

Counting parental contribution - how large sample size makes strong selection weak

Ivan Krukov, Simon Gravel

Abstract

Neutral models of genetic diversity tend to be easier to analyze compared to models including selection. Because lineages are exchangeable in the neutral Wright-Fisher model, for example, the number of lineages that are relevant to the ancestry of a sample at a single locus can only decrease as we go back in time. As a consequence, useful recursion equations can be derived for patterns of polymorphism. By contrast, under negative selection, the number of relevant lineages can increase as we go back in time, and the equivalent recursion equations do not close. Given a large enough sample size, however, the reduction in the number of lineages due to genetic drift is larger than the increase in the number of lineages due to natural selection, and the number of relevant lineages is unlikely to increase. We use this observation to derive asymptotically closed recursion equations for the distribution of allele frequencies. We show that this approach is accurate under strong drift and strong natural selection. We derive several asymptotic results to understand when the sample size is sufficiently large to overcome the influence of selection.

1. Introduction

The calculation of the allele frequency spectrum (*AFS*) is an important tool for the inference of demographic histories and other population genetic parameters. In the absence of selection, the number of parental lineages that contribute to the sample decreases back in time due to coalescent events. This means that the equations are closed with respect to the sample size under neutrality. This has paved the way for moment-based recursions of the allele frequencies (Kimura and Crow, 1964; Ewens, 1972; Donnelly and Kurtz, 1999). Recently, a number of successful methods, including (Gutenkunst et al., 2009; Jouganous et al., 2017; Kamm et al., 2017), have become popular for this purpose. At their core, these methods describe the time-evolution of the *AFS*. The goal of these

approaches is to obtain the probability of observing a given number of derived alleles in a finite sample conditional on the state of the parental lineages.

Despite many successes, considering selection has been problematic within this framework. Under negative selection, the number of parental lineages that contribute to a sample can be larger than the sample size itself, due to selective death events. As a consequence, the equation loses the closure property (Jouganous et al., 2017).

An extension to the Kingman’s coalescent (Kingman, 1982), the ancestral selection graph (*ASG*), (Krone and Neuhauser, 1997) considers the ancestry of a sample in the presence of selection. The *ASG* framework can be used to study the ancestry of highly deleterious alleles (Wakeley, 2009) under certain assumptions. Another possible resolution is to use an uncontrolled jackknife approximation to add extra lineages - the method proposed in (Jouganous et al., 2017).

Unlike the *ASG*, the present method does not need to assume an infinitely large population size, and also explicitly tracks multiple coalescent events, which is important with large sample sizes (Bhaskar et al., 2014). Our approach can also be combined with the jackknife. Unlike the jackknife, our method allows derivation of bounds on the performance of the approximation.

2. Background

Consider the behavior of a single biallelic locus in a haploid Wright-Fisher model with a population of size N . We consider a sample of size n lineages from within the population. We want to know how many parental lineages r have contributed to the present sample from one generation in the past. Going back in time, the random variable \mathcal{R} is the number of lineages that contribute to the present day sample. Figure 1 shows examples of the different models that we consider below – standard coalescent (1A), coalescent with multiple merges (1B), and coalescent with multiple merges and selection (1C).

First, consider a recursion describing the evolution of the expected allele frequency spectrum (*AFS*) $\Phi_n^{(t)}$, in the standard coalescent without selection (Figure 1A). At generation t , the sample consists of n lineages, and at $t - 1$ there are r parental lineages. In what follows, we closely follow the exposition of Jouganous et al. (2017). The standard coalescent model allows at most one event per generation, so we need to consider only two cases: namely $r = n$ if no coalescent occurs, and $r = n - 1$ if a single event took place.

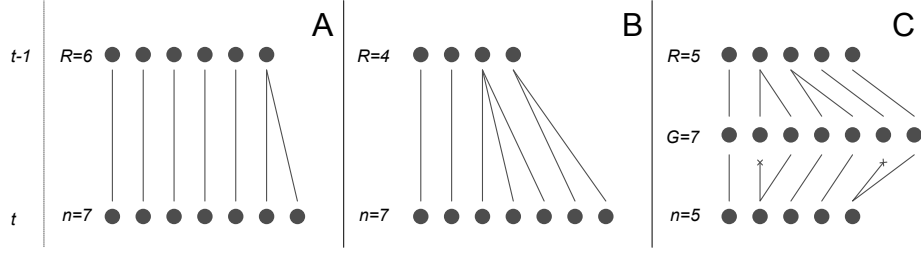


Figure 1: Realizations of sampling models, showing $\mathcal{R} = r$ parental lineages at $t - 1$ contributing to n lineages in present generation t . **A** Standard coalescent - at most 1 coalescent event per generation; $r \in [n - 1, n]$. **B** Coalescent with multiple merges - parents 3 and 4 have 3 and 2 offspring, respectively; $r \in [1, n]$ **C** Including selection - each parent produces a random number of gametes (\mathcal{G}), which then may survive to produce offspring; $r \in [1, n]$, $g \in [n, 2n]$

Without a coalescent event, r parental lineages are sampled randomly into the present generation, so the allele frequencies remain the same: $\Phi_n^{(t)} = \Phi_n^{(t-1)}$. With a single coalescent event, the contribution comes from $r = (n - 1)$ parental lineages: $\Phi_n^{(t)} = \mathcal{D}\Phi_{n-1}^{(t-1)}$, where \mathcal{D} is a sparse $n \times (n - 1)$ matrix describing the effect of drift. We do not use \mathcal{D} in this work (but see (Jouganous et al., 2017)), and instead construct a related matrix in section 3.1. Combining the terms without and with the coalescent event, we have:

$$\Phi_n^{(t)} = \Phi_n^{(t-1)} + \mathcal{D}\Phi_{n-1}^{(t-1)} \quad (1)$$

In other words, the *AFS* of size n at time t can be obtained if we know the *AFS* of sample sizes n and $n - 1$ in the previous generation. This gives rise to a recursion formula for the frequency spectrum that can be solved efficiently (Jouganous et al., 2017). Here we would like to generalize such an approach to cases including natural selection and large sample sizes.

