

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: 5 February 2020)

1 ABSTRACT: Cancers are routinely classified into subtypes according to various fea-
2 tures, including histopathological characteristics and molecular markers. Previous
3 genome-wide association studies have reported heterogeneous associations between
4 loci and cancer subtypes. However, it is not evident what is the optimal modeling
5 strategy for handling correlated tumor features, missing data, and increased degrees-
6 of-freedom in the underlying tests of associations. We propose to test for genetic
7 associations using a mixed-effect two-stage polytomous model score test (MTOPT).
8 In the first stage, a standard polytomous model is used to specify all possible sub-
9 types defined by the cross-classification of the tumor characteristics. In the second
10 stage, the subtype-specific case-control odds ratios are specified using a more parsimonious
11 model based on the case-control odds ratio for a baseline subtype, and the
12 case-case parameters associated with tumor markers. Further, to reduce the degrees-
13 of-freedom, we specify case-case parameters for additional exploratory markers using
14 a random-effect model. We use the Expectation-Maximization (EM) algorithm to
15 account for missing data on tumor markers. Through simulations across a range
16 of realistic scenarios and data from the Polish Breast Cancer Study (PBCS), we
17 show MTOPT outperforms alternative methods for identifying heterogeneous asso-
18 ciations between risk loci and tumor subtypes. The proposed methods have been
19 implemented in a user-friendly and high-speed R statistical package called TOP
20 (<https://github.com/andrewhaoyu/TOP>).

21 KEY WORDS: Cancer subtypes; EM algorithm; Etiologic heterogeneity; Susceptibility

variants; Score tests; Two-stage polytomous model.

23 I. INTRODUCTION

24 Genome-wide association studies (GWAS) have identified hundreds of single nucleotide
25 polymorphisms (SNPs) associated with various cancers ([MacArthur and others, 2016](#)). How-
26 ever, many cancer GWAS have often defined cancer endpoints according to specific anatomic
27 sites, and not according to subtypes of the disease. Many cancers consist of etiologically
28 and clinically heterogeneous subtypes that are defined by multiple correlated tumor charac-
29 teristics. For instance, breast cancer is routinely classified into subtypes defined by tumor
30 expression of estrogen receptor (ER), progesterone receptor (PR), and human epidermal
31 growth factor receptor 2 (HER2) ([Perou and others, 2000](#); [Prat and others, 2015](#)).

32 Increasing numbers of epidemiologic studies with tumor specimens are allowing the char-
33 acterization of cancers at the histological and molecular levels ([Cancer Genome Atlas Re-](#)
34 [search, 2014](#); [Network, 2012](#)), providing tremendous opportunities to investigate for potential
35 distinct etiological pathways between cancer subtypes. For example, a breast cancer ER-
36 negative specific GWAS reported 20 SNPs that were more strongly associated with the risk
37 of developing ER-negative than ER-positive disease ([Milne and others, 2017](#)). Previous
38 studies also suggested traditional breast cancer risk factors, such as age, obesity, and hor-
39 mone therapy use, were heterogeneously associated with the risk of breast cancer subtypes
40 ([Barnard and others, 2015](#)).

41 The most common procedure for testing for associations between risk factors and cancer
42 subtypes is by fitting a standard logistic regression for each subtype versus a control group,
43 then accounting for multiple testing. However, this procedure has several limitations. First,

44 it's common for cancer cases to have missing tumor marker data, leading to many cancer
45 cases with no subtype definition, and often these cases are dropped from the model. Second,
46 the tumor markers that defined the subtypes are commonly highly correlated with each
47 other. Testing each subtype separately without modeling the correlation limits the power
48 of the model. Finally, as the number of tumor markers increases, the number of cancer
49 subtypes dramatically increases, thus the increased degrees of freedom penalizes the power
50 of the model.

51 A two-stage polytomous logistic regression was previously proposed to characterize sub-
52 type heterogeneity of a disease according to the underlying disease characteristics ([Chatter-
53 jee, 2004](#)). The first stage of this method uses a polytomous logistic regression ([Dubin and
54 Pasternack, 1986](#)) to model subtype-specific case-control odds ratios. In the second stage,
55 the subtype-specific case-control odds ratios are decomposed into a case-control odds ratio
56 for a reference subtype, a case-case odds ratio for each tumor characteristic, and higher-order
57 interactions between the tumor characteristics. The two-stage model can reduce the degrees
58 of freedom by constraining some or all of the higher-order interactions to be 0. Moreover,
59 the second stage case-case odds ratios can be interpreted as the measures of etiological
60 heterogeneity for tumor characteristics.

61 Although the two-stage model can improve the power compared to fitting standard logistic
62 regressions for each subtype ([Chatterjee, 2004](#); [Zabor and Begg, 2017](#)), the two-stage model
63 does have notable limitations and has not been widely applied to analyze data on multiple
64 tumor characteristics. First, similar to standard logistic regression, the two-stage model can
65 not handle missing tumor characteristics, which is common in epidemiologic studies. Second,

66 the two-stage model estimation algorithm places high demands on computing power and is
67 therefore not readily applicable to large datasets. Finally, although the two-stage model
68 can reduce the multiple testing burdens compared to traditional methods, as the number
69 of tumor characteristics increases, the two-stage model can still have substantial power loss
70 due to the degrees of freedom penalty.

71 In this paper, we propose a series of computational and statistical innovations to perform
72 computationally scalable and statistically efficient association tests in large cancer GWASs
73 that incorporate tumor characteristic data. Within this two-stage modeling framework, we
74 propose three alternative types of hypotheses for testing genetic associations in the presence
75 of tumor heterogeneity. As the degrees of freedom for the tests can be large in the presence
76 of many tumor characteristics, we propose modeling parameters associated with exploratory
77 tumor characteristics using a random-effect model. We then derive the score tests under the
78 resulting mixed-effect model while taking into account missing data on tumor characteristics
79 using an efficient EM algorithm ([Dempster and others, 1977](#)). All combined, our work
80 represents a conceptually distinct and practically important extension of earlier methods
81 based on mixed-/fixed-effect models ([Lin, 1997](#); [Sun and others, 2013](#); [Wu and others, 2011](#);
82 [Zhang and Lin, 2003](#)) to the novel setting of modeling genetic associations with multiple
83 tumor characteristics.

84 The paper is organized as follows. In Section ??, we describe the proposed three different
85 hypothesis tests, the missing data algorithm, and the score tests. In Section III, we present
86 the simulation results for type I error, power and computation time. In Section IV, the
87 proposed methods are illustrated with applications using data from the Polish Breast Cancer

88 Study (PBCS). In Section V, we discuss the strengths and limitations of the methods and
89 future research directions.

90 II. TWO-STAGE POLYTOMOUS LOGISTIC MODEL

91 The details of the two-stage polytomous logistic model have been described earlier (Chat-
92 terjee, 2004). We briefly summarize them for completeness. Suppose a disease can be clas-
93 sified using K disease characteristics, and each characteristic k can be classified into M_k
94 categories; thus, the disease can be classified into $M \equiv M_1 \times M_2 \cdots \times M_K$ subtypes. For
95 example, breast cancer can be classified into eight subtypes by three tumor characteristics
96 (ER, PR, and HER2), each of which is defined as either positive or negative.

97 Let D_i denote the disease status of subject i in the study such that $D_i \in \{0, 1, 2, \dots, M\}$
98 and $i \in \{1, \dots, N\}$. $D_i = 0$ represents a control, and $D_i = m$ represents a case with
99 disease subtype m . Let G_i be the genotype for subject i , and \mathbf{X}_i be a $P \times 1$ vector of
100 other covariates, where P is the total number of other covariates. In the first stage model,
101 a “saturated” polytomous logistic regression model is constructed as follows:

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

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

104 Because each cancer subtype is defined through a unique combination of the K tumor
105 characteristics, we can always alternatively index the parameters β_m as $\{\beta_{s_1 s_2 \cdots s_K}\}$, where
106 $s_k \in \{0, 1\}$ for binary tumor characteristics, and $s_k \in \{t_1 \leq t_2 \leq \cdots \leq t_{M_k}\}$ for ordinal

107 tumor characteristics with t_1, \dots, t_{M_k} as a set of ordinal scores for M_k different levels. With
108 this new index, the log odds ratios in the first stage can be represented as follows:

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

109 where $\theta^{(0)}$ represents the case-control log odds ratio for a reference disease subtype, $\theta_{k_1}^{(1)}$
110 represents the main effect of k_1 th tumor characteristic, $\theta_{k_1 k_2}^{(2)}$ represents the second order
111 interaction between k_1 th and k_2 th tumor characteristics, and so on. A reference level can
112 be defined for each tumor characteristic, and the reference disease subtype is jointly defined
113 by the combination of the K tumor characteristics.

