

# **Exploring predictive biomarkers from clinical genome-wide association studies via multidimensional hierarchical mixture models**

Takahiro Otani<sup>1</sup>, Hisashi Noma<sup>2</sup>, Shonosuke Sugasawa<sup>3</sup>, Aya Kuchiba<sup>4</sup>,  
Atsushi Goto<sup>5</sup>, Taiki Yamaji<sup>5</sup>, Yuta Kochi<sup>6</sup>, Motoki Iwasaki<sup>5</sup>,  
Shigeyuki Matsui<sup>1</sup>, and Tatsuhiko Tsunoda<sup>7,8</sup>

<sup>1</sup> Department of Biostatistics, Nagoya University Graduate School of Medicine, Nagoya, Aichi, Japan

<sup>2</sup> Department of Data Science, The Institute of Statistical Mathematics, Tachikawa, Tokyo, Japan

<sup>3</sup> Center for Spatial Information Science, The University of Tokyo, Kashiwa, Chiba, Japan

<sup>4</sup> Division of Biostatistical Research, Center for Public Health Sciences, National Cancer Center, Chuo-ku, Tokyo, Japan

<sup>5</sup> Division of Epidemiology, Center for Public Health Sciences, National Cancer Center, Chuo-ku, Tokyo, Japan

<sup>6</sup> Laboratory for Autoimmune Diseases, RIKEN Center for Integrative Medical Sciences, Yokohama, Kanagawa, Japan

<sup>7</sup> Department of Medical Science Mathematics, Medical Research Institute, Tokyo Medical and Dental University, Bunkyo-ku, Tokyo, Japan

<sup>8</sup> Laboratory for Medical Science Mathematics, RIKEN Center for Integrative Medical Sciences, Yokohama, Kanagawa, Japan

Correspondence: Takahiro Otani, Ph.D.

Department of Biostatistics

Nagoya University Graduate School of Medicine

65 Tsurumai-cho, Showa-ku, Nagoya, Aichi 466-8550, Japan

Phone +81-52-744-2487

e-mail: [otani@med.nagoya-u.ac.jp](mailto:otani@med.nagoya-u.ac.jp)

Running title: Exploring biomarkers from clinical GWASs

Conflict of interest: The authors declare no conflict of interest.

## Abstract

1  
2 Although the detection of predictive biomarkers is of particular importance for the  
3 development of accurate molecular diagnostics, conventional statistical analyses based  
4 on gene-by-treatment interaction tests lack sufficient statistical power for this purpose,  
5 especially in large-scale clinical genome-wide studies that require an adjustment for  
6 multiplicity of a huge number of tests. Here we demonstrate an alternative efficient  
7 multi-subgroup screening method using multi-dimensional hierarchical mixture models  
8 developed to overcome this issue, with application to stroke and breast cancer  
9 randomized clinical trials with genomic data. We show that estimated effect size  
10 distributions of single nucleotide polymorphisms (SNPs) associated with outcomes,  
11 which could provide clues for exploring predictive biomarkers, optimizing  
12 individualized treatments, and understanding biological mechanisms of diseases.  
13 Furthermore, using this method we detected three new SNPs that are associated with  
14 blood homocysteine levels, which are strongly associated with the risk of stroke. We  
15 also detected six new SNPs that are associated with progression-free survival in breast  
16 cancer patients.

17 Keywords: predictive biomarker; randomized clinical trial; genome-wide association  
18 study; multidimensional hierarchical mixture model; optimal discovery procedure

19

## 1 **Introduction**

2 The development of accurate molecular diagnostics for choosing the best treatment to  
3 maximize benefits or minimize risks in a particular individual is a crucial issue for the  
4 realization of precision medicine. This will require predicting therapeutic responses, and  
5 to this end, it is of particular importance to efficiently explore predictive biomarkers that  
6 successfully classify patients so that treatment effects differ between subgroups. One of  
7 the most promising study designs for this purpose is the genome-wide association study  
8 (GWAS), as this approach allows for the investigation of medical traits such as drug  
9 metabolism, efficacy, and toxicity<sup>1-3</sup>. In these studies, detection of gene-by-treatment  
10 interactions is one of the crucial issues for developing predictive biomarkers<sup>4</sup>. Despite  
11 the significant effort that has been devoted to GWASs, most published studies have  
12 failed to identify such effective predictive biomarkers.

13 One of the most fundamental problems of gene-by-treatment interaction tests is  
14 their lack of sufficient statistical power. In general, these tests are based on ordinary  
15 regression models and have low statistical power compared to tests for detecting genetic  
16 main effects. The reason is that unreliability is compounded in the interaction term of  
17 the models since the tests assess the difference of magnitudes of gene effects between  
18 treatment groups rather than simply assessing the magnitudes of the gene effects  
19 themselves. A rule of thumb to detect interaction effects has been suggested, whereby  
20 the detection requires a sample size at least four times larger than that required for the  
21 detection of a main effect of comparable magnitude<sup>4,5</sup>. In addition, most GWASs strictly  
22 control the conservative genome-wide significance level ( $p < 5 \times 10^{-8}$ ) for the  
23 interaction tests to adjust multiplicity. These conventional analysis strategies potentially  
24 set serious limits on the value of these studies' outcomes. While we in fact identified a

1 small number of useful biomarkers via these GWASs, the primary purpose of  
2 large-scale data analyses should be to effectively screen for genes that should be further  
3 investigated as candidate biomarkers in individualized precision medicine.

