

Notes on the paper “Reconciling the analysis of IBD and IBS in complex trait studies” and other digressions

Daniel Sorensen

July 25, 2013

Contents

1	Inbreeding as a correlation of uniting gametes	2
2	Estimating relatedness from SNPs	3
2.1	Expected value of \widehat{F}_{12}	3
2.2	Variance of \widehat{F}_{12}	4
3	Estimating relationships between diploid individuals	5
3.1	Expected value of \widehat{A}_{ij}	5
3.2	The special case when $i = j$	7
3.3	The variance of \widehat{A}_{ij}	7
4	A genomic model	7
4.1	The covariance between two individuals	9
5	Marker effects are functions of disequilibria	11
5.1	Example 1: Marker effects fixed and genotypes random	11
5.2	Example 2: Marker effects random and genotypes fixed	12
6	Computing $\rho_{Y\widehat{Y}}$ with limited information	13
6.1	Predicting unobserved genetic value using observed marker information . .	14
6.2	The genomic variance	16
6.3	Cross-validation	16
6.4	Factorising the correlation between Y and \widehat{Y} <i>a la</i> Goddard	17
6.5	Marker effects as fixed parameters - least squares estimation	18
6.5.1	The case of 1 marker	18
6.5.2	The case of m_G markers	20
6.6	Marker effects as random variables - BLUP	22
6.7	Influence of degree of relationship on predictive ability	26

Below are some notes and reflections concerning the important paper by Powell et al. (2010).

The first point is related to the traditional way of calculating probabilities of IBD at a given locus or gametic relationships (between the two gametes of individual i , that gives rise to its inbreeding coefficient F_i , or between gametes of individuals i and j , that results in the gametic relationship between the gametes of these individuals, F_{ij}) from a known pedigree spanning several generations. The classical model assumes the existence of a base population, of size N say, consisting of $2N$ different alleles at a locus. IBD probabilities are estimated with reference to this base population.

Data on genetic markers can be used to estimate probabilities of IBD and often there is no obvious base population. The authors suggest that the base population can be arbitrary, and one choice could be the current generation involving the typed individuals. In this way, the gene frequencies of the observed markers refer to this cohort consisting of the typed individuals. The authors remind us that probabilities of IBD can be interpreted as correlations of uniting gametes (Wright, 1922). By defining the base population as the population of typed individuals, some values of F may be negative. However this has an interpretation in terms of the correlation coefficient; it means that the particular individual is less homozygous than the average.

1 Inbreeding as a correlation of uniting gametes

Let

$$X_i = \begin{cases} 0 & \text{if } G_i = A_2A_2 \\ 1 & \text{if } G_i = A_1A_2 \\ 2 & \text{if } G_i = A_1A_1 \end{cases}$$

and

$$X_{ij} = \begin{cases} 0 & \text{if } j = 1 \text{ and } G_{i1} = A_1, \text{ with probability } (1 - p) \\ 1 & \text{if } j = 2 \text{ and } G_{i2} = A_2, \text{ with probability } p, \end{cases}$$

where G_{ij} is the j th allele ($j = 1, 2$) of individual i . Then,

$$\begin{aligned} X_i &= X_{i1} + X_{i2}, \\ E(X_i) &= 2p, \\ Var(X_i) &= 2p(1 - p)(1 + F_i). \end{aligned} \tag{1}$$

For individual i , let

$$Z = \begin{cases} 0 & \text{if } X_{i1} \neq X_{i2}, \text{ with probability } (1 - F_i) \\ 1 & \text{if } X_{i1} \equiv X_{i2}, \text{ with probability } F_i. \end{cases}$$

The covariance between uniting gametes in individual i is defined as

$$\begin{aligned} Cov(X_{i1}, X_{i2}) &= E[Cov(X_{i1}, X_{i2}) | Z] + Cov[E(X_{i1} | Z), E(X_{i2} | Z)] \\ &= (Cov(X_{i1}, X_{i2}) | Z = 0)(1 - F_i) + (Cov(X_{i1}, X_{i2}) | Z = 1)F_i + 0 \\ &= 0 + Var(X_{ij})F_i = p(1 - p)F_i, \end{aligned} \tag{2}$$

and the correlation between uniting gametes in individual i is

$$\frac{Cov(X_{i1}, X_{i2})}{\sqrt{Var(X_{i1}) Var(X_{i1})}} = \frac{p(1-p) F_i}{p(1-p)} = F_i,$$

the inbreeding coefficient of individual i . If instead one computes the correlation between two arbitrary gametes k and l , from individuals i and j ,

$$\frac{Cov(X_{ik}, X_{jl})}{\sqrt{Var(X_{ik}) Var(X_{jl})}} = \frac{p(1-p) F_{kl}}{p(1-p)} = F_{kl},$$

the gametic relationship between gametes k and l .

2 Estimating relatedness from SNPs

Using the concept of correlation between uniting gametes, one can estimate the gametic relationship between gametes 1 and 2, using

$$\begin{aligned} \hat{F}_{12} &= \frac{Cov(X_{i1}, X_{j2})}{\sqrt{Var(X_{i1}) Var(X_{j2})}} \\ &= \frac{(X_{i1} - E(X_{i1}))(X_{j2} - E(X_{j2}))}{p(1-p)} \\ &= \frac{(X_{i1} - p)(X_{j2} - p)}{p(1-p)}. \end{aligned} \tag{3}$$

We note that the term $(X_{i1} - E(X_{i1}))(X_{j2} - E(X_{j2}))$ is strictly not equal to $Cov(X_{i1}, X_{j2}) = E[(X_{i1} - E(X_{i1}))(X_{j2} - E(X_{j2}))]$. It is best interpreted as a sample covariance with known means. A similar observation holds for several of the expressions below.

2.1 Expected value of \hat{F}_{12}

The expected value of \hat{F}_{12} (given the gene frequency p) is obtained as follows.

$$\begin{aligned} E[(X_{i1} - p)(X_{j2} - p)] &= E(X_{i1}, X_{j2}) - pE(X_{i1}) - pE(X_{j2}) + p^2 \\ &= E(X_{i1}, X_{j2}) - p^2. \end{aligned} \tag{4}$$

The first term on the right hand side is

$$\begin{aligned} E(X_{i1}, X_{j2}) &= Cov(X_{i1}, X_{j2}) + E(X_{i1}) E(X_{j2}) \\ &= Cov(X_{i1}, X_{j2}) + p^2. \end{aligned}$$

Therefore, substituting in (4),

$$\begin{aligned} E[(X_{i1} - p)(X_{j2} - p)] &= Cov(X_{i1}, X_{j2}) \\ &= p(1-p) F_{12} \end{aligned}$$

from (2). Therefore,

$$E\left(\widehat{F}_{12}\right) = \frac{p(1-p)F_{12}}{p(1-p)} = F_{12}.$$

2.2 Variance of \widehat{F}_{12}

The variance of \widehat{F}_{12} (given gene frequency, and assuming $F_{12} = 0$) is obtained as follows.

$$Var\left(\widehat{F}_{12}\right) = \frac{Var\left[(X_{i1} - p)(X_{j2} - p)\right]}{p^2(1-p)^2}.$$

The numerator is

$$\begin{aligned} Var\left[(X_{i1} - p)(X_{j2} - p)\right] &= Var\left[(X_{i1}X_{j2}) - pX_{i1} - pX_{j2} + p^2\right] \\ &= Var(X_{i1}X_{j2}) + 2p^2Var(X_{i1}) - 2Cov(X_{i1}X_{j2}, pX_{i1}) - 2Cov(X_{i1}X_{j2}, pX_{j2}) \\ &\quad + 2Cov(pX_{i1}, pX_{j2}). \end{aligned} \quad (5)$$

The expressions for each term are as follows. Since $F_{12} = 0$, X_{i1} and X_{j2} are independent. Then,

$$\begin{aligned} Var(X_{i1}X_{j2}) &= E(X_{i1}^2)E(X_{j2}^2) - [E(X_{i1})]^2[E(X_{j2})]^2 \\ &= [p(1-p) + p^2]^2 - p^4 \\ &= p^2(1-p^2). \\ 2p^2Var(X_{i1}) &= 2p^3(1-p). \end{aligned}$$

$$\begin{aligned} 2Cov(X_{i1}X_{j2}, pX_{i1}) &= 2[E(X_{i1}^2X_{j2}p) - E(X_{i1}X_{j2})E(X_{i1}p)] \\ &= 2[E(X_{i1}^2)E(X_{j2}p) - E(X_{i1})E(X_{j2})E(X_{i1}p)] \\ &= 2[(p(1-p) + p^2)p^2] - p^4 \\ &= 2p^3(1-p). \end{aligned}$$

$$2Cov(X_{i1}X_{j2}, pX_{j2}) = 2p^3(1-p).$$

$$2Cov(pX_{i1}, pX_{j2}) = 0.$$

Collecting terms,

$$\begin{aligned} Var\left[(X_{i1} - p)(X_{j2} - p)\right] &= p^2(1-p^2) + 2p^3(1-p) - 4p^3(1-p) \\ &= p^2(1-p^2) - 2p^3(1-p) \\ &= [p(1-p)]^2. \end{aligned} \quad (6)$$

