

Notations

In this section, we define the settings of the following methods, and introduce the frequently used identities.

- \mathbf{I}_J denotes a $J \times J$ identity matrix.
- $\mathbf{1}_{I \times J}$ denotes a $I \times J$ matrix whose elements are all 1.
- $\text{MVN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes a multivariate normal distribution with a mean $\boldsymbol{\mu}$ and a variance-covariance matrix $\boldsymbol{\Sigma}$.
- For any $a \times b$ matrix \mathbf{A} , the $b \times a$ matrix \mathbf{A}^\dagger denotes the Moore-Penrose pseudo-inverse of \mathbf{A} .
- For any $a \times b$ matrix \mathbf{A} and $a \times b$ matrix \mathbf{B} , the $a \times b$ matrix $\mathbf{C} = \mathbf{A} \circ \mathbf{B}$ denotes the Hadamard product between \mathbf{A} and \mathbf{B} .
- For any $a \times a$ square matrix \mathbf{A} , the expression $|\mathbf{A}|_+$ denotes the pseudo-determinant of \mathbf{A} . If \mathbf{A} is positive semi-definite, $|\mathbf{A}|_+$ will be computed as the product of non-zero eigen values of \mathbf{A} .
- \mathbf{y} denotes a $n \times 1$ vector of phenotypic values.
- \mathbf{X} denotes a $n \times p$ covariate matrix with full column rank p where $n \geq p$.
- $\boldsymbol{\beta}$ denotes a $p \times 1$ vector of fixed effects for covariates.
- \mathbf{Z}_c denotes a $n \times m_c$ design matrix corresponding to random effects \mathbf{u}_c for family relatedness.
- \mathbf{u}_c denotes a $m_c \times 1$ vector of random effects for family relatedness. We assume $\mathbf{u}_c \sim \text{MVN}(\mathbf{0}, \mathbf{K}_c \sigma_c^2)$ where \mathbf{K}_c is the additive genetic relationship matrix estimated from marker genotype $\tilde{\mathbf{W}}_c$ and σ_c^2 is the additive genetic variance.
- \mathbf{Z}_{r_i} denotes a $n \times m_{r_i}$ design matrix corresponding to random effects \mathbf{u}_{r_i} for each SNP-set.
- \mathbf{u}_{r_i} denotes a $m_{r_i} \times 1$ vector of random effects for each SNP-set. We assume $\mathbf{u}_{r_i} \sim \text{MVN}(\mathbf{0}, \mathbf{K}_{r_i} \sigma_{r_i}^2)$ where \mathbf{K}_{r_i} is the Gram matrix computed from marker genotype in SNP-set of interest $\tilde{\mathbf{W}}_{r_i}$ and $\sigma_{r_i}^2$ is the genetic variance of each SNP-set.

- $\boldsymbol{\epsilon}$ denotes a $n \times 1$ vector of residuals. We assume $\boldsymbol{\epsilon} \sim \text{MVN}(\mathbf{0}, \mathbf{I}_n \sigma_e^2)$ where σ_e^2 is the residual variance.
- \mathbf{V} is a $n \times n$ phenotypic variance-covariance matrix. In this paper, $\mathbf{V} = \mathbf{Z}_c \mathbf{K}_c \mathbf{Z}_c^T \sigma_c^2 + \mathbf{Z}_{r_i} \mathbf{K}_{r_i} \mathbf{Z}_{r_i}^T \sigma_{r_i}^2 + \mathbf{I}_n \sigma_e^2$. Under the null hypothesis $\sigma_{r_i}^2 = 0$ since $\sigma_{r_i}^2$ is the parameter to be tested. \mathbf{V} is assumed to be a full rank (n) matrix.
- \mathbf{H} is a $n \times n$ matrix which satisfies $\mathbf{H} = \mathbf{V} / \sigma_c^2 = \mathbf{Z}_c \mathbf{K}_c \mathbf{Z}_c^T + \mathbf{Z}_{r_i} \mathbf{K}_{r_i} \mathbf{Z}_{r_i}^T \gamma_{r_i} + \mathbf{I}_n \delta_e$ where $\gamma_{r_i} = \sigma_{r_i}^2 / \sigma_c^2$ and $\delta_e = \sigma_e^2 / \sigma_c^2$ to be estimated.
- $\mathbf{S} = \mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is a $n \times n$ symmetric covariance orthogonal projection matrix with rank $n - p$.
- $\mathbf{P} = \mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} = \mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{H}^{-1}$ is a $n \times n$ matrix with rank $n - p$.
- $\mathbf{w} = \mathbf{P} \mathbf{y}$ is a $n \times 1$ vector used for the REML estimation. We also define \mathbf{w}_L , \mathbf{w}_R , \mathbf{P}_L and \mathbf{P}_R by

$$\mathbf{w} = \begin{bmatrix} \mathbf{w}_L \\ \mathbf{w}_R \end{bmatrix} = \begin{bmatrix} \mathbf{P}_L \\ \mathbf{P}_R \end{bmatrix} \mathbf{y} = \mathbf{P} \mathbf{y},$$

where \mathbf{w}_L (or \mathbf{w}_R) stands for a linearly independent (or redundant) part of \mathbf{w} , and \mathbf{P}_L (or \mathbf{P}_R) denotes the corresponding partition of \mathbf{P} . Here, \mathbf{P}_L should satisfy $\mathbf{P}_L^T \mathbf{P}_L = \mathbf{S}$ and $\mathbf{P}_L \mathbf{P}_L^T = \mathbf{I}_{n-p}$ [1, 2, 3]. We also define the eigen decomposition of $\mathbf{S} \mathbf{H} \mathbf{S}$ as

$$\mathbf{S} \mathbf{H} \mathbf{S} = \mathbf{U}_L \boldsymbol{\Lambda}_L \mathbf{U}_L^T,$$

where $\boldsymbol{\Lambda}_L$ is a $(n-p) \times (n-p)$ diagonal matrix whose elements are non-zero eigen values of $\mathbf{S} \mathbf{H} \mathbf{S}$ in the decreasing order, and \mathbf{U}_L is a $n \times (n-p)$ eigen vector matrix whose each eigen vector corresponds to each eigen value. Here, \mathbf{U}_L is the part of the unitary matrix with first $n-p$ columns of eigen vectors, and it satisfies $\mathbf{U}_L \mathbf{U}_L^T = \mathbf{S}$ and $\mathbf{U}_L^T \mathbf{U}_L = \mathbf{I}_{n-p}$ (Proposition 7). Therefore, we can say $\mathbf{P}_L = \mathbf{U}_L^T$.

- $\mathbf{Q} = \mathbf{H}^{-1} - \mathbf{H}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{H}^{-1} \mathbf{P} = \mathbf{P}^T \mathbf{H}^{-1}$ is a $n \times n$ symmetry matrix with rank $n - p$.
- \mathbf{W}_{r_i} is a $m_{r_i} \times M_i$ marker genotype matrix belonging to the i -th SNP-set. Here, M_i is the number of SNPs in the i -th SNP-set.

Multi-kernel linear mixed model

In this study, the alternative model can be written as the multi-kernel mixed model.

$$\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \mathbf{Z}_c \mathbf{u}_c + \mathbf{Z}_{r_i} \mathbf{u}_{r_i} + \boldsymbol{\epsilon} \quad (1)$$

On the other hand, the null model can be written as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}_c\mathbf{u}_c + \boldsymbol{\epsilon} \quad (2)$$

Therefore, we somehow test the null hypothesis $H_0 : \sigma_{r_i}^2 = 0$ for evaluating the significance of the effects of SNP-set of interest.