4 To overcome the lack of statistical power, an alternative effective multi-subgroup  
5 gene screening method<sup>6</sup> using a multidimensional semi-parametric hierarchical mixture  
6 model<sup>7,8</sup> has been developed by Matsui et al. (Figure 1). This method efficiently reveals  
7 the existence of predictive genes that are differently associated with outcomes between  
8 subgroups (treatment and control groups in a randomized clinical trial, for example),  
9 and prognostic genes that are similarly associated with outcomes irrespective of  
10 subgroups. To this end, the method eliminates two types of nuisance factors in  
11 association analysis results: (i) genes that are not associated with outcomes and (ii)  
12 random variation irrespective of association with outcomes. The method achieves these  
13 goals by using the hierarchical mixture model, and reveals the underlying effect size  
14 distribution of genes that are associated with outcomes. The denoised distribution itself  
15 could provide clues for exploring predictive biomarkers, optimizing individualized  
16 treatments, and understanding the biological mechanisms of diseases. One can  
17 demonstrate the existence of predictive/prognostic gene subgroups and their effect sizes  
18 from the distribution rather than by identifying individual genes using interaction tests  
19 with a conservative significance criterion. Furthermore, effect size estimates adjusted  
20 for overestimation error arising from association analyses, the so-called winner's curse  
21 phenomenon<sup>9</sup>, can be obtained based on the distribution. The adjusted estimates would  
22 serve as fundamental information for developing appropriate therapeutic strategies. In  
23 addition, an efficient test can also be developed using the optimal discovery procedure  
24 (ODP)<sup>10,11</sup>, which can provide an optimal ranking of genes as well as the most powerful

1 test for detecting disease-related genes with control of multiplicity, e.g., the false  
2 discovery rate (FDR). Note that the primary purpose of the newly developed method is  
3 to use a different principle than that underlying conventional gene-by-treatment  
4 interaction tests in order to provide an alternative strategy to more effectively identify  
5 predictive biomarkers, one that will overcome the issue of low statistical power and  
6 facilitate precision medicine research from a data analytic perspective.

7 In this article, we demonstrate the important strengths of these newly available  
8 tools by applying them to two large randomized clinical trials, the vitamin intervention  
9 stroke prevention (VISP) trial<sup>12</sup> and the SUCCESS-A trial, which study stroke and  
10 breast cancer, respectively, using genomic data (see Descriptions of GWAS Datasets and  
11 Section A in the Supplementary Notes) to detect single nucleotide polymorphisms  
12 (SNPs) that can be used to predict responses to therapeutics. We present the denoised  
13 effect size distributions of SNPs that are associated with medical outcomes, so as to  
14 assess the existence of predictive SNPs; disease-related SNPs detected by the ODP,  
15 along with their characteristics; and the results of a genomic investigation of these  
16 disease-related SNPs conducted using publicly available tools and databases, for the  
17 purposes of biological investigation and validation of the new method.

## 18 **Materials and methods**

### 19 **Descriptions of GWAS datasets**

20 The two datasets used in this analysis were deposited in the dbGaP database (available  
21 at <https://www.ncbi.nlm.nih.gov/gap>) and derived from the VISP trial (study accession  
22 number: phs000343.v3.p1) and the SUCCESS-A trial (study accession number:  
23 phs000547.v1.p1). For details, see Section A in the Supplementary Notes.

24 The VISP trial was a multi-center, double-blind, randomized, controlled clinical

1 trial that enrolled patients aged 35 or older with homocysteine levels above the 25th  
2 percentile at screening and a non-disabling cerebral infarction within 120 days of  
3 randomization. The trial was designed to determine if daily intake of a multivitamin  
4 tablet containing high-dose folic acid, vitamin B6, and vitamin B12 reduced recurrent  
5 cerebral infarction as well as nonfatal myocardial infarction or mortality. Subjects were  
6 randomly assigned to receive daily doses of the high-dose formulation (treatment group)  
7 or the low-dose formulation (control group). A total of 1533 subjects (760 assigned to  
8 the treatment group and 773 assigned to the control group) with 774670 SNPs passed a  
9 quality control filter (see Section B in the Supplementary Notes). In this study, we used  
10 this dataset to investigate SNPs associated with blood homocysteine levels which are  
11 strongly associated with the risk of stroke. We took as outcome the difference in blood  
12 homocysteine levels between baseline and the first post-baseline measurements, as in  
13 the study of Wakefield et al.<sup>13</sup>. Association tests were conducted using a linear  
14 regression model (see Section C in the Supplementary Notes).

15 The SUCCESS-A trial was a randomized phase III study of treatment response of  
16 early primary breast cancer to adjuvant therapy after surgical resection. The trial was  
17 designed to determine if adjuvant chemotherapy with gemcitabine, an antimetabolite  
18 frequently used in the treatment of pancreatic cancer and other diseases<sup>14</sup>, improved  
19 progression-free survival, overall survival, and toxicity. Subjects were randomly  
20 assigned to chemotherapy with gemcitabine (treatment group) or without gemcitabine  
21 (control group). A total of 3289 subjects (1621 assigned to the treatment group and 1668  
22 assigned to the control group) with 424121 SNPs passed the quality control filter (see  
23 Section B in the Supplementary Notes). In this study, we used this dataset to investigate  
24 SNPs associated with progression-free survival in breast cancer patients. Association

1 tests were conducted using a proportional hazards regression model (see Section C in  
2 the Supplementary Notes).

