
A comparative study on the unified model based multifactor dimensionality reduction methods for identifying gene-gene interactions associated with the survival phenotype

Jung Wun Lee¹ and Seungyeoun Lee^{2*}

*Correspondence: leesy@sejong.ac.kr

¹Department of Statistics, University of Connecticut, Storrs, Connecticut, USA

²Department of Mathematics and Statistics, Sejong University, 209 Neungdong-ro, Gwangjin-gu, 05006, Seoul, Korea

E-mail addresses

JL: leejwegg@gmail.com; SL: leesy@sejong.ac.kr

1

2 Abstract

3 **Background:** For gene-gene interaction analysis, the multifactor dimensionality reduction (MDR) method has been widely employed to
4 reduce multi-levels of gene-gene interactions into high- or low-risk groups using a binary attribute. For the survival phenotype, the Cox-
5 MDR method has been proposed using a martingale residual of a Cox model since Surv-MDR was first proposed using a log-rank test sta-
6 tistic. Recently, the KM-MDR method was proposed using the Kaplan-Meier median survival time as a classifier. All three methods used
7 the cross-validation procedure to identify single nucleotide polymorphism (SNP) using SNP interactions among all possible SNP pairs.
8 Furthermore, these methods require the permutation test to verify the significance of the selected SNP pairs. However, the unified model-
9 based multifactor dimensionality reduction method (UM-MDR) overcomes this shortcoming of MDR by unifying the significance testing
10 with the MDR algorithm within the framework of the regression model. Neither cross-validation nor permutation testing is required to
11 identify SNP by SNP interactions in the UM-MDR method. The UM-MDR method comprises two steps: in the first step, multi-level geno-
12 types are classified into high- or low-risk groups, and an indicator variable for the high-risk group is defined. In the second step, the signifi-
13 cance of the indicator variable of the high-risk group is tested in the regression model included with other adjusting covariates. The Cox-
14 UMMDR method was recently proposed by combining Cox-MDR with UM-MDR to identify gene-gene interactions associated with the
15 survival phenotype. In this study, we propose two simple methods either by combining KM-MDR with UM-MDR, called KM-UMMDR
16 or by modifying Cox-UMMDR by adjusting for the covariate effect in step 1, rather than in step 2, a process called Cox2-UMMDR. The

17 KM-UMMDR method allows the covariate effect to be adjusted for in the regression model of step 2, although KM-MDR cannot adjust for
18 the covariate effect in the classification procedure of step 1. In contrast, Cox2-UMMDR differs from Cox-UMMDR in the sense that the
19 martingale residuals are obtained from a Cox model by adjusting for the covariate effect in step 1 of Cox2-UMMDR whereas Cox-
20 UMMDR adjusts for the covariate effect in the regression model in step 2. We performed simulation studies to compare the power of sev-
21 eral methods such as KM-UMMDR, Cox-UMMDR, Cox2-UMMDR, Cox-MDR, and KM-MDR by considering the effect of covariates
22 and the marginal effect of SNPs. We also analyzed a real example of Korean leukemia patient data for illustration and a short discussion is
23 provided.

24 **Results:** In the simulation study, two different scenarios are considered: the first scenario compares the power of the cases with and with-
25 out the covariate effect. The second scenario is to compare the power of cases with the main effect of SNPs versus without the main effect
26 of SNPs. From the simulation results, Cox-UMMDR performs the best across all scenarios among KM-UMMDR, Cox2-UMMDR, Cox-
27 MDR and KM-MDR. As expected, both Cox-UMMDR and Cox-MDR perform better than KM-UMMDR and KM-MDR when a covariate
28 effect exists because the former adjusts for the covariate effect but the latter cannot. However, Cox2-UMMDR behaves similarly to KM-
29 UMMDR and KM-MDR even though there is a covariate effect. This implies that the covariate effect would be more efficiently adjusted
30 for in the regression model of the second step rather than under the classification procedure of the first step. When there is a main effect of
31 any SNP, Cox-UMMDR, Cox2-UMMDR and KM-UMMDR perform better than Cox-MDR and KM-MDR if the main effects of SNPs are
32 properly adjusted for in the regression model. From the simulation results of two different scenarios, Cox-UMMDR seems to be the most
33 robust when there is either any covariate effect adjusting for or any SNP that has a main effect on the survival phenotype. In addition, the
34 power of all methods decreased as the censoring fraction increased from 0.1 to 0.3, but increased as heritability increased. The power of all
35 methods seems to be greater under $MAF = 0.2$ than under $MAF = 0.4$. For illustration, both KM-UMMDR and Cox2-UMMDR were ap-
36 plied to identify SNP by SNP interactions with the survival phenotype to a real dataset of Korean leukemia patients.

37 **Conclusion:** Both KM-UMMDR and Cox2-UMMDR were easily implemented by combining KM-MDR and Cox-MDR with UM-MDR,
38 respectively, to detect significant gene-gene interactions associated with survival time without cross-validation and permutation testing.
39 The intensive simulation results demonstrate the utility of KM-UMMDR, Cox2-UMMDR and Cox-UMMDR, which outperforms Cox-
40 MDR and KM-MDR when some SNPs with only marginal effects might mask the detection of causal epistasis. In addition, Cox-UMMDR,
41 Cox2-UMMDR and Cox-MDR performed better than KM-UMMDR and KM-MDR when there were potentially confounding covariate
42 effects.

43 **Keywords:** Survival time, Cox model, multifactor dimensionality reduction method, gene-gene interaction, unified model based method,
44 Kaplan-Meier estimate.

45

46 Introduction

47 With the advent of high-throughput genotyping techniques, a large amount of genotype data has been analyzed in genome-wide asso-
48 ciation studies. Among these, single nucleotide polymorphisms (SNPs) can modify many phenotypes, including cancer progression, re-
49 sponses to varying levels of drugs and survival outcomes. Numerous statistical methods for genome-wide association studies (GWAS)
50 have been developed to identify susceptibility genes by analyzing these data for single SNP effects. Since the first published GWAS on
51 age-related macular degeneration [1], the GWAS Catalog now contains 179,364 SNP-trait associations with 120,219 SNPs based on 4,493
52 publications since March 2020 (www.ebi.ac.uk/gwas). However, the problem of missing heritability since only a small proportion of her-
53 itability has been explained for the common and complex human diseases [2, 3, 4].

