

A mixed-model approach for powerful testing of genetic associations with cancer risk incorporating tumor characteristics

Haoyu Zhang,^{1,2} Ni Zhao,¹ Thomas U. Ahearn,² William Wheeler,³ Montserrat
García-Closas,² and Nilanjan Chatterjee^{1,4}

¹*Department of Biostatistics Johns Hopkins Bloomberg SPH, Baltimore, MD 21205,
U.S.A.*

²*National Cancer Institute, Division of Cancer Epidemiology and Genetics,
Rockville, MD 20850, U.S.A.*

³*Information Management Services, Inc., Rockville, MD 20850,
USA*

⁴*Department of Oncology, School of Medicine, Johns Hopkins University, Baltimore,
MD 21205, U.S.A.*

(Dated: 17 October 2018)

ABSTRACT: Cancers are routinely classified into subtypes according to various features, including histo-pathological characteristics and molecular markers. Previous investigations of genetic loci have reported heterogeneous association between loci and cancer subtypes. However, it is not evident what is the optimal modeling strategy for handling correlated tumor features, missing data, and increased degrees-of-freedom in the underlying tests of associations. We propose a score test for genetic associations using a mixed-effect two-stage polytomous model (MTOP). In the first stage, a standard polytomous model is used to specify for all possible subtypes defined by the cross-classification of different markers. In the second stage, the subtype-specific case-control odds ratios are specified using a more parsimonious model based on the case-control odds ratio for a baseline subtype, and the case-case parameters associated with tumor markers. Further, to reduce the degrees-of-freedom, we specify case-case parameters for additional markers using a random-effect model. We use the EM algorithm to account for missing data on tumor markers. The score-test distribution theory is developed by borrowing analogous techniques from group-based association tests. Through analysis of simulations across a wide range of realistic scenarios and data from the Polish Breast Cancer Study (PBCS), we show MTOP substantially outperform several alternative methods for identifying heterogeneous associations between risk loci and tumor subtypes.

KEY WORDS: Two-stage polytomous model; Susceptibility variants; Cancer subtypes; EM algorithm; Score tests; Etiologic heterogeneity.

I. INTRODUCTION

Genome-wide association studies (GWAS) have identified hundreds of single nucleotide polymorphisms (SNPs) associated with various cancers (MacArthur *et al.*, 2016; Visscher *et al.*, 2017); However, many cancer GWAS have often defined cancer endpoints according to specific anatomic sites, and not according to subtypes of the disease. Many cancers consist of etiologically and clinically heterogeneous subtypes that are defined by multiple correlated tumor characteristics, for instance, breast cancer is routinely classified into subtypes defined by tumor expression of estrogen receptor (ER), progesterone receptor (PR), and human epidermal growth factor receptor 2 (HER2) (Curigliano *et al.*, 2017; Perou *et al.*, 2000; Prat *et al.*, 2015).

Increasing number of epidemiologic studies with tumor specimens are allowing the characteristics of cancers at the histological and molecular levels (Cancer Genome Atlas, 2012; Cancer Genome Atlas Research, 2012; 2014). This provides tremendous opportunities to characterize distinct etiological pathways within cancer subtypes. For example, a breast cancer ER-negative specific GWAS reported 20 SNPs that are more strongly associated with risk of developing ER-negative than ER-positive disease (Milne *et al.*, 2017). Previous studies also suggest traditional breast cancer risk factors, such as age, obesity, and hormone therapy use, may also be heterogeneously associated with breast cancer subtypes (Barnard *et al.*, 2015). However, there are complexities when using subtype information to identify distinct risk factor associations, such as missing tumor marker data, the correlation between tumor markers, and the large dimensionality of subtypes.

Polytomous logistic regression is a common approach for analyzing cancer data with information on multiple tumor characteristics (Dubin and Pasternack, 1986; Gortmaker *et al.*, 1994). This method estimates the odds ratio of each cancer subtype compared to the control group, i.e., people without the disease. A major limitation of this approach is that it loses power due to the increased degrees of freedom when there are many different cancer subtypes. A two-stage polytomous logistic regression has been proposed to characterize subtype heterogeneity of a disease using the underlying disease characteristics (Chatterjee, 2004). The first stage of this method uses the polytomous logistic regression to model each subtype specific case-control odds ratio. In the second stage, the subtype-specific case-control odds ratios are decomposed to the case-control odds ratio of a reference subtype, case-case odds ratio of tumor characteristic and higher order interactions between the case-case odds ratio of the tumor characteristics. The two-stage model can reduce the degrees of freedom due to the estimation of subtype specific odds ratio. Moreover, the second stage parameters can be interpreted as the case-case parameters for tumor characteristics.

Although, selected applications have demonstrated the power of the two-stage regression method (Falk *et al.*, 2014; Peters *et al.*, 2004; Sherman *et al.*, 2007; Zabor and Begg, 2017), for several reasons the method has not been widely applied to analyze data on multiple tumor characteristics. First, tumor characteristic data in epidemiologic studies are often incomplete. Second, the two-stage model estimation algorithm places high demands on computing power, and is therefore not readily applicable to large datasets. Finally, as the number of tumor characteristics increases, the method can have substantial power loss due to the increase in the degrees of freedom.

In this paper, we propose a series of computational and statistical innovations to adapt the two-stage model for large scale hypothesis testing in GWAS. We first briefly review the two-stage polytomous model in Section II A. Then in Section II B, we propose to use the two-stage model to test alternative forms of hypotheses for genetic associations in the presence of heterogeneity. And in Section II C, we propose an Expectation-Maximization (EM) algorithm (Dempster *et al.*, 1977) within the two-stage model framework to account for the missing tumor characteristics. In Section II D, we develop a computationally scalable score test for fixed-effect two-stage model, and in Section II E, we introduce a mixed-effect two-stage model to handle potentially large number of exploratory tumor markers minimizing loss of power. We study the type one error and power of the proposed methods on simulated data in Section III. Moreover, we illustrate the methods with two applications using the Polish Breast Cancer Study (PBCS) data in Section IV. Finally we discuss the strengths and limitations of the methods, and future research directions in Section V. The proposed methods are available in a high speed R package called TOP (<https://github.com/andrewhaoyu/TOP>), with all the core functions implemented in C code.

II. METHOD

A. Two-stage polytomous model

Following (Chatterjee, 2004), we first briefly introduce the two-stage model for tumor heterogeneity. Suppose a disease can be classified using K disease characteristics. Assuming each characteristic k can be classified into M_k categories, then the disease can be classified

into $M \equiv M_1 \times M_2 \cdots \times M_K$ subtypes. For example, breast cancer can be classified into eight subtypes by three tumor characteristics (ER, PR, and HER2), each of which is either positive or negative. Note, that we will use this breast cancer example to demonstrate the methods throughout the methods section. Let D_i denote the disease status, taking values in $\{0, 1, 2, \dots, M\}$, of the i th ($i \in 1, \dots, N$) subject in the study. $D_i = 0$ represents a control, and $D_i = m$ represent a subject with disease of subtype m . Let G_i be the genotype for i th subject and \mathbf{X}_i be a $P \times 1$ vector of other covariates we want to adjust for in the model, where P is the total number of other covariates. In the first-stage model, we use the standard “saturated” polytomous logistic regression model

$$Pr(D_i = m | G_i, \mathbf{X}_i) = \frac{\exp(\beta_m G_i + \mathbf{X}_i^T \boldsymbol{\eta}_m)}{1 + \sum_{m=1}^M \exp(\beta_m G_i + \mathbf{X}_i^T \boldsymbol{\eta}_m)}, \quad m \in \{1, 2, \dots, M\}, \quad (1)$$

where β_m and $\boldsymbol{\eta}_m$ are the regression coefficients for the SNP and other covariates for association with the m th subtype.

