

Molecular Evolution

Attribution: The development of the logical arguments building up the Jukes and Cantor and the Kimura models, and measuring sequence divergence with those models, have been taken in large part from

Fundamentals of Molecular Evolution,

Grauer, D and Li, W-H

Sinauer Associates, Mass 2000

Questions for Today

- What is the probability of finding any specific nucleotide at any given site?
- How many mutations have occurred since sequences diverged
- What is the rate of mutation

Bioinformatics *vis à vis* Biology

- We can use theory to estimate the number of mutations
- We need paleontological evidence to determine the passage of time

Question 1

- Examine some models of mutation
- Show how to determine ultimate probabilities of specific nucleotides appearing, given mutation rate(s)
 - Show that specific nucleotide probabilities are stationary in the limit
 - Analytic approach: By calculus to understand the underlying theory
 - Numeric Approach: By Markov Chains to process simple numerical examples
 - Linear Algebraic approach: By transition matrix eigenstructure to handle complex models and complicated matrices

RECALL

- A (purine) binds (2 H bonds) to T (pyrimidine)
- C (pyrimidine) binds (3 H bonds) to G (purine)

Possible mutations:

purine → purine A → G, G → A (transition)

pyrimidine → pyrimidine C → T, T → C (transition)

purine → pyrimidine A → T, A → C, G → T, G → C (transversion)

pyrimidine → purine T → A, C → A, T → G, C → G (transversion)

Note that there are twice as many transversions as transitions possible

The First Question, again:

What is the probability of a specific nucleotide in a specific position?

A Calculus Solution

Nucleotide Substitution

A Simple Model: Jukes and Cantor

Let α = the probability of substitution from one nucleotide to another in one time click. By hypothesis, whether it is a purine-purine, purine -pyrimidine, or pyrimidine-pyrimidine substitution, is not considered in this model; α is the same for all mutations in the Jukes and Cantor model.

If we choose some fixed nucleotide position in the DNA, and begin with a given nucleotide, WLOG, Adenine (A), then

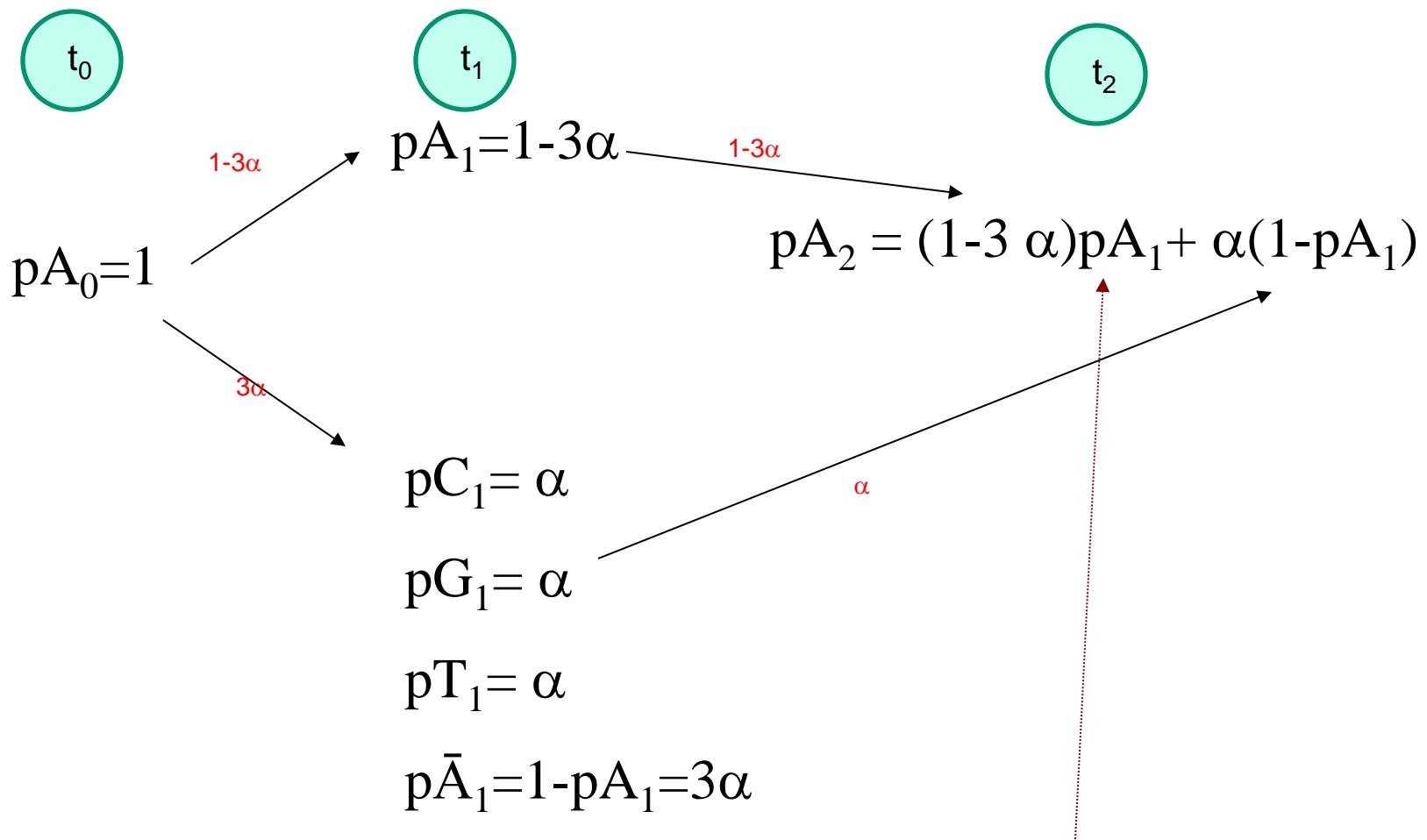
- at the next time click, **the probability of substituting to C,G,or T is α each, or 3α total.**
- The probability of remaining in **A is $1-3\alpha$.**

