

① Consider SNPs in a given gene, run Han's code to get their scores in population p ($p=1, \dots, P$). Define the scores as

$(z_1^{(p)}, z_2^{(p)}, \dots, z_K^{(p)})$, K is the number of SNPs in the gene.

Note: get these scores directly, no need to scale back. Let the covariance matrix be $\Sigma^{(p)}$, (from Han's code, no need to scale it back, just take from the output, This matrix should be the one for the whole block, which includes several genes)

② From $\Sigma^{(p)}$, generate N copies of Z , put them in the matrix

$$Z^{(p)} = \begin{pmatrix} z_{1,1}, & z_{1,2}, & \dots, & z_{1,K} \\ z_{2,1}, & z_{2,2}, & \dots, & z_{2,K} \\ \vdots & \vdots & \ddots & \vdots \\ z_{N,1}, & z_{N,2}, & \dots, & z_{N,K} \end{pmatrix}$$

③ For SNP k , do the following steps, $k=1, \dots, K$.

(a) let $t_{0,k} = \max_p \left\{ \frac{|z_k^{(p)}|}{\sigma_k^{(p)}}, p=1, \dots, P \right\}$, where $\sigma_k^{(p)}$ is $\sqrt{\Sigma_{(k,k)}^{(p)}}$, square root of diagonal term in $\Sigma^{(p)}$.

(b) For each i -th row of $Z^{(p)}$, $p=1, \dots, P$, let

$$t_{i,k} = \max_p \left\{ \frac{|z_{i,k}^{(p)}|}{\sigma_k^{(p)}}, p=1, \dots, P \right\}, \text{ with the}$$

same $\sigma_k^{(p)}$ as in Step (a).

③ Form the following matrix

$$Q = \begin{pmatrix} t_{0,1} & t_{0,2} & \dots & t_{0,k} \\ t_{1,1} & t_{1,2} & \dots & t_{1,k} \\ \vdots & \vdots & \ddots & \vdots \\ t_{N,1} & t_{N,2} & \dots & t_{N,k} \end{pmatrix}$$

④ In Step 3, the max is over the populations that have scores Z for the SNP.

⑤ Convert each element of Q into a p-value. Let the corresponding matrix be R .

More details on Step ⑤ $S^{(p)} = \frac{Z}{\sigma}$ as above.

Suppose $t = \max\{S^{(p)} | p=1, \dots, L\}$, the maximum of L Z -scores. (L can be less than P as some populations do not have that SNP.)

Since $Z^{(p)}, p=1, \dots, L$ are independent standard normal, let X_1, X_2, \dots, X_L be i.i.d. normal random variables.

$$P(\max\{|X_p|, p=1, \dots, L\} < t) = \prod_{p=1}^L P(|X_p| < t) = (\Phi(t))^L$$

wrong! should be $(2\Phi(t) - 1)^L$

so the p-value for t can be written as

$$1 - (\Phi(t))^2 \rightarrow \text{should be } 1 - (2\Phi(t) - 1)^2$$

⑥ Treat R as the matrix for SNP-level p-value for ARTP for one population, and proceed as one population pathway analysis.

Note: The algorithm should be applied to a block at a time in each population. So, we should use a common block for all populations.