54 For the missing heritability problem, many researchers have focused on the challenge of identifying SNP by SNP interactions because
55 complex diseases might be associated with multiple genes and their interactions [3]. Most parametric statistical methods such as logistic
56 regression and ordinary regression models, have difficulty dealing with high-dimensional data because the number of variables exponen-
57 tially increases with higher-order SNP by SNP interactions. One solution to this problem is to collect a large number of samples which
58 yields a robust estimate of interaction effects. However, this is often far from reality because of the expensive sampling cost. An alterna-
59 tive solution is reducing the high dimensionality of multi-genotypes to a very low level, such as one dimension. The multifactor dimen-
60 sionality reduction (MDR) method was proposed by Ritchie *et al.* [5] as a new statistical and computational method to analyze gene-gene
61 interactions in genetic studies. The main principle of MDR is to reduce multi-dimensional genotypes into one-dimensional binary attrib-
62 utes, in which multi-level genotypes of SNPs are classified into either high- or low-risk groups, using a ratio of cases and controls in case-
63 control studies. The MDR algorithm then determines the best pair of SNPs among all possible SNP combinations, yielding the maximum
64 balanced accuracy through a cross-validation procedure. The MDR mechanism can apply higher-order interactions such as two-way, and
65 three-way interactions, because all combinations of multi-way interactions can be reduced to either high or low risk groups using the ap-
66 propriate classification rules.

67 The key algorithm of MDR has been extensively applied to quantitative traits and survival phenotypes because it was originally pro-
68 posed for case-control studies. For a prospective cohort study, the first approach, called Surv-MDR, was proposed by Gui *et al.* [6] fol-
69 lowed by Cox-MDR [7], AFT-MDR [8], and KM-MDR [9]. All these methods follow the same procedure as in the original MDR except
70 for using different classification rules, which are appropriate for the survival phenotype. For example, Surv-MDR uses a log-rank test
71 statistic, whereas Cox-MDR uses a martingale residual of the Cox model, AFT-MDR uses a standardized residual of an accelerated failure
72 time (AFT) model and KM-MDR uses the Kaplan-Meier median survival time. In previous simulation studies [7,8,9], in which the powers
73 of these methods were compared, both Surv-MDR and KM-MDR performed better than both Cox-MDR and AFT-MDR when there was no
74 covariate effect whereas both Cox-MDR and AFT-MDR had superior power than both Surv-MDR and KM-MDR when any covariate ef-
75 fect existed. This is because both Surv-MDR and KM-MDR are nonparametric and cannot adjust for any confounding covariate effect
76 whereas both Cox-MDR and AFT-MDR can adjust for confounding covariates in the frame of the regression model.

77 However, all these methods require a cross-validation procedure to identify the best SNP pairs among all possible combinations of
78 SNPs and computationally intensive permutation testing to check the statistical significance for the identified SNP pairs as performed in the
79 original MDR method. To overcome this shortcoming of MDR, the unified model-based multifactor dimensionality reduction (UM-MDR)
80 method was proposed by unifying significance testing with the MDR algorithm in the frame of the regression model [10]. In the UM-MDR
81 method, multi-level genotypes are classified into high- and low-risk groups, and an indicator variable for the high-risk group is defined in
82 the first step. Then, significance testing is unified in the frame of the regression model by including this indicator variable for the high-risk
83 group as one of covariates with other adjusting covariates. One of the advantages of UM-MDR is that it allows different types of classifica-
84 tion rules in the first step. Thus, a simple approach, called the Cox-UMMDR, was recently proposed by plugging the Cox-MDR into the
85 first step and combining the classified indicator with the significance testing procedure of the UM-MDR.

86 In this study, we propose two different simple methods, called KM-UMMDR and Cox2-UMMDR. The KM-UMMDR method uses the
87 KM-MDR algorithm in the first step and combines it with the second step of UM-MDR. The Cox2-UMMDR modifies the classification
88 step of Cox-UMMDR by allowing the covariate effect to be adjusted in the first step and fitting only one indicator variable for the high-risk
89 group in the second step whereas the adjusting covariates are considered in the regression model in the second step for Cox-UMMDR.
90 Throughout the intensive simulation study, the powers of the five methods— KM-UMMDR, Cox2-UMMDR, Cox-UMMDR, Cox-MDR
91 and KM-MDR— were compared under the two different scenarios mainly focusing on two-way interactions. One scenario compares the
92 power of semi-parametric methods such as Cox-UMMDR, Cox2-UMMDR and Cox-MDR against that of nonparametric methods such as
93 KM-UMMDR and KM-MDR by considering cases with and without a covariate effect. The other scenario compares the power of KM-
94 UMMDR, Cox-UMMDR and Cox2-UMMDR against that of KM-MDR and Cox-MDR by considering cases with and without the main
95 effect of SNPs. In addition, one simulation result is given under a three-way interaction model. In addition, both KM-UMMDR and Cox2-
96 UMMDR were applied to a real dataset of Korean leukemia patients and a short discussion is provided.

97

98 **Methods**

99 As described in Cox-UMMDR, there are two-step procedures for KM-UMMDR and Cox2-UMMDR. In the first step of KM-UMMDR,
100 we classify the multi-genotypes into high- or low-risk groups using the Kaplan-Meier median survival time as in KM-MDR. For the two-
101 way interaction model, all individuals are divided into nine groups with the same genotypes. We then compare the median survival time of
102 each cell with the overall median survival time. If the median survival time of each cell is less than the overall median survival time, then
103 the corresponding cell is classified as high-risk group. Otherwise, it is classified as low-risk group. Once all cells are classified as high(H)
104 or low(L) risk groups, an H/L binary indicator, \mathcal{S} , is defined. In contrast, in the first step of Cox2-UMMDR, we classify the multi-
105 genotypes into high- or low-risk groups using the sign of the sum of martingale residuals within each cell, where the martingale residuals
106 are obtained from a reduced Cox model with covariates such as age, sex and other confounding variables. Cox2-UMMDR differs from
107 Cox-UMMDR in the sense that the covariate effects are adjusted for in step 1 and then only an indicator variable, \mathcal{S} , is fitted in the regres-

108 sion model of step 2. By contrast, Cox-UMMDR does not adjust for the covariate effects in step 1 and fits a Cox model with all adjusting
109 covariates and an H/L indicator, \mathcal{S} , in step 2.

110 Once an indicator, \mathcal{S} , for the high-risk group is defined from the classification procedure in step 1, a similar procedure is implement-
111 ed in step 2 for both KM-UMMDR and Cox2-UMMDR. In step 2, we fit the following two Cox models

$$\lambda(t|S, Z) = \lambda_0(t) \exp(\alpha S + \gamma Z)$$

112 for KM-UMMDR and

$$\lambda(t|S) = \lambda_0(t) \exp(\alpha S)$$

113 for Cox2-UMMDR.

