CTGTGCTGTTGCGTTATAATACTGGGCTCCAAGTACGTGTCCGGTCCCCATGTCTGCGCAAGAATTTAAG GAGCCASIMULATING Sequence data on a phylogenetic tree TATC ATACATTACAACCCTTTTTATTACTACTCGTCAAACAGGTTTGCTTTTTGGTCTGAAATCGAGGCGTCATA GCGATATCAGTACCGGTCGTGAGAGGGAAGTGAGGCCACGGGGTAAAAACAACAACTGGCCCGTTC TGGCTCGTGCGAGTTTACTTGTCGCTTCCACACGCGAGCCGATCTCTTGATAACAGACTAAGAGCAGGCG AACATTTTTTACTGATTTCAGGCAAACGCCATTCTAATTAGGCGGTTAGGTCTGGTGCCATAGACAACG CATACGCATAATCTCCTTGGAAAGCTATTTGGACTGTTGGCGTGGAGCATCGGCGGGGGTAGAAATAGGGC GCTGAGGCGAAAAGGCTCACCTAGATCGTTATCCGACACTTCAATCTCTACCTCGGGTCTATTACCAT CGCTCCTGCCTTAAGCCACAGGACAAGTGCGTTCAATCTCACCCAACTGGAGTAGGAAAATTAGCCC CGGGGTGCCCCCGCAATGCTAGATTTTCAAGACAGATGGGACCATGTCTCCGATATGACAATATTT AAGTCTAAGTGACGGGACGCATTACAACGTATTATAAAATCCATATGTGTCTTCCTATCTGGAGTGGTTG ATTGGCAAGTTGGGTAAGAGTTATTCATGACAGGCGGCCCGGATCCCGCAAATACTCTTGGTGGTATCAC GGCCCGCCCCGGCGTTTAGAACTGTCTATTGAAACCGCTGTTAGAGTAATTATTTTACCGCATATATGAA GAGTCACCCACCTAATGTCTGTGCTGTTGCGTTATAATACTGGGCTCCAAGTACGTGTCCGGTCCCCA GCGCAAGAATTTAAGGAGCCAGCAGTAAGTACTCCGTCTAGTAAAATTCGGGCATAAGTCGGAGG CAAGTAAACGCCCTATCATACATTACAACCCTTTTTATTACTACTCGTCAAACAGGTTTGCTTTTTGG GAAATCGAGGCGTCATAGTTACGCGATATCAGTACCGGTCGTGAGAGGAAGTGAGGCCACGGGGTAA AAACAACAACTGGCCCGTTCTGGCTCGTGCGAGTTTACTTGTCGCTTCCACACGCGAGCCGATCTCTTGA TAACAGACTAAGAGCAGGCGTAAACATTTTTACTGATTTCAGGCAAACGCCATTCTAATTAGGCGGTTAG CTGGTGCCATAGACAACGCATACGCATAATCTCCTTGGAAAGCTATTTGGACTGTTGGCGTGGAGTAT CGGCGGGGTAGAAATAGGGCTCACGGTCTACATTAATGAACTAGTCTTAGCCATACGTGGACGCGGGGGGC GAATCCAACGAACACGGATTGCTGAGGCGAAAAGGCTCACCTAGATCGTTATCCGACACTTCAATCTC CTCGGGTCTATTACCATCCCCGCTCCTGCTTTAACAGGACAAGTGCGTTCAATCTCACCCAACTGGA GGAAAATTAGCCCGAGCCGGGGTGCCCCCGCAATGCTAGATTTTCAAGACAGATGGGACCATGTCTC CGATATGACAATATTTAAGTCTAAGTGACGGGACGCATTACAACGTATTATAAAATCCATATGTGTGTCTTC

- Tree relating samples to each other
- Branches with known length
- Mutation model











Hasegawa-Kishino-Yano model

$$\mathbf{Q} = \{q_{ij}\} = \begin{pmatrix} - & \pi_C & \kappa \pi_G & \pi_T \\ \pi_A & - & \pi_G & \kappa \pi_T \\ \kappa \pi_A & \pi_C & - & \pi_T \\ \pi_A & \kappa \pi_C & \pi_G & - \end{pmatrix} \mu$$



Hasegawa-Kishino-Yano model

$$\mathbf{Q} = \{q_{ij}\} = \begin{pmatrix} - & \pi_C & \kappa \pi_G & \pi_T \\ \pi_A & - & \pi_G & \kappa \pi_T \\ \kappa \pi_A & \pi_C & - & \pi_T \\ \pi_A & \kappa \pi_C & \pi_G & - \end{pmatrix} \mu$$