We first incorporate the effects of large sample size. As previously shown in (Bhaskar et al. (2014) and Nelson et al. (2019)), the coalescent approximation may not be adequate in this setting, since multiple coalescent events can take place within a single generation (Figure 1B). With a slight abuse of notation we denote \mathcal{D}_i as the i^{th} -order drift matrix in which i lineages are lost due to genetic drift. The dimensionality of \mathcal{D}_i is $n \times (n - i)$. For example, \mathcal{D}_2 includes both three-way coalescent and double two-way coalescent. In section 3.1, we demonstrate an efficient dynamic programming algorithm to exhaustively enumerate all the events for a drift-only model.

With multiple coalescent events per generation, (1) becomes:

$$\Phi_n^{(t)} = \Phi_n^{(t-1)} + \sum_{i=1}^n \mathcal{D}_i \Phi_{n-i}^{(t-1)} \quad (2)$$

The equation (2) is still closed in terms of the sample size, since Φ_n^t only depends on $r = (n-i) < n$ parental lineages.

If we now consider selective death events, we must also account for lineages that were not transmitted due to selection. We use the model shown in figure 1C to describe this. Each generation, a random number of $\mathcal{R} = r$ parental lineages produce a large number of gametes, $\mathcal{G} = g$. Then n individuals are formed by randomly sampling gametes, without replacement. The probability of successfully sampling a particular gamete is $1 - xs$, where x is the frequency of the derived allele in the parental generation. In the case of a selective death a new lineage is re-drawn with probability xs . This sampling scheme allow us to consider drift ($r \rightarrow g$) and selection ($g \rightarrow n$) within the same generation as distinct processes. [IK: Should we use \$s\$ or \$s/1 + s\$?](#)

For example, in case of a single selective death event, we have $\Phi_n^{(t)} = \mathcal{S} \Phi_{n+1}^{(t-1)}$, with selection matrix \mathcal{S} . Multiple selection events are possible per generation, but we restrict our attention to the case where each lineage experiences at most one selective death event. This still allows us to consider strong selection, with at most $r \leq 2n$ parental lineages contributing:

$$\Phi_n^{(t)} = \Phi_n^{(t-1)} + \sum_{j=n}^{2n} \sum_{i=1}^n \mathcal{S}_j \mathcal{D}_i \Phi_{n-i+j}^{(t-1)} \quad (3)$$

The closure no longer holds for (3), as up to $2n$ gametes can contribute to the sample. However, the effect of a large sample size counteracts the additional lineages needed due to selection. The opposite effects of drift and selection on the number of lineages relevant to a sample are particularly clear in the context of the size of the ancestral selection graph (*ASG*): in which the number of lineages relevant to a sample can be described as a birth-death process (Krone and Neuhauser, 1997; Wakeley, 2009):

$$n \rightarrow \begin{cases} n+1 & \text{at rate } Ns \frac{n}{2} & \text{(selection)} \\ n-1 & \text{at rate } \frac{n(n-1)}{2} & \text{(coalescence)} \end{cases} \quad (4)$$

The coalescence (drift) term is quadratic with respect to the sample size, while the selection term is linear. The rate of coalescence is higher than the rate of selective deaths if the number of lineages $n > Ns + 1$.

Our goal here is to define recursions generalizing equation (1). For the equations to be closed, we need to ensure that the rate of coalescence is large enough that it overcomes the rate of selection not only on average, but *almost always*.

The rest of the paper is organized into two sections. In the first section, we construct a Markovian recursion to track the number of derived lineages in a large sample from a Wright-Fisher model, similar to (Jouganous et al., 2017; Kamm et al., 2017). In this, we fully account for multiple coalescent events per generation, and show that we restore closure with increasing sample size. In the second part, we derive a number of asymptotic results to get a better understanding of the process. We construct an exact probability distribution for the number of contributing parental lineages, together with several approximations. Importantly, we derive a normal approximation that allows us to calculate a quantile of the sample size where the system is approximately closed.

3. Results

3.1. Markov process construction

We first define a recursion equation for the distribution of allele frequency in a sample of size n from a haploid Wright-Fisher population of size N . Given that the sample in parental generation has j copies of the derived allele, we seek to calculate the probability that the sample will contain i derived copies in the following generation.

In the neutral case, this transition probability can be calculated if we know the transition probabilities in the smaller samples of size $n' \in [1, n - 1]$ (similar to (2)) (Bhaskar et al., 2014). We do not construct the intermediate matrices \mathcal{D}_j of (2) explicitly, but instead calculate a matrix $P((j, n) \rightarrow (i, n))$, giving the probabilities of transitioning from j to i derived alleles in sample of size n within one generation. In brief, we can construct the transition probability matrix $P((j, n) \rightarrow (i, n))$ if we know the matrix for $n - 1$, $P((\cdot, n - 1) \rightarrow (\cdot, n - 1))$, and then use the conditional probabilities for each type of an event. Starting from the base cases for a sample size of 1, we build up a set of square transition probability matrices. By using a dynamic programming algorithm, we can achieve reasonable performance for realistic sample sizes. The full derivation is shown in appendix A.

The case of selection is slightly more complicated, since now we need to consider a larger set of states that can lead to transition from j to i copies in a sample size of n . In particular, there can be a large number of gametes (fig. 1C) that can contribute from n when there is no selection, to a potentially infinite number under strong negative selection. In our calculation, we only consider up to one selective death event per lineage, so the number of contributing gametes is between n and $2n$.

Note that this is analogous to the outer summation boundaries in eq. (3). Again, we do not calculate the matrix \mathcal{S}_i directly, but rather construct a transition probability $Q((j, n) \rightarrow (i, n))$. We can define these transitions recursively if we know the transitions of $Q((\cdot, m) \rightarrow (\cdot, n-1))$. Note that with selection we can have $m \in [1, 2n]$ lineages contribute, whereas we only needed $n' \in [1, n-1]$ in the neutral case.