114 Here, $\lambda_0(t)$ is a baseline hazard function, S is an indicator variable for the high-risk group, Z is the vector coding for the adjusting covari-
115 ates, and α and γ are the corresponding parameters to S and Z , respectively. Because the median survival time is used as a classifier, the
116 covariate effect cannot be adjusted for in step 1 of KM-UMMDR. However, KM-UMMDR adjusted for the covariate effect in step 2 by
117 fitting the Cox model shown above, whereas KM-MDR cannot adjust for the covariate effect. In contrast, Cox2-UMMDR adjusted for the
118 covariate effect in step 1 and fit a Cox model with only an indicator for the high-risk group shown above.

119 We tested the following null hypothesis: $H_0: \alpha = 0$, i.e., whether the corresponding multi-locus is associated with the survival pheno-
120 type after adjusting for the covariate effect. To test the significance of the multi-locus model, a Wald type test statistic, $W = \hat{\alpha}^2 / \widehat{Var}(\hat{\alpha})$,
121 was used, however, its asymptotic distribution followed a non-central chi-square distribution with one degree of freedom and the non-
122 centrality parameter, q , as described in [10]. This is because S is defined as an indicator for the high-risk group after classification is per-
123 formed in step 1. Because the mean of the non-central chi-square distribution is $q + 1$, we first estimate the mean, $\hat{\mu}$, of W under the null
124 distribution and take $\hat{q} = \max(0, \hat{\mu} - 1)$. As described in [10], we permute the trait a few times, for example, 5 or 10 times, and take the
125 sample mean of the statistic W as $\hat{\mu}$. Here we can estimate the non-centrality parameter for each multi-locus model, or we can pool all the
126 statistics and then estimate the common non-centrality parameter for all multi-locus models.

127

128 **Simulation study**

129 Through intensive simulation studies, we compared the power of the five methods— KM-UMMDR, Cox2-UMMDR, Cox-UMMDR, Cox-
130 MDR and KM-MDR—under the two different scenarios. Scenario I compares the power of nonparametric methods with that of semi-
131 parametric methods with and without adjusting for the covariate effect. As described, KM-UMMDR and KM-MDR are nonparametric
132 methods in the sense that the classification procedure is performed by the Kaplan-Meier median survival time in step 1, whereas Cox-MDR,
133 Cox-UMDMR and Cox2-UMMDR are semi-parametric methods because they classify the martingale residuals from a Cox model with
134 adjustment for the covariate effect. However, KM-UMMDR can adjust for covariate effects in the regression model in step 2, and can be
135 regarded as a more improved approach than KM-MDR. Scenario II compares the power of original MDR methods, such as Cox-MDR and

136 KM-MDR, with that of the unified model-based MDR methods, such as KM-UMMDR, Cox2-UMMDR and Cox-UMMDR without and
 137 with adjustment for marginal SNP effects.

138 First we focus on a two-way interaction model and consider two disease-causal SNPs among 10 unlinked diallelic loci with the as-
 139 sumption of Hardy-Weinberg equilibrium and linkage equilibrium. For covariate adjustment, we consider only one covariate associated
 140 with survival time but with no interactions with any SNPs. We generated simulation datasets from different penetrance functions [12],
 141 which define a probabilistic relationship between the high- or low-risk status and SNPs. We then considered 14 different combinations of
 142 two different minor allele frequencies of (0.2, 0.4) and seven different heritability values of (0.01, 0.025, 0.05, 0.1, 0.2, 0.3, and 0.4). For
 143 each of the 14 heritability and minor allele frequency combinations, a total of five models were generated, yielding 70 epistasis models
 144 with various penetrance functions, as described in [12] (Supplemental Table 1).

145 Let f_{ij} be an element from the i^{th} row and the j^{th} column of a penetrance function. Assuming that SNP1 and SNP2 are two disease-
 146 causal SNPs, we obtain the following penetrance function:

$$f_{ij} = P(\text{high risk} | SNP_1 = i, SNP_2 = j)$$

147 We generated data of 400 patients from each of the 70 penetrance models to create one simulated dataset and repeated this procedure
 148 100 times. We generated the survival time from a Cox model, which can be specified as follows:

$$\lambda(t|x, z) = \lambda_0(t) \exp(\alpha x + \gamma z)$$

149 Here x is an indicator variable with value 1 for a high-risk group and 0 for a low-risk group. We set $\alpha = 0.8$, $\gamma = 0.8$ and z as the ad-
 150 justing covariate generated from $N(0, 1)$. In addition, the baseline hazard function follows a Weibull distribution with a shape parameter of
 151 5 and a scale parameter of 2, the censoring time being generated from a uniform distribution, $U(0, c)$ depending on the censoring fraction
 152 (0.1, 0.3).

153 First, we calculated the Type I error to check the validity of the Cox2-UMMDR and KM-UMMDR methods for which data were gen-
 154 erated under the null hypothesis where $H_0 : \alpha = 0$, across various heritability and minor allele frequencies (MAF). We generated 1000
 155 null dataset with eight non-causal SNPs and the Type I error of the selection rate of each SNP pair under the null model became $\frac{1}{\binom{8}{2}} =$
 156 0.0357. We considered five different MAFs as 0.05, 0.1, 0.2, 0.3, and 0.4 and four different censoring fractions (CF) were 0.0, 0.1, 0.3,
 157 and 0.5. Table 1 shows the Type I error of the Cox2-UMMDR and KM-UMMDR methods across various combinations of MAF and CF.
 158 Although most cases showed a conservative trend, it is concluded that the Type I error is well controlled with being less than 0.035.

Table 1: Type I error of Cox2-UMMDR and KM-UMMDR for the combinations of MAF and CF.

MAF	Cox2-UMMDR				KM-UMMDR			
	CF = 0.0	CF = 0.1	CF = 0.3	CF = 0.5	CF = 0.0	CF = 0.1	CF = 0.3	CF = 0.5
0.05	0.025	0.021	0.022	0.019	0.029	0.021	0.026	0.024
0.1	0.023	0.024	0.022	0.021	0.027	0.022	0.024	0.027
0.2	0.024	0.022	0.023	0.029	0.023	0.023	0.027	0.029
0.3	0.024	0.023	0.028	0.027	0.023	0.027	0.028	0.028
0.4	0.026	0.023	0.027	0.021	0.024	0.026	0.024	0.031

162

163 For the power comparison, we consider two different scenarios for the simulation study. In Scenario I, the survival times are gener-
164 ated from the Cox model as follows:

$$\lambda(t|x, z) = \lambda_0(t)\exp(\alpha x + \gamma z)$$

165 where $\alpha = 0.8, \gamma = 0.0$ and $\alpha = 0.8, \gamma = 0.8$, respectively.

166 In scenario II, the survival times were generated from the Cox model as follows:

$$\lambda(t|x, z, SNP3) = \lambda_0(t)\exp(\alpha x + \gamma z + \delta SNP3)$$

167 where $\alpha = 0.8, \gamma = 0.8, \delta = 0.5$ and $\alpha = 0.8, \gamma = 0.0, \delta = 0.5$. Here, δ denotes the marginal main effect of *SNP3* on the hazard
168 rate.