Restricted maximum likelihood (REML)

For the multi-kernel linear mixed model described above, the restricted log likelihood of \mathbf{y} can be regarded as the log likelihood of \mathbf{w}_L , and it can be expressed as a quite simple format by only using \mathbf{Q} .

$$\begin{aligned} l_R(\mathbf{y}; \sigma_c, \gamma_{r_i}, \delta_e) &= -\frac{n-p}{2} \log(2\pi) - \frac{1}{2} \log |\mathbf{P}_L \mathbf{V} \mathbf{P}_L^T| - \frac{1}{2} \mathbf{y} \mathbf{P}_L^T (\mathbf{P}_L \mathbf{V} \mathbf{P}_L^T)^{-1} \mathbf{P}_L \mathbf{y} \\ &= -\frac{n-p}{2} \log(2\pi) + \frac{1}{2} \log |\mathbf{P}^T \mathbf{V}^{-1} \mathbf{P}|_+ - \frac{1}{2} \mathbf{y}^T \mathbf{V}^{-1} \mathbf{P} \mathbf{y} \\ &= -\frac{n-p}{2} \log(2\pi\sigma_c^2) + \frac{1}{2} \log |\mathbf{Q}|_+ - \frac{1}{2\sigma_c^2} \mathbf{y}^T \mathbf{Q} \mathbf{y} \end{aligned} \quad (3)$$

by using Proposition 8 and 9.

In addition, by using Proposition 10, this restricted log likelihood is same as the well known format [4, 5, 6] as follows.

$$\begin{aligned} l_R(\mathbf{y}; \sigma_c, \gamma_{r_i}, \delta_e) &= -\frac{n-p}{2} \log(2\pi\sigma_c^2) + \frac{1}{2} \log |\mathbf{Q}|_+ - \frac{1}{2\sigma_c^2} \mathbf{y}^T \mathbf{Q} \mathbf{y} \\ &= -\frac{n-p}{2} \log(2\pi\sigma_c^2) - \frac{1}{2} \log |\mathbf{H}| + \frac{1}{2} \log |\mathbf{X}^T \mathbf{X}| \\ &\quad - \frac{1}{2} \log |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}| - \frac{1}{2\sigma_c^2} \mathbf{y}^T \mathbf{Q} \mathbf{y} \end{aligned} \quad (4)$$

Plugging $\hat{\sigma}_c^2 = \mathbf{y}^T \mathbf{Q} \mathbf{y} / (n-p)$ into Eq 3, we get

$$l_R(\mathbf{y}; \hat{\sigma}_c, \gamma_{r_i}, \delta_e) = -\frac{n-p}{2} \left\{ \log \left(\frac{2\pi e}{n-p} \right) + \log (\mathbf{y}^T \mathbf{Q} \mathbf{y}) \right\} + \frac{1}{2} \log |\mathbf{Q}|_+ \quad (5)$$

Efficient likelihood ratio test used in RAINBOW

In this section, we describe how to implement computationally efficient algorithm for the likelihood ratio (LR) test [7, 8] for the two kernel linear mixed model.

In this study, we assume the Gram matrix for each SNP-set $\mathbf{Z}_{r_i} \mathbf{K}_{r_i} \mathbf{Z}_{r_i}^T$ is low rank, so $\text{rank}(\mathbf{Z}_{r_i} \mathbf{K}_{r_i} \mathbf{Z}_{r_i}^T) \ll n$. Then, \mathbf{H} can be written as

$$\begin{aligned} \mathbf{H} &= \mathbf{Z}_c \mathbf{K}_c \mathbf{Z}_c^T + \mathbf{Z}_{r_i} \mathbf{K}_{r_i} \mathbf{Z}_{r_i}^T \gamma_{r_i} + \mathbf{I}_n \delta_e \\ &= \mathbf{Z}_c \mathbf{K}_c \mathbf{Z}_c^T + \tilde{\mathbf{W}}_{r_i} \tilde{\boldsymbol{\Gamma}}_{r_i} \tilde{\mathbf{W}}_{r_i}^T + \mathbf{I}_n \delta_e, \end{aligned} \quad (6)$$

where $\tilde{\mathbf{W}}_{r_i}$ is a $n \times k$ matrix and $\tilde{\boldsymbol{\Gamma}}_{r_i}$ is a $k \times k$ square matrix. Here, k is the rank of $\mathbf{Z}_{r_i} \mathbf{K}_{r_i} \mathbf{Z}_{r_i}^T$, so $k \ll n$. Concrete examples of $\tilde{\mathbf{W}}_{r_i}$ and $\tilde{\boldsymbol{\Gamma}}_{r_i}$ will be described later.

Low rank update of \mathbf{Q}

One of the drawbacks of the LR test is a large amount of computation because the LR test requires the maximization of restricted likelihood for each SNP-set. To reduce the computational complexity, [9, 10] proposed the efficient computation of the restricted log likelihood Eq 5 by using low rank update of \mathbf{Q} as follows.

First, the low-rank update of \mathbf{Q} is

$$\begin{aligned}
\mathbf{Q} &= \mathbf{P}^T \mathbf{H}^{-1} \mathbf{P} \\
&= (\mathbf{S} \mathbf{H} \mathbf{S})^\dagger \\
&= \left(\mathbf{S} \left(\mathbf{Z}_c \mathbf{K}_c \mathbf{Z}_c^T + \tilde{\mathbf{W}}_{r_i} \tilde{\mathbf{\Gamma}}_{r_i} \tilde{\mathbf{W}}_{r_i}^T + \mathbf{I}_n \delta_e \right) \mathbf{S} \right)^\dagger \\
&= \left(\mathbf{S} \left(\mathbf{Z}_c \mathbf{K}_c \mathbf{Z}_c^T + \mathbf{I}_n \delta_e \right) \mathbf{S} + \mathbf{S} \tilde{\mathbf{W}}_{r_i} \tilde{\mathbf{\Gamma}}_{r_i} \tilde{\mathbf{W}}_{r_i}^T \mathbf{S} \right)^\dagger \\
&= \left(\mathbf{U}_L \left(\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e \right) \mathbf{U}_L^T + \mathbf{U}_L \mathbf{U}_L^T \tilde{\mathbf{W}}_{r_i} \tilde{\mathbf{\Gamma}}_{r_i} \tilde{\mathbf{W}}_{r_i}^T \mathbf{U}_L \mathbf{U}_L^T \right)^\dagger \\
&= \left(\mathbf{U}_L \left(\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e + \mathbf{U}_L^T \tilde{\mathbf{W}}_{r_i} \tilde{\mathbf{\Gamma}}_{r_i} \tilde{\mathbf{W}}_{r_i}^T \mathbf{U}_L \right) \mathbf{U}_L^T \right)^\dagger \\
&= \mathbf{U}_L \left(\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e + \mathbf{U}_L^T \tilde{\mathbf{W}}_{r_i} \tilde{\mathbf{\Gamma}}_{r_i} \tilde{\mathbf{W}}_{r_i}^T \mathbf{U}_L \right)^{-1} \mathbf{U}_L^T \\
&= \mathbf{U}_L \left(\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e \right)^{-1} \mathbf{U}_L^T \\
&\quad - \mathbf{U}_L \left(\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e \right)^{-1} \mathbf{U}_L^T \tilde{\mathbf{W}}_{r_i} \\
&\quad \cdot \left(\tilde{\mathbf{\Gamma}}_{r_i}^{-1} + \tilde{\mathbf{W}}_{r_i}^T \mathbf{U}_L \left(\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e \right)^{-1} \mathbf{U}_L^T \tilde{\mathbf{W}}_{r_i} \right)^{-1} \\
&\quad \cdot \tilde{\mathbf{W}}_{r_i}^T \mathbf{U}_L \left(\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e \right)^{-1} \mathbf{U}_L^T \\
&= \mathbf{O}_c - \mathbf{O}_c \tilde{\mathbf{W}}_{r_i} \left(\tilde{\mathbf{\Gamma}}_{r_i}^{-1} + \tilde{\mathbf{W}}_{r_i}^T \mathbf{O}_c \tilde{\mathbf{W}}_{r_i} \right)^{-1} \tilde{\mathbf{W}}_{r_i}^T \mathbf{O}_c \tag{7}
\end{aligned}$$