114 The reparameterization in 2 provides a way to decompose the first stage parameters to
115 a lower dimension. We can constrain different main effects or interaction effects to be 0 to
116 specify different second stage models. The first stage and second stage parameters can be
117 linked with a matrix form, $\beta = \mathbf{Z}_G \theta = \mathbf{Z}_G \begin{bmatrix} \theta^{(0)} & \theta_{\mathbf{H}}^T \end{bmatrix}^T$, where $\beta = (\beta_1, \beta_2, \dots, \beta_M)^T$ is a
118 vector of first stage case-control log odds ratios for all the M subtypes, $\theta^{(0)}$ is the case-control
119 log odds ratio for a reference subtype, and $\theta_{\mathbf{H}}$ is a vector containing the main effects and
120 interactions effects in the second stage. We will refer to $\theta_{\mathbf{H}}$ as case-case parameters, and
121 $\theta = (\theta^{(0)}, \theta_{\mathbf{H}}^T)^T$ as the vector of second stage parameters. \mathbf{Z}_G is the second stage design
122 matrix connecting the first stage and second stage parameters. By constraining different
123 second stage main effects or interaction effects to be 0, we can construct different \mathbf{Z}_G to
124 build different two-stage models.

125 Up to now, we have only described second stage decomposition for the regression coef-
126 ficients of \mathbf{G} . The second stage decomposition can also be applied to the other covariates,
127 the details of which are in Supplementary Section 1. We suggest not to perform second

128 stage decomposition on the intercepts parameters of the first stage polytomous model, i.e.,
129 the coefficients of intercepts are saturated, because decomposing the intercepts equates to
130 making assumptions on the prevalence of different cancer subtypes, which can potentially
131 lead to bias. Moving forward, we use \mathbf{Z}_X to denote the second stage design matrix for the
132 other covariates \mathbf{X} , $\boldsymbol{\lambda}$ to denote the second stage parameters for \mathbf{X} , and \mathbf{Z} to denote the
133 second stage design matrix for all the covariates.

134 **A. Hypothesis test under two-stage model**

135 The first stage case-control log odds ratios of subtypes can be decomposed into the second
136 stage case-control log odds ratio of the reference subtype, main effects and interaction effects
137 of tumor characteristics. This decomposition presents multiple options for comprehensively
138 testing for the association between a SNP and cancer subtypes. The first hypothesis test
139 is the global association test, $H_0^A : \boldsymbol{\theta} = [\boldsymbol{\theta}^{(0)} \quad \boldsymbol{\theta}_H^T]^T = [0 \quad \mathbf{0}^T]^T$ versus $H_1^A : \boldsymbol{\theta} \neq \mathbf{0}$, which
140 tests for an overall association between the SNP and the disease. Because $\boldsymbol{\theta} = \mathbf{0}$ implies
141 $\boldsymbol{\beta} = \mathbf{0}$, rejecting this null hypothesis means the SNP is associated with at least one of the
142 subtypes. The null hypothesis can be rejected if the SNP is significantly associated with a
143 similar effect size across all subtypes (i.e. $\boldsymbol{\theta}^{(0)} \neq 0$, $\boldsymbol{\theta}_H = \mathbf{0}$), or if the SNP has heterogeneous
144 effects on different subtypes ($\boldsymbol{\theta}_H \neq \mathbf{0}$).

145 The second hypothesis test is the global heterogeneity test, $H_0^{EH} : \boldsymbol{\theta}_H = \mathbf{0}$ versus $H_1^{EH} :$
146 $\boldsymbol{\theta}_H \neq \mathbf{0}$. This test simultaneously evaluates the etiologic heterogeneity with respect to a
147 SNP and all the tumor characteristics. Rejecting this null hypothesis indicates that the first

148 stage case-control log odds ratios are significantly different between at least two different
149 subtypes.

150 Notably, the global heterogeneity test does not identify which tumor characteristic(s)
151 is/are driving the heterogeneity. To identify the tumor characteristic(s) responsible for
152 observed heterogeneity, we propose the individual tumor marker heterogeneity test, $H_0^{\text{IH}} :$
153 $\theta_{\text{H}(k)} = 0$ versus $H_1^{\text{IH}} : \theta_{\text{H}(k)} \neq 0$, where $\theta_{\text{H}(k)}$ is one of the case-case parameters of $\boldsymbol{\theta}_{\text{H}}$. The
154 case-case parameter ($\theta_{\text{H}(k)}$) provides a measurement of etiological heterogeneity according
155 to a specific tumor characteristic (Begg and Zhang, 1994). In the breast cancer example, we
156 can directly test $H_0^{\text{IH}} : \theta_{\text{ER}}^{(1)} = 0$ versus $H_1^{\text{IH}} : \theta_{\text{ER}}^{(1)} \neq 0$. Rejecting the null hypothesis provides
157 evidence that the case-control log odds ratios of ER+ and ER- subtypes are significantly
158 different.

159 B. EM algorithm accounting for cases with incomplete tumor characteristics

160 In the previous sections, all the tumor characteristics were assumed to have no miss-
161 ing data. However, in epidemiological research, it is very common to have missing tumor
162 characteristics. This problem becomes exacerbated as the number of tumor characteristics
163 grows. Restricting to cases with complete tumor characteristics can reduce statistical power
164 and potentially introduce selection bias. To solve this problem, we propose to use the EM
165 algorithm (Dempster and others, 1977) to find the maximum likelihood estimate (MLE) of
166 the two-stage model, while incorporating all available information from the study. Let \mathbf{T}_{io}
167 be the observed tumor characteristics of subject i , and $Y_{im} = I(D_i = m)$ denote whether
168 the i th subject is disease subtype m . Given \mathbf{T}_{io} , the possible subtypes for subject i , denoted

169 as $\mathcal{Y}_{io} = \{Y_{im} : Y_{im} \text{ that is consistent with } \mathbf{T}_{io}\}$, are within a limited subset of all possible
 170 tumor subtypes. We assume that $(Y_{i1}, Y_{i2}, \dots, Y_{iM}, G_i, \mathbf{X}_i)$ are independently and identically
 171 distributed (i.i.d.), and that the tumor characteristics are missing at random (MAR). Let
 172 $\boldsymbol{\delta} = (\boldsymbol{\theta}^T, \boldsymbol{\lambda}^T)^T$ represent the second stage parameters of both \mathbf{G} and \mathbf{X} . Given the notation,
 173 the E step of them EM algorithm at the v th iteration is

174

$$Y_{im}^E = E(Y_{im}|G_i, \mathbf{X}_i, \mathbf{T}_{io}; \boldsymbol{\delta}^{(v)}) = \frac{Pr(Y_{im} = 1|G_i, \mathbf{X}_i; \boldsymbol{\delta}^{(v)})I(Y_{im} \in \mathcal{Y}_{io})}{\sum_{Y_{im} \in \mathcal{Y}_{io}} Pr(Y_{im} = 1|G_i, \mathbf{X}_i; \boldsymbol{\delta}^{(v)})}, \quad (3)$$

175 where Y_{im}^E is the probability of the i th person to be the m th subtype given his observed
 176 tumor characteristics (\mathbf{T}_{io}) , genotype (G_i) , and other covariates (\mathbf{X}_i) . $I(Y_{im} \in \mathcal{Y}_{io})$ denotes
 177 whether the m th subtype for the i th subject belong to the subsets of possible subtypes given
 178 the observed tumor characteristics. The M step at the v th iteration is

$$\boldsymbol{\delta}^{(v+1)} = \arg \max_{\boldsymbol{\delta}} \sum_{i=1}^N \left[\left(1 - \sum_{m=1}^M Y_{im}^E\right) \log Pr(D_i = 0|G_i, \mathbf{X}_i) + \sum_{m=1}^M Y_{im}^E \log \{Pr(D_i = m|G_i, \mathbf{X}_i)\} \right]. \quad (4)$$