### 3 **Multi-subgroup gene screening method**

4 In this analysis, we used the efficient multi-subgroup gene screening method<sup>6</sup> (Figure 1)  
5 developed to overcome the problem of insufficient power to detect interaction effects  
6 (see Section D in the Supplementary Notes for details). Contrary to standard  
7 gene-by-treatment interaction tests using regression models with interaction terms  
8 between genes and treatments, this method first separates control and treatment groups.  
9 Then, association analyses using regression models without interaction terms are  
10 independently conducted for each group. As a result, summary statistics, i.e., estimated  
11 gene main effect sizes and their standard errors, are obtained for each group. After that,  
12 using the summary statistics, this method reveals the proportion of genes that are  
13 disease related and the underlying effect size distribution of disease-related genes across  
14 treatment and control groups via empirical Bayes estimation under the multidimensional  
15 hierarchical mixture model. Furthermore, posterior probabilities of association and  
16 effect size estimates adjusted for gene selection errors and overestimation for each SNP  
17 are obtained based on the estimated distribution. Finally, based on the optimal gene  
18 ranking and posterior probabilities of association for each gene derived from the fitted  
19 model, disease-related genes are detected by the ODP with control of FDR.

### 20 **Genomic annotation**

21 For biological investigation of detected SNPs and for validation of the new analysis  
22 method, we conducted genomic annotations on LD surrogates of SNPs detected by our  
23 analysis using the publicly available ENCODE (Encyclopedia of DNA Elements) tools.  
24 Genomic annotation on LD surrogates ( $r^2 > 0.6$ , 1000 Genomes Project EUR data) of

1 newly detected SNPs was conducted using HaploReg<sup>15,16</sup> v4.1 (available at  
2 <http://archive.broadinstitute.org/mammals/haploreg/haploreg.php>) and RegulomeDB<sup>17</sup>  
3 v1.1 (available at <http://regulomedb.org/>).

## 4 **Results**

### 5 **Blood homocysteine levels**

6 We conducted ordinary association analyses for control and treatment groups using the  
7 linear regression model and obtained effect size estimates (regression coefficients) and  
8 their standard errors for each SNP (Supplementary Fig. S1a). We then eliminated  
9 nuisance factors in these results using the hierarchical mixture model and obtained the  
10 effect size distribution of SNPs associated with homocysteine levels (Figures 2a and  
11 2c).

12 The estimated distribution identifies the multi-subgroup SNPs that can be  
13 classified as possible prognostic or predictive biomarkers; peaks on the diagonal line in  
14 the first and third quadrants of the distribution correspond to prognostic markers, while  
15 others correspond to predictive markers. The proportion of disease-related SNPs was  
16 estimated to be 0.001. This suggests that 793 SNPs are associated with outcomes; note  
17 that the number of independently associated SNPs should be much smaller since some  
18 SNPs are in linkage disequilibrium (LD). The denoised distribution suggests that almost  
19 all SNPs have small effects on homocysteine levels regardless of control and treatment  
20 assignments, although the large peak is shifted slightly in the negative direction,  
21 corresponding to a decrease in homocysteine levels in the treatment group. The slight  
22 shift of the peak is due to the effect of high-dose administration of multivitamin tablets  
23 and is not related to any genetic properties.

24 In addition to the large peak, there is a small peak that deviates from the low-effect



1 area, suggesting the existence of SNPs that will predict the benefit of high-dose  
2 administration of multivitamin tablets. This peak has a small positive effect on low-dose  
3 administration (effect size of 0.2) but has a large negative effect on high-dose  
4 administration (effect size of -3.3). The location of the peak suggests the existence of  
5 strong interaction effects between particular SNPs and high-dose administration that  
6 drastically decrease homocysteine levels. The marginal distribution for control and  
7 treatment groups clearly shows this difference (Figure 2c). The deviant peak is shown  
8 for the high-dose formulation (treatment) group, while no corresponding peak exists for  
9 the low-dose formulation (control) group.

10 We next detected SNPs associated with homocysteine levels using the ODP (Table  
11 1), based on the optimal ranking via posterior probabilities of association (Figure 3)  
12 obtained from the estimated effect size distribution. The ODP detected five independent  
13 SNPs (FDR<5%) that are associated with outcomes. Of these, rs12631354, rs2367209,  
14 and rs10017302 are newly detected in this analysis while others have previously been  
15 suggested as being associated with homocysteine levels<sup>13</sup>. In particular, rs12631354 and  
16 rs10017302 might significantly change the effect of the administration of multivitamin  
17 tablets. These have strong interaction effects with the high-dose administration of  
18 multivitamin tablets, drastically decreasing homocysteine levels. From the therapeutic  
19 point of view, these two SNPs might be useful predictive biomarkers.

20 The ODP also detected SNPs that have been suggested as being associated with  
21 homocysteine levels<sup>13</sup>, as well as rs3736238, which has been previously reported as the  
22 most statistically significant SNP<sup>13</sup> and which reached genome-wide significance  
23 ( $p < 5 \times 10^{-8}$ ) using the standard regression model with interaction terms (see Section  
24 E in the Supplementary Notes and Supplementary Fig. S2a). According to the adjusted

1 effect size estimates, rs3736238 shows the largest benefit for the high-dose  
2 administration group but demonstrates a harmful effect for the low-dose administration  
3 group (Table 1 and Supplementary Fig. S3a). Although the biological mechanisms of  
4 the harmful effect with low-dose administration are unclear, this result suggests the need  
5 for high-dose administration of multivitamin tablets in stroke patients. Another SNP,  
6 rs1739317, which has also been reported previously<sup>13</sup> and which reached a suggestive  
7 level ( $p < 10^{-6}$ ) using the standard regression model (Supplementary Fig. S2a), was  
8 also detected. As with the newly detected SNPs, this SNP has a strong interaction effect  
9 with the high-dose administration that drastically decreases homocysteine levels and  
10 results in neither benefit nor harm with the low-dose administration. On the other hand,  
11 rs16893296 on chromosome 6, which has been suggested as being associated with  
12 homocysteine levels<sup>13</sup>, was not detected by the ODP.

