

Methods for d-ldsc: Dominance LD-score regression

Duncan Palmer^{1,2,3*}, Wei Zhou^{1,2,3}, Liam Abbott¹,
Nik Baya¹, Claire Churchhouse^{1,2,3}, Cotton Seed^{1,3}, Tim Poterba^{1,3},
Daniel King^{1,3}, Masahiro Kanai^{1,2,3}, Alex Bloemendal¹ and Benjamin Neale^{1,2,3}

¹ Program in Medical and Population Genetics, Broad Institute of MIT and Harvard,
Cambridge, Massachusetts, USA.,

² Analytical and Translational Genetics Unit, Department of Medicine,
Massachusetts General Hospital and Harvard Medical School, Boston, Massachusetts, USA.,

³ Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard,
Cambridge, Massachusetts, USA.

* Correspondence to: duncan.stuart.palmer@gmail.com.

Methods: LD-score dominance extension summary

In order to estimate the contribution of within locus non-additive effects across the genome to phenotypic variation, we extend the infinitesimal model underlying LD-score regression to incorporate a non-additive effect at each site which is uncorrelated from the additive contribution at that site. We then ask: what is the variance of this additional contribution? Letting X^A denote the n samples \times m sites matrix of genotypes (after rescaling to enforce a mean of 0 and a variance of 1), we define X^D to be a re-coding of the genotypes whose columns (X_j^D) are orthonormal to the columns of X^A (X_j^A) under Hardy-Weinberg equilibrium. This extension to the additive infinitesimal model is:

$$y_i = \sum_{j=1}^m X_{i,j}^A \beta_{A_j} + \sum_{j=1}^m X_{i,j}^D \beta_{D_j} + \varepsilon_i, \quad (1)$$

where y_i is a continuous phenotype measured in individual i , β_{A_j} and β_{D_j} are the additive and dominance effect sizes at site j respectively, and ε_i is a noise term. We mean center the genotypes and construct an orthonormal basis using the Gram-Schmidt procedure. The resultant additive and dominance encodings are

$$\frac{1}{\sqrt{2pq}} \begin{bmatrix} -2p \\ 1 - 2p \\ 2 - 2p \end{bmatrix} \text{ and } \begin{bmatrix} -p/q \\ 1 \\ -q/p \end{bmatrix}, \quad (2)$$

respectively. We wish to estimate $h_D^2 := \text{Var}(\sum_{j=1}^m X_{i,j}^D \beta_{D_j})$. By defining additive summary statistics as the marginal effect sizes obtained by regressing y on X_j^A as usual; $\hat{\beta}_{A_j} = \frac{1}{n} (X_j^A)^\top y$, and introducing dominance summary statistics as the analogous marginal associations in the dominance encoding; $\hat{\beta}_{D_j} = \frac{1}{n} (X_j^D)^\top y$, we may then proceed to derive a dominance LD-score equation relating dominance summary statistics to dominance LD-scores:

$$\mathbb{E} [\chi_{D_j}^2] = \frac{nh_D^2}{m} l_j^D + 1 \quad (3)$$

where $l_j^D = \sum_{k=1}^m (r_{j,k}^D)^2$ is the sum of the squared correlations between SNP j and all other SNPs under the dominance encoding, and $\chi_{D_j}^2$ the chi-squared statistic for SNP j under the dominance encoding, $\chi_{D_j}^2 = n\beta_{D_j}^2$.

Detailed derivations of additive LD-score regression and dominance LD-score regression

Additive LD score regression

We first consider the standard LD score regression in which we assume the infinitesimal model: the effect sizes of the genotypes β_j ; $j = 1, 2, \dots, m$ are independent with mean 0 and variance $\frac{h^2}{m}$. This is formulated as

$$y = X\beta + \varepsilon. \quad (4)$$

That is,

$$y_i = \sum_{j=1}^m X_{i,j} \beta_j + \varepsilon_i, \quad (5)$$

where y is a vector of standardised phenotypes, and X is an $n \times m$ matrix of n standardised genotypes both taken from m independent samples (unscaled genotypes $X'_{i,j} \in \{0, 1, 2\}$, subject to rescaling: $X_{i,j} = (X'_{i,j} - 2p_j) / \sqrt{2p_j(1 - p_j)}$, where p_j is the prevalence of the genotype at site j , in the population). Finally we assume that there is are independent uncorrelated noise terms which absorb the effect of the environment.

