

## 1 Supporting Information

### 2 Appendix S1

#### 3 $g_2$ statistics

4 Consider a sufficiently large population of individuals in inbreeding equilibrium. Sample  $N$  individuals out  
 5 of this population, each sequenced at a set of loci  $\{1, \dots, L\}$ . Strong inbreeding increases the dependence  
 6 of loci being homozygous within an individual. Thus, inbreeding is assumed to have an effect on the joint  
 7 distribution of homozygous loci within an individual compared to the marginal assortment of single-locus  
 8 ones. Based on David *et al.* (2007) we recall three different estimators for the second-order heterozygosity  
 9 disequilibrium  $g_2$ , which reflects the excess of joint heterozygous loci relative to their expectation under  
 10 random assortment. For a better comparison to the implementation in R, see the following table:

Estimator	Dataset	Equation	R function
$\hat{g}_2$	small datasets without missing values	eqn 3	
$\hat{g}'_2$	small datasets that include missing values	eqn 5	<code>g2_microsats</code>
$\hat{\underline{g}}_2$	large datasets; missing values do not differ too much across loci (see last paragraph)	eqn 7	<code>g2_snps</code>

#### 12 Notations and Mathematical background

13 Let  $h_i$  denote the true heterozygosity at locus  $i$  in the population. Due to scoring artefacts, the true  
 14 value might differ from the apparent one, which will be denoted here by  $H_i$ . For our sample, define the  
 15 indicator function  $H_{ik}$ ,  $i = 1, \dots, L$ ,  $k = 1, \dots, N$  as follows: Set  $H_{ik} = 1$ , if locus  $i$  is heterozygous in  
 16 individual  $k$  and set  $H_{ik} = 0$  if locus  $i$  is found homozygous.

17 Following Weir & Cockerham (1973), the second-order heterozygosity disequilibrium  $g_2(i, j)$  between loci  
 18  $i$  and  $j$  might be quantified through the identity disequilibrium  $\mathbb{E}[h_i h_j] = \mathbb{E}[h_i] \mathbb{E}[h_j] (1 + g_2(i, j))$ . If one  
 19 assumes independent scoring artefacts across the set of loci, then this identity also holds for the apparent  
 20 heterozygosity, i.e.  $\mathbb{E}[H_i H_j] = \mathbb{E}[H_i] \mathbb{E}[H_j] (1 + g_2(i, j))$ . Commonly,  $g_2$  is assumed to be constant for  
 21 every pair of loci and defined via

$$\mathbb{E}[H_i H_j] = \mathbb{E}[H_i] \mathbb{E}[H_j] (1 + g_2). \quad (\text{eqn 1})$$

22 In practice, tightly linked pairs of loci probably have a higher  $g_2$ , which is why in the following we will give  
 23 an alternative to eqn 1 by averaging out over all locus pairs  $(i, j)$ . To avoid confusions, let us denote the  
 24 'averaged'  $g_2$  by  $\bar{g}_2$ , which will serve for a more robust estimator and goes back to David *et al.* (2007). The

well known decomposition for the variance  $\text{Var}[\sum_i H_i] = \sum_i \text{Var}[H_i] + \sum_i \sum_{j \neq i} \text{Cov}[H_i, H_j]$ , together with

$$\text{Cov}[H_i, H_j] = \mathbb{E}[H_i H_j] - \mathbb{E}[H_i] \mathbb{E}[H_j] = \mathbb{E}[H_i] \mathbb{E}[H_j] g_2$$

leads to an expression of  $\bar{g}_2$  of the form

$$\bar{g}_2 = \frac{\sum_{i=1}^L \sum_{j \neq i} \mathbb{E}[H_i H_j]}{\sum_{i=1}^L \sum_{j \neq i} \mathbb{E}[H_i] \mathbb{E}[H_j]} - 1. \quad (\text{eqn 2})$$

