

Smooth-threshold multivariate genetic prediction incorporating gene-environment interactions

Masao Ueki*,1, Gen Tamiya^{$+,\pm,\$$} and for Alzheimer's Disease Neuroimaging Initiative²

*School of Information and Data Sciences, Nagasaki University, 1-14 Bunkyo-machi, Nagasaki 852-8521, Japan, [†]Tohoku University Graduate School of Medicine, 2-1, Seiryo-machi, Aoba-ku, Sendai, Miyagi, 980-8575, Japan, [‡]Statistical Genetics Team, RIKEN Center for Advanced Intelligence Project, Nihonbashi 1-chome Mitsui Building 15F, 1-4-1 Nihonbashi, Chuo-ku, Tokyo, 103-0027, Japan, [§]Tohoku Medical Megabank Organization, Tohoku University, 2-1, Seiryo-machi, Aoba-ku, Sendai, Miyagi, 980-8573, Japan

ABSTRACT We propose a genetic prediction modeling approach for genome-wide association study (GWAS) data that can include not only marginal gene effects but also gene-environment (GxE) interaction effectsi.e., multiplicative effects of environmental factors with genes rather than merely additive effects of each. The proposed approach is a straightforward extension of our previous multiple-regression-based method, STMGP (smooth-threshold multivariate genetic prediction), with the new feature being that genome-wide test statistics from a GxE interaction analysis are used to weight the corresponding variants. We develop a sim-7 ple univariate regression approximation to the GxE interaction effect that allows a direct fit of the STMGP 8 framework without modification. The sparse nature of our model automatically removes irrelevant predic-9 tors (including variants and GxE combinations), and the model is able to simultaneously incorporate multiple 10 environmental variables. Simulation studies to evaluate the proposed method in comparison with other mod-11 eling approaches demonstrate its superior performance under the presence of GxE interaction effects. We 12 illustrate the usefulness of our prediction model through application to real GWAS data from the Alzheimer's 13 Disease Neuroimaging Initiative (ADNI). 14

KEYWORDS

Genetic prediction Geneenvironment interaction Smooth thresholding Downloaded from https://academic.oup.com/g3journal/advance-article/doi/10.1093/g3journal/jkab278/6343458 %y UnWeBity oPS&utterP California u&eRow 1 a Agguet

INTRODUCTION

15

Although discovery of genetic risk factors for disease is an important goal of genome-wide association studies (GWAS), predicting disease development or related traits is an important task for applying GWAS results in precision medicine. Many researchers have explored algorithms for accurate genetic prediction based on GWAS data with a large number of single nucleotide polymorphisms (SNPs) (Purcell *et al.* 2009; Evans *et al.*2009; Yang *et al.* 2011; Makowsky *et al.* 2013; de Los Campos *et al.*

¹Corresponding author: School of Information and Data Sciences, Nagasaki University, 1-14 Bunkyo-Machi, Nagasaki 852-8521, Japan. E-mail: uekimrsd@nifty.com ²Data used in preparation of this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database (adni.loni.usc.edu). As such, the investigators within the ADNI contributed to the design and implementation of ADNI and/or provided data but did not participate in analysis or writing of this report. A complete listing of ADNI investigators can be found at: http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf 2013; Chatterjee et al. 2013; Dudbridge 2013; Maier et al. 2015; Moser et al. 2015; Vilhjálmsson et al. 2015; Privé et al. 2019), but no model has been found that performs universally well with all data, and performance is highly dependent on the datagenerating mechanism (Cherlin et al. 2018). Popular models are linear in the variants (or SNPs), such as Purcell's gene score (Purcell et al. 2009) and genomic best linear unbiased prediction (BLUP) (Yang et al. 2011). As an alternative, we developed a statistical method for genetic prediction modeling called smooth-threshold multivariate genetic prediction (STMGP) (Ueki and Tamiya 2016), and Takahashi et al. (2020) recently demonstrated that the performance of STMGP was superior to that of other genetic prediction methods for predicting status of depression with actual GWAS data. STMGP is a sparse modeling method based on a multiple linear regression model such as the lasso (Tibshirani 1996) or the elastic net (Zou and Hastie 2005), and it is able to account for the ultrahigh dimensionality of the $p \gg n$ situation by filtering variants based on the corresponding marginal-effect p-values calculated from univariate regressions arising from a genome-wide scan. Sparseness is achieved

© The Author(s) (2021). Published by Oxford University Press on behalf of the Genetics Society of America. This is an Open Access article distributed under the terms of the Creative Commons Attribution-Non-Commercial-NoDerivs licence (<u>http://creativecommons.org/licenses/by-nc-nd/4.0/</u>), which permits non-commercial reproduction and distribution of the work, in any medium, provided the original work is not altered or transformed in any way, and that the work is properly cited. For commercial re-use, please <u>contactjournals.permissions@oup.com</u> 280 220 221

Manuscript compiled: Saturday 31st July, 2021

[🚰] G Genes | Genomes | Genetics

by ignoring irrelevant variants; the corresponding regression coefficient estimates are set to zero as a result of shrinkage based on the strength of the marginal effect through the smooth-threshold estimating equations developed by Ueki (2009). STMGP also automatically tunes the prediction model by a C_p -type model selection criterion (as with the Akaike information criterion (Akaike 1973)), where the tuning parameter corresponds to the cutoff or threshold value for the marginal *p*-values that determines which effects to filter. The proposed C_p -type criterion based on Stein's unbiased risk estimation (SURE, Stein 1981; Ye 1998; Efron 2004) 10 has a closed-form expression and is a computationally efficient 11 alternative to cross-validation that is often used to choose a p-12 value cutoff in the genetic prediction context (Purcell et al. 2009; 13 Warren et al. 2013). 14

Recent advances in data platforms now make it possible 15 to integrate feature variables other than variants, such as those associated with lifestyle, clinical variables, imaging, etc. 17 The simplest integration is to enter everything as an additive 18 term in a multiple linear regression model as implemented in 19 Ueki and Tamiya (2016); Takahashi et al. (2020). While such an 20 additive modeling approach is simple and straightforward, there 21 may be cases where other approaches are more appropriate. 22 23 One example is gene-environment (GxE) interaction, which has 24 received attention recently as one potential candidate to unveil the missing heritability problem (Manolio et al. 2009; Maher 25 2008; Manolio 2013). With GxE interaction, the model to be 26 estimated is no longer simply additive; rather, it involves terms 27 that are multiplicative in the covariates. Many investigations 28 have aimed at discovering genetic factors that contribute to GxE 29 interactions in disease risk (Kraft et al. 2007; Ober and Vercelli 30 2011; Kraft and Aschard 2015; McAllister et al. 2017; Ritchie et al. 31 2017; Kooperberg and LeBlanc 2008; Hamza et al. 2011; Sung et al. 32 2014; Aschard et al. 2012; Khoury 2017; Kraft and Aschard 33 2015; Sung et al. 2016; Gauderman et al. 2017; Moore et al. 2018; 34 Osazuwa-Peters et al. 2020): the approach using GWAS data is 35 sometimes called a genome-wide environment interaction study 36 (GWEIS) (Meijsen et al. 2018; Ueki et al. 2019; Arnau-Soler et al. 37 2019). The need for GxE interactions depends on the data and target traits, but as with variant discovery, it would be beneficial to have a model for genetic prediction also that can incorporate 40 GxE interactions (Aschard 2016). However, currently the number 41 of such studies is very limited, especially with respect to human 42 disease prediction. 43

To address this issue, we present a straightforward extension 44 of our STMGP method to allow incorporation of GxE interaction effects for building a genetic prediction model using large-46 scale genome-wide SNP data in conjunction with environmental 47 variables. The proposed method can incorporate multiple envi-48 ronmental variables. The STMGP method requires as input the 49 marginal association *p*-values from univariate regression models 50 for each individual variant. This requirement implies that GxE 51 interaction can be fit directly in the STMGP framework if it is expressed in a univariate regression model. The standard univariate 53 GxE interaction model for variant *j* in *n* samples is 54

$$y_i = \mu_i + \epsilon_i = \beta_{0i} + \beta_{1i}E_i + \beta_{2i}G_{ij} + \beta_{3i}E_iG_{ij} + \epsilon_i,$$

where i = 1, ..., n. This model contains three terms: E_i , G_{ij} , and E_iG_{ij} . Here, y_i is the response variable, μ_i is the conditional mean of y_i , E_i is the environmental variable, G_{ij} is the *j*th variant (j = 1, ..., p), p is the number of all variants, ϵ_i is the error variable, and β_{0j} , β_{1j} , β_{2j} , and β_{3j} are the corresponding regression coefficients. In general, removing either E_i or G_j will change the regression coefficient estimate of the GxE interaction term (see Appendix for additional discussion). In this sense, the three terms — E_i , G_{ij} , and E_iG_{ij} — are considered one set, meaning that the GxE interaction effects cannot be represented by a univariate model. To overcome this issue, we propose a simple approximation by a univariate regression model (the rationale is given in the "Materials and Methods" section),

$$y_i = \beta_{0j} + \beta_{1j}\widetilde{E}_i + \beta_{3j}\widetilde{E}_iG_{ij} + \epsilon_i,$$

in which \tilde{E}_i is the centered value of E_i , i.e. $\tilde{E}_i = E_i - \bar{E}$ with \bar{E} the sample mean of E_1, \ldots, E_n . In words, $\beta_{2j}G_{ij}$ is simply removed from the standard model and \tilde{E}_i is used instead of E_i . As a result of this approximation, a one-to-one correspondence is made between the regression coefficient β_{3j} and the single predictor variable E_iG_{ij} . Thus, the STMGP method can now incorporate the GxE interaction directly.

MATERIALS AND METHODS

We use vector and matrix notation. Let $y = (y_1, ..., y_n)^T$, $\mu = (\mu_1, ..., \mu_n)^T$, $E = (E_1, ..., E_n)^T$, and $G_j = (G_{1j}, ..., G_{nj})^T$ (j = 1, ..., p). We first briefly explain the STMGP framework (Ueki and Tamiya 2016), then we present our proposed approach.

STMGP framework

Consider the linear multiple regression model, $y = \mu + \epsilon$, where $\mu = X\beta$, $\epsilon = (\epsilon_1, \dots, \epsilon_n)^T$ is the error vector, X is an *nxp*dimensional design matrix, and β is the corresponding vector of p regression coefficients. In application to GWAS data without GxE interactions, we set $X = (G_1, \ldots, G_p)$. Note that p is much larger than *n* in typical GWAS data—i.e. $p \gg n$. Sparse modeling in which some of the regression coefficients are set to zero is often used in GWAS (Hoggart et al. 2008; Ayers and Cordell 2010; Abraham et al. 2013; Lello et al. 2018; Privé et al. 2019). If disease-susceptibility SNPs show relatively large marginal signals, marginal association screening effectively reduces the dimensionality. The polygenic score, including the gene score method (Purcell et al. 2009) and its multivariate generalization (Warren et al. 2013), uses upper-ranked SNPs with marginal association as predictors to build the prediction model. The former uses independent SNPs after pruning on the basis of LD (linkage disequilibrium), which means that LD is not modeled.