Because each cancer subtype m is defined through a unique combination of the K characteristics, we can always alternatively index the parameters β_m as $\{\beta_{s_1 s_2 \dots s_K}\}$, where $s_k \in \{0, 1\}$ for binary tumor characteristics, and $s_k \in \{t_1 \leq t_2 \leq \dots \leq t_{M_k}\}$ for ordinal tumor characteristics with t_1, \dots, t_{M_k} as a set of ordinal scores for M_k different levels. Under the same breast cancer example, originally β_1 could be the coefficient of cancer subtype ER-PR-HER2-. With the new index, β_1 could be written as β_{000} , which means the three tumor characteristics are all negative. With this new index, we can represent the log odds ratio as

$$\beta_{s_1 s_2 \dots s_K} = \theta^{(0)} + \sum_{k_1=1}^K \theta_{k_1}^{(1)} s_{k_1} + \sum_{k_1=1}^K \sum_{k_2 > k_1}^K \theta_{k_1 k_2}^{(2)} (s_{k_1} s_{k_2}) + \dots + \theta_{12 \dots K}^{(K)} (s_1 s_2 \dots s_K). \quad (2)$$

Here $\theta^{(0)}$ represents the standard case-control log odds ratio for a reference disease subtype compared to the control and $\theta_{k_1}^{(1)}$ represents a case-case log odds ratio associated with the levels of k_1 th tumor characteristics after adjusting for other tumor characteristics, and $\theta_{k_1 k_2}^{(2)}$ represent case-case log odds ratios associated with pairwise interactions among the tumor characteristics and so on.

We can represent the Equation 2 into matrix form as

$$\boldsymbol{\beta} = \mathbf{Z}_G \boldsymbol{\theta} = \mathbf{Z}_G \begin{bmatrix} \theta^{(0)} & \boldsymbol{\theta}_H^T \end{bmatrix}^T. \quad (3)$$

Here $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_M)^T$ is a vector of first stage case-control log odds ratios for all the M subtypes, \mathbf{Z}_G is the second stage design matrix, and $\boldsymbol{\theta} = (\theta^{(0)}, \boldsymbol{\theta}_H^T)^T$ is the vector of second stage parameters including the case-control log odds ratio for the reference subtype $\theta^{(0)}$ and all the case-case log odds ratios $\boldsymbol{\theta}_H$. This second stage design matrix connects the first stage case-control log odds ratios for all subtypes to the second stage case-control and case-case log odds ratios. We can build models specifying different second stage matrix by constraining different case-case parameters to be zero in a hierarchical manner.

Up to now, we have only described second stage decomposition for the regression coefficients of \mathbf{G} . We could also apply second stage decomposition on the other covariates. The details of this could be found in Supplementary Section 1. We don't perform any second stage decomposition on regression coefficients of intercepts, since making assumption on the prevalence of different cancer subtypes could potentially yield bias. Moving forward, we use \mathbf{Z}_X to denote the second stage design matrix for the other covariates \mathbf{X} , $\boldsymbol{\lambda}$ to denote the second stage parameters for \mathbf{X} , and \mathbf{Z} to denote the second stage design matrix for all the covariates.

B. Hypothesis test under two-stage model

We can decompose the first stage case-control log odds ratios of all the subtypes into the second stage case-control log odds ratio of a reference subtype and case-case log odds ratios of tumor characteristics through Equation 3. This decomposition presents multiple options for comprehensively testing the association between a SNP and disease subtypes. The first hypothesis test is the global association test,

$$H_0^A : \boldsymbol{\theta} = \begin{bmatrix} \theta^{(0)} & \boldsymbol{\theta}_H^T \end{bmatrix}^T = \begin{bmatrix} 0 & \mathbf{0}^T \end{bmatrix}^T \text{ versus } H_1^A : \boldsymbol{\theta} \neq \mathbf{0}. \quad (4)$$

This tests for an overall association between the SNP and the disease. Because $\boldsymbol{\theta} = \mathbf{0}$ implies $\boldsymbol{\beta} = \mathbf{0}$, rejecting this null hypothesis means the SNP is significantly associated with at least one of the subtypes. The null hypothesis may be rejected if the SNP is significantly associated with a similar effect size across all subtypes (i.e. $\boldsymbol{\theta}^{(0)} \neq 0$, $\boldsymbol{\theta}_H = \mathbf{0}$), or if the SNP has heterogeneous effects on different subtypes ($\boldsymbol{\theta}_H \neq \mathbf{0}$). The second hypothesis test is the global heterogeneity test,

$$H_0^{EH} : \boldsymbol{\theta}_H = \mathbf{0} \text{ versus } H_1^{EH} : \boldsymbol{\theta}_H \neq \mathbf{0}. \quad (5)$$

The global heterogeneity test evaluates for etiologic heterogeneity with respect to a SNP and all tumor characteristics simultaneously. Rejecting this null hypothesis indicates that the first stage case-control log odds ratios of at least two different subtypes are significantly different from each other. Notably, the global heterogeneity test does not identify which tumor characteristic(s) is/are driving the heterogeneity between the subtypes. To identify the tumor characteristic(s) responsible for observed heterogeneity, we propose the individual

tumor marker heterogeneity test,

$$H_0^{\text{IH}} : \theta_{H(k)} = 0 \text{ versus } H_1^{\text{IH}} : \theta_{H(k)} \neq 0, \quad (6)$$

where $\theta_{H(k)}$ is one of the case-case parameters of θ_H . The case-case parameters $\theta_{H(k)}$ provide a measurement of etiological heterogeneity according to a specific tumor characteristic (Begg and Zhang, 1994). Under the breast cancer example, we could directly test $H_0^{\text{IH}} : \theta_{\text{ER}}^{(1)} = 0$ versus $H_1^{\text{IH}} : \theta_{\text{ER}}^{(1)} \neq 0$. In this example, rejecting the null hypothesis provides evidence that the case-control log odds ratios of ER+ and ER- subtypes are significantly different.

C. EM algorithm accounting for cases with incomplete tumor characteristics

In previous sections, we assumed all the tumor characteristics are observed for every case in the study. In epidemiological research it is very common that tumor characteristic data is missing across study participants. This problem becomes exacerbated as the number of analyzed tumor characteristics grows. Restricting to cases with complete tumor characteristics can reduce statistical power and potentially introduce selection bias. To solve this problem, we propose to use the EM algorithm (Dempster *et al.*, 1977) to find the MLE of two-stage model and all available information from the study cases. Let $Y_{im} = I(D_i = m)$ denote whether the i th subject has subtype m and \mathbf{T}_{io} be the observed tumor characteristics status of the i th subject. Given the observed tumor characteristics, the possible subtypes for subject i would be a limited subset of all possible tumor subtypes, which can be denoted as $\mathcal{Y}_{io} = \{Y_{im} : Y_{im} \text{ that is consistent with } \mathbf{T}_{io}\}$. We assume that $(Y_{i1}, Y_{i2}, \dots, Y_{iM}, G_i, \mathbf{X}_i)$ are independently and identically distributed (i.i.d.), and that the tumor characteristics are

missing at random. Given the notation above, the EM algorithm at the v th iteration would be:

E step:

$$Y_{im}^E = E(Y_{im}|G_i, \mathbf{X}_i, \mathbf{T}_{io}; \boldsymbol{\delta}^{(v)}) = \frac{Pr(Y_{im}|G_i, \mathbf{X}_i; \boldsymbol{\delta}^{(v)})}{\sum_{Y_{im} \in \mathcal{Y}_{io}} Pr(Y_{im} = 1|G_i, \mathbf{X}_i; \boldsymbol{\delta}^{(v)})} \quad (7)$$

Where Y_{im}^E is the probability of the i th person to be of the m th subtype given his observed tumor characteristics, genotype and other covariates.