To retain the closure property under selection, the Markov process needs to take $2n$ lineages in the parental generation to n lineages in the present. However, such rectangular transition probability matrix is not suited for our purposes, since we want to describe the behavior of a sample with constant size n . Instead, we only calculate the *truncated* transition probabilities for a $n \times n$ matrix Q . This means that under very strong selection, some transitions will be unaccounted for. However, since the total sum of transition probabilities sums to 1, we can easily calculate the total missing probabilities. As we show in the rest of this work, this missing probability tends to 0 as sample size n increases.

The construction of the full and truncated transition probability matrices is implemented via a dynamic programming approach similar to the neutral case, and is described fully in appendix B. Because we need to account for additional lineages in the selection case, the calculation time is of the order of $O(n^4)$, while it is only $O(n^3)$ for the neutral case. The increase in complexity makes this approach less suitable for large sample size, but we derive several approximations in the following sections.

3.2. Calculation of allele frequency spectra

Once the truncated matrix Q is constructed, it can be used to calculate the allele frequency spectrum. For the infinite sites model at equilibrium, we can calculate the *AFS* Φ as a solution to a linear system:

$$\Phi = \Phi Q + n\mu e_1 \quad (5)$$

where μ is the per-site mutation rate, and e_1 is the first column of the identity matrix of size n . Figure 2 shows the comparison of the *AFS* calculated from Equation (5), the diffusion approximation (Ewens, 2004, eq. 9.23), and the calculation performed in **Moments** (Jouganous et al., 2017). Panel A shows a comparison at $Ns = 50$, with the population size ($N = 2000$), which is substantially larger than the sample size ($n = 200$). There is a small deviation between the approaches at large allele frequencies. At stronger selection coefficients, **Moments** suffers from numerical instability, while the diffusion approximation performs well (not shown).

If the sample size is the same as the population size ($n = N = 200$) (Fig. 2B), the diffusion approximation and **Moments** perform poorly, while our approach remains stable. This is expected, since the diffusion framework does not perform well if multiple coalescent events contribute. Furthermore, if our sample size is the entire population, we expect recursion equations to be closed. To confirm this, we compare our result to the *AFS* calculated from a whole-population haploid Wright-Fisher model, with $N = 200$. (Fig. 2B) shows that our calculation is close to the full Wright-Fisher model. The discrepancy between the curves is due to a difference in the way the selection coefficients are calculated.

3.3. Closure properties

To show the closure properties of Q , we can calculate the total probability that more than n parental lineages contribute to the sample of a given size. By construction, the sum of rows of Q should correspond to the total probability mass that included configurations contribute (Fig. ??). Thus, the probability that some number of configurations are unaccounted for, with j derived alleles in the parental sample, is given by $1 - \sum_{i=0}^n Q_{i,j}$. This probability depends on the number of derived alleles carried by the parental sample: the more derived alleles, the higher the likelihood of a selective event. Figure 3 shows the probability of missing configurations in a sample size of $n = 200$ in the worst-case scenario, with $j = 200$ derived lineages.

Since the expected number of drift events increases quadratically and the number of selective events increases only linearly, the probability that we need additional lineages decreases rapidly with sample sizes.

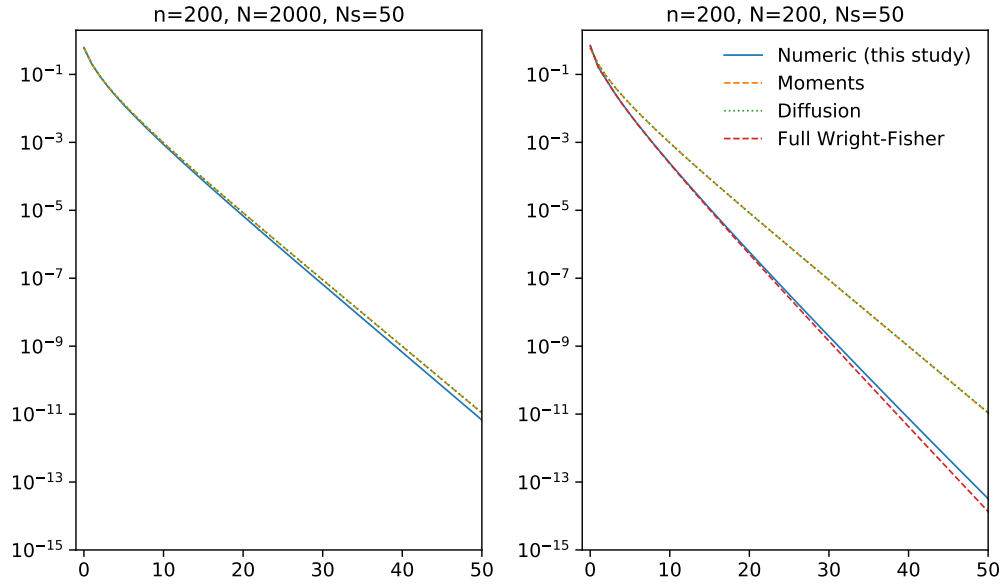


Figure 2: Normalized allele frequency spectra in a sample of size $n = 200$, for highly deleterious alleles ($Ns = -50$). (A) shows the frequency spectrum in a sample from a large population ($N = 2000$), (B) in a small population ($N = 200$). Both panels are truncated at 10^{-15} , to show only moderately high allele frequencies.