In the time click that follows (click 2), there are 2 ways to produce A:

- A Remains as A, where it had been in click 1

OR

- C,T,G (extant in time click 1) mutates to A



$(1-3\alpha)pA_1 = (1-3\alpha)^2$ but we leave it in the expanded form so that we can write a recursion

THE ITERATION

$$pA_{t+1} = (1 - \alpha)pA_t + \alpha(1 - pA_t)$$

Moving and re-arranging terms, we derive a difference equation

$$pA_{t+1} = (1-3\alpha)pA_t + \alpha(1-pA_t)$$

$$pA_{t+1} - pA_t = -3\alpha pA_t + \alpha(1-pA_t)$$

$$\Delta pA_t = -3\alpha pA_t + \alpha(1-pA_t)$$

$$\Delta pA_t = -4\alpha pA_t + \alpha$$

Making a difference equation into a differential equation

$$\Delta pA_t = -4\alpha pA_t + \alpha$$



$$\frac{d pA_t}{dt} = -4\alpha pA_t + \alpha$$

Solving the 1st order linear differential equation

$$pA_t = \int \left(\frac{d pA_t}{dt} \right) = \int (-4\alpha pA_t + \alpha)$$

$$pA_t = \frac{1}{4} + \left(pA_0 - \frac{1}{4} \right) e^{-4\alpha t}$$

We now have an expression for the probability of being in state A at time t

If we start in A (*i.e.* $pA_0=1$):

$$pA_t = \frac{1}{4} + \left(\mathbf{1} - \frac{1}{4} \right) e^{-4\alpha t} = \frac{1}{4} + \left(\frac{3}{4} \right) e^{-4\alpha t}$$

If we start in C,G,T (*i.e.* $pA_0=0$)

$$pA_t = \frac{1}{4} + \left(\mathbf{0} - \frac{1}{4} \right) e^{-4\alpha t} = \frac{1}{4} + \left(-\frac{1}{4} \right) e^{-4\alpha t}$$

Equilibrium

It is easy to see that the probability of ending up in state A asymptotically approaches an equilibrium probability (1/4) after a very large number of clicks. Likewise, the equilibrium probability is approached from a not-A initial condition.

since $\lim_{t \rightarrow \infty} \left[\left(\frac{3}{4} \right) e^{-4\alpha t} \right] = 0$

then $\lim_{t \rightarrow \infty} \left[\frac{1}{4} + \left(\frac{3}{4} \right) e^{-4\alpha t} \right] = \frac{1}{4}$ and, likewise, $\lim_{t \rightarrow \infty} \left[\frac{1}{4} + \left(-\frac{1}{4} \right) e^{-4\alpha t} \right] = \frac{1}{4}$

Stationarity Property

- In the language of stochastic processes, this asymptotically derived probability, $p=0.25$, is called the stationary probability.

2 key stochastic properties to be aware of are stationarity and ergodicity.

- A stationary process has the same moments regardless of which realization is sampled.
- An ergodic process has the same moments regardless of when a single realization is sampled.

