

Rarecoal - Demographic Inference from Rare Mutations

Stephan Schiffels

1 The rarecoal coalescent framework

This model describes a coalescent framework for rare alleles. We define rare alleles roughly by requiring i) the allele count of the derived mutation to be small, typically not larger than 10, and ii) the total number of samples to be much larger, say 100 or more. The idea is to provide a general approach of computing the joint allele frequency spectrum for rare alleles under an arbitrary demographic model under population splits and population size changes. Migration and admixture will be incorporated in the future.

2 Definitions

In the following, we compute the probability to observe a pattern of rare alleles seen across multiple populations, given a demographic model. In the simplest case, a demographic model is tree-like and consists of population split times and constant population sizes in each branch of the tree. Time is counted backwards in time, with $t = 0$ denoting the present and $t > 0$ denoting scaled time in the past. We denote the scaled coalescence rate (scaled inverse population size) in population k at time t by $\lambda_k(t) = N_0/N_k(t)$, where $N_k(t)$ is the population size in population k at time t , and N_0 is a scaling constant which we set to $N_0 = 20000$ for modeling human evolution.

We consider a number of P subpopulations. We define a vector $\mathbf{n} = \{n_k\}$ for $k = 1 \dots P$ summarizing the number of sampled haplotypes in each population. We also define vector $\mathbf{m} = \{m_k\}$ as the set of derived allele counts at a single site in each population. As an example, consider 5 populations with 200 haplotypes sampled in each population, and a rare allele with total allele count 3, with one derived allele seen in population 2 and 2 derived alleles seen in population 3. Then we have $\mathbf{n} = \{200, 200, 200, 200, 200\}$ and $\mathbf{m} = \{0, 1, 2, 0, 0\}$.

Looking back in time, lineages coalesce and migrate, so the numbers of ancestral and derived alleles in the past decrease over time. In theory one needs to consider a very large state space of configurations for this process, with one state for each possible number of ancestral and derived lineages in each population. Here we make a major simplification: While we will consider the full probability distribution over the derived lineages, we will consider only the expected number of ancestral alleles over time. Specifically, we define the expected number of ancestral alleles in population k at time t as $\mathbf{a}(t) = \{a_k(t)\}$. For the derived alleles, we define a state $\mathbf{x} = \{x_k\}$ as a configuration of derived lineages in each population. The probability for state \mathbf{x} at time t is defined by $b(\mathbf{x}, t)$.

3 Coalescence

We now consider the evolution of the two variables $a(t)$ and $b(\mathbf{x}, t)$ through time under the standard coalescent. We first introduce a time discretization. We define time points $t_0 = 0, \dots, t_T$. Here, $t_T = t_{\max}$ should be far enough in the past to make sure that most lineages have coalesced by then with a high probability. We choose a time patterning that is linear in the beginning and crosses over to an exponentially increasing interval width. Specifically, the patterning follows this equation, inspired by the time discretization in [1]:

$$t_i = \alpha \exp \left(\frac{i}{T} \log \left(1 + \frac{t_{\max}}{\alpha} \right) \right) - \alpha. \quad (1)$$

Here, T is the number of time intervals, and α is a parameter that controls the crossover from linear to exponential scale. In practice, we use $\alpha = 0.01$, $t_{\max} = 20$ and $T = 3044$, which are chosen such that

the initial step width equals one generation (in scaled units with $N_0 = 20000$), and the crossover scale is 400 generations.

Given the number of sampled haplotypes in each population n_k , and the observed number of derived alleles m_k in each population, we initialize our variables as follows:

$$a_k(t = 0) = n_k - m_k. \quad (2)$$

for each population k , and

$$b(\mathbf{x}, t = 0) = 1 \text{ if } x_k = m_k \text{ for all } k = 1 \dots P \quad (3)$$

$$b(\mathbf{x}, t = 0) = 0 \text{ otherwise} \quad (4)$$

Under a linear approximation, we can compute the value of \mathbf{a} at a time point $t + \Delta t$, given the value at time t :

$$a_k(t + \Delta t) = a_k(t) \left(1 - \frac{1}{2}(a_k(t) - 1)\lambda_k(t)\Delta t \right). \quad (5)$$

The factor $1/2$ corrects overcounting: any one coalescence takes one of two lineages out, so it should be counted half per participating lineage. We can improve this update equation slightly beyond the linear approximation: In the limit of $\Delta t \rightarrow 0$, equation 5 forms a differential equation which can be solved for finite intervals Δt :

$$a_k(t + \Delta t) = \left(1 + \left(\frac{1}{a_k(t)} - 1 \right) \exp \left(-\frac{1}{2}\lambda_k(t) \times (t + \Delta t) \right) \right)^{-1}. \quad (6)$$

For the derived alleles, we need to update the full probability distribution $b(\mathbf{x}, t)$:

$$\begin{aligned} b(\mathbf{x}, t + \Delta t) = & b(\mathbf{x}, t) \exp \left(- \sum_k \left(\binom{x_k}{2} \lambda_k(t) + x_k a_k(t) \lambda_k(t) \right) \Delta t \right) \\ & + \sum_l b(x_1 \dots (x_l + 1) \dots x_P, t) \left(1 - \exp \left(\binom{x_l + 1}{2} \lambda_l(t) \Delta t \right) \right) \end{aligned} \quad (7)$$

where the first term accounts for the reduction of the probability over time due to derived lineages coalescing among themselves or coalescing with an ancestral lineage, and the second term accounts for the increase from those two processes occurring in states with a higher number of derived lineages. In contrast to the equation for $a(t)$, we cannot solve this as a differential equation and will only use this linear approximation in Δt .

