

File S1

Detailed Theory

Covariance Between Individuals in a Founder Population due to Arbitrary Epistasis

We model each diploid individual t 's phenotype as $Y_t = G_t + \epsilon_t$, where G_t is the effect of the genotype of individual t , and ϵ_t is the residual error arising from noise and the environment, with $\mathbb{E}[\epsilon_t] = 0 \forall t$. We do not model genotype-by-environment interactions here, so $\text{Cov}(G_t, \epsilon_t) = 0 \forall t$.

We model G_t to accommodate any possible interaction effects between loci and to give an orthogonal partitioning of the genetic variance. The model allows every subset of the causal loci L to interact in an arbitrary way to produce an effect, X_{tL} , on the phenotype of individual t . Within each subset L , each possible sequence of alleles across the loci, $s \in S_L$, can have a different effect on the phenotype, β_{Ls} . Therefore, the model is general for any genotype-phenotype map without dominance effects.

For a set of causal loci N , which we assume are in linkage equilibrium,

$$G_t = \sum_{L \subseteq N} X_{tL}; \quad X_{tL} = \sum_{s \in S_L} \beta_{Ls} \prod_{l \in L} x_{tl_s[l]}; \quad (1)$$

where X_\emptyset is the phenotypic mean; $S_L = \{A, T, G, C\}^{|L|}$ is the set of possible sequences of alleles across the loci in L ; $s[l]$ is the allelic state of the locus l in the sequence s ; and

$$x_{tlA} = g_{tlA}^p + g_{tlA}^m - 2f_{lA}, \quad (2)$$

where g_{tlA}^p is an indicator variable for the presence of allele A at locus l on the paternally inherited chromosome of individual t , g_{tlA}^m indicates the equivalent on the maternally inherited chromosome, and f_{lA} is the frequency of the A allele at locus l .

Because $\mathbb{E}[x_{tlk}] = 0 \forall t, l, k$ and the loci are in linkage equilibrium, $\text{Cov}(X_{tL}, X_{tL'}) = 0$ for $L \neq L'$. Proof: without loss of generality, $\exists d \in L \setminus L'$, implying

$$\text{Cov}(X_{tL}, X_{tL'}) = \sum_{s \in S_L} \sum_{s' \in S_{L'}} \beta_{Ls} \beta_{L's'} \quad (3)$$

$$\mathbb{E} \left[\prod_{l \in L \setminus \{d\}} x_{tls[l]} \prod_{l' \in L'} x_{tl's'[l']} \right] \mathbb{E}[x_{tds[d]}] \\ = 0, \text{ as } \mathbb{E}[x_{tds[d]}] = 0. \quad (4)$$

The covariance between arbitrary relatives t and u with kinship coefficient $K_{t,u}$ is therefore

$$\text{Cov}(G_t, G_u) = \sum_{L \subseteq N} \text{Cov}(X_{tL}, X_{uL}); \quad (5)$$

$$= \sum_{s \in S_L} \sum_{s' \in S_{L'}} \beta_{Ls} \beta_{L's'} \prod_{l \in L} \mathbb{E}[x_{tls[l]} x_{uls'[l]}], \quad (6)$$

due to linkage equilibrium between the loci.

$$\mathbb{E}[x_{tls[l]} x_{uls'[l]}] = \sum_{i=m,p} \sum_{j=m,p} \text{Cov}(g_{tls[l]}^i, g_{uls'[l]}^j). \quad (7)$$

$$\text{Cov}(g_{tls[l]}^i, g_{uls'[l]}^j) = \mathbb{E}[g_{tls[l]}^i g_{uls'[l]}^j] - f_{ls[l]} f_{ls'[l]}; \quad (8)$$

$$= \left(\frac{K_{t,u}^{i,j} - K_0}{1 - K_0} \right) f_{ls[l]} (1 - f_{ls'[l]}), \quad (9)$$

when $s[l] = s'[l]$, from the genotypic covariance in a founder population, which is derived in the main text - see (7) and (11).