A more sophisticated model (Kimura)

We know that transitions

purine \rightarrow purine A \rightarrow G, G \rightarrow A

or pyrimidine \rightarrow pyrimidine C \rightarrow T, T \rightarrow C

are more thermodynamically favorable
(hence more probable) than
transversions

purine \rightarrow pyrimidine A \rightarrow T, A \rightarrow C, G \rightarrow T, G \rightarrow C,

or pyrimidine \rightarrow purine T \rightarrow A, C \rightarrow A, T \rightarrow G, C \rightarrow G

The Kimura model captures that notion by
assigning separate probabilities to each

Kimura

- In the Kimura model, we will have a mutational rate α for a transition
- There will be a separate transversion rate β

where presumably $\alpha > \beta$

and where $\alpha + 2\beta$ must be < 1 (mathematical requirement)

and where $\alpha + 2\beta$ should be $\lll 1$ (evolutionary fact)

$$pAA_1 = 1 - \alpha - 2\beta$$

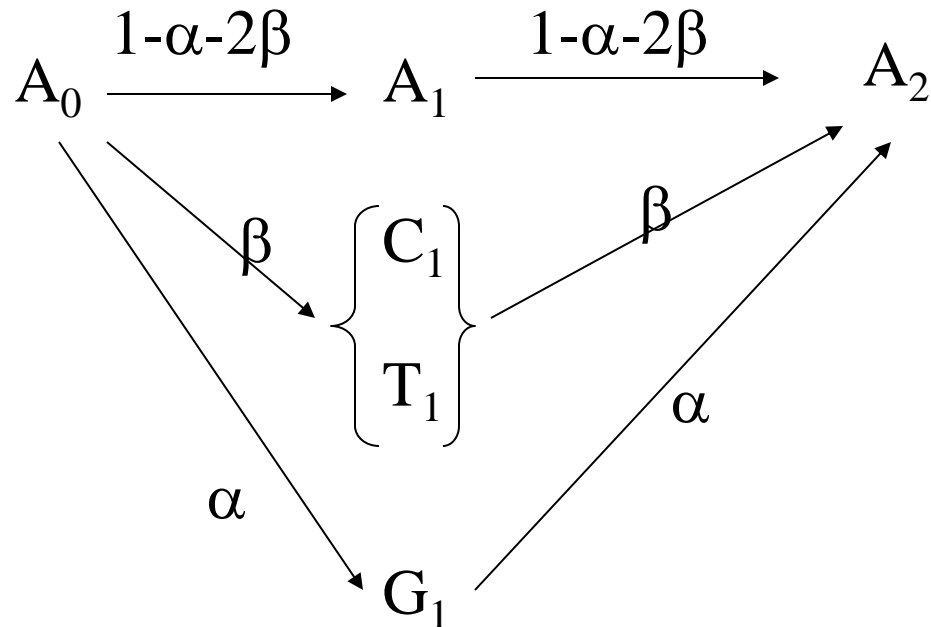
The probability of remaining in state A at click 2 is the sum of all the possible journeys

$$pAA_2 = (1 - \alpha - 2\beta) pAA_1$$

$$+ \beta pTA_1$$

$$+ \beta pCA_1$$

$$+ \alpha pGA_1$$



Again, generalizing

$$pAA_{t+1} = (1 - \alpha - 2\beta) pAA_t + \beta pTA_t + \beta pCA_t + \alpha pGA_t$$

Then moving thru the difference equation to the DE

$$\frac{d pAA_{t+1}}{dt} = -(\alpha + 2\beta) pAA_t + \beta pTA_t + \beta pCA_t + \alpha pGA_t$$

Solution

$$pAA_t = \frac{1}{4} + \frac{1}{4}e^{-\beta t} + \frac{1}{2}e^{-2(\alpha+\beta)t}$$

Equilibrium

Just as in the Jukes-Cantor model, in the limit, this much more complicated Kimura model is also stationary with asymptotic probability of 0.25

Another Viewpoint: A Numeric Solution

We can show this same convergence to stationary probabilities by:

- Creating a Transition Matrix (a Markov Matrix)
- Evolve the Transition Matrix to a vector of stationary probabilities using matrix arithmetic

DEFINITIONS

Markov property

*A next-state stochastic process variable, dependent only on the current state, with the property of "forgetting" all states before, has the **Markov property**.*

Markov process

A continuous stochastic process with the Markov property is called a Markov process. The probability of change from one state to the subsequent state is governed by a Markov propagator.