M step:

$$\boldsymbol{\delta}^{(v+1)} = \arg \max_{\boldsymbol{\delta}} \sum_{i=1}^N \left[\left(1 - \sum_{m=1}^M Y_{im}^E\right) \log Pr(D_i = 0|G_i, \mathbf{X}_i) + \sum_{m=1}^M Y_{im}^E \log \{Pr(D_i = m|G_i, \mathbf{X}_i)\} \right] \quad (8)$$

The M step could be solved through a weighted least square iteration steps and the details of EM algorithm procedure could be found in Supplementary Section 2. The MLE of the second stage parameters (denoted as $\hat{\boldsymbol{\delta}}$) can be obtained when the EM algorithm converges.

Let $\mathbf{Y}_m = (Y_{1m}, \dots, Y_{Nm})^T$, and $\mathbf{Y} = (\mathbf{Y}_1^T, \dots, \mathbf{Y}_M^T)^T$. Let $\mathbf{C} = (\mathbf{G}, \mathbf{X})$ and $\mathbf{C}_M = \mathbf{I}_M \otimes \mathbf{C}$. Following (Louis, 1982), the observed information matrix \mathbf{I} would be:

$$\mathbf{I} = \mathbf{Z}^T \mathbf{C}_M^T \mathbf{W} \mathbf{C}_M^T \mathbf{Z} \quad (9)$$

where the weighted matrix $\mathbf{W} = \mathbf{D} - \mathbf{A}\mathbf{A}^T$, with $\mathbf{D} = \text{diag}(\mathbf{P} - \mathbf{P}_{\text{mis}})$, $\mathbf{P} = E(\mathbf{Y}|\mathbf{C}; \hat{\boldsymbol{\delta}})$, $\mathbf{P}_{\text{mis}} = E(\mathbf{Y}|\mathbf{C}, \mathbf{T}_o; \hat{\boldsymbol{\delta}})$, and $\mathbf{A} = \mathbf{D}(\mathbf{1}_M \otimes \mathbf{I}_N)$. We can construct the Wald test statistics for the global association test, global etiological heterogeneity test, and individual tumor characteristic heterogeneity test using the MLE of corresponding second stage parameters $\hat{\boldsymbol{\theta}}^*$ and covariance matrix $\hat{\boldsymbol{\Sigma}}$:

$$\hat{\boldsymbol{\theta}}^{*T} \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\theta}}^* \sim \chi_l^2, \quad (10)$$

where the degrees of freedom l equals the length of $\hat{\boldsymbol{\theta}}^*$.

D. Fixed effect two-stage polytomous model (FTOP) score test

Although the hypothesis tests could be implemented through the Wald test, estimating the model parameters for all SNPs in the genome is time consuming and computationally intensive. In this section, we develop a score test for the global association test assuming the second stage parameters as fixed. The score test only needs to estimate the second stage parameters of \mathbf{X} under the null hypothesis once, which makes it much more computational efficient than the Wald test.

Let $\mathbf{G}_M = \mathbf{I}_M \otimes \mathbf{G}$, and $\mathbf{X}_M = \mathbf{I}_M \otimes \mathbf{X}$. Under the null hypothesis, $H_0 : \boldsymbol{\theta} = \mathbf{0}$, the score of $\boldsymbol{\theta}$ is $U_{\boldsymbol{\theta}}(\hat{\boldsymbol{\lambda}}) = \mathbf{Z}_G^T \mathbf{G}_M^T (\mathbf{Y} - \mathbf{P}_f)$, where $\mathbf{P}_f = E_{\boldsymbol{\theta}=\mathbf{0}}(\mathbf{Y}|\mathbf{X}; \hat{\boldsymbol{\lambda}})$. The corresponding efficient information matrix is:

$$\tilde{\mathbf{I}} = \mathbf{I}_{\boldsymbol{\theta}\boldsymbol{\theta}} - \mathbf{I}_{\boldsymbol{\theta}\boldsymbol{\lambda}}^T \mathbf{I}_{\boldsymbol{\lambda}\boldsymbol{\lambda}} \mathbf{I}_{\boldsymbol{\lambda}\boldsymbol{\theta}}, \quad (11)$$

where $\mathbf{I}_{\boldsymbol{\theta}\boldsymbol{\theta}} = \mathbf{Z}_G^T \mathbf{G}_M^T \mathbf{W}_f \mathbf{G}_M \mathbf{Z}_G$, $\mathbf{I}_{\boldsymbol{\lambda}\boldsymbol{\lambda}} = \mathbf{Z}_X^T \mathbf{X}_M^T \mathbf{W}_f \mathbf{X}_M \mathbf{Z}_X$, and $\mathbf{I}_{\boldsymbol{\lambda}\boldsymbol{\theta}} = \mathbf{I}_{\boldsymbol{\theta}\boldsymbol{\lambda}}^T = \mathbf{Z}_X^T \mathbf{X}_M^T \mathbf{W}_f \mathbf{G}_M \mathbf{Z}_G$.

The weighted matrix \mathbf{W}_f has the same definition as in Equation 9, but evaluated under the null hypothesis $H_0 : \boldsymbol{\theta} = \mathbf{0}$. The score test statistics $Q_{\boldsymbol{\theta}}$ for fixed-effect two stage model would be:

$$Q_{\boldsymbol{\theta}} = U_{\boldsymbol{\theta}}(\hat{\boldsymbol{\lambda}})^T \tilde{\mathbf{I}}^{-1} U_{\boldsymbol{\theta}}(\hat{\boldsymbol{\lambda}}) \sim \chi_l^2, \quad (12)$$

where the degrees of freedom l equal the length of $U_{\boldsymbol{\theta}}(\hat{\boldsymbol{\lambda}})$.

E. Mixed effect two-stage polytomous model (MTOP) score test

The two-stage model decreases the degrees of freedom compared to the polytomous logistic regression; however, the power gains in the two-stage model can be lost as additional tumor characteristics are added into the model. We further propose a mixed-effect two-stage model by modeling some of the second stage case-case parameters as a random effect. Let $\mathbf{u} = (u_1, \dots, u_s)^T$, where each u_j follows an arbitrary distribution F with mean zero and variance σ^2 . The mixed effect second stage model links the first and second stage parameters via the following:

$$\beta = \mathbf{Z}_f \theta_f + \mathbf{Z}_r \mathbf{u}, \quad (13)$$

where \mathbf{Z}_f is the second stage design matrix of fixed effect, \mathbf{Z}_r is the second stage design matrix of random effect, and θ_f are the fixed-effect second stage parameters. Let $\theta_f = (\theta^{(0)}, \theta_{\text{HH}}^T)^T$, where $\theta^{(0)}$ is the case-control log odds ratio of the reference subtype and θ_{HH} are the fixed case-case parameters. The baseline effect $\theta^{(0)}$ is always kept fixed, since the baseline effect parameter captures the SNP's overall effect on all the cancer subtypes.