$$\varepsilon \sim \mathcal{N}(0, 1 - h^2), \quad (6)$$

Note that $\text{Var}(\varepsilon) = 1 - h^2$ since $\text{Var}(X\beta) = h^2$ by definition of the narrow sense heritability ($\text{Var}\left(\sum_{j=1}^m X_{i,j} \beta_j\right) = \text{Var}(\beta_j) \text{Var}\left(\sum_{j=1}^m X_{i,j}\right) = \frac{h^2}{m} \sum_{j=1}^m \text{Var}(X_{i,j}) = h^2$) and $\text{Var}(y) = 1$ by construction. Under this model, all of the genetic heritability is additive.

We estimate effect sizes

$$\hat{\beta}_j = \frac{1}{n} \sum_{i=1}^n X_{i,j} y_i = \frac{1}{n} X_j^\top y \quad (7)$$

where X_j here are the rows of X . This is the marginal effect of SNP j in the sample, and is simply the estimate of the effect size of SNP j using a linear regression of genotype X_j against the phenotype y ($\hat{\beta}_j = \frac{\sum_{i=1}^n (X_{i,j} - \bar{X}_j)(y_i - \bar{y})}{\sum_{i=1}^n (X_{i,j} - \bar{X}_j)^2} = \frac{\sum_{i=1}^n X_{i,j} y_i}{n}$). Let us define χ_1^2 coefficients for each j as $\chi_j^2 = n \hat{\beta}_j^2$. Note that $n \hat{\beta}_j^2$ is χ_1^2 distributed under the null due to the central limit theorem:

Recall that if X_i are IID random variables with mean μ and variance σ^2 , then

$$\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n (X_i - \mu) \right) \rightarrow \mathcal{N}(0, \sigma^2) \quad (8)$$

in distribution. We have $\mathbb{E}[X_{i,j}y_i] = 0$, $\text{Var}(X_{i,j}y_i) = 1$. Thus $\sqrt{n}(\sum_{i=1}^m X_{i,j}y_i) \rightarrow \mathcal{N}(0, 1)$, so $\hat{\beta}_j \sim \frac{1}{\sqrt{n}}\mathcal{N}(0, 1)$. Now recall that if $Y = \sum_{i=1}^j X_i^2$ where $X_i \sim \mathcal{N}(0, 1)$, then $Y \sim \chi_j^2$. Thus, $n\hat{\beta}_j^2 \sim \chi_1^2$ (a χ^2 distribution with 1 degree of freedom).

Given our effect size estimate and our model in Equation (5), we can substitute into Equation (7):

$$\hat{\beta}_j = \frac{1}{n} \sum_{i=1}^n X_{i,j} \left(\sum_{k=1}^m X_{i,k} \beta_k + \varepsilon_i \right) \quad (9)$$

$$= \sum_{k=1}^m \beta_k \left(\frac{1}{n} \sum_{i=1}^n X_{i,j} X_{i,k} \right) + \tilde{\varepsilon}_j \quad (10)$$

$$= \sum_{k=1}^m \hat{r}_{j,k} \beta_k + \tilde{\varepsilon}_j \quad (11)$$

where $\tilde{\varepsilon}_j = \frac{1}{n} \sum_{i=1}^n X_{i,j} \varepsilon_i$. We now determine the expectation of our χ^2 coefficients to obtain the key formula.

