

EMBO Course in Population Genomics, practicals on Coalescent Theory

Question 1. **Group 1**

Consider a diploid population of size N (with $2N$ chromosomes) that has two alleles, C and c , segregating in i and $2N - i$ copies respectively, at $t = 0$. The population is evolving following the standard Wright-Fisher model.

- (a) Calculate the probability that allele c will be present in k copies in the first generation ($t = 1$).
- (b) What is the probability of observing allele C in $2N$ copies in the first generation ($t = 1$).
- (c) What is the probability that a sample of size $n = 2$ taken from generation t share the same parent at $t - 1$?
- (d) With what probability k offspring at generation t share the same parent at $t - 1$?

Question 2. **Group 2**

- (a) In a haploid population of size $N = 10$, what is the probability that two lineages coalesce exactly 11 generations in the past?
- (b) What is the probability that it takes at least 11 generations for them to coalesce?

Question 3. **Group 3**

The coalescence times T_i in the standard coalescent process are independent and exponentially distributed as

$$f_{T_i}(t) = \binom{i}{2} e^{-\binom{i}{2} t}$$

when time is measured appropriately.

- (a) What is the distribution of T_2 , T_3 , and T_4 ?
- (b) Write down the average and variance of waiting times T_2 , T_3 , and T_4 . Also, write down the value of $E[T_2^2]$, $E[T_3^2]$, and $E[T_4^2]$.
- (c) If one infers that $T_3 > T_4 > T_2$ upon looking at a gene genealogy for a sample of size 4, would this be consistent with the standard coalescent process? Justify your answer.
- (d) Does $\binom{i}{2}$ in the expression given above have some meaning in standard coalescent process (Kingman's coalescent)? If yes, what is it?

Question 4. Group 4

$T_{MRC A}$ is the time to the most recent common ancestor of the entire sample and T_{total} is the total length of all the branches in the genealogy.

- (a) Draw a coalescent tree with sample size $n = 4$ and indicate on the tree what $T_{MRC A}$ and T_{total} correspond to.
- (b) Find $E[T_{MRC A}]$ and $E[T_{total}]$ for a sample of size n . Write down all the steps involved.
- (c) Calculate the value of $E[T_{MRC A}]$ as n goes to infinity.
- (d) Write down the $E[T_{MRC A}]$ and $E[T_{total}]$ when sample size is 4. Also, compute $\frac{E[T_{MRC A}]}{E[T_2]}$ (definition of T_2 is given in the previous question).
- (e) Can $\frac{E[T_{MRC A}]}{E[T_{total}]} > 1$ be true? Justify your answer.

No, we have that $E[T_{total}] < E[T_{MRC A}]$ cannot be true. The branches that are summed to get to $T_{MRC A}$ are included in the sum of the T_{total} calculation.

Question 5. Group 1+2

Consider the data below obtained from $n = 12$ human mtDNA sequences from African populations:

