

Supplementary note

A brief review of the liability threshold model

Consider n independent subjects with genotyped SNPs and their relatives whose relation with the subject and disease status and are only known. Phenotypes of individual i and his relative j are denoted by Y_i and Y_{i_j} , respectively and further define a vector $\mathbf{Y}_i = (Y_i, \mathbf{Y}_{i_j})^t$, where $\mathbf{Y}_{i_j} = (Y_{i_1}, \dots, Y_{i_{n_i}})^t$. Similarly, the liability vector, the genotype vector and any environmental effects are denoted as $\mathbf{L}_i = (L_i, \mathbf{L}_{i_j})^t$, $\mathbf{G}_i = (G_i, \mathbf{G}_{i_j})^t$, and $\mathbf{Z}_i = (Z_i, \mathbf{Z}_{i_j})^t$ in order. Polygenic effects explain phenotypic similarity between family members, and correlations among family members are assumed to be constructed by kinship coefficients. We denote $\pi_{jj'}$ as the kinship coefficient between two relatives j and j' of subject i and d_{i_j} as the inbreeding coefficient for relative j of subject i . The inbreeding coefficient d_{i_j} is a parameter quantifying the departure from HWE and ranges from 0 to 1. We denote the kinship coefficient matrix of subject i 's relatives by

$$\Psi_i^{\text{rel}} = \begin{pmatrix} 1 + d_{i_1} & \cdots & 2\pi_{1n_i} \\ \vdots & \ddots & \vdots \\ 2\pi_{n_i 1} & \cdots & 1 + d_{i_{n_i}} \end{pmatrix},$$

and the corresponding kinship coefficient matrix for both subject i and his/her relatives are defined by Ψ_i . We denote a $w \times w$ dimensional identity matrix by I_w , and w dimensional column vector of which all elements are 0 and 1 by $\mathbf{0}_w$ and $\mathbf{1}_w$ respectively. If we let σ_g^2 and σ_ϵ^2 be variances of polygenic effect and random effect, and Z_i is assumed to include the intercept, we can assume that

$$\mathbf{L}_i = \mathbf{Z}_i \alpha + \mathbf{P}_i + \mathbf{E}_i, \mathbf{P}_i \sim MVN(\mathbf{0}_{n_i+1}, \sigma_g^2 \Psi_i), \mathbf{E}_i \sim MVN(\mathbf{0}_{n_i+1}, \sigma_\epsilon^2 \mathbf{I}_{n_i+1}).$$

Here we assume no main genetic effects (thus, \mathbf{G}_i is not included in the model) and $\sigma_\epsilon^2 = 1$. If the heritability, h^2 , and the prevalence of the disease, q are known, then liability threshold T and σ_g^2 can be evaluated from the standard normal distribution, i.e.,

$$\Phi\left(-\frac{T}{\sqrt{\sigma_g^2 + 1}}\right) = 1 - q.$$

Here we used $q = 0.099$. Under this liability threshold model, the conditional mean (CM) of disease risk can be defined as,

$$CM = E\left(L_i \mid \mathbf{Y}_{i_j} = \mathbf{y}_{i_j}\right) \quad (S1)$$

Note that neither direct genotyped variables nor environmental variables (which is not available for the relatives) were included in our liability model. However, it is straightforward to extend the model in the existence of those variables.

Calculation of the conditional mean

Based on this liability threshold model, we can calculate the conditional mean of L_i , or CM, when conditioning on family history. To make the phenotype in liability scale, let $I_{L_{i_j}} = 1$ if $L_{i_j} \in (T, \infty)$ and 0 otherwise. Then $\mathbf{I}_{L_i} = \left(I_{L_{i_1}}, \dots, I_{L_{i_{n_i}}}\right)^t$ and the CM of for subject i becomes $E\left(L_i \mid \mathbf{I}_{L_{i_j}} = \mathbf{1}_{n_i}\right)$, which can be calculated using the truncated multivariate normal distribution. Here n_i dimensional column vector of which all elements are 1 by $\mathbf{1}_{n_i}$.