For the averaged quantity  $\bar{g}_2$  in eqn 2 one can find an estimator given in eqn (8) in David *et al.* (2007) (corrected for typographical errors) given by

$$\hat{g}_2 = \frac{\sum_{i=1}^L \sum_{j \neq i} \left( \sum_{k=1}^N H_{ik} H_{jk} \right)}{\frac{1}{N-1} \sum_{i=1}^L \sum_{j \neq i} \left( \sum_{k_1=1}^N \sum_{k_2 \neq k_1} H_{ik_1} H_{jk_2} \right)} - 1, \quad (\text{eqn 3})$$

with little bias of order  $1/N$  (see Appendix S1 in (David *et al.*, 2007)). Problematically, real data sets, especially *microsatellites* ones, do have missing values, such that the apparent value  $H_{ik}$  might be unknown for some pairs  $(i, k)$ . In this case define

$$\tilde{H}_{ik} = \begin{cases} 1, & \text{if locus } i \text{ is heterozygous in individual } k \\ 0, & \text{if locus } i \text{ is either homozygous in individual } k, \text{ or unknown,} \end{cases}$$

as well as  $M_{ik} = 1$  if the datum is missing at locus  $i$  in individual  $k$  and  $M_{ik} = 0$  otherwise.

The expected values of  $H_i$  and  $\tilde{H}_i$  are correlated via  $\mathbb{E}[\tilde{H}_i] = (1 - m_i) \mathbb{E}[H_i]$ , where  $m_i := \frac{1}{N} \sum_{k=1}^N M_{ik}$  is the proportion of individuals with missing data at locus  $i$ . Equivalently, for the joint distribution one finds the identity  $\mathbb{E}[\tilde{H}_i \tilde{H}_j] = (1 - m_i - m_j + m_{ij}) \mathbb{E}[H_i H_j]$ , with  $m_{ij}$  being the proportion of individuals with missing values both at loci  $i$  and  $j$ . Note, that  $1 - m_i - m_j + m_{ij}$  is the exact proportion of individuals with non-missing values at both, loci  $i$  and  $j$ . The analogue of eqn 1 now reads

$$\mathbb{E}[\tilde{H}_i \tilde{H}_j] = \frac{(1 - m_i - m_j + m_{ij})}{(1 - m_i)(1 - m_j)} \mathbb{E}[\tilde{H}_i] \mathbb{E}[\tilde{H}_j] (1 + g_2),$$

39 which, with the same procedure as above leads to the more robust averaged parameter

$$\bar{g}_2 = \frac{\sum_{i=1}^L \sum_{j \neq i} \mathbb{E}[\tilde{H}_{ij} \tilde{H}_j]}{\sum_{i=1}^L \sum_{j \neq i} \frac{(1-m_i-m_j+m_{ij})}{(1-m_i)(1-m_j)} \mathbb{E}[\tilde{H}_i] \mathbb{E}[\tilde{H}_j]} - 1,$$

40 where the underline stands for presence of missing data and the overbar again indicates averaging over  
41 all locus pairs. Eqn S1 in the Supplementary Informations in Hoffman *et al.* (2014) provide an estimator  
42 for  $\bar{g}_2$ , which can be rewritten as

$$\hat{g}_2 = \frac{\sum_{i=1}^L \sum_{j \neq i} \frac{1}{N-M_i-M_j+M_{ij}} \left( \sum_{k=1}^N \tilde{H}_{ik} \tilde{H}_{jk} \right)}{\sum_{i=1}^L \sum_{j \neq i} \frac{1}{(N-1)(N-M_i-M_j)+M_i M_j - M_{ij}} \left( \sum_{k_1=1}^N \sum_{k_2 \neq k_1} \tilde{H}_{ik_1} \tilde{H}_{ik_2} \right)} - 1, \quad (\text{eqn 4})$$

43 where the  $M_i, M_{ij} \in \{1, \dots, N\}$ , compared to the  $m_i, m_{ij}$ , now relate to the absolute values of missing  
44 data at some loci in contrast to the relative ones. Eventhough, eqn S1 in Hoffman *et al.* (2014) differs  
45 from (eqn 4) at first glance, it is straightforward to verify their equality.