Markov chain

*A discrete process with the Markov property is called a Markov chain. The probability of change from one state to the subsequent state is called the **transition probability**.*

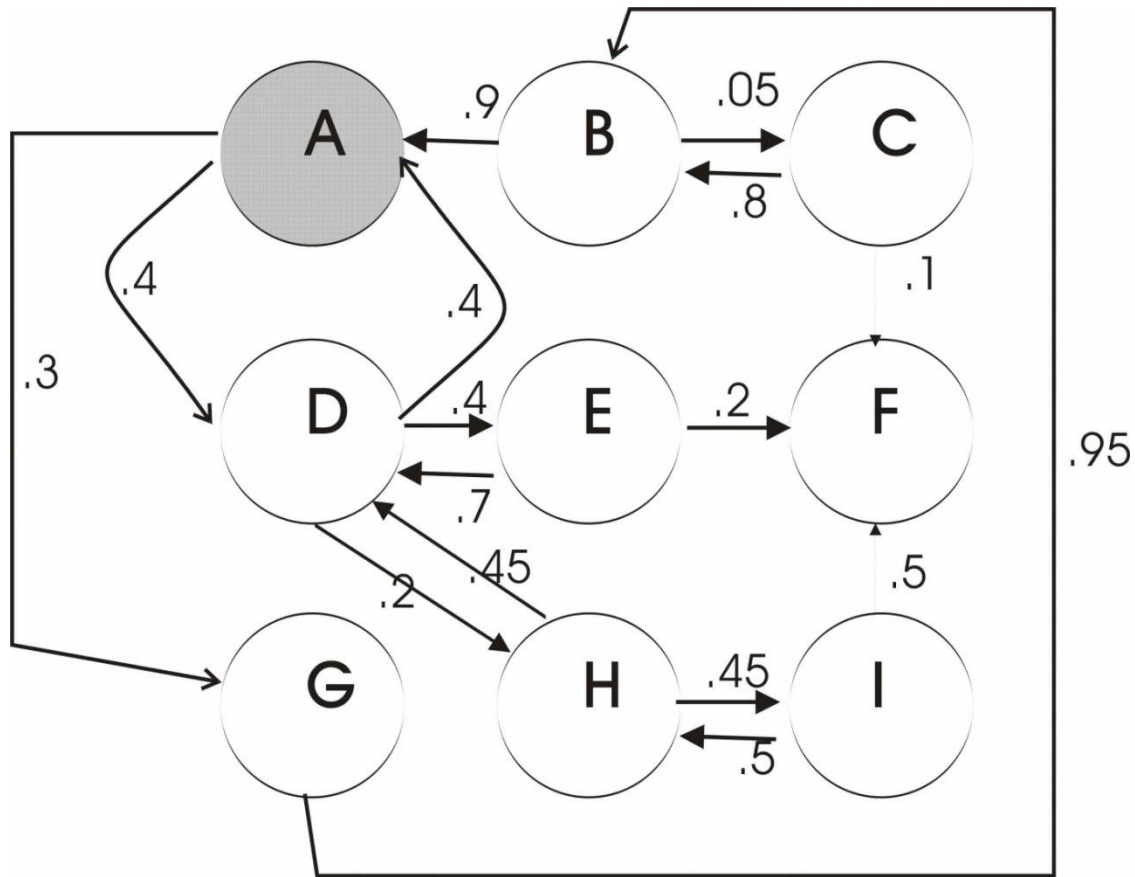
Equilibrium

- We can represent nucleotide change probabilities as transition probabilities in a Markov chain
- Specifically, we write a transition matrix for the Jukes and Cantor model, with each element representing a transition probability from one nucleotide to the next

Properties of Markov Chains

- Absorbing: Possible to get trapped
- Non-absorbing
 - Ergodic (irreducible) [visits everywhere at least once]
 - Regular: There is some power at which all transitions are (positive) nonzero

Markov chain in state A with transition probabilities



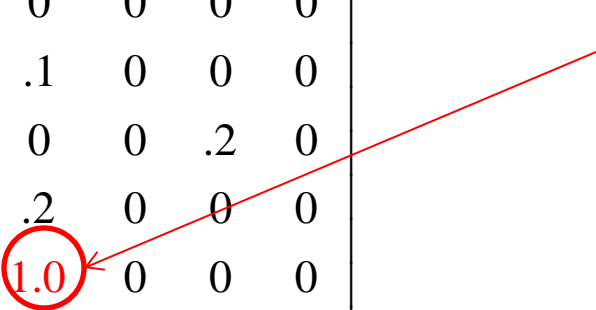
Note that all transition probabilities out of state F are nil.