179 The M step can be solved through a weighted least square iteration. Let $\mathbf{Y}_m = (Y_{1m}, \dots, Y_{Nm})^T$,
 180 and $\mathbf{Y} = (\mathbf{Y}_1^T, \dots, \mathbf{Y}_M^T)^T$. Let $\mathbf{C} = (\mathbf{G}, \mathbf{X})$, and $\mathbf{C}_M = \mathbf{I}_M \otimes \mathbf{C}$. Let $\mathbf{W} = \mathbf{D} - \mathbf{A}\mathbf{A}^T$,
 181 $\mathbf{D} = \text{diag}(\mathbf{P})$, $\mathbf{P} = E(\mathbf{Y}|\mathbf{C}; \boldsymbol{\delta})$, and $\mathbf{A} = \mathbf{D}(\mathbf{1}_M \otimes \mathbf{I}_N)$. During the t th iteration of the
 182 weighted least square, $\mathbf{Y}^{*(t)} = \mathbf{W}^{(t)}(\mathbf{Y}^E - \mathbf{P}^{(t)}) + \mathbf{C}_M \mathbf{Z} \boldsymbol{\delta}^{(t)}$, where $\mathbf{P}^{(t)}$ and $\mathbf{W}^{(t)}$ are re-
 183 spectively defined as \mathbf{P} and \mathbf{W} evaluated at the $\boldsymbol{\delta}^{(t)}$. The weighted least square update is
 184 $\boldsymbol{\delta}^{(t+1)} = (\mathbf{Z}^T \mathbf{C}_M^T \mathbf{W}^{(t)} \mathbf{C}_M \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{C}_M^T \mathbf{Y}^{*(t)}$. As $t \rightarrow \infty$, the weighted least square interaction
 185 converges to $\hat{\boldsymbol{\delta}}^{(v+1)}$, which will be used in next iteration. The EM algorithm will converge
 186 to the MLE of the second stage parameters (denoted as $\hat{\boldsymbol{\delta}}$), and the observed information
 187 matrix \mathbf{I} is $\mathbf{I} = \mathbf{Z}^T \mathbf{C}_M^T (\mathbf{W} - \mathbf{W}_{\text{mis}}) \mathbf{C}_M^T \mathbf{Z}$, where $\mathbf{W}_{\text{mis}} = \mathbf{D}_{\text{mis}} - \mathbf{A}_{\text{mis}} \mathbf{A}_{\text{mis}}^T$, $\mathbf{D}_{\text{mis}} = \text{diag}(\mathbf{P}_{\text{mis}})$,

188 $\mathbf{P}_{\text{mis}} = E(\mathbf{Y}|\mathbf{C}, \mathbf{T}_o; \delta)$, and $\mathbf{A}_{\text{mis}} = \mathbf{D}_{\text{mis}}(\mathbf{1}_M \otimes \mathbf{I}_N)$ (Louis, 1982). More details of the EM
189 algorithm are in Supplementary Section 2.

190 With the MLE of the second stage parameters of G as $\hat{\boldsymbol{\theta}}$, we can construct the the Wald
191 statistics as $\hat{\boldsymbol{\theta}}^{*T} \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\theta}}^* \sim \chi_l^2$ for the global association test, global etiological heterogeneity
192 test, and individual tumor characteristic heterogeneity test using the corresponding second
193 stage parameters and covariance matrix, where the degrees of freedom l equal the length of
194 $\hat{\boldsymbol{\theta}}^*$.

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

196 Although the hypothesis tests can be implemented through the Wald test, estimating
197 the model parameters for all SNPs in the genome is time-consuming and computationally
198 intensive. In this section, we develop a score test for the global association test assuming
199 the second stage parameters to be fixed. The score test only needs to estimate the second
200 stage parameters of \mathbf{X} under the null hypothesis once, making it much more computationally
201 efficient than the Wald test. Moreover, the EM algorithm only needs to be implemented
202 once under the null hypothesis. Since we don't perform any second stage decomposition on
203 the intercept parameters in the first stage polytomous model, the correlations between the
204 tumor characteristics are kept close to the empirical correlations for tumor markers. Most
205 of the imputation power is due to the high correlation between the tumor markers. In the
206 breast cancer example, the correlation between ER and PR is 0.63, between ER and HER2
207 is -0.16, and between PR and HER2 is -0.17 (Supplementary Table 1). Also, The association
208 of \mathbf{X} with the tumor markers can improve the power of the EM algorithm. Since a single

209 SNP \mathbf{G} usually has a small effect, the fact that the effect of individual \mathbf{G} is not incorporated
 210 in the EM algorithm itself doesn't result in much loss of efficiency.

211 Let $\mathbf{G}_M = \mathbf{I}_M \otimes \mathbf{G}$, and $\mathbf{X}_M = \mathbf{I}_M \otimes \mathbf{X}$. Under the null hypothesis, $H_0 : \boldsymbol{\theta} = \mathbf{0}$, let $\hat{\boldsymbol{\lambda}}$ denote
 212 the MLE of $\boldsymbol{\lambda}$ under the null hypothesis. The efficient score of $\boldsymbol{\theta}$ is $U_{\boldsymbol{\theta}}(\hat{\boldsymbol{\lambda}}) = \mathbf{Z}_G^T \mathbf{G}_M^T (\mathbf{Y} - \mathbf{P}_f)$,
 213 where $\mathbf{P}_f = E_{\boldsymbol{\theta}=\mathbf{0}}(\mathbf{Y}|\mathbf{X}; \hat{\boldsymbol{\lambda}})$. Let $\mathbf{W}_f = \mathbf{D}_f - \mathbf{A}_f \mathbf{A}_f^T$, with $\mathbf{P}_f = E_{\boldsymbol{\theta}=\mathbf{0}}(\mathbf{Y}|\mathbf{X}, \mathbf{T}_o; \hat{\boldsymbol{\lambda}})$, $\mathbf{P}_{f,\text{mis}} =$
 214 $E(\mathbf{Y}|\mathbf{X}, \mathbf{T}_o; \hat{\boldsymbol{\lambda}})$, $\mathbf{D}_f = \text{diag}(\mathbf{P}_f - \mathbf{P}_{f,\text{mis}})$ and $\mathbf{A}_f = \mathbf{D}_f(\mathbf{1}_M \otimes \mathbf{I}_N)$. The corresponding efficient
 215 information matrix of $U_{\boldsymbol{\theta}}(\hat{\boldsymbol{\lambda}})$ is

$$\tilde{\mathbf{I}} = \mathbf{I}_{\boldsymbol{\theta}\boldsymbol{\theta}} - \mathbf{I}_{\boldsymbol{\theta}\boldsymbol{\lambda}} \mathbf{I}_{\boldsymbol{\lambda}\boldsymbol{\lambda}}^{-1} \mathbf{I}_{\boldsymbol{\lambda}\boldsymbol{\theta}}, \quad (5)$$

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

217 The score test statistic $Q_{\boldsymbol{\theta}}$ for fixed-effect two stage model is

$$Q_{\boldsymbol{\theta}} = U_{\boldsymbol{\theta}}(\hat{\boldsymbol{\lambda}})^T \tilde{\mathbf{I}}^{-1} U_{\boldsymbol{\theta}}(\hat{\boldsymbol{\lambda}}) \sim \chi_t^2. \quad (6)$$

218 FTOP has the same degrees of freedoms and similar asymptotic power (Yi and Wang, 2011)
 219 as the Wald test. In GWAS which needs to perform millions of tests, FTOP can be first
 220 used to scan the whole genome with global association test, and then select the potential
 221 risk regions. In the selected risk regions, each SNP can be tested for global heterogeneity
 222 and individual tumor characteristic heterogeneity using Wald test.

223 **D. Mixed-effect two-stage polytomous model score test (MTOPT)**

224 The two-stage model decreases the degrees of freedom compared to the polytomous logis-
 225 tic regression. However, the power gains in the two-stage model can be reduced as additional
 226 tumor characteristics are added into the model. We further propose a mixed-effect two-stage

227 model by modeling some of the second stage case-case parameters as random effects. Let
228 $\mathbf{u} = (u_1, \dots, u_s)^T$, where each u_j follows an arbitrary distribution F with mean zero and
229 variance σ^2 . The mixed-effect second stage model links the first and second stage parameters
230 as follows:

$$\boldsymbol{\beta} = \mathbf{Z}_f \boldsymbol{\theta}_f + \mathbf{Z}_r \mathbf{u}, \quad (7)$$

231 where \mathbf{Z}_f is the second stage design matrix of fixed effect, \mathbf{Z}_r is the second stage design matrix
232 of random effect, and $\boldsymbol{\theta}_f$ are the fixed-effect second stage parameters. Let $\boldsymbol{\theta}_f = (\theta^{(0)}, \boldsymbol{\theta}_{\text{fH}}^T)^T$,
233 where $\theta^{(0)}$ is the case-control log odds ratio of the reference subtype, and $\boldsymbol{\theta}_{\text{fH}}$ are the fixed
234 case-case parameters. The baseline effect $\theta^{(0)}$ is always kept fixed, since it captures the
235 SNP's overall effect on all the cancer subtypes.

