

What is a mismatch distribution?

Count the number of site differences between each pair of sequences in a sample, and use the resulting counts to build a histogram. You end up with a "mismatch distribution." The i th entry of the mismatch distribution is the number of pairs of sequences that differ by i sites.

Mismatch Distributions and Population Growth

Alan R. Rogers

October 2, 2015

Partial mtDNA sequences from Asia

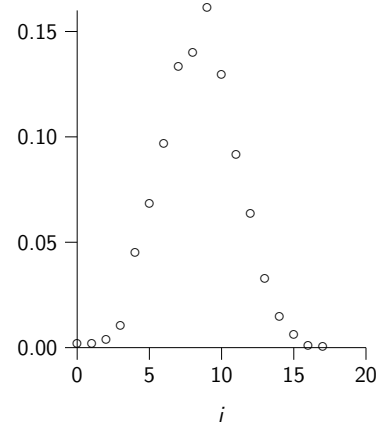
```

1 CATTAGCACCCAAAGCTAAGATTCTAATTTAACTATTCTCTGTTCTTTCATGG...
2 CATTAGCACCCAAAGCTAAGATTCTAATTTAACTATTCTCTGTTCTTTCATGG...
3 CATTAGCACCCAAAGCTAAGATTCTAATTTAACTATTCTCTGTTCTTTCATGG...
4 CATTAGCACCCAAAGCTAAGATTCTAATTTAACTATTCTCTGTTCTTTCATGG...
5 CATTAGCACCCAAAGCTAAGATTCTAATTTAACTATTCTCTGTTCTTTCATGG...
6 CATTAGCACCCAAAGCTAAGATTCTAATTTAACTATTCTCTGTTCTTTCATGG...
7 CATTAGCACCCAAAGCTAAGATTCTAATTTAACTATTCTCTGTTCTTTCATGG...
8 CATTAGCACCCAAAGCTAAGATTCTAATTTAACTATTCTCTGTTCTTTCATGG...
9 CATTAGCACCCAAAGCTAAGATTCTAATTTAACTATTCTCTGTTCTTTCATGG...
10 CATTAGCACCCAAAGCTAAGATTCTAATTTAACTATTCTCTGTTCTTTCATGG...
11 CATTAGCACCCAAAGCTAAGATTCTAATTTAACTATTCTCTGTTCTTTCATGG...
12 CATTAGCACCCAAAGCTAAGATTCTAATTTAACTATTCTCTGTTCTTTCATGG...
13 CATTAGCACCCAAAGCTAAGATTCTAATTTAACTATTCTCTGTTCTTTCATGG...
14 CATTAGCACCCAAAGCTAAGATTCTAATTTAACTATTCTCTGTTCTTTCATGG...
15 CATTAGCACCCAAAGCTAAGATTCTAATTTAACTATTCTCTGTTCTTTCATGG...
.....

```

Mismatch distribution for Asian data

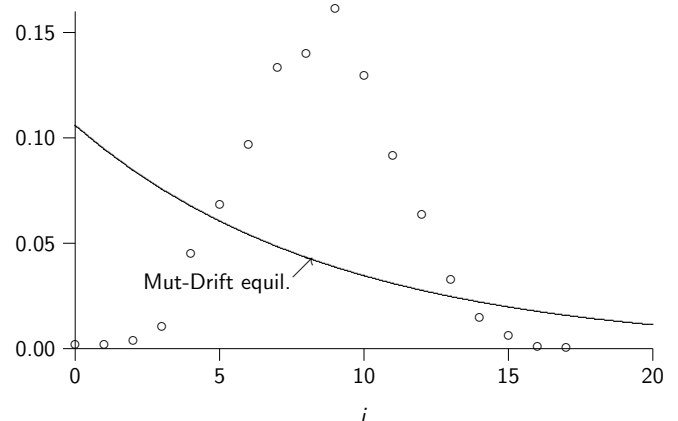
| i | n | i | n |
|-----|-----|-----|-----|
| 0 | 5 | 10 | 379 |
| 1 | 5 | 11 | 268 |
| 2 | 11 | 12 | 186 |
| 3 | 30 | 13 | 95 |
| 4 | 131 | 14 | 43 |
| 5 | 200 | 15 | 17 |
| 6 | 283 | 16 | 2 |
| 7 | 390 | 17 | 1 |
| 8 | 409 | | |
| 9 | 471 | | |



At mutation-drift equilibrium, a random pair of sequences differs by i sites with probability

$$F_i = \left(\frac{1}{\theta + 1} \right) \left(\frac{\theta}{\theta + 1} \right)^i, \quad (i = 0, 1, 2, \dots) \quad (1)$$

mtDNA mismatch distribution doesn't fit equilibrium model



Why does the stationary neutral model fit human data so poorly?

- There are several hypotheses to consider:
1. Sampling error. (Important because the pairs of genes in our sample are correlated.)
 2. Selection.
 3. Failure of infinite sites hypothesis.
 4. Non-random mating.
 5. Variation in population size.
- Work has been done on all of these possibilities, but I will only try to tell you about the last one.

