assignment snippet: how to start with the root

```
def rootgen(rootnode, basefreq, numsites):
  йнн
  this function generates numsites nucleotide
  given a set of probabilities
  йнн
  r = [[0.0, '']]*4
  b = ['A', 'C', 'G', 'T']
  r[0] = [probs[0], b[0]]
  for i in range(len(r)-1):
      r[i+1] = [r[i][0]+basefreq[i+1],b[i+1]]
  for i in range(numsites):
      ran = random.uniform(0.,1.0)
      for ri in r:
          if (ran < ri[0]):
               rootnode.sequence.append(ri[1])
              break
```

Jukes-Cantor (JC) allows for a single parameter and has a transition matrix

$$Q = egin{pmatrix} -rac{3}{4}\mu & rac{1}{4}\mu & rac{1}{4}\mu & rac{1}{4}\mu \ rac{1}{4}\mu & -rac{3}{4}\mu & rac{1}{4}\mu & rac{1}{4}\mu \ rac{1}{4}\mu & rac{1}{4}\mu & -rac{3}{4}\mu & rac{1}{4}\mu \ rac{1}{4}\mu & rac{1}{4}\mu & -rac{3}{4}\mu & rac{1}{4}\mu \ rac{1}{4}\mu & rac{1}{4}\mu & rac{1}{4}\mu & -rac{3}{4}\mu \end{pmatrix}$$

The base frequencies $\pi_A, \pi_C, \pi_G, \pi_T$ are all the same and 0.25. There are only two types of changes possible, either one does not change or one changes. This results in two probabilities:

$$Prob(t)_{ii} = \frac{1}{4} + \frac{3}{4}e^{-\mu t}$$
(19)
$$Prob(t)_{ij} = \frac{1}{4} - \frac{1}{4}e^{-\mu t}$$
(20)



Paul O. Lewis (2017 Woods Hole Workshop in Molecular Evolution)



Maximum likelihood estimation

First 32 nucleotides of the $\psi\eta$ -globin gene of gorilla and orangutan:

gorilla GAAGTCCTTGAGAAATAAACTGCACACACTGG orangutan GGACTCCTTGAGAAATAAACTGCACACACTGG

$$L = \left[\left(\frac{1}{4} \right) \left(\frac{1}{4} + \frac{3}{4} e^{-4\alpha t} \right) \right]^{30} \left[\left(\frac{1}{4} \right) \left(\frac{1}{4} - \frac{1}{4} e^{-4\alpha t} \right) \right]^2$$



number of substitutions = rate × time



existing A changes to a T

Overall substitution rate is 3α , so the expected number of substitutions (v) is

$$v = 3\alpha t$$



On Tuesday, Tracy Heath will introduce models that allow separate estimation of rates and times, but without extra information/constraints, sequence data allow only estimation of the **number** of substitutions.

Evolutionary distances for several common models

Model	Expected no. substitutions: $v = \{r\}t$
JC69	$v = \{3\alpha\} t$
F81	$v = \{2\mu(\pi_R\pi_Y + \pi_A\pi_G + \pi_C\pi_T)\}t$
K80	$v = \{\beta(\kappa + 2)\} t$
HKY85	$v = \{2\mu \left[\pi_R \pi_Y + \kappa (\pi_A \pi_G + \pi_C \pi_T)\right]\} t$

In the formulas above, the overall rate r (in curly brackets) is a function of all parameters in the substitution model.

One substitution model parameter is always determined from the edge length; the others are usually global (i.e. same value applies to all edges).

Likelihood of an unrooted tree

(data shown for only one site)





Brute force approach would be to calculate L_k for all 16 combinations of ancestral states and sum them



Note use of the OR probability rule



Felsenstein, J. 1981. Evolutionary trees from DNA sequences: a maximum likelihood approach. *Journal of Molecular Evolution* **17**:368-376

Algorithm 1 Likelihood downpass algorithm

for all *i* do $h_i^{(q)} \leftarrow 0$ for all *j* do $h_i^{(q)} \leftarrow h_i^{(q)} + p_{ij}^{(q)}g_j^{(q)}$

end for

end for

for all i do $h_i^{(r)} \leftarrow 0$ for all j do $h_i^{(r)} \leftarrow h_i^{(r)} + p_{ij}^{(r)}g_j^{(r)}$ end for end for for all i do

$$g_i^{(p)} \leftarrow h_i^{(q)} h_i^{(r)}$$

end for

Rate Heterogeneity

Green Plant rbc*L*

First 88 amino acids (translation is for Zea mays)