236 The fixed-effect parameters $\boldsymbol{\theta}_{\text{fH}}$ can be used for tumor characters with prior information
237 suggesting that they are a source of heterogeneity, and the random-effect parameters \mathbf{u}
238 can model tumor characteristics with little or no prior information. In the breast cancer
239 example, the baseline parameter ($\theta^{(0)}$) and the main effect of ER (θ_{fH}) can be modeled as
240 fixed effects, since previous evidence indicates ER as a source of breast cancer heterogeneity
241 ([García-Closas and others, 2013](#); [Milne and others, 2017](#)). The main effects of PR and HER2
242 and other potential interactions effects can be modeled as random effects (\mathbf{u}). In the mixed
243 effect two-stage model, the global association test is $H_0^A : \boldsymbol{\theta}_f = \mathbf{0}, \sigma^2 = 0$ versus $H_1^A : \boldsymbol{\theta}_f \neq$
244 $\mathbf{0}$ or $\sigma^2 \neq 0$, and the global etiology heterogeneity test is $H_0^{\text{EH}} : \boldsymbol{\theta}_{\text{fH}} = \mathbf{0}, \sigma^2 = 0$ versus $H_1^{\text{EH}} :$
245 $\boldsymbol{\theta}_{\text{fH}} \neq \mathbf{0}$ or $\sigma^2 \neq 0$.

246 To derive the score statistic for the global null $H_0^A : \boldsymbol{\theta}_f = \mathbf{0}, \sigma^2 = 0$, the common approach
247 is to take the partial derivatives of loglikelihood with respect to $\boldsymbol{\theta}_f$ and σ^2 respectively.

248 However, under the null hypothesis, the score for $\boldsymbol{\theta}_f$ follows a normal distribution, and
 249 for σ^2 follows a mixture of chi-square distribution (Supplementary Section 3). With the
 250 correlation between the two scores, getting the joint distribution between the two becomes
 251 very complicated. Inspired by methods for the rare variants testing ([Sun and others, 2013](#)),
 252 we propose to modify the derivations of score statistic so that two independent scores can
 253 be independent. First for $\boldsymbol{\theta}_f$, the score test statistic $Q_{\boldsymbol{\theta}_f}$ is derived under the global null
 254 hypothesis $H_0^A : \boldsymbol{\theta}_f = \mathbf{0}, \sigma^2 = 0$ as usual. But for σ^2 , the the score statistic Q_{σ^2} is derived
 255 under the null hypothesis $H_0 : \sigma^2 = 0$ without constraining $\boldsymbol{\theta}_f$. Through this procedure,
 256 the two score test statistics ($Q_{\boldsymbol{\theta}_f}$ and Q_{σ^2}) can be proved to be independent (Supplementary
 257 Section 4), and the Fisher's procedure ([Koziol and Perlman, 1978](#)) can be used to combine
 258 the p-value generated from the two independent tests. Similarly to FTOP, the EM algorithm
 259 under the null hypothesis of MTOP can efficiently handle the missing tumor marker problems
 260 given the high correlations between the tumor characteristics. However, since MTOP needs
 261 to estimate $\boldsymbol{\theta}_f$ under the null hypothesis $H_0 : \sigma^2 = 0$ for every single SNP, the computation
 262 speed for MTOP is slower than FTOP.

263 The score statistic of the fixed effect $\boldsymbol{\theta}_f$ under the global null $H_0^A : \boldsymbol{\theta}_f = \mathbf{0}, \sigma^2 = 0$ is

$$Q_{\boldsymbol{\theta}_f} = (\mathbf{Y} - \mathbf{P}_f)^T \mathbf{G}_M \mathbf{Z}_f \tilde{\mathbf{I}}_f^{-1} \mathbf{Z}_f^T \mathbf{G}_M^T (\mathbf{Y} - \mathbf{P}_f) \sim \chi_{l_f}^2, \quad (8)$$

264 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 5, but substi-
 265 tute \mathbf{Z}_G with \mathbf{Z}_f . Under the null hypothesis, $Q_{\boldsymbol{\theta}_f}$ follows a χ^2 distribution with the degrees
 266 of freedom l_f the same as the length of $\boldsymbol{\theta}_f$.

267 To explicitly express Q_{σ^2} , let $\boldsymbol{\tau} = (\boldsymbol{\theta}_f^T, \boldsymbol{\lambda}^T)^T$ be the second stage fixed effect, and \mathbf{Z}_τ is
 268 the corresponding second stage design matrix. The variance component score statistic of σ^2

269 under the null hypothesis $H_0 : \sigma^2 = 0$ without constraining θ_f is as follows:

$$Q_{\sigma^2} = (\mathbf{Y} - \mathbf{P}_r)^T \mathbf{G}_M \mathbf{Z}_r \mathbf{Z}_r^T \mathbf{G}_M^T (\mathbf{Y} - \mathbf{P}_r) \sim \sum_{i=1}^s \rho_i \chi_{i,1}^2, \quad (9)$$

270 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, $H_0 : \sigma^2 = 0$.

271 Under the null hypothesis, Q_{σ^2} follows a mixture of chi square distribution (Supplementary

272 Section 3), where $\chi_{i,1}^2$ i.i.d. follows χ_1^2 . (ρ_1, \dots, ρ_s) are the eigenvalues of $\tilde{\mathbf{I}}_r = \mathbf{I}_{\mathbf{uu}} - \mathbf{I}_{\mathbf{u}\boldsymbol{\tau}}^T \mathbf{I}_{\boldsymbol{\tau}\boldsymbol{\tau}}^{-1} \mathbf{I}_{\boldsymbol{\tau}\mathbf{u}}$,

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

274 where $\mathbf{W}_r = \mathbf{D}_r - \mathbf{A}_r \mathbf{A}_r^T$, with $\mathbf{P}_r = E_{\sigma^2=0}(\mathbf{Y}|\mathbf{G}, \mathbf{X}, \mathbf{T}_o; \hat{\boldsymbol{\tau}})$, $\mathbf{P}_{r,\text{mis}} = E(\mathbf{Y}|\mathbf{G}, \mathbf{X}, \mathbf{T}_o; \hat{\boldsymbol{\tau}})$,

275 $\mathbf{D}_r = \text{diag}(\mathbf{P}_r - \mathbf{P}_{r,\text{mis}})$ and $\mathbf{A}_r = \mathbf{D}_r (\mathbf{1}_M \otimes \mathbf{I}_N)$. The Davies exact method (Davies, 1980) is

276 used here to calculate the p-value of the mixture of chi square distribution.

277 Let $P_{\theta_f} = Pr(Q_{\theta_f} \geq \chi_{l_f}^2)$ and $P_{\sigma^2} = Pr(Q_{\sigma^2} \geq \sum_{i=1}^s \rho_i \chi_{i,1}^2)$ be the p-values of the two

278 independent score statistics. Under the null hypothesis $H_0^A : \theta_f = \mathbf{0}$, $\sigma^2 = 0$, following the

279 Fisher's procedure, $-2 \log(P_{\theta_f}) - 2 \log(P_{\sigma^2})$ follows χ_4^2 ; thus, the p-value of mixed effect

280 two-stage model under the null hypothesis is

$$P_{\text{mix}} = Pr \{ -2 \log(P_{\theta_f}) - 2 \log(P_{\sigma^2}) \geq \chi_4^2 \}. \quad (10)$$

281 The extension of the score statistics of the global etiology heterogeneity test, $H_0^{\text{EH}} : \theta_{\text{EH}} =$

282 $\mathbf{0}$, $\sigma^2 = 0$, can be computed following a similar procedure as the global association test.

283 III. SIMULATION EXPERIMENTS

284 Large scale simulations across a wide range of practical scenarios were conducted to

285 evaluate the type I error (Section III A), statistical power (Section III B), and computation

286 time (Supplementary Section 5) of the fixed-effect and mixed-effect two-stage models. Data

287 were simulated to mimic the PBCS. We simulated four tumor characteristics: ER (positive
288 vs. negative), PR (positive vs. negative), HER2 (positive vs. negative), and grade (ordinal
289 1, 2, 3), which collectively defined $2^3 \times 3 = 24$ breast cancer subtypes.