This is an absorbing Markov chain; transition out of any state will ultimately reach F and become absorbed

Here is the transition matrix for the Markov Chain illustrated in the previous slide

	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>	<i>G</i>	<i>H</i>	<i>I</i>
<i>A</i>	.3	0	0	.4	0	0	.3	0	0
<i>B</i>	.9	.05	.05	0	0	0	0	0	0
<i>C</i>	0	.8	.1	0	0	.1	0	0	0
<i>D</i>	.4	0	0	0	.4	0	0	.2	0
<i>E</i>	0	0	0	.7	.1	.2	0	0	0
<i>F</i>	0	0	0	0	0	1.0	0	0	0
<i>G</i>	0	.95	0	0	0	0	.05	0	0
<i>H</i>	0	0	0	.45	0	0	0	.1	.45
<i>I</i>	0	0	0	0	0	.5	0	.5	0

Note the
absorbing state



Imagine one click of arbitrary time. We note the transition probabilities. In the next click, we operate on the state that is current, oblivious to any previous states, applying the transition probabilities that exist in this moment.

In other words, we multiply the current transition matrix by itself to get the 'new' transition probabilities. (Keep in mind that this multiplication is matrix, multiplication, not element-by-element ordinary multiplication.)

After only 20 time clicks, one can begin to observe the probability of staying in state F (transition $F \rightarrow F$) approaching 1.0

	A	B	C	D	E	F	G	H	I
A	0.1448	0.0528	0.0031	0.1190	0.0556	0.5207	0.0478	0.0383	0.0180
B	0.1513	0.0552	0.0032	0.1245	0.0580	0.4989	0.0500	0.0399	0.0188
C	0.1412	0.0516	0.0030	0.1160	0.0542	0.5325	0.0467	0.0373	0.0175
D	0.1190	0.0434	0.0025	0.0981	0.0456	0.6057	0.0394	0.0314	0.0149
E	0.0972	0.0355	0.0021	0.0798	0.0374	0.6781	0.0321	0.0258	0.0120
F	0	0	0	0	0	1.0000	0	0	0
G	0.1585	0.0579	0.0034	0.1302	0.0608	0.4754	0.0524	0.0419	0.0196
H	0.0861	0.0314	0.0018	0.0707	0.0332	0.7149	0.0284	0.0229	0.0107
I	0.0449	0.0164	0.0010	0.0371	0.0172	0.8511	0.0149	0.0119	0.0056

After 100 time clicks, a steady state is reached, and the process is nearly all absorbed into state F

	A	B	C	D	E	F	G	H	I
A	0.0045	0.0016	0.0001	0.0037	0.0017	0.9852	0.0015	0.0012	0.0006
B	0.0047	0.0017	0.0001	0.0038	0.0018	0.9846	0.0015	0.0012	0.0006
C	0.0043	0.0016	0.0001	0.0036	0.0017	0.9856	0.0014	0.0011	0.0005
D	0.0037	0.0013	0.0001	0.0030	0.0014	0.9879	0.0012	0.0010	0.0005
E	0.0030	0.0011	0.0001	0.0025	0.0011	0.9901	0.0010	0.0008	0.0004
F	0	0	0	0	0	1.0000	0	0	0
G	0.0049	0.0018	0.0001	0.0040	0.0019	0.9838	0.0016	0.0013	0.0006
H	0.0026	0.0010	0.0001	0.0022	0.0010	0.9912	0.0009	0.0007	0.0003
I	0.0014	0.0005	0.0000	0.0011	0.0005	0.9954	0.0005	0.0004	0.0002

The Probability Matrix

- We have given an example, and if you look closely, note that the row values add up to 1.0. These could well represent probabilities, but the sum-to-1 requirement alone is not sufficient for our purposes. We are seeking a structure called a 'Probability* Matrix' with these requirements:
 - Row entries sum to 1.0
 - Non-absorbing
 - Ergodic
 - Regular

*Frequently called a 'Stochastic Matrix'

Suppose we allowed transition out of F to C, E, or I, each with probability 0.01. Then the chain {A, B, C, D, E, F, G, H, I} is **no longer absorbing**. Further, it is **ergodic** (irreducible) because every element is eventually reachable. Yet further, the underlying Markov Chain is **regular**, since there is a power where all elements are positive (non-negative and non-zero) , although this is not obvious at this power of 1).