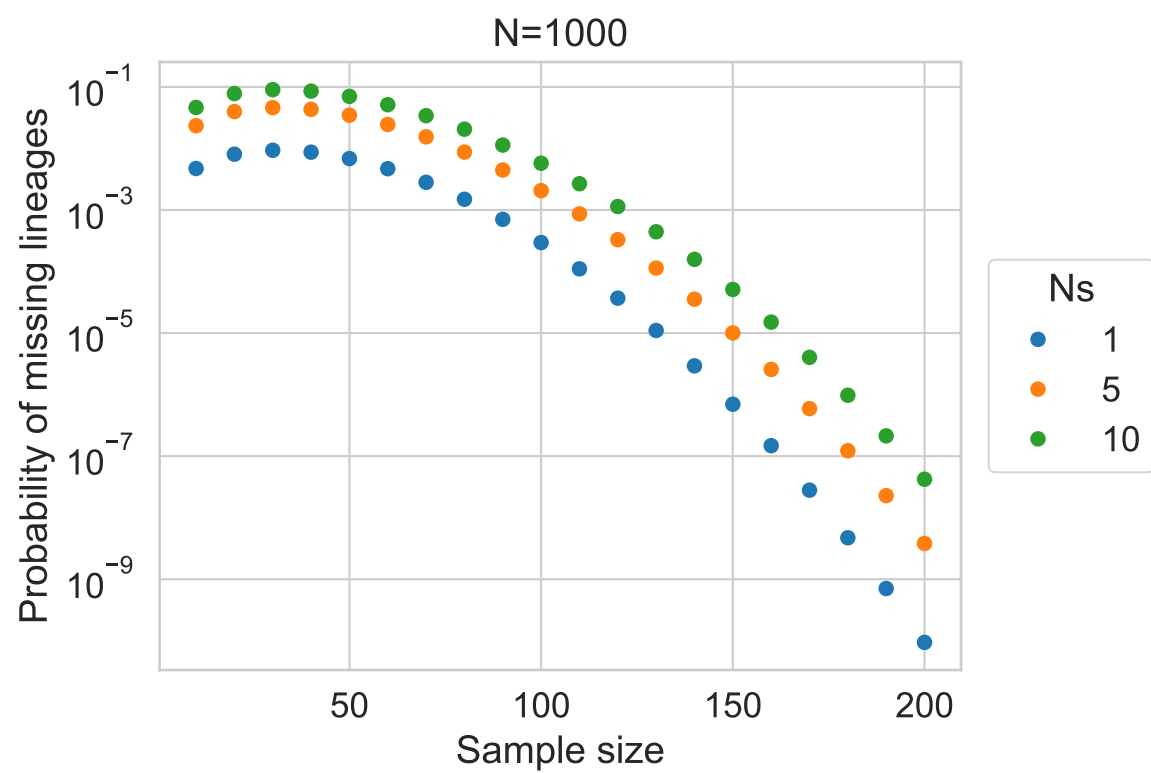


Figure 3: Probability that unaccounted lineages contribute to the transition probabilities. The probabilities are calculated as 1 minus the sum of probabilities for the state where every allele is derived. [IK: Need to keep N consistent](#)

4. Asymptotic closure properties

We now want to determine what sample size is sufficient so that the number of coalescent events due to drift is almost always larger than the number of selection events, such that the system remains closed (3). We derive several approximations to the model proposed in the first section, in order to get a better understanding of this behavior.

As a first order approximation, we consider the mean number of lineages that contribute by selection or drift. Then, we construct a full probability distribution of the number of contributing lineages one generations into the past. Finally, we propose a normal approximation to this distribution, in order to derive a simple quantile function for the number of used lineages.

The upper bound on the number of lineages used is of particular interest - it represents the worst case scenario in terms of extra lineages used by selection. Since the maximum number of lineages will be resampled when all lineages are derived, we will usually assume $x = 1$ in the following calculations. This also allows us to treat the lineages as exchangeable (Wakeley, 2009). Note that in section 3.1, we did not assume exchangeability of lineages, which led to a considerably more complex formulation.

For a given sample size, the probability that r parents have contributed is:

$$Pr(\mathcal{R} = r|n) = \sum_{\mathcal{R}} Pr(\mathcal{R} = r|\mathcal{G} = g)Pr(\mathcal{G} = g|n) \quad (6)$$

Where \mathcal{R} and \mathcal{G} are random variables denoting the number of contributing parents and gametes, respectively (Fig. 1C).

Before deriving the distribution formally, we seek to obtain several approximate results.

4.1. Expected number of lineages used

First, we seek an approximate expression for the expectation of the total number of lineages used. This can be approximated as the sum of expectations of the number of lineages sampled under drift ($E[\mathcal{R}]$) plus the number of lineages rejected by selection ($E[\mathcal{G} - n]$). The expected number of contributing lineages is then $E[\mathcal{R} + \mathcal{G} - n]$. While this is an imprecise approximation, as is essentially assumes that selection and drift are independent, it allows us to derive several closed form results. The number of parents that contribute to n gametes (drift) will be:

$$\hat{E}[\mathcal{R}|n] = N(1 - \left(1 - \frac{1}{N}\right)^n) \quad (7)$$

This expression is not hard to derive intuitively. First, the probability of selecting a particular parent is $\frac{1}{N}$, so the probability of selecting different parents for n individuals is $(1 - \frac{1}{N})^n$. Then one minus this value is the probability that the same parent was picked at least once by any of the n individuals.

For selection, we want to consider the expected number of gametes that are rejected by selection to form a sample size of n , which is the number of extra lineages used by selection. If the probability of rejection is xs , the scheme is described by the negative binomial distribution, where the random variable is the number of failures, given n successes. The expectation of this parameterization of negative binomial is:

$$\hat{E}[\mathcal{G} - n|n] = n \left(\frac{xs}{1 - xs} \right) \quad (8)$$

Then summing the expectations of the two random variables yields:

$$\hat{E}[\mathcal{R}] = \hat{E}[\mathcal{G} - n|n] + \hat{E}[\mathcal{R}|n] \quad (9)$$

$$= N(1 - \left(1 - \frac{1}{N}\right)^n) + n \left(\frac{xs}{1 - xs} \right) \quad (10)$$

$$\underset{N \gg n}{\approx} \frac{nx s}{1 - xs} - \frac{n^2}{2N} \quad (11)$$

The second approximation is made under the assumption that the sample size is much smaller than the population size. We can see that the expected number of lineages sampled will be increased by selection as a linear term. Drift tends to decrease the number of lineages as a quadratic term with respect to the sample size. This is analogous to the results from the ancestral selection graph (Krone and Neuhauser, 1997), eq. (4), but now includes sample size directly.