46 Unfortunately, eqn 4 entails a weighting problem for loci with many missing data. To see this, imagine  
47 a pair of loci  $(i, j)$  with a considerable high fraction of missing values in the sample. In this case, any  
48 individual  $k$  which is heterozygous at both loci is downweighted, whereas in case of no missing data, it  
49 will be fully taken into account. Therefore one can update the formula to

$$\hat{g}'_2 = \frac{\left[ \sum_{i=1}^L \sum_{j \neq i} ((N-1)(N-M_i-M_j) + M_i M_j - M_{ij}) \right] \sum_{i=1}^L \sum_{j \neq i} \left( \sum_{k=1}^N \tilde{H}_{ik} \tilde{H}_{jk} \right)}{\left[ \sum_{i=1}^L \sum_{j \neq i} (N-M_i-M_j+M_{ij}) \right] \left( \sum_{i=1}^L \sum_{j \neq i} \sum_{k_1=1}^N \sum_{k_2 \neq k_1} \tilde{H}_{ik_1} \tilde{H}_{ik_2} \right)} - 1, \quad (\text{eqn 5})$$

50 which now weights with the appropriate number of double non-missing data over all locus pairs (Hardy,  
51 2015). It is indeed eqn 5 (not eqn 4), which is implemented in the RMES software by David *et al.* (2007).  
52 The R package `g2_microsats` is also based on eqn 5.

53 All estimators we have seen so far, require double summation over all loci and are thus unfeasible for  
54 large datasets. To fasten the algorithm, the Supplementary Information in Hoffman *et al.* (2014) provides  
55 another estimator based on the general decomposition of a double sum of the form  $\sum_i \sum_{j \neq i} a_i a_j$  into  
56  $\sum_{i=1}^n \sum_{j \neq i} a_i a_j = \left( \sum_{i=1}^n a_i \right)^2 - \sum_{i=1}^n a_i^2$ , where  $a_i \in \mathbb{R}$  is some parameter and  $n \in \mathbb{N}$ . The estimator  
57 reads

$$\hat{g}_2 = \frac{1 + (\hat{B} - \hat{C})/\hat{A}}{1 + \hat{a}} - 1, \quad (\text{eqn 6})$$

58 where  $\hat{A}$ ,  $\hat{B}$ ,  $\hat{C}$ ,  $\hat{a}$  are encapsulated weighted sums of  $\tilde{H}_{ik}$ :

$$\begin{aligned}\hat{A} &= \frac{N}{N-1} \left[ \left( \sum_{i=1}^L \hat{\mu}_i \right)^2 - \sum_{i=1}^L \hat{\mu}_i^2 \right] - \frac{\hat{J}}{N-1}, & \hat{J} &= \frac{1}{N} \sum_{k=1}^N \tilde{h}_k^2 - \sum_{i=1}^L \hat{\mu}_i, \\ \hat{B} &= \frac{1}{N-1} \left[ \sum_{k=1}^N \tilde{h}_k^2 - \frac{1}{N} \left( \sum_{k=1}^N \tilde{h}_k \right)^2 \right], & \hat{C} &= \frac{N}{N-1} \left[ \sum_{i=1}^L \hat{\mu}_i - \sum_{i=1}^L \hat{\mu}_i^2 \right], \\ \hat{a} &= \frac{\sum_{k=1}^N \hat{M}_k - N \left( \sum_{i=1}^L \hat{\mu}_i \frac{m_i}{1-m_i} \right)^2 + N \sum_{i=1}^L \left( \hat{\mu}_i \frac{m_i}{1-m_i} \right)^2}{(N-1)\hat{A} + \hat{J}}, & \hat{M}_k &= \left( \sum_{i=1}^L \frac{\hat{\mu}_i x_{ik}}{1-m_i} \right)^2 - \sum_{i=1}^L \left( \frac{\hat{\mu}_i x_{ik}}{1-m_i} \right)^2,\end{aligned}$$