Here is the transition matrix:

$$R^1 = \begin{matrix} & \begin{matrix} A & B & C & D & E & F & G & H & I \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \\ G \\ H \\ I \end{matrix} & \left[\begin{array}{cccccccccc} .3 & 0 & 0 & .4 & 0 & 0 & .3 & 0 & 0 \\ .90 & .05 & .05 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & .8 & .1 & 0 & 0 & .1 & 0 & 0 & 0 \\ .4 & 0 & 0 & 0 & .4 & 0 & 0 & .2 & 0 \\ 0 & 0 & 0 & .7 & .1 & .2 & 0 & 0 & 0 \\ 0 & 0 & .01 & 0 & .01 & .97 & 0 & 0 & .01 \\ 0 & .95 & 0 & 0 & 0 & 0 & .05 & 0 & 0 \\ 0 & 0 & 0 & .45 & 0 & 0 & 0 & .1 & .45 \\ 0 & 0 & 0 & 0 & 0 & .5 & 0 & .5 & 0 \end{array} \right] \end{matrix}$$

Markov Transition Matrix

Note that this matrix mimics the transition rules of a Non-Deterministic Finite Automaton

- The Markov property is key
- The requirement for regularity is removed from the NDFA (although it would be bizarre machine design to have non-positive transition probabilities in an NDFA)
- The requirement for ergodicity is also removed

Fundamental Limit Theorem for Markov Chains

If P is a regular transition matrix then

$$\lim_{n \rightarrow \infty} P^n = U$$

where U is a probability matrix with each column a fixed probability and all rows equal

Fundamental Limit Theorem for Markov Chains at Work

Let's put the theorem to work!

Insist that our Markov Matrix is a probability matrix. Since it is already

- Ergodic
- Regular
- non absorbing

just exponentiate it (with matrix multiplication, of course)

The corresponding transition matrix is **converging** to the stationary matrix.

After 157 powers it is stationary up to >4 decimal places.

$$R^{157} = \begin{matrix} & \begin{matrix} A & B & C & D & E & F & G & H & I \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \\ G \\ H \\ I \end{matrix} & \left[\begin{array}{cccccccccc} .0980 & .0383 & .0094 & .0823 & .0438 & .6508 & .0304 & .0292 & .0197 \\ .0980 & .0383 & .0094 & .0823 & .0438 & .6508 & .0304 & .0292 & .0197 \\ .0980 & .0383 & .0094 & .0823 & .0438 & .6508 & .0304 & .0292 & .0197 \\ .0980 & .0383 & .0094 & .0823 & .0438 & .6508 & .0304 & .0292 & .0197 \\ .0980 & .0383 & .0094 & .0823 & .0438 & .6508 & .0304 & .0292 & .0197 \\ .0980 & .0383 & .0094 & .0823 & .0438 & .6508 & .0304 & .0292 & .0197 \\ .0980 & .0383 & .0094 & .0823 & .0438 & .6508 & .0304 & .0292 & .0197 \\ .0980 & .0383 & .0094 & .0823 & .0438 & .6508 & .0304 & .0292 & .0197 \\ .0980 & .0383 & .0094 & .0823 & .0438 & .6508 & .0304 & .0292 & .0197 \end{array} \right] \end{matrix}$$

NOTICE THAT THE STEADY STATE ROWS ARE IDENTICAL AND THE COLUMNS ARE PROBABILITIES. In looking back at the graph, it is not surprising that, in the limit, A is more likely to reach F than any other state, as is borne out by the steady state probability of A→F.

Stationary Probability: An Algebraic Approach to the Eigenstructure of the Transition Matrix

The Perron-Frobenius Theorem motivates the following development:

The stationary transition probabilities are given by the eigenvector corresponding to the eigenvalue with value 1 of the original transition matrix.