Therefore,

$$Var\left(\widehat{F}_{12}\right) = \frac{Var\left[(X_{i1} - p)(X_{j2} - p)\right]}{p^2(1-p)^2} = 1.$$

3 Estimating relationships between diploid individuals

Expression (3) can readily be generalised to compute the coefficient of relationship between diploid individuals. At a given locus, for individual i , the random variable X_i defines the number of copies of the A_1 allele (known as gene content in the genetics literature). It can also be interpreted as a label for the genotype. The coefficient of relationship between individuals i and j is defined as

$$\begin{aligned}\hat{A}_{ij} &= \hat{F}_{ij} = \frac{Cov(X_i, X_j)}{\sqrt{Var(X_i) Var(X_j)}} \\ &= \frac{(X_i - 2p)(X_j - 2p)}{2p(1-p)}.\end{aligned}\tag{7}$$

the coefficient of relationship takes the form of a "sampling" correlation with known mean (gene frequencies) between the genotypes (represented by the random variable X , that can take values 0, 1 or 2). At this point we note that

$$\begin{aligned}Var(X_i) &= Var(X_{i1} + X_{i2}) \\ &= Var(X_{i1}) + Var(X_{i2}) + 2Cov(X_{i1}, X_{i2}) \\ &= p(1-p) + p(1-p) + 2p(1-p)F_i \\ &= 2p(1-p)(1+F_i),\end{aligned}$$

so the term $\sqrt{(1+F_i)(1+F_j)}$ is ignored in the denominator of (7).

3.1 Expected value of \hat{A}_{ij}

Given p the expected value is

$$\begin{aligned}E[(X_i - 2p)(X_j - 2p)] &= E(X_i, X_j) - 2pE(X_i) - 2pE(X_j) + 4p^2 \\ &= E(X_i, X_j) - 4p^2.\end{aligned}$$

The expectation of the first term is

$$\begin{aligned}E(X_i, X_j) &= Cov(X_i, X_j) + E(X_i)E(X_j) \\ &= Cov(X_i, X_j) + 4p^2.\end{aligned}$$

Substituting above gives

$$E[(X_i - 2p)(X_j - 2p)] = Cov(X_i, X_j),$$

as it should. The covariance can be written as

$$\begin{aligned}Cov(X_i, X_j) &= Cov(X_{i1} + X_{i2}, X_{j1} + X_{j2}) \\ &= 4Cov(X_{i1}, X_{j1}).\end{aligned}$$

The covariance term takes the form

$$\begin{aligned}
Cov(X_{i1}, X_{j1}) &= E[Cov(X_{i1}, X_{j1})|Z] + Cov[E(X_{i1}|Z), E(X_{j1}|Z)] \\
&= (Cov(X_{i1}, X_{j1})|Z=0) \Pr(X_{i1} \neq X_{j1}) + (Cov(X_{i1}, X_{j1})|Z=1) \Pr(X_{i1} \equiv X_{j1}) + 0 \\
&= Var(X_{i1}) \Pr(X_{i1} \equiv X_{j1}) \\
&= p(1-p) \Pr(X_{i1} \equiv X_{j1}).
\end{aligned}$$

The coefficient of coancestry between individuals i and j is defined as

$$\varphi_{ij} = \frac{1}{4} [\Pr(X_{i1} \equiv X_{j1}) + \Pr(X_{i1} \equiv X_{j2}) + \Pr(X_{i2} \equiv X_{j1}) + \Pr(X_{i2} \equiv X_{j2})].$$

Then,

$$\begin{aligned}
E[(X_i - 2p)(X_j - 2p)] &= Cov(X_i, X_j) \\
&= 4Cov(X_{i1}, X_{j1}) \\
&= 4p(1-p) \frac{1}{4} \varphi_{ij} \\
&= 2p(1-p) F_{ij},
\end{aligned} \tag{8}$$

where $F_{ij} = A_{ij} = 2\varphi_{ij}$ is the correlation between additive genetic values of related individuals, or additive genetic relationship. Therefore,

$$E(\hat{A}_{ij}) = \frac{2p(1-p) F_{ij}}{2p(1-p)} = A_{ij} = F_{ij}, \tag{9}$$

provided the term $\sqrt{(1+F_i)(1+F_j)}$ is ignored in the denominator of (7).

An important observation is that the coefficient of relationship obtained using marker information provides an "actual or true" relationship. This relationship differs from the pedigree based relationship (even when the number of markers tends to "infinity"), which is based on expected relationships. The expected relationship between full-sibs is 0.5. This expectation is taken over replicated full-sib individuals. But given two full-sibs, their actual relationship will vary around the expected value of 0.5, and the markers disclose this deviation.

3.2 The special case when $i = j$

When $i = j$ (the same individual),

$$\begin{aligned}
E(\hat{A}_i) &= \frac{E(X_i - 2p)^2}{2p(1-p)} \\
&= \frac{E(X_i^2) - 4p^2}{2p(1-p)} \\
&= \frac{Var(X_i) + [E(X_i)]^2 - 4p^2}{2p(1-p)} \\
&= \frac{Var(X_i)}{2p(1-p)} \\
&= \frac{Var(X_{i1} + X_{i2})}{2p(1-p)} \\
&= \frac{2p(1-p)(1 + F_i)}{2p(1-p)} \\
&= (1 + F_i).
\end{aligned}$$

3.3 The variance of \hat{A}_{ij}

Assume $A_{ij} = 0$ and $F_i = F_j = 0$. To compute this variance, first notice that

$$\begin{aligned}
(X_i - 2p)(X_j - 2p) &= [(X_{i1} - p) + (X_{i2} - p)][(X_{j1} - p) + (X_{j2} - p)] \\
&= (X_{i1} - p)(X_{j1} - p) + (X_{i1} - p)(X_{j2} - p) + (X_{i2} - p)(X_{j1} - p) + (X_{i2} - p)(X_{j2} - p).
\end{aligned}$$

Using (6), each of the 4 terms has variance $p^2(1-p)^2$. All covariance terms are of the form

$$\begin{aligned}
Cov[(X_{i1} - p)(X_{j1} - p), (X_{i1} - p)(X_{j2} - p)] &= E[(X_{i1} - p)^2(X_{j1} - p)(X_{j2} - p)] \\
&\quad - E[(X_{i1} - p)(X_{j1} - p)]E[(X_{i1} - p)(X_{j2} - p)] \\
&= E(X_{i1} - p)^2 E(X_{j1} - p) E(X_{j2} - p) - [E(X_{i1} - p)]^2 E(X_{j1} - p) E(X_{j2} - p) \\
&= E(X_{j1} - p) E(X_{j2} - p) Var(X_{i1}) = 0.
\end{aligned}$$

Therefore,

$$Var(\hat{A}_{ij}) = \frac{4p^2(1-p)^2}{4p^2(1-p)^2} = 1.$$

4 A genomic model

Consider the model for the record of individual j

$$y_j = \mu + W_1 b_1 + W_2 b_2 + e_j, \tag{10}$$

where the SNP effects are

$$\begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \sim N \left[\begin{pmatrix} 0 \\ 0 \end{pmatrix}, I\sigma_b^2 \right],$$

and

$$W_i = \frac{X_i - 2p_i}{\sqrt{2p_i(1-p_i)}}, \quad i = 1, 2 \quad (11)$$

which results in $E(W_i) = 0$ and $Var(W_i) = 1$. Let g_M denote the genomic value (captured by the SNPs), equal to

$$\begin{aligned} g_M &= W_1 b_1 + W_2 b_2 \\ &= Wb. \end{aligned}$$

Then, conditionally on W ,

$$\begin{aligned} Var(g_M | W_1, W_2) &= W Var(b) W' \\ &= WW' \sigma_b^2 \\ &= W_1^2 \sigma_b^2 + W_2^2 \sigma_b^2. \end{aligned} \quad (12)$$

and the unconditional variance is

$$Var(g_M) = \sigma_{g_M}^2 = 2\sigma_b^2 \quad (13)$$

known as the genomic variance, equal to the variance per SNP times the number of SNPs. Notice that (12) does not include a contribution due to disequilibrium between marker loci 1 and 2.

