

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.

Recall: Assuming Wright-Fisher, the probability that an allele with i copies in the present generation is found in j copies in the next generation is given by

$$P_{ij} = \binom{2N}{j} p_0^j (1 - p_0)^{2N-j},$$

$p_0 = \frac{i}{2N}$ is the frequency of the allele in the present generation. P_{ij} is a binomial distribution.

(a) Calculate the probability that allele c will be present in k copies in the first generation ($t = 1$).

Answer: For this case: probability of success $p = (1 - \frac{i}{2N})$, therefore, the required probability is:

$$\binom{2N}{k} (1 - \frac{i}{2N})^k (\frac{i}{2N})^{2N-k}.$$

(b) What is the probability of observing allele C in $2N$ copies in the first generation ($t = 1$).

Answer: For this case: probability of success $p = \frac{i}{2N}$, therefore,

required probability = $\binom{2N}{2N} (\frac{i}{2N})^{2N} (1 - \frac{i}{2N})^{2N-2N} = (\frac{i}{2N})^{2N}$.

(c) What is the probability that a sample of size $n = 2$ taken from generation t share the same parent at $t - 1$?

Answer: Required probability = (probability that the 1st progeny chooses a parent from the $2N$ individuals present at $t - 1$)*(probability that the 2nd progeny chooses the same parent present at $t - 1$)*(number of parents present at $t - 1$),

$$\text{Required probability} = \left(\frac{1}{2N}\right) * \left(\frac{1}{2N}\right) * 2N = \frac{1}{2N}.$$

(d) With what probability k offspring at generation t share the same parent at $t - 1$?

Answer: Following the same logic, required probability = $\left(\frac{1}{2N}\right)^{k-1}$.

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?

Answer:

$$(a) P[X = 11] = (1 - p)^{10}p = 0.03486784, \text{ where } p = \frac{1}{10}.$$

$$(b) P[X > 10] = 1 - \sum_{i=0}^9 (1 - p)^i p = (1 - p)^{10} = 0.3486784$$

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 ?

Answer:

$$f_{T_2}(t) = \binom{2}{2} e^{-\binom{2}{2} t} = e^{-t},$$

$$f_{T_3}(t) = \binom{3}{2} e^{-\binom{3}{2} t} = 3 e^{-3 t},$$

$$f_{T_4}(t) = \binom{4}{2} e^{-\binom{4}{2} t} = 6 e^{-6 t}.$$

(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 $E[T_4^2]$.

Answer:

$$E[T_2] = 1, \quad E[T_3] = \frac{1}{3}, \quad E[T_4] = \frac{1}{6}.$$

$$\text{Var}[T_2] = 1, \quad \text{Var}[T_3] = \left(\frac{1}{3}\right)^2 = \frac{1}{9}, \quad \text{Var}[T_4] = \left(\frac{1}{6}\right)^2 = \frac{1}{36}.$$

By definition,

$$E[i^2] = \text{Var}[i] + (E[i])^2, \text{ therefore,}$$

$$E[T_2^2] = 2, \quad E[T_3^2] = \frac{2}{9}, \quad \text{and} \quad E[T_4^2] = \frac{2}{36}.$$

(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.

Answer: Yes, it would be consistent with the standard coalescent process. T_2 , T_3 , and T_4 are the random variables. So, for some genealogies, $T_3 > T_4 > T_2$ holds. Though $E[T_2] > E[T_3] > E[T_4]$ will always be true.

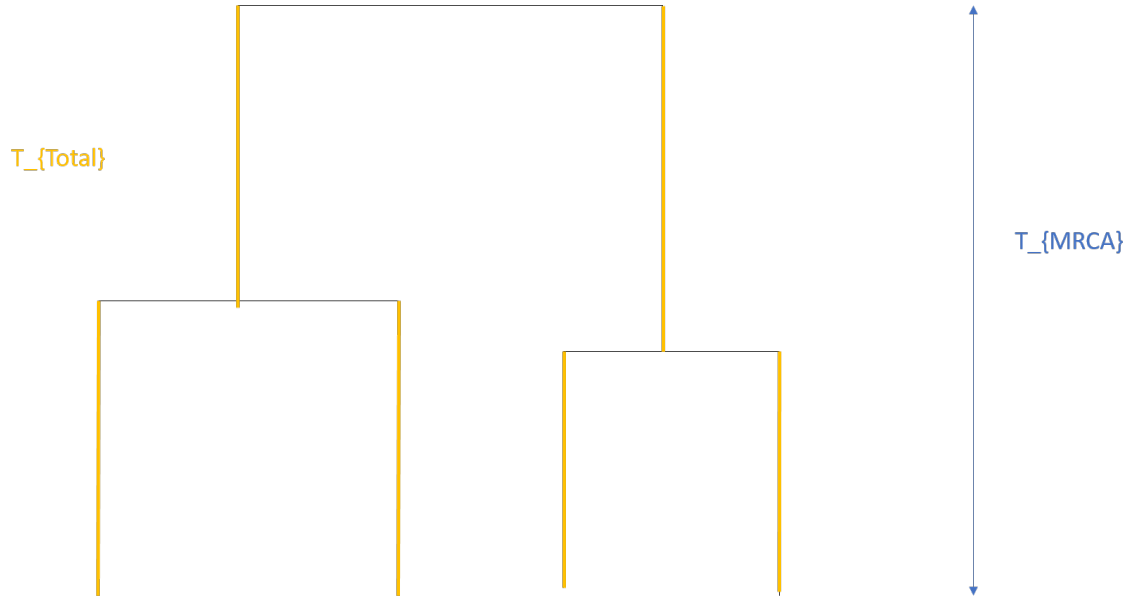
(d) Does $\binom{i}{2}$ in the expression given above have some meaning in standard coalescent process (Kingman's coalescent)? If yes, what is it?