290 In each simulation, genotype data \mathbf{G} was simulated under the Hardy-Weinberg equilib-
291 rium with minor allele frequency (MAF) as 0.25. An additional covariate (\mathbf{X}) was simulated
292 following a standard normal distribution independent of \mathbf{G} . We simulated a multinomial
293 outcome with 25 groups, one for the control group, and the other 24 for different cancer
294 subtypes, using the polytomous logistic regression model as follows:

$$Pr(D_i = m|X_i) = \frac{\exp(\alpha_m + \beta_m G_i + 0.05 X_i)}{1 + \sum_{m=1}^M \exp(\alpha_m + \beta_m G_i + 0.05 X_i)}. \quad (11)$$

295 The effect of \mathbf{X} was set as 0.05 for all subtypes. Using the frequency of the breast cancer
296 subtypes from Breast Cancer Association Consortium (Supplementary Table 2) ([Michailidou](#)
297 [and others, 2017](#)), we computed the corresponding polytomous logistic regression intercept
298 parameters α_m . The case-control ratio was set around 1:1, and the proportions of ER+,
299 PR+ and HER2+ were 0.81, 0.68, and 0.17, respectively. The proportions of grade 1, 2,
300 and 3 were 0.20, 0.48, and 0.32. The missing tumor markers were selected randomly with
301 missing rates of 0.17, 0.25, 0.42, and 0.27 for ER, PR, HER2 and grade, respectively. Under
302 this simulation, approximately 70% cases had at least one missing tumor characteristic.

303 **A. Type I error**

304 We evaluated the type I error of the global association test, global heterogeneity test, and
305 individual tumor marker heterogeneity test under the global null hypothesis. The data were

306 generated by setting $\beta_m = 0$ in Equation 11, where none of the subtypes was associated
307 with genotypes. The total sample size n was set to be 5,000, 50,000 and 100,000. We
308 conducted 2.4×10^7 simulations to evaluate the type I error at $\alpha = 1.0 \times 10^{-4}$, 1.0×10^{-5} ,
309 and 1.0×10^{-6} level.

310 Both MTOP and FTOP were applied with an additive two-stage model by constraining
311 all the interaction terms as 0 in Equation 2. The subtype-specific case-control log ORs were
312 specified into the case-control log OR of a baseline disease subtype (ER- , PR- , HER2-,
313 grade 1) and the main effects associated with the four tumor markers. Furthermore, the
314 MTOP assumed the baseline and ER case-case parameter as fixed effects and the other case-
315 case parameters as random effects. The global association test and global heterogeneity test
316 were implemented using both MTOP and FTOP, but the individual tumor characteristic
317 heterogeneity test could only be implemented with FTOP. For MTOP and FTOP, we re-
318 moved all the subtypes with fewer than 10 cases to avoid potential nonconvergence of the
319 model.

320 Table I presents the estimated type I errors under the global null hypothesis. Both MTOP
321 and FTOP correctly control the type I error, especially for the larger sample sizes. FTOP
322 is conservative with 5,000 subjects, especially for $\alpha = 1.0 \times 10^{-6}$, however, the method is
323 still valid. The well-controlled type I error also shows that removing rare subtypes doesn't
324 bias the estimate, as further demonstrated by additional simulations that are presented in
325 Supplementary Section 6. In the later sections, we generally used the additive second stage
326 structure for both MTOP and FTOP unless otherwise specified.

327 **B. Statistical power**

328 We assessed the statistical power of the proposed methods using various simulation set-
329 tings with sample sizes as 25,000, 50,000, and 100,000. For each setting, we performed
330 2×10^5 simulations to evaluate the power at $\alpha = 5.0 \times 10^{-8}$ level.

331 **1. Global association test**

332 The data were simulated with three different scenarios: I. no heterogeneity between tumor
333 markers, II. heterogeneity according to one tumor marker, and III. heterogeneity according
334 to multiple tumor markers. The disease subtypes were generated through Equation 11.
335 Under scenario I, we set β_m as 0.08 for all the subtypes. For scenarios II and III, β_m was
336 simulated following the additive two-stage model. Under scenarios II, datasets were simu-
337 lated with only ER heterogeneity by setting the case-case parameter for ER as 0.08, and
338 all the other as 0. For scenario III, we simulated a scenario with heterogeneity according to
339 all 4 tumor markers by setting the baseline effect to be 0, the ER case-case parameter to
340 be 0.08, and all the other case-case parameters following a normal distribution with mean
341 0 and variance 4.0×10^{-4} . Under this scenario, all tumor characteristics contributed to the
342 subtype-specific heterogeneity. Moreover, to evaluate different methods under a larger num-
343 ber of tumor characteristics, additional simulations were conducted by adding two additional
344 binary tumor characteristics to the previous four tumor characteristic setting. This defined
345 $2^5 \times 3 = 96$ cancer subtypes. The two additional tumor characteristics were randomly se-
346 lected to be missing with 5% missing rate. Under this setting, around 77% of the cases have

347 at least one tumor characteristic missing. We compared the statistical power to detect the
348 overall association using FTOP, MTOP, standard logistic regression, FTOP with only com-
349 plete data, and polytomous logistic regression. For MTOP, FTOP and polytomous model,
350 we removed all the subtypes with fewer than 10 cases to avoid potential nonconvergence of
351 the model.

352 Overall, MTOP had robust power under all scenarios (Figure 1). Standard logistic re-
353 gression had the highest power when there was no subtype-specific heterogeneity (Scenario
354 I), but suffered from substantial power loss when heterogeneity existed between subtypes.
355 MTOP, followed by FTOP, consistently demonstrated the highest power among the five
356 methods when subtype-specific heterogeneity existed (scenarios II and III). The power gain
357 of MTOP over FTOP ranged from 2% to 49%. The power gain was small when there were
358 four tumor characteristics because the difference in the degrees of freedom between MTOP
359 and FTOP was small. However, with six tumor markers, the power gain of MTOP was
360 more apparent owing to the larger difference in the degrees of freedom between the models.
361 FTOP was the least efficient in scenarios with no or little heterogeneity, such as scenarios
362 I and II, but with increasing heterogeneity, such-as scenario III, the power of MTOP and
363 FTOP were more similar.

364 The simulation study also showed that the incorporation of cases with missing tumor
365 characteristics significantly increased the power of the methods (Figure 1). Under the four
366 tumor markers setting with around 70% incomplete cases, the power gain of FTOP incorpo-
367 rating the missing data algorithm was at least 200% compared to FTOP with only complete
368 data. As expected, under the six tumor markers setting, which resulted in more missing

369 tumor marker data, the power of FTOP with the missing data algorithm was once again
370 significantly higher than FTOP with only complete data. MTOP was the most powerful
371 method when heterogeneity across cancer subtypes was present. Additional power simula-
372 tions with 5,000 subjects are described in Supplementary Section 7.

373 The previous simulations mainly focused on the two-stage model with additive effects.
374 Additional simulations were also implemented with pairwise interactions in the model. We
375 simulated data with β_m following a second stage model that included main effects and
376 pairwise interactions as shown in Equation 2 with the case-case parameter for ER ($\theta_1^{(1)}$) as
377 0.08, the pairwise interaction effect between ER and HER2 ($\theta_{13}^{(2)}$) as 0.04, and all the other
378 parameters as 0. Four methods were evaluated including FTOP with/without pairwise
379 interactions and MTOP with/without pairwise interactions (baseline and ER fixed). FTOP
380 without interaction terms still had high power (Figure 2). However, FTOP with pairwise
381 interaction structure had limited power because of the incorporation of the interaction terms
382 as fixed effects. On the other hand, MTOP with/without pairwise interactions maintained
383 a high power even when there were underlying interaction effects.

384 2. *Global heterogeneity test*

385 Supplementary Figure 3 shows the simulation results for global heterogeneity tests under
386 similar simulation settings as global association tests. MTOP had the highest power when
387 there were heterogeneous associations across the subtypes.

388 **3. Individual tumor marker heterogeneity test**

389 We further evaluated the power of the individual tumor marker heterogeneity test. The
390 data were generated with four tumor characteristics with the ER case-case parameter ($\theta_1^{(1)}$)
391 as 0.08, and all other parameters as 0. ER was randomly selected to be missing with a
392 rate of 0.17, 0.30 and 0.50. We compared two different methods, FTOP with all four tumor
393 characteristics and the polytomous model. The polytomous model was set up to test each
394 marker at a time. In the polytomous model, we removed cases with missing data only on
395 the relevant tumor marker to avoid penalizing the power of the model by removing cases
396 that were missing tumor marker data on the other tumor markers. FTOP with all four
397 tumor characteristics had smaller power compared to the polytomous model in testing the
398 effect of ER (Supplementary Figure 4). Since FTOP included all four tumor characteristics,
399 and the tumor markers were highly correlated, the variability of underlying parameters was
400 larger. However, the type I errors of the polytomous model in testing PR, HER2 and grade
401 were inflated under this case (Supplementary Figure 5). Under this simulation, these three
402 markers had no effect. On the other hand, FTOP controlled the type I error of all the tests.