- (a) Count the number of segregating sites S for the given sampled sequences.

```

Biaka_AF346969/1-1656tata t t a c t a c c a c t g a c a t g a c t t t c c a a a a a a c a c a t a a t t t g a a t c a a c a c a c c a c c c a c a g c c t a a t t a t t a g c a t c a t c c c c t a c t a t t t t t t t a a c c a a a t c a a c a a c a a c c t a t t t a g c t g t t c c c c a a c c t t
Mbenzele_AF346996/1-1656tata t t a c t a c c a c t g a c a t g a c t t t c c a a a a a a c a c a t a a t t t g a a t c a a c a c a c c a c c c a c a g c c t a a t t a t t a g c a t c a t c c c c t a c t a t t t t t t t a a c c a a a t c a a c a a c a a c c t a t t a g c t g t t c c c c a a c c t t
Ibo_AF346987/1-1656tata t t a c t a c c a c t g a c g t g a c t t t c c a a a a a a c a c a t a a t t t g a a t c a a c a c a c c a c c c a c a g c c t a a t t a t t a g c a t c a t c c c c t a c t a t t t t t t t a a c c a a a t c a a c a a c a a c c t a t t a g c t g t t c c c c a a c c t t
Kikuyu_AF346992/1-1656tata t t a c t a c c a c t g a c a t g a c t t t c c a a a a a a c a c a t a a t t t g a a t c a a c a c a c c a c c c a c a g c c t a a t t a t t a g c a t c a t c c c c t a c t a t t t t t t t a a c c a a a t c a a c a a c a a c c t a t t a g c t g t t c c c c a a c c t t
SouthAfrican_AY195766tata t t a c t a c c a c t g a c a t g a c t t t c c a a a a a a c a c a t a a t t t g a a t c a a c a c a c c a c c c a c a g c c t a a t t a t t a g c a t c a t c c c c t a c t a t t t t t t t a a c c a a a t c a a c a a c a a c c t a t t a g c t g t t c c c c a a c c t t
Mauritania_AF381981/1-1656tata t t a c t a c c a c t g a c a t g a c t t t c c a a a a a a c a c a t a a t t t g a a t c a a c a c a c c a c c c a c a g c c t a a t t a t t a g c a t c a t c c c c t a c t a t t t t t t t a a c c a a a t c a a c a a c a a c c t a t t a g c t g t t c c c c a a c c t t
Mandenka_AF346995/1-1656tata t t a c t a c c a c t g a c a t g a c t t t c c a a a a a a c a c a t a a t t t g a a t c a a c a c a c c a c c c a c a g c c t a a t t a t t a g c a t c a t c c c c t a c t a t t t t t t t a a c c a a a t c a a c a a c a a c c t a t t a g c t g t t c c c c a a c c t t
SouthAfrican_AY195785tata t t a c t a c c a c t g a c a t g a c t t t c c a a a a a a c a c a t a a t t t g a a t c a a c a c a c c a c c c a c a g c c t a a t t a t t a g c a t c a t c c c c t a c t a t t t t t t t a a c c a a a t c a a c a a c a a c c t a t t a g c t g t t c c c c a a c c t t
San_AY195788/1-16570tata t a c t a c c a c t g a c a t g a c t t t c c a a a a a a c a c a t a a t t t g a a t c a a c a c a c c a c c c a c a g c c t a a t t a t t a g c a t c a t c c c c t a c t a t t t t t t t a a c c a a a t c a a c a a c a a c c t a t t a g c t g t t c c c c a a c c t t
Effik_AF346976/1-1656tata t t a c t a c c a c t g a c a t g a c t t t c c a a a a a a c a c a t a a t t t g a a t c a a c a c a c c a c c c a c a g c c t a a t t a t t a g c a t c a t c c c c t a c t a t t t t t t t a a c c a a a t c a a c a a c a a c c t a t t a g c t g t t c c c c t a a c c t t
Effik_AF346977/1-1656tata t t a c t a c c a c t g a c a t g a c t t t c c a a a a a a c a c a t a a t t t g a a t c a a c a c a c c a c c c a c a g c c t a a t t a t t a g c a t c a t c c c c t a c t a t t t t t t t a a c c a a a t c a a c a a c a a c c t a t t a g c t g t t c c c c a a c c t t
SouthAfrican_AY195776tata t t a c t a c c a c t g a c a t g a c t t t c c a a a a a a c a c a t a a t t t g a a t c a a c a c a c c a c c c a c a g c c t a a t t a t t a g c a t c a t c c c c t a c t a t t t t t t t a a c c a a a t c a a c a a c a a c c t a t t a g c t g t t c c c c a a c c t t

```

Figure 1: Snapshot of $L = 141$ multiple aligned mtDNA nucleotides from 12 individuals. Nucleotides are highlighted in blue whenever they match the base reported in the reference genome.

