010)	-0.186 (0.006)	-0.216 (0.006)	0.068 (0.003)	0.077 (0.003)
Recessive	Rare	2.0	0.985 (0.004)	0.985 (0.004)	0.025 (0.011)	0.035 (0.013)	-0.688 (0.006)	-0.689 (0.005)	0.481 (0.006)	0.481 (0.006)
		2.5	1.000^e (0.000)	0.995 (0.002)	0.025 (0.011)	0.025 (0.011)	-0.916 (0.000)	-0.914 (0.002)	0.840 (0.000)	0.836 (0.003)
		3.0	0.985 (0.004)	0.980 (0.004)	0.025 (0.011)	0.030 (0.012)	-1.089 (0.005)	-1.088 (0.006)	1.192 (0.009)	1.189 (0.009)
		3.5	1.000 (0.000)	0.990 (0.003)	0.000 (0.000)	0.005 (0.005)	-1.253 (0.000)	-1.248 (0.003)	1.569 (0.000)	1.561 (0.006)
	Common	2.0	0.975 (0.005)	0.950 (0.007)	0.005 (0.005)	0.020 (0.010)	-0.685 (0.004)	-0.684 (0.004)	0.472 (0.004)	0.470 (0.004)
		2.5	0.890 (0.010)	0.890 (0.010)	0.015 (0.009)	0.025 (0.013)	-0.876 (0.008)	-0.877 (0.008)	0.781 (0.012)	0.782 (0.012)
		3.0	0.860 (0.011)	0.845 (0.011)	0.015 (0.009)	0.030 (0.012)	-1.051 (0.008)	-1.044 (0.009)	1.120 (0.016)	1.107 (0.017)
		3.5	0.740 (0.014)	0.715 (0.014)	0.010 (0.007)	0.020 (0.010)	-1.156 (0.012)	-1.146 (0.012)	1.365 (0.025)	1.343 (0.026)

^aAdditive analysis means an additive model is adopted in the analysis when the underlying genetic mode is not.

^bFalse negatives are counts out of 1 and false positives are counts out of 6.

^cMode refers to underlying genetic mode of data; Freq refers to frequency of risk haplotype; OR refers to the odds ratio of the risk haplotype.

^dMeasures in this table are found by averaging over 1,000 simulated data sets; standard errors (SE) are shown in parentheses.

^eBolded measures (pro vs. retro) are statistically significantly different at 0.05 level (i.e. the ±2SE intervals do NOT overlap).

(Table IV). While having comparable FP counts, the FN count of aRetro is less than that of aPro although the differences are not always significant (i.e. the intervals of mean ±2SEs overlap). Similar results are observed for the estimation measures. The bias and MSE of aRetro are similar to or lower than that of aPro.

When the data are incorrectly analyzed additively (Table V), the FP count is similar for both methods, and stays roughly at the same level as the correct analysis. However, both methods suffer from an increase in FNs, bias, and MSE when compared to the performance of the correct analysis. It appears that the impact of using an incorrect genetic model is larger on aRetro than on aPro. For example, the aRetro FN count of retaining a common dominant risk haplotype with OR = 1.7 increased from 0.206 to 0.301 under the additive analysis, while the FN count of aPro increased from 0.263 to 0.291, with the SEs of

all figures around 0.014. The aRetro bias in this setting increases from |−0.145| to |−0.242|, while the bias of aPro increases from |−0.161| to |−0.225|. As a result, aPro performs worse than aRetro under correct analysis (Table IV) but is comparable or slightly better than aRetro under the additive analysis (Table V).

Recessive genetic mode. Under a recessive genetic mode, aRetro clearly outperforms aPro when the data are analyzed under the correct model (Table IV). Both methods have comparable FP counts, but the FN count of aRetro is significantly less than that of aPro. When the risk haplotype is rare, the FN count of aPro is almost 1, indicating that aPro has almost no power to retain the risk haplotype in the model under this scenario. When the risk haplotype is common, the FN count of aRetro is 75% to 17% smaller than the FN count of aPro. The lack of power of aPro also manifests in significantly more bias on the effect estimates.

The bias on effect estimates from aRetro is 84% to 50% of the bias from aPro. The MSE on the effect estimates from aRetro is 90% to 26% of the MSE from aPro in most cases.

When this data are incorrectly analyzed additively (Table V), the performance of each method suffers from an increase in FNs, bias, and MSE. The magnitude of the performance loss due to incorrect modeling is more severe than what was observed under a dominant mode and is much more severe for aRetro than aPro. For example, the aRetro FN count of retaining a common recessive risk haplotype with OR = 2.5 increased from 0.350 to 0.890 under the additive analysis, while the FN count of aPro increased from 0.735 to 0.890, with the SEs of all figures

around 0.014. The aRetro bias in this setting increases from |−0.329| to |−0.877|, while the bias of aPro increases from |−0.549| to |−0.876|. Consequently, while aRetro exhibits absolute superiority over the aPro method (Table IV) under a correct analysis, it becomes comparable to or slightly better than aPro under the additive analysis (Table V).

SIMULATION II FOR HAPLOTYPE DISTRIBUTION 1

Simulation II examines the performance of aPro and aRetro when two haplotypes are associated with the disease (Tables VI and VII). Under each genetic mode, the patterns

TABLE VI. BIC-penalized results of Sim II (two-risk haplotypes) for Distribution 1—correct analysis^a

Mode ^c	Freq		Model selection results				Parameter estimation results			
			False negatives ^b		False positives		Bias		MSE	
			Pro	Retro	Pro	Retro	Pro	Retro	Pro	Retro
Additive	R/R	R1 ^d	1.548 ^e (0.020)	1.536 (0.021)	0.029 (0.006)	0.031 (0.006)	−0.443 (0.008)	−0.444 (0.008)	0.267 (0.002)	0.266 (0.003)
		R2					−0.327 (0.009)	−0.324 (0.009)	0.193 (0.004)	0.190 (0.004)
	R/C	R1	0.596 (0.020)	0.588 (0.020)	0.089 (0.010)	0.098 (0.011)	−0.250 (0.006)	−0.249 (0.006)	0.156 (0.003)	0.155 (0.002)
		C1					−0.109 (0.010)	−0.109 (0.010)	0.048 (0.004)	0.046 (0.004)
	C/C	C1	0.204 (0.017)	0.205 (0.017)	0.095 (0.010)	0.090 (0.010)	−0.113 (0.006)	−0.115 (0.006)	0.053 (0.002)	0.054 (0.002)
		C2					−0.101 (0.006)	−0.103 (0.006)	0.042 (0.003)	0.042 (0.003)
Dominant	R/R	R1	1.628 (0.020)	1.566 (0.021)	0.036 (0.006)	0.039 (0.006)	−0.441 (0.009)	−0.433 (0.009)	0.270 (0.003)	0.261 (0.003)
		R2					−0.371 (0.009)	−0.352 (0.009)	0.214 (0.004)	0.202 (0.004)
	R/C	R1	0.741 (0.022)	0.704 (0.022)	0.077 (0.009)	0.084 (0.010)	−0.253 (0.008)	−0.252 (0.007)	0.167 (0.003)	0.162 (0.003)
		C1					−0.142 (0.010)	−0.133 (0.010)	0.077 (0.004)	0.071 (0.004)
	C/C	C1	0.477 (0.024)	0.433 (0.021)	0.087 (0.010)	0.114 (0.011)	−0.170 (0.007)	−0.149 (0.006)	0.092 (0.003)	0.074 (0.003)
		C2					−0.152 (0.008)	−0.130 (0.007)	0.076 (0.004)	0.059 (0.003)
Recessive	R/R	R1	2.000^f (0.000)	1.865 (0.024)	0.005 (0.005)	0.020 (0.010)	−0.916 (0.000)	−0.916 (0.000)	0.840 (0.000)	0.840 (0.000)
		R2					−0.916 (0.000)	−0.731 (0.034)	0.840 (0.000)	0.766 (0.015)
	R/C	R1	1.775 (0.030)	1.215 (0.046)	0.010 (0.007)	0.030 (0.014)	−0.916 (0.000)	−0.684 (0.038)	0.840 (0.000)	0.749 (0.019)
		C1					−0.618 (0.041)	−0.367 (0.033)	0.720 (0.026)	0.354 (0.027)
	C/C	C1	1.490 (0.054)	0.595 (0.053)	0.000 (0.000)	0.050 (0.015)	−0.662 (0.043)	−0.319 (0.028)	0.717 (0.029)	0.322 (0.025)
		C2					−0.547 (0.037)	−0.292 (0.033)	0.664 (0.022)	0.244 (0.027)

^aCorrect analysis means the genetic model adopted in the analysis is the same as the underlying genetic mode.

^bFalse negatives are counts out of 2 and false positives are counts out of 5.

^cMode refers to underlying genetic mode of data; Freq refers to frequency of risk haplotype.