403 Overall, for the global test for association and the global test for heterogeneity, when there
404 was no heterogeneity, the standard logistic regression was the most powerful method. How-
405 ever, in the presence of subtype heterogeneity, MTOP was the most powerful method, and
406 MTOP had stable power even with a large number of pairwise interactions terms included.

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

408 We applied our proposed methods to the PBCS, a population-based breast cancer case-
409 control study conducted in Poland between 2000 and 2003 ([García-Closas and others, 2006](#)).
410 The study consisted of 2,078 cases of histologically or cytologically confirmed invasive breast
411 cancer and 2,219 women without a history of breast cancer at enrollment. Information on
412 ER, PR, and grade were available from pathology records ([García-Closas and others, 2006](#)),
413 and information on HER2 was available from immunohistochemical staining of tissue mi-
414 croarray blocks ([Yang and others, 2007](#)). We used genome-wide genotyping data to compare
415 MTOP, FTOP, standard logistic regression, and polytomous logistic regression to detect
416 SNPs associated with breast cancer risk.

417 Supplementary Table 4 presents the sample size of the tumor characteristics. The four
418 tumor characteristics defined 24 mutually exclusive breast cancer subtypes. Subtypes with
419 less than 10 cases were excluded, leaving 17 subtypes in the analysis. Both MTOP and FTOP
420 used the additive second stage design. Besides, we modeled the baseline and ER case-case
421 parameters as fixed effects in MTOP, and all other effects as random effects. We put ER as
422 a fixed effect because of the previously reported heterogeneity in genetic association by ER
423 ([García-Closas and others, 2013](#); [Milne and others, 2017](#)). Genotype imputation was done
424 using IMPUTE2 based on 1000 Genomes Project as reference ([Michailidou and others, 2017](#);
425 [Milne and others, 2017](#)). In total, 7,017,694 common variants on 22 auto chromosomes with
426 $MAF \geq 5\%$ were included in the analysis. In all the models, we adjusted for age and the
427 first four genetic principal components to account for population stratification.

428 As Figure 3 shows, MTOP, FTOP and standard logistic regression all identified a known
429 susceptibility variant in the FGFR2 locus on chromosome 10 (Michailidou *and others*, 2017),
430 with the most significant SNP being rs11200014 ($P < 5.0 \times 10^{-8}$). Further, both MTOP and
431 FTOP identified a second known susceptibility locus on chromosome 11 (CCND1) (Michaili-
432 dou *and others*, 2017), with the most significant SNP in both models being rs78540526 (P
433 $< 5.0 \times 10^{-8}$). The individual heterogeneity test of this SNP showed evidence for heterogene-
434 ity by ER ($P=0.011$) and grade ($P=0.024$). Notably, the CCND1 locus was not genome-wide
435 significant in standard logistic regression or polytomous models. The type I error of the four
436 methods was well-controlled (Supplementary Figure 6).

437 Additional sensitivity analysis of MTOP was implemented by specifying baseline, ER
438 and grade as fixed effects, and PR and HER2 as random effects (Supplementary Figure
439 7). The results for MTOP with grade as fixed vs. random effect were similar. We also
440 implemented MTOP and FTOP incorporating pairwise interactions in the second stage
441 model (Supplementary Figure 8-9). With pairwise interactions, both MTOP and FTOP
442 detected FGFR2 and CCND1 with the genome-wide significant threshold. However, the
443 P-value of FTOP with pairwise interactions was less significant compared to FTOP without
444 these interaction terms (for rs11200014, $P = 4.3 \times 10^{-8}$ vs. $P = 1.0 \times 10^{-9}$; for rs78540526,
445 $P = 2.7 \times 10^{-10}$ vs. $P = 8.1 \times 10^{-12}$). The P-value of MTOP with pairwise interactions
446 was also less significant compared to MTOP without interaction terms (for rs11200014,
447 $P = 1.0 \times 10^{-9}$ vs. $P = 2.2 \times 10^{-10}$; for rs78540526, $P = 1.7 \times 10^{-11}$ vs. $P = 1.8 \times 10^{-12}$). In
448 both scenarios with pairwise interactions parameter included, however, the power loss was
449 smaller.

450 Next, we compared the ability of MTOP and standard logistic regressions to detect 178
451 previously identified breast cancer susceptibility loci (*Michailidou and others, 2017*). For
452 eight of the 178 loci, the MTOP global association test p-value was more than ten fold lower
453 compared to the standard logistic regression p-value (Table II). In the MTOP model, these
454 eight loci all had significant global heterogeneity tests ($P < 0.05$). Confirming these results,
455 in a previous analysis applying MTOP to 106,571 breast cancer cases and 95,762 controls,
456 these eight loci were reported to have significant global heterogeneity (*Ahearn and others,*
457 *2019*).

458 V. DISCUSSION

459 We present a series of novel methods for performing genetic association testing for cancer
460 outcomes accounting for potential heterogeneity across subtypes. These methods efficiently
461 account for multiple testing, correlations between markers, and missing tumor data. Un-
462 der the model framework, we develop two computationally efficient score tests, FTOP and
463 MTOP, which model the underlying heterogeneity parameters in terms of fixed effects or
464 mixed effects, respectively. We demonstrate these methods have greater statistical power in
465 the presence of subtype heterogeneity than either standard or polytomous logistic regression
466 analysis.

467 Several methods have been proposed to study the etiological heterogeneity of cancer
468 subtypes (*Chatterjee, 2004; Rosner and others, 2013; Wang and others, 2015*). A recent
469 review showed the well-controlled type I error and good statistical power of the two-stage
470 model (*Zabor and Begg, 2017*). However, previous two-stage models haven't accounted for

471 missing tumor markers, which is a common problem in epidemiological studies. We show that
472 by incorporating the EM algorithm into the two-stage model we can take advantage of all
473 available information and substantially increase the statistical power (Figure 1). Moreover,
474 the newly proposed mixed effect model can mitigate the degrees of freedom penalty caused
475 by analyzing many tumor characteristics. In a recent large breast cancer GWAS analysis
476 with 106,571 cases and 95,762 controls, the newly developed methods MTOP and FTOP
477 have identified 16 novel loci (Zhang *and others*, 2019).

478 Incorporating missing tumor characteristics based on the proposed EM algorithm requires
479 the assumption of MAR, i.e. the mechanism of missing of the individual tumor characteris-
480 tics can depend only on other observed tumor characteristics and covariates, but not on the
481 unobserved missing value themselves. For the analysis of tumor heterogeneity, information
482 on aggressive types of tumors may be systematically missing. If the missing tumor char-
483 acteristics are important determinants of aggressiveness, then the underlying assumption
484 is violated. In general, dealing with non-ignorable missing data is a complex problem and
485 certain sensitivity analyses can be performed to explore the degree of bias (Little and Rubin,
486 2019). In the context of genetic association testing, non-ignorable missingness can lead to
487 inflated type I error only if the missingness mechanism itself is related to the genetic variant.
488 Further research is merited to explore the complex effects of non-ignorable missingness in
489 type I error and power of the proposed tests.

490 The computation time of MTOP is greater than FTOP (Supplementary Section 5). To
491 construct the score tests in FTOP, the coefficients of covariates need to be estimated once
492 under the null hypothesis, while in MTOP they need to be estimated for every SNP. The

493 computational complexity of FTOP is $O(NM^2P^2)$, with P as the number of other covariates
494 **X**. For MTOP, the computational complexity is $O(NM^2P^2lk)$, where l and k are respectively
495 the numbers of iteration required for weighted least square and EM algorithm to converge.

496 Currently, we only implement the linear kernel in MTOP, but other common kernels that
497 capture the similarity between tumor characteristics can be used in the future. If there is
498 prior knowledge about the overlapping genetic architecture across different tumor subtypes,
499 this will help to choose the kernel function, and improve the power of the methods.

500 The proposed methods have been implemented in a user-friendly and high-speed R statis-
501 tical package called TOP (<https://github.com/andrewhaoyu/TOP>), which includes all the
502 core functions implemented in C code.