On the other hand, if (10) had been parameterised as

$$\begin{aligned} y_j &= \mu + X_1 b_1 + X_2 b_2 + e_j \\ &= \mu + b_1 X_1 + b_2 X_2 + e_j \\ &= \mu + b' X + e_j \\ &= \mu + g_M + e_j, \end{aligned}$$

and one conditions on the b 's instead,

$$\begin{aligned} Var(g_M | b_1, b_2) &= Var(bX | b) \\ &= b Var(X) b' \\ &= b' \begin{pmatrix} Var(X_1) & Cov(X_1 X_2) \\ Cov(X_1 X_2) & Var(X_2) \end{pmatrix} b \\ &= \begin{pmatrix} b_1 & b_2 \end{pmatrix} \begin{pmatrix} 2p_1(1-p_1) & 2D_{12} \\ 2D_{12} & p_2(1-p_2) \end{pmatrix} \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \\ &= 2b_1^2 p_1(1-p_1) + 2b_2^2 p_2(1-p_2) + 4b_1 b_2 D_{12}, \end{aligned}$$

where $2D_{12} = Cov(X_1, X_2)$. Here the X 's are treated as random and the inference is conditional on the b 's, and this incorporates the contribution due to disequilibrium. Unconditionally with respect to the b 's, the variance is of course equal to (13), because the b 's are a priori uncorrelated.

4.1 The covariance between two individuals

Consider the records of two individuals

$$\begin{aligned} y_i &= \mu + W_{1i}b_1 + W_{2i}b_2 + e_i, \\ y_j &= \mu + W_{1j}b_1 + W_{2j}b_2 + e_j. \end{aligned}$$

The covariance between their genomic breeding values is

$$\begin{aligned} \text{Cov}(g_{Mi}, g_{Mj} | W_1, W_2) &= \text{Cov}(W_{1i}b_1 + W_{2i}b_2, W_{1j}b_1 + W_{2j}b_2) \\ &= W_{1i}W_{1j}\sigma_b^2 + W_{2i}W_{2j}\sigma_b^2 \end{aligned}$$

and the other terms vanish because $\text{Cov}(b_1, b_2) = 0$. The unconditional covariance is

$$\begin{aligned} \text{Cov}(g_{Mi}, g_{Mj}) &= \text{Cov}[E(g_{Mi} | W_1, W_2), E(g_{Mj} | W_1, W_2)] + E[\text{Cov}(g_{Mi}, g_{Mj} | W_1, W_2)] \\ &= E[\text{Cov}(g_{Mi}, g_{Mj} | W_1, W_2)] \\ &= \sigma_b^2 E[W_{1i}W_{1j} + W_{2i}W_{2j}] \\ &= \sigma_b^2 (A_{ij,1} + A_{ij,2}) \\ &= 2\sigma_b^2 A_{ij} \\ &= \sigma_{g_M}^2 A_{ij}, \end{aligned}$$

where $A_{ij} = \frac{1}{2}(A_{ij,1} + A_{ij,2})$, the average (over markers) of the additive genetic relationship between individuals i and j .

In general, if the model in matrix notation is

$$\begin{aligned} \mathbf{y} &= \mathbf{1}\mu + \mathbf{W}\mathbf{b} + \mathbf{e} \\ &= \mathbf{1}\mu + \mathbf{g}_M + \mathbf{e}, \end{aligned}$$

with

$$\mathbf{b} \sim N(\mathbf{0}, \mathbf{I}\sigma_b^2)$$

then

$$\begin{aligned} \text{Var}(\mathbf{g}_M) &= \mathbf{W}\mathbf{W}'\sigma_b^2 \\ &= \frac{1}{m_G}\mathbf{W}\mathbf{W}'\sigma_{g_M}^2. \end{aligned}$$

The term $\frac{1}{m_G}\mathbf{W}\mathbf{W}'$ is the average (over SNPs) *realised* additive genetic relationship among individuals and

$$\sigma_{g_M}^2 = m_G\sigma_b^2 \quad (14)$$

is the sum of the variances of all (m_G) SNPs.

It is instructive to make quite explicit the structure of $\mathbf{W}\mathbf{W}'$. For two individuals and two markers it takes the form

$$\begin{aligned} \mathbf{W}\mathbf{W}' &= \begin{bmatrix} W_{1i} & W_{2i} \\ W_{1j} & W_{2j} \end{bmatrix} \begin{bmatrix} W_{1i} & W_{1j} \\ W_{2i} & W_{2j} \end{bmatrix} \\ &= \begin{bmatrix} W_{1i}W_{1i} + W_{2i}W_{2i} & W_{1i}W_{1j} + W_{2i}W_{2j} \\ W_{1j}W_{1i} + W_{2j}W_{2i} & W_{1j}W_{1j} + W_{2j}W_{2j} \end{bmatrix}. \end{aligned}$$

This shows that disequilibrium among loci cannot be accounted for when the inference is conditional on W .

On the other hand, if the model is parameterised as

$$\begin{aligned} y_i &= \mu + b_1 X_{1i} + b_2 X_{2i} + e_i \\ y_j &= \mu + b_1 X_{1j} + b_2 X_{2j} + e_j. \end{aligned}$$

In matrix notation,

$$y = 1\mu + bX + e,$$

where

$$b = \begin{bmatrix} b_1 & b_2 & 0 & 0 \\ 0 & 0 & b_1 & b_2 \end{bmatrix}$$

and

$$X' = [X_{1i}, X_{2i}, X_{1j}, X_{2j}].$$

and the X 's are treated as random with variance-covariance structure

$$\begin{bmatrix} Var(X_{1i}) & Cov(X_{1i}, X_{2i}) & Cov(X_{1i}, X_{1j}) & Cov(X_{1i}, X_{2j}) \\ & Var(X_{2i}) & Cov(X_{2i}, X_{1j}) & Cov(X_{2i}, X_{2j}) \\ & & Var(X_{1j}) & Cov(X_{1j}, X_{2j}) \\ & & & Var(X_{2j}) \end{bmatrix},$$

then the 2×2 off-diagonal block defines covariances between the two individuals. These covariances include terms involving the same locus, like $Cov(X_{1i}, X_{1j})$ and $Cov(X_{2i}, X_{2j})$, and terms involving different loci, like $Cov(X_{1i}, X_{2j})$ and $Cov(X_{2i}, X_{1j})$. The within individual terms include inbreeding contributions in the diagonals, and disequilibrium contributions in the off-diagonals.

The covariance between i and j takes the form

$$\begin{aligned} & Cov(X_{1i}b_1 + X_{2i}b_2, X_{1j}b_1 + X_{2j}b_2 | b_1, b_2) \\ &= b_1^2 Cov(X_{1i}, X_{1j}) + b_2^2 Cov(X_{2i}, X_{2j}) + b_1 b_2 Cov(X_{1i}, X_{2j}) + b_2 b_1 Cov(X_{2i}, X_{1j}). \end{aligned}$$

From (8), the within locus covariances between individuals i and j are

$$Cov(X_{ki}, X_{kj}) = 2p_k(1 - p_k) A_{ij,k}, \quad k = 1, 2,$$

where $A_{ij,k}$ is the additive genetic relationship between individuals i and j at locus k computed using the SNP markers. The between loci covariances between individuals i and j involve each 4 terms, each one of the form

$$D_{1i,2j} = \Pr(X_{11i} = 1, X_{21j} = 1) - \Pr(X_{11i} = 1) \Pr(X_{21j} = 1).$$

The conditional covariance between the genomic breeding values of individuals i and j is then

$$\begin{aligned} Cov(X_{1i}b_1 + X_{2i}b_2, X_{1j}b_1 + X_{2j}b_2 | b_1, b_2) &= 2b_1^2 p_1(1 - p_1) A_{ij,1} + 2b_2^2 p_2(1 - p_2) A_{ij,2} \\ &\quad + 4b_1 b_2 D_{1i,2j} + 4b_2 b_1 D_{1j,2i}. \end{aligned}$$

5 Marker effects are functions of disequilibria

5.1 Example 1: Marker effects fixed and genotypes random

Consider the task of predicting an unobserved genotypic value for a quantitative trait, $g = X_a a$, based on observed genotypic information from two markers X_1, X_2 . Here, X_a is a label for the unobserved genotype at the causal locus, X_1 and X_2 are labels for the observed marker genotypes, a is the effect of the causal locus. Assume that the variance-covariance structure of $[X_a, X_1, X_2]$ is

$$\begin{bmatrix} \sigma_{aa} & \sigma_{a1} & \sigma_{a2} \\ \sigma_{a1} & \sigma_{11} & \sigma_{12} \\ \sigma_{a2} & \sigma_{12} & \sigma_{22} \end{bmatrix}.$$