^dR refers to a rare risk haplotype and C refers to a common risk haplotype.

^eMeasures in this table are found by averaging over 1,000 simulated data sets; standard errors (SE) are shown in parentheses.

^fBolded measures (pro vs. retro) are statistically significantly different at 0.05 level (i.e. the ±2SE intervals do NOT overlap).

TABLE VII. BIC-penalized results of Sim II (two-risk haplotypes) for Distribution 1—additive analysis^a

Mode ^c	Freq		Model selection results				Parameter estimation results			
			False negatives ^b		False positives		Bias		MSE	
			Pro	Retro	Pro	Retro	Pro	Retro	Pro	Retro
Dominant	R/R	R1 ^d	1.653 ^e (0.019)	1.660 (0.019)	0.028 (0.005)	0.029 (0.005)	-0.448 (0.008)	-0.452 (0.008)	0.268 (0.002)	0.265 (0.002)
		R2					-0.391 (0.008)	-0.397 (0.008)	0.218 (0.004)	0.219 (0.004)
	R/C	R1	0.790 (0.024)	0.789 (0.024)	0.079 (0.009)	0.083 (0.010)	-0.261 (0.007)	-0.265 (0.007)	0.163 (0.004)	0.160 (0.004)
		C1					-0.216 (0.01)	-0.233 (0.009)	0.094 (0.004)	0.098 (0.004)
	C/C	C1	0.594 (0.026)	0.603 (0.026)	0.085 (0.010)	0.093 (0.010)	-0.243 (0.007)	-0.255 (0.006)	0.110 (0.004)	0.113 (0.004)
		C2					-0.225 (0.007)	-0.237 (0.007)	0.096 (0.004)	0.098 (0.004)
Recessive	R/R	R1	2.000 (0.000)	2.000 (0.000)	0.015 (0.009)	0.025 (0.011)	-0.916 (0.000)	-0.916 (0.000)	0.840 (0.000)	0.840 (0.000)
		R2					-0.916 (0.000)	-0.916 (0.000)	0.840 (0.000)	0.840 (0.000)
	R/C	R1	1.910 (0.021)	1.910 (0.021)	0.025 (0.011)	0.030 (0.012)	-0.916 (0.000)	-0.916 (0.000)	0.843 (0.011)	0.844 (0.011)
		C1					-0.885 (0.005)	-0.884 (0.005)	0.795 (0.009)	0.794 (0.009)
	C/C	C1	1.875 (0.028)	1.865 (0.029)	0.015 (0.009)	0.025 (0.011)	-0.895 (0.006)	-0.894 (0.006)	0.809 (0.009)	0.807 (0.009)
		C2					-0.895 (0.006)	-0.894 (0.006)	0.808 (0.009)	0.806 (0.009)

^aAdditive analysis means an additive model is adopted in the analysis when the underlying genetic mode is not.

^bFalse negatives are counts out of 2 and false positives are counts out of 5.

^cMode refers to underlying genetic mode of data; Freq refers to frequency of risk haplotype.

^dR refers to a rare risk haplotype and C refers to a common risk haplotype.

^eMeasures in this table are found by averaging over 1,000 simulated data sets; standard errors (SE) are shown in parentheses.

of results observed in Simulation I remain the same in Simulation II. First, when the risk haplotypes act additively on disease susceptibility, the performances of aPro and aRetro are comparable for FNs, FPs, bias and MSE. Second, the performance of aRetro is better than that of aPro under a dominant mode with correct analysis. The gain brought by aRetro is similar across all three simulation scenarios but is rarely significantly different from aPro. However, when analyzing the data with an additive model, all measures increase for both methods. The performance loss is more severe in aRetro, resulting in a comparable performance of aPro and aRetro. Finally, under a recessive mode analyzed correctly, aRetro has significantly fewer FNs, slightly more FPs, and significantly smaller bias/MSE when compared to aPro for almost all simulation scenarios. However, when recessively acting haplotypes are analyzed using an additive genetic model, the performance of each procedure suffers, especially aRetro. Both methods essentially lose their power to detect the two-risk haplotypes (e.g. having FN counts close to 2 as no-risk haplotypes were retained in the model) and yield sizable biases/MSE.

SIMULATION I AND II FOR HAPLOTYPE DISTRIBUTION 2

When comparing results between the two haplotype distributions for Simulations I and II, all measures were

typically higher for Distribution 2. Like for the Null Simulation, these results are not unexpected because Distribution 2 has a larger dimension, which means more parameters needed to be estimated in the analysis. Increasing the number of parameters and using the same amount of data for estimation can decrease power (i.e. increase the FN count) and increase bias. Although the magnitudes of model selection and estimation measures differ between the two haplotype distributions, when comparing the relative performance of aPro and aRetro, the pattern of results observed in the first haplotype distribution is similar to that in the second haplotype distribution for both Simulations I and II (Tables VIII–XI). The relative performance of aPro and aRetro depends on both the underlying genetic mode of the risk haplotypes and the genetic model adopted in the analysis. When the haplotypes associated with disease risk act additively, the two procedures perform comparably with respect to model selection and estimation measures. Under a dominant mode, aRetro performs slightly better than aPro and substantially better under a recessive mode. When these data are analyzed using an additive model, both procedures suffer from a loss in power and an increase in bias/MSE. The impact of imposing the incorrect genetic model is more severe for aRetro than for aPro, and the performance gain of aRetro is lost.

TABLE VIII. BIC-penalized results of Sim I (one-risk haplotype) for Distribution 2—correct analysis^a

Mode ^c	Freq	OR	Model selection results				Parameter estimation results			
			False Negatives ^b		False Positives		Bias		MSE	
			Pro	Retro	Pro	Retro	Pro	Retro	Pro	Retro
Additive	Rare	1.3	0.965 ^d (0.006)	0.960 (0.006)	0.065 (0.017)	0.035 (0.013)	-0.245 (0.006)	-0.244 (0.006)	0.068 (0.000)	0.068 (0.000)
		1.5	0.885 (0.010)	0.865 (0.011)	0.090 (0.021)	0.045 (0.015)	-0.338 (0.014)	-0.329 (0.014)	0.151 (0.003)	0.148 (0.003)
		1.7	0.675 (0.015)	0.685 (0.015)	0.080 (0.02)	0.070 (0.021)	-0.349 (0.020)	-0.351 (0.020)	0.198 (0.009)	0.201 (0.009)
		2.0	0.290 (0.014)	0.290 (0.014)	0.140 (0.027)	0.080 (0.022)	-0.250 (0.023)	-0.250 (0.023)	0.169 (0.015)	0.168 (0.015)
	Common	1.3	0.715 (0.014)	0.710 (0.014)	0.085 (0.021)	0.075 (0.020)	-0.229 (0.009)	-0.221 (0.010)	0.067 (0.001)	0.067 (0.001)
		1.5	0.385^e (0.015)	0.315 (0.015)	0.050 (0.017)	0.050 (0.017)	-0.311 (0.013)	-0.286 (0.014)	0.130 (0.005)	0.121 (0.005)
		1.7	0.125 (0.010)	0.120 (0.010)	0.065 (0.020)	0.080 (0.023)	-0.443 (0.013)	-0.382 (0.016)	0.233 (0.008)	0.199 (0.009)
		2.0	0.005 (0.002)	0.010 (0.003)	0.090 (0.025)	0.160 (0.034)	-0.646 (0.013)	-0.583 (0.018)	0.450 (0.009)	0.404 (0.012)
Dominant	Rare	1.3	0.955 (0.007)	0.930 (0.008)	0.085 (0.021)	0.040 (0.016)	-0.234 (0.010)	-0.219 (0.012)	0.073 (0.003)	0.076 (0.004)
		1.5	0.895 (0.010)	0.880 (0.010)	0.085 (0.023)	0.045 (0.016)	-0.343 (0.013)	-0.336 (0.014)	0.152 (0.003)	0.151 (0.003)
		1.7	0.770 (0.013)	0.690 (0.015)	0.125 (0.023)	0.110 (0.024)	-0.388 (0.019)	-0.346 (0.021)	0.224 (0.008)	0.203 (0.009)
		2.0	0.510 (0.016)	0.425 (0.016)	0.125 (0.023)	0.080 (0.020)	-0.361 (0.026)	-0.322 (0.025)	0.260 (0.016)	0.225 (0.016)
	Common	1.3	0.815 (0.012)	0.800 (0.013)	0.030 (0.012)	0.035 (0.013)	-0.240 (0.007)	-0.234 (0.008)	0.067 (0.001)	0.068 (0.001)
		1.5	0.545 (0.016)	0.520 (0.016)	0.070 (0.019)	0.110 (0.025)	-0.296 (0.015)	-0.273 (0.016)	0.132 (0.004)	0.123 (0.005)
		1.7	0.270 (0.014)	0.235 (0.013)	0.075 (0.019)	0.100 (0.025)	-0.376 (0.018)	-0.339 (0.019)	0.205 (0.008)	0.185 (0.009)
		2.0	0.020 (0.004)	0.010 (0.003)	0.110 (0.025)	0.130 (0.026)	-0.509 (0.021)	-0.441 (0.024)	0.349 (0.014)	0.304 (0.016)
Recessive	Rare	2.0	1.000 (0.000)	0.930 (0.008)	0.050 (0.015)	0.035 (0.013)	-0.693 (0.000)	-0.610 (0.022)	0.480 (0.000)	0.467 (0.005)
		2.5	1.000 (0.000)	0.810 (0.012)	0.050 (0.015)	0.075 (0.019)	-0.916 (0.000)	-0.660 (0.038)	0.840 (0.000)	0.728 (0.018)
		3.0	1.000 (0.000)	0.810 (0.012)	0.070 (0.018)	0.025 (0.011)	-1.099 (0.000)	-0.859 (0.036)	1.207 (0.000)	0.995 (0.031)
		3.5	1.000 (0.000)	0.685 (0.015)	0.045 (0.015)	0.055 (0.018)	-1.253 (0.000)	-0.841 (0.045)	1.569 (0.000)	1.103 (0.049)
	Common	2.0	0.935 (0.008)	0.820 (0.012)	0.025 (0.011)	0.060 (0.020)	-0.693 (0.000)	-0.539 (0.024)	0.481 (0.000)	0.404 (0.012)
		2.5	0.880 (0.010)	0.610 (0.015)	0.010 (0.007)	0.070 (0.019)	-0.856 (0.021)	-0.555 (0.034)	0.824 (0.007)	0.531 (0.028)
		3.0	0.700 (0.014)	0.465 (0.016)	0.020 (0.010)	0.060 (0.018)	-0.932 (0.033)	-0.548 (0.038)	1.087 (0.024)	0.590 (0.041)
		3.5	0.585 (0.016)	0.405 (0.016)	0.000 (0.000)	0.100 (0.024)	-0.974 (0.041)	-0.590 (0.041)	1.284 (0.041)	0.690 (0.052)

