

APPENDIX A: Proof of the unbiasedness of the r_s^2 estimate for a two population structured sample when loci are unlinked.

Let S denote the random Bernoulli variable which equals 1 if the sampled individual comes from the first population and 0 otherwise. Let t denote the probability of $S = 1$ then,

$$Esp(S) = t \tag{A1}$$

$$Var(S) = Esp(S) - (Esp(S))^2 = t(1-t) \tag{A2}$$

Let X^l denote the dummy random variable equaling 1 when an individual carries the A allele at locus l and 0 otherwise, then

$$\begin{aligned} Esp(X^l) &= Esp(S=0)Esp(X^l | S=0) + Esp(S=1)Esp(X^l | S=1) \\ &= (1-t)Esp(X^l | S=0) + tEsp(X^l | S=1) \end{aligned} \tag{A3}$$

Using equations (A1) and (A3), we get

$$\begin{aligned} Cov(X^l, S) &= Esp(X^l S) - Esp(X^l)Esp(S) \\ &= Esp(S=1)Esp(X^l | S=1) - t(1-t)Esp(X^l | S=0) - t^2Esp(X^l | S=1) \\ &= t(1-t)(Esp(X^l | S=0) - Esp(X^l | S=1)) \end{aligned} \tag{A4}$$

Using equations (A2) and (A4) we obtain

$$\begin{aligned} Cov(X^l, S)Var(S)Cov(S, X^m) &= t(1-t)(Esp(X^l | S=0) - Esp(X^l | S=1)) \\ &\quad (Esp(X^m | S=0) - Esp(X^m | S=1)) \end{aligned} \tag{A5}$$

On an other hand, theoretical developments of $Cov(X^l, X^m)$ give

$$\begin{aligned} Cov(X^l, X^m) &= Esp(X^l X^m) - Esp(X^l)Esp(X^m) \\ &= (1-t)Esp(X^l X^m | S=0) + tEsp(X^l X^m | S=1) - Esp(X^l)Esp(X^m) \\ &= (1-t)Cov(X^l, X^m | S=0) + tCov(X^l, X^m | S=1) - Esp(X^l)Esp(X^m) \\ &\quad + (1-t)Esp(X^l | S=0)Esp(X^m | S=0) + tEsp(X^l | S=1)Esp(X^m | S=1) \\ &= (1-t)Cov(X^l, X^m | S=0) + tCov(X^l, X^m | S=1) \\ &\quad + t(1-t)(Esp(X^l | S=0) - Esp(X^l | S=1))(Esp(X^m | S=0) - Esp(X^m | S=1)) \end{aligned} \tag{A6}$$

Thus, using equations (A5) and (A6) we obtain that

$$\begin{aligned} \text{Cov}(X^l | S, X^m | S) &= \text{Cov}(X^l, X^m) - \text{Cov}(X^l, S)\text{Var}(S)\text{Cov}(S, X^m) \\ &= (1-t)\text{Cov}(X^l, X^m | S = 0) + t\text{Cov}(X^l, X^m | S = 1) \end{aligned}$$

which equals 0 when the loci are unlinked.

APPENDIX B: Proof that the r_s^2 measure is the proportional factor to apply to sample size in order to achieve the same power of structure corrected association test at a SNP locus in linkage disequilibrium with the causal locus, as at the causal locus itself.

Let a trait, observed on a sample of size N , be explained by a causal locus l in the following linear model

$$Y = \mathbb{1}_N \mu + S\beta + X^l \theta^l + \varepsilon$$

where $Y = (y_1, \dots, y_i, \dots, y_N)^T$ is the vector of observed trait values, ε is the residual vector that is assumed to have expectation 0 and variance σ^2 , and (μ, β, θ^l) are the parameters for the mean, the structure effect and the causal locus effect, respectively.

The association t-test is equal to

$$t^l = \frac{\hat{\theta}^l}{\sqrt{\text{Var}(\hat{\theta}^l)}}$$

with

$$\text{Var}(\hat{\theta}^l) = \hat{\sigma}^2 ([\tilde{S}, \tilde{X}^l]^T [\tilde{S}, \tilde{X}^l])_{2,2}^{-1}$$

where $\hat{\sigma}^2$ is the estimate of σ^2 , \tilde{S} and \tilde{X}^l are the centered S and X^l matrices, respectively, and $_{2,2}$ denotes the second diagonal block of the matrix.