We now want to ask when the expected number of contributing lineages is less than the sample size:

$$\hat{E}[\mathcal{R} - \mathcal{G} - n] < n$$

$$\frac{nx s}{1 - x s} - \frac{n^2}{2N} < n \quad (12)$$

$$n \geq \frac{2Nxs}{1 - xs}$$

$$\approx 2Nxs \quad (13)$$

This gives a simple expression for the sample size where drift overcomes selection: $n \geq 2Nxs$. Figure 4 shows this for several selection coefficients, assuming the entirety of the sample is derived in a population of $N = 1,000$. The Y axis shows the fraction of contributing parental lineages to the sample size, $\frac{r}{n}$. Above the horizontal line $\frac{r}{n} > 1$, selection dominates. Below, drift reduces the number of used lineages. The intercept of the line with $\frac{r}{n} = 1$ is the critical sample size, which is well-approximated by $2Ns$ (we assume that every lineage is derived, $x = 1$).

Using the same equation, we can track the expected number of used parental lineages back in time, which we denote as n_{t-1} :

$$n_{t-1} = \frac{n_t x s}{1 - x s} - \frac{n_t^2}{2N} \quad (14)$$

We solve this recurrence going back in time 10,000 generations, producing figure 5. The equilibrium point is well-approximated by $2Ns$, shown as a dashed line here. This is the same as solving equation 12 explicitly. Non-withstanding of the starting sample size, we converge to the equilibrium relatively quickly.

4.2. Distribution of number of contributing lineages

We now construct a probability distribution of the number of contributing lineages one generation into the past.

The number of parental lineages used by drift can be modelled by the modified occupancy (Arfwedson) distribution (Wakeley, 2009; O'Neill, 2019; Johnson et al., 2005). This is given by:

$$P(\mathcal{R} = r | \mathcal{G} = g) = \frac{S_2(g, r) N!}{(N - r)! N^g} \quad (15)$$

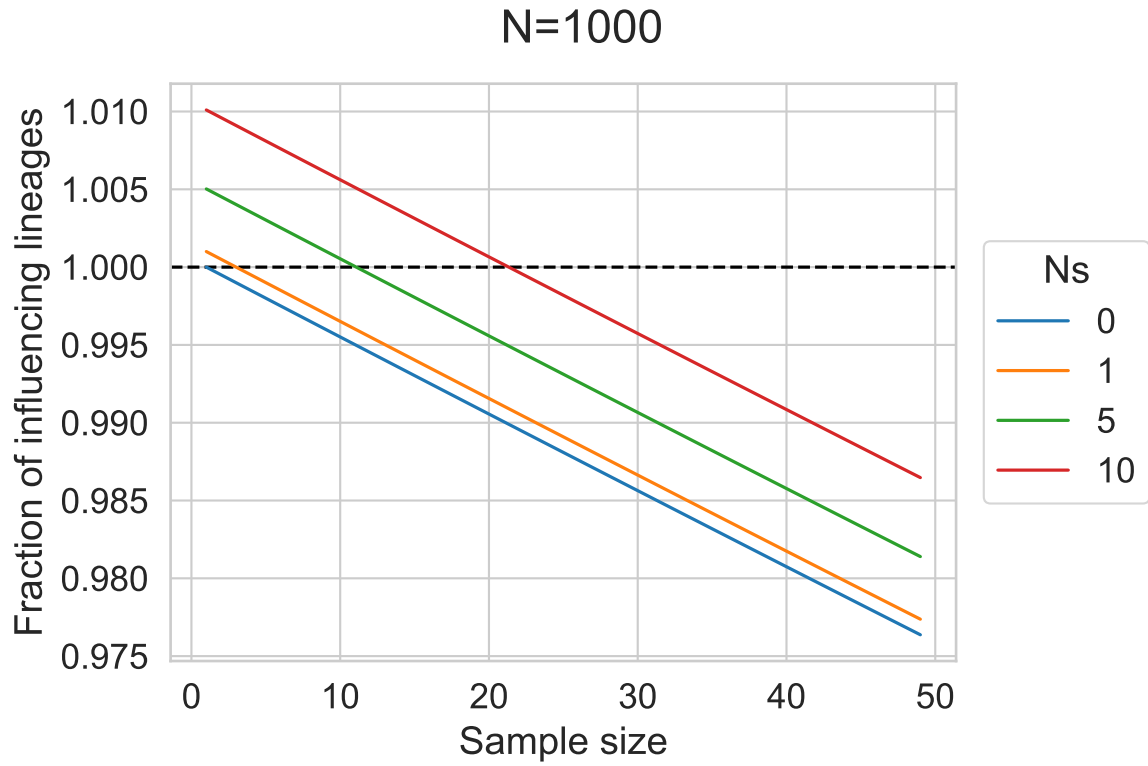


Figure 4: Critical sample size for different selection coefficients. The Y axis shows the fraction of parental lineages over the sample size, $\frac{r}{n}$, each line corresponds to a different selection coefficient. Above $\frac{r}{n} \geq 1$, selection dominates, below – drift. The critical sample size, where the expected number of parental contributing lineages is smaller than the sample size is well-approximated by $2Ns$.