$$\mathbb{E}[\chi_j^2] = \mathbb{E} \left[n \hat{\beta}_j^2 \right] \quad (12)$$

$$= \mathbb{E} \left[n \left(\sum_{k=1}^m \hat{r}_{j,k} \beta_k + \tilde{\varepsilon}_j \right)^2 \right] \quad (13)$$

$$= \mathbb{E} \left[n \sum_{k=1}^m \sum_{l=1}^m \hat{r}_{j,k} \hat{r}_{j,l} \beta_k \beta_l \right] + \underbrace{\mathbb{E} \left[2n \sum_{k=1}^m \hat{r}_{j,k} \tilde{\varepsilon}_j \right]}_{=0} + \mathbb{E} \left[n \tilde{\varepsilon}_j^2 \right] \quad (14)$$

$$= \underbrace{n \mathbb{E} \left[\sum_{k \neq l} \hat{r}_{j,k} \hat{r}_{j,l} \beta_k \beta_l \right]}_{=0} + \mathbb{E} \left[n \sum_{k=1}^m \hat{r}_{j,k}^2 \beta_k^2 \right] + n \mathbb{E} \left[\tilde{\varepsilon}_j^2 \right] \quad (15)$$

$$= n \left(\sum_{k=1}^m \mathbb{E} \left[\hat{r}_{j,k}^2 \right] \mathbb{E} \left[\beta_k^2 \right] \right) + n \mathbb{E} \left[\tilde{\varepsilon}_j^2 \right] \quad (16)$$

by independence and linearity of expectations. Since β_k is standardised, $\text{Var}(\beta_k) = \mathbb{E}[\beta_k^2] = \frac{h^2}{m}$. Also, in an unstructured sample, we can make the approximation $\mathbb{E}[\hat{r}_{j,k}^2] \approx r_{j,k}^2 + \frac{1}{n}$. This approximation comes from noting that \hat{r} is unbiased as we know the population variances for X_j ; $j \in \{1, 2, \dots, m\}$ are 1. Therefore, setting $X = X_j$, $Y = X_k$ to make notation simpler

and using subscripts i and l to reiterate that summation is over samples,

$$\mathbb{E} [\hat{r}_{X,Y}^2] = \mathbb{E} \left[\left(\frac{1}{n} \sum_{i=1}^n X_i Y_i \right)^2 \right] \quad (17)$$

$$= \frac{1}{n^2} \mathbb{E} \left[\sum_{l=1}^n \sum_{i=1}^n X_i Y_i X_l Y_l \right] \quad (18)$$

$$= \frac{1}{n^2} \sum_{i,l: i \neq l}^n \mathbb{E} [X_i Y_i] \mathbb{E} [X_l Y_l] + \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [X_i^2 Y_i^2] \quad (19)$$

$$= \frac{1}{n^2} (n^2 - n) r_{X,Y}^2 + \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [X_i^2 Y_i^2]. \quad (20)$$

since $\mathbb{E} [X_i Y_i] = r_{X,Y}$ and the genotypes of individuals i and l are independent in an unstructured population. Now, the majority of X and Y will be far apart and so approximately independent. Therefore,

$$\mathbb{E} [\hat{r}_{X,Y}^2] \approx \left(1 - \frac{1}{n} \right) r_{X,Y}^2 + \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [X_i^2] \mathbb{E} [Y_i^2] = \left(1 - \frac{1}{n} \right) r_{X,Y}^2 + \frac{1}{n} \quad (21)$$

$$\approx r_{X,Y}^2 + \frac{1}{n}. \quad (22)$$

Substituting this approximation into Equation (16), we have

$$\mathbb{E} [\chi_j^2] \approx \frac{nh^2}{m} \sum_{k=1}^m \left(r_{j,k}^2 + \frac{1}{n} \right) + \mathbb{E} [\tilde{\varepsilon}_j^2]. \quad (23)$$