^aCorrect analysis means the genetic model adopted in the analysis is the same as the underlying genetic mode.

^bFalse negatives are counts out of 1 and false positives are counts out of 9.

^cMode refers to underlying genetic mode of data; Freq refers to frequency of risk haplotype; OR refers to the odds ratio of the risk haplotype.

^dMeasures in this table are found by averaging over 1,000 simulated data sets; standard errors (SE) are shown in parentheses.

^eBolded measures (pro vs. retro) are statistically significantly different at 0.05 level (i.e. the $\pm 2SE$ intervals do NOT overlap).

TABLE IX. BIC-penalized results of Sim I (one-risk haplotype) for Distribution 2—additive analysis^a

Mode ^c	Freq	OR	Model selection results				Parameter estimation results			
			False negatives ^b		False positives		Bias		MSE	
			Pro	Retro	Pro	Retro	Pro	Retro	Pro	Retro
Dominant	Rare	1.3	0.950 ^d (0.007)	0.955 (0.007)	0.095 (0.023)	0.045 (0.015)	-0.232 (0.010)	-0.235 (0.009)	0.073 (0.003)	0.072 (0.002)
		1.5	0.900 (0.009)	0.900 (0.009)	0.070 (0.022)	0.045 (0.016)	-0.351 (0.012)	-0.351 (0.012)	0.151 (0.003)	0.152 (0.003)
		1.7	0.795 (0.013)	0.790 (0.013)	0.130 (0.024)	0.075 (0.019)	-0.406 (0.018)	-0.408 (0.017)	0.229 (0.008)	0.227 (0.008)
		2.0	0.505 (0.016)	0.550 (0.016)	0.100 (0.021)	0.070 (0.019)	-0.390 (0.024)	-0.414 (0.023)	0.263 (0.016)	0.281 (0.016)
	Common	1.3	0.820 (0.012)	0.805 (0.013)	0.035 (0.013)	0.025 (0.011)	-0.243 (0.006)	-0.245 (0.006)	0.067 (0.001)	0.067 (0.001)
		1.5	0.565 (0.016)	0.570 (0.016)	0.075 (0.019)	0.075 (0.019)	-0.314 (0.013)	-0.308 (0.013)	0.133 (0.004)	0.129 (0.005)
		1.7	0.275 (0.014)	0.270 (0.014)	0.075 (0.020)	0.100 (0.024)	-0.397 (0.016)	-0.385 (0.016)	0.208 (0.008)	0.200 (0.009)
		2.0	0.015 (0.004)	0.020 (0.004)	0.090 (0.025)	0.100 (0.026)	-0.539 (0.018)	-0.478 (0.019)	0.358 (0.014)	0.304 (0.015)
Recessive	Rare	2.0	1.000 (0.000)	1.000 (0.000)	0.075 (0.020)	0.025 (0.011)	-0.693 (0.000)	-0.693 (0.000)	0.480 (0.000)	0.480 (0.000)
		2.5	0.995 (0.002)	0.995 (0.002)	0.105 (0.026)	0.055 (0.020)	-0.912 (0.004)	-0.912 (0.004)	0.835 (0.004)	0.835 (0.004)
		3.0	0.995 (0.002)	0.990 (0.003)	0.080^e (0.019)	0.010 (0.007)	-1.096 (0.002)	-1.094 (0.003)	1.203 (0.004)	1.199 (0.006)
		3.5	0.985 (0.004)	0.985 (0.004)	0.075 (0.020)	0.035 (0.015)	-1.244 (0.005)	-1.242 (0.006)	1.552 (0.010)	1.551 (0.011)
	Common	2.0	0.930 (0.008)	0.950 (0.007)	0.035 (0.013)	0.015 (0.009)	-0.691 (0.002)	-0.691 (0.002)	0.478 (0.002)	0.478 (0.002)
		2.5	0.890 (0.010)	0.890 (0.010)	0.045 (0.015)	0.055 (0.016)	-0.903 (0.005)	-0.895 (0.007)	0.821 (0.007)	0.809 (0.009)
		3.0	0.790 (0.013)	0.805 (0.013)	0.040 (0.014)	0.020 (0.010)	-1.087 (0.005)	-1.078 (0.006)	1.186 (0.009)	1.171 (0.011)
		3.5	0.730 (0.014)	0.715 (0.014)	0.020 (0.010)	0.030 (0.014)	-1.230 (0.007)	-1.216 (0.009)	1.523 (0.014)	1.495 (0.018)

^aAdditive analysis means an additive model is adopted in the analysis when the underlying genetic mode is not.

^bFalse negatives are counts out of 1 and false positives are counts out of 9.

^cMode refers to underlying genetic mode of data; Freq refers to frequency of risk haplotype; OR refers to the odds ratio of the risk haplotype.

^dMeasures in this table are found by averaging over 1,000 simulated data sets; standard errors (SE) are shown in parentheses.

^eBolded measures (pro vs. retro) are statistically significantly different at 0.05 level (i.e. the ±2SE intervals do NOT overlap).

ANALYSIS OF GAW15 RA DATA

Rheumatoid arthritis (RA) is a condition that is believed to be influenced by a number of genetic and environmental factors. Genetic Analysis Workshop (GAW) 15 designed a series of studies to investigate the genetic risk factors of this condition. One of these studies focused on a dense panel of 2,300 SNPs from a 10-megabase region of chromosome 18q. This region of the genome was chosen because it had shown prior evidence of linkage to RA [Browning and Thomas, 2007]. The panel of SNPs was genotyped for 460 cases and 460 controls collected by the North American Rheumatoid Arthritis Consortium. Before implementing the prospective and retrospective analyses discussed in this article, the quality control filters described by Browning and Thomas [2007] were applied to the data, and five SNPs in Hardy-Weinberg

disequilibrium and one case with no genotype information were removed from the data set prior to analysis. We analyzed this data using the four methods discussed in this article—unpenalized prospective likelihood, prospective likelihood coupled with the ALASSO penalty, unpenalized retrospective likelihood, and retrospective likelihood coupled with the ALASSO penalty. The unpenalized methods were fit using haplo.glm in R and HAPSTAT in Linux for the prospective and retrospective likelihoods, respectively. The penalized methods were fit using the LSA code posted on the corresponding author’s website.