### 13 **Breast cancer**

14 As with the analysis of homocysteine levels, we conducted ordinary association  
15 analyses for control and treatment groups using the proportional hazards regression  
16 model and obtained estimated regression coefficients (log hazard ratios) and their  
17 standard errors for each SNP (Supplementary Fig. S1b). We then applied the  
18 hierarchical mixture model to the analysis results and obtained the effect size  
19 distribution of SNPs associated with progression-free survival in breast cancer patients  
20 (Figures 2b and 2d). The proportion of disease-related SNPs was estimated to be 0.002;  
21 although some of these SNPs are in LD, this finding suggests that 903 SNPs are  
22 associated with outcomes.

23 The estimated distribution is more spread out in the gemcitabine (treatment) group,  
24 although since the peak of the distribution is centered almost at 0, almost all

1 disease-related SNPs have small effects on progression-free survival in patients  
2 irrespective of treatment assignments. The difference in spread of effect size  
3 distributions between the control group and the treatment group is clearly revealed by  
4 the marginal distribution (Figure 2d). The spread in the treatment group suggests the  
5 existence of “beneficial” and “harmful” SNPs, i.e., some SNPs might increase survival  
6 rate with gemcitabine administration while others might do the opposite. These  
7 estimates suggest that careful evaluations of effect sizes for each predictive SNP are  
8 necessary to develop appropriate therapeutic strategies.

9 We next applied the ODP to detect SNPs associated with progression-free survival  
10 in breast cancer patients based on the posterior probability of association (Figure 4). The  
11 ODP detected new six SNPs (FDR<5%) associated with progression-free survival in  
12 breast cancer patients (Table 1). Since the existence of harmful SNPs was suggested by  
13 the denoised distribution, we assessed the effect sizes of these SNPs adjusted for errors  
14 arising from the association analysis. The adjusted effect size estimates (Table 1 and  
15 Supplementary Fig. S3b) suggest that rs6712299 and rs17367673 have beneficial effects  
16 on progression-free survival with gemcitabine administration while rs4690351 and  
17 rs12449931 have harmful effects. Although further investigations are necessary, these  
18 SNPs might be useful predictive biomarkers to determine whether or not gemcitabine  
19 treatment should be conducted. On the other hand, rs12620133 and rs7311993 have  
20 beneficial effects in both the control and treatment groups. These SNPs might be used as  
21 prognostic biomarkers for developing risk-stratification systems.

## 22 **External validation**

23 We conducted genomic annotations on LD surrogates of SNPs detected by our analysis.  
24 For homocysteine level analysis, three SNPs (rs12631354, rs2367209, and rs10017302)

1 were used as queries, while for breast cancer analysis, six SNPs (rs12620133,  
2 rs6712299, rs4690351, rs7311993, rs12449931, and rs17367673) were used. Summaries  
3 of the results are presented in Supplementary Tables S1 to S9.

4 For blood homocysteine levels, we queried the three newly detected SNPs and  
5 obtained lists of LD surrogates with their characteristics (Supplementary Tables S1 to  
6 S3). As a high-LD surrogate of rs12631354, we identified rs4450813 ( $r^2 = 1$  based on  
7 1000 Genomes Project European population), which has an expression quantitative trait  
8 locus (eQTL) effect on the *RYK* gene in the liver<sup>18</sup> (Supplementary Table S1). This  
9 result seems biologically plausible because homocysteine is metabolized in the liver.  
10 For rs2367209, we identified the LD surrogate rs4679904 ( $r^2 = 0.66$ ), which has an  
11 eQTL effect on the *AF038199* gene in liver tissue<sup>19</sup> and is associated with primary  
12 biliary cirrhosis<sup>20</sup> (Supplementary Table S2). The association between rs4679904 and  
13 primary biliary cirrhosis is strongly supported by existing evidence and is recorded in  
14 the National Human Genome Research Institute (NHGRI) GWAS catalog<sup>21</sup> (available at  
15 <https://www.genome.gov/gwastudies/>). This result suggests the existence of an actual  
16 regulatory variant, despite the fact that the effect size of rs2367209 on the high-dose  
17 administration of multivitamin tablets is comparably small (Table 1), and further  
18 investigation should be conducted. According to epigenomic information for another  
19 high-LD surrogate, rs1879797 ( $r^2 = 0.84$ ), there is a cluster of active transcription start  
20 site/enhancer in digestive. This result might indicate the contribution of the variant to  
21 digestion of homocysteine. For rs10017302, the high-LD surrogates rs2126029  
22 ( $r^2 = 0.94$ ) and rs1460781 ( $r^2 = 0.94$ ) were found, both of which have an eQTL effect  
23 for the *GPM6A* gene in peripheral blood monocytes<sup>22</sup> (Supplementary Table S3). Since  
24 the purpose of this study was to investigate associations with the transcriptome of

1 circulating monocytes, a key cell type involved in immunity-related diseases and  
2 atherosclerosis, the SNPs might consistently be associated with cardiovascular diseases.  
3 These results support the biological validity of the testing method under the hierarchical  
4 mixture models.