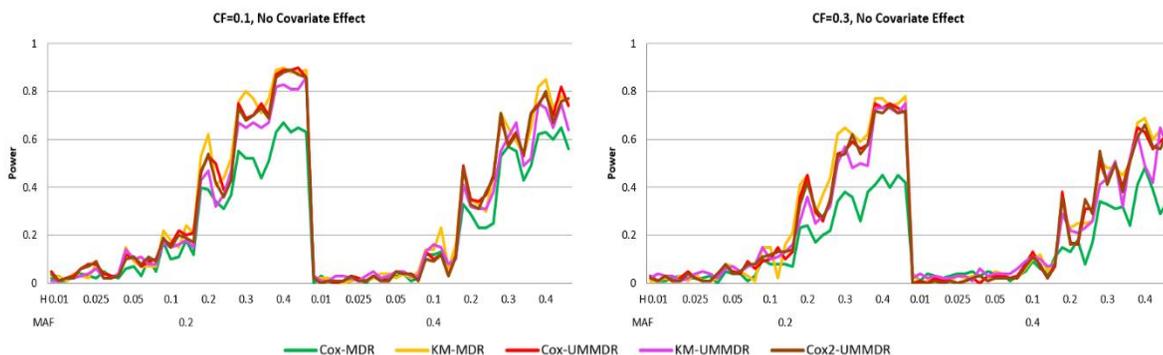
169 In addition, we consider a three-way interaction model to compare the power of the five methods, similar to the two-way interaction
170 model. For the three-way interaction model, we use penetrance functions, $P(\text{high risk}|AAbbCc) = 0.2, P(\text{high risk}|AaBbcc) =$
171 $0.2, P(\text{high risk}|aaBBcc) = 0.2, P(\text{high risk}|aaBbCc) = 0.2, P(\text{high risk}|AabbCc) = 0.2, P(\text{high risk}|aabbCC) =$
172 0.2 as suggested by Ritchie et al. [5]. Here, we assume that *SNP1*, *SNP2* and *SNP3* are the three disease-causal SNPs and each has dial-
173 lelic locus of (**A, a**), (**B, b**), (**C, c**). Similar to the two-way interaction models, we generated data of 400 patients from the penetrance
174 model to create one simulated dataset and repeated this procedure 100 times. The survival time was generated from a Cox model specified
175 as follows:

$$\lambda(t|x, z, SNP4) = \lambda_0(t)\exp(\alpha x + \gamma z + \delta SNP4)$$

176 Here, α was set to be 1.2, and γ and δ had four different combinations (i) $[\gamma, \delta] = [1.2, 0.5]$, (ii) $[\gamma, \delta] = [1.2, 0.0]$, (iii) $[\gamma, \delta] = [0.0,$
177 $0.5]$ and (iv) $[\gamma, \delta] = [0.0, 0.0]$

178 First, we consider the simulation result for the two-way interaction model under Scenario I. Figure 1 shows the power curve of the five
179 methods under $\alpha = 0.8, \gamma = 0.0$ where there is no covariate effect and the CFs are 0.1 and 0.3. The x-axis represents the heritability,
180 which is $H=(0.01, 0.025, 0.05, 0.1, 0.2, 0.3, 0.4)$ and $MAF=(0.2, 0.4)$, whereas the y-axis represents the power of the five methods across
181 70 different penetrance models [11].

182

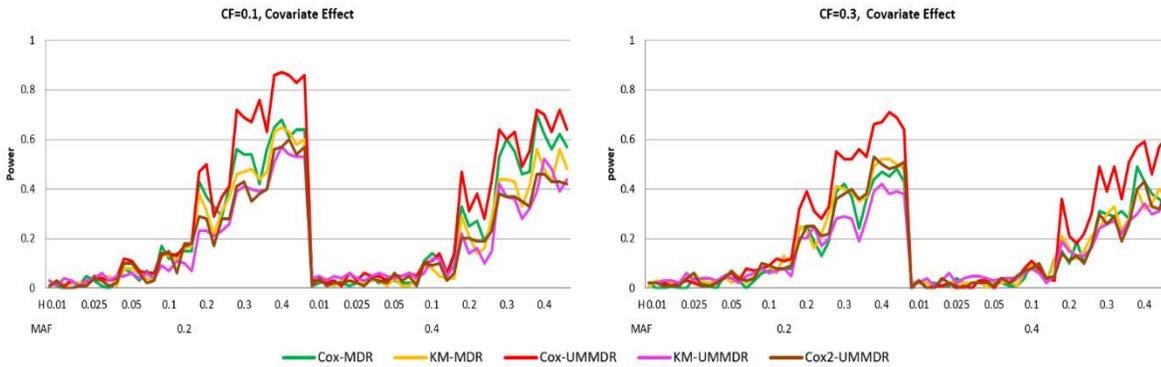


195
196
197

Figure 1: Power curves of the five methods under $\alpha=0.8, \gamma=0.0$

198 As shown in Figure 1, the power decreases as the CF increases whereas the power is greater when MAF=0.2 than when MAF=0.4. Under
199 the no covariate effect model, KM-UMMDR, KM-MDR, Cox-UMMDR, and Cox2-UMMDR performed similarly, and their powers in-
200 creased as heritability increased, except for Cox-MDR. It is expected that nonparametric approaches, such as KM-UMMDR and KM-
201 MDR, are more powerful than semi-parametric approaches, such as Cox-UMMDR, Cox2-UMMDR, and Cox-MDR. However, Cox-
202 UMMDR and Cox2-UMMDR seem to be very close to KM-UMMDR and KM-MDR in terms of the power curve.