Here, we use Proposition 6, 5, 2, 7 and the Woodbury identity, and we define the eigen decomposition of $\mathbf{S} \mathbf{Z}_c \mathbf{K}_c \mathbf{Z}_c^T \mathbf{S} = \mathbf{U}_L \mathbf{\Lambda}_c \mathbf{U}_L^T$ where $\mathbf{\Lambda}_c$ is a $n - p \times n - p$ diagonal matrix and $n \times n$ matrix $\mathbf{O}_c = \mathbf{U}_L \left(\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e \right)^{-1} \mathbf{U}_L^T$ to shorten the notation. Assuming that the eigen decomposition of $\mathbf{S} \left(\mathbf{Z}_c \mathbf{K}_c \mathbf{Z}_c^T + \mathbf{I}_{n-p} \delta_e \right) \mathbf{S}$ has been pre-computed, the additional computation required will be $\mathbf{O}_c \tilde{\mathbf{W}}_{r_i}$, an $O(n^2 k)$ operation.

Update of the squared form $\mathbf{y}^T \mathbf{Q} \mathbf{y}$

Since we know the derivation of the low-rank update of \mathbf{Q} (Eq 7) now, we can plug this into the squared form.

$$\mathbf{y}^T \mathbf{Q} \mathbf{y} = \mathbf{y}^T \mathbf{O}_c \mathbf{y} - \mathbf{y}^T \mathbf{O}_c \tilde{\mathbf{W}}_{r_i} \left(\tilde{\mathbf{\Gamma}}_{r_i}^{-1} + \tilde{\mathbf{W}}_{r_i}^T \mathbf{O}_c \tilde{\mathbf{W}}_{r_i} \right)^{-1} \tilde{\mathbf{W}}_{r_i}^T \mathbf{O}_c \mathbf{y} \tag{8}$$

Update of the determinant $\log|\mathbf{Q}|_+$

Since the update matrix is not necessarily positive semi-definite, we have to slightly modify Eq 7 to avoid numerical instability.

$$\begin{aligned}
\log|\mathbf{Q}|_+ &= -\log|\mathbf{U}_L^T \mathbf{H} \mathbf{U}_L| \\
&= -\log\left|\mathbf{I}_{n-p}\delta_e + \mathbf{\Lambda}_c + \mathbf{U}_L^T \tilde{\mathbf{W}}_{r_i} \tilde{\mathbf{\Gamma}}_{r_i} \tilde{\mathbf{W}}_{r_i}^T \mathbf{U}_L\right| \\
&= -\log|\mathbf{I}_{n-p}\delta_e + \mathbf{\Lambda}_c + \mathbf{A}\mathbf{B}| \\
&= -\log\left(\left|\mathbf{I}_{n-p}\delta_e + \mathbf{\Lambda}_c\right| \cdot \left|\mathbf{I}_{n-p} + (\mathbf{I}_{n-p}\delta_e + \mathbf{\Lambda}_c)^{-1} \mathbf{A}\mathbf{B}\right|\right) \\
&= -\log|\mathbf{I}_{n-p}\delta_e + \mathbf{\Lambda}_c| - \log\left|\mathbf{I}_k + \mathbf{B}(\mathbf{I}_{n-p}\delta_e + \mathbf{\Lambda}_c)^{-1} \mathbf{A}\right| \quad (9)
\end{aligned}$$

Here, we used Proposition 5, 6, 7 and the Sylvester's determinant identity, and we also define a $n-p \times k$ matrix $\mathbf{A} = \mathbf{U}_L^T \tilde{\mathbf{W}}_{r_i} \tilde{\mathbf{\Gamma}}_{r_i}$ and a $k \times n-p$ matrix $\mathbf{B} = \tilde{\mathbf{W}}_{r_i}^T \mathbf{U}_L$. When calculating the second term of Eq 9, we have to be careful about the possibility that $\mathbf{I}_k + \mathbf{B}(\mathbf{I}_{n-p}\delta_e + \mathbf{\Lambda}_c)^{-1} \mathbf{A}$ is not necessarily positive semi-definite. Since the matrix \mathbf{Q} is positive semi-definite, it is possible to calculate the second term of Eq 9, however, the matrix $\mathbf{I}_k + \mathbf{B}(\mathbf{I}_{n-p}\delta_e + \mathbf{\Lambda}_c)^{-1} \mathbf{A}$ has an even numbers of negative eigen values, so instead of simply computing the determinant of this matrix, we have to avoid taking logarithms of such negative eigen values.

Estimation of variance components

Now we can derive the efficient computation of the restricted log likelihood from Eq 5, Eq 8 and Eq 9, we can optimize γ_{r_i} and δ_e over maximization of Eq 5 by using L-BFGS optimization method as we introduce in the paper [11].

Then, we calculated the weighted variance-covariance matrix by γ_{r_i} , and reestimate variance components by using EMMA (efficient mixed model association) or GEMMA (genome wide efficient mixed model association) [6, 12].

Discussion on the identity of $\tilde{\mathbf{W}}_{r_i}$ and $\tilde{\mathbf{\Gamma}}_{r_i}$

In this subsection, we discuss on what $\tilde{\mathbf{W}}_{r_i}$ and $\tilde{\mathbf{\Gamma}}_{r_i}$ correspond to depending on the kind of the kernels for \mathbf{K}_{r_i} .

For the case where \mathbf{K}_{r_i} is the linear kernel

Here, we discuss the case where \mathbf{K}_{r_i} is calculated as the linear kernel of the marker genotype \mathbf{W}_{r_i} belonging to the i -th SNP-set.

For example, if we assume \mathbf{K}_{r_i} is the additive genetic matrix of \mathbf{W}_{r_i} , $\tilde{\mathbf{W}}_{r_i}$ corresponds to

$$\tilde{\mathbf{W}}_{r_i} = \frac{\mathbf{Z}_{r_i} (\mathbf{W}_{r_i} + \mathbf{1}_{m_{r_i} \times M_i} - 2 \cdot \mathbf{\Phi})}{2 \cdot \sum_{m=1}^{M_i} p_m (1 - p_m)}, \quad (10)$$

where p_m is the allele frequency of the 1 allele at marker m and $\mathbf{\Phi}$ is a $m_{r_i} \times M_i$ matrix whose m th column equals to $\mathbf{\Phi} = p_m \cdot \mathbf{1}_{m_{r_i} \times 1}$ [13].

On the other hand, $\tilde{\mathbf{\Gamma}}_{r_i}$ corresponds to

$$\tilde{\mathbf{\Gamma}}_{r_i} = \mathbf{I}_{M_i} \gamma_{r_i}, \quad (11)$$

so in this case, $\tilde{\mathbf{\Gamma}}_{r_i}$ is a $M_i \times M_i$ diagonal matrix. Here, k in Eq 7 equals to the number of SNPs in the i -th SNP-set, M_i .

For the case where \mathbf{K}_{r_i} is the exponential or gaussian kernel

Here, we discuss the case where \mathbf{K}_{r_i} is calculated as the exponential or gaussian kernel of the marker genotype \mathbf{W}_{r_i} belonging to the i -th SNP-set.