To calculate the CM, we assume a liability vector, \mathbf{L} (hereafter we dropped the subscripts for simplicity) follows n -dimensional multivariate normal distribution and the probability density function (*pdf*) is,

$$f(\mathbf{L}) = (2\pi)^{-\frac{n}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} \mathbf{L}^t \boldsymbol{\Sigma}^{-1} \mathbf{L}\right) \quad (S2)$$

where $\boldsymbol{\Sigma} = \text{cov}(\mathbf{L})$. Correspondingly, the conditional *pdf* of \mathbf{L} becomes

$$f_{\alpha}(\mathbf{L}) = \begin{cases} \frac{1}{\alpha} f(\mathbf{L}), & \text{for } (T, \infty) \\ 0, & \text{otherwise} \end{cases} \quad (\text{S3})$$

where $\alpha = \text{Pr}(\mathbf{L} \in (T, \infty))$, the fraction after truncation.

The moment generating approach can be applied to the truncated density in Eq (3) and provides the n -dimensional truncated *mgf* (Manjunath and Wilhelm, 2012),

$$m(\mathbf{t}) = \frac{\exp\left(\frac{\mathbf{t}^t \boldsymbol{\Sigma} \mathbf{t}}{2}\right)}{\alpha (2\pi)^{n/2} |\boldsymbol{\Sigma}|^{1/2}} \int_A \exp\left(-\frac{1}{2} \mathbf{L}^t \boldsymbol{\Sigma} \mathbf{L}\right) d\mathbf{L} \quad (\text{S4})$$

If we let $(\boldsymbol{\Sigma})_{jk} = \sigma_{jk}$ and $F_k(x)$ be the marginal *pdf* of L_k , the CM for subject i can be obtained by

$$\mu_i = \left. \frac{\partial m(\mathbf{t})}{\partial t_i} \right|_{\mathbf{t}=0} = \sum_{k=1}^n \sigma_{ik} F_k^* \quad (\text{S5})$$

where

$$F_k^* = \begin{cases} F_k(T) - F_k(\infty) & \text{if } y_k = 1 \\ F_k(-\infty) - F_k(T) & \text{otherwise} \end{cases} \quad (\text{S6})$$

The F_k can be similarly derived as was done in (Manjunath and Wilhelm, 2012). First, we partition \mathbf{L} into two parts, L_i and \mathbf{L}_{i_j} , and then \mathbf{L} can be rewritten as

$$\mathbf{L} = \begin{pmatrix} L_i \\ \mathbf{L}_{i_j} \end{pmatrix} \sim \text{MVN} \left(\begin{pmatrix} 0 \\ \mathbf{0}_n \end{pmatrix}, \begin{pmatrix} 1 & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix} \right) \quad (\text{S7})$$

If we denote the lower and upper truncated bound of \mathbf{L} respectively as \mathbf{a} and \mathbf{b} , *i.e.*,

$\mathbf{a} < \mathbf{L} < \mathbf{b}$ where $\mathbf{a} = (a_i, \mathbf{a}_{i_j})^t$, $\mathbf{b} = (b_i, \mathbf{b}_{i_j})^t$. Then the truncated normal

distribution function is

$$f_{\alpha}(\mathbf{L}_{i_j}, L_i = x) = \alpha^{-1} f(L_i = x) f(\mathbf{L}_{i_j} | L_i = x) I_{\alpha}. \quad (\text{S8})$$

Because a conditional distribution of a normal distribution is also normally distributed, $L_{ij}|L_i = x$ is normally distributed with $E(L_{ij}|L_i = x) = \Sigma_{12}x$ and $Var(L_{ij}|L_i = x) = \Sigma_{22} - \Sigma_{12}\Sigma_{21}^t$. Accordingly, the multivariate marginal pdf of L_{ij} becomes

$$F_{L_{ij}}(x) = \int_{a_{ij}}^{b_{ij}} \alpha^{-1} f(L_i = x) f(L_{ij}|L_i = x) dL_{ij}. \quad (S8)$$

The integral is over L_{ij} and thus $F_{L_{ij}}(x)$ can be rewritten as $\alpha^{-1} f(L_i = x) \int_{a_{ij}}^{b_{ij}} f(L_{ij}|L_i = x) dL_{ij}$, which can be readily computed by using conventional statistical software. For this purpose, we used the `pmvnorm()` function in `mtvnorm` R package.