The STMGP method (Ueki and Tamiya 2016) is a variant of the multivariate gene score method (Warren *et al.* 2013), which is essentially the multiple regression model for the upper-ranked SNPs, and it accounts for correlations among SNPs by not including LD-based pruning. Let $T_j(y, X)$ denote a test statistic for marginal association that takes a nonnegative value. Examples of $T_j(y, X)$ include the squared Pearson's correlation and the *F* statistic. Let t > 0 be a cutoff value for $T_j(y, X)$ defining inclusion of SNPs. The cutoff value *t* corresponds to a quantile of the null distribution of $T_j(y, X)$, as in hypothesis testing. The linear multiple regression after marginal association screening uses X_j satisfying $T_j(y, X) > t$ in the model. Without loss of generality, assume that a large value of $T_j(y, X)$ indicates stronger marginal association. Multiple regression after marginal association screening can be expressed by

$$\begin{split} \hat{\mu} &= X\hat{\beta}, \\ \hat{\beta} &= \begin{pmatrix} \hat{\beta}_{\mathcal{A}} \\ \hat{\beta}_{\mathcal{A}^{c}} \end{pmatrix} = \begin{pmatrix} (X_{\mathcal{A}}^{T}X_{\mathcal{A}})^{-1}X_{\mathcal{A}}^{T}y \\ 0 \end{pmatrix}, \\ \mathcal{A} &= \{j: T_{j}(y, X) > t\}, \end{split}$$
(1)

where $X_{\mathcal{A}} = (X_j)_{j \in \mathcal{A}}$ and \mathcal{A}^c indicates the complement set of \mathcal{A} . Note that the above procedure is similar to sure independence screening (Fan and Lv 2008), which uses predictor variables that are upper-ranked in marginal association analyses. The procedure (1) is feasible for $p \gg n$ data and is useful in building a predictive model. In view of the normal equations, it can be seen that $\hat{\beta}$ in (1) satisfies, for j = 1, ..., p,

$$(1 - \hat{D}_j) \{ X_j^T (X\hat{\beta} - y) \} + \hat{D}_j \hat{\beta}_j = 0,$$
(2)

8 or, in vector form,

$$(I_p - \hat{D}) \{ X^T (X\hat{\beta} - y) \} + \hat{D}\hat{\beta} = 0,$$

where $\hat{D}_i = 1\{T_i(y, X) \le t\}$, where $1\{\cdot\}$ denotes the indicator 9 function, $\hat{D} = \text{diag}(\hat{D}_i : j)$, and I_p is the *p*-dimensional identity 10 matrix. Obviously, for $j \in A^c$, $\hat{D}_j = 1$ and (2) reduces to $\hat{\beta}_j = 0$, 11 i.e. a sparse solution; for $j \in A$, $\hat{D}_j = 0$ and the above normal 12 equations reduce to $X_A^T(X_A \hat{\beta}_A - y) = 0$ because $\hat{\beta}_{A^c} = 0$. These 13 are the normal equations for an ordinary least squares regression 14 with design matrix X_A . The resulting prediction process forms $\hat{\mu}(y) = X_A \hat{\beta}_A = X_A (X_A^T X_A)^{-1} X_A^T y$, which is discontinuous in *y* 15 16 due to the thresholding induced by \hat{D}_i . 17

The main innovative idea in STMGP is to replace the discontinuous thresholding \hat{D}_j in (2) with a smooth thresholding using the smooth-threshold estimating equations (STEE) proposed by Ueki (2009). Following Ueki (2009), $\hat{D}_j = 1\{T_j(y, X) \le t\}$ is replaced by an adaptive lasso smooth-thresholding function

$$\check{D}_{j} = \min[1, \{t/T_{j}(y, X)\}^{\frac{1+\gamma}{2}}],$$
(3)

where $\gamma > 0$ is a tuning parameter. This smooth-thresholding 23 function is chosen so as to be identical to the adaptive lasso es-24 timator under the simplest least squares regression of $y = \beta + \epsilon$ 25 (Ueki 2009). If $T_j(y, X) \leq t$ (or $j \in A^c$), $\dot{D}_j = 1$, producing a zero-26 valued regression coefficient; if $T_i(y, X) > t$ (or $i \in A$), $\check{D}_i < 1$ 27 producing a nonzero regression coefficient. Therefore, the condi-28 tion for a sparse solution with D_i is the same as that with D_i . Note 29 that \check{D}_i is monotonically decreasing in $T_i(y, X)$, so regression coef-30 ficients having large $T_i(y, X)$ are penalized to a lesser extent than 31 those having small $T_i(y, X)$. 32

For a given screening cutoff value t > 0, which gives a SNP set $A = \{j : T_j(y, X) > t\}$, the estimates of the *p* regression coefficients are

$$\begin{split} \check{\boldsymbol{\beta}} &= \begin{pmatrix} \check{\boldsymbol{\beta}}_{\mathcal{A}} \\ \check{\boldsymbol{\beta}}_{\mathcal{A}^{c}} \end{pmatrix} \\ &= \begin{pmatrix} \{(I_{|\mathcal{A}|} - \check{\boldsymbol{D}}_{\mathcal{A}})(X_{\mathcal{A}}^{T}X_{\mathcal{A}} + \lambda I_{|\mathcal{A}|}) + \tau\check{\boldsymbol{D}}_{\mathcal{A}}\}^{-1}(I_{|\mathcal{A}|} - \check{\boldsymbol{D}}_{\mathcal{A}})X_{\mathcal{A}}^{T}y \\ & 0 \end{pmatrix} \end{split}$$

$$(4)$$

where $|\mathcal{A}|$ is the cardinality of \mathcal{A} . The non-negative tuning parameters γ and τ are set to 1 and $n/\sqrt{\log n}$, respectively, following previous studies (Ueki and Tamiya 2016; Takahashi *et al.*) 2020), and $\lambda > 0$ is a small constant to avoid singularity of $X_{\mathcal{A}}^T X_{\mathcal{A}}$. The corresponding prediction of y_i is then $\check{\mu}_i(y) = X_i^T \check{\beta}$, where \check{D}_j is an adaptive lasso smooth-thresholding function defined as $\check{D}_j = \min[1, \{t/T_j(y, X)\}^{\frac{1+\gamma}{2}}]$. Since $\check{D}_j = 1$ if and only if $T_j(y, X) \leq t$, the screened set \mathcal{A} with \check{D}_j is the same as that with $\hat{D}_j = 1\{T_j(y, X) \leq t\}$. It can be seen that \check{D}_j replaces the discontinuous screening process \hat{D}_j by a continuous function. As a result, $\check{\mu}_i(y)$ turns out to be continuous in y, enabling stable model selection (Breiman 1996).

According to Ueki (2009); Ueki and Tamiya (2016), the regression coefficients for the screened set in (4) can equivalently be considered as the solution of the generalized ridge regression with loss $||y - X_{\mathcal{A}}\beta_{\mathcal{A}}||^2 + \sum_{j \in \mathcal{A}} \beta_j^2 W_j$, in which $W_j = \lambda + \tau \check{D}_j / (1 - \chi)^2 + \chi \check{D}_j /$ \check{D}_i). The ridge weight for each predictor variable, W_i , represents the uncertainty of the marginal association screening. If the marginal association is very weak, $D_i \approx 1$ and W_i is large, and the corresponding regression coefficient is strongly shrunken towards zero. If the marginal association is strong, $D_i \approx 0$ and $W_i \approx \lambda$, and the corresponding regression coefficient is less penalized. Continuity due to the smooth thresholding also allows computation of a C_{v} -type model selection criterion using SURE. The C_p -type criterion enables a computationally efficient choice of optimal *p*-value cutoff from the perspective of model selection. Details are provided in the Supplementary Material of Ueki and Tamiya (2016). We now outline the STMGP algorithm for $X = (G_1, ..., G_p)$.

Outline of the STMGP algorithm

- Step 1. Perform single-SNP association analysis for *p* SNPs with a univariate model for each SNP.
- Step 2. Retain SNPs whose single-SNP association *p*-value is less than α_{max} .
- Step 3. Fix $\gamma = 1$ and $\tau = n/\sqrt{\log n}$, and select an optimal α from candidate values in $[\alpha_{\min}, \alpha_{\max}]$ by minimizing the C_p -type criterion:

$$C(\alpha) = \sum_{i=1}^{n} \{y_i - \check{\mu}_i(\alpha)\}^2 + 2\hat{\sigma}^2 \text{GDF}(\alpha).$$

Step 4. Compute $\check{\beta}$ in (4) by using the selected α in Step 3.

Here, $\mu_i(\alpha)$ denotes the predicted value for the *i*th subject at the *p*-value threshold α corresponding to the test statistic threshold *t*; α_{max} is the maximum *p*-value in the search, which is set to make the expected number of screened SNPs to be on the order of *n* in practice; $\hat{\sigma}^2$ is an error variance estimate; and GDF(α) denotes the generalized degrees of freedom (Ye 1998; Efron 2004). The univariate model for the *j*th variant G_j (j = 1, ..., p) in Step 1 is

$$\mu_{01} = 1_n \beta_{0j} + G_j \beta_{1j}. \tag{5}$$

Step 3 outputs estimates of regression coefficients, $\check{\beta}_0, \check{\beta}_1, \ldots, \check{\beta}_p$, for the intercept and each variant, which allows computation of the prediction model in an additive form. Some of the regression coefficients $\check{\beta}_1, \ldots, \check{\beta}_p$ can be exactly zero (i.e. sparsity). The predicted value for a new individual who has variants $(G_j^*)_{j=1,\ldots,p}$ can be calculated as $\check{\beta}_0 + \sum_{j=1}^p G_j^* \check{\beta}_j$. The above method assumes a linear regression model for a quantitative phenotype. For a binary phenotype, a logistic regression model is used.

Incorporating GxE interactions with univariate regression ap proximation

³ In what follows, we describe our procedure to incorporate GxE

⁴ interactions into the STMGP framework. Consider the standard ⁵ GxE interaction model for the *j*th variant G_j and an environmental

6 variable *E*,

$$\mu_{0123} = 1_n \beta_{0j} + E \beta_{1j} + G_j \beta_{2j} + (G_j \circ E) \beta_{3j}, \tag{6}$$

where \circ denotes the Hadamard product—i.e. the *i*th element of 7 $(G_i \circ E)$ is given by $G_{ii}E_i$. As seen in Steps 1 and 2 of the STMGP algorithm, because the STMGP framework requires input of mul-9 10 tiple predictors that pass a marginal association *p*-value threshold 11 from each univariate regression model, the above GxE interaction model does not directly fit the STMGP framework due to there 12 being two regression coefficients — β_{2i} and β_{3i} — that associate 13 with G_i . For example, if β_{2i} is highly significant but β_{3i} is not, 14 it is uncertain whether we may include only G_j , because β_{2j} dif-15 fers from the regression coefficient of G_i in the univariate regres-16 sion model without interaction term $(G_i \circ E)$. In contrast, if β_{3i} 17 is highly significant but β_{2i} is not, then it is unclear whether we need $(G_i \circ E)$ only, for the same reason. Furthermore, including 19 both $(G_i \circ E)$ and G_i might reduce predictive power by increasing 20 the number of predictors included: in other words, the curse of 21 dimensionality. 22

We propose a simple approximation to the above GxE inter-23 action model by using a univariate regression model to elimi-24 nate these complications. To this end, we assume independence 25 between E and each G_i . Such assumption is sometimes made 26 in the literature on GxE interaction (Chatterjee and Carroll 2005; 27 Mukherjee and Chatterjee 2007), and it is reasonable for many real 28 GWAS data as the majority of variants have small marginal effects 29 on environmental factors. Our proposed method (the main result) 30 is simply to use the following univariate regression model instead 31 of (6): 32

$$\mu_{013} = 1_n \beta_{0j} + \tilde{E} \beta_{1j} + (G_j \circ \tilde{E}) \beta_{3j}, \tag{7}$$

in which \tilde{E} is the centered E as defined previously. In the Ap-33 pendix we show that, under independence between G_i and E, the 34 least squares estimate of the regression coefficient of $(G_i \circ E)$ in 35 (6) is approximated by that of $(G_i \circ \widetilde{E})$ in (7). This implies a one-36 to-one correspondence between the effects of the regression coef-37 ficient of $(G_i \circ E)$ in (6) and that of the single predictor $(G_i \circ \tilde{E})$. 38 As a consequence, the STMGP framework can be directly applied 39 by setting the following design matrix with 2*p* predictors: 40

$$X = (G_1, \ldots, G_p, G_1 \circ \widetilde{E}, \ldots, G_p \circ \widetilde{E}).$$

If we have *m* environmental variables, E_1, \ldots, E_m , we may set

$$X = (G_1, \ldots, G_p, G_1 \circ \widetilde{E}_1, \ldots, G_p \circ \widetilde{E}_1, \ldots, G_1 \circ \widetilde{E}_m, \ldots, G_p \circ \widetilde{E}_m),$$

which has (1 + m)p predictors. To implement this proposal, we simply include an additional procedure into Steps 1 and 2 above. The following is the modification to include *m* environmental variables.