203



210

217
218
219

Figure 2: Power curves of the five methods under $\alpha=0.8, \gamma=0.8$

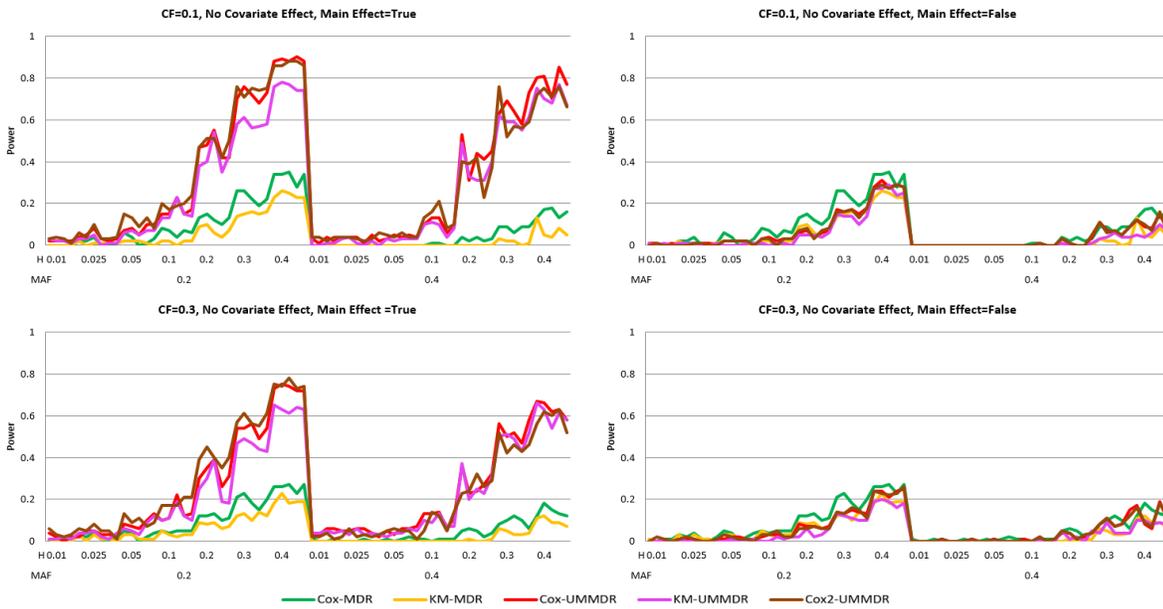
220 Figure 2 shows the power curve of the five methods under $\alpha = 0.8, \gamma = 0.8$, where there is a confounding covariate effect. As
221 shown in Figure 2, Cox-UMMDR has the highest power, however, the performance of Cox-MDR is similar to those of KM-UMMDR,
222 KM-MDR, and Cox2-UMMDR, except when MAF = 0.4, and the CF is 0.1. It is interesting that Cox2-UMMDR behaves almost the same
223 as KM-UMMDR, although Cox2-UMMDR adjusted for the covariate effects in step 1. This implies that the effect of covariates should be
224 adjusted for in step 2 rather than in step 1 because Cox-UMMDR is more powerful than Cox2-UMMDR. This is because the difference
225 between Cox2-UMMDR and Cox-UMMDR depends on the adjustment of the effect of the covariates.

226 Second, we consider the simulation result for the two-way interaction model under Scenario II. Because a non-causal SNP effect exists
227 in a true model, we compare two different cases where the main effect of SNPs is adjusted for (main = True) and not adjusted for (main =
228 False). Figure 3 shows the power curve of the five methods under $\alpha = 0.8, \gamma = 0.0, \delta = 0.5$, where there is no covariate effect, but
229 there exists a non-causal SNP main effect. As shown in Figure 3, the powers of all methods were very low when the main effect of SNPs is
230 not adjusted for regardless of any combination of MAF, heritability, and CF. However, when the main effect of SNPs is adjusted for, the
231 powers of Cox-UMMDR, Cox2-UMMDR and KM-UMMDR are relatively higher than those of Cox-MDR and KM-MDR. This implies
232 that adjusting for the main effects of causal SNPs helps to identify the interaction effect of these SNPs although there exists a non-causal

233 SNP main effect. In summary, three methods based on UM-MDR are more powerful than two methods based on the original MDR method
 234 when there exists an SNP main effect and it is properly adjusted for.

235 Figure 4 shows the power curve of five methods under $\alpha = 0.8$, $\gamma = 0.8$, $\delta = 0.5$, where there exists a non-causal SNP main effect
 236 as well as a confounding covariate effect. As shown in Figure 3, the power of all methods is very low when the main effect of an SNP is
 237 not adjusted for (main = False). However, when the main effect of an SNP is adjusted for (main = True), the power of Cox-UMMDR is
 238 substantially greater than that of Cox2-UMMDR and KM-UMMDR. As mentioned previously, it is more efficient that the effect of con-
 239 founding covariates should be adjusted for in step 2 than in step 1. Both Cox2-UMMDR and KM-UMMDR were very similar when the CF
 240 was 0.1, whereas Cox2-UMMDR had a relatively greater power than KM-UMMDR when the CF was 0.3. Similar to Figure 3, the powers
 241 of Cox-MDR and KM-MDR were noticeably lower than those of Cox-UMMDR, Cox2-UMMDR, and KM-UMMDR.

242



205

264
 265
 266
 267

Figure 3: Power curves of the five methods under $\alpha=0.8$, $\gamma=0.0$, $\delta=0.5$

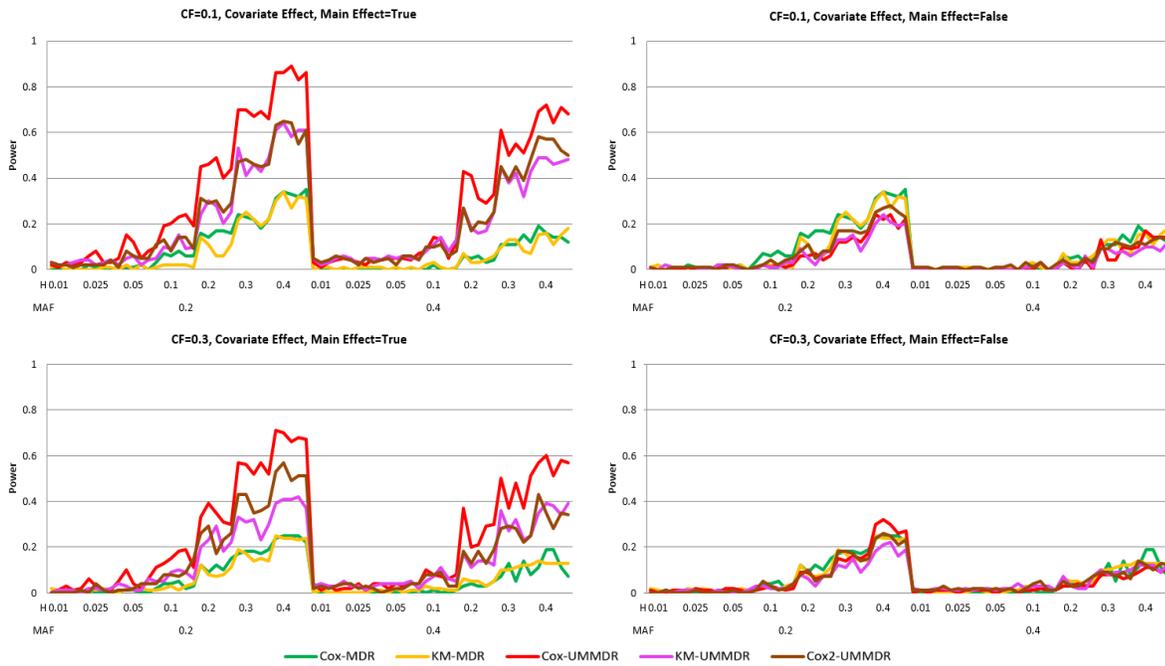


Figure 4 : Power curves of the five methods under $\alpha=0.8, \gamma=0.8, \delta=0.5$

Finally, we consider a three-way interaction model that has only one type of penetrance function: $P(\text{high risk}|AAbbcc) = 0.2$, $P(\text{high risk}|aabbCC) = 0.2$, as suggested by Ritchie et al. [5]. Similar to the simulation study for the two-way interaction, we considered four different combinations of the covariate effect and SNP main effect: γ (1.2 and 0.0) and δ (0.5 and 0.0).