5 We also investigated SNPs that were associated with progression-free survival in  
6 breast cancer patients. As with the stroke trial, we queried six independent SNPs (Table  
7 1) detected by our analysis and obtained lists of LD surrogates with their characteristics  
8 (Supplementary Tables S4 to S9). Some of these indicated the existence of biological  
9 mechanisms and might be truly associated with outcomes. Genomic annotation  
10 identified rs3821340 ( $r^2 = 0.88$ ), which has eQTL effects for the *AC073464.7*,  
11 *CYP4F32P*, and *ZNF514* genes in the pancreas<sup>23</sup>, as a surrogate of rs6712299  
12 (Supplementary Table S5), as well as rs4690439 ( $r^2 = 0.64$ ), which has an eQTL effect  
13 for the *WDR17* gene in the ovary<sup>23</sup>, as a surrogate of rs4690351 (Supplementary Table  
14 S6). Although existing pharmacogenomic/pharmacogenetic genome-wide studies have  
15 identified SNPs that are associated with responses to gemcitabine<sup>24–27</sup> and their results  
16 are recorded in the NHGRI GWAS catalog<sup>21</sup>, the newly detected SNPs did not match  
17 any records, partly because the previous studies focused on drug responses in pancreatic  
18 cancer, neutropenia, and/or leucopenia rather than breast cancer. Epigenomic  
19 information also suggests an association between some of the detected SNPs and breast  
20 cancer; enhancer activities in breast variant human mammary epithelial cells (vHMEC)  
21 exist for rs4690351, rs7311993, and rs17367673. Although there is still no strong  
22 evidence of an association with gemcitabine response and further validations are  
23 necessary, these results suggest the possibility of associations between these SNPs and  
24 breast cancer.

## 1 **Discussion**

2 In this analysis, we demonstrated an efficient multi-subgroup gene screening method  
3 using hierarchical mixture models and the ODP, with applications to molecular data  
4 from randomized clinical trials of stroke and breast cancer to detect predictive  
5 biomarkers. We found three new SNPs that were associated with blood homocysteine  
6 levels, and six new SNPs that were associated with progression-free survival in breast  
7 cancer patients. These SNPs have not been reported by existing GWASs.

8 This new method can more effectively detect predictive disease-related SNPs than  
9 conventional association tests that use regression models with interaction terms. For  
10 comparison, we conducted association tests based on standard regression models with  
11 interaction terms for each trial (see Section E in the Supplementary Notes). In the  
12 association analysis for homocysteine levels, only one SNP, rs3736238, reached a  
13 genome-wide significance level of  $p < 5 \times 10^{-8}$ , while rs16893296 and rs1739317  
14 reached the suggestive level of  $p < 10^{-6}$  (Supplementary Fig. S2a). Three other  
15 predictive SNPs detected by our analysis did not reach the suggestive level. On the  
16 other hand, rs16893296, which has been suggested as being associated with  
17 homocysteine levels<sup>13</sup>, was not detected by our analysis. No signal peak corresponding  
18 to this SNP was found in the estimated effect size distribution, and the posterior  
19 probability of association was estimated as only 13%. Also, the adjusted effect size  
20 estimates suggest that this SNP does not have a large effect on homocysteine levels  
21 (Supplementary Fig. S3a). This result may indicate the possibility of a false positive in  
22 the existing study or a false negative in our analysis, and further validation is necessary.  
23 In the breast cancer trial, no SNPs reached the suggestive level (Supplementary Fig.  
24 S2b), while our analysis detected six independent SNPs consisting of four predictive

1 SNPs and two prognostic SNPs. We also conducted a comparison of the number of  
2 detected SNPs between the standard association tests and the new method under  
3 specified FDR levels (see Section E in the Supplementary Notes). The new method  
4 detected more SNPs associated with homocysteine levels than the standard method  
5 under the same FDR levels, and also, unlike the standard method, effectively detected  
6 several SNPs associated with progression-free survival in breast cancer patients  
7 (Supplementary Table S10). Note that categorization of SNPs as predictive markers or  
8 prognostic markers was conducted in a subjective manner. Basically, SNPs that have  
9 different effect sizes between control and treatment groups would be categorized as  
10 predictive markers, while others would be categorized as prognostic markers. However,  
11 a specific criterion to categorize SNPs will be subjectively determined (see Matsui et  
12 al.<sup>6</sup>, for example). In addition, we assessed the performance of the ODP through a  
13 simulation study based on the two clinical trials (see Section F in the Supplementary  
14 Notes). The ODP detected larger numbers of significant SNPs with controlling FDR  
15 accurately, compared with the conventional methods (Supplementary Tables S11 and  
16 S12).

17 As demonstrated in this analysis, the denoised distribution can be used to explore  
18 the existence of predictive biomarkers and identify the best therapeutic strategies using  
19 a different approach than the ordinary gene identification scheme that uses  
20 gene-by-treatment interaction tests. For example, because the daily intake of high-dose  
21 multivitamin tablets has no serious harmful effects, such as increasing homocysteine  
22 levels, it can be administered to stroke patients according to the obtained distribution  
23 with no risk (Figures 2a and 2c). Also, in the low-dose formulation (control) group, all  
24 disease-related SNPs have effect sizes of nearly zero, corresponding to no impact on

1 homocysteine levels, and no peaks in the distribution deviate either positively or  
2 negatively (Figures 2c). This result means that low-dose administration does not  
3 maximize response or minimize side effects. From the therapeutic point of view, these  
4 findings suggest that high-dose multivitamin tablets should be administered to stroke  
5 patients irrespective of their genetic variations. On the other hand, the use of  
6 gemcitabine requires careful assessment of patients' genetic characteristics because the  
7 estimated distribution suggests that the drug might cause serious side effects in breast  
8 cancer patients with particular genetic variants and might decrease survival rates  
9 (Figures 2b and 2d). From a therapeutic point of view, these results indicate that  
10 individualized gemcitabine administration is necessary to improve survival rates and  
11 avoid side effects. Note that although this analysis strategy can reveal the existence of  
12 predictive/prognostic biomarkers as shown in the above examples, follow-up studies are  
13 necessary to definitively identify all of these markers, although some of them were  
14 successfully detected by our analysis using the ODP.

