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

Today

- Markov chains
 - Review
 - Stationary distributions
- Models of sequence evolution
 - Nucleotide substitution models
 - (Amino acids in about 2 weeks)

Markov chain properties

In this course, we consider *finite, discrete, time-homogeneous* Markov chains:

- Number of states finite
- Independent variable is discrete
- *Time homogeneous:* The transmission matrix does not change over time.

that are

- *irreducible:* every state may be reached from every other state
- aperiodic:

There is no state that can only be visited multiples of m time steps, where m>1

Steady state behavior:

A finite, irreducible, aperiodic Markov chain has a unique stationary distribution, φ^* , such that

$$\varphi^* = \varphi^* P$$

Further,

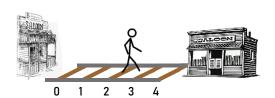
$$Q = \lim_{n \to \infty} P^n$$

has a solution:

$$Q = \begin{bmatrix} \varphi_1^* & \cdots & \varphi_S^* \\ \vdots & \ddots & \vdots \\ \varphi_1^* & \cdots & \varphi_S^* \end{bmatrix}$$

where $\varphi_1^* \dots \varphi_s^*$ is the limiting and stationary disribution.

Random walk with absorbing boundaries

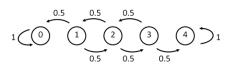


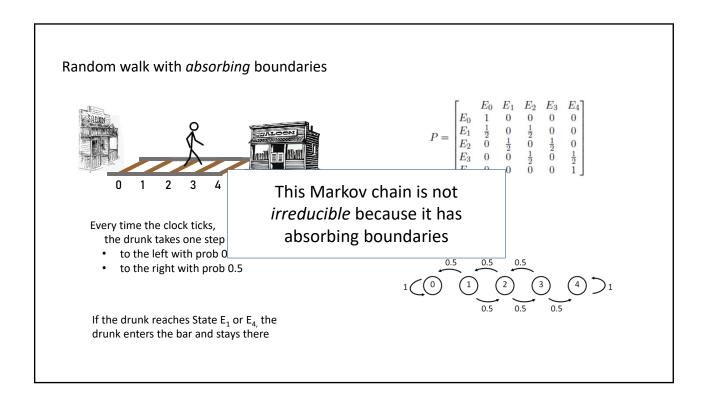
 $P = \begin{bmatrix} E_0 & E_1 & E_2 & E_3 & E_4 \\ E_0 & 1 & 0 & 0 & 0 & 0 \\ E_1 & \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 \\ E_2 & 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ E_3 & 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ E_4 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$

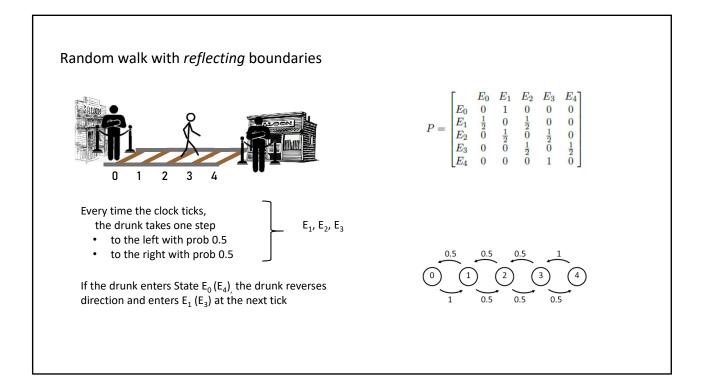
Every time the clock ticks, the drunk takes one step

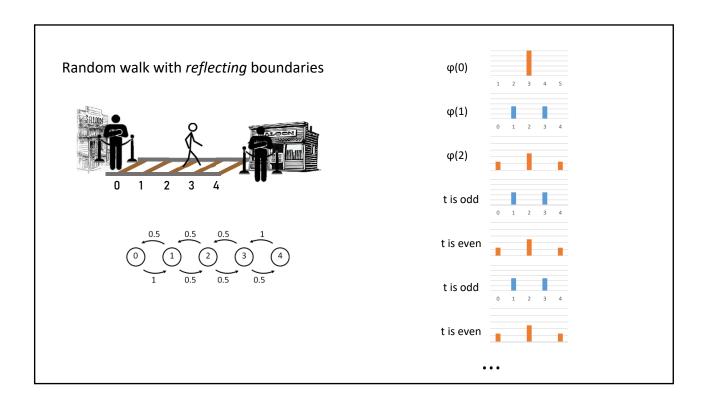
- to the left with prob 0.5
- to the right with prob 0.5

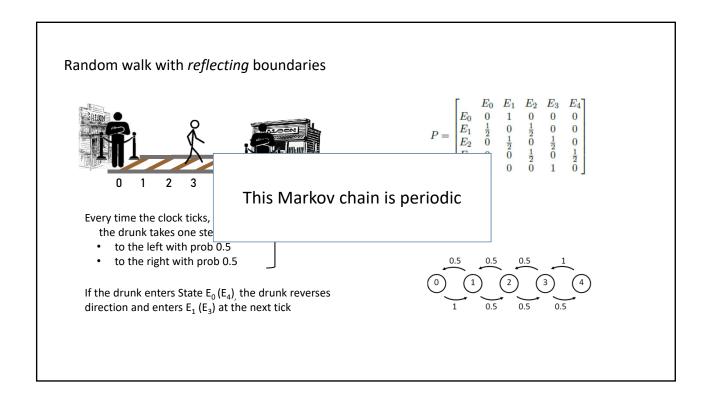
If the drunk reaches State E₁ or E₄, the drunk enters the bar and stays there