46 Steps 1 and 2 of STMGP algorithm modified to incorporate GxE in-47 teractions with m environmental variables E_1, \ldots, E_m

Step 1': Perform single-SNP association analysis for each of the p SNPs with a univariate model for each variant, and perform

- ⁵⁰ SNP $x\tilde{E}_k$ interaction analysis for each of the *p* SNPs and \tilde{E}_k
- with the model (7) (k = 1, ..., m), where $\tilde{E}_k = E_k \bar{E}_k 1_n$ with
- \bar{E}_k the sample mean of E_k .

y

Step 2': Screen (retain) SNPs on the basis of single-SNP association *p*-values, and screen SNP–environmental variable pairs on the basis of SNPx \tilde{E}_k interaction *p*-values (k = 1, ..., m) at α_{max} .

The above steps are easily performed with PLINK (Purcell *et al.* 2007; Chang *et al.* 2015), as follows. Prepare the centered environmental variable in a covariate file, say environment.cov. Then, the PLINK command option is --linear --covar environment.cov --interaction --parameters 1,2,3

--tests 1,3. It is also possible to include additional covariates. We have implemented the above algorithm in our STMGP package. We have also implemented a prediction model for binary traits with a logistic regression model based on the method developed in Ueki and Tamiya (2016).

Simulation study

To examine the performance of the proposed method, we conducted simulation studies based on real SNP-GWAS data analogous to those of Takahashi et al. (2020). We used an ADNI-GWAS dataset obtained from the publicly available ADNI database (adni.loni.usc.edu). The ADNI was launched in 2003 as a publicprivate partnership led by Principal Investigator Michael W. Weiner, MD. The goal of the ADNI has been to test whether serial magnetic resonance imaging (MRI), positron emission tomography (PET), other biological markers, and clinical and neuropsychological assessments can be combined to measure the progression of mild cognitive impairment (MCI) and early Alzheimer's disease (AD). For up-to-date information, see www.adni-info.org. The ADNI is an ongoing, longitudinal study with the primary purpose being to explore the association of genetic and neuroimaging information with late-onset Alzheimer's disease (LOAD). The study investigators recruited subjects older than 65 years of age comprising about 400 subjects with mild cognitive impairment (MCI), about 200 subjects with Alzheimer's disease (AD), and about 200 healthy controls. Each subject was followed for at least 3 years. During the study period, the subjects were assessed with magnetic resonance imaging (MRI) measures and psychiatric evaluation to determine the diagnostic status at each time point.

The ADNI-GWAS data were obtained from 818 DNA samples of ADNI1 participants by using the Illumina Human 610-Quad genotyping array (Shen *et al.* 2014). The data initially included 620,901 SNPs. We included the *apolipoprotein E (APOE)* SNPs rs429358 and rs7412 in our analysis. We used data from 684 non-Hispanic Caucasian samples after we excluded one pair showing cryptic relatedness (revealed by the PLINK pairwise $\hat{\pi}$ statistic being greater than 0.125) (Purcell *et al.* 2007), and we excluded subjects whose reported sex did not match the sex inferred from X-chromosome SNPs. We then applied further quality control measures by excluding SNPs with missing genotype rate > 0.1, Hardy–Weinberg equilibrium test *p*-value < 10⁻⁶, and MAF < 5%; the total number of remaining SNPs was 528,984, which is the value of *p* for this analysis.

For the 684 individuals, given that the above real genotype data remain fixed, we artificially generated a quantitative trait, which was used as a target variable to be predicted. We also simulated two environmental variables (sex, E_1 , and years of education, E_2) as follows. E_1 was generated from a Bernoulli distribution with success probability 0.5. E_2 was generated from a standard normal distribution. Both variables were standardized to have mean zero and variance 1 in the generated sample. First, we denote by p_0 the number of causal variants for the main effects of genes, GxE_1 effects, and GxE_2 effects; note that the p_0

variants of each type are not the same. The corresponding $3p_0$ regression coefficients, β_i^* ($j = 1, ..., 3p_0$), were generated from 2 pre-specified distributions. Specifically, the first p_0 regression co-3 efficients were generated independently and identically from a normal, NEG2 (normal-exponential-gamma with shape parameter 2), or Laplace distribution with mean zero and variance h_C^2 ; the second p_0 regression coefficients were generated independently and identically from a normal, NEG2, or Laplace distribution with 8 mean zero and variance $h_{G \times E_1}^2$; the remaining p_0 regression co-9 efficients were generated independently and identically from a 10 normal, NEG2, or Laplace distribution with mean zero and vari-11 ance $h_{G \times E_2}^2$. Next, we randomly selected $3p_0$ causal variants, 12 $G_1^*, \ldots, G_{3p_0}^*$, from among the *p* SNPs, (G_1, \ldots, G_p) . The first p_0 13 variants $(G_1^*, \ldots, G_{p_0}^*)$ had a nonzero gene main effect, the second p_0 variants $(G_{1+p_0}^*, \ldots, G_{2p_0}^*)$ had a nonzero GxE interaction effect 14 15 with E_1 , and the remaining p_0 variants $(G^*_{1+2p_0}, \ldots, G^*_{3p_0})$ had a 16 nonzero GxE interaction effect with E_2 . 17

Then, the conditional mean was set as

$$\begin{split} \mu_{\text{true}} &= \frac{1}{\sqrt{p_0}} \sum_{j=1}^{p_0} \widetilde{G}_j^* \beta_j^* + \frac{1}{\sqrt{p_0}} \sum_{j=1+p_0}^{2p_0} (\widetilde{G_j^* \circ E_1}) \beta_j^* \\ &+ \frac{1}{\sqrt{p_0}} \sum_{j=1+2p_0}^{3p_0} (\widetilde{G_j^* \circ E_2}) \beta_j^*, \end{split}$$

in which \widetilde{G}_{i}^{*} , $(\widetilde{G}_{i}^{*} \circ E_{1})$, and $(\widetilde{G}_{i}^{*} \circ E_{2})$ denote the corresponding 18 terms standardized to have mean zero and variance one. Fi-19 nally, a quantitative trait was generated as $y = \mu_{true} + \epsilon$, where 20 ϵ is an independently and identically distributed normal ran-21 dom variable with mean zero and variance $1 - h_G^2 - h_{G \times E_1}^2 - h_{G \times E_1}^2$ 22 23 similarly, $\frac{1}{\sqrt{p_0}} \sum_{j=1+p_0}^{2p_0} (\widetilde{G_j^* \circ E_1}) \beta_j^*$ and $\frac{1}{\sqrt{p_0}} \sum_{j=1+2p_0}^{3p_0} (\widetilde{G_j^* \circ E_2}) \beta_j^*$ 25 have mean zero and variance $h_{G \times E_1}^2$ and $h_{G \times E_2}^2$, respectively. Also 26 note that the three terms in μ_{true} and ϵ are mutually indepen-27 dent. Thus, y has mean zero and variance 1, and the triplet 28 $h^2 = (h_G^2, h_{G \times E_1}^2, h_{G \times E_2}^2)$ is referred to as heritability throughout 29 this paper. We considered a total of eight scenarios for h^2 . First, 30 we considered (0.3, 0, 0), (0.6, 0, 0), (0, 0.3, 0), and (0, 0.6, 0), where 31 the first and second are scenarios with gene effect without GxE in-32 teractions, and the third and fourth are scenarios with GxE inter-33 actions only for E_1 . Then we considered four additional scenarios: 34 (0, 0.15, 0.15), (0, 0.3, 0.3), (0, 0, 0.3), (0, 0, 0.6), where the first and 35 second are scenarios with GxE interactions both for E_1 and E_2 , 36 and the third and fourth are scenarios with GxE interactions only 37 for E_2 . 38

We used cross-validation to evaluate the prediction models. 39 The data were randomly divided into two parts: 20% for training data and the remaining 80% for test data. The training dataset 41 was used to build prediction models, and then the prediction ac-42 curacy of each model was evaluated on the basis of how well the 43 simulated quantitative traits in the test dataset were predicted by 11 the trained model. We used the prediction correlation coefficient 45 (PCC) to measure the prediction accuracy. The above procedure 46 was repeated 100 times. We note that the $3p_0$ causal SNPs and true regression coefficients differed for each replicate.

We also considered simulations for prediction of binary traits. A binary trait was generated by dichotomizing the quantitative trait on the basis of whether or not its value exceeded $\Phi^{-1}(0.7)$, in which Φ^{-1} is the standard normal quantile function. With a binary trait, the prediction accuracy of each model was evaluated by the area under the receiver operating characteristic curve (AUC).

Comparisons among prediction models We compared the proposed extension of the STMGP method with other prediction models. We included the usual STMGP without GxE interaction as a competitor; specifically, the STMGP models compared were the STMGP without environmental variables, STMGP with environmental variable E_1 , STMGP with environmental variable E_2 , and STMGP with both environmental variables E_1 and E_2 .

We also compared the proposed STMGP extension with other prediction models based on genomic BLUP. Specifically, we considered the following four genomic BLUP models,

$$\mu_b = 1_n \beta_0 + E_1 \beta_{1,1} + E_2 \beta_{2,1} + \sum_{j=1}^p \widetilde{G}_j \beta_{j,2}, \tag{8}$$

$$\mu_{bge1} = 1_n \beta_0 + E_1 \beta_{1,1} + E_2 \beta_{2,1} + \sum_{j=1}^p \widetilde{G}_j \beta_{j,2} + \sum_{j=1}^p (\widetilde{G}_j \circ E_1) \beta_{j,3},$$
(9)

$$\mu_{bge2} = 1_n \beta_0 + E_1 \beta_{1,1} + E_2 \beta_{2,1} + \sum_{j=1}^p \widetilde{G}_j \beta_{j,2} + \sum_{j=1}^p (\widetilde{G}_j \circ E_2) \beta_{j,3},$$
(10)

$$\mu_{bge12} = 1_n \beta_0 + E_1 \beta_{1,1} + E_2 \beta_{2,1} + \sum_{j=1}^p \widetilde{G}_j \beta_{j,2} + \sum_{j=1}^p (\widetilde{G}_j \circ \bar{E}_{12}) \beta_{j,3},$$
(11)

where $\bar{E}_{12} = (E_1 + E_2)/2$, β_0 and β_1 are fixed effects, and $\beta_{j,2}$ and $\beta_{j,3}$ are random effects that are independently distributed as $N(0, \sigma_G^2)$ and $N(0, \sigma_{G \times E}^2)$, respectively. Similar BLUP models have been considered in previous studies (Moore *et al.* 2018; e Sousa *et al.* 2017). We constructed the prediction model by BLUP implemented in the BGEE package for R (Granato *et al.* 2018) by using the BGEE function with options ite=20000, burn=1000, and thin=3.

Application to prediction of real traits

We applied the proposed extension of the STMGP to the prediction of real traits. All variables were obtained from the ADNIMERGE package for R. We considered four cognitive scores as target traits for prediction: FAQ (Functional Assessment Questionnaire), CDRSB (Clinical Dementia Rating Sum of Boxes), MMSE (Mini-Mental State Examination), and ADAS11 (the 11-item ADAS-cog [Alzheimer's Disease Assessment Scale-Cognitive Subscale]). We used SEX and EDU (years of education) as environmental variables. We also considered two additional covariates, AGE and APOE4 genotype. The latter is a known risk allele for AD development. As with the above simulations, we evaluated prediction accuracy via cross-validation.