15       Although we demonstrated the two dimensional hierarchical mixture model to  
16 analyse molecular data consisting of two subgroups in this study, models with three or  
17 more dimensions can also be developed to explore higher order interactions. The  
18 existence of such interactions can be evaluated by the multidimensional models defining  
19 subgroups according to combinations of specific values of multiple variables.

20       The denoised distribution and the estimated proportion of disease-related SNPs can  
21 also be used for designing future medical genomics studies to identify predictive  
22 biomarkers. In particular, using the estimated results, we can obtain required sample  
23 sizes to find predictive biomarkers under a specified power and  $FDR^7$ . Furthermore, the  
24 estimated proportions of disease-related SNPs suggest the existence of other underlying



1 markers that might have smaller effect sizes. Although efforts to obtain adequate  
2 numbers of samples are naturally crucial, it is equally important to develop more  
3 efficient association tests on gene-by-treatment interactions, as demonstrated in this  
4 analysis, as these tests would be serve as a realistic approach to discovering predictive  
5 markers.

## 6 **Acknowledgments**

7 This work was supported by CREST, Japan Science and Technology Agency  
8 (JPMJCR1412), the Practical Research for Innovative Cancer Control (17ck0106266  
9 since 2017) from the Japan Agency for Medical Research and Development, and JSPS  
10 KAKENHI Grant Numbers JP16H06299 and JP17H01557.

## 11 **Conflicts of interest**

12 The authors declare no conflict of interest.

13

14 Supplementary information is available at European Journal of Human Genetics'  
15 website.

16

## 17 **References**

- 18 1 Daly AK. Genome-wide association studies in pharmacogenomics. *Nat Rev*  
19 *Genet* 2010; **11**: 241–246.
- 20 2 Motsinger-Reif AA, Jorgenson E, Relling MV *et al*. Genome-wide association  
21 studies in pharmacogenomics: successes and lessons. *Pharmacogenet Genomics*  
22 2013; **23**: 383–394.
- 23 3 Wheeler HE, Maitland ML, Dolan ME, Cox NJ, Ratain MJ. Cancer  
24 pharmacogenomics: strategies and challenges. *Nat Rev Genet* 2012; **14**: 23–34.

- 1 4 Thomas D. Gene-environment-wide association studies: emerging approaches.  
2 *Nat Rev Genet* 2010; **11**: 259–72.
- 3 5 Smith PG, Day NE. The design of case-control studies: the influence of  
4 confounding and interaction effects. *Int J Epidemiol* 1984; **13**: 356–365.
- 5 6 Matsui S, Noma H, Qu P *et al.* Multi-subgroup gene screening using  
6 semi-parametric hierarchical mixture models and the optimal discovery  
7 procedure: application to a randomized clinical trial in multiple myeloma.  
8 *Biometrics* 2017. doi:10.1111/biom.12716.
- 9 7 Matsui S, Noma H. Estimating effect sizes of differentially expressed genes for  
10 power and sample-size assessments in microarray experiments. *Biometrics* 2011;  
11 **67**: 1225–1235.
- 12 8 Nishino J, Kochi Y, Shigemizu D *et al.* Empirical Bayes estimation of  
13 semi-parametric hierarchical mixture models for unbiased characterization of  
14 polygenic disease architectures. *bioRxiv*,  
15 <http://biorxiv.org/lookup/doi/101101/080945> 2016. doi:10.1101/080945.
- 16 9 Ferguson JP, Cho JH, Yang C, Zhao H. Empirical Bayes correction for the  
17 Winner’s Curse in genetic association studies. *Genet Epidemiol* 2013; **37**: 60–68.
- 18 10 Storey JD, Dai JY, Leek JT. The optimal discovery procedure for large-scale  
19 significance testing, with applications to comparative microarray experiments.  
20 *Biostatistics* 2007; **8**: 414–432.
- 21 11 Noma H, Matsui S. The optimal discovery procedure in multiple significance  
22 testing: an empirical Bayes approach. *Stat Med* 2012; **31**: 165–176.
- 23 12 Spence JD, Howard VJ, Chambless LE *et al.* Vitamin Intervention for Stroke  
24 Prevention (VISP) trial: rationale and design. *Neuroepidemiology* 2001; **20**: 16–

- 1 25.
- 2 13 Wakefield J, Skrivankova V, Hsu F-C, Sale M, Heagerty P. Detecting signals in  
3 pharmacogenomic genome-wide association studies. *Pharmacogenomics J* 2014;  
4 **14**: 309–15.
- 5 14 Soo RA, Yong W-P, Innocenti F. Systemic therapies for pancreatic cancer - the  
6 role of pharmacogenetics. *Curr Drug Targets* 2012; **13**: 811–828.
- 7 15 Ward LD, Kellis M. HaploReg: a resource for exploring chromatin states,  
8 conservation, and regulatory motif alterations within sets of genetically linked  
9 variants. *Nucleic Acids Res* 2012; **40**: 930–934.
- 10 16 Ward LD, Kellis M. HaploReg v4: systematic mining of putative causal variants,  
11 cell types, regulators and target genes for human complex traits and disease.  
12 *Nucleic Acids Res* 2016; **44**: D877–D881.
- 13 17 Boyle AP, Hong EL, Hariharan M *et al.* Annotation of functional variation in  
14 personal genomes using RegulomeDB. *Genome Res* 2012; **22**: 1790–1797.
- 15 18 Innocenti F, Cooper GM, Stanaway IB *et al.* Identification, replication, and  
16 functional fine-mapping of expression quantitative trait loci in primary human  
17 liver tissue. *PLoS Genet* 2011; **7**: e1002078.
- 18 19 Greenawalt DM, Dobrin R, Chudin E *et al.* A survey of the genetics of stomach,  
19 liver, and adipose gene expression from a morbidly obese cohort. *Genome Res*  
20 2011; **21**: 1008–1016.
- 21 20 Hirschfield GM, Liu X, Xu C *et al.* Primary biliary cirrhosis associated with  
22 HLA, IL12A, and IL12RB2 variants. *N Engl J Med* 2009; **360**: 2544–2555.
- 23 21 Welter D, MacArthur J, Morales J *et al.* The NHGRI GWAS Catalog, a curated  
24 resource of SNP-trait associations. *Nucleic Acids Res* 2014; **42**: D1001–D1006.