The fixed case-case parameters θ_{HH} can be used for the tumor characters with prior information suggesting that they are a source of heterogeneity. And the random effect case-case parameters \mathbf{u} can be used for tumor characteristics with little or no prior information to suggest that they are a source of heterogeneity. Under the breast cancer example, the baseline parameter ($\theta^{(0)}$) and the case-case parameter for ER (θ_{HH}) could be modeled fixed effects, since previous evidence indicates ER as a source breast cancer heterogeneity ([Garcia-](#)

Closas *et al.*, 2013). And the case-case parameters of PR and HER2 can be modeled as random effect (\mathbf{u}).

Under the mixed effect two-stage model, the global association test would be:

$$H_0^A : \boldsymbol{\theta}_f = \mathbf{0}, \sigma^2 = 0 \text{ versus } H_1^A : \boldsymbol{\theta}_f \neq \mathbf{0} \text{ or } \sigma^2 \neq 0, \quad (14)$$

And the global etiology heterogeneity test would be:

$$H_0^{EH} : \boldsymbol{\theta}_{fH} = \mathbf{0}, \sigma^2 = 0 \text{ versus } H_1^{EH} : \boldsymbol{\theta}_{fH} \neq \mathbf{0} \text{ or } \sigma^2 \neq 0. \quad (15)$$

We derive the corresponding score statistics and associated distribution under two-stage model by drawing parallels from recent studies on association tests for groups of rare variants using kernel machine regression methodology (Lin, 1997; Sun *et al.*, 2013; Wu *et al.*, 2011; Zhang and Lin, 2003). The score statistics of fixed effect $\boldsymbol{\theta}_f$ under the global null $H_0^A : \boldsymbol{\theta}_f = \mathbf{0}, \sigma^2 = 0$ would be:

$$Q_{\boldsymbol{\theta}_f} = (\mathbf{Y} - \mathbf{P}_f)^T \mathbf{G}_M \mathbf{Z}_f \tilde{\mathbf{I}}_f^{-1} \mathbf{Z}_f^T \mathbf{G}_M^T (\mathbf{Y} - \mathbf{P}_f) \sim \chi_{l_f}^2, \quad (16)$$

where $\mathbf{P}_f = E_{\boldsymbol{\theta}_f=\mathbf{0}, \sigma^2=0}(\mathbf{Y}|\mathbf{X}; \hat{\boldsymbol{\lambda}})$. Here $\tilde{\mathbf{I}}_f$ has the same definition as Equation 11, but substitute \mathbf{Z}_G with \mathbf{Z}_f . Under the null hypothesis, $Q_{\boldsymbol{\theta}_f}$ follows a χ^2 distribution, and the degrees of freedom l_f is the same as the length of $\boldsymbol{\theta}_f$.

Let $\boldsymbol{\tau} = (\boldsymbol{\theta}_f^T, \boldsymbol{\lambda}^T)^T$ be the second stage fixed effect, and \mathbf{Z}_τ is the corresponding second stage design matrix. The variance component score statistics of σ^2 under the null hypothesis: $H_0 : \sigma^2 = 0$ without constraining $\boldsymbol{\theta}_f$ would be:

$$Q_{\sigma^2} = (\mathbf{Y} - \mathbf{P}_r)^T \mathbf{G}_M \mathbf{Z}_r \mathbf{Z}_r^T \mathbf{G}_M^T (\mathbf{Y} - \mathbf{P}_r) \sim \sum_{i=1}^s \rho_i \chi_{i,1}^2, \quad (17)$$

where $\mathbf{P}_r = E_{\sigma^2=0}(\mathbf{Y}|\mathbf{G}, \mathbf{X}; \hat{\boldsymbol{\tau}})$, and $\hat{\boldsymbol{\tau}}$ is the MLE under the null hypothesis: $H_0 : \sigma^2 = 0$.

Under the null hypothesis, Q_{σ^2} follows a mixture of chi square distribution, where $\chi_{i,1}^2$

i.i.d. follows χ_1^2 . (ρ_1, \dots, ρ_s) are the eigenvalues of $\tilde{\mathbf{I}}_r = \mathbf{I}_{\mathbf{uu}} - \mathbf{I}_{\mathbf{u}\boldsymbol{\tau}}^T \mathbf{I}_{\boldsymbol{\tau}\boldsymbol{\tau}}^{-1} \mathbf{I}_{\boldsymbol{\tau}\mathbf{u}}$, with $\mathbf{I}_{\mathbf{uu}} =$

$\mathbf{Z}_r^T \mathbf{G}_M^T \mathbf{W}_r \mathbf{G}_M \mathbf{Z}_r$, $\mathbf{I}_{\boldsymbol{\tau}\boldsymbol{\tau}} = \mathbf{Z}_{\boldsymbol{\tau}}^T \mathbf{C}_M^T \mathbf{W}_r \mathbf{C}_M \mathbf{Z}_{\boldsymbol{\tau}}$ and $\mathbf{I}_{\boldsymbol{\tau}\mathbf{u}} = \mathbf{I}_{\mathbf{u}\boldsymbol{\tau}}^T = \mathbf{Z}_{\boldsymbol{\tau}}^T \mathbf{C}_M^T \mathbf{W}_r \mathbf{G}_M \mathbf{Z}_r$. The weighted

matrix \mathbf{W}_r has the same definition as the one used for Equation 9, but evaluated under

the null hypothesis $H_0 : \sigma^2 = 0$. The Davies exact method (Davies, 1980) is used here to

calculate the p-value of the mixture of chi square distribution. The details of the derivation

of Q_{σ^2} are in Supplementary Section 3.

Following similar logic as (Sun et al., 2013), we prove that $Q_{\boldsymbol{\theta}_f}$ and Q_{σ^2} are independent

with each other (see proof in Supplementary Section 4). We use Fisher's procedure (Koziol

and Perlman, 1978) of to combined the p-value coming out from the two independent tests.

Let $P_{\boldsymbol{\theta}_f} = Pr(Q_{\boldsymbol{\theta}_f} \geq \chi_{l_f}^2)$ and $P_{\sigma^2} = Pr(Q_{\sigma^2} \geq \sum_{i=1}^s \rho_i \chi_{i,1}^2)$. Under the null hypothesis

$H_0^A : \boldsymbol{\theta}_f = \mathbf{0}, \sigma^2 = 0$, $-2 \log(P_{\boldsymbol{\theta}_f}) - 2 \log(P_{\sigma^2})$ follows χ_4^2 . Then the p-value P_{mix} of mixed

effect two-stage model score test under the null hypothesis would be:

$$P_{\text{mix}} = Pr \{ -2 \log(P_{\boldsymbol{\theta}_f}) - 2 \log(P_{\sigma^2}) \geq \chi_4^2 \}. \quad (18)$$

The extension of the score statistics to global test for etiology heterogeneity, $H_0^{\text{EH}} : \boldsymbol{\theta}_{\text{EH}} =$

$\mathbf{0}, \sigma^2 = 0$, would be straightforward.

III. SIMULATION EXPERIMENTS

In this section, large scale simulations across a wide range of practical scenarios are

conducted to evaluate the type I error and power of the fixed effect and mixed effect two-stage

models. Data were simulated to mimic the PBCS. Four tumor characteristics were simulated: ER (positive vs. negative), PR (positive vs negative), HER2 (positive vs. negative) and grade (ordinal 1, 2, 3). This defined a total of $2^3 \times 3 = 24$ breast cancer subtypes.