Since $X_{i,j}$ and ε_i are independent, the $X_{i,j}$ are standardised, and $\mathbb{E} [\tilde{\varepsilon}_j] = \mathbb{E} [\varepsilon_i] = 0$,

$$\mathbb{E} [\tilde{\varepsilon}_j^2] = \text{Var} (\tilde{\varepsilon}_j) = \text{Var} \left(\frac{1}{n} \sum_{i=1}^n X_{i,j} \varepsilon_i \right) \quad (24)$$

$$= \frac{1}{n^2} \sum_{i=1}^n \text{Var} (X_{i,j}) \text{Var} (\varepsilon_i) \quad (25)$$

$$= \frac{1}{n^2} \sum_{i=1}^n \text{Var} (\varepsilon_i) \quad (26)$$

$$= \frac{1}{n} (1 - h^2). \quad (27)$$

Thus,

$$\mathbb{E} [\chi_j^2] \approx \frac{nh^2}{m} \sum_{k=1}^m r_{j,k}^2 + h^2 + 1 - h^2. \quad (28)$$

Finally, letting $l_j := \sum_{k=1}^m r_{j,k}^2$; dubbed the LD score of SNP j , we obtain the LD score regression formula:

$$\mathbb{E} [\chi_j^2] = \frac{nh^2}{m} l_j + 1. \quad (29)$$

LD score regression with dominance

For a more general model, we now allow a dominance term to affect the phenotype y .

$$y = X\beta + \varepsilon \quad (30)$$

$$= X^A \beta_A + X^D \beta_D + \varepsilon \quad (31)$$

$$y_i = \sum_{j=1}^m X_{i,j}^A \beta_{A_j} + \sum_{j=1}^m X_{i,j}^D \beta_{D_j} + \varepsilon_i \quad (32)$$

where now X is the concatenation of two matrices; X^A and X^D (similarly, $\beta = (\beta_A, \beta_D)$). X^A encodes the genotype (in exactly the same way as X does in Equation (5)), and X^D encodes a dominance contribution to the phenotype y . We will pick an zero mean-centered orthonormal encoding for X^A and X^D . This encoding isolates the pure dominance and additive contributions of the genotypes. We start with the basis:

$$X^{0'} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \quad X^{A'} = \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix}, \quad X^{D'} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \quad (33)$$

where the first, second and third entries are the contributions of homozygous reference, heterozygous, and homozygous variant respectively. We then apply the Gram-Schmidt process (a method to orthonormalise a set of vectors in an inner product space), where our inner product is defined as

$$\langle X, Y \rangle = q^2 X_0 Y_0 + 2pq X_1 Y_1 + p^2 X_2 Y_2. \quad (34)$$

The subscript here denotes the genotype of the entry: (hom ref, het, hom var) for (0,1,2). We use this inner product as we wish to weight the contributions of the genotypes according to their prevalence in the population, assuming Hardy-Weinberg equilibrium. The inclusion of $[1, 1, 1]^T$ will ensure that our resultant orthonormal basis has mean 0. Proceeding with Gram-Schmidt:

$$X^0 = \frac{1}{\|X^{0'}\|} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \quad (35)$$

$$X^A = \frac{X^{A'} - \langle X^{A'}, X^0 \rangle X^0}{\|X^{A'} - \langle X^{A'}, X^0 \rangle X^0\|} = \frac{1}{\sqrt{2pq}} \left(\begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix} - (2pq + 2q^2) \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right) \quad (36)$$

$$= \frac{1}{\sqrt{2pq}} \begin{bmatrix} -2p \\ 1 - 2p \\ 2 - 2p \end{bmatrix} \quad (37)$$

$$= \frac{1}{\sqrt{2pq}} \begin{bmatrix} -2p \\ q - p \\ 2q \end{bmatrix}, \quad (38)$$

$$X^D = \frac{X^{D'} - \langle X^{D'}, X^0 \rangle X^0 - \langle X^{D'}, X^A \rangle X^A}{\|X^{D'} - \langle X^{D'}, X^0 \rangle X^0 - \langle X^{D'}, X^A \rangle X^A\|} = \frac{\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} - 2pq \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - (q - p) \begin{bmatrix} -2p \\ q - p \\ 2q \end{bmatrix}}{\|X^{D'} - \langle X^{D'}, X^0 \rangle X^0 - \langle X^{D'}, X^A \rangle X^A\|} \quad (39)$$