Table 2 shows the power comparison of the five methods across various combinations of γ , δ , CF, and adjustment of the main effect. When there exists a covariate effect (that is, $\gamma = 1.2$), Cox-MDR, Cox-UMMDR and Cox2-UMMDR have relatively greater power than KM-MDR and KM-UMMDR. When there is an SNP main effect (that is, $\delta = 0.5$), the powers of KM-UMMDR, Cox-UMMDR, and Cox2-UMMDR were relatively greater when it is properly adjusted for than when it is not adjusted for. Interestingly, the powers of Cox-MDR and KM-MDR were higher than those of KM-UMMDR, Cox-UMMDR, and Cox2-UMMDR when both the covariate effect and SNP main effect exist (that is, $\gamma = 1.2, \delta = 0.5$). This implies that the main effect of the SNP may be reflected as any covariate effect. However, the power of Cox-MDR hardly decreases when only the main effect of SNP (that is, $\gamma = 0.0, \delta = 0.5$) is considered and is properly adjusted for. Note that the power of KM-MDR, Cox-UMMDR, and Cox2-UMMDR were similar higher to that of Cox-MDR when there was neither a covariate effect nor the SNP main effect (that is, $\gamma = \delta = 0.0$).

288
289

Table 2: Power of the five methods under 3-way interaction model

[γ, δ]	[1.2, 0.5]				[1.2, 0.0]		[0.0, 0.5]				[0.0, 0.0]	
	True		False		False		True		False		False	
main effect												
CF	0.10	0.30	0.10	0.30	0.10	0.30	0.10	0.30	0.10	0.30	0.10	0.30
Cox-MDR	0.60	0.33	0.60	0.33	0.70	0.40	0.10	0.02	0.10	0.02	0.60	0.34
KM-MDR	0.41	0.26	0.41	0.26	0.43	0.27	0.35	0.22	0.25	0.17	0.84	0.64
KM-UMMDR	0.34	0.26	0.10	0.05	0.45	0.32	0.28	0.21	0.10	0.03	0.74	0.62
Cox-UMMDR	0.34	0.26	0.10	0.05	0.71	0.59	0.35	0.30	0.17	0.06	0.90	0.70
Cox2-UMMDR	0.46	0.30	0.10	0.05	0.58	0.39	0.41	0.25	0.17	0.05	0.90	0.70

290

291 **Real data analysis**

292 We analyzed a real dataset consisting of 97 Korean AML patients with demographic information of age and sex, and genetic information of
293 139 SNPs. At the end of the study, 40 deaths occurred, and 57 patients were alive. We applied the two proposed methods of KM-
294 UMMDR and Cox2-UMMDR to detect SNP-SNP interactions associated with survival time by adjusting for age and sex.

295 To consider the marginal effect of SNP, we first fitted a univariate Cox model with each SNP by adjusting for age and sex. We found
296 that 21 SNPs have significant marginal effects on survival time. To summarize these marginal effects of 21 SNPs, we performed the prin-
297 ciple component analysis and took two principal components (PC) as a covariate, which accounted for 78% of the variation. We considered
298 four different models in identifying gene-gene interactions by KM-UMMDR and Cox2-UMMDR as follows:

299 (1) PC unadjusted and main effects of SNP1 and SNP2 unadjusted:

300

$$\lambda(t|S, age, sex) = \lambda_0(t) \exp(\beta S + \gamma_1 age + \gamma_2 sex)$$

301

302 (2) PC adjusted and main effects of SNP1 and SNP2 unadjusted:

303

$$\lambda(t|S, age, sex, PC_1, PC_2) = \lambda_0(t) \exp(\beta S + \gamma_1 age + \gamma_2 sex + \delta_1 PC_1 + \delta_2 PC_2)$$

304

305 (3) PC unadjusted and main effects of SNP1 and SNP2 adjusted:

306

$$\lambda(t|S, age, sex, SNP1, SNP2) = \lambda_0(t) \exp(\beta S + \gamma_1 age + \gamma_2 sex + \theta_1 SNP1 + \theta_2 SNP2)$$

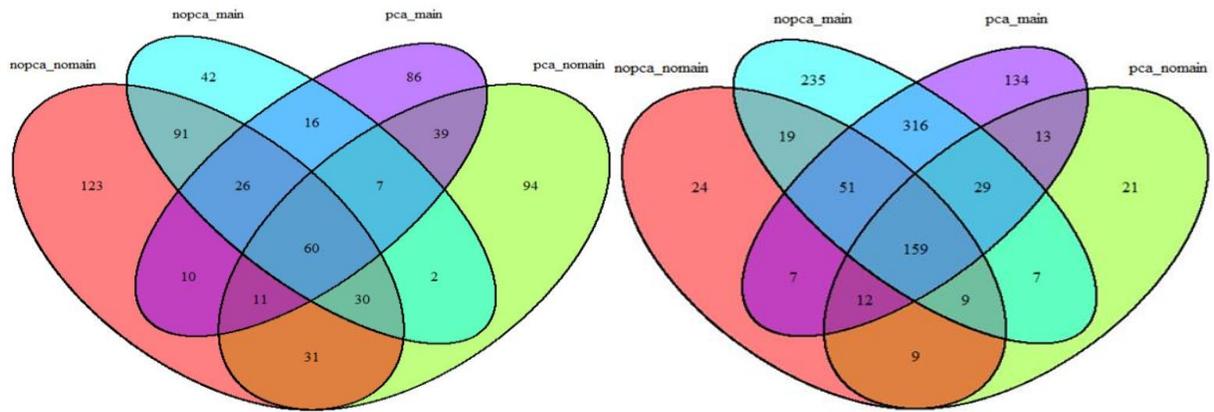
307

308 (4) PC adjusted and main effects of SNP1 and SNP2 adjusted:

309

$$\lambda(t|S, age, sex, PC_1, PC_2, SNP1, SNP2) = \lambda_0(t) \exp(\beta S + \gamma_1 age + \gamma_2 sex + \delta_1 PC_1 + \delta_2 PC_2 + \theta_1 SNP1 + \theta_2 SNP2)$$

310



311
 312 Figure 5: Venn Diagram for the number of significant 2-way SNP pairs identified by Cox2-UMMDR(left) and KM-UMMDR(right)
 313

314 We obtained a list of significant SNP pairs from each model, and summarized their numbers using a Venn diagram. Figure 5 shows
 315 the Venn Diagrams for the number of significant two-way SNP pairs identified by Cox2-UMMDR (left) and KM-UMMDR (right). The
 316 number of significant SNP pairs differed by these two methods under the four different models. We show that 159 significant SNP pairs
 317 were overlapped by four different models under KM-UMMDR whereas only 60 significant SNP pairs were overlapped by four different
 318 models under Cox2-UMMDR. In addition, although the number of SNP pairs varied across the four models, 40 SNP pairs were found to
 319 be commonly significant by both Cox2-UMMDR and KM-UMMDR.