Coalescent theory in a population of varying size

At any given time, t , the hazard of a coalescent event is

$$h_i(t) = \frac{i(i-1)}{4N(t)}$$

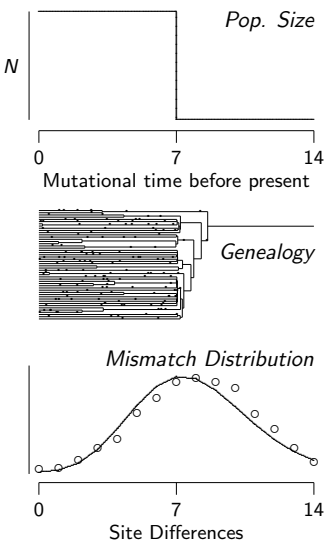
But $N(t)$ is no longer constant.

$$E[t_i] \neq 1/h_i$$

We need computer simulations.

Principles

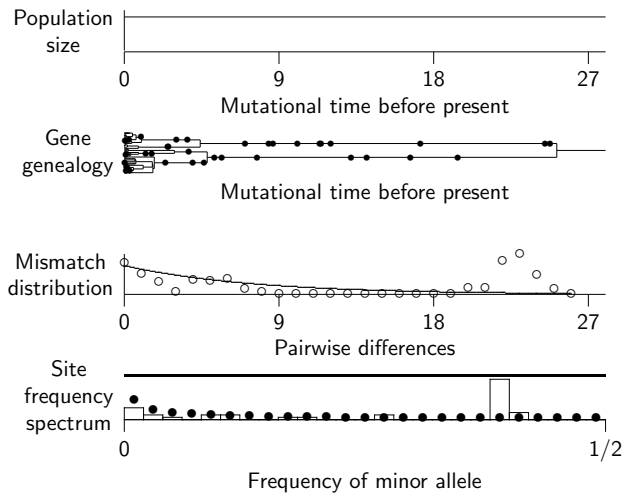
- Coalescent intervals tend to be long in those parts of the tree where
- ▶ there are only a few lines of descent
 - ▶ the population size is large



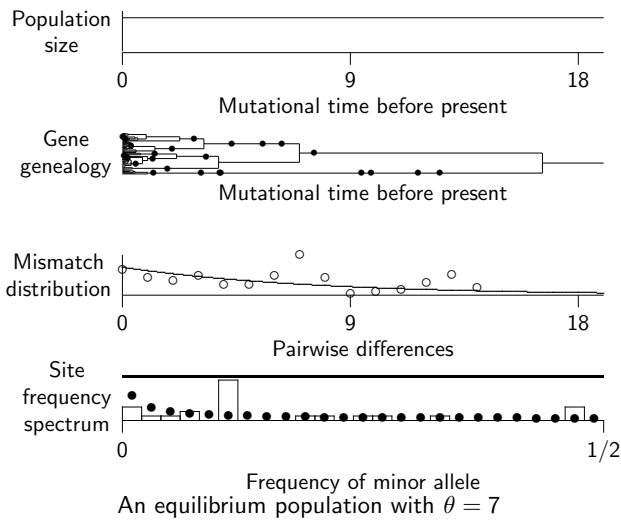
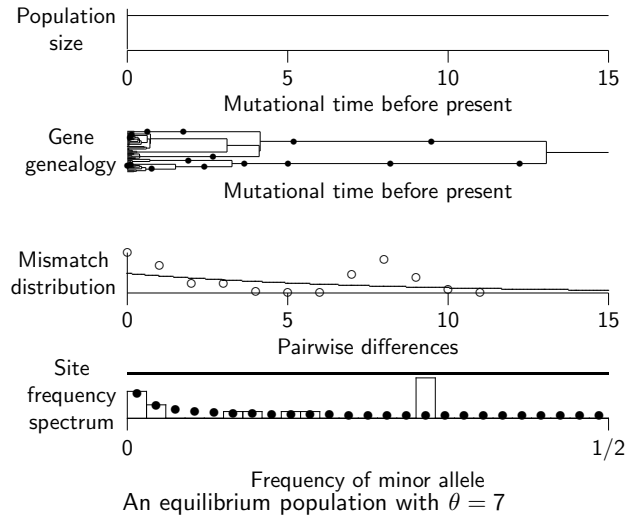
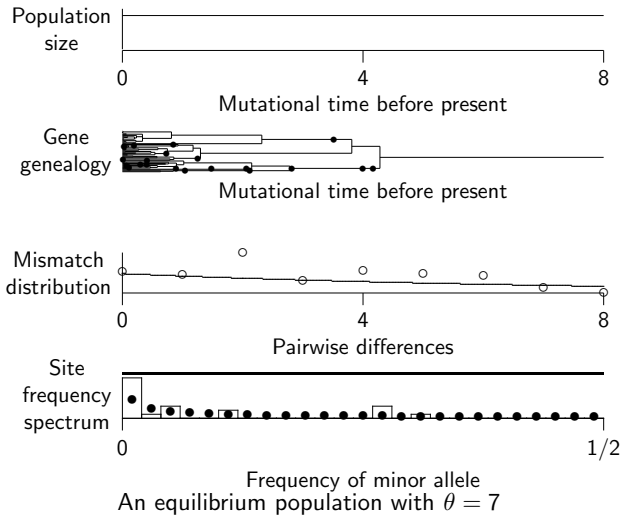
Effect of a population explosion

Middle: genealogy of 50 individuals; dots are mutations.
 Bottom: \circ = simulated data, line = theory.