Then, given a , the best linear predictor of $X_a a$ takes the form

$$\hat{g} = g(X_1, X_2) = aE(X_a) + b_1(X_1 - E(X_1)) + b_2(X_2 - E(X_2)), \quad (15)$$

where, labeling $\hat{g}_i = b_i(X_i - E(X_i))$,

$$b_1 = a \frac{\sqrt{\sigma_{aa}}}{\sqrt{\sigma_{11}}} \left(\frac{r_{g, \hat{g}_1} - r_{g, \hat{g}_2} r_{\hat{g}_1, \hat{g}_2}}{1 - r_{\hat{g}_1, \hat{g}_2}^2} \right), \quad (16a)$$

$$b_2 = a \frac{\sqrt{\sigma_{aa}}}{\sqrt{\sigma_{22}}} \left(\frac{r_{g, \hat{g}_2} - r_{g, \hat{g}_1} r_{\hat{g}_1, \hat{g}_2}}{1 - r_{\hat{g}_1, \hat{g}_2}^2} \right) \quad (16b)$$

are the marker effects. In these expressions

$$\sigma_{ii} = 2p_i(1 - p_i), \quad i = a, 1, 2$$

and the correlation between the genotypic value at the causal locus, g , and $g_{M_i} = X_i b_i$, the genomic value of marker i take the form

$$\begin{aligned} r_{g, g_{M_i}} &= \frac{\text{Cov}(X_a a, X_i b_i)}{\sqrt{\text{Var}(X_a a) \text{Var}(X_i b_i)}} = \frac{\text{Cov}(X_a, X_i)}{\sqrt{\text{Var}(X_a) \text{Var}(X_i)}} \\ &= \frac{\sigma_{ai}}{\sqrt{\sigma_{aa} \sigma_{ii}}} \\ &= \frac{D_{ai}}{\sqrt{p_a(1 - p_a) p_i(1 - p_i)}} \quad i = 1, 2 \end{aligned} \quad (17)$$

which is a correlation of gene frequencies involving the causal locus and the marker. The term $\text{Cov}(X_a, X_i) = D_{ai}$ is the linkage disequilibrium (covariance of gene frequencies) involving the causal locus and the marker. Similarly, for the two marker loci,

$$\begin{aligned} r_{12} &= \frac{\text{Cov}(X_1, X_2)}{\sqrt{\text{Var}(X_1) \text{Var}(X_2)}} \\ &= \frac{D_{12}}{\sqrt{\sigma_{11} \sigma_{22}}}. \end{aligned}$$

In the model, the predictor of the genetic value g is \hat{g} . The correlation between these two quantities is

$$\begin{aligned}\rho_{g,\hat{g}} &= \frac{Cov(g, \hat{g})}{\sqrt{Var(g) Var(\hat{g})}} \\ &= \frac{Cov(X_a a, X_1 b_1 + X_2 b_2)}{\sqrt{Var(X_a a) Var(X_1 b_1 + X_2 b_2)}} \\ &= \frac{ab_1 Cov(X_a X_1) + ab_2 Cov(X_a X_2)}{\sqrt{a^2 Var(X_a) [b_1^2 Var(X_1) + b_2^2 Var(X_2) + 2b_1 b_2 Cov(X_1, X_2)]}}.\end{aligned}\quad (18)$$

In this approach, the X 's are treated as random and the partial regression coefficients b_1 and b_2 are parameters, assumed known in the context of best linear prediction. The point is that marker effects b_1 and b_2 in expressions (16) are functions of linkage disequilibria involving the causal locus and the markers. Also, unconditionally with respect to a , marker effects are correlated.

5.2 Example 2: Marker effects random and genotypes fixed

Assume that the model for the datum y (with mean 0) is

$$\mathbf{y} = W_a a + W_b b + e \quad (19)$$

where the W 's, defined in (11), are labels for the causal locus, whose effect is a , and for the marker locus, whose effect is b . Assume that

$$E(a, b, e) = (0, 0, 0)$$

and

$$Var(a, b, e) = diag(\sigma_a^2, \sigma_b^2, \sigma_e^2).$$

Write (19) as

$$y = Wu + e \quad (20)$$

where

$$W = (W_a, W_b), \quad u' = (a, b).$$

Then the posterior mean of u is the solution to

$$(W'W + D^{-1})\hat{u} = Z'y, \quad (21)$$

where

$$D^{-1} = \begin{bmatrix} \frac{\sigma_e^2}{\sigma_a^2} & 0 \\ 0 & \frac{\sigma_e^2}{\sigma_b^2} \end{bmatrix}.$$

The solution is

$$\begin{bmatrix} \hat{a} \\ \hat{b} \end{bmatrix} = \frac{1}{\frac{\sigma_e^2}{\sigma_a^2} \frac{\sigma_e^2}{\sigma_b^2} W_a W_a W_b W_b - (W_a W_b)^2} \begin{bmatrix} W_b W_b \frac{\sigma_e^2}{\sigma_b^2} & -W_a W_b \\ -W_a W_b & W_a W_a \frac{\sigma_e^2}{\sigma_a^2} \end{bmatrix} \begin{bmatrix} W_a \\ W_b \end{bmatrix} y. \quad (22)$$

The posterior mean of the marker effect is

$$\hat{b} = \frac{W_a W_b \left(\frac{\sigma_e^2}{\sigma_a^2} - 1 \right) W_b}{\frac{\sigma_e^2}{\sigma_a^2} \frac{\sigma_e^2}{\sigma_b^2} W_a W_a W_b W_b - (W_a W_b)^2} y$$

which indicates that it is a function of the disequilibrium between the marker and the causal locus ($W_a W_b$).

6 Computing $\rho_{Y\hat{Y}}$ with limited information

Here we address how imperfect knowledge of marker effects and ignorance about the true genetic model affect $\rho_{Y\hat{Y}}$ in a cross-validation scenario.

We first review briefly some basic results of best linear prediction. Consider two random variables Y and X of finite non-zero variance. A linear function $k + \beta X$ predicts Y with mean squared error $E(y - k + \beta X)^2$. This is minimised with $k = E(Y) - \beta E(X)$ and $\beta = Cov(Y, X) / Var(X)$. The best linear predictor of Y is

$$\hat{Y} = E(Y) + \beta(X - E(X)) \quad (23)$$

with

$$\begin{aligned} E(\hat{Y}) &= E_X [E(\hat{Y}|X)] \\ &= E(Y). \end{aligned}$$

One may write

$$Y = \hat{Y} + (Y - \hat{Y}) = \hat{Y} + e \quad (24)$$

where e , the prediction error, is uncorrelated with \hat{Y} . Then it is easily verified that

$$\begin{aligned} E(e) &= E(Y - \hat{Y}) = E(Y) - E(\hat{Y}) = 0, \\ Var(\hat{Y}) &= \frac{[Cov(Y, X)]^2}{Var(X)} = \beta^2 Var(X) = \rho_{YX}^2 Var(Y), \\ Cov(Y, \hat{Y}) &= Var(\hat{Y}), \text{ and therefore, } \beta_{Y, \hat{Y}} = 1, \\ Var(e) &= Var(Y) - \frac{[Cov(Y, X)]^2}{Var(X)} = (1 - \rho_{YX}^2) Var(Y), \\ Cov(\hat{Y}, e) &= 0, \\ \rho_{YX} &= \rho_{Y\hat{Y}}, \\ \rho_{Y\hat{Y}} &= \sqrt{\frac{Var(\hat{Y})}{Var(Y)}}. \end{aligned}$$

6.1 Predicting unobserved genetic value using observed marker information

In this section it is assumed that marker effects (labelled as b) are known parameters. Consider first the case of predicting $X_a a = g$ with one marker W_1 genotype only. Here X_a is a label for the genotype at the unobserved causal locus, a the substitution effect at the causal locus, and W_1 is a label for the genotype at the observed marker. Then, given a , the best linear predictor of g is

$$\hat{g} = aE(X_a) + b_1(W_1 - E(W_1)).$$

The decomposition $g = \hat{g} + e$ results in

$$Var(g) = \rho_{g\hat{g}}^2 Var(g) + (1 - \rho_{g\hat{g}}^2) Var(g).$$

The proportion of variance of the causal locus explained by the marker X_1 , given a , is

$$\rho_{g\hat{g}}^2 Var(g) = 2a^2 p_a (1 - p_a) r_{a1}^2 \quad (25)$$

where $Var(X_a a | a) = 2a^2 p_a (1 - p_a)$ and

$$\begin{aligned} \rho_{g\hat{g}} &= r_{a1} = \frac{Cov(X_a a, b_1 W_1)}{\sqrt{Var(X_a a | a) Var(b_1 W_1 | b_1)}} \\ &= \frac{D_{a1}}{\sqrt{p_a (1 - p_a) p_1 (1 - p_1)}}. \end{aligned}$$

where $Var(b_1 W_1 | b_1) = 2b_1^2 p_1 (1 - p_1)$. Expression (25) shows that if the causal locus is rare, the proportion of variance explained by the marker is very small. As an example, if the frequency of the causal locus is 0.1, that of the SNP is 0.4, and assuming that $D_{a1} = -0.1 \times 0.4 = -0.04$, then $r_{a1}^2 = 0.074$, and the variance explained by the marker is $0.013a^2$. Even when $r_{a1}^2 = 1$, if the frequency of the causal locus is 0.1, the variance explained by the marker is $0.18a^2$.