Note that if we denote the number of relatives for subject i by n_i the convergence rate of the MC algorithm for subject i and for all subjects are $O(\sqrt{n_i})$ and $O(\sqrt{n_1} + \dots + \sqrt{n_n})$, respectively (<http://arxiv.org/pdf/1206.5387.pdf>). Thus, computational intensity is positively related with the number of relatives and subjects and if the same size is large and many relatives' phenotypes are known, it can be computationally intensive.

Supplementary Figures and Tables

Figure S1. MDS plot of two datasets.

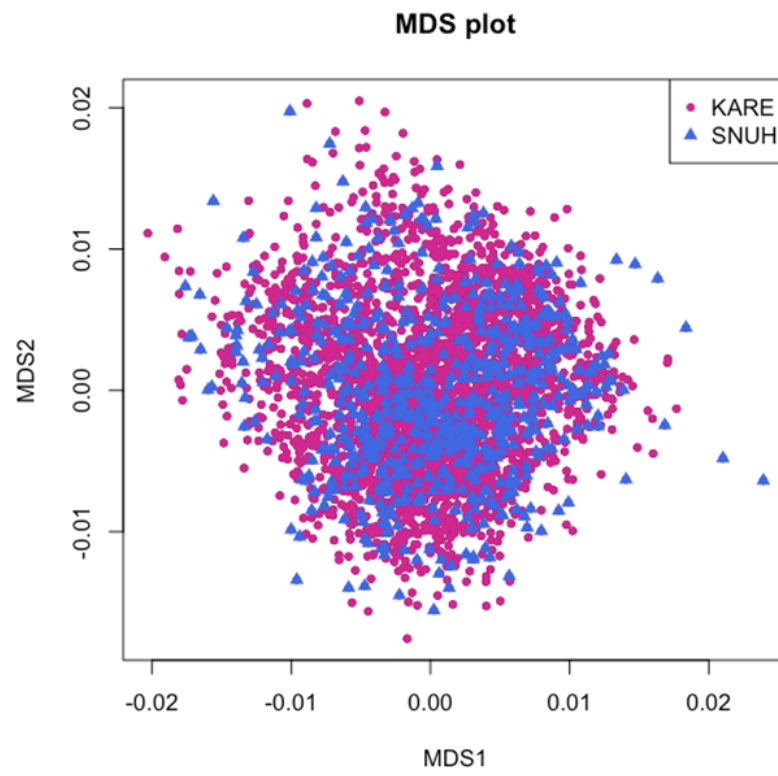


Figure S2. MAF scatter plot. SNPs MAFs (~300k) in SNUH and KARE datasets are plotted.