First, we calculated the Euclidean distance matrix \mathbf{D}_{r_i} from the marker genotype \mathbf{W}_{r_i} belonging to the i -th SNP-set. Then for the case where \mathbf{K}_{r_i} is the exponential kernel, \mathbf{K}_{r_i} is calculated as

$$\mathbf{K}_{r_i} = \exp\left(-\frac{h_{r_i}\mathbf{D}_{r_i}}{\sqrt{M_i}}\right), \quad (12)$$

where h_{r_i} is a hyperparameter calculated as the inverse of median of the off-diagonal elements of $\mathbf{D}_{r_i}^2/M_i$ for the default setting of RAINBOW. To scale the distance matrix, the distance matrix is divided by $\sqrt{M_i}$ in the exponential. Similarly, for the case where \mathbf{K}_{r_i} is the gaussian kernel, \mathbf{K}_{r_i} is calculated as

$$\mathbf{K}_{r_i} = \exp\left(-\frac{h_{r_i}\mathbf{D}_{r_i}^2}{M_i}\right), \quad (13)$$

where the term in the exponential is also deivide by M_i to fit its scale to the linear kernel.

In this case, we cannot apply the decomposition of \mathbf{K}_{r_i} as seen in the linear kernel case, however, it is assumed that the rank of \mathbf{K}_{r_i} is still much smaller than the number of observations, n . This is because, since there should be strong linkage disequilibrium between SNPs in each SNP-set and many accessions share the same SNPs in that SNP-set, the number of genotypes in that SNP-set m_{r_i} may be much smaller than n . Therefore, in such case,

$$\tilde{\mathbf{W}}_{r_i} = \mathbf{Z}_{r_i}, \quad (14)$$

$$\tilde{\mathbf{\Gamma}}_{r_i} = \mathbf{K}_{r_i}\gamma_{r_i} \quad (15)$$

Here, k in Eq 7 equals to the number of SNPs in the genotypes in the i th SNP-set, m_{r_i} . We should be careful about that $\tilde{\mathbf{\Gamma}}_{r_i}$ is not a diagonal matrix in this case.

Efficient likelihood ratio test for dominance and epistatic effects

In this section, we extend the efficient LR test for the one random effect of each SNP-set to that for multiple random effects of the SNP-set. In this case, the multi-kernel mixed model will be

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}_c\mathbf{u}_c + \sum_{l=1}^L \mathbf{Z}_{r_i,l}\mathbf{u}_{r_i,l} + \boldsymbol{\epsilon}, \quad (16)$$

where $\mathbf{u}_{r_i,l}$ is the l th random effect of the i th SNP-set and $\mathbf{Z}_{r_i,l}$ is a $n \times m_{r_i,l}$ design matrix which correspond to $\mathbf{u}_{r_i,l}$. Examples of multiple random effects are additive effects, dominance effects and epistatic effects, and the

test for the significance of these effects will be described later.
The restricted log likelihood is still Eq 5, however, here

$$\begin{aligned}\mathbf{H} &= \mathbf{Z}_c \mathbf{K}_c \mathbf{Z}_c^T + \sum_{l=1}^L \mathbf{z}_{r_i,l} \mathbf{K}_{r_i,l} \mathbf{z}_{r_i,l}^T \gamma_{r_i,l} + \mathbf{I}_n \delta_e \\ &= \mathbf{Z}_c \mathbf{K}_c \mathbf{Z}_c^T + \sum_{l=1}^L \tilde{\mathbf{W}}_{r_i,l} \tilde{\mathbf{\Gamma}}_{r_i,l} \tilde{\mathbf{W}}_{r_i,l}^T + \mathbf{I}_n \delta_e,\end{aligned}\quad (17)$$

where $\mathbf{K}_{r_i,l}$, $\gamma_{r_i,l}$, $\tilde{\mathbf{W}}_{r_i,l}$ and $\tilde{\mathbf{\Gamma}}_{r_i,l}$ are the extensions of \mathbf{K}_{r_i} , γ_{r_i} , $\tilde{\mathbf{W}}_{r_i}$ and $\tilde{\mathbf{\Gamma}}_{r_i}$ for the l th random effect.

Low rank update of $\mathbf{Q}^{(L)}$

In this case, the low rank update of \mathbf{Q} is realized by extending Eq 7 to

$$\begin{aligned}\mathbf{Q}^{(L)} &= \mathbf{P}^T \mathbf{H}^{-1} \mathbf{P} \\ &= (\mathbf{S} \mathbf{H} \mathbf{S})^\dagger \\ &= \left(\mathbf{S} \left(\mathbf{Z}_c \mathbf{K}_c \mathbf{Z}_c^T + \sum_{l=1}^L \tilde{\mathbf{W}}_{r_i,l} \tilde{\mathbf{\Gamma}}_{r_i,l} \tilde{\mathbf{W}}_{r_i,l}^T + \mathbf{I}_n \delta_e \right) \mathbf{S} \right)^\dagger \\ &= \left(\mathbf{S} \left(\mathbf{Z}_c \mathbf{K}_c \mathbf{Z}_c^T + \mathbf{I}_n \delta_e \right) \mathbf{S} + \sum_{l=1}^L \mathbf{S} \tilde{\mathbf{W}}_{r_i,l} \tilde{\mathbf{\Gamma}}_{r_i,l} \tilde{\mathbf{W}}_{r_i,l}^T \mathbf{S} \right)^\dagger \\ &= \left(\mathbf{U}_L \left(\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e \right) \mathbf{U}_L^T + \sum_{l=1}^L \mathbf{U}_L \mathbf{U}_L^T \tilde{\mathbf{W}}_{r_i,l} \tilde{\mathbf{\Gamma}}_{r_i,l} \tilde{\mathbf{W}}_{r_i,l}^T \mathbf{U}_L \mathbf{U}_L^T \right)^\dagger \\ &= \left(\mathbf{U}_L \left(\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e + \sum_{l=1}^L \mathbf{U}_L^T \tilde{\mathbf{W}}_{r_i,l} \tilde{\mathbf{\Gamma}}_{r_i,l} \tilde{\mathbf{W}}_{r_i,l}^T \mathbf{U}_L \right) \mathbf{U}_L^T \right)^\dagger \\ &= \mathbf{U}_L \left(\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e + \sum_{l=1}^L \mathbf{U}_L^T \tilde{\mathbf{W}}_{r_i,l} \tilde{\mathbf{\Gamma}}_{r_i,l} \tilde{\mathbf{W}}_{r_i,l}^T \mathbf{U}_L \right)^{-1} \mathbf{U}_L^T \\ &= \mathbf{U}_L \left(\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e + \sum_{l=1}^{L-1} \tilde{\mathbf{W}}_{r_i,l} \tilde{\mathbf{\Gamma}}_{r_i,l} \tilde{\mathbf{W}}_{r_i,l}^T \right)^{-1} \mathbf{U}_L^T \\ &\quad - \mathbf{U}_L \left(\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e + \sum_{l=1}^{L-1} \tilde{\mathbf{W}}_{r_i,l} \tilde{\mathbf{\Gamma}}_{r_i,l} \tilde{\mathbf{W}}_{r_i,l}^T \right)^{-1} \mathbf{U}_L^T \tilde{\mathbf{W}}_{r_i,L} \\ &\quad \cdot \left(\tilde{\mathbf{\Gamma}}_{r_i,L}^{-1} + \tilde{\mathbf{W}}_{r_i,L}^T \mathbf{U}_L \left(\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e + \sum_{l=1}^{L-1} \tilde{\mathbf{W}}_{r_i,l} \tilde{\mathbf{\Gamma}}_{r_i,l} \tilde{\mathbf{W}}_{r_i,l}^T \right)^{-1} \mathbf{U}_L^T \tilde{\mathbf{W}}_{r_i,L} \right)^{-1} \\ &\quad \cdot \tilde{\mathbf{W}}_{r_i,L}^T \mathbf{U}_L \left(\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e + \sum_{l=1}^{L-1} \tilde{\mathbf{W}}_{r_i,l} \tilde{\mathbf{\Gamma}}_{r_i,l} \tilde{\mathbf{W}}_{r_i,l}^T \right)^{-1} \mathbf{U}_L^T \\ &= \mathbf{Q}^{(L-1)} - \mathbf{Q}^{(L-1)} \tilde{\mathbf{W}}_{r_i,L} \left(\tilde{\mathbf{\Gamma}}_{r_i,L}^{-1} + \tilde{\mathbf{W}}_{r_i,L}^T \mathbf{Q}^{(L-1)} \tilde{\mathbf{W}}_{r_i,L} \right)^{-1} \tilde{\mathbf{W}}_{r_i,L}^T \mathbf{Q}^{(L-1)}\end{aligned}\quad (18)$$