320 Among these SNP pairs, we chose the best two SNP pairs that yielded the H/L indicator and checked their p-value from the fitted Cox
 321 model, adjusting for age and sex in step 2. The SNP pair (“rs580032”, “rs1960207”) were present in the list of the best SNP pairs obtained
 322 by both Cox2-UMMDR and KM-UMMDR, and a very significant p-value less than 0.001 was obtained.

323

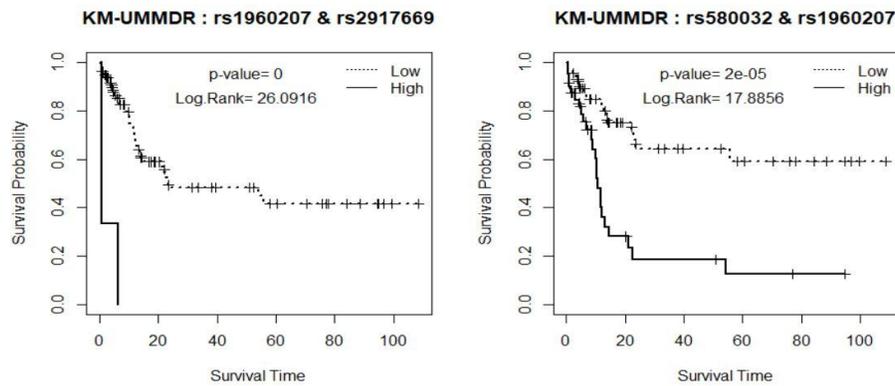
324 **Table 3. Significance test for the interaction effects of top two SNP pairs identified by Cox2-UMMDR and KM-UMMDR**

Method	SNP1	SNP2	P-value
KM UM-MDR	rs1960207	rs2917669	< 0.001
	rs580032	rs1960207	< 0.001
Cox2 UM-MDR	rs580032	rs1960207	< 0.001
	rs17170153	rs1801133	< 0.001

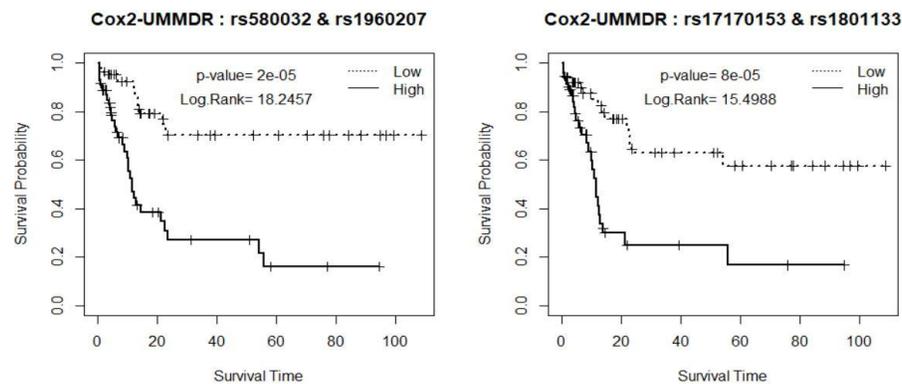
325

326 In addition, we provide the Kaplan-Meier survival plots of these two groups that were divided with respect to the H/L indicator at-
 327 tributed to the corresponding SNP pairs, as shown in Figures 6 and 7. In other words, both ‘high-risk’ and ‘low-risk’ groups were defined
 328 via step 1 classification using the best SNP pairs obtained by KM-MDR and the modified Cox-MDR, respectively. It is shown that the
 329 identified SNP pairs could separate the two survival plots significantly, which implies that there may be two-way interaction effects associ-
 330 ated with the survival phenotype. As shown in Figure 6, the SNP pair (“rs1960207”, “rs2917669”) was present in the extremely sparse
 331 high-risk group, which consists of only 3 individuals and the others belong to the low-risk group. Even though there is a large imbalance

332 between the high- and low-risk groups, the survival curve for the high-risk group was significantly lower than that for the low-risk group.
 333 The SNP pair (“rs580032”, “rs1960207”) was selected as one of the best top two pairs by both KM-UMMDR and Cox2-UMMDR. How-
 334 ever, the two KM-plots were not the same and the corresponding log-rank test statistics were also different at 17.8856 and 18.2457, respec-
 335 tively. This implies that KM-UMMDR and Cox2-UMMDR classify the high- and low-risk groups differently in step 1 because the former
 336 uses the median survival time and the latter uses the martingale residual, respectively. In fact, one cell among all combinations of two SNPs
 337 (“rs580032”, “rs1960207”) is classified as high-risk by KM-UMMDR but as low-risk by Cox2-UMMDR. However, the corresponding cell
 338 has only three individuals and provides very similar values to the log-rank test statistics.
 339



340
 341 Figure 6: Kaplan-Meier survival plots attributed by SNP pairs from KM-UMMDR
 342



343
 344 Figure 7: Kaplan-Meier survival plots attributed by SNP pairs from Cox2-UMMDR
 345
 346

347 Discussion

348 We proposed two simple methods, KM-UMMDR and Cox2-UMMDR, by adopting the classification rules from KM-MDR and the modi-
 349 fied Cox-MDR. These methods are extensions of KM-MDR and Cox-MDR to UM-MDR, similar to Cox-UMMDR. In this study, we fo-

350 cused on the comparison of the power of the five different methods such as KM-UMMDR, Cox2-UMMDR, Cox-UMMDR, Cox-MDR,
351 and KM-MDR across various scenarios.

352 Among those, both Cox-MDR and KM-MDR need cross-validation and permutation testing to identify significant interaction models
353 whereas Cox-UMMDR, KM-UMMDR, and Cox2-UMMDR provide the significance of the interaction model in the frame of the regression
354 model without any intensive computing. In contrast, both KM-MDR and KM-UMMDR are nonparametric approaches in which any co-
355 variate effect cannot be adjusted for in the classification procedure, whereas Cox-MDR, Cox-UMMDR, and Cox2-UMMDR can adjust for
356 the covariate effect. However, the covariate effect is considered in the procedure of testing for the high-risk group classified by KM-
357 UMMDR in the regression model in step 2. Therefore, KM-UMMDR is more flexible than KM-MDR in the sense that the former can
358 adjust for the covariate effect and substantially reduce the computing time.