$$= \frac{\begin{bmatrix} -2p^2 \\ 2pq \\ -2q^2 \end{bmatrix}}{\|X^{D'} - \langle X^{D'}, X^0 \rangle X^0 - \langle X^{D'}, X^A \rangle X^A\|} \quad (40)$$

$$= \frac{1}{\sqrt{q^2 p^4 + 2p^3 q^3 + p^2 q^4}} \begin{bmatrix} -p^2 \\ pq \\ -q^2 \end{bmatrix} \quad (41)$$

$$= \frac{1}{pq} \begin{bmatrix} -p^2 \\ pq \\ -q^2 \end{bmatrix}. \quad (42)$$

Each site has additive and dominance effect sizes (β_A and β_D) which act along X^A and X^D respectively. β_A and β_D are drawn from distributions with mean 0 and variances $\frac{h_A^2}{m}$, $\frac{h_D^2}{m}$ respectively. We assume $\beta_j = (\beta_{A_j}, \beta_{D_j})$ are independent across sites, but may have dependencies at a given locus. Let's determine the variance of $X_i \beta (= X_i^A \beta_A + X_i^D \beta_D)$, where here, X_i denotes

the i^{th} row of the matrix X .

$$\text{Var}(X_i\beta) = \text{Var}\left(\sum_{j=1}^m X_{i,j}\beta_j\right) \quad (43)$$

$$= \mathbb{E}\left[\sum_{j=1}^m \sum_{l=1}^m X_{i,j}X_{i,l}\beta_j\beta_l\right] \quad (44)$$

$$= \mathbb{E}\left[\sum_{j=1}^m \sum_{l=1}^m X_{i,j}^A X_{i,l}^A \beta_{A_j} \beta_{A_l} + X_{i,j}^A X_{i,l}^D \beta_{A_j} \beta_{D_l} + X_{i,j}^D X_{i,l}^A \beta_{D_j} \beta_{A_l} + X_{i,j}^D X_{i,l}^D \beta_{D_j} \beta_{D_l}\right]. \quad (45)$$

We may remove all cross terms over sites as we assume that β_{A_j}, β_{D_j} are independent across sites, and also independent of the genetic data encoded in X^A and X^D . Thus,

$$\mathbb{E}[X_i\beta] = \mathbb{E}\left[\sum_{j=1}^m X_{i,j}^A{}^2 \beta_{A_j}^2 + X_{i,j}^A X_{i,j}^D \beta_{A_j} \beta_{D_j} + X_{i,j}^D X_{i,j}^A \beta_{D_j} \beta_{A_j} + X_{i,j}^D{}^2 \beta_{D_j}^2\right]. \quad (46)$$

Now we may use that $\mathbb{E}[X_{i,j}^A X_{i,j}^D] = \mathbb{E}[X_{i,j}^A] \mathbb{E}[X_{i,j}^D] = 0$, as X^A and X^D are orthogonal by construction. Thus, we are left with

$$\text{Var}(X_i\beta) = \mathbb{E}\left[\sum_{j=1}^m X_{i,j}^A{}^2 \beta_{A_j}^2\right] + \mathbb{E}\left[\sum_{j=1}^m X_{i,j}^D{}^2 \beta_{D_j}^2\right] \quad (47)$$

$$= \sum_{j=1}^m \mathbb{E}[X_{i,j}^A{}^2] \mathbb{E}[\beta_{A_j}^2] + \sum_{j=1}^m \mathbb{E}[X_{i,j}^D{}^2] \mathbb{E}[\beta_{D_j}^2] \quad (48)$$

$$= \frac{h_A^2}{m} \sum_{j=1}^m 1 + \frac{h_D^2}{m} \sum_{j=1}^m 1 = h_A^2 + h_D^2. \quad (49)$$

So $\text{Var}(X\beta) = h_A^2 + h_D^2 = h^2$. As in the standard infinitesimal model, the noise term ε has mean 0 and variance $1 - h^2$. $h_A^2 \leq h_A^2 + h_D^2 = h^2$. In particular, equalities hold when the dominance terms don't contribute to the variance.

