Announcements

Problem Set 2 due Friday at midnight.

First exam:

- Tuesday, Sept 30, 7:30pm 9:30pm
- WH 5403
 - This exam is closed book. You may bring two pages (or one page, front and back) of your own notes.
 - The exam covers material covered in Lectures 1 through 9.
 - You will not need a calculator

Models of sequence evolution

• Nucleotide substitution models

Last Tuesday

• Applications of substitution models

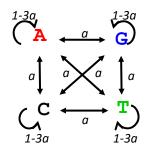
Today

Properties of DNA substitution models

- State space: $\{E_1 = A, E_2 = G, E_3 = C, E_4 = T\}$
- States are fully connected
- Transition probabilities: substitution frequencies (a, b, c, d ...)
- Implicitly also specifies stationary base frequencies: ϕ^* = (p_A, p_G, p_C, p_T)

3

Jukes-Cantor model (1969)



	A	С	G	T
A	1-3a	а	а	а
С	а	1-3a	а	а
G	а	а	1-3a	а
T	а	а	а	1-3a

Assumptions:

• All substitutions have equal probability

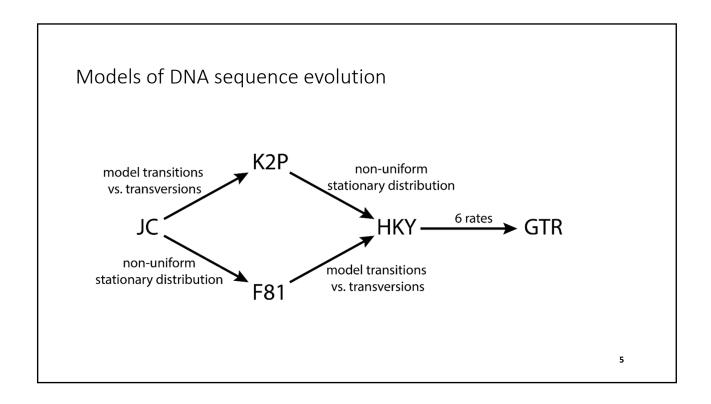
• Base frequencies are equal

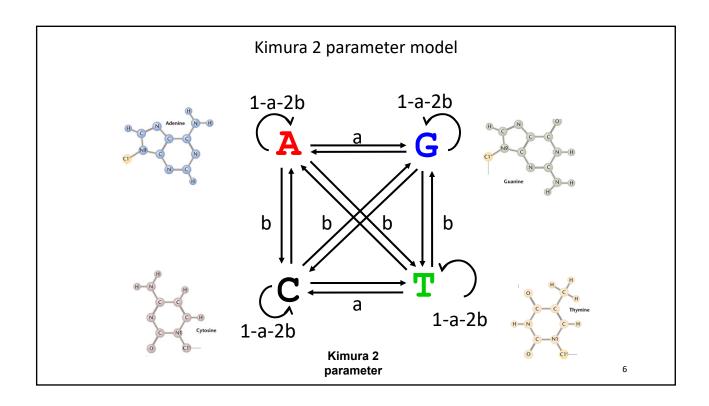
p(**A**)=0.25

p(G)=0.25

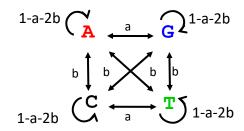
p(C)=0.25

p(T)=0.25





Kimura 2 parameter model (K2P) (1980)



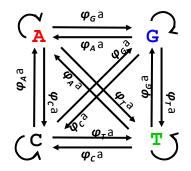
	A	С	G	T
A	1-a-2b	b	а	b
С	b	1-a-2b	b	а
G	а	b	1-a-2b	b
T	b	а	b	1-a-2b

Assumptions

- Transitions and transversions have different probabilities
- Base frequencies are equal

p(A)=0.25 p(G)=0.25 p(C)=0.25p(T)=0.25

Felsenstein (1981)

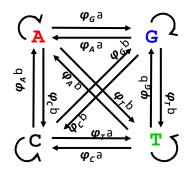


 $p(\mathbf{A}) = \boldsymbol{\varphi}_{\mathbf{A}}$ $p(\mathbf{G}) = \boldsymbol{\varphi}_{\mathbf{G}}$ $p(\mathbf{C}) = \boldsymbol{\varphi}_{\mathbf{C}}$ $p(\mathbf{T}) = \boldsymbol{\varphi}_{\mathbf{T}}$

- All substitutions have the same base rate
- Unequal base frequencies

 $p(A)\neq p(G)\neq p(C)\neq p(T)$

Hasegawa, Kishino & Yano (HKY) (1985)



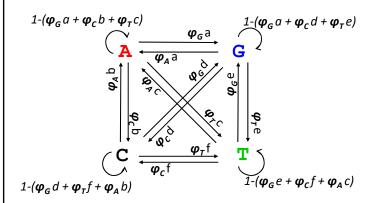
 $p(\mathbf{A}) = \boldsymbol{\varphi}_{\mathbf{A}}$ $p(\mathbf{G}) = \boldsymbol{\varphi}_{\mathbf{G}}$ $p(\mathbf{C}) = \boldsymbol{\varphi}_{\mathbf{C}}$ $p(\mathbf{T}) = \boldsymbol{\varphi}_{\mathbf{T}}$

- Transitions and transversions have different probabilities
- Unequal base frequencies

$$p(A)\neq p(G)\neq p(C)\neq p(T)$$

9

General Time Reversible model



- All six pairs have different substitution frequencies
- Unequal base frequencies $p(A) \neq p(G) \neq p(C) \neq p(T)$

Models of sequence evolution

• Nucleotide substitution models

Last Tuesday

- Applications of substitution models
 - Simulation
 - Estimating the probability of observing x aligned with y
 - · Estimating rates of evolution
 - Phylogenetic inference
 - Correcting for multiple substitutions at the same site

Next Tuesday

Deriving scoring matrices, S[x,y]

October

Questions to ask:

Given a site evolving according to Jukes Cantor with parameter a, what is the probability of observing x aligned with y?

...ATGCGAGGACTCXCAGTGA...
...ATGTGAGGTCTCYCAATGA...

Given an alignment of s_1 and s_2 with m observed mismatches, how many substitutions occured since the divergence of s_1 and s_2 ?

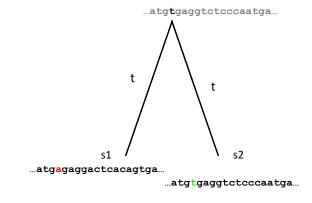
...CACATACGAAGATACGAACGAGC... ,..CAGATAGGAAGAGACGATCTAGC...

n nucleotides with *m* mismatches

We apply these models to the following scenario

Sequences s_1 and s_2 are DNA sequences of length n

- s_1 and s_2 have been diverging from a common ancestor for t million years (MY) according to <u>a given model with parameters a, b, c, ...</u>
 - This framework does not model indels
 - · Assumes site independence
 - · Parameters are estimated from data



We apply these models to the following scenario

Sequences s_1 and s_2 are DNA sequences of length n

- s_1 and s_2 have been diverging from a common ancestor for t million years (MY) according to the **Jukes Cantor (JC) model with parameter** α
 - This framework does not model indels
 - · Assumes site independence
 - All substitutions are equally probable
 - Parameter a is estimated from data