With many markers, the decomposition (24) still holds, but the expression for the (multiple) correlation between g and \hat{g} is less transparent. However the result $\rho_{g\hat{g}} = [Var(\hat{g}) / Var(g)]^{1/2}$ holds for any number of markers.

We consider now a model discussed in Goddard (2009), that assumes that each marker genotype is correlated with one causal genotype, but that pairs of markers and causal genotypes, are independent. The model assumes that there is an equal number of markers and causal genotypes. For simplicity we consider here 2 pairs, where each pair consists of a marker genotype and a causal genotype. The causal genotypic value of an individual is

$$g = X_1 a_1 + X_2 a_2$$

where X_i is a label (random variable) for the genotype at the i th causal locus, and a_i is the effect of the i th causal genotype. The conditional variance is

$$Var(g|a_1, a_2) = a_1^2 \sigma_{X_1}^2 + a_2^2 \sigma_{X_2}^2$$

where $Var(X_i) = \sigma_{X_i}^2$. The best linear predictor of g is

$$\hat{g} = E(g) + b_1(W_1 - E(W_1)) + b_2(W_2 - E(W_2)) \quad (26)$$

where W_i is a label (random variable) for the genotype at the i th marker locus, and b_i is the effect of the i th marker genotype. The derivation of (26) uses standard best linear prediction theory and is sketched at the end of this subsection. In (26),

$$\begin{aligned} b_i &= \frac{Cov(g, W_i|a_1, a_2)}{Var(W_i)} \\ &= a_i \frac{\sigma_{ii}}{\sigma_{W_i}^2}, \quad i = 1, 2, \end{aligned} \quad (27)$$

where $\sigma_{ii} = Cov(X_i, W_i)$, and $\sigma_{W_i}^2 = Var(W_i)$, $i = 1, 2$. Then

$$\begin{aligned} Var(\hat{g}) &= a_1^2 \frac{(\sigma_{11})^2}{\sigma_{W_1}^2} + a_2^2 \frac{(\sigma_{22})^2}{\sigma_{W_2}^2} \\ &= a_1^2 \sigma_{X_1}^2 r_{11}^2 + a_2^2 \sigma_{X_2}^2 r_{22}^2 \end{aligned} \quad (28)$$

where r_{ii}^2 is the (squared) coefficient of linkage disequilibrium between the i th pair of causal and marker genotypes. The covariance is

$$Cov(g, \hat{g}) = Var(\hat{g}),$$

and the squared correlation

$$\begin{aligned} \rho_{g, \hat{g}}^2 &= \frac{Var(\hat{g})}{Var(g)} \\ &= \frac{a_1^2 \sigma_{X_1}^2 r_{11}^2 + a_2^2 \sigma_{X_2}^2 r_{22}^2}{a_1^2 \sigma_{X_1}^2 + a_2^2 \sigma_{X_2}^2} \end{aligned} \quad (29)$$

The denominator is the total (causal) genetic variance of the trait. Therefore $\rho_{g, \hat{g}}^2 = 1$ when the marker genotype is the causal genotype. (Compare the above with (18), derived under a different genetic model). Similarly,

$$\begin{aligned} \rho_{y, \hat{g}}^2 &= \frac{Var(\hat{y})}{Var(y)} \\ &= \frac{a_1^2 \sigma_{X_1}^2 r_{11}^2 + a_2^2 \sigma_{X_2}^2 r_{22}^2}{a_1^2 \sigma_{X_1}^2 + a_2^2 \sigma_{X_2}^2 + \sigma_e^2} \frac{Var(g)}{Var(g)} \\ &= \rho_{g, \hat{g}}^2 h^2 \\ &= h_G^2, \end{aligned} \quad (30)$$

the genomic heritability.

NOTE

The derivation of (26) is as follows.

6.2 The genomic variance

Typically one computes genomic values for each individual. Given the above model, with m independent pairs of marker-causal locus genotypes, the genomic values are defined as

$$\begin{aligned} g &= \sum_{i=1}^m b_i W_i \\ &= \sum_{i=1}^m a_i \frac{\sigma_{ii}}{\sigma_{W_i}^2} W_i \end{aligned}$$

where $a_i r_{ii} \frac{\sigma_{X_i}}{\sigma_{W_i}} = b_i$, the regression of the causal locus genotypic value on marker genotype, given the substitution effect of the causal locus a_i . Then the genomic variance is

$$\begin{aligned} \text{Var}(g|b) &= \sum_{i=1}^m b_i^2 \text{Var}(W_i) \\ &= \sum_{i=1}^m a_i^2 \sigma_{X_i}^2 r_{ii}^2. \end{aligned} \quad (31)$$

When the markers are causal, $r_{ii}^2 = 1$ and $\text{Var}(g|b) = \sum_{i=1}^m a_i^2 \sigma_{X_i}^2$ is the additive genetic variance of the trait, assuming linkage equilibrium. This expression provides one explanation for the so-called *missing heritability*.

6.3 Cross-validation

The following problem is considered below. Assume data $Y_i^t = bX_i + e_i$, $i = 1, 2, \dots, N$ from which an estimate \hat{b} of b is obtained. These are referred to as the *training data* Y_i^t . The genomic value is now defined as

$$g_M = bX_i.$$

An independent (given X) data (the *test data*) labelled Y_i , $i = 1, 2, \dots$, are to be predicted using $\hat{Y}_i = \hat{g}_M = \hat{b}X_i$ and the first objective is to obtain an expression for the (marginal, with respect to \hat{b} and X) correlation between g_M and \hat{Y}_i , $\rho_{g_M, \hat{Y}}$. Later we consider the correlation between Y_i and \hat{Y}_i , $\rho_{Y, \hat{Y}}$ and show that (Goddard, 2009)

$$\rho_{Y, \hat{Y}} = \rho_{Y, g_M} \rho_{g_M, \hat{Y}} \quad (32)$$

provided that Y is conditionally independent of \hat{Y} , given g_M . The correlation between a datum from the test data and its genomic value is of the form

$$\rho_{Y, g_M} = \frac{\text{Cov}(Y, g_M)}{\sqrt{\text{Var}(Y) \text{Var}(g_M)}}$$

and

$$\text{Cov}(Y, g_M) = \text{Cov}(g, g_M) \quad (33)$$

where g is the true, unobserved genetic value. Here we are expressing the observed phenotype in the validating data as

$$Y_i = g_i + \epsilon_i$$

where g_i is the unobserved genetic value of the causal genotype of individual i . As shown above in (29), $Cov(g, g_M)$ is a function of the degree of linkage disequilibrium between the markers and unobserved causal factors and of the effect of the causal locus on the trait

Two modelling scenarios are considered. In the first one, SNP effects are treated as unknown, fixed parameters, and are estimated using least squares. In the second scenario, SNP effects are random variables and are predicted using BLUP. Before doing this we show the derivation of (32).