We do not know μ and we substitute it with the expected number of nucleotide substitutions per site. For example we could assume that the mean rate of substitutions is 1.0

$$\mu = \frac{1}{\sum_{i \in \{A,C,G,T\}} \pi_i q_{ii}}$$

- Transition/Transversion
- Stationary distribution

Let us consider a specific example of a rate matrix, with all of the parameters of the model taking specific values. For example, if we use the HKY85 model and fix the parameters to $\kappa = 5$, $\pi_A = 0.4$, $\pi_C = 0.3$, $\pi_G = 0.2$, and $\pi_T = 0.1$, we get the following matrix of instantaneous rates

$$\mathbf{Q} = \{q_{ij}\} = \begin{pmatrix} -0.886 & 0.190 & 0.633 & 0.063 \\ 0.253 & -0.696 & 0.127 & 0.316 \\ 1.266 & 0.190 & -1.519 & 0.063 \\ 0.253 & 0.949 & 0.127 & -1.329 \end{pmatrix}$$

Let us consider a specific example of a rate matrix, with all of the parameters of the model taking specific values. For example, if we use the HKY85 model and fix the parameters to $\kappa = 5$, $\pi_A = 0.4$, $\pi_C = 0.3$, $\pi_G = 0.2$, and $\pi_T = 0.1$, we get the following matrix of instantaneous rates

$$\mathbf{Q} = \{q_{ij}\} = \begin{pmatrix} -0.886 & 0.190 & 0.633 & 0.063 \\ 0.253 & -0.696 & 0.127 & 0.316 \\ 1.266 & 0.190 & -1.519 & 0.063 \\ 0.253 & 0.949 & 0.127 & -1.329 \end{pmatrix}$$

Focusing on G $Q = \{q_{ij}\} = \begin{pmatrix} & & & & & & & \\ & & & & & & & \\ 1.266 & 0.190 & -1.519 & 0.063 \\ & & & & & & & \\ & & & & & & & \\ \end{pmatrix}$

Figure 2: Simulation under the HKY85 substitution process. A single realization of the substitution process under the HKY85 model when $\kappa = 5$, $\pi_A = 0.4$, $\pi_C = 0.3$, $\pi_G = 0.2$, and $\pi_T = 0.1$. The length of the branch is v = 0.5 and the starting nucleotide is G (light gray). A, The process starts in nucleotide G. B, The first change is 0.152 units up the branch. C, the change is from G to A (dark gray). The time at which the next change occurs exceeds the total branch length, so the process ends in state \mathbf{A}

We choose the next nucleotide

 $G \to A: \frac{1.266}{1.519} = 0.833, \quad G \to C: \frac{0.190}{1.519} = 0.125, \quad G \to T: \frac{0.063}{1.519} = 0.042$

In the example we picked an A, now we continue with A: pick time, pick next nucleotide



Figure 2: Simulation under the HKY85 substitution process. A single realization of the substitution process under the HKY85 model when $\kappa = 5$, $\pi_A = 0.4$, $\pi_C = 0.3$, $\pi_G = 0.2$, and $\pi_T = 0.1$. The length of the branch is v = 0.5 and the starting nucleotide is G (light gray). A, The process starts in nucleotide G. B, The first change is 0.152 units up the branch. C, the change is from G to A (dark gray). The time at which the next change occurs exceeds the total branch length, so the process ends in state \mathbf{A} .

The following table summarizes the results of 100 simulations, each of which started with the nucleotide G:

	Starting	E	nding	Number of	of	
	Nucleotid	e Nuc	eleotide	Replicate	\mathbf{s}	
-	G		А	27		
Starting with G	G		С	10		
	G		G	59		
	G		Т	4		
	Ending			Ending		
			Nucleotide			
		А	. C	G	Т	
Starting with A,C,G,T	A	0.6	0.1	3 0.20	0.00	
Start	ing C	0.1	3 0.7	0 0.07	0.10	
Nucleo	otide G	0.2	0.1	0 0.59	0.04	
	Т	0.1	2 0.3	0 0.08	0.50	

$$\mathbf{P}(v) = e^{\mathbf{Q}v}$$

For the case we have been simulating, the exact transition probabilities (to four decimal places) are

$$\mathbf{P}(0.5) = \{p_{ij}(0.5)\} = \begin{pmatrix} 0.7079 & 0.0813 & 0.1835 & 0.0271 \\ 0.1085 & 0.7377 & 0.0542 & 0.0995 \\ 0.3670 & 0.0813 & 0.5244 & 0.0271 \\ 0.1085 & 0.2985 & 0.0542 & 0.5387 \end{pmatrix}$$