We wish to determine the extent to which $h^2 > h_A^2$. To achieve this we can simply compare the sums of the first order terms of two LD score regression models. The first under the standard infinitesimal model as outlined above, and the second according the model introduced in Equation (30). It remains to derive the analogue of Equation (29), which can be obtained in a similar manner to the original LD score regression model.

Again, we estimate effect sizes but this time regress our dominance encoded SNPs against the phenotypes as well:

$$\hat{\beta}_{A_j} = \frac{1}{n} \sum_{i=1}^n X_{i,j}^A y_i = \frac{1}{n} X_j^{A\top} y; \quad \hat{\beta}_{D_j} = \frac{1}{n} \sum_{i=1}^n X_{i,j}^D y_i = \frac{1}{n} X_j^{D\top} y. \quad (50)$$

This time, we substitute y_i using Equation (30):

$$\hat{\beta}_{A_j} = \frac{1}{n} \sum_{i=1}^n X_{i,j}^A \left(\sum_{k=1}^m X_{i,k}^A \beta_{A_k} + \sum_{k=1}^m X_{i,k}^D \beta_{D_k} + \varepsilon_i \right) \quad (51)$$

$$= \sum_{k=1}^m \beta_{A_k} \frac{1}{n} \sum_{i=1}^n X_{i,j}^A X_{i,k}^A + \sum_{k=1}^m \beta_{D_k} \frac{1}{n} \sum_{i=1}^n X_{i,j}^A X_{i,k}^D + \frac{1}{n} \sum_{i=1}^n X_{i,j}^A \varepsilon_i \quad (52)$$

$$= \sum_{k=1}^m \hat{r}_{j,k}^{AA} \beta_{A_k} + \sum_{k=1}^m \hat{r}_{j,k}^{AD} \beta_{D_k} + \frac{1}{n} \sum_{i=1}^n X_{i,j}^A \varepsilon_i. \quad (53)$$

Similarly,

$$\hat{\beta}_{D_j} = \sum_{k=1}^m \hat{r}_{j,k}^{DA} \beta_{A_k} + \sum_{k=1}^m \hat{r}_{j,k}^{DD} \beta_{D_k} + \frac{1}{n} \sum_{i=1}^n X_{i,j}^D \varepsilon_i. \quad (54)$$

Setting $\tilde{\varepsilon}_j = \frac{1}{n} \sum_{i=1}^n X_{i,j}^A \varepsilon_i$, we now determine the expectation of the χ^2 statistic,

$$\chi_{A_j}^2 := n \hat{\beta}_{A_j}. \quad (55)$$

(Obtaining the expectation of $\chi_{D_j}^2 := n \hat{\beta}_{D_j}$ is analogous).

$$\mathbb{E} [\chi_{A_j}^2] = \mathbb{E} [n \beta_{A_j}^2] = n \mathbb{E} \left[\left(\sum_{k=1}^m \hat{r}_{j,k}^{AA} \beta_{A_k} + \sum_{k=1}^m \hat{r}_{j,k}^{AD} \beta_{D_k} + \tilde{\varepsilon}_j \right)^2 \right] \quad (56)$$

$$= n \mathbb{E} \left[\sum_{k=1}^m \sum_{l=1}^m (\hat{r}_{j,k}^{AA} \beta_{A_k} \hat{r}_{j,l}^{AA} \beta_{A_l} + \hat{r}_{j,k}^{AA} \beta_{A_k} \hat{r}_{j,l}^{AD} \beta_{D_l} + \right. \quad (57)$$

$$\left. \hat{r}_{j,k}^{AD} \beta_{D_k} \hat{r}_{j,l}^{AA} \beta_{A_l} + \hat{r}_{j,k}^{AD} \beta_{D_k} \hat{r}_{j,l}^{AD} \beta_{D_l} \right) + \underbrace{\tilde{\varepsilon}_j \left(\sum_{k=1}^m \hat{r}_{j,k}^{AA} \beta_{A_k} + \hat{r}_{j,k}^{AD} \beta_{D_k} \right)}_{=0} + \tilde{\varepsilon}_j^2 \right]$$