When $s[l] \neq s'[l]$, assuming no mutation, $\mathbb{E}[g_{tls[l]}^i g_{uls'[l]}^j]$ is only non-zero when the haplotypes are not IBD, because the alleles are different. Given that the haplotypes are not IBD, $\mathbb{E}[g_{tls[l]}^i g_{uls'[l]}^j]$ is the probability that t inherits allele $s[l]$ from one of the A ancestral haplotypes and u inherits allele $s'[l]$ from one of the other $A - 1$ ancestral haplotypes. Therefore, if $c_{ls'[l]}$ is the number of ancestral haplotypes carrying the $s'[l]$ allele,

$$\mathbb{E}[g_{tls[l]}^i g_{uls'[l]}^j] = (1 - K_{t,u}^{i,j}) f_{ls[l]} \frac{c_{ls'[l]}}{A - 1} = \frac{(1 - K_{t,u}^{i,j})}{1 - K_0} f_{ls[l]} f_{ls'[l]}. \quad (10)$$

Therefore,

$$\text{Cov}(g_{t|s[l]}^i, g_{u|s'[l]}^j) = -f_{l|s[l]} f_{l|s'[l]} \left(\frac{K_{t,u}^{i,j} - K_0}{1 - K_0} \right). \quad (11)$$

Therefore,

$$\mathbb{E}[x_{t|s[l]} x_{u|s'[l]}] = 2 \left(\frac{K_{t,u} - K_0}{1 - K_0} \right) \xi_{l:s[l],s'[l]}, \quad (12)$$

where $\xi_{l:s[l],s'[l]} = -2f_{l|s[l]}f_{l|s'[l]}$ is the covariance between $x_{t|s[l]}$ and $x_{t|s'[l]}$ for the distinct alleles $s[l]$ and $s'[l]$ in an outbred population with the same allele frequencies. $\xi_{l:s[l],s[l]} = 2f_{l|s[l]}(1 - f_{l|s[l]})$ is the variance of $x_{t|s[l]}$ in an outbred population. The outbred allele count variances and covariances are equivalent to those for a multinomial distribution with two trials and with event probabilities equal to the allele frequencies at the locus.

Therefore,

$$\text{Cov}(X_{tL}, X_{tL'}) = 2^{|L|} \left(\frac{K_{t,u} - K_0}{1 - K_0} \right)^{|L|} \xi_L, \quad (13)$$

where ξ_L is the variance of X_L in an outbred population, and is equal to

$$\sum_{s \in S_L} \sum_{s' \in S_L} \beta_{Ls} \beta_{Ls'} \prod_{l \in L} \xi_{l:s[l],s'[l]}. \quad (14)$$

Therefore,

$$\text{Cov}(G_t, G_u) = \sum_{\tau=1}^{|N|} 2^\tau \left(\frac{K_{t,u} - K_0}{1 - K_0} \right)^\tau v_\tau, \quad (15)$$

where v_τ is the variance from genetic interactions between τ loci in an outbred population, and is the sum of ξ_L over all subsets L of size τ .

Using the fact that $K_{tt} = (1 + F_t)/2$, where F_t is the inbreeding coefficient of individual t , and setting $t = u$ gives

$$\text{Var}(G_t) = \sum_{\tau=1}^{|N|} \left(1 + \frac{F_t - K_0}{1 - K_0} \right)^\tau v_\tau. \quad (16)$$

The population variance is derived by the law of total variance,

$$\text{Var}(G) = \mathbb{E}_t[\text{Var}(G_t)] + \text{Var}_t(\mathbb{E}[G_t]). \quad (17)$$

Because there is no dominance, the phenotypic mean does not change with inbreeding. In a random-mating

population, the mean inbreeding coefficient is equal to the mean kinship coefficient: $\mathbb{E}_t(F_t) = K_0$. Therefore,

$$\begin{aligned} \text{Var}(G) = \mathbb{E}_t[\text{Var}(G_t)] &= v_1 + \left(1 + \frac{\text{Var}(F_t)}{(1 - K_0)^2}\right) v_2 + \\ &\sum_{\tau=3}^{|N|} \left(1 + \binom{\tau}{2} \frac{\text{Var}(F_t)}{(1 - K_0)^2} + \sum_{i=3}^{\tau} \binom{\tau}{i} \frac{\mathbb{E}[(F_t - K_0)^i]}{(1 - K_0)^i}\right) v_{\tau}. \end{aligned} \quad (18)$$

Dominance Variance in a Founder Population

We consider a phenotype whose genetic contribution is determined by two bi-allelic loci in linkage equilibrium that have non-zero dominance deviations as well as an interaction between their additive effects. The phenotype of an individual s , Y_s , is the sum of the additive contributions of the loci, a_s , the interaction between the additive effects, $(a \times a)_s$, and the sum of the dominance deviations of the loci d_s , giving

$$Y_s = a_s + (a \times a)_s + d_s + \epsilon_s; \text{ where } d_s = \delta_1 \gamma_{s1m} \gamma_{s1p} + \delta_2 \gamma_{s2m} \gamma_{s2p}, \quad (19)$$

and $\gamma_{sim} = g_{si}^m - f_i$ is the mean normalised indicator variable for the presence of the minor allele at locus i , with frequency f_i , on the maternal chromosome, and $\gamma_{sip} = g_{si}^p - f_i$ is the corresponding variable for the paternal chromosome.