Here, we use Proposition 6, 5, 2, 7 and the Woodbury identity, and we define $\mathbf{Q}^{(L)}$ as the \mathbf{Q} matrix for the L random effects of the SNP-set, so

$\mathbf{Q}^{(L)} = \mathbf{U}_L \left(\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e + \sum_{l=1}^L \mathbf{U}_L^T \tilde{\mathbf{W}}_{r_i,l} \tilde{\mathbf{\Gamma}}_{r_i,l} \tilde{\mathbf{W}}_{r_i,l}^T \mathbf{U}_L \right)^{-1} \mathbf{U}_L$. To compute $\mathbf{Q}^{(L)}$, we set $\mathbf{Q}^{(0)} = \mathbf{U}_L (\mathbf{\Lambda}_c + \mathbf{I}_{n-p} \delta_e)$ first, and repeat Eq 18 for L steps. The computation of $\mathbf{Q}^{(L)}$ will be an $O(n^2 \sum_{l=1}^L k_l)$ operation and it is much efficient if all $k_l \ll n$ where k_l is the rank of $\tilde{\mathbf{W}}_{r_i,l}$.

Update of the determinant $\log |\mathbf{Q}^{(L)}|_+$

We also extend Eq 9 to

$$\begin{aligned}
\log |\mathbf{Q}^{(L)}|_+ &= -\log |\mathbf{U}_L^T \mathbf{H} \mathbf{U}_L| \\
&= -\log \left| \mathbf{I}_{n-p} \delta_e + \mathbf{\Lambda}_c + \sum_{l=1}^L \mathbf{U}_L^T \tilde{\mathbf{W}}_{r_i,l} \tilde{\mathbf{\Gamma}}_{r_i,l} \tilde{\mathbf{W}}_{r_i,l}^T \mathbf{U}_L \right| \\
&= -\log \left| \mathbf{I}_{n-p} \delta_e + \mathbf{\Lambda}_c + \sum_{l=1}^{L-1} \mathbf{A}^{(l)} \mathbf{B}^{(l)} + \mathbf{A}^{(L)} \mathbf{B}^{(L)} \right| \\
&= -\log \left| \mathbf{I}_{n-p} \delta_e + \mathbf{\Lambda}_c + \sum_{l=1}^{L-1} \mathbf{A}^{(l)} \mathbf{B}^{(l)} \right| \\
&\quad - \log \left| \mathbf{I}_{n-p} + \left(\mathbf{I}_{n-p} \delta_e + \mathbf{\Lambda}_c + \sum_{l=1}^{L-1} \mathbf{A}^{(l)} \mathbf{B}^{(l)} \right)^{-1} \mathbf{A}^{(L)} \mathbf{B}^{(L)} \right| \\
&= \log |\mathbf{Q}^{(L-1)}|_+ - \log |\mathbf{I}_{k_L} + \mathbf{B}^{(L)} \mathbf{Q}^{(L-1)} \mathbf{A}^{(L)}| \quad (19)
\end{aligned}$$

Here, we used Proposition 5, 6, 7 and the Sylvester's determinant identity, and we also define a $n-p \times k_L$ matrix $\mathbf{A}^{(L)} = \mathbf{U}_L^T \tilde{\mathbf{W}}_{r_i,L} \tilde{\mathbf{\Gamma}}_{r_i,L}$ and a $k \times n-p$ matrix $\mathbf{B}^{(L)} = \tilde{\mathbf{W}}_{r_i,L}^T \mathbf{U}_L$. Computation of $\log |\mathbf{Q}^{(L)}|_+$ can be realized by repeating Eq 19 after setting $\log |\mathbf{Q}^{(0)}|_+ = -\log |\mathbf{I}_{n-p} \delta_e + \mathbf{\Lambda}_c|$. Here, when calculating the second term of Eq 19, the notes mentioned for Eq 9 still exists in this case.

The LR test for dominance and epistatic effects

To test the dominance or epistatic effects of each SNP-set, it should be assumed that the term of additive effects of each SNP-set is included both in the null model and the alternative model. Therefore we should assume the model such as Eq 16 for the alternative model in this case.

If we define $\mathbf{K}_{r_i,d}$ is the dominance genetic matrix of \mathbf{W}_{r_i} , $\tilde{\mathbf{W}}_{r_i,d}$ ($\tilde{\mathbf{W}}_{r_i}$ for the dominance effects) corresponds to

$$\tilde{\mathbf{W}}_{r_i,d} = \frac{\mathbf{Z}_{r_i} (\mathbf{1}_{m_{r_i} \times M_i} - \|\mathbf{W}_{r_i}\|)}{2 \cdot \sum_{m=1}^{M_i} p_m (1-p_m) (1-p_m (1-p_m))}, \quad (20)$$

and $\tilde{\mathbf{\Gamma}}_{r_i,d}$ is

$$\tilde{\mathbf{\Gamma}}_{r_i,d} = \mathbf{I}_{M_i} \gamma_{r_i,d}, \quad (21)$$

where $\gamma_{r_i,d}$ is the weight for the dominance effect of i_{th} SNP-set to be estimated [14]. Therefore, to test the significance of dominance effects, compare the restricted log likelihood of the model including both additive and dominance effects such as Eq 16 ($l = 2$) with that of the model including only additive effects such as Eq 1.

For the random effects of epistatic effects between two SNP-sets consisting of the same number of SNPs M , $\tilde{\mathbf{W}}_{r_{ij},aa}$, $\tilde{\mathbf{W}}_{r_{ij},ad}$, $\tilde{\mathbf{W}}_{r_{ij},da}$ and $\tilde{\mathbf{W}}_{r_{ij},dd}$ ($\tilde{\mathbf{W}}_{r_{ij}}$ for additive \times additive, additive \times dominance, dominance \times additive and dominance \times dominance epistatic effects between the i_{th} and the j_{th} SNP-sets of interest respectively) correspond to

$$\tilde{\mathbf{W}}_{r_{ij},aa} = \tilde{\mathbf{W}}_{r_i,a} \circ \tilde{\mathbf{W}}_{r_j,a}, \quad (22)$$

$$\tilde{\mathbf{W}}_{r_{ij},ad} = \tilde{\mathbf{W}}_{r_i,a} \circ \tilde{\mathbf{W}}_{r_j,d}, \quad (23)$$

$$\tilde{\mathbf{W}}_{r_{ij},da} = \tilde{\mathbf{W}}_{r_i,d} \circ \tilde{\mathbf{W}}_{r_j,a}, \quad (24)$$