In each case control simulation, genotype data \mathbf{G} was assumed to be under Hardy-Weinberg equilibrium in the underlying population with a minor allele frequency (MAF) of 0.25. An additional covariate, \mathbf{X} , was simulated as a standard normal distribution independent of \mathbf{G} . We used polytomous logistic regression model as Equation 19 to simulate a multinomial outcome with 25 groups, one for the control subjects, and the other 24 for different cancer subtypes.

$$Pr(D_i = m|X_i) = \frac{\exp(\alpha_m + \beta_m G_i + 0.05 X_i)}{1 + \sum_{m=1}^M \exp(\alpha_m + \beta_m G_i + 0.05 X_i)}, \quad (19)$$

where β_m is the log OR of G for m th subtype v.s. control. The effect of covariate \mathbf{X} was set as 0.05 across all the subtypes. By using the frequency of 24 breast cancer subtypes estimated from Breast Cancer Association Consortium data (Supplementary Table 1), we computed the corresponding polytomous logistic regression intercept parameters α_m . The cases and controls ratio was set to be around 1:1, and the proportion of ER+, PR+ and HER2+ were set to be 0.81, 0.68, and 0.17, respectively. The proportion of grade 1, 2, and 3 were 0.20, 0.48, and 0.32, respectively. The missing tumor markers were randomly selected and the missing rate of ER, PR, HER2, and grade were set to be 0.17, 0.25, 0.42, and 0.27, respectively. Under this simulation setting, around 70% breast cancer cases had at least one missing tumor characteristic.

A. Type I error

In this subsection, we evaluated the type I error of global tests for association, global tests for heterogeneity, and individual heterogeneity test for the tumor characteristics under the global null hypothesis. We assumed $\beta_m = 0$ in Equation 19, where none of the subtypes is associated with genotypes. The total sample size n was set to be 5,000, 50,000, and 100,000. And 2.4×10^7 simulations were conducted to evaluate the type I error at $\alpha = 10^{-4}, 10^{-5}$ and, 10^{-6} level.

We applied both MTOP and FTOP with an additive second stage design structure as in Equation 20, where the subtype-specific case-controls log ORs were specified into the case-control log OR of a baseline disease subtype (ER- , PR- , HER2-, grade 1) and case-case parameters associated with the four tumor markers. All of the second stage interactions parameters were constrained to be 0. Furthermore, the MTOP assumed the baseline parameter and the ER case-case parameters as fixed effects, and the PR, HER2, and grade case-case parameters as random effects.

$$\beta_{s_1 s_2 \dots s_K} = \theta^{(0)} + \sum_{k_1=1}^4 \theta_{k_1}^{(1)} s_{k_1}. \quad (20)$$

Table I presents the simulated estimated type I error under the global null hypothesis. As expected, the type I error for both MTOP and FTOP tended to be lower for the simulated sample size of 5000, but with larger samples sizes all tests report nearly correct type I error, demonstrating the validity of MTOP and FTOP.

B. Statistical power

In this subsection, we present the statistical power of MTOP and FTOP under three different scenarios using our breast cancer example: I. no heterogeneity between tumor markers, II. heterogeneity according to one tumor marker, and III. heterogeneity according to multiple tumor markers. We generated the subtypes through Equation 19. Under the scenario I, we set β_m as 0.05 for all the subtypes, thus no heterogeneity of the ORs between **G** and the subtypes. For scenarios II and III, it was assumed that β_m followed the additive second stage structure as in Equation 20. Under scenarios II, we simulated a situation with only ER heterogeneity by setting the baseline effect $\theta^{(0)}$ to be 0, the case-case parameter for ER $\theta_1^{(1)}$ was set to be 0.05, and the PR, HER2, and grade case-case parameters to be 0. For scenario III, we simulated a situation with heterogeneity according to all 4 tumor markers by setting the baseline effect $\theta^{(0)}$ set to be 0, the ER $\theta_1^{(1)}$ case-case parameter to be 0.05, and all the other three case-case parameters were set to follow a normal distribution with mean 0 and variance 4×10^{-4} . Under this scenario, all the tumor characteristics contribute subtype-specific heterogeneity. The total sample size n was set to be 25,000, 50,000, and 100,000. We performed 10^5 simulations were conducted to evaluate the power at $\alpha < 10^{-3}$ level.

We compared the statistical power to detect the genetic association between MTOP, FTOP, a standard logistic regression, polytomous logistic regression, and a two-stage model that only uses cases with complete tumor characteristics. The same additive second stage structure as Section III A was used for MTOP and FTOP. When we applied MTOP, FTOP

and polytomous model, we removed all the subtypes with fewer than 10 cases to avoid
unconvergence of the model. All the methods were set to test the overall association between
G and the risk of the cancer.

To evaluate the different methods under a larger number of tumor characteristics, we
added two additional binary tumor characteristics to the previous breast cancer example.
This defined a total of $2^5 \times 3 = 96$ cancer subtypes. Similar to the four tumor characteristics
simulations, we generated subtypes by polytomous model as in Equation 19, and simulated
data under three different scenarios: I. no heterogeneity II. one tumor marker drove the
heterogeneity, and III. multiple tumor markers driving tumor heterogeneity. Total sample
size was set to be 25,000, 50,000, and 100,000. The two additional tumor characteristics
were randomly selected to missing with 5% missing rate. Under this setting, around 77% of
the cases have at least one tumor characteristics missing. We performed 10^5 simulations to
evaluate the power at $\alpha < 10^{-3}$ level.

Figure 1 shows the power comparison between the five methods under different scenarios.
Overall, MTOP had robust power under all the heterogeneity scenarios. Under scenario I
with no subtype-specific heterogeneity, standard logistic regression had the highest power,
but suffered from substantial power loss when heterogeneity existed between subtypes. When
heterogeneity was introduced in scenarios II and III, MTOP, followed by FTOP, consistently
demonstrated the highest power among the five methods. The higher power observed in
MTOP, relative to FTOP, ranged from 102% to 168%. Under the scenarios with four tumor
characteristics the difference in degrees of freedom between MTOP and FTOP were small,
therefore MTOP had only a slight power advantage. However, with six tumor markers,

the differences in degrees of freedom between MTOP and FTOP becomes more apparent, as does the greater power of MTOP. FTOP is least efficient in scenarios of none or little heterogeneity, such as scenarios I and II, but with increasing sources of heterogeneity, such-as scenario III, the power of MTOP and FTOP are more similar.

Simulation study also shows that incorporation of cases with missing tumor characteristics significantly increased the power of the methods. Under the four tumor markers setting with around 70% incomplete cases, the power of MTOP was between 202% to 905% greater compared to the original two-stage model using only complete data. Under the six tumor markers setting with around 77% incomplete cases, the two-stage model with only complete data lost more power compared to the four tumor markers setting; however, MTOP maintained similar power.

Overall, in the scenario of no heterogeneity the standard logistic regression demonstrated the most powerful. However, in the presence of subtype heterogeneity, MTOP was the most powerful method. The polytomous model had the lowest power across all of the settings.