(b) In population genetics, one quantity of interest is the parameter $\theta = 4N_e\mu$ (N_e is the effective population size and μ is the mutation rate per individual/generation). Compute two common estimators of θ , $\hat{\theta}_S = \frac{S}{L \sum_{i=1}^{n-1} \frac{1}{i}}$ (here S is the number of segregating sites, n is the sample size) and $\hat{\theta}_\pi = \frac{1}{L \binom{n}{2}} \sum_{i < j} d_{ij}$ (d_{ij} is the number of pairwise differences between sequence i and j) from these African sequences .

(c) What is your expectation for these estimators for other human populations? Should they be larger or smaller? Explain your answer.

Question 6. **Group 3+4** Consider the data below, where $n=10$ haploid sequences were simulated. Each line represents an individual and each column a segregating site S . The ancestral allele is denoted by "0" and the derived allele is denoted by "1." We observe a total of $S=20$ segregating sites.

```
00001000010000000000
000000000100100010011
000000000100100010011
000000000000101110001
000000000000101110001
100000000000101010001
01010111001010011001
100000000000101010001
00100000100100010111
000000000000101110001
```

(a) What is the frequency of the **derived** allele in each segregating site?

(b) Compute the SFS for the data above, using the allele frequencies you just computed.

The Figure 2 below shows two Site Frequency Spectra (SFS). The red distribution is the expected SFS for a given value of θ for a population of constant size, not subdivided (geographically or otherwise), and where mutations have no effect in fitness. Mutation rate per base pair is constant and loci are not linked. The blue distribution is the observed SFS

for a population simulated under neutrality using *ms* (Hudson 1990).

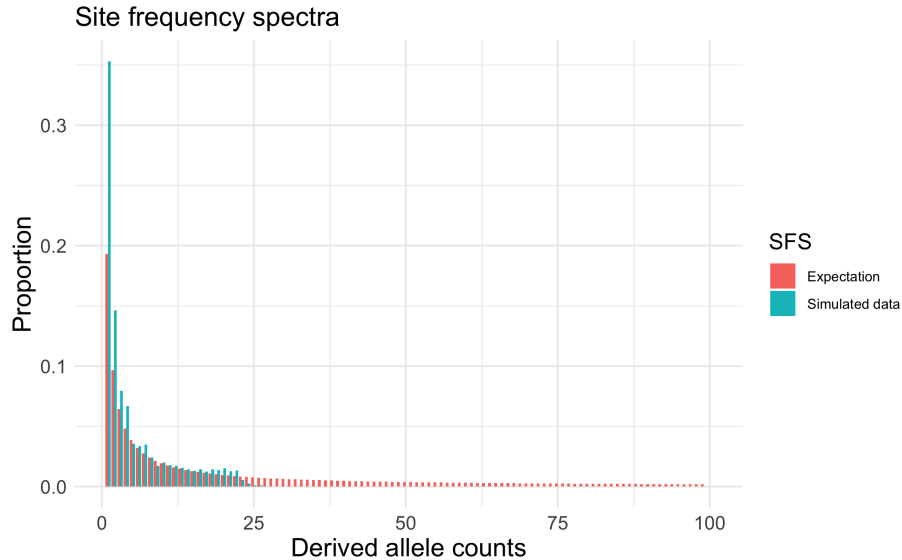


Figure 2: Site frequency spectra for two scenarios: a simulated dataset, and the expected SFS for a population with constant population size.

(c) For the red curve (expectation with constant N_e), we see 800 singletons. What is the expected number of doubletons?

(d) We can see that the SFS for the simulations (blue) departs from the expectation (red). Besides some statistical noise, what is the most striking difference between these two distributions? Please use maximum two sentences to answer this question.

(e) Given that the same value of θ was used to generate both distributions, can you think of an explanation for the observed differences? Please use maximum three sentences to answer this question.

Question 7. **Group 1+2+3+4** Write a small program to simulate allele frequencies through time under a Wright-Fisher model for the simplest case, e.g. neutral alleles, constant population size. Ideally, your program should output the allele frequencies through time.