- 1 22 Zeller T, Wild P, Szymczak S *et al.* Genetics and beyond – the transcriptome of  
2 human monocytes and disease susceptibility. *PLoS One* 2010; **5**: e10693.
- 3 23 GTEx Consortium. The Genotype-Tissue Expression (GTEx) pilot analysis:  
4 multitissue gene regulation in humans. *Science* 2015; **348**: 648–60.
- 5 24 Innocenti F, Owzar K, Cox NL *et al.* A genome-wide association study of overall  
6 survival in pancreatic cancer patients treated with gemcitabine in CALGB 80303.  
7 *Clin Cancer Res* 2012; **18**: 577–584.
- 8 25 Li L, Fridley BL, Kalari K *et al.* Discovery of genetic biomarkers contributing to  
9 variation in drug response of cytidine analogues using human lymphoblastoid  
10 cell lines. *BMC Genomics* 2014; **15**: 93.
- 11 26 Low SK, Chung S, Takahashi A *et al.* Genome-wide association study of  
12 chemotherapeutic agent-induced severe neutropenia/leucopenia for patients in  
13 Biobank Japan. *Cancer Sci* 2013; **104**: 1074–1082.
- 14 27 Kiyotani K, Uno S, Mushiroda T *et al.* A genome-wide association study  
15 identifies four genetic markers for hematological toxicities in cancer patients  
16 receiving gemcitabine therapy. *Pharmacogenet Genomics* 2012; **22**: 229–235.  
17  
18

1 **Titles and legends to figures**

2 Figure 1. Workflow of the gene screening method using multidimensional hierarchical  
3 mixture models. The top panels show the genotypic and phenotypic data of control and  
4 treatment groups. The sample sizes of control and treatment groups are denoted by  
5  $N_0$  and  $N_1$ , and  $M$  is the total number of SNPs. The middle panels show the summary  
6 statistics consisting of estimated effect sizes  $b^{(0)}, b^{(1)}$  and their standard errors for the  
7 main effect of SNPs derived from association analyses using regression models without  
8 interaction terms. The bottom panels show the denoised effect size distribution of SNPs  
9 that are associated with outcomes, posterior probabilities of association, and adjusted  
10 effect size estimates for each SNP. Significant SNPs are detected by the optimal  
11 discovery procedure based on the distribution and posterior probabilities.

12

13 Figure 2. Effect size distributions of SNPs associated with blood homocysteine levels or  
14 progression-free survival in breast cancer patients. (a and b) Two-dimensional  
15 distributions of homocysteine levels (a) and breast cancer (b). The  $x$  axis represents the  
16 effect size for the control group and the  $y$  axis represents the effect size for the treatment  
17 group. (c and d) Marginal distributions of homocysteine levels (c) and breast cancer (d).  
18 The  $x$  axis represents the effect size and the  $y$  axis represents the probability density, and  
19 distributions marginalized by the control group and the treatment group are plotted.

20

21

1 Figure 3. Plots of posterior probabilities of association with blood homocysteine levels  
2 for each SNP based on the estimated effect size distributions. Posterior probabilities for  
3 each SNP ( $y$  axis) are plotted by chromosomal position ( $x$  axis) in a similar way as with  
4 a Manhattan plot. Red points with rsIDs denote the probabilities of significant SNPs  
5 detected by the ODP ( $FDR < 5\%$ ) and are not in LD.

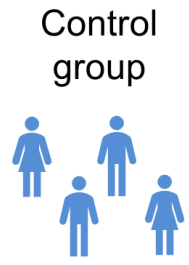
6

7 Figure 4. Plots of posterior probabilities of association with progression-free survival in  
8 breast cancer patients for each SNP based on the estimated effect size distributions.  
9 Posterior probabilities for each SNP ( $y$  axis) are plotted by chromosomal position ( $x$   
10 axis) in a similar way as with a Manhattan plot. Red points with rsIDs denote the  
11 probabilities of significant SNPs detected by the ODP ( $FDR < 5\%$ ) and are not in LD.

Table 1. New significant loci of blood homocysteine levels and breast cancer detected through the optimal discovery procedure.