Based on the findings of Browning and Thomas, we analyzed two blocks of linkage disequilibrium (LD) around the region that showed significant association (SNPs 1631 and 1632). The two regions we investigated were (a) SNPs 1626 to 1632, a block of SNPs that are in

TABLE X. BIC-penalized results of Sim II (two-risk haplotypes) for Distribution 2—correct analysis^a

Mode ^c	Freq		Model selection results				Parameter estimation results			
			False negatives ^b		False positives		Bias		MSE	
			Pro	Retro	Pro	Retro	Pro	Retro	Pro	Retro
Additive	R/R	R1 ^d	1.040 ^e (0.056)	1.025 (0.056)	0.195 (0.037)	0.140 (0.028)	-0.318 (0.021)	-0.301 (0.022)	0.190 (0.009)	0.184 (0.009)
		R2					-0.227 (0.019)	-0.219 (0.019)	0.125 (0.009)	0.122 (0.009)
	R/C	R1	1.030 (0.033)	0.985 (0.054)	0.120 (0.028)	0.120 (0.028)	-0.470 (0.013)	-0.422 (0.017)	0.257 (0.006)	0.234 (0.007)
		C1					-0.417 (0.015)	-0.333 (0.018)	0.219 (0.008)	0.176 (0.009)
	C/C	C1	0.960 (0.037)	0.970 (0.055)	0.095 (0.024)	0.130 (0.029)	-0.467 (0.013)	-0.425 (0.015)	0.251 (0.006)	0.228 (0.007)
		C2					-0.452 (0.013)	-0.416 (0.015)	0.235 (0.007)	0.217 (0.008)
Dominant	R/R	R1	1.120 (0.055)	0.975 (0.054)	0.170 (0.028)	0.130 (0.026)	-0.319 (0.021)	-0.279 (0.021)	0.188 (0.009)	0.166 (0.009)
		R2					-0.248 (0.021)	-0.229 (0.020)	0.151 (0.009)	0.135 (0.009)
	R/C	R1	1.290^f (0.042)	1.060 (0.052)	0.165 (0.029)	0.210 (0.037)	-0.414 (0.019)	-0.337 (0.021)	0.240 (0.007)	0.197 (0.009)
		C1					-0.391 (0.017)	-0.349 (0.018)	0.210 (0.009)	0.184 (0.009)
	C/C	C1	1.180 (0.046)	1.130 (0.056)	0.075 (0.020)	0.095 (0.023)	-0.456 (0.013)	-0.411 (0.015)	0.243 (0.007)	0.212 (0.008)
		C2					-0.434 (0.014)	-0.379 (0.016)	0.229 (0.008)	0.194 (0.009)
Recessive	R/R	R1	2.000 (0.000)	1.610 (0.044)	0.030 (0.012)	0.050 (0.015)	-0.916 (0.000)	-0.749 (0.032)	0.840 (0.000)	0.763 (0.018)
		R2					-0.916 (0.000)	-0.642 (0.034)	0.840 (0.000)	0.642 (0.025)
	R/C	R1	1.880 (0.023)	1.400 (0.045)	0.020 (0.010)	0.045 (0.015)	-0.916 (0.000)	-0.736 (0.032)	0.840 (0.000)	0.746 (0.017)
		C1					-0.846 (0.024)	-0.517 (0.033)	0.827 (0.009)	0.487 (0.028)
	C/C	C1	1.725 (0.034)	1.330 (0.052)	0.010 (0.007)	0.085 (0.022)	-0.885 (0.019)	-0.632 (0.031)	0.854 (0.014)	0.593 (0.026)
		C2					-0.874 (0.017)	-0.604 (0.031)	0.824 (0.007)	0.560 (0.027)

^aCorrect analysis means the genetic model adopted in the analysis is the same as the underlying genetic mode.

^bFalse negatives are counts out of 2 and false positives are counts out of 8.

^cMode refers to underlying genetic mode of data; Freq refers to frequency of risk haplotype.

^dR refers to a rare risk haplotype and C refers to a common risk haplotype.

^eMeasures in this table are found by averaging over 1,000 simulated data sets; standard errors (SE) are shown in parentheses.

^fBolded measures (pro vs. retro) are statistically significantly different at 0.05 level (i.e. the $\pm 2SE$ intervals do NOT overlap).

extremely high LD with SNPs 1631 and 1632 and (b) SNPs 1621 to 1630, a haplotype block [Barrett et al., 2005] that is next to SNPs 1631 and 1632 and located in gene CCBE1. High LD was defined as having a pair-wise $D' \geq 0.99$; D' is a measure of LD that scales the raw difference between the expected and observed haplotype frequencies of a pair of SNPs so that it is between -1 and 1. For each block of SNPs, we considered additive, dominant, and recessive models for haplotype effects. Examining the values of the fitted likelihoods for each method, it appears that the haplotypes in these blocks do not act recessively; so only the results of the additive and dominant analyses are shown. The results of our analysis (coefficient

estimates for all methods and unadjusted p -values for the unpenalized methods) can be found in Tables XII and XIII.

In their analysis, Browning and Thomas located a significant haplotype sequence "21" at SNPs 1631 and 1632. When analyzing the LD block containing these SNPs (a), all four methods considered in this article were able to find the same signal. This result can be seen when examining the small p -value and the non-zero coefficient of the haplotype "1121121" for the unpenalized and penalized methods, respectively (Table XII). All the four methods finding the same result is not surprising and is supported by the results of our simulation studies. Based

TABLE XI. BIC-penalized results of Sim II (two-risk haplotypes) for Distribution 2—additive analysis^a

Mode ^c	Freq		Model selection results				Parameter estimation results			
			False negatives ^b		False positives		Bias		MSE	
			Pro	Retro	Pro	Retro	Pro	Retro	Pro	Retro
Dominant	R/R	R1 ^d	1.140 ^e (0.055)	1.135 (0.056)	0.160 (0.028)	0.110 (0.025)	-0.337 (0.020)	-0.348 (0.019)	0.192 (0.009)	0.195 (0.009)
		R2					-0.286 (0.019)	-0.289 (0.018)	0.153 (0.009)	0.148 (0.009)
	R/C	R1	1.425^f (0.040)	1.120 (0.055)	0.155 (0.032)	0.195 (0.039)	-0.443 (0.016)	-0.402 (0.018)	0.249 (0.006)	0.225 (0.008)
		C1					-0.413 (0.015)	-0.381 (0.015)	0.213 (0.008)	0.192 (0.009)
	C/C	C1	1.135 (0.044)	1.140 (0.050)	0.085 (0.020)	0.125 (0.025)	-0.469 (0.011)	-0.458 (0.012)	0.246 (0.006)	0.238 (0.007)
		C2					-0.458 (0.012)	-0.453 (0.012)	0.241 (0.007)	0.236 (0.007)
Recessive	R/R	R1	1.980 (0.010)	1.960 (0.014)	0.050 (0.015)	0.025 (0.011)	-0.911 (0.004)	-0.911 (0.004)	0.833 (0.005)	0.833 (0.005)
		R2					-0.913 (0.003)	-0.905 (0.005)	0.834 (0.004)	0.823 (0.007)
	R/C	R1	1.890 (0.023)	1.910 (0.017)	0.045 (0.015)	0.025 (0.011)	-0.913 (0.003)	-0.913 (0.003)	0.836 (0.004)	0.836 (0.004)
		C1					-0.902 (0.005)	-0.895 (0.006)	0.819 (0.007)	0.809 (0.009)
	C/C	C1	1.730 (0.032)	1.740 (0.020)	0.045 (0.016)	0.050 (0.017)	-0.910 (0.003)	-0.907 (0.004)	0.830 (0.005)	0.826 (0.006)
		C2					-0.906 (0.005)	-0.901 (0.006)	0.824 (0.007)	0.818 (0.008)

^aAdditive analysis means an additive model is adopted in the analysis when the underlying genetic mode is not.

^bFalse negatives are counts out of 2 and false positives are counts out of 8.

^cMode refers to underlying genetic mode of data; Freq refers to frequency of risk haplotype.

^dR refers to a rare risk haplotype and C refers to a common risk haplotype.

^eMeasures in this table are found by averaging over 1,000 simulated datasets; standard errors (SE) are shown in parentheses.

^fBolded measures (pro vs. retro) are statistically significantly different at 0.05 level (i.e. the $\pm 2SE$ intervals do NOT overlap).

TABLE XII. Analysis of GAW15-RA data SNPs 1626 to 1632^a

Model ^b	Haplotype	Prospective analysis			Retrospective analysis		
		Unpenalized (<i>p</i> -value)	BIC-penalized	AIC-penalized	Unpenalized (<i>p</i> -value)	BIC-penalized	AIC-penalized
Additive	1121111	0.016 (0.883)	0.000	0.000	0.028 (0.784)	0.000	0.000
	1121112	0.069 (0.883)	0.000	0.000	0.076 (0.871)	0.000	0.000
	1121121	-2.028 (0.000)	-2.023	-2.023	-1.979 (0.000)	-1.978	-1.978
	2212222	-0.016 (0.919)	0.000	0.000	-0.016 (0.922)	0.000	0.000
Dominant	1121111	-0.039 (0.797)	0.000	0.000	0.126 (0.335)	0.000	0.000
	1121112	0.043 (0.927)	0.000	0.000	0.110 (0.814)	0.000	0.000
	1121121	-2.022 (0.000)	-2.018	-2.018	-1.967 (0.000)	-1.861	-1.861
	2212222	-0.026 (0.880)	0.000	0.000	0.018 (0.914)	0.000	0.000

^aSNPs in high LD (i.e. $D' \geq 0.99$) with SNPs 1631 and 1632, a significant region identified by Browning and Thomas [2007].