Linear Algebra Mini Review

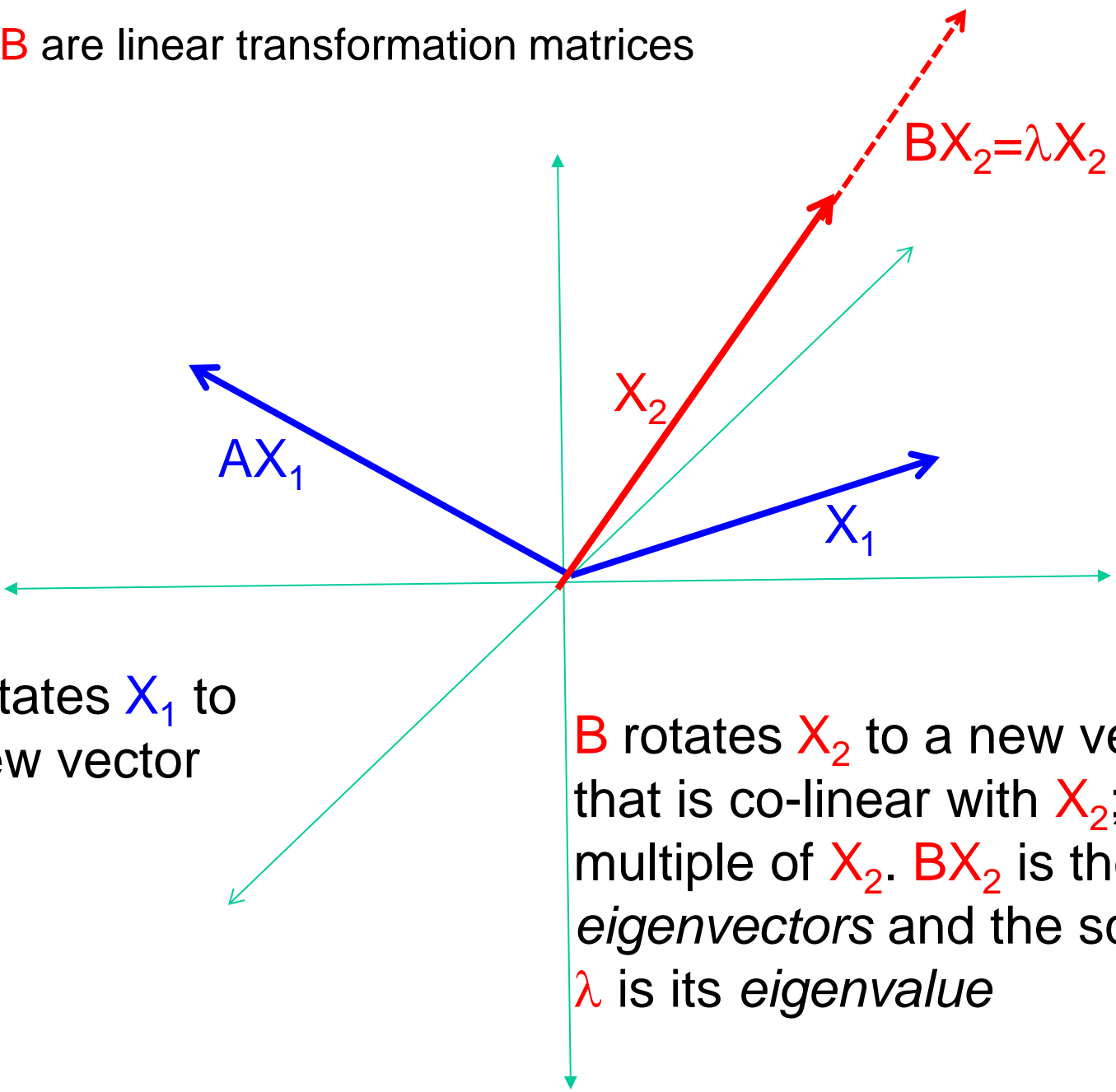
Let A be an $n \times n$ matrix.

λ is an eigenvalue of A if \exists non-zero vector X in \mathcal{R}^n such that

$$AX = \lambda X$$

If X is a non-zero vector satisfying the above, X is an ***eigenvector***. There can be more than one solution; there are as many solutions (eigenvectors) as the order (n) of the matrix, although the solutions may not be distinct.

A and B are linear transformation matrices



A rotates X_1 to a new vector AX_1

B rotates X_2 to a new vector BX_2 that is co-linear with X_2 ; it is a scalar multiple of X_2 . BX_2 is the one of the *eigenvectors* and the scalar multiple λ is its *eigenvalue*

When you do this rotation, you need to know two things about where you end up:

1. What vectors emerge from the rotation (the eigenvectors)
2. What scale factors emerge from the rotation (the eigenvalues)

Given $\mathbf{AX}=\lambda\mathbf{X}$, which represents a linear system of equations.

Then, re-writing this:

$\mathbf{AX}-\lambda\mathbf{IX}=\mathbf{0}$ where \mathbf{I} is the identity matrix

For this linear system to have a non-trivial solution, its determinant must be zero.

$$\text{So, } \begin{vmatrix} a_{11} - \lambda & a_{12} \\ a_{21} & a_{22} - \lambda \end{vmatrix} = 0 \quad \text{for a system of rank 2}$$

Solving, we get $(a_{11} - \lambda)(a_{22} - \lambda) - (a_{12})(a_{21}) = 0$

resulting in a quadratic (in a matrix of rank 2) equation in λ . This produces two values for λ in a system of rank 2, λ_1 and λ_2 .