The additive and additive-by-additive components are orthogonal, as shown in the main text. The additive-by-additive and the dominance components are orthogonal, because the dominance deviation of each locus is uncorrelated with the additive effect of the other locus. Inbreeding, however, induces a correlation between the additive effect and dominance deviation at a locus, implying that

$$\text{Var}(Y_s) = \text{Var}(a_s) + \text{Var}((a \times a)_s) + \text{Var}(d_s) + \text{Cov}(a_s, d_s). \quad (20)$$

The additive and additive-by-additive variance components are as derived in the main text and supplement.

$\text{Var}(d_s)$ is derived by applying the law of total variance to $d_{si} = \delta_i \gamma_{sim} \gamma_{sip}$, the contribution of locus $i \in \{1, 2\}$ to d_s .

$$\text{Var}(d_{si}) = \mathbb{E}_{I_{si}}[\text{Var}(d_{si}|I_{si})] + \text{Var}_{I_{si}}(\mathbb{E}[d_{si}|I_{si}]), \quad (21)$$

where I_{si} is the indicator variable for whether individual s is inbred or not at locus i . The conditional expectation

of d_{si} depends on

$$\mathbb{E}[\gamma_{sim}\gamma_{sip}|I_{si}] = I_{si}(f_i(1-f_i) - \mathbb{E}[\gamma_{sim}\gamma_{sip}|I_{si}=0]) + \mathbb{E}[\gamma_{sim}\gamma_{sip}|I_{si}=0]. \quad (22)$$

Using the expression for the genotypic covariance in a founder population derived in the main text,

$$\mathbb{E}[\gamma_{sim}\gamma_{sip}|I_{si}=0] = \frac{-K_0}{1-K_0} f_i(1-f_i), \quad (23)$$

where K_0 is the mean inbreeding (and kinship) coefficient. Therefore,

$$f_i(1-f_i) - \mathbb{E}[\gamma_{sim}\gamma_{sip}|I_{si}=0] = \frac{f_i(1-f_i)}{1-K_0}, \quad (24)$$

and therefore

$$\text{Var}_{I_{si}}(\mathbb{E}[d_{si}|I_{si}]) = \frac{F_s(1-F_s)}{(1-K_0)^2} \mu_{hi}^2, \quad (25)$$

where $\mu_{hi} = \delta_i f_i(1-f_i)$ is the inbreeding depression at locus i .

We now calculate $\mathbb{E}_{I_{si}}[\text{Var}(d_{si}|I_{si})]$. First, in the inbred case,

$$\text{Var}(d_{si}|I_{si}=1) = \delta_i^2 f_i(1-f_i)(1-2f_i)^2 = v_{hi}, \quad (26)$$

which is the dominance variance at locus i in the homozygous population. When there is no inbreeding at locus i ,

$$\text{Var}(d_{si}|I_{si}=0) = \delta_i^2 (\mathbb{E}[\gamma_{sim}^2 \gamma_{sip}^2 | I_{si}=0] - \mathbb{E}[\gamma_{sim}\gamma_{sip}|I_{si}=0]^2) \quad (27)$$

By expanding the squares in $\gamma_{sim}^2 \gamma_{sip}^2$ and using the result for the genotypic covariance when there is no IBD sharing, it can be shown that

$$\mathbb{E}[\gamma_{sim}^2 \gamma_{sip}^2 | I_{si}=0] = f_i^2(1-f_i)^2 - \frac{K_0}{1-K_0} (1-2f_i)^2 f_i(1-f_i). \quad (28)$$

Using the result for genotypic covariance when there is no IBD again gives

$$\mathbb{E}[\gamma_{sim}\gamma_{sip}|I_{si}=0]^2 = \frac{K_0^2}{(1-K_0)^2} f_i^2(1-f_i)^2. \quad (29)$$

Therefore,

$$\text{Var}(d_{si}|I_{si} = 0) = \mu_{hi}^2 \left(1 - \frac{K_0^2}{(1 - K_0)^2} - \frac{K_0}{1 - K_0} \frac{(1 - 2f_i)^2}{f_i(1 - f_i)} \right) = v_{\delta_i}, \quad (30)$$

which we have defined to be v_{δ_i} .