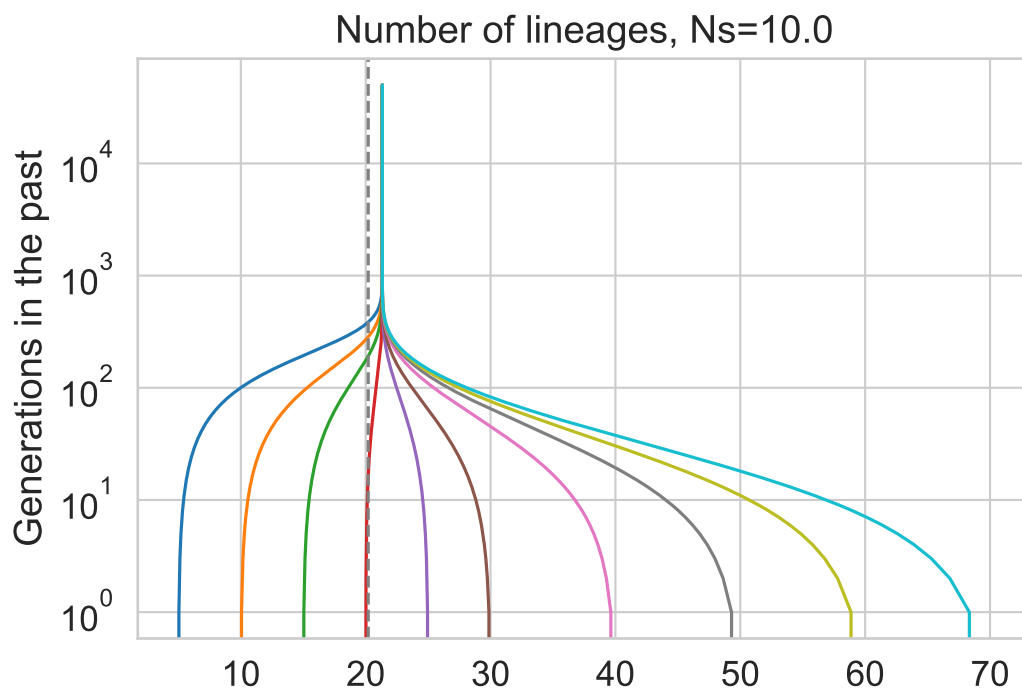


Figure 5: Expected number of contributing parental lineages back in time. Starting at a given sample size, the number of contributions is tracked with align 14. The Y axis shows time on a logarithmic scale, X axis is the sample size. $N = 1000$, $Ns = 10$.

where $S_2(g, r)$ is a Stirling number of the second kind, which is the number of ways to partition g objects into r categories (g gametes produced by r parents). A typical statement of the occupancy distribution is that we have N urns and g colored balls, and we want to know the probability that exactly r of the urns will be occupied (see Johnson et al. (2005) section 10.4 for a thorough treatment). In our case, N is the population size, urns correspond to the parents, colored balls to gametes. Note that the under drift, the number of parents will be smaller or equal to the number of gametes $r \leq g$. The expectation of this distribution is given by equation (7).

The occupancy distribution is not simple to evaluate, but good performance can be achieved by pre-computing a table of reduced occupancy numbers, using the algorithm of O'Neill (2019).

As stated before, the number of lineages sampled under selection is described with a negative binomial distribution. Unlike 8, however, we are now looking for the total number of lineages sampled, not simply the number of failed trials. In this parameterization, the probability of the negative binomial is given by:

$$P(\mathcal{G} = g|n) = \binom{g-1}{n-1} (1-xs)^n (xs)^{g-n} \quad (16)$$

Here, the number of gametes can be larger than the sample size $n \leq g$, if selection is present ($s < 0$).

Combining the two distributions together through 6, we get:

$$Pr(\mathcal{R} = r|n) = \sum_{g=1}^{\infty} \frac{S_2(g, r)N!}{(N-r)!N^g} \binom{g-1}{n-1} (1-xs)^n (xs)^{g-n} \quad (17)$$

Unfortunately, this distribution does not have a simple analytical form. In certain parameter regimes, this can be approximated by the normal distribution (Johnson et al., 2005; O'Neill, 2019), which we describe in the next section.

Figure 6 shows the distribution of the number of contributing parental lineages for several selection coefficients for a sample $n = 20$. In the absence of selection, the distribution has zero probability above $n = 20$, as no extra lineages can be sampled. As the strength of selection is increased, we begin requiring larger number of lineages. At the equilibrium point ($Ns = 10$, (12)), the distribution is symmetric.

We note that at the critical sample size, the probability that we will have a sufficient number of lineages is only 50%. In order to guarantee that drift will out-pace selection, we can calculate

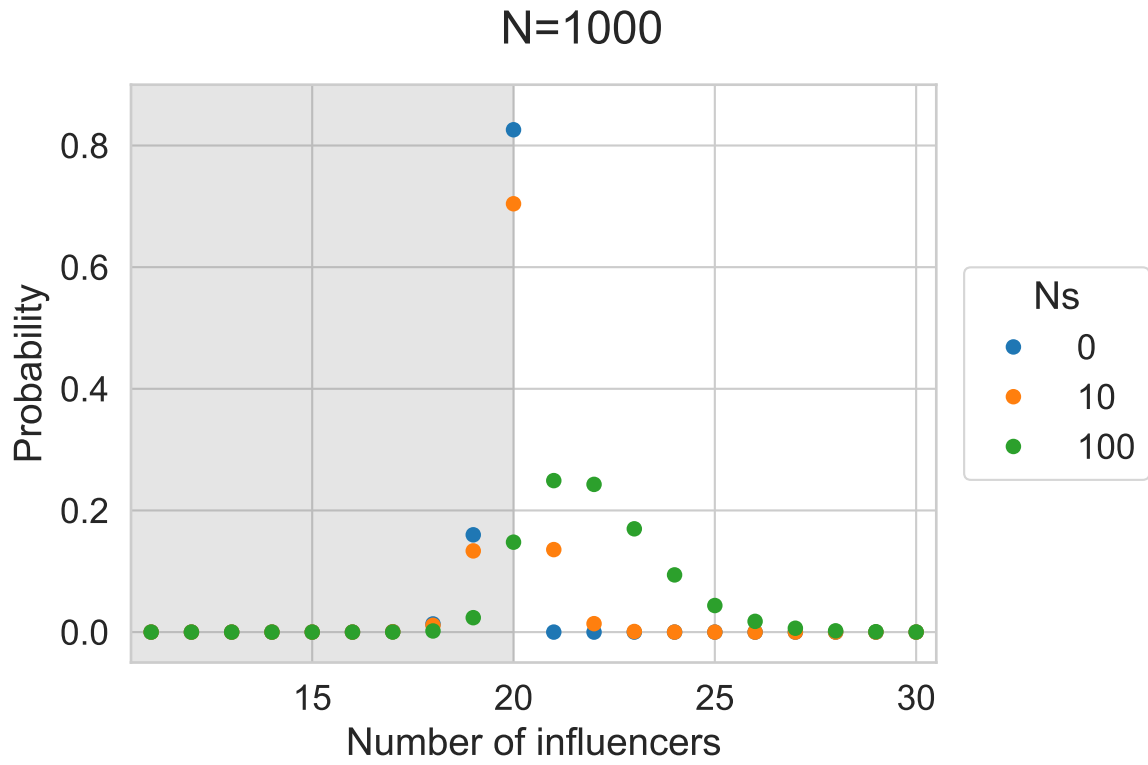


Figure 6: The distribution of the number of parental contributing lineages one generation into the past ($n = 20$, $N = 1000$). Shaded area shows the drift-dominated regime, where the number of lineages is smaller than the sample size.