These λ are the characteristic roots, or *eigenvalues*, of the system. There will be an *eigenvector* corresponding to each eigenvalue

Example

For $A = \begin{pmatrix} 1 & 1 \\ -2 & 4 \end{pmatrix}$ find

$$\begin{pmatrix} 1 & 1 \\ -2 & 4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \lambda \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

where λ is the eigenvalue

$$\begin{aligned}x_1 + x_2 &= \lambda x_1 \\ -2x_1 + 4x_2 &= \lambda x_2\end{aligned}$$

$$\begin{vmatrix} \lambda - 1 & -1 \\ 2 & \lambda - 4 \end{vmatrix} = 0$$

$$(\lambda - 1)(\lambda - 4) + 2 = 0$$

$$\lambda_1 = 2, \lambda_2 = 3$$

Eigenvalues may not be distinct

Finding the eigenvectors for λ_1

For $\lambda_1 = 2$

$$AX = 2X$$

$$\begin{pmatrix} 1 & 1 \\ -2 & 4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 2 \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

so

$$x_1 = x_2,$$

$x_2 = \text{any real number}$

is the eigenvector associated with λ_1

Linear Algebra and Markov Models

Review: There are three prerequisites in order to make classify a transition matrix as a probability matrix and to make it consistent with our biology

- The rows of the transition matrix are probability vectors (column entries in each row sum to 1.0)
- The Matrix is regular (there is at least one power of the matrix which has only positive entries)
- The chain is ergodic (irreducible) and aperiodic

Back to the Question:

What is the probability of a specific nucleotide at any given site?

Stationary Probability

If V is the steady state vector, then it stands to reason that $\mathbf{M}\mathbf{V}=\mathbf{V}$ where \mathbf{M} is the transition matrix

This is the same thing as $\mathbf{M}\mathbf{V}=\lambda\mathbf{V}$ where $\lambda=1$

Then it further stands to reason that the steady state vector of \mathbf{M} is the eigenvector \mathbf{V} corresponding to the eigenvalue $\lambda=1$ which is a probability vector (remember there is an infinitude of parallel eigenvectors)

The Jukes-Cantor Model of Nucleotide Substitution Probabilities

- Assumes that all nucleotides are equally likely to substitute for each other
- α is a probability reflecting how often a nucleotide substitutes in 1 generation

$$\begin{matrix} & A & C & G & T \\ A & \left(\begin{array}{cccc} 1-3\alpha & \alpha & \alpha & \alpha \\ \alpha & 1-3\alpha & \alpha & \alpha \\ \alpha & \alpha & 1-3\alpha & \alpha \\ \alpha & \alpha & \alpha & 1-3\alpha \end{array} \right) \\ C & & & & \\ G & & & & \\ T & & & & \end{matrix}$$

Jukes-Cantor Model

The eigenvalues are $1, (1 - 4\alpha), (1 - 4\alpha), (1 - 4\alpha)$

The corresponding eigenvectors are

$$\begin{pmatrix} r \\ r \\ r \\ r \end{pmatrix}, \begin{pmatrix} -r \\ 0 \\ 0 \\ r \end{pmatrix}, \begin{pmatrix} -r \\ 0 \\ r \\ 0 \end{pmatrix}, \begin{pmatrix} -r \\ r \\ 0 \\ 0 \end{pmatrix} \quad \text{Choose the eigenvector} \quad \begin{pmatrix} r \\ r \\ r \\ r \end{pmatrix}$$

corresponding to the eigenvalue 1

Jukes-Cantor Model

Because this is a probability vector, the coefficients must sum to 1. Because the elements are equal, the coefficients must be .25 each.

Thus, the specific eigenvector we require must be:

$$\begin{pmatrix} .25 \\ .25 \\ .25 \\ .25 \end{pmatrix}$$

IMPLICATION:

When our Jukes_Cantor model reaches equilibrium, the probabilities of transitions from each of the nucleotides are equal.

More Complicated Models Kimura

transitions (purine \rightarrow purine, pyrimidine \rightarrow pyrimidine)
[A \rightarrow G , C \rightarrow T with frequency α]

more likely than

transversions (purine \rightarrow pyrimidine, pyrimidine \rightarrow purine)
[A \rightarrow C , G \rightarrow T with frequency β]

NOTE simpler versions assume
freq A \rightarrow G = freq G \rightarrow A *etc*

Kimura Transition Matrix

$$\begin{matrix} & A & C & G & T \\ A & 1 - \alpha - 2\beta & \beta & \alpha & \beta \\ C & \beta & 1 - \alpha - 2\beta & \beta & \alpha \\ G & \alpha & \beta & 1 - \alpha - 2\beta & \beta \\ T & \beta & \alpha & \beta & 1 - \alpha - 2\beta \end{matrix}$$

Here, again, the eigenvector corresponding to the eigenvalue of 1 would represent the stationary transition probabilities. It is left as an exercise to evaluate the stationary transition probabilities in this model.