503 **VI. SUPPLEMENTARY MATERIALS**

504 In Supplementary Section 1, the two-stage model is generalized to multivariates. In Sup-
505 plementary Section 2-3, the details of the EM algorithm and the variance component score
506 statistic are respectively presented. In Supplementary Section 4, Q_{θ_i} and Q_{σ^2} are proved to
507 be independent. In Supplementary Section 5, computation time simulations are presented.
508 In Supplementary Section 6, the simulations to evaluate the bias of the estimates are shown.
509 In Supplementary Section 7, simulations with 5,000 subjects are presented. Supplementary
510 Table 1 shows the correlations of ER, PR, HER2, and grade. Supplementary Table 2 presents
511 the frequencies of the joint distribution of ER, PR, HER2, and grade. Supplementary Table
512 3 shows the simulation results to evaluate bias. Supplementary Table 4 presents the sam-
513 ple size of tumor characteristics in PBCS. Supplementary Figure 1 shows the computation

514 time simulations results. Supplementary Figure 2 presents the power analysis of the global
515 association test with 5,000 subjects. Supplementary Figure 3 presents global heterogeneity
516 test simulation results. Supplementary 4-5 respectively present the power and type I error
517 simulations results of individual tumor marker heterogeneity test. Supplementary Figure 6
518 is the QQ plot of GWAS with PBCS. Supplementary Figure 7 shows the GWAS with PBCS
519 using MTOP with ER and grade as fixed effects. Supplementary Figure 8-9 respectively
520 present the GWAS with PBCS using MTOP/FTOP with pairwise interactions.

521 **ACKNOWLEDGEMENTS**

522 This work was supported by funds from the NCI Intramural Research Program, Bloomberg
523 Distinguished Professorship endowment, and NHGRI (1R01 HG010480-01). The simulation
524 experiments and data analysis were implemented using the high performance computation
525 Biowulf cluster at National Institutes of Health, USA.

526

527 AHEARN, T. U. *and others.* (2019). Common breast cancer risk loci predispose to distinct
528 tumor subtypes. *bioRxiv*, 733402.

529 BARNARD, M. E., BOEKE, C. E. AND TAMIMI, R. M. (2015). Established breast cancer
530 risk factors and risk of intrinsic tumor subtypes. *Biochimica et Biophysica Acta (BBA)-*
531 *Reviews on Cancer* **1856**(1), 73–85.

532 BEGG, C. B. AND ZHANG, Z. F. (1994). Statistical analysis of molecular epidemiology
533 studies employing case-series. *Cancer Epidemiology and Prevention Biomarkers* **3**(2), 173–

534 175.

535 CANCER GENOME ATLAS RESEARCH, NETWORK. (2014). Comprehensive molecular pro-
536 filing of lung adenocarcinoma. *Nature* **511**(7511), 543–50.

537 CHATTERJEE, N. (2004). A two-stage regression model for epidemiological studies with
538 multivariate disease classification data. *Journal of the American Statistical Associa-*
539 *tion* **99**(465), 127–138.

540 DAVIES, R. B. (1980). The distribution of a linear combination of χ^2 random variables.
541 *Journal of the Royal Statistical Society. Series C (Applied Statistics)* **29**(3), 323–333.

542 DEMPSTER, A. P., LAIRD, N. M. AND RUBIN, D. B. (1977). Maximum likelihood from
543 incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B*
544 *(Methodological)* **39**(1), 1–22.

545 DUBIN, N. AND PASTERNAK, B. S. (1986). Risk assessment for case-control subgroups by
546 polychotomous logistic regression. *American journal of epidemiology* **123**(6), 1101–1117.

547 GARCÍA-CLOSAS, M. *and others.* (2006). Established breast cancer risk factors by clinically
548 important tumour characteristics. *British journal of cancer* **95**(1), 123.

549 GARCÍA-CLOSAS, M. *and others.* (2013). Genome-wide association studies identify four er
550 negative-specific breast cancer risk loci. *Nature genetics* **45**(4), 392.

551 KOZIOL, J. A. AND PERLMAN, M. D. (1978). Combining independent chi-squared tests.
552 *Journal of the American Statistical Association* **73**(364), 753–763.

553 LIN, X. (1997). Variance component testing in generalised linear models with random
554 effects. *Biometrika* **84**(2), 309–326.

- 555 LITTLE, R. J. AND RUBIN, D. B. (2019). *Statistical analysis with missing data*, Volume
556 793. John Wiley & Sons.
- 557 LOUIS, T. A. (1982). Finding the observed information matrix when using the em algorithm.
558 *Journal of the Royal Statistical Society: Series B (Methodological)* **44**(2), 226–233.
- 559 MACARTHUR, J. *and others*. (2016). The new nhgri-ebi catalog of published genome-wide
560 association studies (gwas catalog). *Nucleic acids research* **45**(D1), D896–D901.
- 561 MICHAILED, K. *and others*. (2017). Association analysis identifies 65 new breast cancer
562 risk loci. *Nature* **551**(7678), 92–94.
- 563 MILNE, R. L. *and others*. (2017). Identification of ten variants associated with risk of
564 estrogen-receptor-negative breast cancer. *Nature genetics* **49**(12), 1767.
- 565 NETWORK, CANCER GENOME ATLAS. (2012). Comprehensive molecular portraits of hu-
566 man breast tumours. *Nature* **490**(7418), 61–70.
- 567 PEROU, C. M. *and others*. (2000). Molecular portraits of human breast tumours. *Na-*
568 *ture* **406**(6797), 747–52.
- 569 PRAT, A. *and others*. (2015). Clinical implications of the intrinsic molecular subtypes of
570 breast cancer. *The Breast* **24**, S26–S35.
- 571 ROSNER, B. *and others*. (2013). Breast cancer risk prediction with heterogeneous risk pro-
572 files according to breast cancer tumor markers. *American Journal of Epidemiology* **178**(2),
573 296–308.
- 574 SUN, J., ZHENG, Y. AND HSU, L. (2013). A Unified Mixed-Effects Model for Rare-Variant
575 Association in Sequencing Studies. *Genetic Epidemiology* **37**(4), 334–344.

- 576 WANG, M., KUCHIBA, A. AND OGINO, S. (2015). A meta-regression method for studying
577 etiological heterogeneity across disease subtypes classified by multiple biomarkers. *Amer-*
578 *ican journal of epidemiology* **182**(3), 263–270.
- 579 WU, M. C. *and others.* (2011). Rare-variant association testing for sequencing data with
580 the sequence kernel association test. *American Journal of Human Genetics* **89**(1), 82–93.
- 581 YANG, X. R. *and others.* (2007). Differences in risk factors for breast cancer molecular sub-
582 types in a population-based study. *Cancer Epidemiology and Prevention Biomarkers* **16**(3),
583 439–443.
- 584 YI, Y. AND WANG, X. (2011). Comparison of wald, score, and likelihood ratio tests for
585 response adaptive designs. *Journal of Statistical Theory and Applications* **10**(4), 553–569.
- 586 ZABOR, E. C. AND BEGG, C. B. (2017). A comparison of statistical methods for the study
587 of etiologic heterogeneity. *Statistics in medicine* **36**(25), 4050–4060.
- 588 ZHANG, D. AND LIN, X. (2003). Hypothesis testing in semiparametric additive mixed
589 models. *Biostatistics* **4**(1), 57–74.
- 590 ZHANG, H. *and others.* (2019). Genome-wide association study identifies 32 novel breast
591 cancer susceptibility loci from overall and subtype-specific analyses. *bioRxiv*, 778605.

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. 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 | 3.1×10^{-2} | 9.5×10^{-1} | 9.5×10^{-3} |
| rs10816625 | 9 | 110,837,073 | .06 | 5.0×10^{-2} | 9.8×10^{-1} | 2.2×10^{-2} |
| rs7904519 | 10 | 114,773,927 | .46 | 6.5×10^{-2} | 8.5×10^{-1} | 3.1×10^{-2} |
| rs554219 | 11 | 69,331,642 | .13 | 7.3×10^{-11} | 1.4×10^{-7} | 5.1×10^{-6} |
| rs11820646 | 11 | 129,461,171 | .40 | 1.5×10^{-2} | 8.6×10^{-1} | 4.5×10^{-3} |
| rs2236007 | 14 | 37,132,769 | .21 | 2.1×10^{-3} | 1.9×10^{-1} | 3.5×10^{-3} |
| rs1436904 | 18 | 24,570,667 | .40 | 7.2×10^{-4} | 6.6×10^{-2} | 9.7×10^{-4} |
| rs1436904 | 22 | 29,121,087 | .01 | 9.8×10^{-3} | 1.6×10^{-1} | 2.3×10^{-2} |