First, we randomly divided the 684 individuals into five groups of roughly equal size. Then, one of the five groups was selected as the test set and the remaining groups were used as the training set. Consequently, by repeating this with each group in turn acting as the test set, we had five different test/training sample combinations (i.e. 5-fold cross-validation). For each of the five combinations, we built a prediction model based on the training set and predicted each trait value for the test set with the constructed prediction model.

For each training set, we used 528,984 SNPs as predictors as in the above simulation studies. The prediction models we

compared were STMGP with SEX as the environmental variable,
STMGP with EDU as the environmental variable, and STMGP
with SEX and EDU both as environmental variables. BLUP-based
prediction models are (8)–(11). Since the target traits are cognitive scores, we additionally studied regression models including
APOE4 genotype interaction without other variants; specifically,
we considered the following models without GWAS data:

$$\mu_{l0} = 1_n \beta_0 + SEX\beta_{1,1} + EDU\beta_{2,1},$$
(12)
$$\mu_l = 1_n \beta_0 + SEX\beta_{1,1} + EDU\beta_{2,1} + AGE\beta_{3,1} + APOE4\beta_{4,1},$$
(13)

$$\mu_{lge1} = 1_n \beta_0 + SEX\beta_{1,1} + EDU\beta_{2,1} + AGE\beta_{3,1} + APOE4\beta_{4,1} + APOE4 \circ SEX\beta_{5,1},$$
(14)
$$\mu_{lge1} = 1_1 \beta_2 + SEX\beta_{2,1} + EDU\beta_{2,1} + AGE\beta_{3,1} + APOE4\beta_{4,1} (14)$$

$$\mu_{lge2} = I_n \beta_0 + SEX\beta_{1,1} + EDU\beta_{2,1} + AGE\beta_{3,1} + APOE4\beta_{4,1} + APOE4 \circ EDU\beta_{5,1},$$
(15)

$$\mu_{lge12} = 1_n \beta_0 + SEX\beta_{1,1} + EDU\beta_{2,1} + AGE\beta_{3,1} + APOE4\beta_{4,1} + APOE4 \circ SEX\beta_{5,1} + APOE4 \circ EDU\beta_{6,1}.$$
(16)

Prediction accuracy was evaluated with PCC, which compares the
 predicted value with the actual trait in the test set.

RESULTS

11 Simulation results

Results of the quantitative trait simulation are shown in Figures 1, 2, and Supplementary Figure S1, where each cell exhibits mean PCC and the number of causal variants is $p_0 = 100$, 1000, and 500, respectively.

The first and second scenarios for h^2 , (0.3, 0, 0) and (0.6, 0, 0), 16 are those with gene effects but no GxE interactions. From Fig-17 ures 1 and 2, and Supplementary Figure S1, all methods showed 18 a higher predictive power in the latter scenario than in the former 19 scenario due to the larger heritability. The four STMGP meth-20 ods resulted in comparable predictive power, implying that the 21 inclusion of GxE interactions had virtually no effect on predic-22 tive power, which is a reasonable result because no GxE inter-23 action effects were assumed in the data-generating model. The 24 BLUP models had lower predictive power than the STMGP meth-25 ods, which is also reasonable because only a small proportion 26 of variants was assumed to be causal and the BLUP models do 27 not carry out variable selection. Indeed, by comparing Figures 1 28 and 2 and Supplementary Figure S1, it can be seen that an in-29 crease in the number of causal variants made the difference be-30 tween the STMGP and BLUP methods smaller. The difference in 31 effect size distribution had a non-negligible impact on predictive 32 power. While the BLUP methods assume a normal distribution, 33 the STMGP methods do not rely on the effect size distribution, 34 and the STMGP methods had much higher predictive power than 35 the BLUP methods, in particular, when the effect size distribution was non-normal. The difference between the STMGP and 37 BLUP methods was pronounced under the NEG2 distribution, 38 which has the heaviest tails among the three effect-size distribu-39 tions compared. A similar result was observed in the simulation 40 studies of Takahashi et al. (2020). 41

The third and fourth scenarios for h^2 , (0, 0.3, 0) and (0, 0.6, 0), are those with GxE interactions only for E_1 . As in the scenarios for $h^2 = (0.3, 0, 0)$ and (0.6, 0, 0), all prediction models gave higher predictive power in the latter scenario than in the former scenario. Unlike the scenarios with no GxE interactions $h^2 = (0.3, 0, 0)$ and

(0.6, 0, 0), the STMGP methods incorporating GxE interaction effects had higher predictive power than the STMGP method without GxE interactions. For example, in scenario $h^2 = (0, 0.6, 0)$ under a normal effect-size distribution, the STMGP without GxE interaction produced mean PCC 0.36 (standard deviation 0.26), while the STMGP with GxE interaction on variable E_1 resulted in mean PCC 0.41 (standard deviation 0.22). On the other hand, the STMGP with GxE interaction on variable E₂ resulted in mean PCC 0.37 (standard deviation 0.26), which is comparable with STMGP without GxE interaction. This is reasonable since no GxE interaction effect on variable E_2 was assumed. The STMGP with GxE interaction on both E_1 and E_2 gave mean PCC 0.41 (standard deviation 0.23), a predictive power comparable to that of STMGP with GxE interaction on variable E_1 . Total heritability and the difference in effect size distribution had a similar impact on predictive power in scenarios (0.3, 0, 0) and (0.6, 0, 0). For $p_0 = 100$ and the larger heritability scenario, $h^2 = (0, 0.6, 0)$, or under the NEG2 distribution, STMGP with GxE interaction on variable E_1 tended to produce higher predictive power than the BLUP methods, which is perhaps due to the fact that only a small proportion of variants was assumed to be causal. In the other cases among the third and fourth scenarios (any distribution with other than (0, 0.6, 0) and $p_0 = 100$, or $p_0 = 100$ and NEG2 with any heritability [(0, 0.3, 0) or (0, 0.6, 0)]), the STMGP methods did not always perform better than the BLUP methods.

Results of the additional four scenarios are shown in Supplementary Figures S3, S4, and S5. The first and second scenarios for h^2 , (0,0.15,0.15) and (0,0.3,0.3), are the scenarios with GxE interactions both for E_1 and E_2 . Unlike the scenarios (0,0.3,0) and (0,0.6,0), all three STMGP methods with GxE interaction had comparably higher predictive power than STGMP without GxE interaction. This is reasonable as GxE interaction was assumed for both variables, E_1 and E_2 . The third and fourth scenarios for h^2 , (0,0,0.3) and (0,0,0.6), are those with GxE interactions only for E_2 . The results were similar to those for (0,0.3,0) and (0,0.6,0), in which the role of E_2 was replaced by E_1 .

Results of the binary trait simulation are shown in Figures 3 and 4, and Supplementary Figure S2, in which each cell exhibits the mean AUC. The results were consistent overall with the results of the quantitative trait simulation, but differences in predictive power between methods were smaller than with the quantitative trait simulation.

Prediction of real quantitative trait

Results of predicting the four cognitive scores - FAQ, CDRSB, MMSE, and ADAS11 — as target traits are shown in Table 1, which convey the five PCCs from 5-fold cross-validation. Generally, the prediction accuracy differed across the four traits. By comparing 10 with l, lge1, lge2, and lge12, which correspond to formulae (12)-(16), we see that inclusion of the APOE4 genotype (without genome-wide variants) gave much higher predictive power. However, the observed comparable prediction ability among models l, lge1, lge2, and lge12 implies that the inclusion of an interaction between APOE4 and either SEX or EDU did not impact predictive power. The BLUP methods, s, sge1, sge2, and sge12, resulted in performance that was comparable to those of l, lge1, lge2, and lge12, and did not show any extremely distinctive behavior. Similarly, the STMGP methods did not behave much differently from the other methods, but STMGP with a GxE interaction with EDU (sge2) tended to show slightly higher predictive power and improved upon the STMGP without GxE interaction. In particular, for prediction of FAQ, STMGP with a GxE interac-

(0.6,0,0)_Normal -	0.22	0.19	0.21	0.19	0.05	0.04	0.05	0.04
(0.6,0,0)_NEG2 -	0.41	0.4	0.42	0.41	0.05	0.02	0.04	0.03
(0.6,0,0)_Laplace -	0.36	0.34	0.36	0.34	0.06	0.04	0.05	0.05
(0.3,0,0)_Normal -	0.06	0.05	0.05	0.03	0.03	0.03	0.03	0.03
(0.3,0,0)_NEG2 -	0.16	0.15	0.17	0.16	0.02	0	0.01	0.01
(0.3,0,0)_Laplace -	0.13	0.12	0.13	0.12	0.03	0.02	0.02	0.02
(0,0.6,0)_Normal -	0.36	0.41	0.37	0.41	0.38	0.39	0.38	0.39
(0,0.6,0)_NEG2 -	0.34	0.49	0.35	0.49	0.37	0.38	0.37	0.37
(0,0.6,0)_Laplace -	0.34	0.45	0.35	0.45	0.36	0.37	0.37	0.37
(0,0.3,0)_Normal -	0.24	0.25	0.25	0.26	0.26	0.27	0.27	0.27
(0,0.3,0)_NEG2 -	0.23	0.27	0.23	0.27	0.26	0.26	0.26	0.26
(0,0.3,0)_Laplace -	0.25	0.26	0.26	0.28	0.28	0.28	0.28	0.28
	s	sge1	sge2	sge12	bg	bge1	bge2	bge12

Figure 1 Quantitative trait simulation with $p_0 = 100$. Average predictive correlation coefficient (PCC) for eight models. For each scenario (shown in rows), high values are highlighted in red and low values in white. s: STMGP with E_1 and E_2 as covariates; sge1: STMGP with E_1 and E_2 as covariates and E_1 as environmental variable for GxE interaction; sge2: STMGP with E_1 and E_2 as covariates and E_2 as environmental variable for GxE interaction; sge12: STMGP with E_1 and E_2 as covariates, and E_1 and E_2 as environmental variables for GxE interaction; bg: BLUP with E_1 and E_2 as covariates; bge1: BLUP with E_1 and E_2 as covariates and E_1 as environmental variable for GxE interaction; bge2: BLUP with E_1 and E_2 as covariates and E_2 as environmental variable for GxE interaction; bge12: BLUP with E_1 and E_2 as covariates, and average of E_1 and E_2 as environmental variable for GxE interaction. Scenarios are denoted as $(h_G^2, h_{G \times E_1}^2, h_{G \times E_2}^2)$ _dist, where dist means effect size distribution: Normal, NEG2, or Laplace.

(0.6,0,0)_Normal -	0.54	0.53	0.53	0.52	0.51	0.5	0.51	0.51
(0.6,0,0)_NEG2 -	0.63	0.61	0.62	0.61	0.51	0.5	0.51	0.51
(0.6,0,0)_Laplace -	0.59	0.58	0.59	0.58	0.53	0.52	0.52	0.52
(0.3,0,0)_Normal -	0.5	0.49	0.49	0.49	0.5	0.5	0.51	0.51
(0.3,0,0)_NEG2 -	0.54	0.54	0.54	0.54	0.5	0.5	0.5	0.5
(0.3,0,0)_Laplace -	0.53	0.53	0.52	0.52	0.51	0.51	0.51	0.51
(0,0.6,0)_Normal -	0.67	0.68	0.67	0.68	0.67	0.68	0.68	0.68
(0,0.6,0)_NEG2 -	0.67	0.7	0.66	0.7	0.66	0.67	0.66	0.66
(0,0.6,0)_Laplace -	0.66	0.68	0.66	0.68	0.66	0.67	0.66	0.67
(0,0.3,0)_Normal -	0.62	0.62	0.61	0.62	0.62	0.62	0.62	0.62
(0,0.3,0)_NEG2 -	0.61	0.63	0.62	0.63	0.61	0.62	0.61	0.61
(0,0.3,0)_Laplace -	0.62	0.62	0.62	0.62	0.62	0.63	0.62	0.62
	s	sge1	sge2	sge12	bg	bge1	bge2	bge12

Figure 3 Binary trait simulation with $p_0 = 100$. Average area under the ROC curve (AUC) is shown for eight models. For each scenario (in rows), high values are highlighted in red and low values in white. s: STMGP with E_1 and E_2 as covariates; sge1: STMGP with E_1 and E_2 as covariates and E_1 as environmental variable for GxE interaction; sge2: STMGP with E_1 and E_2 as covariates and E_2 as environmental variable for GxE interaction; sge12: STMGP with E_1 and E_2 as covariates, and E_1 and E_2 as environmental variables for GxE interaction; bg: BLUP with E_1 and E_2 as covariates; bge1: BLUP with E_1 and E_2 as covariates and E_1 as environmental variable for GxE interaction; bge2: BLUP with E_1 and E_2 as covariates and E_2 as environmental variable for GxE interaction; bge12: BLUP with E_1 and E_2 as covariates, and average of E_1 and E_2 as environmental variable for GxE interaction. Scenarios are denoted as $(h_G^2, h_{G \times E_1}^2, h_{G \times E_2}^2)$ _dist, where dist means effect size distribution: Normal, NEG2, or Laplace.