4 Population Splits

Population splits forward in time are joins backward in time. We consider a population join backward in time from population l into population k . For the ancestral lineages, this means that after the join, population k contains the sum of lineages from population k and l :

$$a'_k(t) = a_k(t) + a_l(t) \quad (8)$$

$$a'_l(t) = 0 \quad (9)$$

where the primed variable marks the variable after the event, which will then be used as the basis for the next coalescence update.

For the derived lineages, we need to sum probabilities in the correct way. We first define a transition function that changes a state before the join to new states after the join:

$$\mathbf{x}' = J(\mathbf{x}), \quad (10)$$

where

$$J((\dots x_k \dots x_l \dots)) = (\dots (x_k + x_l) \dots 0 \dots) \quad (11)$$

We can then define the join itself as a sum over all states before the join that give rise to the same state after the join:

$$b'(\mathbf{x}', t) = \sum_{\mathbf{x}, J(\mathbf{x})=\mathbf{x}'} b(\mathbf{x}, t) \quad (12)$$

5 The likelihood of a configuration of rare alleles

Eventually we want to compute the probability for a given configuration (\mathbf{n}, \mathbf{m}) observed in the present. This probability is equal to the probability that a) all derived lineages coalesce before any of them coalesces to any ancestral-allele lineage, and b) that a mutation occurred on the single lineage ancestral to all derived lineages.

We define a singleton state \mathbf{s}^k to be the state in which only $x_k = 1$ and $x_l = 0$ for $l \neq k$. We accumulate the total probability for a single derived lineage:

$$d(t + \Delta t) = d(t) + \sum_k b(\mathbf{s}^k) \Delta t. \quad (13)$$

Then the likelihood of the configuration under the model is

$$L(\mathbf{n}, \mathbf{m}) = \mu d(t_{\max}) \prod_{k=1}^P \binom{n_k}{m_k}, \quad (14)$$

which is the total probability of a mutation occurring on a single derived lineage, times the number of ways that \mathbf{m} derived alleles can be drawn from a pool of \mathbf{n} samples. Note that $d(t_{\max})$ depends on \mathbf{n}, \mathbf{m} and the demographic parameters, which we have omitted for brevity so far.

6 Parameter estimation

The above framework presents a way to efficiently compute the probability of observing a distribution of rare alleles, \mathbf{m} for a large number of samples \mathbf{n} in multiple subpopulations, given a demographic model. We can summarize the full data as a histogram of rare allele configurations. We denote the i th allele configuration by \mathbf{m}_i and the number of times that this configuration is seen in the data by $N(\mathbf{m}_i)$. We then write

$$\mathcal{L}(\{N(\mathbf{m}_i)\}|\Theta) = \prod_i L(\mathbf{m}_i|\Theta)^{N(\mathbf{m}_i)}, \quad (15)$$

where we have introduced a meta-parameter Θ that summarizes the entire model specification (population split times and branch population sizes), and we have made the dependency of L (eq. 14) on Θ explicit. For brevity we have omitted the sample sizes \mathbf{n} . For numerical purpose, we always consider the logarithm of this:

$$\log \mathcal{L}(\{N(\mathbf{m}_i)\}|\Theta) = \sum_i N(\mathbf{m}_i) \log L(\mathbf{m}_i|\Theta). \quad (16)$$

The sum in equation 16 comprises all possible configurations in the genome, in principle. In practice, we only explicitly compute it for configurations between allele count 1 and 4, and replace the rest of the counts with a bulk probability:

$$\log \mathcal{L}(\{N(\mathbf{m}_i)\}|\Theta) = \sum_i I(\text{AC}(i)) N(\mathbf{m}_i) \log L(\mathbf{m}_i|\Theta) + N_{\text{other}} \log L_{\text{other}}(\Theta), \quad (17)$$

where the indicator function $I(\text{AC}(i))$ gives 1 if the allele count is between 1 and 4, and 0 otherwise. The bulk count N_{other} simply counts up sites with either no variant or variants with allele count larger than 4. The bulk probability is simply:

$$L_{\text{other}}(\Theta) = 1 - \sum_i (1 - I(\text{AC}(i)) L(\mathbf{m}_i|\Theta)), \quad (18)$$

With a given population tree and a given histogram of allele configuration counts $N(\mathbf{m}_i)$, we implemented numerical optimizations over the parameters Θ to find the maximum likelihood parameters, and MCMC to estimate the posterior distributions for all parameters given the data. We usually first search for the maximum with the optimization method, which is much faster than MCMC, and then use MCMC to explore the distribution around that maximum.

7 Implementation

We implemented this method in the Haskell programming language as a program called “rarecoal”, available from github at <https://github.com/stschiff/rarecoal.hs>.

References

- [1] Li, H., and Durbin, R. (2011). Inference of human population history from individual whole-genome sequences. <http://doi.org/10.1038/nature10231>
- [2] Staab, P. R., Zhu, S., Metzler, D., and Lunter, G. (2015). scrm: efficiently simulating long sequences using the approximated coalescent with recombination. *Bioinformatics*, btu861. <http://doi.org/10.1093/bioinformatics/btu861>
- [3] 1000 Genomes Project. (2015). A global reference for human genetic variation. *Nature* (in Revision).