IV. APPLICATION TO THE POLISH BREAST CANCER STUDY (PBCS)

In this section, we used data from the PBCS, a population-based breast cancer case-control study conducted in Poland between 2000 and 2003 ([García-Closas *et al.*, 2006](#)). The study population consisted of 2,078 cases of histologically or cytologically confirmed invasive breast cancer, and 2,219 women without a history of breast cancer at enrollment. Tumor characteristic information on ER, PR, and grade were available from pathology records ([García-Closas *et al.*, 2006](#)) and information on HER2 status was available from

immunohistochemical staining of tissue microarray blocks (Yang *et al.*, 2007). We used genome-wide genotyping data to compare MTOP, FTOP, standard logistic regression, and polytomous logistic regression to detect SNPs associated with breast cancer risk. Table II presents the sample size distribution of the tumor characteristics. Combinations of the four tumor characteristics define 24, mutually exclusive breast cancer subtypes. Subtypes with less than 10 cases were excluded, leaving 17 subtypes that were evaluated. Both MTOP and FTOP used additive second stage design structure as in Equation 20. In addition, MTOP assumed the baseline parameter and the ER case-case parameter as fixed effects, and the case-case parameters of PR, HER2 and grade as random effects. We put ER as a fixed effect because of the previously reported heterogeneity of genetic association by ER status (Garcia-Closas *et al.*, 2013). Genotype imputation was done using IMPUTE2 based on 1000 Genomes Project dataset as reference panel (Michailidou *et al.*, 2017; Milne *et al.*, 2017). In total, 7,017,694 common variants on 22 auto chromosomes with $MAF \geq 5\%$ were included in the analysis. In the four models, we adjusted for age, and the first four genetic principal components to account for population stratification bias.

Figure 2 shows the Manhattan plot of genome-wide associations analysis with PBCS using the four different methods. MTOP, FTOP and standard logistic regression identified a known susceptibility variant in the FGFR2 locus on chromosome 10 (Michailidou *et al.*, 2013), with the most significant SNP being rs11200014 ($P < 5 \times 10^{-8}$). Further, both MTOP and FTOP identified a second known susceptibility locus on chromosome 11 (CCND1) (Michailidou *et al.*, 2017), with the most significant SNP in both models being rs78540526 ($P < 5 \times 10^{-8}$). The individual heterogeneity test of this SNP showed evidence

for heterogeneity by ER ($P=0.011$) and grade ($P=0.024$). Notably, the CCND1 locus was not genome-wide significant in standard logistic regression or polytomous models. The QQ plots of the four GWAS analysis can be found in Supplementary Figure 1.

Next, we compared the ability of the same MTOP and standard logistic regressions to detect 178 previously identified breast cancer susceptibility loci (Michailidou *et al.*, 2017). As shown in Table III, for eight of the 178 loci, the MTOP global association test p value was more than ten fold lower compared to the standard logistic regression p value. In the MTOP model, these eight loci all had significant global heterogeneity tests ($P < 0.05$). Confirming these results, in a previous analysis applying MTOP to 106,571 breast cancer cases and 95,762 controls, these eight loci were reported to have significant global heterogeneity (Ahearn *et al.*, 2017).

V. DISCUSSION

We present novel methods for performing genetic association testing for cancer outcomes accounting for potential heterogeneity across subtypes defined by multiple, correlated, tumor characteristics. These methods efficiently account for multiple testing, correlation between markers, and missing tumor data. We demonstrated that MTOP has greater statistical power in the presence of subtype heterogeneity than either standard logistic regression or polytomous logistic regression. Moreover, we show that the EM algorithm is an efficient method for handling missing data and substantially increases statistical power. Furthermore, we developed a publicly available R package called TOP (two-stage polytomous logistic

regression, <https://github.com/andrewhaoyu/TOP>), which includes all the core functions implemented in C code.

Several statistical methods have been proposed to study the etiological heterogeneity of cancer subtypes (Chatterjee, 2004; Rosner *et al.*, 2013; Wang *et al.*, 2015). A recent review showed the well controlled type-one error and good statistical power of two-stage model among these methods (Zabor and Begg, 2017). However, previous two-stage models have not accounted for missing tumor marker, a common problem of epidemiological studies. We show that by incorporating the EM algorithm into the two-stage model we can take advantage of all available information and make substantial gains in statistical power (as shown in Figure 1). Moreover, we show that modeling some of the second stage parameters as random effects that follow an arbitrary distribution with mean 0 and variance σ^2 is an efficient method to mitigate the degrees of freedom penalty caused by analyzing a large number of tumor characteristics.

Notably, the computation time of MTOP is markedly greater than FTOP due to estimating the coefficients of covariates. To construct the score tests in FTOP, the coefficients of covariates need to be estimated once under the null hypothesis, while for MTOP they need to be estimated for every SNP. The computational complexity of FTOP is $O(NMS_\theta)$, with S_θ as the number of second stage parameters of \mathbf{G} ; whereas the computational complexity for MTOP is $O(N^3M^2P^2S_\tau)$, with S_τ as the number of fixed effect second stage parameters of \mathbf{G} and \mathbf{X} .

We parallel the recent studies on rare genetic association tests using kernel machine regression methods to develop MTOP (Lin, 1997; Sun *et al.*, 2013; Wu *et al.*, 2011; Zhang

and Lin, 2003). Currently, we have only implemented the linear kernel in MTOP, but other kernel functions that capture the similarity between tumor characteristics could be implemented in the future. If there is prior knowledge about the genetic architecture of different tumor subtypes, this could help to choose the kernel function and improve the power of the methods.

In conclusion, we have proposed an efficient and systematic approach for incorporating tumor characteristics information to identify genetic associations in the presence of subtype heterogeneity. The methods leverage all available tumor information and have robust statistical power. We have limited our demonstration of the benefit of these methods to analyzing the association between genetic variants and breast cancer subtypes; however these methods can easily be applied to the analysis of other non-genetic risk factors and/or other endpoints characterized by subtypes. The proposed methods have been implemented in a user-friendly and high-speed R statistical package called TOP (<https://github.com/andrewhaoyu/TOP>).

VI. SUPPLEMENTARY MATERIAL

In Supplementary Section 1, we describe generalizing two-stage polytomous model to multiple variates with different second stage design matrix. In Supplementary Section 2, we derive the EM algorithm under two-stage model. In Supplementary Section 3, we derive the variance component score statistics in two-stage model. In Supplementary Section 4, we prove the independence between Q_{θ_f} and Q_{σ^2} . The Supplementary Table 1 contains the 24 breast cancer subtypes frequency estimated from Breast Cancer Association Consortium data. The Supplementary Figure 1 is the QQ plot of GWAS with PBCS.

ACKNOWLEDGMENTS

This work was supported by funds from the NCI Intramural Research Program and Bloomberg Distinguished Professorship endowment. The simulation experiments and data analysis were implemented using the high performance computation Biowulf cluster at National Institutes of Health, USA.

Ahearn, T. *et al.* (2017). “Novel analysis incorporating multiple tumor characteristics provide evidence of highly heterogeneous associations for known breast cancer risk loci,” The American Society of Human Genetics poster .

Barnard, M. E., Boeke, C. E., and Tamimi, R. M. (2015). “Established breast cancer risk factors and risk of intrinsic tumor subtypes,” *Biochim Biophys Acta* **1856**(1), 73–85.

Begg, C. B., and Zhang, Z. (1994). “Statistical analysis of molecular epidemiology studies employing case-series,” *Cancer Epidemiology and Prevention Biomarkers* **3**(2), 173–175.

Cancer Genome Atlas, N. (2012). “Comprehensive molecular portraits of human breast tumours,” *Nature* **490**(7418), 61–70.

Cancer Genome Atlas Research, N. (2012). “Comprehensive genomic characterization of squamous cell lung cancers,” *Nature* **489**(7417), 519–25.