Combining the results gives

$$\text{Var}(d_{si}) = (1 - F_s)v_{\delta_i} + F_s v_{hi} + \frac{F_s(1 - F_s)}{(1 - K_0)^2} \mu_{hi}^2, \quad (31)$$

where F_s is the inbreeding coefficient of individual s . Summing across the loci gives

$$\text{Var}(d_s) = (1 - F_s)v_{\delta} + F_s v_h + \frac{F_s(1 - F_s)}{(1 - K_0)^2} SS_{\mu_h}, \quad (32)$$

where $v_{\delta} = v_{\delta 1} + v_{\delta 2}$, $v_h = v_{h1} + v_{h2}$, and $SS_{\mu_h} = \mu_{h1}^2 + \mu_{h2}^2$ is the sum of the squared inbreeding depressions at the loci.

In a founder population, the maternal and paternal alleles are not independent, which implies that

$$\text{Cov}(\gamma_{s1m}, \gamma_{s1m}\gamma_{s1p}) = \mathbb{E}[\gamma_{s1m}^2 \gamma_{s1p}] \neq 0. \quad (33)$$

This implies that there is covariance between an individual's additive effect and dominance deviation, depending on their inbreeding coefficient. The above expectation can be evaluated by conditioning on whether or not individual s is inbred or not, giving

$$\text{Cov}(\gamma_{s1m}, \gamma_{s1m}\gamma_{s1p}) = f_1(1 - f_1)(1 - 2f_1) \frac{F_s - K_0}{1 - K_0}. \quad (34)$$

Summing the contributions of the four possible covariances within the locus and summing across loci gives

$$\text{Cov}(a_s, d_s) = 4 \frac{F_s - K_0}{1 - K_0} C_{a,d}, \quad (35)$$

where $C_{a,d} = \sum_{i=1}^2 \beta_i \delta_i f_i (1 - f_i) (1 - 2f_i)$ parameterises the strength of the covariance between additive and dominance effects.

Combining these results with those from the main text:

$$\text{Var}(Y_s) = \sum_{\tau=1}^2 \left(1 + \frac{F_t - K_0}{1 - K_0}\right)^\tau v_\tau + (1 - F_s)v_\delta + \quad (36)$$

$$4 \frac{F_s - K_0}{1 - K_0} C_{a,d} + F_s v_h + \frac{F_s(1 - F_s)}{(1 - K_0)^2} SS_{\mu_h} + \sigma_\epsilon^2; \quad (37)$$

where v_1 and v_2 are as defined in the main text. v_δ is the covariance between two individuals' dominance deviations conditional on both alleles of one individual being IBD to distinct alleles of the other individual, implying that neither individual is inbred. v_δ differs slightly from the dominance variance in an infinite, outbred population, where $v_{\delta 1} = \mu_{h1}^2$. Similarly, $v_\delta \approx SS_{\mu_h}$. These difference will be very small apart from for populations descending from a very small number of founders (large K_0), such as in certain cross designs.

The variance in the population is, by the Law of Total Variance,

$$\text{Var}(Y) = \mathbb{E}_s[\text{Var}(Y_s)] + \text{Var}_s(\mathbb{E}[Y_s]). \quad (38)$$

Because the mean inbreeding coefficient is K_0 in an outbred population,

$$\mathbb{E}_s[\text{Var}(Y_s)] = v_1 + \left(1 + \frac{\text{Var}(F)}{(1 - K_0)^2}\right) v_2 + (1 - K_0)v_\delta + K_0 v_h + \frac{K_0(1 - K_0) - \text{Var}(F)}{(1 - K_0)^2} SS_{\mu_h} + \sigma_\epsilon^2. \quad (39)$$

The expectation of Y_s depends only on the expectation of d_s , the others being zero.

$$\mathbb{E}[d_{s1}] = \frac{F_s}{1 - K_0} \mu_{h1} + C, \quad (40)$$

where C is a constant that does not depend on s . Therefore,

$$\mathbb{E}[d_s] = \frac{F_s}{1 - K_0} \mu_h + 2C, \quad (41)$$

where $\mu_h = \mu_{h1} + \mu_{h2}$ is the inbreeding depression of the phenotype. Therefore

$$\text{Var}_s(\mathbb{E}[Y_s]) = \frac{\text{Var}(F)}{(1 - K_0)^2} \mu_h^2 \quad (42)$$