^bModel refers to genetic model adopted in the analysis.

TABLE XIII. Analysis of GAW15-RA data SNPs 1621 to 1630^a

Model ^b	Haplotype	Prospective analysis			Retrospective analysis		
		Unpenalized (<i>p</i> -value)	BIC-penalized	AIC-penalized	Unpenalized (<i>p</i> -value)	BIC-penalized	AIC-penalized
Additive	122121211	0.175 (0.096)	0.000	0.120	0.173 (0.098)	0.000	0.152
	211212212	0.098 (0.539)	0.000	0.000	0.070 (0.656)	0.000	0.000
	221121211	-0.368 (0.355)	0.000	0.000	-0.451 (0.224)	0.000	-0.407
	222121211	-0.762 (0.100)	0.000	-0.522	-0.767 (0.094)	0.000	-0.729
Dominant	122121211	0.283 (0.042)	0.000	0.225	0.246 (0.052)	0.000	0.226
	211212212	0.128 (0.461)	0.000	0.000	0.072 (0.660)	0.000	0.000
	221121211	-0.356 (0.371)	0.000	0.000	-0.446 (0.231)	0.000	-0.398
	222121211	-0.739 (0.111)	0.000	-0.506	-0.768 (0.094)	0.000	-0.727

^aA haplotype block that is next to SNPs 1631 and 1632 and is in gene CCBE1.

^bModel refers to genetic model adopted in the analysis.

on our simulations, we found that the difference between a prospective and retrospective approach is often negligible when the haplotype effect size is large under additive and dominant genetic modes.

Browning and Thomas also found that individuals with the haplotype “21” at SNPs 1631 and 1632 also share the sequence “11211” at SNPs 1626 to 1630. When analyzing the LD block (a) defined by SNPs 1626 to 1632, we replicated the exact haplotype using the article’s four methods. Motivated by this finding, we analyzed the second LD block (b), which is a 10-SNP haplotype block covering SNPs 1626 to 1630 and is located in gene CCBE1. The unpenalized methods (at significance level 0.01) and penalized methods tuned with BIC did not find any significant signals. However, in practice, the significance level is often set more liberally at 0.05; at these levels, one haplotype (“122121211”) sits at the border of significance. More interestingly, the penalized method tuned with AIC identified some of the haplotypes in the sample with the shared sequence “11211” under a prospective analysis and all of them under a retrospective analysis. Although it is unknown whether these findings are true signals or false positives, it is promising that penalized methods are able to detect other haplotypes containing the sequence “11211”, and that a retrospective penalized approach is able to find all haplotypes with the sequence “11211” significant.

DISCUSSION

Like other haplotype-based methods developed to assess haplotype-phenotype association, the success of penalized regression methods depends on the underlying data likelihood. For unpenalized methods, prospective analyses are valid but less efficient than their retrospective counterparts for hypothesis testing, and can result in substantial bias when estimating the haplotype effects. Based on our simulation studies, the same can be said for

penalized methods, which combine testing and estimation into one procedure. We found that the impact of using a prospective likelihood in the analysis depends on the underlying genetic mode of the associated genetic variant and the genetic model adopted in the analysis. When the genetic mode of the haplotypes is known and the correct inheritance model is imposed, using a prospective analysis in place of the more appropriate retrospective analysis is detrimental when the associated haplotypes act dominantly or recessively with respect to disease risk. These results agree with the findings for non-penalized likelihood methods [Satten and Epstein, 2004]. Because the genetic mode of a genetic variant is usually unknown, researchers often analyze the data additively. When dominant or recessive data are analyzed under an additive genetic model, the performance of the prospective and retrospective analyses become comparable. Both methods suffer from an increase in FNs and bias for using an incorrect genetic model. The retrospective analysis appears to be more sensitive to model misspecification and exhibits a larger degree of performance loss, thus making its performance gain over the prospective analysis negligible or slight.

While our simulations focused on penalized methods using the ALASSO penalty with the prospective likelihood of `haplo.glm` and the retrospective likelihood of `HAPSTAT`, we hope our findings can provide insight when coupling other penalized approaches with a prospective or a retrospective likelihood for case-control studies. If the main consideration is the relative performance of the retrospective versus prospective penalized method, then our results suggest that the negative impact of developing haplotype-based penalized methods based on a prospective likelihood for case-control data is non-trivial only when the risk haplotypes act non-additively and the correct genetic model is adopted in the analysis. However, we think a more appropriate way to summarize our findings is to note that a careful haplotype-based penalized analysis of case-control data requires the use

TABLE XIV. Unpenalized results of Sim I (one-risk haplotype) for Distribution 1—correct analysis^a

Mode ^c	Freq	OR	Model selection results				Parameter estimation results			
			False negatives ^b		False positives		Bias		MSE	
			Pro	Retro	Pro	Retro	Pro	Retro	Pro	Retro
Additive	Rare	1.3	0.934 ^d (0.008)	0.934 (0.008)	0.030 (0.006)	0.028 (0.005)	0.013 (0.007)	0.010 (0.007)	0.049 (0.002)	0.048 (0.002)
		1.5	0.865^e (0.012)	0.785 (0.011)	0.050 (0.008)	0.050 (0.007)	0.012 (0.007)	0.005 (0.008)	0.050 (0.001)	0.048 (0.002)
		1.7	0.677 (0.016)	0.680 (0.016)	0.040 (0.007)	0.042 (0.007)	0.005 (0.007)	-0.003 (0.007)	0.047 (0.002)	0.046 (0.002)
		2.0	0.249 (0.013)	0.253 (0.013)	0.032 (0.006)	0.028 (0.005)	0.021 (0.007)	0.013 (0.007)	0.045 (0.002)	0.044 (0.002)
	Common	1.3	0.823 (0.012)	0.820 (0.012)	0.026 (0.005)	0.024 (0.005)	0.011 (0.005)	0.009 (0.005)	0.024 (0.001)	0.023 (0.001)
		1.5	0.575 (0.01)	0.550 (0.011)	0.035 (0.006)	0.035 (0.005)	-0.004 (0.005)	-0.009 (0.005)	0.022 (0.002)	0.022 (0.001)
		1.7	0.168 (0.012)	0.167 (0.012)	0.033 (0.006)	0.032 (0.006)	0.009 (0.005)	0.003 (0.005)	0.023 (0.001)	0.022 (0.001)
		2.0	0.012 (0.003)	0.011 (0.003)	0.028 (0.005)	0.029 (0.005)	0.001 (0.005)	-0.009 (0.005)	0.021 (0.001)	0.020 (0.001)
Dominant	Rare	1.3	0.938 (0.008)	0.930 (0.008)	0.034 (0.006)	0.030 (0.005)	0.011 (0.008)	0.007 (0.008)	0.062 (0.003)	0.058 (0.003)
		1.5	0.850 (0.011)	0.820 (0.012)	0.020 (0.006)	0.025 (0.007)	0.024 (0.007)	0.008 (0.008)	0.054 (0.003)	0.048 (0.002)
		1.7	0.678 (0.015)	0.656 (0.015)	0.028 (0.005)	0.025 (0.005)	0.016 (0.007)	0.003 (0.007)	0.052 (0.002)	0.047 (0.002)
		2.0	0.332 (0.015)	0.292 (0.014)	0.036 (0.006)	0.035 (0.006)	0.013 (0.007)	-0.005 (0.007)	0.053 (0.003)	0.048 (0.002)
	Common	1.3	0.888 (0.010)	0.857 (0.011)	0.026 (0.005)	0.025 (0.005)	-0.004 (0.006)	-0.01 (0.005)	0.032 (0.001)	0.028 (0.001)
		1.5	0.670 (0.011)	0.600 (0.011)	0.045 (0.005)	0.030 (0.006)	-0.011 (0.007)	-0.022 (0.006)	0.028 (0.001)	0.023 (0.001)
		1.7	0.328 (0.015)	0.262 (0.014)	0.024 (0.005)	0.023 (0.005)	0.012 (0.006)	-0.002 (0.005)	0.031 (0.001)	0.026 (0.001)
		2.0	0.080 (0.009)	0.052 (0.007)	0.036 (0.007)	0.040 (0.007)	0.016 (0.006)	-0.006 (0.005)	0.031 (0.002)	0.025 (0.001)
Recessive	Rare	2.0	1.000 (0.000)	0.920 (0.019)	0.010 (0.007)	0.020 (0.010)	-0.256 (0.056)	-0.328 (0.106)	0.686 (0.066)	0.537 (1.128)
		2.5	1.000 (0.000)	0.865 (0.024)	0.010 (0.007)	0.005 (0.005)	-0.270 (0.053)	-0.333 (0.094)	0.622 (0.062)	0.868 (1.019)
		3.0	1.000 (0.000)	0.800 (0.030)	0.010 (0.007)	0.015 (0.009)	-0.240 (0.054)	-0.218 (0.045)	0.641 (0.071)	0.457 (0.049)
		3.5	1.000 (0.000)	0.730 (0.032)	0.010 (0.007)	0.010 (0.007)	-0.386 (0.053)	-0.271 (0.041)	0.711 (0.075)	0.404 (0.046)
	Common	2.0	0.945 (0.016)	0.780 (0.033)	0.000 (0.000)	0.025 (0.011)	0.016 (0.029)	-0.080 (0.021)	0.165 (0.017)	0.092 (0.010)
		2.5	0.755 (0.032)	0.395 (0.034)	0.000 (0.000)	0.015 (0.009)	0.096 (0.037)	-0.086 (0.021)	0.281 (0.033)	0.095 (0.010)
		3.0	0.485 (0.035)	0.105 (0.022)	0.010 (0.007)	0.015 (0.009)	0.083 (0.028)	-0.076 (0.017)	0.163 (0.015)	0.065 (0.008)
		3.5	0.340 (0.034)	0.055 (0.016)	0.015 (0.009)	0.010 (0.007)	0.061 (0.034)	-0.122 (0.017)	0.229 (0.034)	0.073 (0.008)