Answer: This quantity represents the rate of coalescence events. Since only two lineages merge at the same in the standard coalescence, it is equal to $\binom{i}{2}$ the number of pairs among i lineages.

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_{MRCA} and T_{total} correspond to.



(b) Find $E[T_{MRCA}]$ and $E[T_{total}]$ for a sample of size n . Write down all the steps involved.

Answer:

$$T_{MRCA} = \sum_{i=2}^n T_i$$

where n is the sample size, since T_i are independent random variables, therefore,

$$E[T_{MRCA}] = \sum_{i=2}^n E[T_i],$$

$$E[T_{MRCA}] = \sum_{i=2}^n \frac{2}{i(i-1)} = 2 \sum_{i=2}^n \left(\frac{1}{i-1} - \frac{1}{i} \right),$$

$$E[T_{MRCA}] = 2 \left(\frac{1}{1} - \frac{1}{2} + \frac{1}{2} - \frac{1}{3} + \frac{1}{3} - \dots - \frac{1}{n-1} + \frac{1}{n-1} - \frac{1}{n} \right) = 2 \left(1 - \frac{1}{n} \right).$$

$$T_{total} = \sum_{i=2}^n i T_i,$$

using $E[cX] = c E[X]$,

$$\begin{aligned} E[T_{total}] &= \sum_{i=2}^n i E[T_i], \\ E[T_{total}] &= \sum_{i=2}^n \frac{2 i}{i(i-1)} = \sum_{i=2}^n \frac{2}{(i-1)}. \end{aligned}$$

(c) Calculate the value of $E[T_{MRCA}]$ as n goes to infinity.

Answer:

$$\lim_{n \rightarrow \infty} E[T_{MRCA}] = 2$$

(d) Write down the $E[T_{MRCA}]$ and $E[T_{total}]$ when sample size is 4. Also, compute $\frac{E[T_{MRCA}]}{E[T_2]}$ (definition of T_2 is given in the previous question).

Answer:

$$E[T_{MRCA}] = 2 \left(1 - \frac{1}{n}\right) = 2 \left(1 - \frac{1}{4}\right) = \frac{3}{2} = 1.5,$$

$$E[T_{total}] = \sum_{i=2}^4 \frac{2}{(i-1)} = \frac{11}{3} = 3.67.$$

$$\frac{E[T_{MRCA}]}{E[T_2]} = \frac{1.5}{1} = 1.5.$$

(e) Can $\frac{E[T_{MRCA}]}{E[T_{total}]} > 1$ be true ? Justify your answer.

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

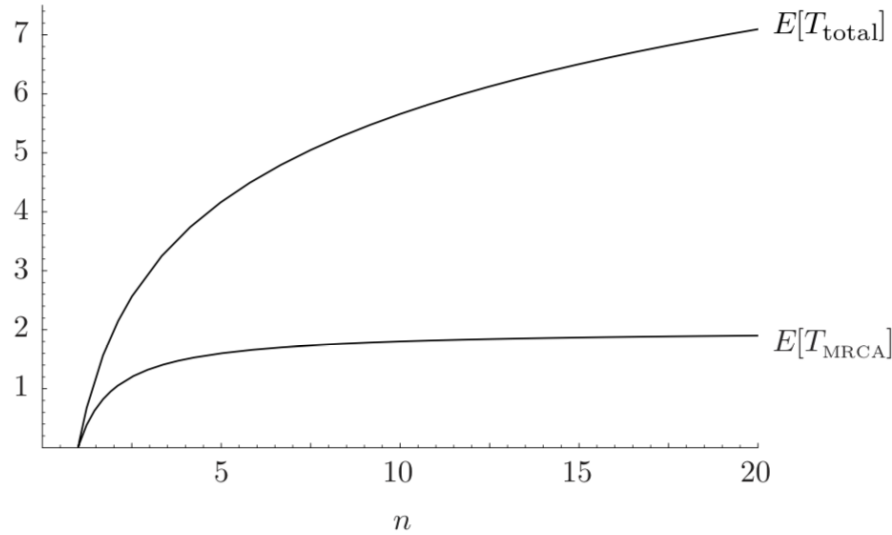


Figure 1: Reference: An Introduction to Coalescence Theory by John Wakeley.