Therefore,

$$\text{Var}(Y) = v_1 + \left(1 + \frac{\text{Var}(F)}{(1 - K_0)^2}\right) v_2 + (1 - K_0)v_\delta + K_0v_h + \frac{K_0}{(1 - K_0)}SS_{\mu_h} + \frac{\text{Var}(F)}{(1 - K_0)^2}(\mu_h^2 - SS_{\mu_h}) + \sigma_\epsilon^2. \quad (43)$$

Asymptotic Variance of Fitting a Quadratic

We extend the analogy of fitting a quadratic to derive an analytic approximation of the standard error of the estimator of the variance from pairwise interactions. We imagine fitting the off diagonal elements of the covariance matrix as a quadratic function of $R_{s,t} = 2(K_{s,t} - K_0)/(1 - K_0)$, with normal error:

$$\Sigma_{st} \sim N(v_1R_{st} + v_2R_{st}^2, \sigma^2), \quad (44)$$

for all $\eta = N(N - 1)/2$ pairs st . This assumes that the off diagonal elements of the covariance matrix are independent, which may be problematic for samples which contain large sets of closely related individuals. The homoscedasticity assumption could also be problematic when there are many levels of relatedness present in the sample.

We invert the information matrix to obtain the asymptotic error of the maximum likelihood estimator of v_2 . We note that number of pairs, η , scales quadratically with the sample size, justifying the use of the asymptotic error for even moderate sample sizes.

The log likelihood is

$$l = \frac{\eta}{2} \log\left(\frac{\sigma^2}{2\pi}\right) - \sum_{i=1}^{\eta} \frac{(\Sigma_i - v_1R_i - v_2R_i^2)^2}{\sigma^2} \quad (45)$$

If we define $e_i = \Sigma_i - v_1R_i - v_2R_i^2$ to be the i^{th} residual, then the matrix of the second derivatives of the log-likelihood is

$$H = -\sigma^{-2} \begin{bmatrix} S_{R^2} & S_{R^3} & \sigma^{-2} \sum_{i=1}^{\eta} e_i R_i \\ S_{R^3} & S_{R^4} & \sigma^{-2} \sum_{i=1}^{\eta} e_i R_i^2 \\ \sigma^{-2} \sum_{i=1}^{\eta} e_i R_i & \sigma^{-2} \sum_{i=1}^{\eta} e_i R_i^2 & -2\sigma^{-4} \sum_{i=1}^{\eta} e_i^2 \end{bmatrix}, \quad (46)$$

where $S_{R^c} = \sum_{i=1}^{\eta} R_i^c$.

We now take the negative expectation of H to obtain the Fisher information matrix. Because $R_i = 2(K_i - K_0)/(1 - K_0)$, $\mathbb{E}[R] = 0$. Therefore $\mathbb{E}[S_{R^c}] = \eta\mu_c$, where μ_c is the c^{th} central moment of the distribution of R .

Therefore, the information matrix is

$$\mathbb{I} = \eta\sigma^{-2} \begin{bmatrix} \text{Var}(R) & \mu_3 & 0 \\ \mu_3 & \mu_4 & 0 \\ 0 & 0 & \sigma^{-2} \end{bmatrix}. \quad (47)$$

Using elimination to invert the matrix gives

$$\mathbb{I}^{-1} = \frac{\sigma^2}{\eta} \begin{bmatrix} \frac{1}{\text{Var}(R)} + \frac{\mu_3^2}{\mu_4 - \mu_3/\text{Var}(R)} & \frac{-\mu_3}{\text{Var}(R)\mu_4 - \mu_3^2} & 0 \\ \frac{-\mu_3}{\text{Var}(R)\mu_4 - \mu_3^2} & \left(\mu_4 - \frac{\mu_3^2}{\text{Var}(R)}\right)^{-1} & 0 \\ 0 & 0 & \sigma^2 \end{bmatrix} \quad (48)$$

This implies that the asymptotic standard error of the estimate of the variance from pairwise interactions is

$$\frac{\sigma}{\sqrt{\eta(\mu_4 - \mu_3^2/\text{Var}(R))}} \quad (49)$$

If the phenotype has been normalised to have variance one, then σ should be approximately 1.