Stationary Distribution

$\mathbf{P}(0.00) =$	$ \left(\begin{array}{c} 1.000\\ 0.000\\ 0.000\\ 0.000 \end{array}\right) $	$0.000 \\ 1.000 \\ 0.000 \\ 0.000$	$0.000 \\ 0.000 \\ 1.000 \\ 0.000$	$\begin{array}{c} 0.000\\ 0.000\\ 0.000\\ 1.000 \end{array} \right)$	$\mathbf{P}(0.01) = \begin{pmatrix} 0.991\\ 0.003\\ 0.013\\ 0.003 \end{pmatrix}$	$\begin{array}{c} 0.002 \\ 0.993 \\ 0.002 \\ 0.009 \end{array}$	$\begin{array}{c} 0.006 \\ 0.001 \\ 0.985 \\ 0.001 \end{array}$	$\begin{array}{c} 0.001 \\ 0.003 \\ 0.001 \\ 0.987 \end{array} \right)$
P(0.10) =	$\left(\begin{array}{c} 0.919\\ 0.024\\ 0.113\\ 0.025\end{array}\right)$	$\begin{array}{c} 0.018 \\ 0.934 \\ 0.018 \\ 0.086 \end{array}$	$0.056 \\ 0.012 \\ 0.863 \\ 0.012$	$\begin{pmatrix} 0.006 \\ 0.029 \\ 0.006 \\ 0.877 \end{pmatrix}$	$\mathbf{P}(0.50) = \begin{pmatrix} 0.708\\ 0.106\\ 0.367\\ 0.109 \end{pmatrix}$	$\begin{array}{c} 0.081 \\ 0.738 \\ 0.081 \\ 0.299 \end{array}$	$0.184 \\ 0.054 \\ 0.524 \\ 0.054$	$\left(\begin{array}{c} 0.027 \\ 0.100 \\ 0.027 \\ 0.539 \end{array} \right)$
P(1.00) =	$ \left(\begin{array}{c} 0.580\\ 0.188\\ 0.464\\ 0.188 \end{array}\right) $	$\begin{array}{c} 0.141 \\ 0.587 \\ 0.141 \\ 0.394 \end{array}$	$\begin{array}{c} 0.232 \\ 0.094 \\ 0.348 \\ 0.094 \end{array}$	$\begin{pmatrix} 0.047 \\ 0.131 \\ 0.047 \\ 0.324 \end{pmatrix}$	$\mathbf{P}(5.00) = \begin{pmatrix} 0.411\\ 0.383\\ 0.411\\ 0.383 \end{pmatrix}$	$\begin{array}{c} 0.287 \\ 0.319 \\ 0.287 \\ 0.319 \end{array}$	$\begin{array}{c} 0.206 \\ 0.192 \\ 0.206 \\ 0.192 \end{array}$	$\left(\begin{array}{c} 0.096 \\ 0.106 \\ 0.096 \\ 0.107 \end{array} \right)$
P(10.0) =	$\left(\begin{array}{c} 0.401 \\ 0.399 \\ 0.401 \\ 0.399 \end{array}\right)$	$\begin{array}{c} 0.299 \\ 0.301 \\ 0.299 \\ 0.301 \end{array}$	$\begin{array}{c} 0.200 \\ 0.199 \\ 0.200 \\ 0.199 \end{array}$	$\begin{array}{c} 0.099 \\ 0.100 \\ 0.099 \\ 0.100 \end{array} \right)$	$\mathbf{P}(100) = \begin{pmatrix} 0.400\\ 0.400\\ 0.400\\ 0.400 \end{pmatrix}$	$\begin{array}{c} 0.300 \\ 0.300 \\ 0.300 \\ 0.300 \end{array}$	$0.200 \\ 0.200 \\ 0.200 \\ 0.200$	$\begin{pmatrix} 0.100 \\ 0.100 \\ 0.100 \\ 0.100 \end{pmatrix}$