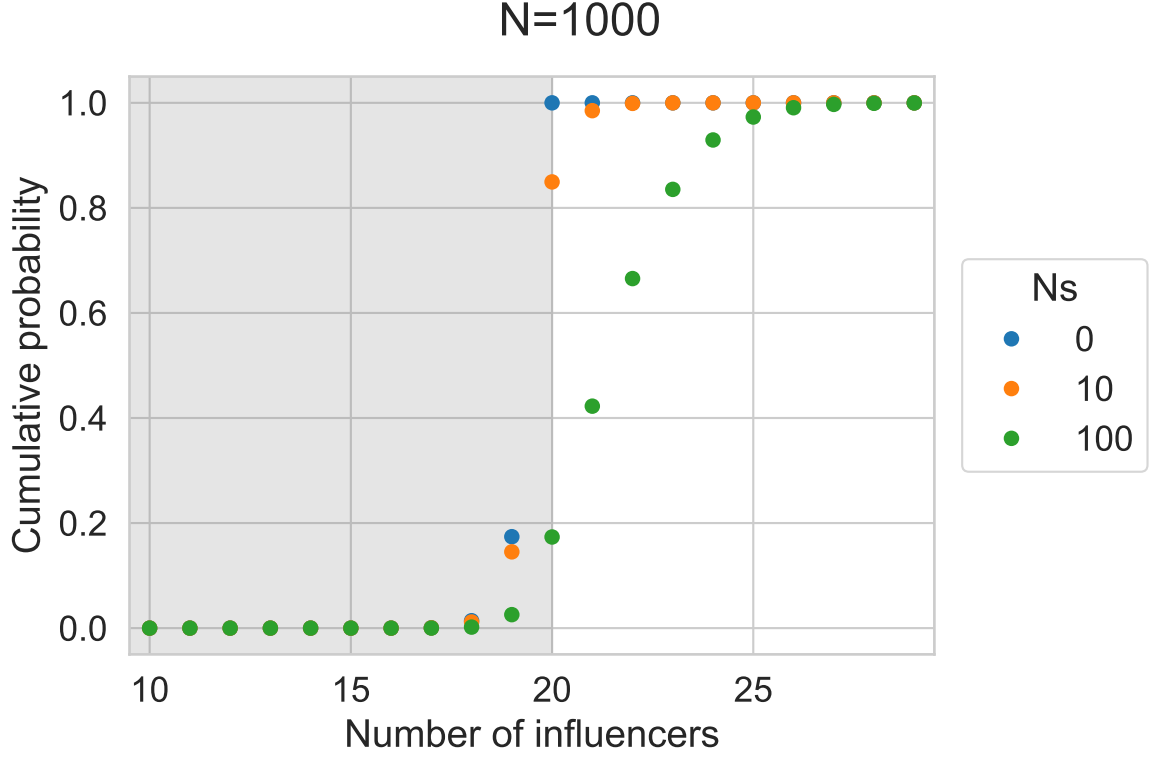


Figure 7: The cumulative distribution of the number of parental contributing lineages one generation into the past ($n = 20$, $N = 1000$). Shaded area shows the drift-dominated regime, where the number of lineages is smaller than the sample size. [IK: This should be a two-panel with the previous figure](#)

the cumulative distribution - Figure 7. This shows that a sample size in which the majority of lineages are accounted for can be substantially larger than the critical sample size of equation (12). To derive a convenient expression, we turn to the normal approximation in the next section.

4.3. Normal approximation

Finally, we can construct a normal approximation to the distribution of the number of contributing lineages. The occupancy distribution is approximated by the normal (O'Neill, 2019) when $n \ll N$. Likewise, the number of failures (eq. (8)) before a given number of successes, can be approximated by the normal distribution. In the case of large population size, as required by the approximation of the occupancy by the normal, we can approximate the total number of contributing lineages as the sum of lineages contributed by the two distributions. The random variable

which is a sum of two normally-distributed random variables is also normal, with $\mu = \mu_1 + \mu_2$ and $\sigma^2 = \sigma_1^2 + \sigma_2^2$. By combining the required expectations and variance, we find that the normal approximation then has the form:

$$Pr(\mathcal{R} = r|n) \approx \mathcal{N}(\mu = [(sn)/(1-s) + N(1 - (1 - 1/N)^n)], \quad (18)$$

$$\sigma = \sqrt{N \left((N-1) \left(1 - \frac{2}{N}\right)^n + \left(1 - \frac{1}{N}\right)^n - N \left(1 - \frac{1}{N}\right)^{2n} \right) + \frac{ns}{(1-s)^2}} \quad (19)$$

Figure 8 shows the quantiles of the normal approximation. We see that up to 99% of the lineages will be contained within the sample of 200 with $Ns = 20$. Larger percentiles will require larger sample sizes.

5. Conclusion

In this work we show that with the increasing sample size, the effect of drift overcomes the effect of selection. As a result, it is possible to construct asymptotically closed solutions to coalescent with selection, provided the sample size is sufficiently large.

The sample size where the expected number of extra lineages required by selection is less than the sample size is well approximated by $2Nxs$. However, the sample size that guarantees that almost no extra lineages are required is considerably larger (7).

Using this observation, we can construct a Markov model that describes the number of derived alleles in the sample. With a sufficiently large sample size, such Markov chains are closed, and can be used for the calculation of the allele frequency spectra with strong selection.

As a future direction, we want to combine the jackknife approximation (Jouganous et al., 2017) with the results presented here. Since the jackknife is an uncontrolled approximation, the current results provide a more sensible approach. However, we can still employ the jackknife in the cases where extra lineages are still required in the current approach. This has the potential of further improving the accuracy of the model and computational efficiency.