SNP	Chr.	Position	A <sub>12</sub>	MAF	Effect size (95% CI)		Adjusted effect size (95% CI)		Posterior probability	Gene	Distance
					Control	Treatment	Control	Treatment			
<b>Blood homocysteine levels</b>											
rs12631354	3	134195045	CT	0.01	-0.07 (-1.41 to 1.28)	-3.90 (-5.24 to -2.56)	0.21 (0.19 to 0.42)	-3.30 (-3.46 to -1.55)	0.989	<i>RYK</i>	Intragenic
rs2367209	3	160681097	TG	0.27	0.14 (-0.28 to 0.56)	-0.99 (-1.40 to -0.58)	0.07 (-0.23 to 0.15)	-0.51 (-0.54 to -0.47)	0.871	<i>ARL14</i>	2.7 kb 3'
rs10017302	4	114563671	CT	0.02	0.52 (-0.85 to 1.89)	-3.63 (-4.89 to -2.37)	0.22 (0.20 to 0.69)	-3.30 (-3.63 to -1.17)	0.990	<i>UGT8</i>	35 kb 5'
rs1739317	6	24947873	TC	0.03	0.25 (-0.61 to 1.11)	-2.86 (-3.74 to -1.98)	0.22 (0.19 to 0.24)	-3.28 (-3.38 to -0.66)	>0.999	<i>FAM65B</i>	Intragenic
rs3736238	17	28881308	TC	0.01	1.88 (0.31 to 3.44)	-4.75 (-6.46 to -3.04)	1.35 (0.28 to 1.80)	-4.22 (-4.61 to -2.05)	0.982	<i>FLOT2</i>	Intragenic
<b>Breast cancer</b>											
rs12620133	2	3879585	AC	0.20	0.39 (0.14 to 0.64)	0.54 (0.28 to 0.80)	0.22 (0.09 to 0.36)	0.33 (0.20 to 0.47)	0.938	<i>DCDC2C</i>	32 kb 3'
rs6712299	2	95369055	CA	0.40	0.11 (-0.10 to 0.32)	0.58 (0.34 to 0.81)	0.07 (-0.07 to 0.22)	0.40 (0.28 to 0.50)	0.924	<i>KCNIP3</i>	Intragenic
rs4690351	4	176498757	GA	0.17	-0.05 (-0.35 to 0.25)	-0.69 (-0.95 to -0.44)	0.01 (-0.10 to 0.11)	-0.42 (-0.59 to -0.33)	0.963	<i>SPCS3</i>	167 kb 3'
rs7311993	12	11778652	AG	0.22	0.40 (0.14 to 0.66)	0.58 (0.33 to 0.84)	0.22 (0.08 to 0.36)	0.36 (0.25 to 0.51)	0.971	<i>ETV6</i>	Intragenic
rs12449931	17	79279996	GA	0.42	0.03 (-0.19 to 0.25)	-0.69 (-0.94 to -0.43)	0.03 (-0.07 to 0.14)	-0.44 (-0.59 to -0.33)	0.965	<i>RBFOX3</i>	Intragenic
rs17367673	19	46894238	AG	0.17	-0.14 (-0.41 to 0.14)	0.64 (0.39 to 0.89)	-0.06 (-0.18 to 0.08)	0.41 (0.31 to 0.53)	0.942	<i>ARHGAP35</i>	24 kb 5'

Effect size estimates and their 95% confidential intervals were obtained from association analyses under linear regression models for homocysteine levels and proportional hazards models for breast cancer. Effect size estimates for breast cancer correspond to log hazard ratios. Adjusted effect size estimates were calculated based on the denoised effect size distributions under the hierarchical mixture models. Positions based on hg38 were obtained from dbSNP build 141. Allele frequencies are based on 1000 Genomes Project European population and gene annotations are based on GENCODE version 13. Chr., chromosome. A<sub>12</sub>, reference and alternate alleles.



Control group

$N_0$ individuals	$M$ SNPs					Trait
	$rs_1$	$rs_2$	...	$rs_M$		
$id_1^{(0)}$	1	0	...	0	-0.4	
$id_2^{(0)}$	0	2	...	1	1.1	
$\vdots$	$\vdots$	$\vdots$		$\vdots$	$\vdots$	
$id_{N_c}^{(0)}$	2	0	...	1	0.6	

Treatment group



$N_1$ individuals	$M$ SNPs					Trait
	$rs_1$	$rs_2$	...	$rs_M$		
$id_1^{(1)}$	2	1	...	0	-0.8	
$id_2^{(1)}$	0	0	...	2	0.8	
$\vdots$	$\vdots$	$\vdots$		$\vdots$	$\vdots$	
$id_{N_t}^{(1)}$	0	2	...	0	1.3	

Association analysis

Association analysis

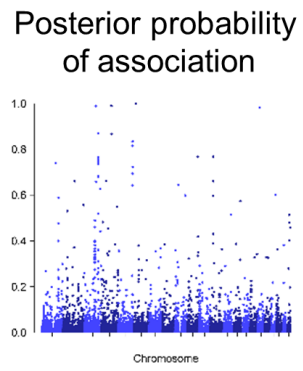
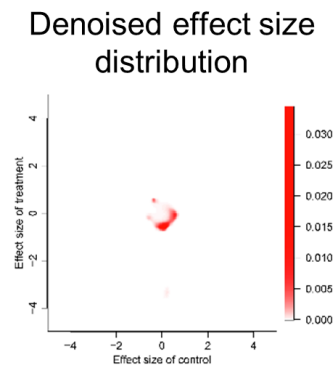
Summary statistics

$M$ SNPs	$b^{(0)}$	s.e.
$rs_1$	-0.2	0.3
$rs_2$	0.4	0.1
$\vdots$	$\vdots$	$\vdots$
$rs_M$	0.3	0.5

Summary statistics

$M$ SNPs	$b^{(1)}$	s.e.
$rs_1$	-0.3	0.6
$rs_2$	0.3	0.2
$\vdots$	$\vdots$	$\vdots$
$rs_M$	0.6	0.3

Empirical Bayes estimation



Adjusted effect size estimates

	$w^{(0)}$	$w^{(1)}$
$rs_1$	-0.1	-0.2
$rs_2$	0.3	0.2
$\vdots$	$\vdots$	$\vdots$
$rs_M$	0.2	0.5

Optimal discovery procedure

Significant SNPs