6.4 Factorising the correlation between Y and \hat{Y} *a la* Goddard

Assume that (Y, \hat{Y}, g_M) are multivariate normally distributed with covariance structure

$$Var(Y, \hat{Y}, g_M) = \begin{bmatrix} \sigma_Y^2 & \sigma_{Y, \hat{Y}} & \sigma_{Y, g_M} \\ \sigma_{Y, \hat{Y}} & \sigma_{\hat{Y}}^2 & \sigma_{\hat{Y}, g_M} \\ \sigma_{Y, g_M} & \sigma_{\hat{Y}, g_M} & \sigma_{g_M}^2 \end{bmatrix}.$$

Then

$$\begin{aligned} Var(Y, \hat{Y} | g_M) &= \begin{bmatrix} \sigma_Y^2 & \sigma_{Y, \hat{Y}} \\ \sigma_{Y, \hat{Y}} & \sigma_{\hat{Y}}^2 \end{bmatrix} - \begin{bmatrix} \sigma_{Y, g_M} \\ \sigma_{\hat{Y}, g_M} \end{bmatrix} \begin{bmatrix} \sigma_{Y, g_M} & \sigma_{\hat{Y}, g_M} \end{bmatrix} \frac{1}{\sigma_{g_M}^2} \\ &= \begin{bmatrix} \sigma_Y^2 & \sigma_{Y, \hat{Y}} \\ \sigma_{Y, \hat{Y}} & \sigma_{\hat{Y}}^2 \end{bmatrix} - \frac{1}{\sigma_{g_M}^2} \begin{bmatrix} (\sigma_{Y, g_M})^2 & \sigma_{Y, g_M} \sigma_{\hat{Y}, g_M} \\ \sigma_{Y, g_M} \sigma_{\hat{Y}, g_M} & (\sigma_{\hat{Y}, g_M})^2 \end{bmatrix}. \end{aligned}$$

The off-diagonal term is

$$Cov(Y, \hat{Y} | g_M) = \sigma_{Y, \hat{Y}} - \frac{1}{\sigma_{g_M}^2} \sigma_{Y, g_M} \sigma_{\hat{Y}, g_M}.$$

The diagonal terms are

$$\begin{aligned} Var(Y | g_M) &= \sigma_Y^2 - \frac{(\sigma_{Y, g_M})^2}{\sigma_{g_M}^2} = \sigma_Y^2 (1 - \rho_{Y, g_M}^2), \\ Var(\hat{Y} | g_M) &= \sigma_{\hat{Y}}^2 (1 - \rho_{\hat{Y}, g_M}^2). \end{aligned}$$

The conditional correlation is then

$$\begin{aligned} \rho_{Y, \hat{Y} | g_M} &= \frac{\sigma_{Y, \hat{Y}}}{\sigma_Y \sigma_{\hat{Y}} \sqrt{(1 - \rho_{Y, g_M}^2) (1 - \rho_{\hat{Y}, g_M}^2)}} - \frac{\frac{1}{\sigma_{g_M}^2} \sigma_{Y, g_M} \sigma_{\hat{Y}, g_M}}{\sigma_Y \sigma_{\hat{Y}} \sqrt{(1 - \rho_{Y, g_M}^2) (1 - \rho_{\hat{Y}, g_M}^2)}} \\ &= \frac{\rho_{Y, \hat{Y}}}{\sqrt{(1 - \rho_{Y, g_M}^2) (1 - \rho_{\hat{Y}, g_M}^2)}} - \frac{\rho_{Y, g_M} \rho_{\hat{Y}, g_M}}{\sqrt{(1 - \rho_{Y, g_M}^2) (1 - \rho_{\hat{Y}, g_M}^2)}}. \end{aligned} \quad (34)$$

If

$$p\left(Y, \hat{Y}|g_M\right) = p\left(Y|g_M\right) p\left(\hat{Y}|g_M\right)$$

then $\rho_{Y, \hat{Y}|g_M} = 0$ and (34) shows that

$$\rho_{Y, \hat{Y}} = \rho_{Y, g_M} \rho_{\hat{Y}, g_M}.$$

6.5 Marker effects as fixed parameters - least squares estimation

6.5.1 The case of 1 marker

To address this question, consider first the single marker model

$$Y_i = bX_i + e_i, \quad i = 1, 2, \dots, N, \quad (35)$$

with $\text{Var}(Y_i|X_i) = \sigma_e^2$ and where N is the number of individuals. In (35), $X_i \sim (0, 2p(1-p))$; thus, in contrast with previous notation, X_i represents now a centered (but not scaled) random variable that labels the marker genotype and b is treated as a known parameter. The genomic value is

$$g_M = bX_i \quad (36)$$

and

$$\text{Var}(g_M) = \sigma_{g_M}^2 = b^2 2p(1-p)$$

is the genomic variance.

The variance of Y_i given b , is

$$\text{Var}(Y_i|b) = \sigma_{g_M}^2 + \sigma_e^2.$$

Assume that b is estimated by least squares using test data Y^t , and that the predictor is (we use here \hat{Y} for the predictor, not to be confused with (23))

$$\hat{Y}_i = \hat{b}X_i, \quad (37)$$

where \hat{b} is the least squares estimator which has the distribution

$$\hat{b} \sim N\left(b, \frac{\sigma_e^2}{N2p(1-p)}\right) \quad (38)$$

and N is the size of the training data. The denominator arises by writing $\text{Var}(\hat{b}|X) = \frac{\sigma_e^2}{\sum_{i=1}^N (X_i - \bar{X})^2}$, and expressing $\sum_{i=1}^N (X_i - \bar{X})^2 = N\widehat{\text{Var}}(X_i) = N2p(1-p)$, assuming p is known.

To obtain an expression for $\rho_{g_M, \hat{Y}} = \frac{\text{Cov}(g_M, \hat{Y})}{\sqrt{\text{Var}(g_M)\text{Var}(\hat{Y})}}$ we proceed as follows. Firstly,

$$\begin{aligned} \text{Cov}(g_M, \hat{Y}|X) &= \text{Cov}(bX, \hat{b}X|X) \\ &= X^2 \text{Cov}(b, \hat{b}) = 0 \end{aligned}$$

because b is a parameter. The marginal covariance is

$$\begin{aligned}
Cov(g_M, \hat{Y}) &= E[Cov(g_M, \hat{Y}|X)] + Cov[E(g_M|X), E(\hat{Y}|X)] \\
&= 0 + Cov_X[bX, bX] \\
&= b^2 Var(X) = \sigma_{g_M}^2.
\end{aligned} \tag{39}$$

NOTE: $E(\hat{Y}|X) = E(\hat{b}X|X) = XE(\hat{b}) = bX$.

The same result is obtained (of course) if instead one conditions on \hat{b} . Thus

$$\begin{aligned}
Cov(g_M, \hat{Y}|\hat{b}) &= Cov(bX, \hat{b}X|\hat{b}) \\
&= \hat{b}\hat{b}Var(X)
\end{aligned}$$

and

$$\begin{aligned}
Cov(g_M, \hat{Y}) &= E[Cov(g_M, \hat{Y}|\hat{b})] + Cov[E(g_M|\hat{b}), E(\hat{Y}|\hat{b})] \\
&= E[\hat{b}\hat{b}Var(X)] \\
&= b^2 Var(X) = \sigma_{g_M}^2.
\end{aligned}$$

and the second term on the right hand side of the first line is zero because $E(\hat{Y}|\hat{b}) = E(\hat{b}X|\hat{b}) = \hat{b}E(X) = 0$.

Also,

$$Var(\hat{Y}|\hat{b}) = \hat{b}^2 Var(X).$$

Notice that if b is known without error, $Var(\hat{Y}|\hat{b}) = Var(\hat{Y}|b) = b^2 Var(X) = \sigma_{g_M}^2$ and

$$\begin{aligned}
\rho_{g_M \hat{Y}}^2 &= \frac{(b^2 Var(X))^2}{Var(X) b^2 Var(g_M)} \\
&= \frac{b^2 Var(X)}{Var(g_M)} = 1
\end{aligned}$$

When b is estimated from the data, the unconditional variance is

$$\begin{aligned}
Var(\hat{Y}) &= E[Var(\hat{Y}|\hat{b})] + Var[E(\hat{Y}|\hat{b})] \\
&= E[Var(\hat{Y}|\hat{b})] \\
&= Var(X) E(\hat{b}^2) \\
&= Var(X) (Var(\hat{b}) + b^2) \\
&= \sigma_{g_M}^2 + \frac{\sigma_e^2}{N}.
\end{aligned}$$

Then,

$$\begin{aligned}
\rho_{g_M, \hat{Y}}^2 &= \frac{(\sigma_{g_M}^2)^2}{\sigma_{g_M}^2 \left(\sigma_{g_M}^2 + \frac{\sigma_e^2}{N} \right)} \\
&= \frac{\sigma_{g_M}^2}{\sigma_{g_M}^2 + \frac{\sigma_e^2}{N}} \\
&= \frac{g^2}{\frac{(1-g^2)}{N} + g^2},
\end{aligned} \tag{40}$$

where $1 - g^2 = \sigma_e^2 / \text{Var}(Y)$. When N is large, $\rho_{g_M, \hat{Y}}^2 \approx 1$.

6.5.2 The case of m_G markers

The extended model assumes that m_G marker genotypes are available. The linear structure for a datum is

$$Y_{ij} = \sum_{j=1}^{m_G} b_j X_{ij} + e_i; \quad j = 1, 2, \dots, m_G; \quad i = 1, 2, \dots, N,$$

where as before, the b 's are treated as fixed but unknown parameters. The genomic value is

$$g_{M_i} = \sum_{j=1}^{m_G} b_j X_{ij}.$$

The variance of Y_{ij} given all the b 's is

$$\begin{aligned}
\text{Var}(Y_{ij}|b) &= \sum_{j=1}^{m_G} b_j^2 \text{Var}(X_{ij}) + \sum_{i < j} 2b_i b_j \text{Cov}(X_i, X_j) + \sigma_e^2 \\
&= \sigma_{g_M}^2 + \sigma_e^2 = \sigma^2.
\end{aligned}$$

At this point, we note that here with the b 's treated as fixed effects,

$$\sigma_{g_M}^2 = \sum_{j=1}^{m_G} b_j^2 \text{Var}(X_{ij}) + \sum_{i < j} 2b_i b_j \text{Cov}(X_i, X_j) \tag{41}$$

the genomic variance has a different interpretation than in (13). In the latter, the b 's are random, and the variance is marginalised with respect to b and W . Here it is only marginalised over the X 's.