Again, we may remove all cross terms over sites due to the independence of the β s.

$$= n \mathbb{E} \left[\sum_{k=1}^m \hat{r}_{j,k}^{AA^2} \beta_{A_k}^2 \right] + 2n \mathbb{E} \left[\sum_{k=1}^m \hat{r}_{j,k}^{AA} \beta_{A_k} \hat{r}_{j,k}^{AD} \beta_{D_k} \right] + \quad (58)$$

$$n \mathbb{E} \left[\sum_{k=1}^m \hat{r}_{j,k}^{AD^2} \beta_{D_k}^2 \right] + n \mathbb{E} [\tilde{\varepsilon}_j^2]$$

We may replace the $\widehat{r}_{j,k}^{AA^2}$ and $\widehat{r}_{j,k}^{AD^2}$ as before, but we also have a collection of cross terms (i.e. $\widehat{r}_{j,k}^{AA} \beta_{A_k} \widehat{r}_{j,k}^{AD} \beta_{D_k}$) to consider. Let's determine the expectation of this summation of cross terms.

$$\mathbb{E} \left[\sum_{k=1}^m \beta_{A_k} \beta_{D_k} \widehat{r}_{j,k}^{AA} \widehat{r}_{j,k}^{AD} \right] = \sum_{k=1}^m \mathbb{E} [\beta_{A_k} \beta_{D_k}] \mathbb{E} [\widehat{r}_{j,k}^{AA} \widehat{r}_{j,k}^{AD}]. \quad (59)$$

There may be dependency between the additive and dominance effect sizes at a locus. So let's turn our attention to $\mathbb{E} [\widehat{r}_{j,k}^{AA} \widehat{r}_{j,k}^{AD}]$:

$$\mathbb{E} [\widehat{r}_{j,k}^{AA} \widehat{r}_{j,k}^{AD}] = \mathbb{E} \left[\frac{1}{n^2} \sum_{i=1}^n \sum_{l=1}^n X_{i,j}^A X_{i,k}^A X_{l,j}^A X_{l,k}^D \right] \quad (60)$$

$$= \frac{1}{n^2} \sum_{i,l: i \neq l} \mathbb{E} [X_{i,j}^A X_{i,k}^A X_{l,j}^A X_{l,k}^D] + \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [X_{i,j}^{A^2} X_{i,k}^A X_{i,k}^D] \quad (61)$$

$$= \frac{1}{n^2} (n^2 - n) r_{j,k}^{AA} r_{j,k}^{AD} + \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [X_{i,j}^{A^2} X_{i,k}^A X_{i,k}^D]. \quad (62)$$

Since, j and k are far apart for most SNPs, we make the approximation

$$\mathbb{E} [\widehat{r}_{j,k}^{AA} \widehat{r}_{j,k}^{AD}] \approx \left(1 - \frac{1}{n}\right) r_{j,k}^{AA} r_{j,k}^{AD} + \sum_{i=1}^n \mathbb{E} [X_{i,j}^{A^2}] \mathbb{E} [X_{i,k}^A X_{i,k}^D]. \quad (63)$$

Note that $X_{i,j}^D$ and $X_{i,j}^A$ are independent for all j by construction. So,

$$\mathbb{E} [\widehat{r}_{j,k}^{AA} \widehat{r}_{j,k}^{AD}] \approx r_{j,k}^{AA} r_{j,k}^{AD}. \quad (64)$$

Plugging in this approximation, and the approximation in Equation (22) into Equation (58), we have

$$\mathbb{E} [\chi_{A_j}^2] \approx n \underbrace{\sum_{k=1}^m \left(r_{j,k}^{AA^2} + \frac{1}{n} \right) \mathbb{E} [\beta_{A_k}^2]}_{\text{by Equation (22)}} + n \underbrace{\sum_{k=1}^m \left(r_{j,k}^{AD^2} + \frac{1}{n} \right) \mathbb{E} [\beta_{D_k}^2]}_{\text{by Equation (22)}} + \quad (65)$$

