

② I've only done a) as b) is not required for Assignment 2.

The simulator code is on page 6.

The tree topology is generated in linear time (see function `simulate-tree`). The distribution of edge lengths is different for exponential, increasing N . The simulation is done by inverting the cumulative distribution function.

Formally let $d_T = d_T^{(k)}$ be the length of the edge starting at some generation T ($\Rightarrow N_T = e^{-d_T N_0}$) // for k individuals.

$$P(d_T = 1) = \frac{\binom{k}{2}}{N_T}, \quad P(d_T = 2) = \frac{\binom{k}{2}}{N_T e^{-d}} \left(1 - \frac{\binom{k}{2}}{N_T e^{-d}}\right)$$

$$\Rightarrow P(d_T = l+1) = \frac{\binom{k}{2}}{N_T e^{-d}} \prod_{s=1}^l \left(1 - \frac{\binom{k}{2}}{N_T e^{-ds}}\right).$$

The simulation idea comes from the fact that

$$\frac{P(d_T = l+1)}{P(d_T = l)} = e^{-d} \left(1 - \frac{\binom{k}{2}}{e^{-dk} N_T}\right) \quad \text{— a very simple close form that allows one to go from } P(d_T = l) \text{ to } P(d_T = l+1) \text{ in short time.}$$

The mutations are simulated by Poisson as was described on the lecture.

Simulator saves the tree in newick format and SNP matrix as `np.savetxt` result.

On page 7 I also saved SNP matrix and a tree for $n=10$.

CSE 280A

Assignment 2

Andrey Bzikalae

List 5 || ②