The predictor is

$$\hat{Y}_{ij} = \sum_{j=1}^{m_G} \hat{b}_j X_{ij}.$$

Now

$$\begin{aligned} Cov \left(g_{M_i}, \hat{Y}_i | \hat{b} \right) &= Cov \left(\sum_{j=1}^{m_G} b_j X_{ij}, \sum_{j=1}^{m_G} \hat{b}_j X_{ij} | \hat{b} \right) \\ &= \sum_{j=1}^{m_G} \hat{b}_j b_j Var (X_j) + \sum_{k,j} b_k \hat{b}_j Cov (X_k, X_j) \end{aligned}$$

and unconditionally with respect to \hat{b} ,

$$\begin{aligned} Cov \left(g_{M_i}, \hat{Y}_i \right) &= \sum_{j=1}^{m_G} b_j^2 Var (X_j) + \sum_{k < j} 2b_k b_j Cov (X_k, X_j) \\ &= \sigma_{g_M}^2, \end{aligned} \tag{42}$$

the total variance due to the markers. Similarly,

$$\begin{aligned} Var \left(\hat{Y} | \hat{b} \right) &= Var \left(\sum_{j=1}^{m_G} \hat{b}_j X_j \right) \\ &= \sum_{j=1}^{m_G} \hat{b}_j^2 Var (X_j) + \sum_{i < j} 2\hat{b}_i \hat{b}_j Cov (X_i, X_j), \end{aligned}$$

and the marginal variance is

$$\begin{aligned} Var \left(\hat{Y} \right) &= E \left[Var \left(\sum_j \hat{b}_j X_j \right) \right] + Var \left[E \left(\hat{Y} | \hat{b} \right) \right] \\ &= \sum_{j=1}^{m_G} Var (X_j) E \left(\hat{b}_j^2 \right) + \sum_{i < j} 2Cov (X_i, X_j) E \left(\hat{b}_i \hat{b}_j \right) \\ &= \sum_{j=1}^{m_G} Var (X_j) \left(Var \left(\hat{b}_j \right) + b_j^2 \right) + \sum_{i < j} 2Cov (X_i, X_j) \left(Cov \left(\hat{b}_i, \hat{b}_j \right) + b_i b_j \right) \\ &= \sigma_{g_M}^2 + \sum_j Var (X_j) Var \left(\hat{b}_j \right) + \sum_{i < j} 2Cov (X_i, X_j) Cov \left(\hat{b}_i, \hat{b}_j \right). \end{aligned} \tag{43}$$

The second term in the first line is zero because $E \left(\hat{Y} | \hat{b} \right) = 0$. Also in the 4th line

$$\begin{aligned} \sigma_{g_M}^2 &= \sum_{j=1}^{m_G} Var (X_j) b_j^2 + \sum_{i < j} 2Cov (X_i, X_j) b_i b_j \\ &= \sum_{j=1}^{m_G} 2p_j (1 - p_j) b_j^2 + \sum_{i < j} 2D_{ij} b_i b_j \end{aligned}$$

which differs from the definition given before in (13) as mentioned above.

As an approximation, we assume

$$Var \left(\hat{b}_j \right) = \frac{\sigma_e^2}{2Np_j(1 - p_j)}$$

and then

$$\begin{aligned} \sum_{j=1}^{m_G} Var (X_j) Var \left(\hat{b}_j \right) &= \sum_{j=1}^{m_G} 2p_j (1 - p_j) \frac{\sigma_e^2}{2Np_j(1 - p_j)} \\ &= \frac{m_G \sigma_e^2}{N}. \end{aligned}$$

We shall also ignore the third term in (43), which is equivalent to assuming that the m_G markers represent an *effective number of markers* whose correlation is zero. Therefore, ignoring the third term $\sum_{i < j} 2Cov(X_i, X_j) Cov(\hat{b}_i, \hat{b}_j)$

$$\begin{aligned} Var(\hat{Y}) &\approx \sigma_{g_M}^2 + \sum_j Var(X_j) Var(\hat{b}_j) \\ &= \sigma_{g_M}^2 + \frac{m_G \sigma_e^2}{N}, \end{aligned}$$

and, approximately

$$\begin{aligned} \rho_{g_M, \hat{Y}}^2 &\approx \frac{[Cov(g_{M_i}, \hat{Y}_i)]^2}{Var(g_{M_i}) Var(Var \hat{Y}_i)} \\ &= \frac{(\sigma_{g_M}^2)^2}{\sigma_{g_M}^2 \left(\sigma_{g_M}^2 + \frac{m_G \sigma_e^2}{N} \right)} \\ &= \frac{\sigma_{g_M}^2}{\frac{m_G \sigma_e^2}{N} + \sigma_{g_M}^2} \\ &= \frac{g^2}{\frac{m_G(1-g^2)}{N} + g^2}. \end{aligned} \tag{44}$$

This expression reduces to (40) when $m_G = 1$. The point is that depending on m_G/N , the squared correlation between observed and predicted values may be significantly smaller than the proportion of variance due to the markers, g^2 .

Incidentally, (44) produces

$$m_G = \frac{Ng^2(1 - \rho_{g_M, \hat{Y}}^2)}{\rho_{g_M, \hat{Y}}^2(1 - g^2)}. \tag{45}$$

6.6 Marker effects as random variables - BLUP

When W and b are treated as random, the genomic variance conditional on the observed genotypes is

$$\begin{aligned} Var(w'_i b | w_i) &= Var(g_{Mi} | w_i) \\ &= w'_i w_i \sigma_b^2 \\ &= \sum_{i=1}^{m_G} w_i^2 \sigma_b^2, \end{aligned}$$

where w'_i is the i th row (individual) of $W = \{w_{ij}\}$, $j = 1, 2, \dots, m_G$ and m_G is the observed number of marker genotypes. Unconditionally,

$$Var(g_M) = \sigma_{g_M}^2 = m_G \sigma_b^2$$

as indicated in (14) and different from (41). Given $\sigma_{g_M}^2$, the larger m_G the smaller the variance of the SNP effect σ_b^2 . The genomic heritability or the proportion of the total variance explained by the SNP genotypes is

$$g^2 = \frac{\sigma_{g_M}^2}{\sigma^2},$$

where $\sigma^2 = \sigma_{g_M}^2 + \sigma_e^2$ is the marginal variance of Y .

As above an expression for the squared correlation between observed data g_M and predicted data \hat{Y} , $\rho_{g_M, \hat{Y}}^2$ is derived.

The model is

$$y = Wb + e, \quad (46)$$

y is $(N \times 1)$, W is $(N \times m_G)$ and b is $(m_G \times 1)$. Further,

$$\begin{aligned} y|W, b &\sim N(Wb, I\sigma_e^2), \\ b &\sim N(0, I\sigma_b^2) \end{aligned}$$

and the variance components σ_e^2 , σ_b^2 are assumed “known”. Although inferences are conditional on $W = \{w_{ij}\}$, since the elements of W are centered and scaled as in (11),

$$\begin{aligned} w_{ij} &\sim (0, 1), \\ \text{Cov}(w_{ij}, w_{ik}) &= E(w_{ij}, w_{ik}) \\ &= r_{jk} = D_{jk} \end{aligned}$$

where D_{jk} is the linkage disequilibrium parameter between SNP j and k .

