

Review of de Meeus and Nous (2023): *A new and almost perfectly accurate approximation of the effective population size of dioecious populations: comparisons with former estimates and detailed proofs*

This manuscript makes three contributions to the theoretical population genetics of dioecious populations. First, the authors derive several novel expressions for the eigenvalue effective population size of a randomly-mating population conforming to a two-sex version of the Wright-Fisher model. They begin by deriving a system of recurrence equations for the probability of identity-by-descent between pairs of alleles taken either from a single randomly sampled individual or from two distinct randomly sampled individuals. They then show that the equilibrium effective population size can be calculated from the leading eigenvalue of this recurrence, which allows them to derive an exact expression for N_e in terms of the numbers of females N_f and males N_m (their equation 13). They also give a somewhat simpler but very accurate approximation for N_e (their equation 15). Both the exact result and its approximation are compared with several alternative expressions for the eigenvalue effective population size of a dioecious population and shown to be more accurate, especially at very low population sizes (Figure 1). Secondly, the authors use these results to formulate two novel estimators of N_e based on Wright's F_{IS} (equations (19) and (21)) and they show that these outperform several existing F_{IS} -based estimators of N_e when F_{IS} is equal to its expected value (Figure 2). On the other hand, when F_{IS} is estimated from sequence data (as would usually be the case), simulations suggest that all of the various estimators of N_e considered in this paper are significantly biased. The third contribution comes in the form of an extended series of appendices in which the authors provide detailed derivations of a number of expressions for the inbreeding, eigenvalue and coalescent effective population sizes that have been suggested by other researchers, sometimes without explicit derivation.

All three of these contributions have some value. In particular, one strength of the manuscript is that the authors provide detailed, step-by-step derivations of all of the main results, both new and old, making it easier for the reader to follow the assumptions and algebra leading to the various equations that appear in the text. Nonetheless, I found some of the authors' claims about effective population size unclear or questionable. These and other concerns and suggestions are discussed below.

(1) The most pressing question that I have concerns the scope and robustness of the authors' main results. Genetic drift can be influenced by several factors, including (a) the reproductive system, (b) within-generation reproductive variance, (c) life history (e.g., overlapping vs. non-overlapping generations; iteroparity vs. semelparity), (d) demographic history (e.g., changes in population size and/or sex ratio), (e) population structure, and (f) selection at linked sites (background selection and selective sweeps). This paper focuses on the impact of dioecy and sex ratio on genetic drift but it largely ignores all of the other factors mentioned above, leaving me wondering about the biological relevance and validity of the authors' results, especially equations (13), (15), (19) and (21). Are there any real populations that closely conform to the two-sex Wright-Fisher model used to derive these results and, if not, what are we actually estimating if we apply equations (19) and (21) to genetic data?

To be specific, I would be interested in seeing how the results given in these four equations are affected by the following complications. First, how do these expressions change if the within-generation reproductive variance (which could differ between males and females) differs from that obtained under the multinomial sampling assumed by the Wright-Fisher model? The authors do refer to the possibility of defining effective numbers of breeders, but it isn't obvious to me that their results on the eigenvalue effective population size will remain valid if N_f and N_m are simply replaced by effective numbers of

female and male breeders, especially if male and female reproductive success are correlated. Perhaps a dioecious version of the exchangeable Cannings model could be used to investigate this complication?

A second complication that I think merits attention is the impact of population size and sex ratio fluctuations on these results. The authors observe that the discrepancies between their expressions for the eigenvalue effective population size and those obtained by earlier authors are proportionately greatest when the population size is small. However, small populations are also probably more strongly affected by demographic stochasticity, which can lead to proportionately larger fluctuations both in total population size and in sex ratio. The expressions given in equations (13) - (21) were derived by assuming that the population is at equilibrium (in some sense), which likely requires a time average over multiple generations, during which the numbers of adult females and males may fluctuate. How will this impact the formulas given in equation (13) and (15), and do the estimators derived in equations (19) and (21) remain valid if N_f and N_m are replaced by their harmonic means or some other appropriate averages? Alternatively, is it possible to derive comparable results for the single generation eigenvalue effective population size (which can then fluctuate across generations) without insisting on equilibrium? (For example, we can define a time-dependent coalescent effective population size in terms of the instantaneous pairwise coalescent rate, which can then be estimated using Bayesian skyline estimators.)

If it isn't possible to perform additional theoretical or simulation-based studies of these complications within the scope of the current manuscript, then I think that the authors should at least acknowledge that the scope and applicability of their results may be limited in practice.