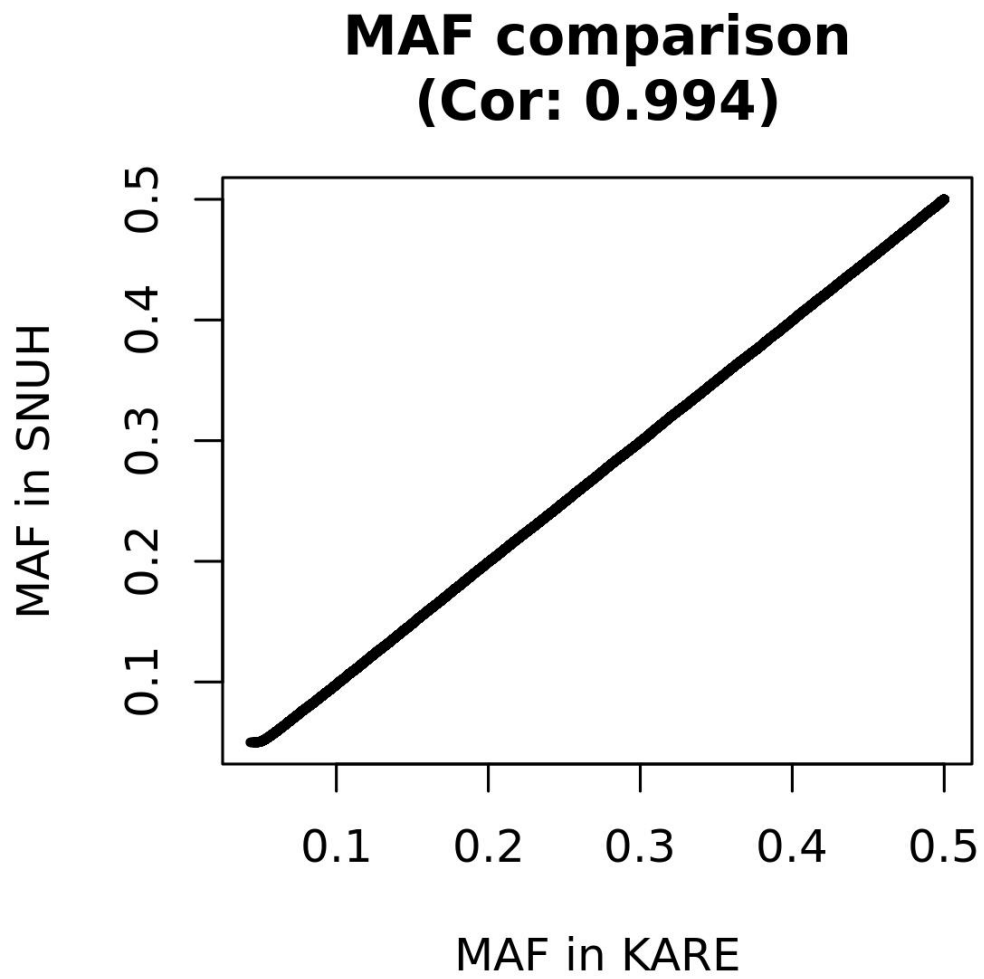


Figure S3. Characteristics of selected SNPs. To find the most effective set of SNPs, we selected SNPs based on the p-value obtained from the logistic regression and the BLUP obtained by the mixed effects model. Since the selected set of SNPs should be applied in penalized regression, we expected that the selection procedure would be more effective if the set of SNPs was uniformly distributed across the genome. Toward this end, we divided the whole genome into 3,233 windows of size 5M and counted the frequency of SNPs in each window. With a varying number of SNPs (0.1k – 20k).