$$\tilde{\mathbf{W}}_{r_{ij},dd} = \tilde{\mathbf{W}}_{r_i,d} \circ \tilde{\mathbf{W}}_{r_j,d}, \quad (25)$$

where $\tilde{\mathbf{W}}_{r_i,a}$ is defined by Eq 10 and $\tilde{\mathbf{W}}_{r_j,d}$ is defined by Eq 20 [15, 14]. On the other hand, corresponding $\tilde{\mathbf{\Gamma}}_{r_{ij},aa}$, $\tilde{\mathbf{\Gamma}}_{r_{ij},ad}$, $\tilde{\mathbf{\Gamma}}_{r_{ij},da}$ and $\tilde{\mathbf{\Gamma}}_{r_{ij},dd}$ are

$$\tilde{\mathbf{\Gamma}}_{r_{ij},aa} = \mathbf{I}_M \gamma_{r_{ij},aa}, \quad (26)$$

$$\tilde{\mathbf{\Gamma}}_{r_{ij},ad} = \mathbf{I}_M \gamma_{r_{ij},ad}, \quad (27)$$

$$\tilde{\mathbf{\Gamma}}_{r_{ij},da} = \mathbf{I}_M \gamma_{r_{ij},da}, \quad (28)$$

$$\tilde{\mathbf{\Gamma}}_{r_{ij},dd} = \mathbf{I}_M \gamma_{r_{ij},dd}, \quad (29)$$

where $\gamma_{r_{ij},aa}$, $\gamma_{r_{ij},ad}$, $\gamma_{r_{ij},da}$ and $\gamma_{r_{ij},dd}$ are weights for each epistatic effects to be estimated. In this case, to test the significance of these epistatic effects, compare the restricted log likelihood of the model including additive and dominance effects of two SNP-set and four epistatic effects such as Eq 16 ($l = 8$) with that of the model including only additive and dominance effects of two SNP-set such as Eq 16 ($l = 4$).

Propositions

Proposition 1. $\mathbf{SX} = \mathbf{0}$, $\mathbf{PX} = \mathbf{0}$ and $\mathbf{U}_L^T \mathbf{X} = \mathbf{0}$.

Proof.

$$\mathbf{SX} = \mathbf{X} - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} = \mathbf{0},$$

$$\mathbf{PX} = \mathbf{X} - \mathbf{X} (\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{H}^{-1} \mathbf{X} = \mathbf{0},$$

and because $\mathbf{PX} = \mathbf{0}$, $\mathbf{P}_L \mathbf{X} = \mathbf{0}$. Therefore, $\mathbf{U}_L^T \mathbf{X} = \mathbf{0}$. □

Proposition 2. $\mathbf{S}^2 = \mathbf{S}$.

Proof.

$$\begin{aligned} \mathbf{S}^2 &= \left(\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right) \left(\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right), \\ &= \mathbf{I}_n - 2 \cdot \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ &= \mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ &= \mathbf{S} \end{aligned}$$

This characteristic is called idempotent. □

Proposition 3. $\mathbf{PS} = \mathbf{P}$ and $\mathbf{SP} = \mathbf{S}$.

Proof.

$$\begin{aligned} \mathbf{PS} &= \left(\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \right) \left(\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right) \\ &= \mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ &\quad + \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ &= \mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \\ &= \mathbf{P}, \end{aligned}$$

and

$$\begin{aligned} \mathbf{SP} &= \left(\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right) \left(\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \right) \\ &= \mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ &\quad + \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \\ &= \mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ &= \mathbf{S} \end{aligned}$$

□

Proposition 4. $\mathbf{P}^2 = \mathbf{P}$.

Proof.

$$\begin{aligned} \mathbf{P}^2 &= \left(\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \right) \left(\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \right) \\ &= \mathbf{I}_n - 2 \times \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \\ &\quad + \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \\ &= \mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \\ &= \mathbf{P} \end{aligned}$$

□

Proposition 5. $\mathbf{P}^T \mathbf{H}^{-1} \mathbf{P} = (\mathbf{SHS})^\dagger$.

Proof.

From Proposition 3 and 4,

$$\begin{aligned} (\mathbf{SHS}) (\mathbf{P}^T \mathbf{H}^{-1} \mathbf{P}) (\mathbf{SHS}) &= \mathbf{SHP}^T \mathbf{H}^{-1} \mathbf{PHS} \\ &= \mathbf{SPPHS} \\ &= \mathbf{SHS}, \end{aligned}$$

and

$$\begin{aligned} (\mathbf{P}^T \mathbf{H}^{-1} \mathbf{P}) (\mathbf{SHS}) (\mathbf{P}^T \mathbf{H}^{-1} \mathbf{P}) &= \mathbf{P}^T \mathbf{H}^{-1} \mathbf{PHP}^T \mathbf{H}^{-1} \mathbf{P} \\ &= \mathbf{P}^T \mathbf{H}^{-1} \mathbf{PPP} \\ &= \mathbf{P}^T \mathbf{H}^{-1} \mathbf{P}. \end{aligned}$$

Moreover, the product of $\mathbf{P}^T \mathbf{H}^{-1} \mathbf{P}$ and \mathbf{SHS} is commutative and becomes a Hermitian (in this case, real symmetric) matrix.

$$(\mathbf{SHS}) (\mathbf{P}^T \mathbf{H}^{-1} \mathbf{P}) = (\mathbf{P}^T \mathbf{H}^{-1} \mathbf{P}) (\mathbf{SHS}) = \mathbf{S}$$

Therefore, from the definition of Moore-Penrose pseudo-inverse matrix,

$$\mathbf{P}^T \mathbf{H}^{-1} \mathbf{P} = (\mathbf{SHS})^\dagger = \mathbf{U}_L \mathbf{\Lambda}_L^{-1} \mathbf{U}_L^T$$

□

Proposition 6. $\mathbf{Q} = \mathbf{P}^T \mathbf{H}^{-1} \mathbf{P}$.

Proof.

By using Proposition 4,

$$\mathbf{P}^T \mathbf{H}^{-1} \mathbf{P} = \mathbf{H}^{-1} \mathbf{PP} = \mathbf{H}^{-1} \mathbf{P} = \mathbf{Q}$$

□

Proposition 7. $\mathbf{U}_L \mathbf{U}_L^T = \mathbf{S}$.

Proof.

By using Proposition 5,

$$\begin{aligned} \mathbf{S} &= (\mathbf{SHS}) (\mathbf{P}^T \mathbf{H}^{-1} \mathbf{P}) = (\mathbf{SHS}) (\mathbf{SHS})^\dagger \\ &= \mathbf{U}_L \mathbf{\Lambda}_L \mathbf{U}_L^T \mathbf{U}_L \mathbf{\Lambda}_L^{-1} \mathbf{U}_L^T \\ &= \mathbf{U}_L \mathbf{U}_L^T \end{aligned}$$

□

Proposition 8. $\log |\mathbf{P}_L \mathbf{V} \mathbf{P}_L^T| = -\log |\mathbf{P}^T \mathbf{V}^{-1} \mathbf{P}|_+$.

Proof.