^aCorrect analysis means the genetic model adopted in the analysis is the same as the underlying genetic mode.

^bFalse negatives are counts out of 1 and false positives are counts out of 6.

^cMode refers to underlying genetic mode of data; Freq refers to frequency of risk haplotype; OR refers to the odds ratio of the risk haplotype.

^dMeasures in this table are found by averaging over 1,000 simulated data sets; standard errors (SE) are shown in parentheses.

^eBolded measures (pro vs. retro) are statistically significantly different at 0.05 level (i.e. the $\pm 2SE$ intervals do NOT overlap).

of a retrospective likelihood and the correct genetic mode. In practice, a major concern about using retrospective likelihoods is that they are difficult to implement. When using penalized likelihood methods, optimizing the retrospective likelihood can become even more intractable when the penalty term is incorporated. Use of the exact likelihood coupled with the EM algorithm due to the unknown phase, with the penalization, and tuning on a dense grid of tuning parameters, requires a computationally intensive iterative procedure.

To overcome the computational burden, the least squares objective function provides a promising alternative for implementing penalized retrospective methods. By using the least squares objective function, the need to directly optimize the penalized retrospective likelihood is bypassed. Instead, the unpenalized likelihood is optimized once for the centering value in the objective function;

hence, implementing prospective and retrospective penalized methods have similar computational costs. The spared computational efforts can be put into exploring and identifying the correct genetic mode for potential risk haplotypes.

Penalized likelihood methods can have higher power than unpenalized methods in detecting important haplotypic factors [Guo and Lin, 2009]. Our simulations agree with their findings and reveal that the methods considered here can have better power to identify risk haplotypes (i.e. fewer FNs) than the unpenalized version (e.g. Tables XIV and XV vs. Tables IV and V, respectively). However, the power enjoyed by penalized likelihood methods comes at the expense of obtaining effect estimates with relatively larger bias than their unpenalized likelihood counterparts. As observed in our simulations, the bias on the effect estimates obtained by the penalized method can remain

TABLE XV. Unpenalized results of Sim I (one-risk haplotype) for Distribution 1—additive analysis^a

Mode ^c	Freq	OR	Model selection results				Parameter estimation results			
			False negatives ^b		False positives		Bias		MSE	
			Pro	Retro	Pro	Retro	Pro	Retro	Pro	Retro
Dominant	Rare	1.3	0.954 ^d (0.007)	0.953 (0.007)	0.036 (0.006)	0.035 (0.006)	-0.008 (0.007)	-0.014 (0.007)	0.055 (0.002)	0.052 (0.002)
		1.5	0.850 (0.010)	0.855 (0.009)	0.025 (0.006)	0.030 (0.006)	-0.013 (0.008)	-0.023 (0.007)	0.048 (0.002)	0.046 (0.003)
		1.7	0.700 (0.015)	0.687 (0.015)	0.029 (0.005)	0.032 (0.006)	-0.019 (0.007)	-0.036 (0.007)	0.048 (0.002)	0.046 (0.002)
		2.0	0.367 (0.015)	0.361 (0.015)	0.031 (0.006)	0.033 (0.006)	-0.032 (0.007)	-0.059 (0.007)	0.049 (0.002)	0.047 (0.002)
	Common	1.3	0.904 (0.009)	0.908 (0.009)	0.022 (0.005)	0.022 (0.005)	-0.047 (0.005)	-0.052 (0.005)	0.028 (0.001)	0.027 (0.001)
		1.5	0.735 (0.010)	0.750 (0.009)	0.035 (0.005)	0.030 (0.006)	-0.074 (0.006)	-0.088 (0.005)	0.028 (0.001)	0.028 (0.002)
		1.7	0.386 (0.015)	0.414 (0.016)	0.026 (0.006)	0.031 (0.006)	-0.076^e (0.005)	-0.098 (0.005)	0.031 (0.001)	0.032 (0.001)
		2.0	0.120 (0.010)	0.135 (0.011)	0.038 (0.007)	0.037 (0.006)	-0.096 (0.005)	-0.134 (0.005)	0.035 (0.002)	0.040 (0.002)
Recessive	Rare	2.0	0.990 (0.007)	0.990 (0.007)	0.040 (0.014)	0.050 (0.015)	-0.671 (0.015)	-0.665 (0.016)	0.497 (0.021)	0.491 (0.021)
		2.5	0.995 (0.005)	0.990 (0.007)	0.020 (0.010)	0.030 (0.012)	-0.871 (0.016)	-0.865 (0.017)	0.810 (0.028)	0.805 (0.028)
		3.0	0.970 (0.012)	0.965 (0.013)	0.040 (0.014)	0.040 (0.014)	-0.995 (0.017)	-0.989 (0.018)	1.045 (0.034)	1.039 (0.036)
		3.5	1.000 (0.000)	0.985 (-0.009)	0.010 (-0.007)	0.015 (-0.011)	-1.148 (-0.014)	-1.14 (-0.015)	1.357 (-0.032)	1.345 (-0.035)
		2.0	0.985 (0.009)	0.980 (0.010)	0.010 (0.007)	0.030 (0.012)	-0.560 (0.010)	-0.557 (0.011)	0.335 (0.012)	0.332 (0.013)
	Common	2.5	0.910 (0.020)	0.890 (0.022)	0.015 (0.009)	0.025 (0.011)	-0.722 (0.010)	-0.708 (0.011)	0.542 (0.015)	0.526 (0.016)
		3.0	0.865 (0.024)	0.835 (0.026)	0.025 (0.011)	0.040 (0.014)	-0.867 (0.010)	-0.845 (0.011)	0.771 (0.018)	0.739 (0.020)
		3.5	0.775 (0.030)	0.725 (0.032)	0.025 (0.011)	0.025 (0.013)	-0.979 (0.010)	-0.947 (0.011)	0.978 (0.020)	0.920 (0.021)

^aAdditive analysis means an additive model is adopted in the analysis when the underlying genetic mode is not.

^bFalse negatives are counts out of 1 and false positives are counts out of 6.

^cMode refers to underlying genetic mode of data; Freq refers to frequency of risk haplotype; OR refers to the odds ratio of the risk haplotype.

^dMeasures in this table are found by averaging over 1,000 simulated data sets; standard errors (SE) are shown in parentheses.

^eBolded measures (pro vs. retro) are statistically significantly different at 0.05 level (i.e. the ±2SE intervals do NOT overlap).

sizeable, particularly when there is very little power to detect the effect. Part of this is due to the fact that non-detection of the effect implies that the estimated

coefficient is zero. In contrast, for the MLE, the estimate is still non-zero even though it is not significantly different from zero.