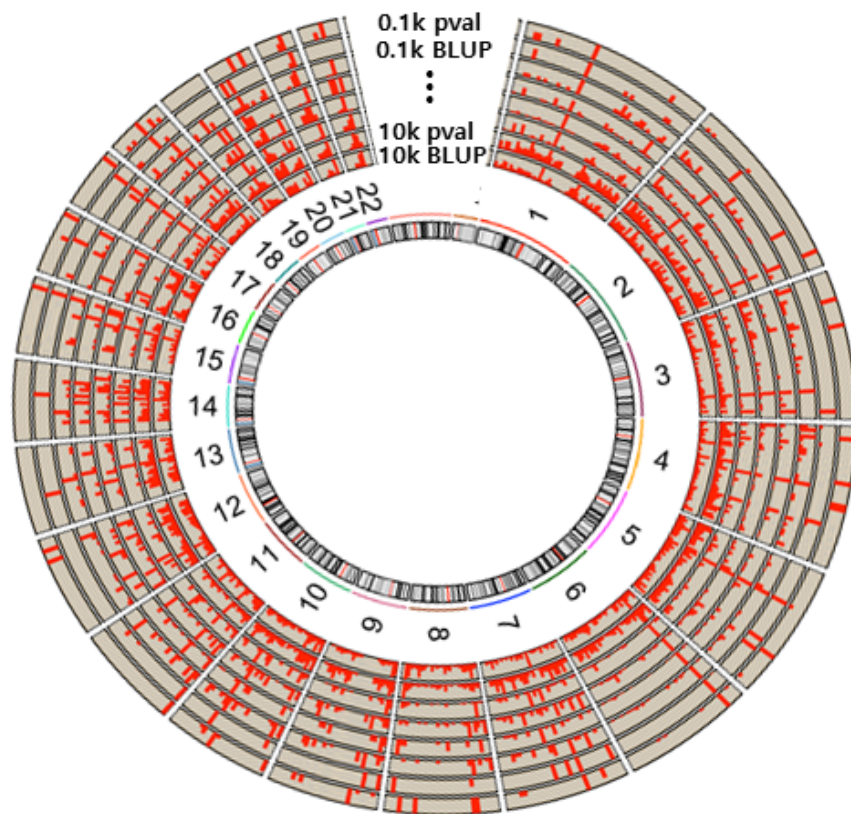


Figure S4. Proportion of variation explained by each variable in the final model.

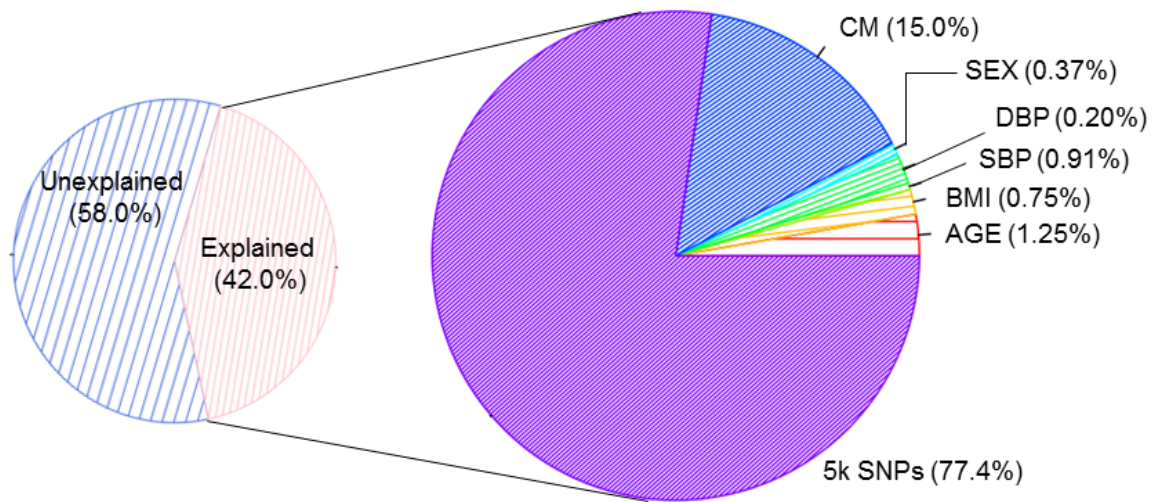


Figure S5. Proportion of variation explained by each variable in the final model without CM variable. For five clinical variables (age, sex, BMI, SBP, DBP) except CM variable, the individual proportions of the variation are shown, whereas the variation explained by the 5,000 SNPs is shown according to their summed proportion.