6. Appendix

IK: TODO

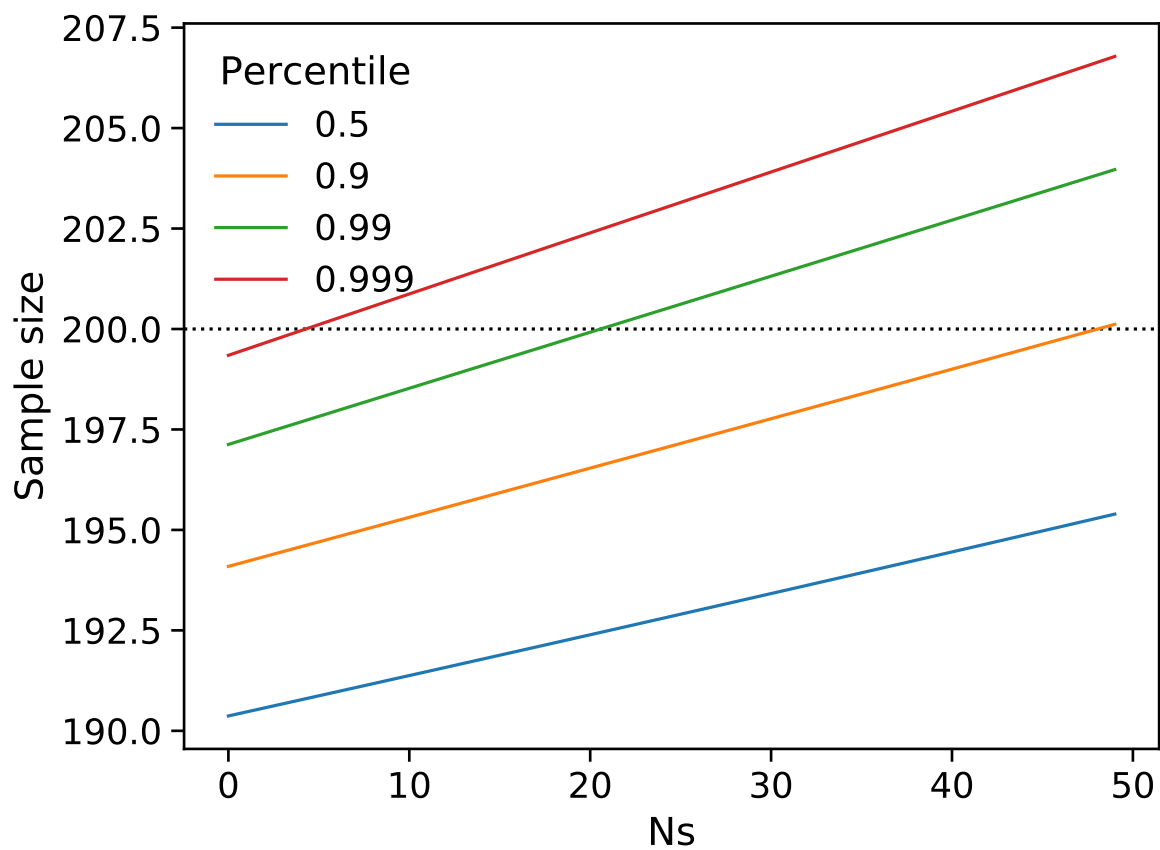


Figure 8: The quantile function of the closure of the sample. Each line corresponds to different percentile of the normal approximation. Black dashed line shows the reference sample size $n = 200$.

References

- Bhaskar, A., Clark, A.G., Song, Y.S., 2014. Distortion of genealogical properties when the sample is very large. *Proceedings of the National Academy of Sciences* 111, 2385–2390. doi:10.1073/pnas.1322709111.
- Donnelly, P., Kurtz, T.G., 1999. Particle representations for measure-valued population models. *The Annals of Probability* 27, 166–205. doi:10.1214/aop/1022677258.
- Ewens, W.J., 1972. The sampling theory of selectively neutral alleles. *Theoretical Population Biology* 3, 87–112. doi:10.1016/0040-5809(72)90035-4.
- Ewens, W.J., 2004. *Mathematical Population Genetics: I. Theoretical Introduction..* volume 27 of *Interdisciplinary Applied Mathematics*. 2 ed., Springer New York, New York. OCLC: 958522782.
- Gutenkunst, R.N., Hernandez, R.D., Williamson, S.H., Bustamante, C.D., 2009. Inferring the joint demographic history of multiple populations from multidimensional snp frequency data. *PLOS Genetics* 5, e1000695. doi:10.1371/journal.pgen.1000695.
- Johnson, N., Kemp, A., Kotz, S., 2005. Occupancy distributions, in: *Univariate Discrete Distributions*. 3 ed.. John Wiley & Sons, Ltd. Wiley Series in Probability and Statistics.
- Jouganous, J., Long, W., Ragsdale, A.P., Gravel, S., 2017. Inferring the joint demographic history of multiple populations: Beyond the diffusion approximation. *Genetics* 206, 1549–1567. doi:10.1534/genetics.117.200493.
- Kamm, J.A., Terhorst, J., Song, Y.S., 2017. Efficient computation of the joint sample frequency spectra for multiple populations. *Journal of Computational and Graphical Statistics* 26, 182–194. doi:10.1080/10618600.2016.1159212.
- Kimura, M., Crow, J.F., 1964. The number of alleles that can be maintained in a finite population. *Genetics* 49, 725–738.
- Kingman, J.F.C., 1982. The coalescent. *Stochastic Processes and their Applications* 13, 235–248. doi:10.1016/0304-4149(82)90011-4.
- Krone, S.M., Neuhauser, C., 1997. Ancestral processes with selection. *Theoretical Population Biology* 51, 210–237. doi:10.1006/tpbi.1997.1299.

- Nelson, D., Kelleher, J., Ragsdale, A.P., McVean, G., Gravel, S., 2019. Coupling wright-fisher and coalescent dynamics for realistic simulation of population-scale datasets. *bioRxiv* , 674440doi:10.1101/674440.
- O'Neill, B., 2019. The classical occupancy distribution: Computation and approximation. *The American Statistician* , 1–12doi:10.1080/00031305.2019.1699445.
- Wakeley, J., 2009. *Coalescent Theory - an Introduction*. W. H. Freeman, New York.