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

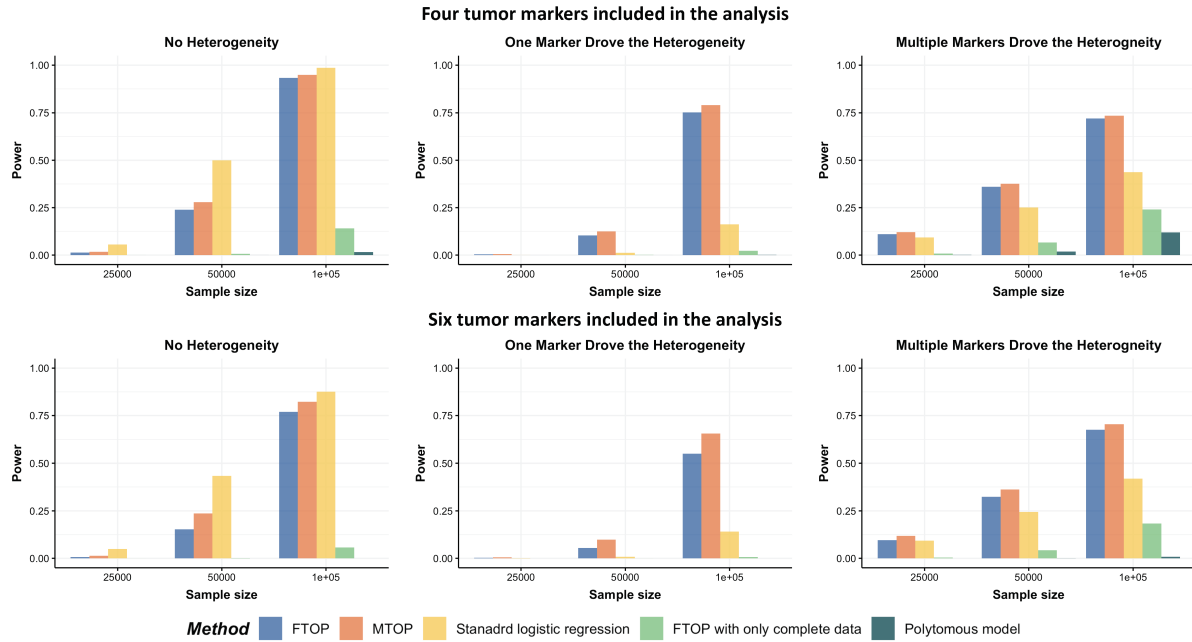


FIG. 1. Power comparison among MTOp, FTOP, standard logistic regression, two-stage model with only complete data and polytomous model with 2×10^5 random samples. For the three figures in the first row, four tumor markers were included in the analysis. Three binary tumor marker and one ordinal tumor marker defined 24 cancer subtypes. Around 70% cases would be incomplete. For the three figures in the second row, two extra binary tumor markers were included in the analysis. The six tumor markers defined 96 subtypes. Around 77% cases would be incomplete. The power was estimated by controlling the type one error $\alpha < 5.0 \times 10^{-8}$.

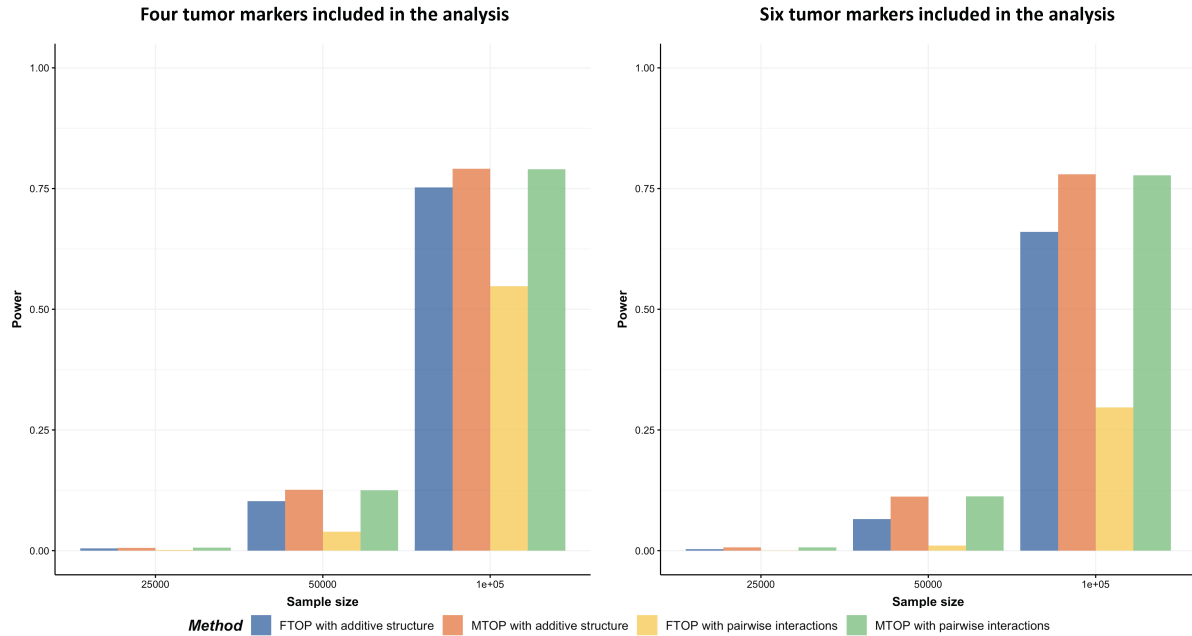


FIG. 2. Power comparison of global association test with pairwise interactions. Four methods were evaluated, including FTOP with additive structure, MTOP with additive structure (ER fixed), FTOP with pairwise interactions and MTOP with pairwise interactions (ER fixed). For the three figures in the first row, four tumor markers were included in the analysis. Three binary tumor marker and one ordinal tumor marker defined 24 cancer subtypes. Around 70% cases were incomplete. For the three figures in the second row, two extra binary tumor markers were included in the analysis. The six tumor markers defined 96 subtypes. Around 77% cases were incomplete. The total sample size was 25,000, 50,000 and 100,000. We generated 2×10^5 random replicates. The power was estimated by controlling the type one error $\alpha < 5.0 \times 10^{-8}$.

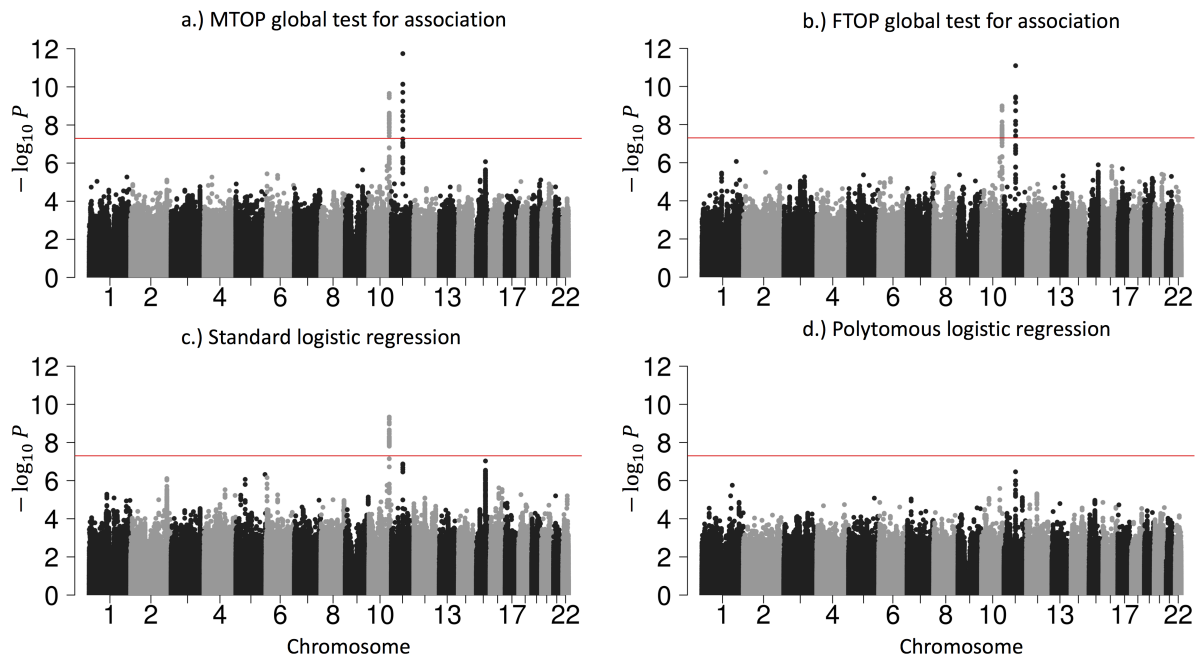


FIG. 3. Manhattan plot of genome-wide association analysis with PBCS using four different methods. PBCS have 2,078 invasive breast cancer and 2,219 controls. In total, 7,017,694 SNPs on 22 auto chromosomes with MAF more than 5% were included in the analysis. ER, PR, HER2 and grade were used to define breast cancer subtypes.