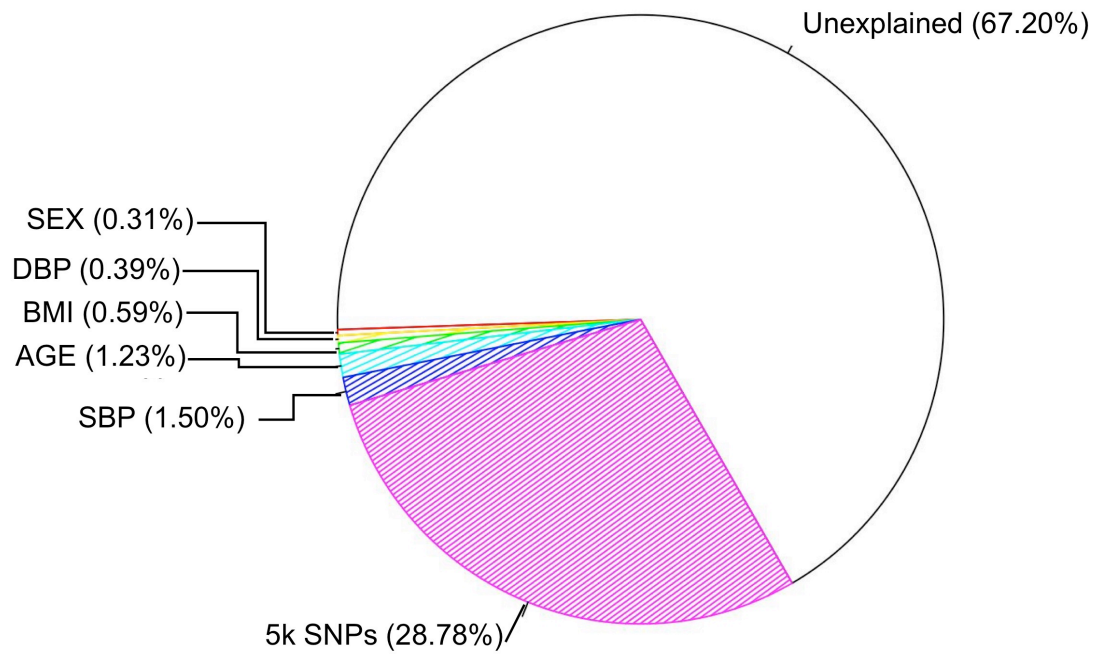


Table S1 Model comparison with different family history measures (AUC)

Family History	# of SNPs	RIDGE	LASSO	EN	SCAD	T.RIDGE
None	100	0.611 (0.023)	0.602 (0.025)	0.602 (0.023)	0.585 (0.023)	0.612 (0.024)
	500	0.614 (0.029)	0.600 (0.031)	0.600 (0.028)	0.594 (0.031)	0.614 (0.030)
	1,000	0.626 (0.029)	0.611 (0.029)	0.611 (0.031)	0.601 (0.032)	0.628 (0.029)
	5,000	0.689 (0.032)	0.647 (0.030)	0.647 (0.030)	-	0.689 (0.031)
	10,000	0.672 (0.028)	0.626 (0.029)	0.626 (0.030)	-	0.672 (0.030)
	20,000	0.674 (0.030)	0.639 (0.031)	0.639 (0.029)	-	0.674 (0.031)
	Weighted Mean	100	0.669 (0.028)	0.605 (0.030)	0.605 (0.030)	-
500		0.617 (0.028)	0.602 (0.026)	0.602 (0.026)	-	0.617 (0.020)
1,000		0.629 (0.028)	0.615 (0.026)	0.615 (0.026)	-	0.630 (0.031)
5,000		0.692 (0.028)	0.650 (0.024)	0.650 (0.024)	-	0.692 (0.012)
10,000		0.676 (0.032)	0.630 (0.035)	0.630 (0.035)	-	0.676 (0.012)
20,000		0.683 (0.033)	0.647 (0.037)	0.647 (0.037)	-	0.683 (0.012)
Conditional Mean		100	0.669 (0.023)	0.659 (0.024)	0.659 (0.024)	0.643 (0.021)
	500	0.659 (0.030)	0.642 (0.031)	0.642 (0.031)	0.639 (0.030)	0.659 (0.031)
	1,000	0.670 (0.029)	0.651 (0.029)	0.651 (0.029)	0.645 (0.028)	0.670 (0.030)
	5,000	0.736 (0.030)	0.691 (0.031)	0.691 (0.031)	-	0.736 (0.027)
	10,000	0.721 (0.029)	0.673 (0.030)	0.673 (0.029)	-	0.721 (0.031)
	20,000	0.725 (0.034)	0.689 (0.031)	0.689 (0.031)	-	0.725 (0.030)