59 with  $\hat{\mu}_i := \frac{1}{N} \sum_{i=1}^N \tilde{H}_{ik}$ ,  $x_{ik} := M_{ik}$  and  $h_k := \sum_{k=1}^L \tilde{H}_{ik}$ . A straightforward calculation leads to a  
60 simplification of eqn 6 that uses less normalisation steps, represented as

$$\hat{g}_2 = \frac{\hat{D}}{1+\hat{E}} - 1, \quad (\text{eqn 7})$$

61 with

$$\hat{D} := (N-1) \frac{\sum_{k=1}^N (\tilde{H}_{.k})^2 - \tilde{H}_{..}}{(\tilde{H}_{..})^2 - \sum_{i=1}^L (\tilde{H}_{i.})^2 - \sum_{k=1}^N (\tilde{H}_{.k})^2 + \tilde{H}_{..}}$$

62 and

$$\hat{E} := \frac{\frac{1}{N} \sum_{k=1}^N \left[ \left( \sum_{i=1}^L \tilde{H}_{i.} \frac{M_{ik}}{1-m_i} \right)^2 - \sum_{i=1}^L \left( \tilde{H}_{i.} \frac{M_{ik}}{1-m_i} \right)^2 \right] + \sum_{i=1}^L \left( \tilde{H}_{i.} \frac{m_i}{1-m_i} \right)^2 - \left( \sum_{i=1}^L \tilde{H}_{i.} \frac{m_i}{1-m_i} \right)^2}{(\tilde{H}_{..})^2 - \sum_{i=1}^L (\tilde{H}_{i.})^2},$$

63 where we used the common notation for the marginalisation over all individuals or respectively over all  
64 loci; more precisely

$$\tilde{H}_{.k} = \sum_{i=1}^L \tilde{H}_{ik}, \quad \tilde{H}_{i.} = \sum_{k=1}^N \tilde{H}_{ik}, \quad \tilde{H}_{..} = \sum_{i=1}^L \sum_{k=1}^N \tilde{H}_{ik}.$$

65 As mentioned in Hoffman *et al.* (2014), the estimator in eqn 6 and thus also the one in eqn 7 are subject  
66 to the assumption, that for each pair  $(i, j)$ , the term  $\frac{m_{ij}-m_i \cdot m_j}{(1-m_i)(1-m_j)}$  can be approximated by the average  
67 over all pairs of loci. Thus,  $\hat{g}_2$  only serves as an estimator for  $g_2$ , if in the underlying data set, the missing  
68 values between pairs do not differ greatly in frequency (which could potentially occur if data quality is  
69 very poor for certain individuals).

70 **References**

- 71 David, P., Pujol, B., Viard, F., Castella, V. & Goudet, J. (2007) Reliable selfing rate estimates from imperfect population  
72 genetic data. *Molecular Ecology*, **16**, 2474–2487.
- 73 Hardy, O.J. (2015) Population genetics of autopolyploids under a mixed mating model and the estimation of selfing rate.  
74 *Molecular ecology resources*.
- 75 Hoffman, J.I., Simpson, F., David, P., Rijks, J.M., Kuiken, T., Thorne, M.A.S., Lacy, R.C. & Dasmahapatra, K.K. (2014)  
76 High-throughput sequencing reveals inbreeding depression in a natural population. *Proceedings of the National Academy  
77 of Sciences*, **111**, 3775–3780.
- 78 Weir, B. & Cockerham, C.C. (1973) Mixed self and random mating at two loci. *Genetical research*, **21**, 247–262.