	(green alga; land plant lineage)	AAAGATTACAGATTAACTTACTACTACTCCTGAGTATAAAACTAAAGATACTGACATTTTAGCTGCATTTCGTGTAACTCCA
Chiorella	(green alga)	CC.T
Volvox	(green alga)	TC.T
Conocephalum	(liverwort)	
Bazzania	(moss)	
Anthoceros	(hornwort)	
Osmunda	(fern)	TCGCCTG.GCGTG.ACAA.GC
Lycopodium	(club "moss")	.GGAC.TC.TC.TTGCACTC.GAAA.GT
Ginkgo	(gymnosperm; Ginkgo biloba)	
Picea	(gymnosperm; spruce)	GT
Iris	(flowering plant)	
Asplenium	(fern; spleenwort)	TCC.GTCCCCACGCCTCGATCGA.GC
Nicotiana	(flowering plant; tobacco)	GAGTCCCGTAGACAT
		.AAA
	.G.A. A. A. T. T. .A. G. T. T. .A. G. T. T. .GG. .G. T. .GG. .G. .G. .GG. .G. .G. .G. G. A. .G. .G. .G. G. <td< th=""><th>A. A. A. T. A. A. A. T. T. T. T. A. C.T. T. T. T. T. C. C. G. G. C. C. G. T. C. C. G. T. C. T. T.</th></td<>	A. A. A. T. A. A. A. T. T. T. T. A. C.T. T. T. T. T. C. C. G. G. C. C. G. T. C. C. G. T. C. T.

Paul O. Lewis (2017 Woods Hole Workshop in Molecular Evolution)

Site-specific rates

Each defined subset (e.g. gene, codon position) has its own relative rate

Subset 1	Subset 2
r_1 applies to subset 1 (e.g. sites 1 - 1000)	<i>r</i> ₂ applies to subset 2 (e.g. sites 1001-2000)
Relative rates have mean 1:	More generally:
$\frac{r_1 + r_2}{2} = 1$	$r_1 p(r_1) + r_2 p(r_2) = 1$

Site-specific rates

$L = \Pr(D_1|r_1) \cdots \Pr(D_{1000}|r_1) \ \Pr(D_{1001}|r_2) \cdots \Pr(D_{2000}|r_2)$

Gene 1

Gene 2





 $r_2 = 0.8$

Site-specific rates

JC69 transition probabilities that would be used for every site if rate *homogeneity* were assumed:



Site specific rates

JC69 transition probabilities that would be used for sites in gene 1:

$$P_{ii}(t) = \frac{1}{4} + \frac{3}{4}e^{-4r_1\alpha t}$$
$$P_{ij}(t) = \frac{1}{4} - \frac{1}{4}e^{-4r_1\alpha t}$$

JC69 transition probabilities that would be used for sites in gene 2:

$$P_{ii}(t) = \frac{1}{4} + \frac{3}{4}e^{-4r_2\alpha t}$$
$$P_{ij}(t) = \frac{1}{4} - \frac{1}{4}e^{-4r_2\alpha t}$$



Site-specific Approach



Good: costs less: need to buy just one coat for every person

<u>Bad:</u> every person in a group has to wear the same size coat, so the fit will be poor for some people if they are much bigger or smaller than the average size for the group in which they have been placed

Mixture Models

All relative rates applied to every site



$L_i = \Pr(D_i | r_1) \Pr(r_1) + \Pr(D_i | r_2) \Pr(r_2)$

Common examples { Invariable sites (I) model Discrete Gamma (G) model



Mixture Model Approach



<u>Good:</u> every person experiences better fit because they can choose the size coat that fits best <u>Bad:</u> costs more because two coats much be provided for each person

Invariable Sites Model

A fraction p_{invar} of sites are assumed to be invariable (i.e. rate = 0.0)



$$L_i = \Pr(D_i|r_1)p_{\text{invar}} + \Pr(D_i|r_2)(1-p_{\text{invar}})$$

 $r_1 = 0.0$ $r_2 = \frac{1}{1 - p_{\text{invaries}}}$

Allows for the possibility that any given site could be variable or invariable

Reeves, J. H. 1992. Heterogeneity in the substitution process of amino acid sites of proteins coded for by mitochondrial DNA. Journal of Molecular Evolution 35:17-31.

Invariable sites model

If site *i* is a *constant* site, both terms will contribute to the site likelihood:



$$L_i = \Pr(D_i|0.0)p_{invar} + \Pr(D_i|r_2)(1-p_{invar})$$

If site *i* is a *variable* site, there is no way to explain the data with a zero rate, so the first term is zero:



$$L_i = \underline{\Pr(D_i|0.0)p_{\text{invar}}} + \Pr(D_i|r_2)(1-p_{\text{invar}})$$

Discrete Gamma Model

No relative rate is exactly 0.0, and all are equally probable



site i

 $L = (\frac{1}{4}) \Pr(D_i | r_1) + (\frac{1}{4}) \Pr(D_i | r_2) + (\frac{1}{4}) \Pr(D_i | r_3) + (\frac{1}{4}) \Pr(D_i | r_4)$

Relative rates are constrained to a discrete gamma distribution Number of rate categories can vary (4 used here)

Yang, Z. 1993. Maximum-likelihood estimation of phylogeny from DNA sequences when substitution rates differ over sites. Molecular Biology and Evolution 10:1396-1401.

Yang, Z. 1994. Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: approximate methods. Journal of Molecular Evolution 39:306-314.

Paul O. Lewis (2017 Woods Hole Workshop in Molecular Evolution)

Relative rates in 4-category case



Gamma distributions