Stationary Distribution

$\mathbf{P}(0.00) =$	$ \left(\begin{array}{c} 1.000\\ 0.000\\ 0.000\\ 0.000 \end{array}\right) $	$0.000 \\ 1.000 \\ 0.000 \\ 0.000$	$0.000 \\ 0.000 \\ 1.000 \\ 0.000$	$\left(\begin{array}{c} 0.000\\ 0.000\\ 0.000\\ 1.000\end{array}\right)$	$\mathbf{P}(0.01) = \begin{pmatrix} 0.991\\ 0.003\\ 0.013\\ 0.003 \end{pmatrix}$	$\begin{array}{c} 0.002 \\ 0.993 \\ 0.002 \\ 0.009 \end{array}$	$\begin{array}{c} 0.006 \\ 0.001 \\ 0.985 \\ 0.001 \end{array}$	$\begin{array}{c} 0.001 \\ 0.003 \\ 0.001 \\ 0.987 \end{array} \right)$
P(0.10) =	$\left(\begin{array}{c} 0.919\\ 0.024\\ 0.113\\ 0.025\end{array}\right)$	$\begin{array}{c} 0.018 \\ 0.934 \\ 0.018 \\ 0.086 \end{array}$	$0.056 \\ 0.012 \\ 0.863 \\ 0.012$	$\begin{array}{c} 0.006 \\ 0.029 \\ 0.006 \\ 0.877 \end{array} \right)$	$\mathbf{P}(0.50) = \begin{pmatrix} 0.708\\ 0.106\\ 0.367\\ 0.109 \end{pmatrix}$	$\begin{array}{c} 0.081 \\ 0.738 \\ 0.081 \\ 0.299 \end{array}$	$\begin{array}{c} 0.184 \\ 0.054 \\ 0.524 \\ 0.054 \end{array}$	$\left(\begin{array}{c} 0.027 \\ 0.100 \\ 0.027 \\ 0.539 \end{array} \right)$
P(1.00) =	$\left(\begin{array}{c} 0.580\\ 0.188\\ 0.464\\ 0.188\end{array}\right)$	$\begin{array}{c} 0.141 \\ 0.587 \\ 0.141 \\ 0.394 \end{array}$	$\begin{array}{c} 0.232 \\ 0.094 \\ 0.348 \\ 0.094 \end{array}$	$\begin{pmatrix} 0.047 \\ 0.131 \\ 0.047 \\ 0.324 \end{pmatrix}$	$\mathbf{P}(5.00) = \begin{pmatrix} 0.411\\ 0.383\\ 0.411\\ 0.383 \end{pmatrix}$	$\begin{array}{c} 0.287 \\ 0.319 \\ 0.287 \\ 0.319 \end{array}$	$\begin{array}{c} 0.206 \\ 0.192 \\ 0.206 \\ 0.192 \end{array}$	$\left(\begin{array}{c} 0.096 \\ 0.106 \\ 0.096 \\ 0.107 \end{array} \right)$
P(10.0) =	$\left(\begin{array}{c} 0.401 \\ 0.399 \\ 0.401 \\ 0.399 \end{array} \right)$	$\begin{array}{c} 0.299 \\ 0.301 \\ 0.299 \\ 0.301 \end{array}$	$\begin{array}{c} 0.200 \\ 0.199 \\ 0.200 \\ 0.199 \end{array}$	$\begin{array}{c} 0.099 \\ 0.100 \\ 0.099 \\ 0.100 \end{array} \right)$	$\mathbf{P}(100) = \begin{pmatrix} 0.400\\ 0.400\\ 0.400\\ 0.400 \end{pmatrix}$	$\begin{array}{c} 0.300 \\ 0.300 \\ 0.300 \\ 0.300 \end{array}$	$0.200 \\ 0.200 \\ 0.200 \\ 0.200$	$\begin{array}{c} 0.100 \\ 0.100 \\ 0.100 \\ 0.100 \end{array} \right)$

We used HKY85 model with these parameters to $\varkappa = 5$, $\pi_A = 0.4$, $\pi_C = 0.3$, $\pi_G = 0.2$, and $\pi_T = 0.1$

Simulating data

$$v_{2} = 0.1 \qquad v_{3} = 0.1 \qquad v_{4} = 0.2 \\ v_{1} = 0.3 \qquad v_{5} = 0.1 \qquad v_{6} = 0.1$$

Figure 4: The model tree for the simulations. We simulate data on this tree. The branch lengths are denoted v_i .

$$\mathbf{Q} = \{q_{ij}\} = \begin{pmatrix} -0.886 & 0.190 & 0.633 & 0.063 \\ 0.253 & -0.696 & 0.127 & 0.316 \\ 1.266 & 0.190 & -1.519 & 0.063 \\ 0.253 & 0.949 & 0.127 & -1.329 \end{pmatrix}$$

- we use our exponential random number generator, and add a tail to the tree and start with an A
- 2. No long tail, but pick root state from base frequencies
- 3. We have only three different branch length, we calculate the probability matrix and us that to generate the nucleotides at the nodes
- 4. We generate all site patterns and then pick among them (this will be discussed when know about likelihood.