(0.6,0,0)_Normal -	0.02	0.01	0.01	0.01	0.05	0.03	0.04	0.03
(0.6,0,0)_NEG2 -	0.16	0.15	0.16	0.15	0.05	0.04	0.05	0.05
(0.6,0,0)_Laplace	0.03	0.01	0.02	0.01	0.05	0.04	0.04	0.04
(0.3,0,0)_Normal -	0.02	0.01	0.01	0.01	0.01	0.01	0.01	0.01
(0.3,0,0)_NEG2 -	0.05	0.05	0.05	0.05	0.03	0.02	0.03	0.03
(0.3,0,0)_Laplace	-0.01	0	0.01	0	0.02	0.02	0.02	0.02
(0,0.6,0)_Normal -	0.33	0.34	0.34	0.34	0.36	0.37	0.36	0.37
(0,0.6,0)_NEG2 -	0.36	0.41	0.36	0.41	0.38	0.39	0.38	0.38
(0,0.6,0)_Laplace -	0.29	0.29	0.3	0.3	0.32	0.33	0.32	0.33
(0,0.3,0)_Normal -	0.24	0.23	0.25	0.24	0.25	0.26	0.25	0.26
(0,0.3,0)_NEG2 -	0.24	0.25	0.25	0.26	0.27	0.28	0.27	0.28
(0,0.3,0)_Laplace -	0.18	0.18	0.19	0.19	0.22	0.23	0.22	0.23
	e l	scel	see2	sne12	ba	bae1	bae2	bre12

Figure 2 Quantitative trait simulation with $p_0 = 1000$. Average predictive correlation coefficient (PCC) for eight models. See Figure 1 for explanation of scenarios (shown in rows).

(0.6,0,0)_Normal -	0.49	0.49	0.5	0.5	0.52	0.51	0.52	0.51
(0.6,0,0)_NEG2 -	0.54	0.54	0.54	0.53	0.52	0.51	0.51	0.52
(0.6,0,0)_Laplace -	0.51	0.51	0.51	0.5	0.52	0.52	0.52	0.52
(0.3,0,0)_Normal -	0.5	0.5	0.49	0.49	0.5	0.5	0.51	0.5
(0.3,0,0)_NEG2 -	0.51	0.51	0.51	0.51	0.51	0.5	0.51	0.51
(0.3,0,0)_Laplace -	0.51	0.5	0.51	0.51	0.51	0.51	0.51	0.51
(0,0.6,0)_Normal -	0.66	0.66	0.66	0.66	0.66	0.66	0.65	0.66
(0,0.6,0)_NEG2 -	0.67	0.68	0.66	0.68	0.68	0.68	0.68	0.68
(0,0.6,0)_Laplace -	0.64	0.64	0.65	0.65	0.65	0.65	0.65	0.65
(0,0.3,0)_Normal -	0.62	0.61	0.61	0.62	0.61	0.61	0.61	0.61
(0,0.3,0)_NEG2 -	0.61	0.62	0.62	0.62	0.62	0.62	0.62	0.62
(0,0.3,0)_Laplace -	0.6	0.6	0.6	0.6	0.6	0.6	0.6	0.6
	s	sge1	sge2	sge12	bg	bge1	bge2	bge12

Figure 4 Binary trait simulation with $p_0 = 1000$. Average area under the ROC curve (AUC) for eight models. See Figure 3 for explanation of scenarios (shown in rows).

tion with EDU (sge2) gave the highest mean PCC (0.22; standard
deviation 0.07) among the methods. However, the differences
among models were small: for example, the second best mean
PCC was 0.21 for l, lge1, lge2, bg, bge1, bge2, and the mean PCC
for the STMGP without GxE interaction was is 0.20 with standard
deviation 0.08. On the other hand, the STMGPs with GxE interaction with SEX (sge1) or with both SEX and EDU (sge12) produced
lower or more variable prediction results.

The above results indicate the possibility that incorporating GXE interactions leads to improved predictive performance. Of course, whether the predictive performance is improved or not depends on the choice of environmental variable, which was also observed in the simulation studies.

Finally, we checked the validity of the proposed univariate re-14 gression approximation in the real data application. Supplemen-15 tary Figures S9-S16 show the accuracy of the proposed approximation, where each figure gives a scatter plot matrix of *p*-values 17 associated with the GxE interaction term $G_i \circ E$ from models (6) 18 and (7) with environmental variables either centered or not. Since 19 centering of environmental variable *E* does not change the model 20 (6), we only compared three *p*-values: model (6), model (7) with 21 centered *E*, and model (7) with non-centered *E*. Among the fig-22 23 ures, Supplementary Figures S9, S11, S13, and S15 show the p-24 values associated with GxE interaction for SEX as the environmental variable, and Supplementary Figures S10, S12, S14, and 25 S16 show the *p*-values associated with GxE interaction for EDU as 26 the environmental variable. In all figures, the $-\log_{10} p$ -values for 27 the GxE interaction term in the approximate univariate regression 28 (i.e. with no gene main effect) using a centered environmental 29 variable were highly correlated (> 0.99) with the $-\log_{10} p$ -values 30 for the GxE interaction term in the interaction model having a 31 gene main effect. On the other hand, with a non-centered envi-32 ronmental variable the same sets of $-\log_{10} p$ -values for the GxE 33 interaction terms were either less correlated (correlation around 34 0.65 for SEX as E) or uncorrelated (< 0.02 for EDU). These re-35 sults confirm the validity of the proposed univariate regression 36 approximation. 37

38 DISCUSSION

In this paper, we presented a procedure to incorporate GxE in-39 teraction effects into our previously developed genetic modeling 40 approach, the STMGP method. Since the STMGP method relies 41 on univariate regression to screen for high-dimensional predictors, we developed a univariate regression approximation to the 43 GxE interaction model so that the STMGP framework can be di-44 rectly applied without modification. The approximation is simply 45 to use "centered" environmental variables and remove gene main 46 effect terms from the standard GxE interaction regression model. 47 Simulation studies and real data analysis showed that incorporat-48 ing GxE interactions may improve the performance of the STMGP, but, as expected, its effectiveness depends to a great extent on the 50 underlying genetic structure. 51

An important point to note is that genome-wide GxE inter-52 action analysis is more sensitive to model misspecification than 53 marginal association analysis, as pointed out by Voorman et al. 54 (2011); Almli et al. (2014); Ueki et al. (2019). Since the model mis-55 specification issue applies to all GxE interaction analyses, special 56 care should be taken in modeling GxE interaction, such as se-57 lection of the environmental variable. We recommend using the 58 check statistic proposed by Ueki et al. (2019) before performing a 59 GxE interaction analysis; this enables prediction of problematic 60 behavior in the GxE interaction analysis without having to per-61

form the actual genome-wide scan.

Most of the existing genetic prediction models treat genetic data separately from non-genetic data. While the widely used additive models to combine genetic and non-genetic data are simple and easy to handle, there must be situations where non-additive models, such as models with GxE interactions, improve upon the additive models. However, studies have reported low power of GxE interaction analysis (Kraft *et al.* 2007). Nevertheless, analogous to the relationship between an association study and prediction modeling, the goal is not to discover GxE interactions but to have a better prediction model. Low statistical power is not necessarily a severe issue in this context: GxE interactions, even if not genome-wide significant, may be useful in helping to improve predictive power.

DATA AVAILABILITY

All data necessary to reproduce the conclusions are fully presented in the paper. The authors do not have ownership of the data used; the data obtained were collected and are owned by the Alzheimer's Disease Neuroimaging Initiative (ADNI). Researchers may request and access the data through the ADNI website (http://adni.loni.usc.edu/). The authors had no special access privileges to use these data. A computer program for the method proposed in this paper is available from the R package stmgp (version 1.0.4).

ACKNOWLEDGEMENTS

Data collection and sharing for this project were funded by the Alzheimer's Disease Neuroimaging Initiative (ADNI) (National Institutes of Health Grant U01 AG024904) and the DOD ADNI (Department of Defense award number W81XWH-12-2-0012). ADNI is funded by the National Institute on Aging, by the National Institute of Biomedical Imaging and Bioengineering, and through generous contributions from the following: Abb-Vie, Alzheimer's Association; Alzheimer's Drug Discovery Foundation; Araclon Biotech; BioClinica, Inc.; Biogen; Bristol-Myers Squibb Company; CereSpir, Inc.; Eisai Inc.; Elan Pharmaceuticals, Inc.; EliLilly and Company; EuroImmun; F. Hoffmann-La Roche Ltd and its affiliated company Genentech, Inc.; Fujirebio; GE Healthcare; IXICO Ltd.; Janssen Alzheimer Immunotherapy Research & Development, LLC.; Johnson & Johnson Pharmaceutical Research & Development LLC.; Lumosity; Lundbeck; Merck & Co., Inc.; Meso Scale Diagnostics, LLC.; NeuroRx Research; Neurotrack Technologies; Novartis Pharmaceuticals Corporation; Pfizer Inc.; Piramal Imaging; Servier; Takeda Pharmaceutical Company; and Transition Therapeutics. The Canadian Institutes of Health Research is providing funds to support ADNI clinical sites in Canada. Private sector contributions are facilitated by the Foundation for the National Institutes of Health (www.fnih.org). The grantee organization is the Northern California Institute for Research and Education, and the study is coordinated by the Alzheimer's Disease Cooperative Study at the University of California, San Diego. ADNI data are disseminated by the Laboratory for Neuro Imaging at the University of Southern California. Computations for the present work were performed by using the facilities of the Institute of Statistical Mathematics. The authors are grateful for the helpful comments by two referees, associate editor, and Dr. John Cologne.