Conclusion: $E[T_{total}] > E[T_{MRCA}]$ **will always be true.**

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.

Answer: We find $S = 7$ SNPs or segregating sites.

(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 .

Answer: $S = 7$, $\sum_{i < j} d_{ij} = 98$, $n = 12$, $L = 141$; $\hat{\theta}_S = 0.0164$ and $\hat{\theta}_\pi = 0.0105$.

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

Biaka_AF346969/1-1656 t a t a 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 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
Mbenzele_AF346996/1- t a t a 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 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
lbo_AF346987/1-16565 t a t a 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 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_AF346992/1-165 t a t a t t g t 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 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_AY195766 t a t a 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 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 t a t a 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 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- t a t a 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 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_AY195785 t a t a 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 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_AY195788/1-16570 t a t a 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 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_AF346976/1-1656 t a t a 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 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_AF346977/1-1656 t a t a 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 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_AY195776 t a t a 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 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 2: 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.

Answer: We expect lower values for non-African populations due to their lower effective population size after the out of Africa bottleneck. The quantity $\theta = 4N_e\mu$ can be seen as a measure of variation in the population.

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
00000000100100010011
00000000100100010011
000000000000101110001
000000000000101110001
000000000000101010001
```

01010111001010011001
10000000000101010001
00100000100100010111
00000000000101110001

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

Answer: 2 1 1 1 1 1 1 3 1 1 8 1 5 3 9 1 1 3 9

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

Answer: We consider the frequencies of the derived allele (1) and count the number of derived alleles at each site to get the following:

Derived allele count	1	2	3	4	5	6	7	8	9
# of sites	12	1	3	0	1	0	0	1	2

Table 1: Unfolded SFS

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).

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

Answer: $E[S_1] = \theta = 800$. $E[S_2] = \theta/2 = 400$. We expect to see 400 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.

Answer: The most striking difference is that the blue curve has more than 50% singletons than the red curve. In other words, there seems to be an excess of singletons in the simu-

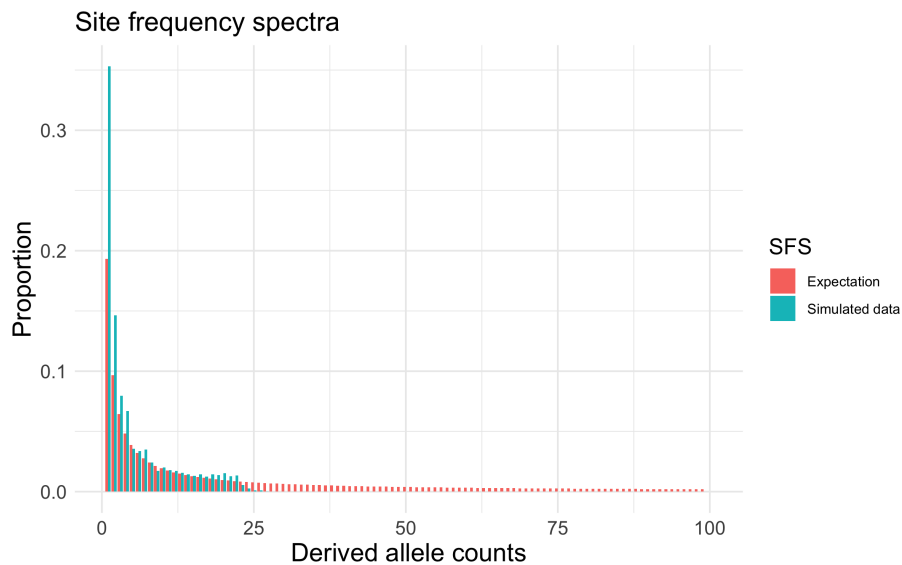


Figure 3: Site frequency spectra for two scenarios: a simulated dataset, and the expected SFS for a population with constant population size.

lations relative to the expectation.

(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.

Answer: The explanation is that the simulated population went through some sort of demographic change. Here, the population was simulated with a bottleneck and exponential growth.

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.

Code used by Andy Clark in 2018.

WFsimsim.r simple simulations of the Wright-Fisher model

Initialize popsize (N), number of samples (samp), number of
generations (ngen) and starting frequency (startfreq)

```

# Note this is for a haploid population of size N.
N<-200
nsamp<-8
ngen<-200
startfreq<-0.5
pcur<-matrix(1:N)

#Wright-Fisher is simply recurrent binomial sampling over generations

x<-rbinom(nsamp,N,startfreq)
p<-x/N
pcur<-p

for (i in 1:ngen){
x<-rbinom(nsamp,N,p)
p<-x/N
pcur<-rbind(pcur,p)
}

#Now plot these trajectories
gen<-seq(1,ngen+1)
for (k in 1:nsamp){
plot(gen,pcur[,k],type="l",ylim=c(0,1))
par(new=TRUE)
}

```