(2) As the authors acknowledge in their introduction, there are several formal definitions of effective population size which do not coincide in general and so it is somewhat misleading to speak of *the* effective population size without specifying which concept is in use. Although it may be acceptable to use the wording 'the effective population size' where the concept is clear by context, I think that this language should be avoided otherwise. For this reason, I would encourage the authors to change the wording of the title and the abstract so that they explicitly refer to the eigenvalue effective population size that is the main focus of this manuscript.

(3) I would encourage the authors to include a more detailed and explicit description of each of the models being studied in the paper. For example, I think that the subsection titled 'The general model of a dioecious pangamic population' should begin with a detailed description of the model that is used to derive equations (7) and (8). This would include the fact that (i) we are considering a diploid locus; (ii) that the numbers of adult females and adult males participating in reproduction is constant from generation to generation; (iii) that the genotype of each individual alive in generation $t + 1$ is determined by independently sampling a single allele uniformly at random and with replacement from the N_f adult females alive in generation t and then doing the same from the N_m males alive in generation t ; (iv) that the maternal and paternal alleles are sampled independently, etc. I think that detailed, explicit descriptions of the biological models used to derive the theoretical results given in the paper are at least as important as the detailed, explicit descriptions that the authors give of the algebraic transformations that they apply to these results.

(4) The authors note that their estimates of N_e are negative and therefore not biologically meaningful whenever F_{IS} is estimated to be positive. To address this problem, they recommend excluding loci at which F_{IS} is estimated to be positive when estimating N_e . This strikes me as *ad hoc* and statistically unsound. Perhaps a better approach would be to estimate N_e directly from the observed

and expected heterozygosities using either maximum likelihood or Bayesian estimation. For example, for the small populations that seem to be most relevant to the concerns of this manuscript, it may not be too computationally challenging to use either approximate Bayesian computation or MCMC to estimate the posterior distribution of (N_e, N_m) (and any nuisance parameters such as mutation rates). One could then easily estimate the posterior distribution of N_e using equation (13). What is the value of introducing yet another statistically questionable method-of-moments type estimator when one can perform maximum likelihood or Bayesian analysis?

(5) It might be useful to mention that the F_{IS} -values used to estimate N_e should be estimated in reproductively mature individuals. In small populations these statistics could fluctuate significantly within generations due to random survival from birth to reproductive maturity.

(6) The equations presented in lines 1185-1211 in Appendix 6 are not quite correct. In particular, beginning on line 1185, there is a series of identities which contain an infinite series ($t = 1, \dots, \infty$) on the left-hand side and a finite series (also over t) on the right-hand side. To fix this, you need to take limits (as $t \rightarrow \infty$) on the right-hand side. Provided that all of the eigenvalues λ_i have modulus less than 1, these limits will exist and be finite. Alternatively, the desired result can be derived as follows. Defining

$$S_i \equiv \sum_{t=1}^{\infty} t \lambda_i^{t-1}$$

we can write

$$\begin{aligned} \lambda S_i &= \sum_{t=1}^{\infty} t \lambda_i^t \\ &= \sum_{t=1}^{\infty} (t+1-1) \lambda_i^t \\ &= \sum_{t=1}^{\infty} (t+1) \lambda_i^t - \sum_{t=1}^{\infty} \lambda_i^t \\ &= S_i - 1 - \frac{\lambda_i}{1 - \lambda_i}. \end{aligned}$$

Here we have used the fact that the second term appearing in the third line is a geometric series, with $\sum_{t=1}^{\infty} \lambda_i^t = \lambda_i / (1 - \lambda_i)$. We can solve for S_i , obtaining

$$S_i = \frac{1}{(1 - \lambda_i)^2}$$

We can also derive this result by noticing that S_i is the term-by-term derivative with respect to t of a geometric series, i.e.,

$$\begin{aligned}
S_i &= \sum_{t=0}^{\infty} \frac{d}{dt} \lambda_i^t \\
&= \frac{d}{dt} \left(\sum_{t=0}^{\infty} \lambda_i^t \right) \\
&= \frac{d}{dt} \left(\frac{1}{1 - \lambda_i} \right) \\
&= \frac{1}{(1 - \lambda_i)^2}.
\end{aligned}$$

Jay Taylor (29 March 2023)

Minor corrections:

line 55: “Many species have separate sexes. Several authors have investigated the impact that dioecy and sex ratio have on effective population size.”

line 58: “leads to an approximation that appears closer”

line 60: “We also propose another estimator of”

line 68: “The effective population size of a dioecious population has been defined in different ways.”

line 75: “for the eigenvalue effective population size N_e ”

line 86: “Another consequence is that”

line 99: No comma is required after Balloux (2004)

line 121: Why is it necessary to assume that the number of matings is very large?

line 235: “The reasons for this discrepancy between these two sets of equations are unclear due to the lack of details in Balloux’ paper.”

line 245: “This bias is very small when $N_e > 10$.”