Trait ^a	Data ^b	10 ^c	lq	lge1 ^e	lge2 ^f	lge12 ^g	s ^h	sge1 ⁱ	sge2 ^j	sge12 ^k	bg [/]	bge1 ^m	bge2 ⁿ	bge12 ⁰
FAQ	CV 1	0.07	0.16	0.15	0.17	0.16	0.11	-0.01	0.15	0.05	0.14	0.13	0.13	0.12
	CV 2	0.17	0.35	0.33	0.36	0.34	0.26	0.24	0.32	0.31	0.32	0.35	0.33	0.33
	CV 3	0.19	0.15	0.15	0.16	0.16	0.19	0.13	0.21	0.15	0.17	0.15	0.18	0.17
	CV 4	0.01	0.26	0.26	0.27	0.27	0.31	0.18	0.24	0.19	0.23	0.28	0.25	0.23
	CV 5	0.08	0.16	0.16	0.10	0.09	0.15	0.14	0.17	0.15	0.17	0.14	0.17	0.15
	mean	0.10	0.21	0.21	0.21	0.20	0.20	0.14	0.22	0.17	0.21	0.21	0.21	0.20
	sd	0.08	0.09	0.08	0.10	0.10	0.08	0.09	0.07	0.09	0.07	0.10	0.08	0.09
CDRSB	CV 1	0.07	0.13	0.13	0.12	0.12	0.21	0.18	0.22	0.17	0.12	0.13	0.10	0.11
	CV 2	0.16	0.38	0.37	0.36	0.35	0.33	0.28	0.33	0.30	0.34	0.36	0.34	0.33
	CV 3	0.22	0.26	0.26	0.26	0.26	0.28	0.26	0.26	0.25	0.25	0.25	0.26	0.27
	CV 4	0.10	0.37	0.37	0.37	0.37	0.44	0.36	0.41	0.31	0.36	0.39	0.37	0.36
	CV 5	0.19	0.27	0.26	0.25	0.22	0.27	0.25	0.28	0.27	0.27	0.25	0.27	0.27
	mean	0.15	0.28	0.27	0.27	0.26	0.31	0.27	0.30	0.26	0.27	0.27	0.27	0.27
	sd	0.06	0.10	0.10	0.10	0.10	0.08	0.06	0.07	0.06	0.09	0.10	0.10	0.10
MMSE	CV 1	0.10	0.27	0.25	0.26	0.25	0.13	0.21	0.18	0.16	0.22	0.23	0.23	0.22
	CV 2	0.19	0.34	0.33	0.33	0.32	0.30	0.33	0.33	0.33	0.29	0.30	0.31	0.30
	CV 3	0.30	0.35	0.35	0.35	0.35	0.28	0.26	0.34	0.35	0.37	0.38	0.36	0.36
	CV 4	0.27	0.35	0.35	0.35	0.36	0.35	0.34	0.39	0.37	0.36	0.37	0.36	0.37
	CV 5	0.17	0.28	0.26	0.28	0.25	0.25	0.23	0.26	0.22	0.29	0.28	0.29	0.27
	mean	0.21	0.32	0.31	0.31	0.31	0.26	0.27	0.30	0.29	0.31	0.31	0.31	0.30
	sd	0.08	0.04	0.05	0.04	0.05	0.08	0.06	0.08	0.09	0.06	0.06	0.05	0.06
ADAS11	CV 1	0.12	0.31	0.32	0.30	0.31	0.30	0.28	0.29	0.26	0.29	0.29	0.28	0.27
	CV 2	0.17	0.30	0.30	0.30	0.30	0.22	0.23	0.24	0.22	0.28	0.27	0.28	0.29
	CV 3	0.15	0.29	0.30	0.29	0.30	0.22	0.26	0.24	0.26	0.29	0.29	0.29	0.29
	CV 4	0.11	0.36	0.36	0.35	0.35	0.29	0.29	0.37	0.29	0.37	0.38	0.35	0.36
	CV 5	0.22	0.34	0.32	0.33	0.32	0.30	0.28	0.34	0.24	0.33	0.31	0.32	0.31
	mean	0.15	0.32	0.32	0.31	0.32	0.27	0.27	0.30	0.25	0.31	0.31	0.30	0.30
	sd	0.04	0.03	0.03	0.03	0.02	0.04	0.02	0.06	0.03	0.04	0.04	0.03	0.03

Table 1 Results of predicting four quantitative traits, FAQ, CDRSB, MMSE, and ADAS11

^a Prediction of each target trait is evaluated by the prediction correlation coefficient (PCC) from 5-fold cross-validation.

^b Data used to calculate PCC (CV 1 – CV 5 denote each cross-validated dataset from 5-fold cross-validation) for each model are shown in row together with mean and standard deviation (sd).

^c Linear regression with SEX and EDU as predictors.

^d Linear regression with SEX, EDU, AGE, and APOE4 as predictors.

^e Linear regression with SEX, EDU, AGE, APOE4, and APOE4xSEX as predictors.

^{*t*} Linear regression with SEX, EDU, AGE, APOE4, and APOE4xEDU as predictors.

^g Linear regression with SEX, EDU, AGE, APOE4, APOE4xSEX, and APOE4xEDU as predictors.

^h STMGP with SEX, EDU, AGE, and APOE4 as covariates.

¹ STMGP with SEX, EDU, AGE, and APOE4 as covariates, and SEX as environmental variable for GxE interaction.

^{*j*} STMGP with SEX, EDU, AGE, and APOE4 as covariates, and EDU as environmental variable for GxE interaction.

^k STMGP with SEX, EDU, AGE, and APOE4 as covariates, and AGE and EDU as environmental variables for GxE interaction.

¹ BLUP with SEX, EDU, AGE, and APOE4 as covariates.

^m BLUP with SEX, EDU, AGE, and APOE4 as covariates, and SEX as environmental variable for GxE interaction.

^{*n*} BLUP with SEX, EDU, AGE, and APOE4 as covariates, and EDU as environmental variable for GxE interaction.

^o BLUP with SEX, EDU, AGE, and APOE4 as covariates, and average of AGE and EDU as environmental variable for GxE interaction.

FUNDING

² This work was supported by JSPS KAKENHI Grant Numbers
 ³ 20K11723 and 20H00576.

4 COMPETING INTERESTS

⁵ The authors declare that there are no conflict of interests.

6 LITERATURE CITED

Abraham, G., A. Kowalczyk, J. Zobel, and M. Inouye, 2013 Perfor mance and robustness of penalized and unpenalized methods

- for genetic prediction of complex human disease. Genetic Epidemiology 37: 184–195.
- Akaike, H., 1973 Information theory and an extension of the max-
- imum likelihood principle. In Proceedings of the 2nd International
 Symposium on Information Theory, Petrov, B. N and and Caski, F.
- 14 (eds.), pp. 267–281, Budapest, Akadimiai Kiado.
- ¹⁵ Almli, L. M., R. Duncan, H. Feng, D. Ghosh, E. B. Binder, *et al.*,
 ¹⁶ 2014 Correcting systematic inflation in genetic association tests
- that consider interaction effects. JAMA Psychiatry 71: 1392–
 1399.
- Arnau-Soler, A., , E. Macdonald-Dunlop, M. J. Adams, T.-K.
 Clarke, *et al.*, 2019 Genome-wide by environment interaction studies of depressive symptoms and psychosocial stress in UK
- ²² biobank and generation scotland. Translational Psychiatry 9.
- Aschard, H., 2016 A perspective on interaction effects in genetic
 association studies. Genetic Epidemiology 40: 678–88.
- Aschard, H., S. Lutz, B. Maus, E. J. Duell, T. E. Fingerlin, *et al.*, 2012
 Challenges and opportunities in genome-wide environmental interaction (GWEI) studies. Human Genetics **131**: 1591–1613.
- Ayers, K. and H. Cordell, 2010 Snp selection in genome-wide and
 candidate gene studies via penalized logistic regression. Genetic Epidemiology 34: 879–891.
- Breiman, L., 1996 Heuristics of instability and stabilization in model selection. Annals of Statistics **24**: 2350–2383.
- ³³ Chang, C. C., C. C. Chow, L. C. Tellier, S. Vattikuti, S. M. Purcell,
 et al., 2015 Second-generation PLINK: rising to the challenge of
 ³⁵ larger and richer datasets. GigaScience 4.
- ³⁶ Chatterjee, N. and R. J. Carroll, 2005 Semiparametric maximum
 ³⁷ likelihood estimation exploiting gene-environment indepen ³⁸ dence in case-control studies. Biometrika **92**: 399–418.
- 39 Chatterjee, N., B. Wheeler, J. Sampson, P. Hartge, S. Chanock,
- *et al.*, 2013 Projecting the performance of risk prediction based
 on polygenic analyses of genome-wide association studies. Nat
 Genetics 45: 400–5.
- ⁴³ Cherlin, S., D. Plant, J. C. Taylor, M. Colombo, A. Spiliopoulou,
 et al., 2018 Prediction of treatment response in rheumatoid
 ⁴⁵ arthritis patients using genome-wide SNP data. Genetic Epi ⁴⁶ demiology **42**: 754–771.
- de Los Campos, G., A. Vazquez, R. Fernando, Y. Klimentidis, and
 D. Sorensen, 2013 Prediction of complex human traits using
- the genomic best linear unbiased predictor. PLOS Genetics 9:
 e1003608.
- ⁵¹ Dudbridge, F., 2013 Power and predictive accuracy of polygenic ⁵² risk scores. PLOS Genetics **9**: e1003348.
- e Sousa, M. B., J. Cuevas, E. G. de Oliveira Couto, P. Pérez-Rodríguez, D. Jarquín, *et al.*, 2017 Genomic-enabled prediction
 in maize using kernel models with genotype x environment interaction. G3:Genes | Genomes | Genetics 7: 1995–2014.
- ⁵⁷ Efron, B., 2004 The estimation of prediction error: covariance
- penalties and ccross-validation. Journal of the American Statis tical Association 99: 619–632.

- Evans, D., P. Visscher, and N. Wray, 2009 Harnessing the information contained within genome-wide association studies to improve individual prediction of complex disease risk. Human Molecular Genetics **18**: 3525–3531.
- Fan, J. and J. Lv, 2008 Sure independence screening for ultrahigh dimensional feature space (with discussion). Journal of the Royal Statistical Society: Series B (Statistical Methodology) **70**: 849–911.
- Gauderman, W. J., B. Mukherjee, H. Aschard, L. Hsu, J. P. Lewinger, *et al.*, 2017 Update on the state of the science for analytical methods for gene-environment interactions. American Journal of Epidemiology 186: 762–770.
- Granato, I., J. Cuevas, F. Luna-Vázquez, J. Crossa, O. Montesinos-López, et al., 2018 BGGE: A new package for genomic-enabled prediction incorporating genotype x environment interaction models. G3:Genes | Genomes | Genetics 8: 3039–3047.
- Hamza, T. H., H. Chen, E. M. Hill-Burns, S. L. Rhodes, J. Montimurro, *et al.*, 2011 Genome-wide gene-environment study identifies glutamate receptor gene GRIN2a as a parkinson's disease modifier gene via interaction with coffee. PLOS Genetics 7: e1002237.
- Hoggart, C., J. Whittaker, M. Iorio, and D. Balding, 2008 Simultaneous analysis of all snps in genome-wide and re-sequencing association studies. PLOS Genetics 4: e1000130.
- Khoury, M. J., 2017 Editorial: emergence of gene-environment interaction analysis in epidemiologic research. American Journal of Epidemiology 186: 751–752.
- Kooperberg, C. and M. LeBlanc, 2008 Increasing the power of identifying gene x gene interactions in genome-wide association studies. Genetic Epidemiology **32**: 255–263.
- Kraft, P. and H. Aschard, 2015 Finding the missing geneenvironment interactions. European Journal of Epidemiology 30: 353–355.
- Kraft, P., Y. C. Yen, D. O. Stram, J. Morrison, and W. J. Gauderman, 2007 Exploiting gene-environment interaction to detect genetic associations. Human Heredity 63: 111–119.
- Lello, L., S. G. Avery, L. Tellier, A. I. Vazquez, G. de los Campos, et al., 2018 Accurate genomic prediction of human height. Genetics 210: 477–497.
- Maher, B., 2008 Personal genomes: The case of the missing heritability. Nature **456**: 18–21.
- Maier, R., G. Moser, G.-B. Chen, S. Ripke, W. Coryell, et al., 2015 Joint analysis of psychiatric disorders increases accuracy of risk prediction for schizophrenia, bipolar disorder, and major depressive disorder. The American Journal of Human Genetics 96: 283–294.
- Makowsky, R., N. Pajewski, Y. Klimentidis, A. Vazquez, C. Duarte, *et al.*, 2013 Power and predictive accuracy of polygenic risk scores. PLOS Genetics 9: e1003348.
- Manolio, T., 2013 Bringing genome-wide association findings into clinical use. Nature Reviews Genetics **14**: 549–558.
- Manolio, T. A., F. S. Collins, N. J. Cox, and D. B. Goldstein, 2009 Finding the missing heritability of complex diseases. Nature **461**: 747–753.
- McAllister, K., L. E. Mechanic, C. Amos, H. Aschard, I. A. Blair, *et al.*, 2017 Current challenges and new opportunities for geneenvironment interaction studies of complex diseases. American Journal of Epidemiology **186**: 753–761.
- Meijsen, J. J., A. Campbell, C. Hayward, D. J. Porteous, I. J. Deary, *et al.*, 2018 Phenotypic and genetic analysis of cognitive performance in major depressive disorder in the generation scotland: Scottish family health study. Translational Psychiatry **8**.