BLUP of b is

$$\hat{b} = E(b|y) = (W'W + Ik)^{-1} W'y \quad (47)$$

with $k = \sigma_e^2/\sigma_b^2$. The prediction error variance is

$$\begin{aligned} \text{Var}(\hat{b} - b) &= \text{Var}(b|y) \quad (\text{not a function of } y) \\ &= E[\text{Var}(b|y)] \\ &= \sigma_e^2 (W'W + Ik)^{-1}. \end{aligned} \quad (48)$$

In general,

$$\begin{aligned} \text{Var}(b) &= E[\text{Var}(b|y)] + \text{Var}[E(b|y)] \\ &= \text{Var}(\hat{b} - b) + \text{Var}(\hat{b}). \end{aligned} \quad (49)$$

As shown below, $\text{Var}(\hat{b}_i - b_i) \approx (N + k)^{-1} \sigma_e^2$ and

$$\text{Var}(\hat{b}_i) \approx \sigma_b^2 - \frac{\sigma_e^2}{N + k}. \quad (50)$$

To obtain an expression for $\rho_{g_M, \hat{Y}}^2$ we need $Var(g_M) = \sigma_{g_M}$, $Var(\hat{Y})$ and $Cov(g_M, \hat{Y})$. The marginal (with respect to b and W) variance of the data is

$$\begin{aligned} Var(Y) &= m_G \sigma_b^2 + \sigma_e^2 \\ &= \sigma_{g_M}^2 + \sigma_e^2 = \sigma^2. \end{aligned}$$

The predicted value of Y is

$$\hat{Y}_i = \sum_{j=1}^{m_G} w_{ij} \hat{b}_j.$$

The conditional expectation is

$$E(\hat{Y}_i | \hat{b}) = \sum_{j=1}^{m_G} \hat{b}_j E(w_{ij}) = 0$$

and the conditional variance

$$Var(\hat{Y}_i | \hat{b}) = \sum_{j=1}^{m_G} \hat{b}_j^2 + \sum_{k < l} 2 \hat{b}_k \hat{b}_l D_{kl}.$$

Unconditionally,

$$\begin{aligned} Var(\hat{Y}_i) &= E[Var(\hat{Y}_i | \hat{b})] \\ &= \sum_{j=1}^{m_G} Var(\hat{b}_j) + 2 \sum_{k < l} r_{kl} Cov(\hat{b}_k, \hat{b}_l) \\ &= \sum_{j=1}^{m_G} Var(\hat{b}_j) + 2 \sum_{k < l} \\ &\approx \sum_{j=1}^{m_G} \left[\sigma_b^2 - \frac{\sigma_e^2}{N + k} \right] \quad \text{ignoring the second term above} \\ &= m_G \sigma_b^2 - \frac{m_G \sigma_e^2}{N + k} \\ &= \sigma_{g_M}^2 - \frac{m_G (1 - g^2) \sigma^2}{N + \frac{(1 - g^2) \sigma^2}{\sigma_b^2}} \\ &= \sigma_{g_M}^2 - \frac{(1 - g^2) \sigma^2}{\frac{N}{m_G} + \frac{(1 - g^2) \sigma^2}{m_G \sigma_b^2}} \\ &= \sigma_{g_M}^2 - \frac{(1 - g^2) \sigma^2}{\frac{N}{m_G} + \frac{(1 - g^2)}{g^2}} \\ &= \sigma_{g_M}^2 \frac{g^2}{g^2 + (1 - g^2) \frac{m_G}{N}}. \end{aligned} \tag{51}$$

The squared correlation is

$$\begin{aligned}
\rho_{g_M, \hat{Y}}^2 &= \frac{Cov(g_M, \hat{Y})^2}{Var(g_M) Var(\hat{Y})} \\
&= \frac{Var(\hat{Y})}{Var(g_M)} \\
&= \frac{g^2}{g^2 + (1 - g^2) \frac{m_G}{N}}
\end{aligned} \tag{52}$$

because $Cov(g_M, \hat{Y}) = Var(\hat{Y})$. Notice that (52) is equal to (44), which seems a strange result despite the different definitions of $\sigma_{g_M}^2$.

NOTE 1

$$\begin{aligned}
Cov(\hat{Y}_i, g'_{Mi}|w_i) &= Cov(w'_i \hat{b}, \hat{b}' w_i | w_i) \\
&= Cov\left[w'_i \hat{b}, (\hat{b} - (\hat{b} - b))' w_i\right] \\
&= w'_i \left[Var(\hat{b}) - Cov(\hat{b}, (\hat{b} - b)') \right] w_i \\
&= w'_i Var(\hat{b}) w_i \\
&= Var(\hat{Y}_i | w_i).
\end{aligned}$$

NOTE 2 - For two records and three SNPs, (46) is

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} + e$$

and

$$W'W + Ik = \begin{bmatrix} w_{11}w_{11} + w_{21}w_{21} + k & w_{11}w_{12} + w_{21}w_{22} & w_{11}w_{13} + w_{21}w_{23} \\ w_{11}w_{12} + w_{21}w_{22} & w_{12}w_{12} + w_{22}w_{22} + k & w_{12}w_{13} + w_{22}w_{23} \\ w_{11}w_{13} + w_{21}w_{23} & w_{12}w_{13} + w_{22}w_{23} & w_{13}w_{13} + w_{22}w_{23} + k \end{bmatrix}. \tag{53}$$

The j th diagonal element has the form

$$\sum_{i=1}^N w_{ij}^2 + k \tag{54}$$

where the subscript i refers to individual and j to the SNP. Since $E(w_{ij}^2) = 1$ and $E(w_{ij}) = 0$, $Var(w_{ij}) = 1$ and we write (54) as $NVar(w_{ij}) + k = N + k$. The off-diagonal terms are function of disequilibria (covariances) involving different SNP genotypes. Parameter k is typically a very large number, because the variance per SNP σ_b^2 is very small relative to σ_e^2 . The diagonal element of the inverse of (53) will be approximated by $(N + k)^{-1}$ and therefore,

$$Var(\hat{b} - b) \approx I(N + k)^{-1} \sigma_e^2. \quad (55)$$

6.7 Influence of degree of relationship on predictive ability

We now consider how the degree of relationship between the training and validating data affect the predictive ability of the genomic model. To formalise this, we consider

$$Cov(\hat{Y}, g'_M | W).$$

Firstly, BLUP of b is computed using the training data $Y_t = W_t b + e$, where as before

$$\begin{aligned} b &\sim N(0, I\sigma_b^2) \\ e &\sim N(0, I\sigma_e^2). \end{aligned}$$

Then

$$\begin{aligned} \hat{b} &= C_t^{-1} W_t' Y_t \\ &= C_t^{-1} W_t' W_t b + C_t^{-1} W_t' e, \end{aligned}$$

where

$$C_t = W_t' W_t + Ik, \quad k = \sigma_e^2 / \sigma_b^2.$$

In these expressions, W_t of order $N \times m$, has elements $W_{t,ij}$, $i = 1, 2, \dots, m$ labelling markers and $j = 1, 2, \dots, N$ labelling individuals in the training data. The predicted phenotype in the validating data is

$$\begin{aligned} \hat{Y} &= W_v \hat{b} \\ &= W_v C_t^{-1} W_t' W_t b + W_v C_t^{-1} W_t' e. \end{aligned}$$

The unobserved genomic effect in the validating data is

$$g_M = W_v b$$

and

$$\begin{aligned} Cov(\hat{Y}, g'_M | W) &= Cov(W_v C_t^{-1} W_t' W_t b + W_v C_t^{-1} W_t' e, b' W_v' | W) \\ &= W_v C_t^{-1} W_t' W_t W_v' \sigma_b^2 \end{aligned} \quad (56)$$

because $Cov(b, e') = 0$. The covariance is a function of the degree of genomic relationship between the training and the validating datasets, via the terms $W_t W_v'$. For the i th phenotypic value, $Cov(\hat{Y}_i, g_{M_i} | W) = W_{v,i} C_t^{-1} W_t' W_t W_{v,i}' \sigma_b^2$, where $W_{v,i}$ is the i th row of matrix W_v . In this case

$$W_t W_{v,i}' = \begin{bmatrix} W_{t,11} W_{v,1i} + W_{t,21} W_{v,2i} + \cdots + W_{t,m1} W_{v,mi} \\ W_{t,12} W_{v,1i} + W_{t,22} W_{v,2i} + \cdots + W_{t,m2} W_{v,mi} \\ \vdots \\ W_{t,1N} W_{v,1i} + W_{t,2N} W_{v,2i} + \cdots + W_{t,mN} W_{v,mi} \end{bmatrix}$$

an $N \times 1$ column vector where the j th row is the sum over the m markers of terms of the form $W_{t,kj} W_{v,ki}$. This term corresponds to individual j in the training data and individual i in the validating data. The elements in this vector describe the degree of marker captured relationship between individual i in the validating data, and all the N individuals in the training data. Therefore, dividing row j say of $W_t W_{v,i}'$ by the number of markers m , yields an estimate of the average (over markers) relationship between individual j in the training data and i in the validating data.

References

- Goddard, M. E. (2009). Genomic selection: prediction of accuracy and maximisation of long term response. *Genetica* 136, 245–252.
- Powell, J. E., P. M. Visscher, and M. E. Goddard (2010). Reconciling the analysis of IBD and IBS in complex trait studies. *Nature Reviews, Genetics* 11, 800–805.
- Wright, S. (1922). Coefficients of inbreeding and relationship. *American Naturalist* 56, 330–338.