Figure 1: Please specify the actual sex ratio(s) used to obtain the results shown in the figure on the left (uneven sex ratio). Also, why does purple curve have such a jagged appearance in this figure and why, in particular, does it appear to jump back and forth when $N_e > 7$?

line 309: What does very big mean in this setting? $N_e > 20$? $N_e > 100$?

line 311: “in the opposite direction”

Figure 2 legend: add a space before the equation numbers, e.g., Eq 5.

line 380: “It is worth recalling that the F_{IS} -based estimate given in Equation (21) assumes an even sex ratio.”

line 396: “large variances”

line 397: “will have a large impact on F_{IS} -based estimates of N_e .”

line 408: “In addition to the fact that it is generally preferable to work with the most accurate equation, these results are likely to be especially pertinent for certain types of biological systems that are able to persist for extended periods despite having very small effective population sizes.”

line 414: “a female enters a brood cell, which she caps, where she feeds on the bee larva and then gives birth to a haploid male, which later mates with its”

line 422: “may not be rare in dioecious parasitic”

line 439: “different kinds of loci”

line 453: “were positive and therefore could not be used to estimate N_e .”

line 480: “Simulation studies could be used to identify an estimator that more accurately approximates the eigenvalue effective population size of genotyped populations.” As discussed above, I would argue that there is no such thing as the ‘real effective population size’.

line 487: “may lead to underestimates.”

line 493: I don’t know what you mean by “important” population sizes.

line 721: “To compute the inverse of a 3x3 matrix”

line 788: “an infinite collection of eigenvalue that all satisfy”

line 802: “Computing matrix powers is difficult except for diagonal matrices.”

line 843: “Consequently, we can use equation (A3-3) to calculate the power of any diagonalizable square matrix.”

line 845: “We can now derive some other properties of eigenvalue-eigenvector pairs (eigenpairs).”

line 874: You should remark that equation (A3-6) only applies when S is invertible.

line 971: Castle-Weinberg (capitalization)

Appendix 5: Perhaps it would be clearer to write H_{di} and H_{mon} in place of H_{obs} and H_{exp} . To me, the notation H_{obs} and H_{exp} suggests that we are working with the observed and the expected heterozygosity, which is not the case here. Instead, H_{obs} is the expected heterozygosity in a dioecious

population, while H_{exp} is the expected heterozygosity in a monoecious population.

line 1024: “The variance of a difference between two uncorrelated (e.g., independent) random variables”

line 1099: “the probabilities that the same allele is sampled twice, either in one individual or in two distinct individuals from the same population. Let u be the mutation rate per generation.”

line 1100: Do equations (A6-1) assume the infinite allele model?

Equation (A6-3): There is an extra 1 in the expression shown in $A_{2,2}$.

line 1131: “the second term in these equations will increase with t , albeit at a diminishing rate, while the first term will decrease with t .”

line 1153: “the probability of pairwise coalescence”

line 156: To be consistent you should capitalize the entries in the vector shown in (A6-8).

lines 1158-1160: Please clarify what you mean by the coalescent probabilities here. In most scenarios commonly considered in biology, the probability that two lineages will eventually coalesce in the past is 1.

line 1181: Since A is a 2x2 matrix, shouldn't i be restricted to the values 1 and 2?

line 1368: This identity (the chain rule) is incorrect as written. It should be replaced by the expression $(f \circ g)'(x) = f'(g(x))g'(x)$.

line 1372: Replace derivable by the word differentiable.

line 1461: “Assuming that $u \ll 1$ ”

line 1468: I disagree that the effective number of breeders is in general equal to the exact number of reproducing adults. As discussed above, we need to account for within-generation variance in reproductive success which depends not only on the number of adults that reproduce but also on the number of offspring that survive to adulthood.

line 1493: “At equilibrium, the vector of genetic identities satisfies the equation”

lines 1509-1525: I think that this paragraph needs to be reworded, e.g., “The recursion for the identity between individuals can be determined by conditioning on the ancestry of the sampled pair in the previous generation. One possibility is that the two sampled individuals are sibs, i.e., they share the same parents, which is true with probability $1/(N/2)$. In this case, with probability $1/2$, the two alleles will have come from the same parent, in which case they are equally likely to be derived from a single parental allele or from both parental alleles. In the former case, the sampled alleles are necessarily IBD, whereas in the latter case, the probability that they are IBD is $Q_{I(t-1)}$. Alternatively, with probability $1/2$, each sampled allele may have come from a different parent, in which case the probability that they are IBD is $Q_{S(t-1)}$. The second possibility, which has probability $1 - 1/(N/2)$,

is that the two sampled individuals are not sibs, in which case the probability that the sampled alleles are IBD is $Q_{S(t-1)}$.