First,

$$\mathbf{P}_L \mathbf{V} \mathbf{P}_L^T = \mathbf{U}_L^T \mathbf{H} \mathbf{U}_L \sigma_c^2 = \mathbf{\Lambda}_L \sigma_c^2,$$

because

$$\mathbf{U}_L \mathbf{U}_L^T \mathbf{H} \mathbf{U}_L \mathbf{U}_L^T = \mathbf{S} \mathbf{H} \mathbf{S} = \mathbf{U}_L \mathbf{\Lambda}_L \mathbf{U}_L^T,$$

from Proposition 7. Then, by using Proposition 5

$$\begin{aligned} \log |\mathbf{\Lambda}_L| &= -\log |\mathbf{\Lambda}_L^{-1}| \\ &= -\log |\mathbf{U}_L \mathbf{\Lambda}_L^{-1} \mathbf{U}_L^T|_+ \\ &= -\log |\mathbf{P}^T \mathbf{H}^{-1} \mathbf{P}|_+ \end{aligned}$$

Therefore,

$$\log |\mathbf{P}_L \mathbf{V} \mathbf{P}_L^T| = \log |\mathbf{\Lambda}_L \sigma_c^2| = -\log |\mathbf{P}^T \mathbf{V}^{-1} \mathbf{P}|_+$$

□

Proposition 9. $\mathbf{P}_L^T (\mathbf{P}_L \mathbf{V} \mathbf{P}_L^T)^{-1} \mathbf{P}_L = \mathbf{V}^{-1} \mathbf{P}$.

Proof. By using Proposition 8 and 6,

$$\begin{aligned} \mathbf{P}_L^T (\mathbf{P}_L \mathbf{V} \mathbf{P}_L^T)^{-1} \mathbf{P}_L &= \mathbf{U}_L (\mathbf{U}_L^T \mathbf{V} \mathbf{U}_L)^{-1} \mathbf{U}_L^T \\ &= \mathbf{U}_L (\mathbf{\Lambda}_L \sigma_c^2)^{-1} \mathbf{U}_L^T \\ &= \mathbf{P}^T \mathbf{H}^{-1} \mathbf{P} / \sigma_c^2 \\ &= \mathbf{Q} / \sigma_c^2 = \mathbf{V}^{-1} \mathbf{P} \end{aligned}$$

□

Proposition 10. $\frac{1}{2} \log |\mathbf{Q}|_+ = -\frac{1}{2} \log |\mathbf{H}| + \frac{1}{2} \log |\mathbf{X}^T \mathbf{X}| - \frac{1}{2} \log |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}|$.

Proof. By using Proposition 6, 5 and 7,

$$\begin{aligned}
|\mathbf{Q}|_+ &= |\mathbf{P}^T \mathbf{H}^{-1} \mathbf{P}|_+ \\
&= |\mathbf{S} \mathbf{H} \mathbf{S}|_+^{-1} \\
&= |\mathbf{\Lambda}_L|^{-1} \cdot |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}| \cdot |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}|^{-1} \\
&= |\mathbf{U}_L^T \mathbf{H} \mathbf{U}_L|^{-1} \cdot \left| (\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{H}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X})^{-1} \right|^{-1} \cdot |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}|^{-1} \\
&= |\mathbf{U}_L^T \mathbf{H} \mathbf{U}_L|^{-1} \cdot \left| (\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{H}^{-1} \mathbf{H} \mathbf{H}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X})^{-1} \right|^{-1} \cdot |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}|^{-1} \\
&= \left| \begin{bmatrix} \mathbf{U}_L^T \mathbf{H} \mathbf{U}_L & \mathbf{0} \\ \mathbf{0} & (\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{H}^{-1} \mathbf{H} \mathbf{H}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X})^{-1} \end{bmatrix} \right|^{-1} \cdot |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}|^{-1}
\end{aligned}$$

Here, we define $\mathbf{M} = (\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{H}^{-1}$ to shorten the notation.

$$\begin{aligned}
&= \left| \begin{bmatrix} \mathbf{U}_L^T \mathbf{H} \mathbf{U}_L & \mathbf{U}_L^T \mathbf{H} \mathbf{M}^T \\ \mathbf{M} \mathbf{H} \mathbf{U}_L & \mathbf{M} \mathbf{H} \mathbf{M}^T \end{bmatrix} \right|^{-1} \cdot |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}|^{-1} \\
&\mathbf{M} \mathbf{H} \mathbf{U}_L = \mathbf{0} \text{ because } \mathbf{X}^T \mathbf{U}_L = \mathbf{0} \text{ (Proposition 1).}
\end{aligned}$$

$$\begin{aligned}
&= \left| \begin{bmatrix} \mathbf{U}_L^T \\ \mathbf{M} \end{bmatrix} \mathbf{H} \begin{bmatrix} \mathbf{U}_L^T \\ \mathbf{M} \end{bmatrix}^T \right|^{-1} \cdot |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}|^{-1} \\
&= |\mathbf{H}|^{-1} \cdot \left| \begin{bmatrix} \mathbf{U}_L^T \\ \mathbf{M} \end{bmatrix} \begin{bmatrix} \mathbf{U}_L^T \\ \mathbf{M} \end{bmatrix}^T \right|^{-1} \cdot |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}|^{-1} \\
&= |\mathbf{H}|^{-1} \cdot \left| \begin{bmatrix} \mathbf{U}_L^T \mathbf{U}_L & \mathbf{U}_L^T \mathbf{M}^T \\ \mathbf{M} \mathbf{U}_L & \mathbf{M} \mathbf{M}^T \end{bmatrix} \right|^{-1} \cdot |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}|^{-1} \\
&\text{Using the well-known formula for the determinant of block-matrix,} \\
&= |\mathbf{H}|^{-1} \cdot |\mathbf{U}_L^T \mathbf{U}_L|^{-1} \cdot \left| \mathbf{M} \mathbf{M}^T - \mathbf{M} \mathbf{U}_L (\mathbf{U}_L^T \mathbf{U}_L)^{-1} \mathbf{U}_L^T \mathbf{M}^T \right|^{-1} \cdot |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}|^{-1} \\
&= |\mathbf{H}|^{-1} \cdot \left| \mathbf{M} \mathbf{M}^T - \mathbf{M} \mathbf{S} \mathbf{M}^T \right|^{-1} \cdot |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}|^{-1} \\
&= |\mathbf{H}|^{-1} \cdot \left| \mathbf{M} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{M}^T \right|^{-1} \cdot |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}|^{-1} \\
&= |\mathbf{H}|^{-1} \cdot |\mathbf{X}^T \mathbf{X}| \cdot |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}|^{-1}
\end{aligned}$$

Therefore,

$$\frac{1}{2} \log |\mathbf{Q}|_+ = -\frac{1}{2} \log |\mathbf{H}| + \frac{1}{2} \log |\mathbf{X}^T \mathbf{X}| - \frac{1}{2} \log |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}|$$

If we plus the both sides of the equation by $-(n-p) \log(\sigma_c^2)/2$, we get

$$\frac{1}{2} \log |\mathbf{Q}/\sigma_c^2|_+ = -\frac{1}{2} \log |\mathbf{V}| + \frac{1}{2} \log |\mathbf{X}^T \mathbf{X}| - \frac{1}{2} \log |\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}|,$$

which is also the frequently used form of the formula. \square

Details about the material and simulations

Details about the simulation of phenotype data

As described in the main text, the model as Eq 30 was used for the simulation of phenotypic values.

$$\mathbf{y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \mathbf{X}_3\beta_3 + \mathbf{Z}\mathbf{u} + \mathbf{e}, \quad (30)$$

where \mathbf{y} is the vector of simulated phenotypic values of 414 accessions, \mathbf{X}_1 , \mathbf{X}_2 and \mathbf{X}_3 correspond to three quantitative trait nucleotides (QTNs) scored as -1, 0 or 1 (hereinafter, referred to as “QTN1”, “QTN2” and “QTN3” respectively), β_1 , β_2 and β_3 are scalars representing the effects of the three QTNs, \mathbf{u} is the vector of polygenetic effects, and \mathbf{e} is the vector of the residuals. Here, we will explain the details about the last two terms.