Table S2 Model comparison with different SNP filtering criteria (without CM variable)

CRITERIA	# of SNPs	RIDGE (SD)	LASSO (SD)	EN (SD)	SCAD (SD)	T.RIDGE (SD)
P-value	100	0.642 (0.024)	0.637 (0.022)	0.637 (0.022)	0.616 (0.026)	0.641 (0.021)
	500	0.640 (0.032)	0.626 (0.031)	0.626 (0.031)	0.608 (0.032)	0.640 (0.031)
	1,000	0.640 (0.027)	0.624 (0.028)	0.624 (0.028)	0.608 (0.029)	0.640 (0.026)
	5,000	0.660 (0.029)	0.635 (0.029)	0.635 (0.029)	-	0.660 (0.027)
	10,000	0.668 (0.028)	0.640 (0.030)	0.640 (0.030)	-	0.668 (0.030)
	20,000	0.674 (0.031)	0.640 (0.026)	0.640 (0.026)	-	0.674 (0.028)
BLUP	100	0.611 (0.023)	0.602 (0.025)	0.602 (0.023)	0.585 (0.023)	0.612 (0.024)
	500	0.614 (0.029)	0.600 (0.031)	0.600 (0.028)	0.594 (0.031)	0.614 (0.030)
	1,000	0.626 (0.029)	0.611 (0.029)	0.611 (0.031)	0.601 (0.032)	0.626 (0.029)
	5,000	0.689 (0.032)	0.647 (0.030)	0.647 (0.031)	-	0.689 (0.031)
	10,000	0.672 (0.028)	0.626 (0.029)	0.626 (0.030)	-	0.672 (0.030)
	20,000	0.674 (0.030)	0.639 (0.031)	0.639 (0.029)	-	0.674 (0.031)

Table S3 Model comparison with different SNP filtering criteria (with CM variable)

CRITERIA	# of SNPs	RIDGE (SD)	LASSO (SD)	EN (SD)	SCAD (SD)	T.RIDGE (SD)
P-value	100	0.693 (0.025)	0.687 (0.024)	0.688 (0.023)	0.676 (0.025)	0.693 (0.023)
	500	0.687 (0.031)	0.672 (0.030)	0.672 (0.031)	0.665 (0.032)	0.687 (0.034)
	1,000	0.685 (0.028)	0.669 (0.029)	0.669 (0.031)	0.664 (0.029)	0.685 (0.029)
	5,000	0.709 (0.030)	0.687 (0.031)	0.687 (0.027)	-	0.709 (0.031)
	10,000	0.717 (0.029)	0.690 (0.027)	0.690 (0.029)	-	0.717 (0.028)
	20,000	0.721 (0.030)	0.689 (0.031)	0.689 (0.027)	-	0.721 (0.030)
BLUP	100	0.669 (0.023)	0.659 (0.024)	0.659 (0.024)	0.643 (0.021)	0.669 (0.021)
	500	0.659 (0.030)	0.642 (0.031)	0.642 (0.031)	0.639 (0.030)	0.659 (0.031)
	1,000	0.670 (0.029)	0.651 (0.029)	0.651 (0.029)	0.645 (0.028)	0.670 (0.030)
	5,000	0.736 (0.030)	0.691 (0.031)	0.691 (0.031)	-	0.736 (0.027)
	10,000	0.721 (0.029)	0.673 (0.030)	0.673 (0.029)	-	0.721 (0.031)
	20,000	0.725 (0.034)	0.689 (0.031)	0.689 (0.032)	-	0.725 (0.030)