TABLE XVI. AIC-penalized results of Sim I (one-risk haplotype) for Distribution 1—correct analysis^a

Mode ^c	Freq	OR	Model selection results				Parameter estimation results			
			False negatives ^b		False positives		Bias		MSE	
			Pro	Retro	Pro	Retro	Pro	Retro	Pro	Retro
Additive	Rare	1.3	0.605 ^d (0.015)	0.600 (0.015)	1.010 (0.089)	1.090 (0.094)	-0.110 (0.015)	-0.106 (0.015)	0.059 (0.003)	0.058 (0.003)
		1.5	0.320 (0.015)	0.330 (0.015)	1.095 (0.088)	1.145 (0.089)	-0.106 (0.018)	-0.109 (0.018)	0.074 (0.005)	0.074 (0.005)
		1.7	0.100 (0.009)	0.105 (0.010)	0.975 (0.078)	0.990 (0.077)	-0.079 (0.016)	-0.082 (0.016)	0.060 (0.006)	0.060 (0.006)
		2.0	0.025 (0.005)	0.025 (0.005)	1.000 (0.082)	1.000 (0.081)	-0.070 (0.016)	-0.067 (0.016)	0.055 (0.006)	0.055 (0.006)
	Common	1.3	0.350 (0.015)	0.325 (0.015)	1.010 (0.082)	1.050 (0.082)	-0.077 (0.011)	-0.073 (0.011)	0.032 (0.002)	0.030 (0.002)
		1.5	0.080 (0.009)	0.070 (0.008)	1.125 (0.085)	1.130 (0.082)	-0.060 (0.011)	-0.062 (0.011)	0.030 (0.003)	0.030 (0.003)
		1.7	0.015 (0.004)	0.005 (0.002)	1.060 (0.084)	1.035 (0.083)	-0.016 (0.010)	-0.018 (0.010)	0.022 (0.003)	0.020 (0.002)
		2.0	0.000 (0.000)	0.000 (0.000)	0.910 (0.077)	0.900 (0.076)	-0.041 (0.010)	-0.046 (0.010)	0.022 (0.002)	0.021 (0.002)
Dominant	Rare	1.3	0.630 (0.015)	0.625 (0.015)	0.885 (0.082)	0.940 (0.086)	-0.116 (0.015)	-0.117 (0.015)	0.057 (0.003)	0.057 (0.003)
		1.5	0.355 (0.015)	0.345 (0.015)	0.855 (0.080)	0.855 (0.079)	-0.128 (0.018)	-0.122 (0.018)	0.078 (0.006)	0.076 (0.005)
		1.7	0.175 (0.012)	0.170 (0.012)	0.905 (0.072)	0.945 (0.079)	-0.102 (0.020)	-0.109 (0.019)	0.088 (0.008)	0.083 (0.007)
		2.0	0.025 (0.005)	0.020 (0.004)	1.080 (0.083)	1.075 (0.086)	-0.070 (0.017)	-0.070 (0.016)	0.064 (0.007)	0.058 (0.006)
	Common	1.3	0.435 (0.016)	0.425 (0.016)	0.950 (0.078)	0.875 (0.075)	-0.087 (0.012)	-0.097 (0.012)	0.038 (0.002)	0.036 (0.002)
		1.5	0.105 (0.010)	0.115 (0.010)	1.045 (0.080)	1.065 (0.081)	-0.085 (0.013)	-0.089 (0.012)	0.040 (0.004)	0.038 (0.004)
		1.7	0.015 (0.004)	0.015 (0.004)	1.010 (0.084)	0.940 (0.079)	-0.033 (0.013)	-0.037 (0.012)	0.034 (0.004)	0.029 (0.003)
		2.0	0.000 (0.000)	0.000 (0.000)	0.935 (0.083)	0.915 (0.088)	-0.018 (0.011)	-0.027 (0.011)	0.026 (0.002)	0.024 (0.002)
Recessive	Rare	2.0	0.900^e (0.009)	0.570 (0.016)	0.450 (0.049)	0.545 (0.057)	-0.563 (0.031)	-0.233 (0.041)	0.506 (0.015)	0.383 (0.017)
		2.5	0.835 (0.012)	0.505 (0.016)	0.515 (0.049)	0.565 (0.052)	-0.717 (0.036)	-0.386 (0.042)	0.771 (0.031)	0.507 (0.027)
		3.0	0.765 (0.013)	0.360 (0.015)	0.500 (0.052)	0.480 (0.048)	-0.779 (0.042)	-0.381 (0.043)	0.965 (0.032)	0.511 (0.038)
		3.5	0.670 (0.015)	0.330 (0.015)	0.500 (0.050)	0.565 (0.056)	-0.804 (0.048)	-0.403 (0.048)	1.101 (0.048)	0.617 (0.049)
	Common	2.0	0.435 (0.016)	0.270 (0.014)	0.380 (0.043)	0.450 (0.047)	-0.202 (0.035)	-0.176 (0.026)	0.290 (0.019)	0.162 (0.014)
		2.5	0.180 (0.012)	0.060 (0.008)	0.425 (0.047)	0.430 (0.048)	-0.130 (0.035)	-0.121 (0.023)	0.259 (0.023)	0.117 (0.014)
		3.0	0.095 (0.009)	0.025 (0.005)	0.505 (0.051)	0.435 (0.047)	-0.135 (0.036)	-0.164 (0.020)	0.273 (0.026)	0.107 (0.015)
		3.5	0.020 (0.004)	0.000 (0.000)	0.375 (0.041)	0.400 (0.041)	-0.056 (0.031)	-0.101 (0.018)	0.193 (0.024)	0.072 (0.008)

^aCorrect analysis means the genetic model adopted in the analysis is the same as the underlying genetic mode.

^bFalse negatives are counts out of 1 and false positives are counts out of 6.

^cMode refers to underlying genetic mode of data; Freq refers to frequency of risk haplotype; OR refers to the odds ratio of the risk haplotype.

^dMeasures in this table are found by averaging over 1,000 simulated data sets; standard errors (SE) are shown in parentheses.

^eBolded measures (pro vs. retro) are statistically significantly different at 0.05 level (i.e. the ±2SE intervals do NOT overlap).

The power and bias is also affected by the model selection criterion used to select the tuning parameter. In our analysis, we used BIC to choose the final model because it can achieve consistent model selection [Shao, 1997; Yang, 2005], and although estimation accuracy was of interest, our primary goal was to identify the true model structure. To achieve selection consistency, BIC penalizes degrees of freedom more heavily, which can place a larger amount of shrinkage on the effect estimates and increase their bias. Alternatively, AIC could be used to select the final model. (See Tables XVI and XVII for the AIC-based results for Tables IV and V, respectively.) This selection criterion targets prediction error rather than finding the true model structure [Shao, 1997; Yang, 2005] and imposes a smaller penalty on the degrees of freedom. As a result, models selected via AIC will incur less shrinkage on the effect estimates, which can decrease their bias, but also

increases the chance of including non-important predictors in the final model. In our analyses, when tuning with AIC instead of BIC, the FN count stays relatively similar, while the FP count often increases and the bias/MSE decreases. This is a direct result of placing less shrinkage/penalty on the final estimates.

Throughout this paper, we assume that the sample size is larger than the number of predictors ($n > p$), which is consistent with typical case-control data sets. However, if the investigator wishes to consider the situation in which the number of candidate haplotypes exceeds the sample size ($n < p$), the analysis can be extended by replacing the initial estimates β_j in the adaptive LASSO penalty with an L_2 penalized solution, i.e. ridge regression. The initial ridge step would then need to be tuned; we suggest the use of AIC, so that less shrinkage occurs in the initial estimate. However, the use of the least squares objective function

TABLE XVII. AIC-penalized results of Sim I (one-risk haplotype) for Distribution 1—additive analysis^a

Mode ^c	Freq	OR	Model selection results				Parameter estimation results					
			False negatives ^b		False positives		Bias		MSE			
			Pro	Retro	Pro	Retro	Pro	Retro	Pro	Retro		
Dominant	Rare	1.3	0.640 ^d (0.015)	0.635 (0.015)	0.870 (0.078)	0.860 (0.079)	-0.128 (0.014)	-0.132 (0.013)	0.055 (0.003)	0.053 (0.003)		
		1.5	0.385 (0.015)	0.395 (0.015)	0.815 (0.081)	0.840 (0.077)	-0.145 (0.017)	-0.154 (0.017)	0.081 (0.005)	0.081 (0.005)		
		1.7	0.215 (0.013)	0.210 (0.013)	0.895 (0.069)	0.920 (0.072)	-0.141 (0.019)	-0.155 (0.019)	0.092 (0.008)	0.092 (0.008)		
		2	0.030 (0.005)	0.035 (0.006)	1.115 (0.082)	1.095 (0.082)	-0.106 (0.017)	-0.132 (0.016)	0.065 (0.007)	0.067 (0.007)		
	Common	1.3	0.460 (0.016)	0.470 (0.016)	0.840 (0.072)	0.815 (0.071)	-0.121 (0.011)	-0.128 (0.010)	0.037 (0.002)	0.038 (0.002)		
		1.5	0.170 (0.012)	0.170 (0.012)	1.120 (0.084)	1.100 (0.085)	-0.142 (0.012)	-0.153 (0.011)	0.049 (0.004)	0.049 (0.004)		
		1.7	0.025 (0.005)	0.025 (0.005)	0.980 (0.081)	1.040 (0.088)	-0.116 (0.011)	-0.134 (0.011)	0.039 (0.004)	0.041 (0.004)		
		2	0.005^e (0.002)	0.000 (0.000)	0.885 (0.083)	1.015 (0.092)	-0.118 (0.011)	-0.153 (0.010)	0.037 (0.004)	0.043 (0.004)		
		Recessive	Rare	2	0.795 (0.013)	0.800 (0.013)	0.405 (0.047)	0.775 (0.072)	-0.653 (0.011)	-0.657 (0.010)	0.449 (0.013)	0.453 (0.013)
				2.5	0.845 (0.011)	0.840 (0.012)	0.455 (0.047)	0.765 (0.069)	-0.883 (0.009)	-0.887 (0.009)	0.796 (0.015)	0.803 (0.015)
3	0.805 (0.013)			0.780 (0.013)	0.555 (0.056)	0.930 (0.085)	-1.047 (0.011)	-1.041 (0.011)	1.121 (0.022)	1.108 (0.021)		
Common	3.5	0.780 (0.013)	0.785 (0.013)	0.455 (0.050)	0.815 (0.080)	-1.198 (0.012)	-1.195 (0.012)	1.462 (0.026)	1.456 (0.026)			
	2	0.695 (0.015)	0.700 (0.014)	0.480 (0.050)	0.900 (0.078)	-0.623 (0.009)	-0.622 (0.009)	0.405 (0.011)	0.405 (0.010)			
	2.5	0.500 (0.016)	0.460 (0.016)	0.540 (0.053)	1.050 (0.082)	-0.783 (0.011)	-0.776 (0.011)	0.635 (0.016)	0.624 (0.016)			
	3	0.430 (0.016)	0.420 (0.016)	0.535 (0.049)	1.020 (0.085)	-0.944 (0.011)	-0.933 (0.012)	0.915 (0.020)	0.897 (0.021)			
	3.5	0.240 (0.014)	0.210 (0.013)	0.590 (0.052)	1.070 (0.080)	-1.036 (0.011)	-1.015 (0.011)	1.095 (0.022)	1.055 (0.023)			