$$\begin{aligned} & \sum_{k=1}^m \mathbb{E} [\beta_{A_k} \beta_{D_k}] \underbrace{r_{j,k}^{AA} r_{j,k}^{AD}}_{\text{by Equation (64)}} + n \mathbb{E} [\widetilde{\varepsilon}_j^2] \\ &= \frac{nh_A^2}{m} \sum_{k=1}^m r_{j,k}^{AA^2} + \frac{nh_D^2}{m} \sum_{k=1}^m r_{j,k}^{AD^2} + h_A^2 + h_D^2 + \sum_{k=1}^m \mathbb{E} [\beta_{A_k} \beta_{D_k}] r_{j,k}^{AA} r_{j,k}^{AD} + \frac{n(1-h^2)}{n} \end{aligned} \quad (66)$$

$$= \frac{nh_A^2}{m} \sum_{k=1}^m r_{j,k}^{AA^2} + \frac{nh_D^2}{m} \sum_{k=1}^m r_{j,k}^{AD^2} + \sum_{k=1}^m \mathbb{E} [\beta_{A_k} \beta_{D_k}] r_{j,k}^{AA} r_{j,k}^{AD} + 1. \quad (67)$$

j/k	0	1
0	$1 - p_j - p_k$	$p_j - p_{j,k}$
1	$p_k - p_{j,k}$	$p_{j,k}$

Table S1

Similarly,

$$\mathbb{E} [\chi_{D_j}^2] = \frac{nh_A^2}{m} \sum_{k=1}^m r_{j,k}^{DA^2} + \frac{nh_D^2}{m} \sum_{k=1}^m r_{j,k}^{DD^2} + \sum_{k=1}^m \mathbb{E} [\beta_{A_k} \beta_{D_k}] r_{j,k}^{DD} r_{j,k}^{AD} + 1. \quad (68)$$

We can investigate the contribution of each of the correlation terms analytically. Let p_j and p_k denote the probability of observing the alternate genotype at sites j and k respectively, and $p_{j,k}$ denote the probability of observing the alternate genotype at both sites j and k .

This then leads to the probability mass function of linked genotypes in offspring (in the absence of inbreeding) by multiplying the relevant entries of Table S1. We may then determine the correlation coefficients (r^{AA} , r^{AD} , r^{DA} , r^{DD}):

$$r^{QR} = \mathbb{E} [X_Q^\top X_R] = \sum_{g_j, g_k \in \{0,1,2\}} P(g_j, g_k) X_j^Q(g_j) X_k^R(g_k); \quad Q, R \in \{A, D\}. \quad (69)$$

where

$$X_j^A = \frac{1}{\sqrt{2p_j q_j}} \begin{bmatrix} -2p_j \\ q_j - p_j \\ 2q_j \end{bmatrix}, \quad X_j^D = \frac{1}{p_j q_j} \begin{bmatrix} -p_j^2 \\ p_j - q_j \\ -q_j^2 \end{bmatrix}. \quad (70)$$

Here, using an abuse of notation, we mean the encoding of the genotypes (homozygous reference, heterozygous, homozygous variant) at site j .

After cancelling, we obtain:

$$r^{AD} = 0, \quad r^{DA} = 0, \quad (71)$$

$$r^{AA} = \frac{p_{j,k} - p_j p_k}{\sqrt{p_j p_k (1 - p_j) (1 - p_k)}}, \quad r^{DD} = r^{AA^2} = \frac{(p_{j,k} - p_j p_k)^2}{p_j p_k (1 - p_j) (1 - p_k)}. \quad (72)$$

Equations (67) and (68) then simplify to

$$\mathbb{E} [\chi_{A_j}^2] = \frac{nh_A^2}{m} l_j^A + 1; \quad (73)$$

$$\mathbb{E} [\chi_{D_j}^2] = \frac{nh_D^2}{m} l_j^D + 1, \quad (74)$$

where $l_j^A = \sum_{k=1}^m r_{j,k}^{AA^2}$ and $l_j^D = \sum_{k=1}^m r_{j,k}^{DD^2}$.