- ¹ Moore, R., , F. P. Casale, M. J. Bonder, D. Horta, *et al.*, 2018 ² A linear mixed-model approach to study multivariate gene-
- environment interactions. Nature Genetics **51**: 180–186.
- ⁴ Moser, G., S. H. Lee, B. J. Hayes, M. E. Goddard, N. R. Wray, et al.,
- 5 2015 Simultaneous discovery, estimation and prediction anal-
- ysis of complex traits using a bayesian mixture model. PLOS
 Genetics 11: e1004969.
- Mukherjee, B. and N. Chatterjee, 2007 Exploiting gene environment independence for analysis of case-control stud-
- ies: An empirical bayes-type shrinkage estimator to trade-off
 between bias and efficiency. Biometrics 64: 685–694.
- ¹² Ober, C. and D. Vercelli, 2011 Gene–environment interactions in
- human disease: nuisance or opportunity? Trends in Genetics
 27: 107–115.
- 15 Osazuwa-Peters, O. L., R. J. Waken, K. L. Schwander, Y. J. Sung,
- P. S. de Vries, *et al.*, 2020 Identifying blood pressure loci whose
 effects are modulated by multiple lifestyle exposures. Genetic
 Epidemiology 44: 629–641.
- Privé, F., H. Aschard, and M. G. B. Blum, 2019 Efficient imple mentation of penalized regression for genetic risk prediction.
 Genetics 212: 65–74.
- Purcell, S., B. Neale, K. Todd-Brown, L. Thomas, M. Ferreira,
 et al., 2007 Plink: a tool set for whole-genome association and
- population-based linkage analyses. The American Journal of
 Human Genetics 81: 559–75.
- Purcell, S., N. Wray, J. Stone, P. Visscher, M. O 'Donovan, *et al.*, 2009 Common polygenic variation contributes to risk of schizophrenia and bipolar disorder. Nature 460: 748–52.
- Ritchie, M. D., J. R. Davis, H. Aschard, A. Battle, D. Conti, *et al.*,
 2017 Incorporation of biological knowledge into the study of
 gene-environment interactions. American Journal of Epidemi ology 186: 771–777.
- Shen, L., P. Thompson, S. Potkin, L. Bertram, L. Farrer, *et al.*,
 2014 Genetic analysis of quantitative phenotypes in ad and mci:
 imaging, cognition and biomarkers. Brain Imaging and Behav-
- ior 8: 183–207.
 Stein, C., 1981 Estimation of the mean of a multivariate normal distribution. Annals of Statistics 9: 1135–1151.
- ³⁹ Sung, Y. J., L. de las Fuentes, K. L. Schwander, J. Simino, and D. C.
- Rao, 2014 Gene–smoking interactions identify several novel
 blood pressure loci in the framingham heart study. American
 Journal of Hypertension 28: 343–354.
- ⁴³ Sung, Y. J., T. W. Winkler, A. K. Manning, H. Aschard, V. Gudna-
- son, *et al.*, 2016 An empirical comparison of joint and stratified
 frameworks for studying g x e interactions: Systolic blood pres-
- sure and smoking in the CHARGE gene-lifestyle interactions
 working group. Genetic Epidemiology 40: 404–415.
- ⁴⁸ Takahashi, Y., M. Ueki, G. Tamiya, S. Ogishima, K. Kinoshita, *et al.*,
- 2020 Machine learning for effectively avoiding overfitting is a
 crucial strategy for the genetic prediction of polygenic psychi atric phenotypes. Translational Psychiatry 10.
- Takane, Y. and H. Yanai, 1999 On oblique projectors. Linear Algebra and its Applications 289: 297–310.
- ⁵⁴ Tibshirani, R., 1996 Regression shrinkage and selection via the lasso. Journal of the Royal Statistical Society: Series B (Methodological) 59: 267–289
- ⁵⁶ ological) **58**: 267–288.
- ⁵⁷ Ueki, M., 2009 A note on automatic variable selection using
 smooth-threshold estimating equation. Biometrika 96: 1005–
 1011.
- ⁶⁰ Ueki, M., M. Fujii, and G. Tamiya, 2019 Quick assessment for sys-
- tematic test statistic inflation/deflation due to null model mis specifications in genome-wide environment interaction studies.
- ² specifications in genome-wide environment interaction studies.

PLOS ONE 14: e0219825.

- Ueki, M. and Y. Kawasaki, 2013 Multiple choice from competing regression models under multicollinearity based on standardized update. Computational Statistics & Data Analysis **63**: 31–41.
- Ueki, M. and G. Tamiya, 2016 Smooth-threshold multivariate genetic prediction with unbiased model selection. Genetic Epidemiology **40**: 233–243.
- Vilhjálmsson, B. J., J. Yang, H. K. Finucane, A. Gusev, S. Lindström, et al., 2015 Modeling linkage disequilibrium increases accuracy of polygenic risk scores. The American Journal of Human Genetics 97: 576–592.
- Voorman, A., T. Lumley, B. McKnight, and K. Rice, 2011 Behavior of qq-plots and genomic control in studies of gene-environment interaction. PLOS ONE **6**: e19416.
- Warren, H., J. Casas, A. Hingorani, F. Dudbridge, and J. Whittaker, 2013 Genetic prediction of quantitative lipid traits: comparing shrinkage models to gene scores. Genetic Epidemiology 38: 72– 83.
- Yang, J., S. H. Lee, M. E. Goddard, and P. M. Visscher, 2011 GCTA: A tool for genome-wide complex trait analysis. The American Journal of Human Genetics 88: 76–82.
- Ye, J., 1998 On measuring and correcting the effects of data mining and model selection. Journal of the American Statistical Association **93**: 120–131.
- Zou, H. and T. Hastie, 2005 Regularization and variable selection via the elastic net. Journal of the Royal Statistical Society: Series B (Statistical Methodology) **67**: 301–320.

22)

APPENDIX

2 Derivation of the univariate regression approximation

³ We consider the GxE interaction model $y_i = \beta_{0j} + E_i\beta_{1j} + G_{ij}\beta_{2j} + E_iG_{ij}\beta_{3j} + \epsilon_i$, where the ϵ_i are independently and identically distributed with mean zero and variance σ_0^2 . Here, we assume that ⁶ G_{ij} and E_i are independent, and that each is independently and ⁷ identically distributed for i = 1, ..., n. We also assume that ⁸ σ_0^2 , $v_{G_j} = \text{Var}(G_{ij})$, and $v_E = \text{Var}(E_i)$ are finite. Let $P_X = X(X^TX)^{-1}X^T$ be the projection matrix onto the column space of ¹⁰ X, and let $Q_X = I_n - P_X$. Then, for *n*-dimensional one-vector 1_n , ¹¹ the operator $Q_{1_n} = I_n - P_{1_n}$ gives centering to have mean zero. ¹² Let $\tilde{E} = Q_{1_n}E$ and $\tilde{G}_j = Q_{1_n}G_j$. Then, for large *n*,

$$n^{-1}\widetilde{G}_j^T\widetilde{E} = \mathbb{E}(n^{-1}\widetilde{G}_j^T\widetilde{E}) + O_p\{\operatorname{Var}(n^{-1}\widetilde{G}_j^T\widetilde{E})^{1/2}\} = O_p(n^{-1/2}).$$
(17)

¹³ Finally, let $W_j = E \circ G_j$ and $\widetilde{W}_j = \widetilde{E} \circ G_j$.

Note that $P_X = P_{n^{-1/2}X}$ for any given matrix X. Thus, $Q_X = Q_{n^{-1/2}X}$ also holds. By Ueki and Kawasaki (2013); Ueki *et al.* (2019), the least squares estimate of regression coefficient β_{3j} in the model $\mu = \mu(E, G_j) = 1_n\beta_{0j} + E\beta_{1j} + G_j\beta_{2j} + (G_j \circ E)\beta_{3j}$, model (6), is

$$\widehat{\beta}_{3j}^{0123}(E,G_j) = \frac{y^T Q_{(1_n,E,G_j)} W_j}{||Q_{(1_n,E,G_j)} W_j||^2} = \frac{n^{-1} y^T Q_{(1_n,E,G_j)} W_j}{n^{-1} ||Q_{(1_n,E,G_j)} W_j||^2}.$$
 (18)

- ¹⁹ Similarly, the least squares estimate of regression coefficient β_{3j} in
- ²⁰ the model $\mu = \mu(E, G_j) = 1_n \beta_{0j} + E \beta_{1j} + G_j \beta_{2j} + (G_j \circ E) \beta_{3j}$ is

$$\widehat{\beta}_{3j}^{013}(E,G_j) = \frac{y^T Q_{(1_n,E)} W_j}{||Q_{(1_n,E)} W_j||^2} = \frac{n^{-1} y^T Q_{(1_n,E)} W_j}{n^{-1} ||Q_{(1_n,E)} W_j||^2}.$$
 (19)

We utilize the decomposition of a projection matrix or blockwise formula (Takane and Yanai 1999, Lemma 3 (iii)), $P_{(A,B)} = P_A + P_{QAB}$ for two matrixes A and B. Note that $P_A P_{QAB} = P_{QA} B P_A = O$ since $Q_A A = O$. Then, $P_{(1_n, E)} = P_{1_n} + P_{\widetilde{E}}$. Using this, and by the blockwise formula again, we have $P_{(1_n, E, G_j)} = P_{1_n} + P_{(\widetilde{E}, \widetilde{G}_j)} = P_{1_n} + P_{\widetilde{E}} + P_{Q_{\widetilde{E}} \widetilde{G}_j} = P_{(1_n, E)} + P_{Q_{\widetilde{E}} \widetilde{G}_j}$. Thus,

$$Q_{(1_n,E,G_j)} = Q_{(1_n,E)} - P_{Q_{\widetilde{E}}\widetilde{G}_j},$$

and applying this identity to (18),

$$\widehat{\beta}_{3j}^{0123}(E,G_{j}) = \frac{n^{-1}y^{T}Q_{(1_{n},E,G_{j})}W_{j}}{n^{-1}W_{j}^{T}Q_{(1_{n},E,G_{j})}W_{j}} \\
= \frac{n^{-1}y^{T}\{Q_{(1_{n},E)} - P_{Q_{\tilde{E}}\tilde{G}_{j}}\}W_{j}}{n^{-1}W_{j}^{T}\{Q_{(1_{n},E)} - P_{Q_{\tilde{E}}\tilde{G}_{j}}\}W_{j}} \\
= \frac{n^{-1}y^{T}Q_{(1_{n},E)}W_{j} - n^{-1}y^{T}P_{Q_{\tilde{E}}\tilde{G}_{j}}W_{j}}{n^{-1}||Q_{(1_{n},E)}W_{j}||^{2} - n^{-1}W_{j}^{T}P_{Q_{\tilde{E}}\tilde{G}_{j}}W_{j}}, \quad (20)$$

which differs from (19) unless $n^{-1}W_j^T P_{Q_{\tilde{E}}\tilde{G}_j}W_j$ and $n^{-1}W_j^T P_{Q_{\tilde{E}}\tilde{G}_j}W_j$ are both negligible. Let $\overline{G}_j = Q_{\tilde{E}}\tilde{G}_j$. The second term of the numerator of (20) can be written as

$$n^{-1}y^{T}P_{\overline{G}_{j}}W_{j} = n^{-1}y^{T}\overline{G}_{j}(\overline{G}_{j}^{T}\overline{G}_{j})^{-1}\overline{G}_{j}^{T}W_{j}$$
$$= (n^{-1}y^{T}\overline{G}_{j})(n^{-1}\overline{G}_{j}^{T}\overline{G}_{j})^{-1}(n^{-1}W_{j}^{T}\overline{G}_{j})$$