Cancer Genome Atlas Research, N. (2014). “Comprehensive molecular profiling of lung adenocarcinoma,” *Nature* **511**(7511), 543–50.

- Chatterjee, N. (2004). “A two-stage regression model for epidemiological studies with multivariate disease classification data,” *Journal of the American Statistical Association* **99**(465), 127–138.
- Curigliano, G., Burstein, H. J., P Winer, E., Gnant, M., Dubsy, P., Loibl, S., Colleoni, M., Regan, M. M., Piccart-Gebhart, M., Senn, H.-J. *et al.* (2017). “De-escalating and escalating treatments for early-stage breast cancer: the st. gallen international expert consensus conference on the primary therapy of early breast cancer 2017,” *Annals of Oncology* **28**(8), 1700–1712.
- Davies, R. B. (1980). “Algorithm as 155: The distribution of a linear combination of χ^2 random variables,” *Journal of the Royal Statistical Society. Series C (Applied Statistics)* **29**(3), 323–333.
- Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). “Maximum Likelihood from Incomplete Data via the EM Algorithm,” *Journal of the Royal Statistical Society. Series B* **39**(1), 1–38.
- Dubin, N., and Pasternack, B. S. (1986). “Risk assessment for case-control subgroups by polychotomous logistic regression,” *American journal of epidemiology* **123**(6), 1101–1117.
- Falk, R. T. *et al.* (2014). “Alcohol and risk of breast cancer in postmenopausal women: an analysis of etiological heterogeneity by multiple tumor characteristics,” *American journal of epidemiology* **180**(7), 705–717.
- García-Closas, M. *et al.* (2006). “Established breast cancer risk factors by clinically important tumour characteristics,” *British Journal of Cancer* **95**(1), 123–129.

Garcia-Closas, M. *et al.* (2013). “Genome-wide association studies identify four er negative-specific breast cancer risk loci,” *Nat Genet* **45**(4), 392–8, 398e1–2.

Gortmaker, S. L., Hosmer, D. W., and Lemeshow, S. (1994). “Applied Logistic Regression.,” *Contemporary Sociology* .

Koziol, J. A., and Perlman, M. D. (1978). “Combining independent chi-squared tests,” *Journal of the American Statistical Association* **73**(364), 753–763.

Lin, X. (1997). “Variance component testing in generalised linear models with random effects,” *Biometrika* **84**(2), 309–326.

Louis, T. A. (1982). “Finding the observed information matrix when using the em algorithm,” *Journal of the Royal Statistical Society. Series B (Methodological)* 226–233.

MacArthur, J. *et al.* (2016). “The new nhgri-ebi catalog of published genome-wide association studies (gwas catalog),” *Nucleic acids research* **45**(D1), D896–D901.

Michailidou, K. *et al.* (2013). “Large-scale genotyping identifies 41 new loci associated with breast cancer risk,” *Nature Genetics* .

Michailidou, K. *et al.* (2017). “Association analysis identifies 65 new breast cancer risk loci,” *Nature* **551**(7678), 92–94.

Milne, R. L. *et al.* (2017). “Identification of ten variants associated with risk of estrogen-receptor-negative breast cancer,” *Nature genetics* **49**(12), 1767.

Perou, C. M. *et al.* (2000). “Molecular portraits of human breast tumours,” *Nature* **406**(6797), 747–52.

Peters, U. *et al.* (2004). “Calcium intake and colorectal adenoma in a US colorectal cancer early detection program.,” *The American journal of clinical nutrition* .

Prat, A. *et al.* (2015). “Clinical implications of the intrinsic molecular subtypes of breast cancer,” *Breast* .

Rosner, B. *et al.* (2013). “Breast cancer risk prediction with heterogeneous risk profiles according to breast cancer tumor markers,” *American Journal of Epidemiology* **178**(2), 296–308.

Sherman, M. E. *et al.* (2007). “Variation in breast cancer hormone receptor and HER2 levels by etiologic factors: A population-based analysis,” *International Journal of Cancer* .

Sun, J., Zheng, Y., and Hsu, L. (2013). “A Unified Mixed-Effects Model for Rare-Variant Association in Sequencing Studies,” *Genetic Epidemiology* **37**(4), 334–344.

Visscher, P. M. *et al.* (2017). “10 years of gwas discovery: Biology, function, and translation,” *Am J Hum Genet* **101**(1), 5–22.

Wang, M., Kuchiba, A., and Ogino, S. (2015). “A meta-regression method for studying etiological heterogeneity across disease subtypes classified by multiple biomarkers,” *American Journal of Epidemiology* .

Wu, M. C. *et al.* (2011). “Rare-variant association testing for sequencing data with the sequence kernel association test,” *American Journal of Human Genetics* **89**(1), 82–93.

Yang, X. R. *et al.* (2007). “Differences in risk factors for breast cancer molecular subtypes in a population-based study,” *Cancer Epidemiology and Prevention Biomarkers* **16**(3), 439–443.

Zabor, E. C., and Begg, C. B. (2017). “A comparison of statistical methods for the study of etiologic heterogeneity,” *Statistics in medicine* **36**(25), 4050–4060.

517 Zhang, D., and Lin, X. (2003). “Hypothesis testing in semiparametric additive mixed mod-
518 els,” *Biostatistics* 4(1), 57–74.

TABLE I. Type one error estimates of MTOP, FTOP with 2.4×10^7 randomly simulated samples.

Global test for association and global test for heterogeneity were applied with FTOP and MTOP.

Heterogeneity test for a tumor marker was applied with only FTOP. All of the type error rates are divided by the α level.

Interested tests	Total sample size	MTOP			FTOP		
		$\alpha = 10^{-4}$	$\alpha = 10^{-5}$	$\alpha = 10^{-6}$	$\alpha = 10^{-4}$	$\alpha = 10^{-5}$	$\alpha = 10^{-6}$
Global association test	5,000	.99	.97	.88	.91	.91	.67
	50,000	.98	1.0	1.0	.99	1.0	.93
	100,000	1.0	.94	1.0	1.0	1.0	1.0
Global heterogeneity test	5,000	1.0	.97	.89	.92	.85	.55
	50,000	1.0	1.0	1.0	1.0	1.0	1.0
	100,000	1.0	.94	1.0	1.0	.98	.97
Heterogeneity test for a tumor marker	5,000				.92	.93	.76
	50,000				.98	.97	1.0
	100,000				1.0	.97	1.0

TABLE II. Sample size of four tumor characteristics in Polish Breast Cancer Study

	ER	PR	HER2	Grade	
Positive	1316	1056	1246	Grade 1	356
Negative	594	847	254	Grade 2	968
Missing	168	157	578	Grade 3	554
				Missing	200

TABLE III. Analysis results of previously identified susceptibility loci. For the listed eight loci, MTOP global association test p value decreased more than ten fold compared to the standard logistic regression p value. All of the loci are significant in global heterogeneity test ($P < 0.05$).

SNP	Chr. ^a	Position	MAF ^b	Global association p	Standard analysis p	Global heterogeneity p
rs4973768	3	27,416,013	.47	3.12×10^{-2}	9.53×10^{-1}	9.48×10^{-3}
rs10816625	9	110,837,073	.06	4.98×10^{-2}	9.79×10^{-1}	2.22×10^{-2}
rs7904519	10	114,773,927	.46	6.51×10^{-2}	8.48×10^{-1}	3.07×10^{-2}
rs554219	11	69,331,642	.13	7.34×10^{-11}	1.42×10^{-7}	5.13×10^{-6}
rs11820646	11	129,461,171	.40	1.48×10^{-2}	8.62×10^{-1}	4.53×10^{-3}
rs2236007	14	37,132,769	.21	2.10×10^{-3}	1.93×10^{-1}	3.49×10^{-3}
rs1436904	18	24,570,667	.40	7.17×10^{-4}	6.61×10^{-2}	9.69×10^{-4}
rs1436904	22	29,121,087	.01	9.83×10^{-3}	1.61×10^{-1}	2.32×10^{-2}

^aChr. chromosome. ^b MAF, minor allele frequency.