The polygenetic effects \mathbf{u} in Eq 31 were sampled from the multivariate normal distribution whose variance-covariance matrix was the additive genetic relationship matrix \mathbf{A} , and were normalized so that their variance was equal to that of three QTN effects.

$$\mathbf{u} \sim \text{MVN}(\mathbf{0}, \mathbf{G}) \quad (31)$$

$\mathbf{G} = \mathbf{A}\sigma_A^2$ is the genetic covariance matrix, and the additive genetic variance σ_A^2 is automatically determined from the relationship with the heritability. In this study, the additive numerator relationship matrix was estimated based on the marker genotype of 112,630 SNPs using the “A.mat” function of the R package “rrBLUP” version 4.6 [16].

The residual \mathbf{e} in Eq 32 was sampled identically and independently from a normal distribution, and then was normalized so that the narrow-sense heritability was equal to 0.6.

$$\mathbf{e} \sim \text{MVN}(\mathbf{0}, \mathbf{I}\sigma_e^2), \quad (32)$$

where \mathbf{I} is an identity matrix, and σ_e^2 is the residual variance determined so that the heritability is equal to 0.6.

Details about the evaluation of RAINBOW

Details of four methods

Here, we will explain the details of four methods for the evaluation of the RAINBOW performance.

Single-SNP GWAS The single-SNP GWAS model was similar to Eq 2, but this method regarded each SNP as a fixed effect and tested the significance of each SNP effect one by one [17]. In this study, we used the method called “P3D”, which first estimated σ_c^2 and σ_e^2 by solving the mixed model without SNP effects (Eq 2) by REML, and then tested each marker using the information of the estimated variance $\hat{\sigma}_c^2$ and $\hat{\sigma}_e^2$ [6, 18, 19]. We used the “GWAS” function of the R package “rrBLUP” version 4.6 [16].

HGF method HGF tests the significance of each haplotype block effect. In this method, all accessions were divided into k_{r_i} groups based on the haplotypes. The haplotype groups were then represented by dummy variables. The significance of the haplotype effect was tested by regarding a matrix of dummy variables as fixed effects [20]. In this study, the k-medoids method (the “pam” function of the R package “cluster” version 2.0.6) [21] or UPGMA (unweighted pair group method with arithmetic mean, “hclust” and “cutree” function of the R package “stats”) [22] were employed to group haplotypes. We set the number of groups k_{r_i} as 2,3 or 4. Therefore, we applied 6 different HGF models in total. To apply HGF models, we modified the R package “rrBLUP” version 4.6 to enable the use of a matrix of dummy variables as fixed effects [16].

SKAT The SKAT is a SNP-set GWAS method, and is widely used in human genetics [23]. This method was invented to detect rare alleles, but the SKAT model does not include the random effects term that accounts for family relatedness, which may cause false positives for GWAS in materials with a strong population structure or family relatedness. In this study, we used the “SKAT” function of the R package “SKAT” version 1.3.2.1 and performed SNP-set GWAS by regarding each haplotype block as each SNP-set.

RAINBOW The methods for RAINBOW were described in the main text. In this study, we performed haplotype-based GWAS for RAINBOW by regarding each haplotype block as each SNP-set. We used the linear kernel of each haplotype block for the Gram matrix \mathbf{K}_{r_i} .

References

- [1] Patterson HD, Thompson R. Recovery of inter-block information when block sizes are unequal. *Biometrika*. 1971;58(3):545–554.
- [2] Harville DA. Bayesian Inference for Variance Components Using Only Error Contrasts. *Biometrika*. 1974;61(2):383–385.
- [3] Sorensen D, Gianola D. Likelihood, Bayesian, and MCMC Methods in Quantitative Genetics. Heidelberg, New York: Springer; 2002.
- [4] Croebel RR, Searle SR. Restricted Maximum Likelihood (REML) Estimation of Variance Components in the Mixed Model. *Technometrics*. 1976;18(1):31–38.
- [5] Welham SJ, Thompson R. Likelihood ratio tests for fixed model terms using residual maximum likelihood. *J R Stat Soc Ser B Stat Methodol*. 1997;59(3):701–714.
- [6] Kang HM, Zaitlen NA, Wade CM, Kirby A, Heckerman D, Daly MJ, et al. Efficient Control of Population Structure in Model Organism Association Mapping. *Genetics*. 2008;178(3):1709–1723.
- [7] Neyman J, Pearson ES. On the Use and Interpretation of Certain Test Criteria for Purposes of Statistical Inference. *Biometrika*. 1928;20A(1-2):175–240.
- [8] Wilks SS. The Large-Sample Distribution of the Likelihood Ratio for Testing Composite Hypotheses. *Ann Math Stat*. 1938;9(1):60–62.
- [9] Listgarten J, Lippert C, Kang EY, Xiang J, Kadie CM, Heckerman D. A powerful and efficient set test for genetic markers that handles confounders. *Bioinformatics*. 2013;29(12):1526–1533.
- [10] Lippert C, Xiang J, Horta D, Widmer C, Kadie C, Heckerman D, et al. Greater power and computational efficiency for kernel-based association testing of sets of genetic variants. *Bioinformatics*. 2014;30(22):3206–3214.
- [11] Byrd R, Lu P, Nocedal J, Zhu C. A Limited Memory Algorithm for Bound Constrained Optimization. *SIAM Journal on Scientific Computing*. 1995;16(5):1190–1208.
- [12] Zhou X, Stephens M. Genome-wide efficient mixed-model analysis for association studies. *Nat Genet*. 2012;44(7):821–824.
- [13] Endelman JB, Jannink JL. Shrinkage Estimation of the Realized Relationship Matrix. *G3 (Bethesda)*. 2012;2(11):1405–1413.
- [14] Su G, Christensen OF, Ostersen T, Henryon M, Lund MS. Estimating Additive and Non-Additive Genetic Variances and Predicting Genetic Merits Using Genome-Wide Dense Single Nucleotide Polymorphism Markers. *PLoS One*. 2012;7(9):1–7.
- [15] Henderson CR. Best Linear Unbiased Prediction of Nonadditive Genetic Merits in Noninbred Populations. *J Anim Sci*. 1985;60(1):111–117.

- [16] Endelman JB. Ridge Regression and Other Kernels for Genomic Selection with R Package rrBLUP. *Plant Genome J.* 2011;4(3):250.
- [17] Yu J, Pressoir G, Briggs WH, Vroh Bi I, Yamasaki M, Doebley JF, et al. A unified mixed-model method for association mapping that accounts for multiple levels of relatedness. *Nat Genet.* 2006;38(2):203–8.
- [18] Kang HM, Sul JH, Service SK, Zaitlen NA, Kong SY, Freimer NB, et al. Variance component model to account for sample structure in genome-wide association studies. *Nat Genet.* 2010;42(4):348–354.
- [19] Zhang Z, Ersoz E, Lai CQ, Todhunter RJ, Tiwari HK, Gore MA, et al. Mixed linear model approach adapted for genome-wide association studies. *Nat Genet.* 2010;42(4):355–360.
- [20] Yano K, Yamamoto E, Aya K, Takeuchi H, Lo PC, Hu L, et al. Genome-wide association study using whole-genome sequencing rapidly identifies new genes influencing agronomic traits in rice. *Nat Genet.* 2016;48(8):927–934.
- [21] Maechler M, Rousseeuw P, Struyf A, Hubert M, Hornik K. *cluster: Cluster Analysis Basics and Extensions*; 2018.
- [22] R Core Team. *R: A Language and Environment for Statistical Computing*; 2019. Available from: <https://www.R-project.org/>.
- [23] Wu MC, Lee S, Cai T, Li Y, Boehnke M, Lin X. Rare-variant association testing for sequencing data with the sequence kernel association test. *Am J Hum Genet.* 2011;89(1):82–93.