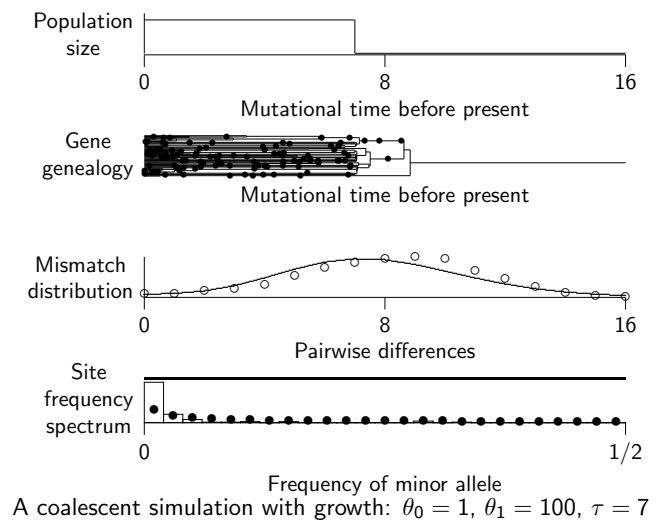
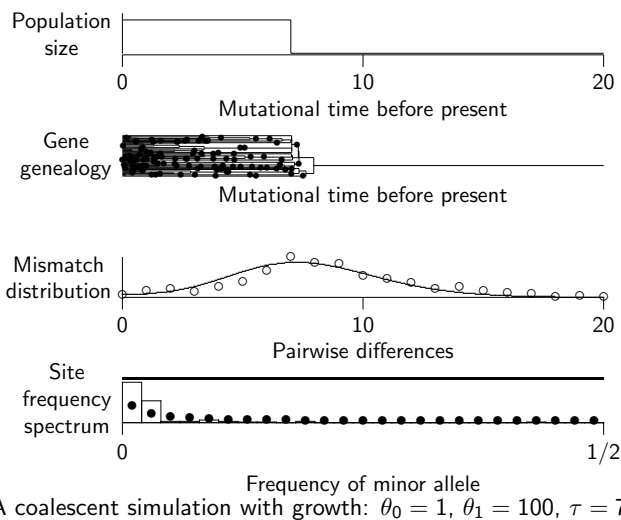
Simulations of stationary populations

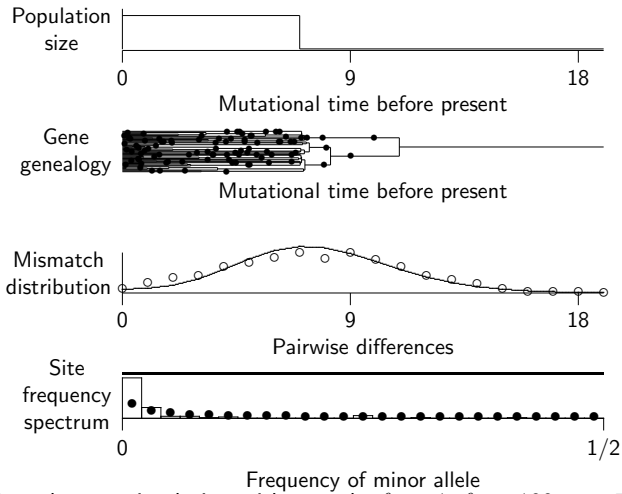


An equilibrium population with $\theta = 7$

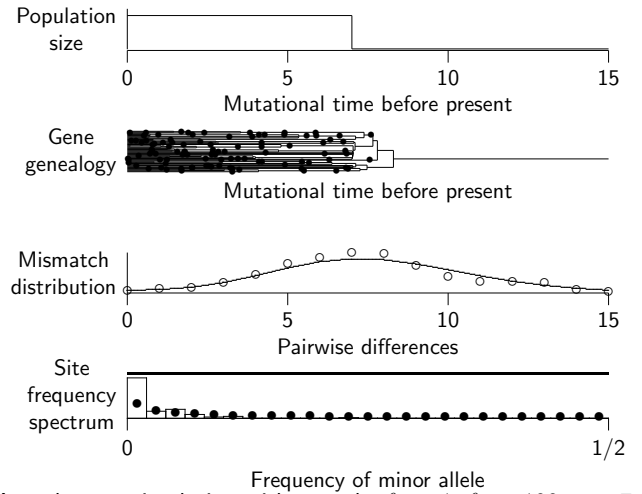


Simulations of expanded populations





A coalescent simulation with growth: $\theta_0 = 1$, $\theta_1 = 100$, $\tau = 7$



A coalescent simulation with growth: $\theta_0 = 1$, $\theta_1 = 100$, $\tau = 7$

Model of sudden growth fit to Asian data