By definition of the sample variance-covariance matrix, we get

$$[\tilde{S}, \tilde{X}^l]^T [\tilde{S}, \tilde{X}^l] = (N-1) \begin{pmatrix} \Sigma_{S,S} & \Sigma_{S,X^l} \\ \Sigma_{X^l,S} & \Sigma_{X^l,X^l} \end{pmatrix}$$

The inversion of the block matrix gives

$$\begin{pmatrix} \Sigma_{S,S} & \Sigma_{S,X^l} \\ \Sigma_{X^l,S} & \Sigma_{X^l,X^l} \end{pmatrix}_{2,2}^{-1} = \frac{1}{(N-1)(\Sigma_{X^l,X^l} - \Sigma_{X^l,S} \Sigma_{S,S}^{-1} \Sigma_{S,X^l})}$$

So, asymptotically we get that t^l is Gaussian with variance 1 and expectation equal to $\sqrt{(N-1)(\Sigma_{X^l, X^l} - \Sigma_{X^l, S} \Sigma_{S, S}^{-1} \Sigma_{S, X^l})} \theta^l / \sigma$.

The t-test at the SNP locus m is $t^m = \frac{\hat{\theta}^m}{\sqrt{\text{Var}(\hat{\theta}^m)}}$ with $\text{Var}(\hat{\theta}^m) = \hat{\sigma}^2 ([\tilde{S}, \tilde{X}^m]^T [\tilde{S}, \tilde{X}^m])_{2,2}^{-1}$.

To find the expectation of t^m under the causal model, which is the correct one for the expectation of the data Y , it is necessary to calculate the second block of

$$([\tilde{S}, \tilde{X}^m]^T)^{-1} [\tilde{S}, \tilde{X}^m] [\tilde{S}, \tilde{X}^m]^T \text{Esp}(Y) = \begin{pmatrix} \Sigma_{S, S} & \Sigma_{S, X^m} \\ \Sigma_{X^m, S} & \Sigma_{X^m, X^m} \end{pmatrix}^{-1} \begin{pmatrix} \Sigma_{S, S} & \Sigma_{S, X^l} \\ \Sigma_{X^m, S} & \Sigma_{X^m, X^l} \end{pmatrix} \begin{pmatrix} s \\ \theta^l \end{pmatrix}$$

This second block is equal to $(\Sigma_{X^m, X^m} - \Sigma_{X^m, S} \Sigma_{S, S}^{-1} \Sigma_{S, X^m})^{-1} (\Sigma_{X^m, X^l} - \Sigma_{X^m, S} \Sigma_{S, S}^{-1} \Sigma_{S, X^l}) \theta^l$

We find, asymptotically,

$$\text{Esp}(t^m) \approx \sqrt{N-1} \frac{\Sigma_{X^m, X^l} - \Sigma_{X^m, S} \Sigma_{S, S}^{-1} \Sigma_{S, X^l}}{\sqrt{\Sigma_{X^m, X^m} - \Sigma_{X^m, S} \Sigma_{S, S}^{-1} \Sigma_{S, X^m}}} \theta^l / \sigma = \sqrt{\hat{r}_S^2(l, m)} \text{Esp}(t^l)$$

Thus, t^m is asymptotically Gaussian with variance 1 and expectation equal to $\sqrt{\hat{r}_S^2(l, m)} \text{Esp}(t^l)$. This finishes the first part of the proof showing that $\hat{r}_S^2(l, m)$ is the reducing power factor between the causal locus and a SNP locus in linkage disequilibrium.

Now, let us suppose that we get a N^m sample at the SNP locus and a N^l sample at the causal locus. Using that asymptotically $\hat{r}_S^2(l, m) \approx r_S^2(l, m)$ and that $\Sigma_{X^l, X^l} - \Sigma_{X^l, S} \Sigma_{S, S}^{-1} \Sigma_{S, X^l} \approx \text{Var}(X^l | S)$, we get

$$\text{Esp}(t^m) \approx \sqrt{r_S^2(l, m)} \sqrt{(N^m - 1) \text{Var}(X^l | S)} \theta^l / \sigma = \sqrt{r_S^2(l, m)} \sqrt{\frac{N^m - 1}{N^l - 1}} \text{Esp}(t^l)$$

Then, we obtain that the sample size has to be increased by a factor equal to $1/r_S^2(l, m)$ to achieve the same power at the SNP locus, compared to the power at the causal locus.