To begin with, by (17) the left, middle, and right terms reduce to

$$n^{-1}y^{T}\overline{G}_{j} = n^{-1}y^{T}Q_{n^{-1/2}\widetilde{E}}\widetilde{G}_{j} = n^{-1}y^{T}\widetilde{G}_{j} - \frac{(n^{-1}y^{T}\widetilde{E})(n^{-1}\widetilde{G}_{j}^{T}\widetilde{E})}{||n^{-1/2}\widetilde{E}||^{2}}$$
$$= n^{-1}y^{T}\widetilde{G}_{j} + o_{p}(1), \qquad (21)$$

$$n^{-1}\overline{G}_{j}^{T}\overline{G}_{j} = n^{-1}\widetilde{G}_{j}^{T}Q_{n^{-1/2}\widetilde{E}}\widetilde{G}_{j} = n^{-1}\widetilde{G}_{j}^{T}\widetilde{G}_{j} - \frac{(n^{-1}G_{j}^{T}E)^{2}}{||n^{-1/2}\widetilde{E}||^{2}}$$
$$= n^{-1}\widetilde{G}_{j}^{T}\widetilde{G}_{j} + o_{p}(1), \qquad ($$

$$n^{-1}W_{j}^{T}\overline{G}_{j} = n^{-1}W_{j}^{T}Q_{n^{-1/2}\widetilde{E}}\widetilde{G}_{j} = n^{-1}W_{j}^{T}\widetilde{G}_{j} - \frac{(n^{-1}W_{j}^{T}\widetilde{E})(n^{-1}\widetilde{G}_{j}^{T}\widetilde{E})}{||n^{-1/2}\widetilde{E}||^{2}}$$

= $n^{-1}W_{j}^{T}\widetilde{G}_{j} + o_{p}(1),$ (23)

respectively. Combining (21)–(23), the numerator of (20) reduces to

$$n^{-1}y^{T}Q_{(1_{n},E)}W_{j} - n^{-1}y^{T}P_{Q_{\tilde{E}}\widetilde{G}_{j}}W_{j}$$

= $n^{-1}y^{T}Q_{(1_{n},E)}W_{j} - \frac{(n^{-1}y^{T}\widetilde{G}_{j})(n^{-1}W_{j}^{T}\widetilde{G}_{j})}{n^{-1}\widetilde{G}_{j}^{T}\widetilde{G}_{j}} + o_{p}(1).$ (24)

By analogous calculations, the denominator of (20) reduces to

$$n^{-1} ||Q_{(1_n,E)}W_j||^2 - n^{-1} W_j^T P_{Q_{\tilde{E}}\tilde{G}_j}W_j$$

$$= n^{-1} ||Q_{(1_n,E)}W_j||^2 - \frac{(n^{-1} W_j^T \tilde{G}_j)^2}{(n^{-1} W_j^T \tilde{G}_j)^2} + o_T(1)$$
(26)

$$= n^{-1} ||Q_{(1_n,E)}W_j||^2 - \frac{(n - W_j G_j)}{n^{-1}\widetilde{G}_j^T \widetilde{G}_j} + o_p(1).$$
 (2)

Substituting (24) and (26) into (20),

$$\widehat{\beta}_{3j}^{0123}(E,G_j) = \frac{n^{-1}y^T Q_{(1_n,E)} W_j - \frac{(n^{-1}y^T G_j)(n^{-1}W_j^T G_j)}{n^{-1} \widetilde{G}_j^T \widetilde{G}_j}}{n^{-1} ||Q_{(1_n,E)} W_j||^2 - \frac{(n^{-1}W_j^T \widetilde{G}_j)^2}{n^{-1} \widetilde{G}_j^T \widetilde{G}_j}} + o_p(1).$$
(27)

This approximates (19) if $(n^{-1}y^T \tilde{G}_j)(n^{-1}W_j^T \tilde{G}_j)$ and $(n^{-1}W_j^T \tilde{G}_j)^2$ are both negligible, which, however, might not be true in general.

Instead, we consider the case where *E* is replaced by $\tilde{E} = Q_{1n}E = E - \tilde{E}1_n$ in (18). In this case, the estimate of regression coefficient (19) is

$$\widehat{\beta}_{3j}^{013}(\widetilde{E},G_j) = \frac{y^T Q_{(1_n,\widetilde{E})} \widetilde{W}_j}{||Q_{(1_n,\widetilde{E})} \widetilde{W}_j||^2} = \frac{n^{-1} y^T Q_{(1_n,\widetilde{E})} \widetilde{W}_j}{n^{-1} ||Q_{(1_n,\widetilde{E})} \widetilde{W}_j||^2}, \quad (28)$$

and the corresponding model is $\mu = \mu(\tilde{E}, G_j) = 1_n \beta_{0j} + \tilde{E} \beta_{1j} + G_j \beta_{2j} + (G_j \circ \tilde{E}) \beta_{3j}$ (i.e. model (7)). By an argument analogous to that which leads to (27),

$$\begin{split} \widehat{\beta}_{3j}^{0123}(\widetilde{E},G_j) &= \frac{y^T Q_{(1_n,\widetilde{E},G_j)} \widetilde{W}_j}{||Q_{(1_n,\widetilde{E},G_j)} \widetilde{W}_j||^2} = \frac{n^{-1} y^T Q_{(1_n,\widetilde{E},G_j)} \widetilde{W}_j}{n^{-1} ||Q_{(1_n,\widetilde{E},G_j)} \widetilde{W}_j||^2} \\ &= \frac{n^{-1} y^T Q_{(1_n,\widetilde{E})} \widetilde{W}_j - \frac{(n^{-1} y^T \widetilde{G}_j)(n^{-1} \widetilde{W}_j^T \widetilde{G}_j)}{n^{-1} \widetilde{G}_j^T \widetilde{G}_j}}{n^{-1} ||Q_{(1_n,\widetilde{E})} \widetilde{W}_j||^2 - \frac{(n^{-1} \widetilde{W}_j^T \widetilde{G}_j)^2}{n^{-1} \widetilde{G}_j^T \widetilde{G}_j}} + o_p(1). \end{split}$$

(29)

29

Here we focus on the quantity

$$n^{-1}\widetilde{W}_{j}^{T}\widetilde{G}_{j} = n^{-1}\sum_{i=1}^{n}\widetilde{E}_{i}G_{ij}\widetilde{G}_{ij}$$
$$= \mathbb{E}\left(n^{-1}\sum_{i=1}^{n}\widetilde{E}_{i}G_{ij}\widetilde{G}_{ij}\right) + O_{p}\left\{\operatorname{Var}\left(n^{-1}\sum_{i=1}^{n}\widetilde{E}_{i}G_{ij}\widetilde{G}_{ij}\right)^{1/2}\right\}$$

By the independence between *E* and G_i , 1

$$\mathbf{E}\left(n^{-1}\sum_{i=1}^{n}\widetilde{E}_{i}G_{ij}\widetilde{G}_{ij}\right) = n^{-1}\sum_{i=1}^{n}\mathbf{E}(\widetilde{E}_{i})\mathbf{E}(G_{ij}\widetilde{G}_{ij}) = 0,$$

where the last identity is due to the fact that $E(\tilde{E}_i) = E(E_i - \bar{E}) =$ 2 3

$$n^{-1}\widetilde{W}_i^T\widetilde{G}_i = O_p(n^{-1/2})$$

and by substituting the above into (29),

$$\widehat{\beta}_{3j}^{0123}(\widetilde{E}, G_j) = \frac{n^{-1}y^T Q_{(1_n, \widetilde{E})} \widetilde{W}_j + (n^{-1}y^T \widetilde{G}_j) O_p(n^{-1/2})}{n^{-1} ||Q_{(1_n, \widetilde{E})} \widetilde{W}_j||^2} + o_p(1).$$
(30)

This representation reveals that, if $n^{-1}y^T Q_{(1_n,\widetilde{E})}\widetilde{W}_j$ dominates $(n^{-1}y^T \widetilde{G}_j)n^{-1/2}$, $\widehat{\beta}_{3j}^{0123}(\widetilde{E}, G_j)$ (eq. (30)) is approximated by $\widehat{\beta}_{3i}^{013}(\widetilde{E}, G_j)$ (eq. (28)). In other words, the approximation breaks down only if $n^{-1/2}(n^{-1}y^T \widetilde{G}_i)$ cannot be ignored in comparison to $n^{-1}y^T Q_{(1, \tilde{E})} \widetilde{W}_j$ for large *n*, which is the case when the *j*th variant has a large marginal effect on y while the GxE interaction effect is weak or absent. Such variants should in principle be captured 10 by the marginal association scan. The proposed algorithm thus 11 implements the marginal association scan in addition to the GxE 12 interaction scan, which avoids missing variants that have strong 13 marginal effects. Supplementary Figures S9-S16 confirm that the 14 approximation works well in practice with real data, in which 15 we can see the importance of centering *E* (see "Prediction of real 16 quantitative trait" section). 17

Invariance of regression coefficient estimate for GxE interaction 18 19

Here we show that the least squares estimate of regression coeffi-20 cient β_{3i} in the model $\mu = \mu(E, G_i) = 1_n \beta_{0i} + E \beta_{1i} + G_i \beta_{2i} + (G_i \circ$ 21 E) β_{3j} , model (6), is invariant if E is replaced by $E^a = E - a \mathbf{1}_n$ 22 and/or G_i is replaced by $G_i^b = G_i - b\mathbf{1}_n$ for any scalar values a 23

and
$$b$$
. Recall (18),

$$\widehat{\beta}_{3j}^{0123}(E,G_j) = \frac{y^{T}Q_{(1_n,E,G_j)}W_j}{||Q_{(1_n,E,G_j)}W_j||^2},$$

where $W_i = (E \circ G_i)$. Therefore,

$$\widehat{\beta}_{3j}^{0123}(E^{a},G_{j}^{b}) = \frac{y^{T}Q_{(1_{n},E^{a},G_{j}^{b})}W_{j}^{a,b}}{||Q_{(1_{n},E^{a},G_{j}^{b})}W_{j}^{a,b}||^{2}},$$
(31)

where $W_i^{a,b} = (E^a \circ G_i^b) = (E - a \mathbf{1}_n) \circ (G_j - b \mathbf{1}_n) = W_j - W_j$

²⁷
$$bE - aG_j + ab1_n$$
. Note that $Q_{(1_n, E^a, G_j^b)} = I_n - P_{(1_n, E^a, G_j^b)} = I_2$
²⁸ $P_{(1_n, E - a1_n, G_j - b1_n)} = I_n - P_{(1_n, E, G_j)} = Q_{(1_n, E, G_j)}$, Hence,

$$Q_{(1_n, E^a, G_j^b)}W_j^{a, b} = Q_{(1_n, E, G_j)}(W_j - bE - aG_j + ab1_n) = Q_{(1_n, E, G_j)}W_j,$$

in which the second identity is due to the fact that bE, aG_i , and $ab1_n$ are included in the linear span by $(1_n, E, G_i)$. Therefore, by (31), for any scalar values *a* and *b*, the following identity holds:

$$\widehat{\beta}_{3j}^{0123}(E^a, G_j^b) = \frac{y^T Q_{(1_n, E, G_j)} W_j}{||Q_{(1_n, E, G_j)} W_j||^2} = \widehat{\beta}_{3j}^{0123}(E, G_j).$$
(32)

It is noteworthy that the invariance is essentially due to the involvement of both *E* and G_i , so it is not guaranteed to hold in the absence of either of the two terms.

Genes | Genomes | Genetics