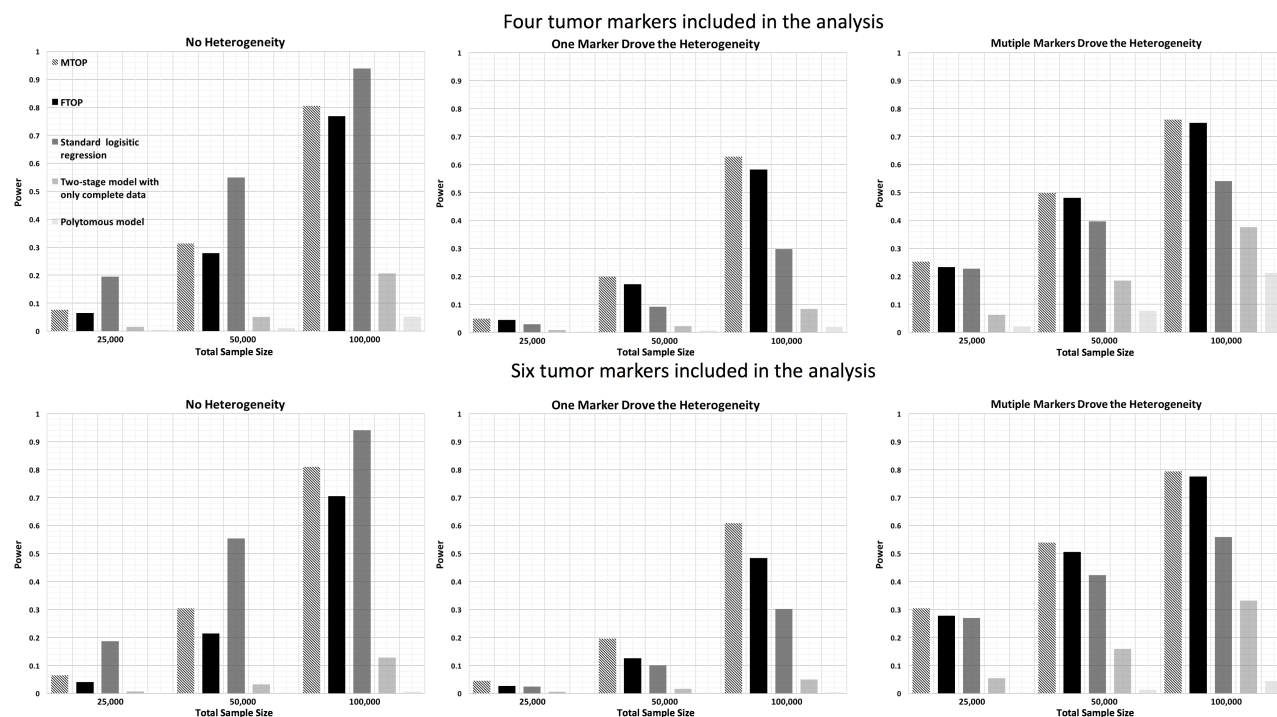


FIG. 1. Power comparison among MTOP, FTOP, standard logistic regression, two-stage model with only complete data and polytomous model with 10^5 random samples. In first setting, four tumor markers were included in the analysis. Three binary tumor marker and one ordinal tumor marker defined 24 cancer subtypes. The missing rate for the four markers were set as 0.17, 0.25, 0.42, and 0.27 respectively. Around 70% cases would be incomplete. The total sample size was set as 25,000, 50,000 and 100,000. The case control ratio was 1:1. Under second setting, two extra binary tumor markers were included in the analysis. The six tumor markers defined 96 subtypes. The missing rate of the two extra markers were 0.05. Around 77% cases would be incomplete. The total sample size was set as 25,000, 50,000 and 100,000. The power was estimated by controlling the type one error $\alpha < 10^{-3}$.

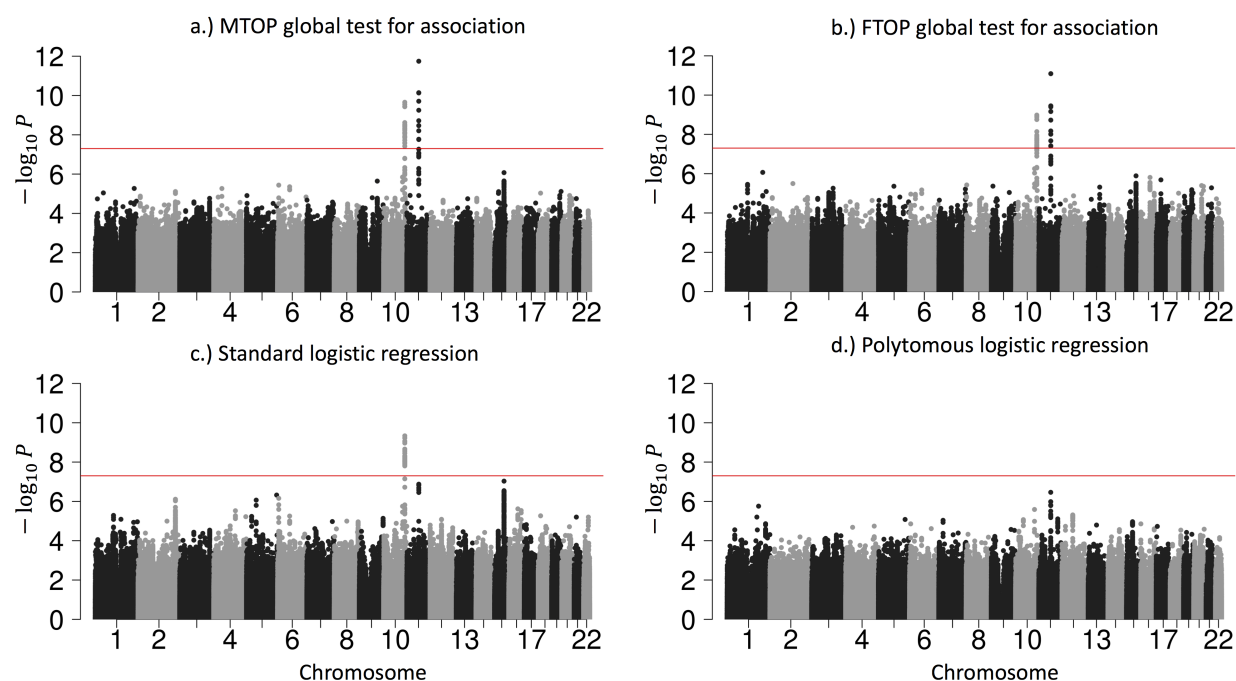


FIG. 2. Manhattan plot of genome-wide association analysis with PBCS using four different methods. PBCS have 2,078 invasive breast cancer and 2,219 controls. In total, 7,017,694 SNPs on 22 auto chromosomes with MAF more than 5% were included in the analysis. ER, PR, HER2 and grade were used to define breast cancer subtypes.