^aAdditive analysis means an additive model is adopted in the analysis when the underlying genetic mode is not.

^bFalse negatives are counts out of 1 and false positives are counts out of 6.

^cMode refers to underlying genetic mode of data; Freq refers to frequency of risk haplotype; OR refers to the odds ratio of the risk haplotype.

^dMeasures in this table are found by averaging over 1,000 simulated data sets; standard errors (SE) are shown in parentheses.

^eBolded measures (pro vs. retro) are statistically significantly different at 0.05 level (i.e. the $\pm 2SE$ intervals do NOT overlap).

would no longer be possible, as the maximum likelihood estimator is no longer necessarily unique.

ACKNOWLEDGMENTS

The authors appreciate the GAW15 RA data set from the Genetic Analysis Workshop. Support for GAW15 data was provided from NIH grants 5R01-HL049609-14, 1R01-AG021917-01A1, the University of Minnesota, and the Minnesota Supercomputing Institute, and the GAW grant, R01-GM031575 and AR44422. The authors thank the reviewers for the helpful and constructive feedback that greatly improved the manuscript. They also thank Dr. Danyu Lin, Tammy Bailey, and Chris Smith for their generous help with HAPSTAT and perl. M.L.K. was supported by NIH T32GM081057 and NIH R01 MH084022-01. H.D.B. was supported by NSF DMS-0705968, NIH R01 MH084022-01, and NIH 1P01-CA142538-01. J.Y.T. was supported by NIH R01 MH084022-01 and NIH 1P01-CA142538-01.

REFERENCES

- Arlot S, Celisse A. 2010. A survey of cross-validation procedures for model selection. *Stat Surv* 4:40–79.
- Balding DJ. 2006. A tutorial on statistical methods for population association studies. *Nat Rev Genet* 7:781–791.
- Barrett JC, Fry B, Maller J, Daly MJ. 2005. Haploview: analysis and visualization of LD and haplotype maps. *Bioinformatics* 21: 263–265.
- Browning S, Thomas J. 2007. Multilocus analysis of GAW15 NARAC chromosome 18 case-control data. *BMC Proc* 1:S1–S11.
- Chen YH, Kao JT. 2006. Multinomial logistic regression approach to haplotype association analysis in population-based case-control studies. *BMC Genet* 7:43.
- Chen YH, Chatterjee N, Carroll RJ. 2009. Shrinkage estimators for robust and efficient inference in haplotype-based case-control studies. *J Am Stat Assoc* 104:220–233.
- Clark AG. 2004. The role of haplotypes in candidate-gene studies. *Genet Epidemiol* 27:321–333.
- Cordell HJ. 2006. Estimation and testing of genotype and haplotype effects in case-control studies: comparison of weighted regression and multiple imputation procedures. *Genet Epidemiol* 30:259–275.
- de Bakker PW, Yelensky R, Pe'er I, Gabriel S, Daly MJ, Altshuler D. 2005. Efficiency and power in genetic association studies. *Nat Genet* 37:1217–1223.
- Epstein MP, Satten GA. 2003. Inference on haplotype effects in case-control studies using unphased genotype data. *Am J Hum Genet* 73:1316–1329.
- French B, Lumley T, Monks SA, Rice KM, Hindorff LA, Reiner AP, Psaty BM. 2006. Simple estimates of haplotype relative risks in case-control data. *Genet Epidemiol* 30:485–494.
- Guo W, Lin S. 2009. Generalized linear modeling with regularization for detecting common disease rare haplotype association. *Genet Epidemiol* 33:308–316.
- Hastie T, Tibshirani R, Friedman J. 2009. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd edition. Berlin: Springer.
- Lake SL, Lyon H, Tantisira K, Silverman EK, Weiss ST, Laird NM, Schaid DJ. 2003. Estimation and tests of haplotype-environment interaction when linkage phase is ambiguous. *Hum Hered* 55: 56–65.
- Li Y, Sung WK, Liu JJ. 2007. Association mapping via regularized regression analysis of single nucleotide polymorphism haplotypes in variable-sized sliding windows. *Am J Hum Genet* 80: 705–715.
- Lin DY, Huang BE. 2008. The use of inferred haplotypes in downstream analysis. *Am J Hum Genet* 80:577–579.
- Lin DY, Zeng D. 2006. Likelihood-based inference on haplotype effects in genetic association studies. *J Am Stat Assoc* 101:89–104.
- Lin DY, Zeng D, Milikan R. 2005. Maximum likelihood estimation of haplotype effects and haplotype-environment interactions in association studies. *Genet Epidemiol* 29:299–312.
- Prentice RL, Pyke R. 1979. Logistic disease incidence models and case-control studies. *Biometrika* 66:403–412.
- Satten GA, Epstein MP. 2004. Comparison of prospective and retrospective methods for haplotype inference in case-control studies. *Genet Epidemiol* 27:192–201.
- Schaid DJ. 2004. Evaluating associations of haplotypes with traits. *Genet Epidemiol* 27:348–364.
- Shao J. 1997. An asymptotic theory for linear model selection. *Stat Sin* 7:221–264.
- Souverein OW, Zwiderman AH, Tanck MW. 2006. Estimating haplotype effects on dichotomous outcome for unphased genotype data using a weighted penalized log-likelihood approach. *Hum Hered* 61:104–110.
- Souverein OW, Zwiderman AH, Jukema JW, Tanck MW. 2008. Estimating effects of rare haplotypes on failure times using a penalized Cox proportional hazards regression model. *BMC Genet* 9:9.
- Stram DO, Leigh Pearce C, Bretsky P, Freedman M, Hirschhorn JN, Altshuler D, Kolonel LN, Henderson BE, Thomas DC. 2003. Modeling and E-M estimation of haplotype specific relative risks from genotype data for a case-control study of unrelated individuals. *Hum Hered* 55:179–190.
- Tanck MW, Klerkx AH, Jukema JW, De Knijff P, Kastelein JJ, Zwiderman AH. 2003. Estimation of multilocus haplotype effects using weighted penalized log-likelihood: analysis of five sequence variations at the cholesteryl ester transfer protein gene locus. *Ann Hum Genet* 67:175–184.
- Tibshirani R. 1996. Regression shrinkage and selection via the LASSO. *J R Stat Soc* 58:267–288.
- Tzeng JY, Bondell HD. 2010. A comprehensive approach to haplotype-specific analysis by penalized likelihood. *Eur J Hum Genet* 18: 95–103.
- Tzeng JY, Lu W, Farmen MW, Liu Y, Sullivan PF. 2010. Haplotype-based pharmacogenetic analysis for longitudinal quantitative traits in the presence of dropout. *J Biopharm Stat* 20:334–350.
- Wang H, Leng C. 2007. Unified LASSO estimation by least squares approximation. *J Am Stat Assoc* 102:1039–1048.
- Yang Y. 2005. Can the strengths of AIC and BIC be shared? A conflict between model identification and regression estimation. *Biometrika* 92:937–950.
- Zaitlen M, Kang H, Eskin E, Halperin E. 2007. Leveraging the HapMap correlation structure in association studies. *Am J Hum Genet* 80: 683–691.
- Zou H. 2006. The adaptive LASSO and its oracle properties. *J Am Stat Assoc* 101:1418–1429.