359 Comparing Cox2-UMMDR with Cox-UMMDR, the simulation results show that Cox-UMMDR is more powerful than Cox2-
360 UMMDR, which implies that the effect of covariates should be adjusted for in the regression model rather than in the classification proce-
361 dure. Furthermore, when there exists any SNP main effect, it should be properly adjusted for in the regression model because the powers
362 of Cox-UMMDR, KM-UMMDR, and Cox2-UMMDR are very low when it is not adjusted for. Throughout the simulation results, KM-
363 UMMDR, Cox2-UMMDR, and Cox-UMMDR were more powerful than Cox-MDR and KM-MDR when there were any main effects of
364 SNPs.

365 We performed a simulation study for only one three-way interaction model because few higher-order interaction models are available
366 for simulation studies. Although only one model is provided in the simulation study, the power trend seems to be similar to that of the two-
367 way interaction model. Focusing on the two-way interaction model, Cox-UMMDR is the most powerful, except for a few cases among the
368 five methods whereas KM-UMMDR and Cox2-UMMDR have reasonable power in most cases.

369 In this study, we only consider a Cox regression model in step 2 but it is possible to fit any other regression model such as an AFT
370 model when the proportional hazard assumption is not valid for some covariates. For example, multi-genotypes are classified into high-
371 and low- risk groups by KM-MDR in step 1 and fit an AFT regression model with an indicator of the high-risk group and the other adjust-
372 ing covariates in step 2. By testing the significance of the indicator for the high-risk group, we can show that SNPs have a high-order inter-
373 action effect on the survival phenotype.

374 **Declarations**

375 **Ethics approval and consent to participate**

376 Not applicable

377 **Consent for publication**

378 Not applicable

379 **Availability of data and materials**

380 The real data of Korean leukemia patients is not available. The R-program for KM-UMMDR and Cox2-UMMDR is available at
381 <http://github.com/leesy80/Seungylee/coxumMDR>

382 **Competing interests**

383 The authors declare that they have no competing interests

384 **Funding**

385 This research was supported by the Basic Science Research Program through the National Research Foundation (NRF) funded by the Ministry of Science,
386 ICT & Future Planning (2019R1F1A1062005) and (2017M3A9C4065964).

387 Conflict of Interest: none declared.

388 **Authors' contributions**

389 SYL and JWL conceived the study and designed the simulation study. JWL and SYL wrote the manuscript and JWL implemented the simulation program
390 and analysed a real data. JWL and SYL read the paper and approved the final manuscript.

391 **Author details**

392 ¹Department of Statistics, University of Connecticut, Storrs, Connecticut, USA

393 ²Department of Mathematics and Statistics, Sejong University, 209 Neungdong-ro, Gwangjin-gu, 05006, Seoul, Korea

394

395

396

397

398

399

400

401

402

403

404

405

References

1. Klein RJ, Zeiss C, Chew EY, Tsai JY, Sackler RS, Haynes C, Henning AK, SanGiovanni JP, Mane SM, Mayne ST, Bracken MB, Ferris FL, Ott J, Narnstable C and Hoh J. Complement factor H polymorphism in age-related macular degeneration. *Science*. 2005; 308(5720):385-389.
2. Moore JH, Williams SW. New strategies for identifying gene-gene interactions in hypertension. *Ann. Med.* 2002; 34:88-95.
3. Manolio TA. Genome-wide association studies and assessment of the risk of disease. *N. Engl. J. Med.* 2010;363: 166-176.
4. Eichler EE, Flint J, Gibson G, Kong A, Leal SM, Moore, JH, and Nadeau JH. Missing heritability and strategies for finding the underlying causes of complex disease. *Nat. Rev. Genet.* 2010; 11(6):446-450.
5. Ritchie MD, Hahn LW, Roodi N, Bailey LR, Dupont WD, Parl FF, Moore JH. Multifactor-dimensionality reduction reveals high-order interactions among estrogen-metabolism genes in sporadic breast cancer. *Am. J. Hum. Genet.* 2001; 69: 138-147.
6. Gui J, Moore JH, Kelsey KT, Marsit CJ, Karagas MR, Andrew AS. A novel survival multifactor dimensionality reduction method for detecting gene-gene interactions with application to bladder cancer prognosis. *Hum. Genet.* 2011; 129:101-110.
7. Lee SY, Kwon MS, Oh JM, Park T. Gene-gene interaction analysis for the survival phenotype based on the Cox model. *Bioinformatics.* 2012; 28: i582-i588.
8. Oh JS, Lee SY. An extension of multifactor dimensionality reduction method for detecting gene-gene interactions with the survival time. *J. Korean Data & Information Science Soc.* 2014; 25(5):1-11.
9. Park M, Lee JW, Park T, Lee SY. Gene-gene interaction analysis for the survival phenotype based on the Kaplan-Meier median estimate. *BioMed Research International.* 2020; Article ID 5282345:1-10.
10. Yu W, Lee, SY, Park T. A unified model based multifactor dimensionality reduction framework for detecting gene-gene interactions. *Bioinformatics.* 2016; 32: i605-i610.
11. Lee SY, Son DH, Kim YK, Yu W, Park T. Unified cox model based multifactor dimensionality reduction method for gene-gene interaction analysis of the survival phenotype. *BioData Mining.* 2018; 11(27):1-13.
12. Velez DR, White BC, Motsinger, AA, Bush WS, Ritchie MD, Williams SM *et al.* A balanced accuracy metric for epistasis modeling in imbalanced datasets using multifactor dimensionality reduction. *Genet. Epidemiol.* 2007; 31:306-315.

Figure 1: Power curves of five methods under $\alpha=0.8, \gamma=0.0$

1 **Figure 2:** Power curves of five MDR methods under $\alpha=0.8, \gamma=0.8$

2

3 **Figure 3:** Power curves of five MDR methods under $\alpha=0.8, \gamma=0.0, \delta=0.5$

4

5 **Figure 4:** Power curves of five MDR methods under $\alpha=0.8, \gamma=0.8, \delta=0.5$

Figure 5: Venn Diagram for the number of significant 2-way SNP pairs identified by Cox2-UMMDR(left) and KM-UMMDR(right)

Figure 6: Kaplan-Meier survival plots attributed by SNP pairs from KM-UMMDR

Figure 7: Kaplan-Meier survival plots attributed by SNP pairs from Cox2-UMMDR