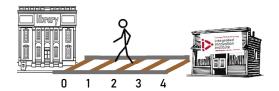








A third random walk



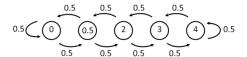
$$P = \begin{bmatrix} E_0 & E_1 & E_2 & E_3 & E_4 \\ E_0 & \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 \\ E_1 & \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 \\ E_2 & 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ E_3 & 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ E_4 & 0 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

Every time the clock ticks, the drunk takes one step

- to the left with prob 0.5
- to the right with prob 0.5

If the drunk enters State E_0 (E_4), the drunk

- rests with prob 0.5
- reverses direction and enters E₁ (E₃) with prob 0.5



A third random walk



Campublishing
Integrated
Integrated
Integrated
Institute

 $P = \begin{bmatrix} E_0 & E_1 & E_2 & E_3 & E_4 \\ E_0 & \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 \\ E_1 & \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 \\ E_2 & 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix}$

Every time the clock ticks,

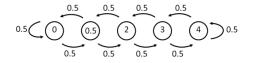
• to the left with prob

• to the right with prob 0.5

This Markov chain is *irreducible* and *aperiodic*. It has a unique stationary distribution.

If the drunk enters State $E_0(E_4)$, the drunk

- rests with prob 0.5
- reverses direction and enters E₁ (E₃) with prob 0.5



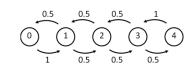
Random walk with absorbing boundaries

- Not irreducible because it has absorbing boundaries
- Does not have a unique stationary distribution.

$1 \underbrace{\begin{array}{c} 0.5 \\ 0.5 \\ 0.5 \\ 0.5 \\ 0.5 \\ \end{array}}_{0.5} \underbrace{\begin{array}{c} 0.5 \\ 0.5 \\ 0.5 \\ \end{array}}_{0.5} \underbrace{\begin{array}{c} 4 \\ 0.5 \\ \end{array}}_{1}$

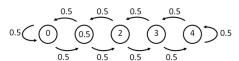
Random walk with reflecting boundaries

- This Markov chain is *periodic*
- It has a unique stationary distribution.
- It does not have a limiting distribution



Random walk with neither absorbing nor reflecting boundaries

- This Markov chain has a *unique stationary* distribution.
- It has a *limiting distribution* which is the same as the *stationary distribution*.



Today

- Announcements
- Markov chains
 - Review
 - Stationary distributions
- Models of sequence evolution
 - Nucleotide substitution models
 - (Amino acids in about 2 weeks)

Properties of DNA substitution models

1-
$$(a+b+c)$$
 1- $(a+c+e)$

A

a

G

T

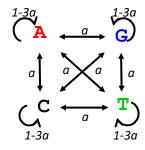
1- $(b+c+f)$

1- $(a+c+e)$

- 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^* = ($p_{_{\rm A}}$, $p_{_{\rm G}}$, $p_{_{\rm C}}$, $p_{_{\rm T}}$)

18

Jukes-Cantor model (1969)



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

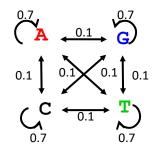
p(T)=0.25

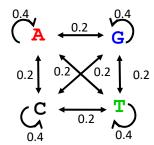
Assumptions:

- · All substitutions have equal probability
- Base frequencies are equal

Two Jukes Cantor models with different rates

All Jukes Cantor models have a single rate parameter, a, 0 < a < 1/3Different instances of the JC model can have different rates. Rates are typically learned from data



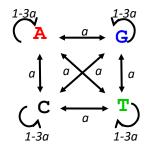


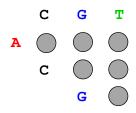
The model on the right changes twice as fast as the model on the left.

In both models, all substitutions are equally probable

20

Three representations of the Jukes Cantor model





	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

More nucleotide substitution models

Jukes Cantor

- Uniform substitution probabilities
- Uniform base frequencies

Substitution models can be extended by allowing

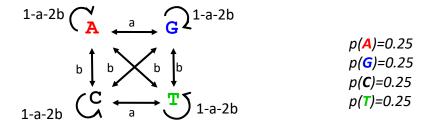
- different substitution probabilities for different base pairs
- non-uniform base frequencies

or both

22

A more complex model different probabilities for transitions and transversions 1-a-2b 1-a-2b 1-a-2b Kimura 2 parameter

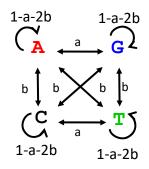
Kimura 2 parameter model (K2P) (1980)



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

24

Three representations of the Kimura 2-parameter model

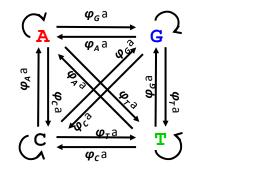




	A	С	G	T
A	1-a-2b	b	а	b
С	b	1-a-2b	b	а
G	а	b	1-a-2b	b
T	b	а	b	1-a-2b

Felsenstein (1981)

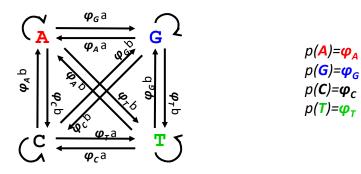
 $p(\mathbf{A}) = \boldsymbol{\varphi}_{\mathbf{A}}$ $p(\mathbf{G}) = \boldsymbol{\varphi}_{\mathbf{G}}$ $p(\mathbf{C}) = \boldsymbol{\varphi}_{\mathbf{C}}$ $p(\mathbf{T}) = \boldsymbol{\varphi}_{\mathbf{T}}$



- All substitutions have the same base rate
- Unequal base frequencies $p(A) \neq p(G) \neq p(C) \neq p(T)$

26

Hasegawa, Kishino & Yano (HKY) (1985)



- Transitions and transversions have different probabilities
- Unequal base frequencies $p(A) \neq p(G) \neq p(T)$

General Time Reversible model

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