We choose the next nucleotide

 $G \to A: \frac{1.266}{1.519} = 0.833, \quad G \to C: \frac{0.190}{1.519} = 0.125, \quad G \to T: \frac{0.063}{1.519} = 0.042$

The third simulation method does away with the need to generate exponential random variables. It takes advantage of our knowledge of the stationary distribution (as does the second method), but also takes advantage of our ability to calculate transition probabilities. There are only three different lengths of branches on our model tree (0.1, 0.2, and 0.3). The transition probabilities are

$$\mathbf{P}(0.10) = \begin{pmatrix} 0.919 & 0.018 & 0.056 & 0.006\\ 0.024 & 0.934 & 0.012 & 0.029\\ 0.113 & 0.018 & 0.863 & 0.006\\ 0.025 & 0.086 & 0.012 & 0.877 \end{pmatrix} \mathbf{P}(0.20) = \begin{pmatrix} 0.851 & 0.035 & 0.100 & 0.011\\ 0.047 & 0.876 & 0.023 & 0.052\\ 0.201 & 0.035 & 0.750 & 0.011\\ 0.047 & 0.156 & 0.023 & 0.771 \end{pmatrix}$$
$$\mathbf{P}(0.30) = \begin{pmatrix} 0.795 & 0.051 & 0.135 & 0.017\\ 0.069 & 0.824 & 0.034 & 0.071\\ 0.270 & 0.051 & 0.659 & 0.017\\ 0.069 & 0.214 & 0.034 & 0.681 \end{pmatrix}$$

Instead of drawing exponential random variables, and generating the process continuously across the entire tree, our simulation jumps from node to node on the tree. First, we generate the nucleotide at the root of the tree (the first split leading to all four taxa) by drawing from the stationary distribution. Then, we use the transition probabilities to simulate from one end to the other of each branch on the tree. We start from the root of the tree, and simulate up the tree to progressively higher branches until we have simulated a nucleotide at each tip of the tree.

Algorithm to simulate data on a tree

Required: a tree with branch lengths and a mutation transition rate matrix

$$\mathbf{Q} = \{q_{ij}\} = \begin{pmatrix} -0.886 & 0.190 & 0.633 & 0.063 \\ 0.253 & -0.696 & 0.127 & 0.316 \\ 1.266 & 0.190 & -1.519 & 0.063 \\ 0.253 & 0.949 & 0.127 & -1.329 \end{pmatrix}$$

1. $\tau = 0$; nuc = {}

2. Do forever

3. find λ

if you are at **G** use the diagonal value row of q_G : 0.886

- 4. draw random number r between 0 and 1 $r_1 = 0.134$
- 5. calculate $t = -\log_e(r_1)\frac{1}{\lambda}$ e.g. $t = -\log(0.134)\frac{1}{0.886} = 0.98521$ 6. $\tau = \tau + t$ $\tau = 0 + 0.98521$
- 7. if $\tau < v$

return current nucleotide nuc

- 8. calculate change of nucleotide cumulative sum $s = [0, \frac{q_{GA}}{q_{GG}}, \frac{q_{GA}}{q_{GG}} + \frac{q_{GC}}{q_{GG}}, 1.0]$ $s = [0, \frac{1.266}{1.519}, \frac{1.266}{1.519} + \frac{0.190}{1.519}, 1.0]$
- 9. draw random number r between 0 and 1
- 10. pick interval in which r_2 lays it is in the interval (0.833, 0.958] and thus nuc=**C**

11. goto 2

$$v_{1} = 0.3$$

$$v_{5} = 0.1$$

$$v_{6} = 0.1$$

$$v_{6} = 0.1$$

 $r_2 = 0.912$

Here is a second round of the example above: we are at ${\bf C}$ now

- 1. do forever 2. $\lambda = 0.696$ using row q_C 3. $r_1 = 0.449$ 4. $t = -\log_e(0.449)\frac{1}{0.696} = 0.49964$ 5. $\tau = 0.98521 + 0.49964 = 1.48485$ 6. if $(\tau = 0.98521) < 10.0$ return current nucleotide nuc 7. $s = [0, \frac{0.253}{0.696}, \frac{0.253}{0.696} + \frac{0.127}{0.696}, 1.0]$ 8. $r_2 = 0.191$ 9. pick interval in which r_2 lays it is in the interval (0, 0.363] and thus nuc=A
- 10. goto 2

