

String Algorithms and Data Structures

Markov Chains

CS 199-225

November 29, 2022

Brad Solomon



UNIVERSITY OF
ILLINOIS
URBANA - CHAMPAIGN

Department of Computer Science

Learning Objectives

Introduce Markov Chains

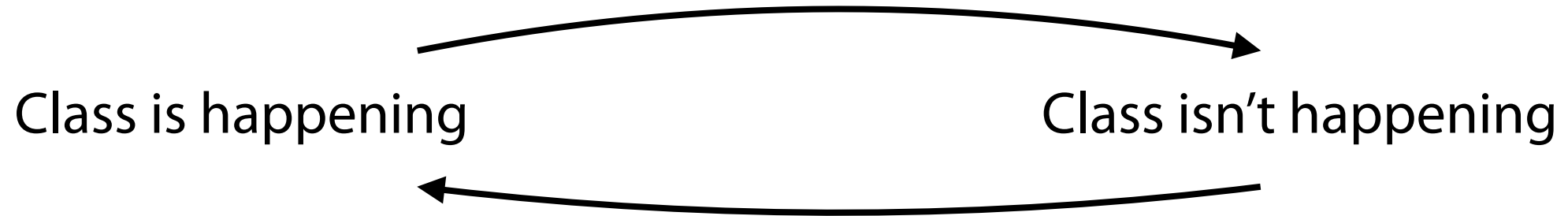
Define and determine stationary states

Identify common Markov Chain irregularities

Introduce Hidden Markov Models

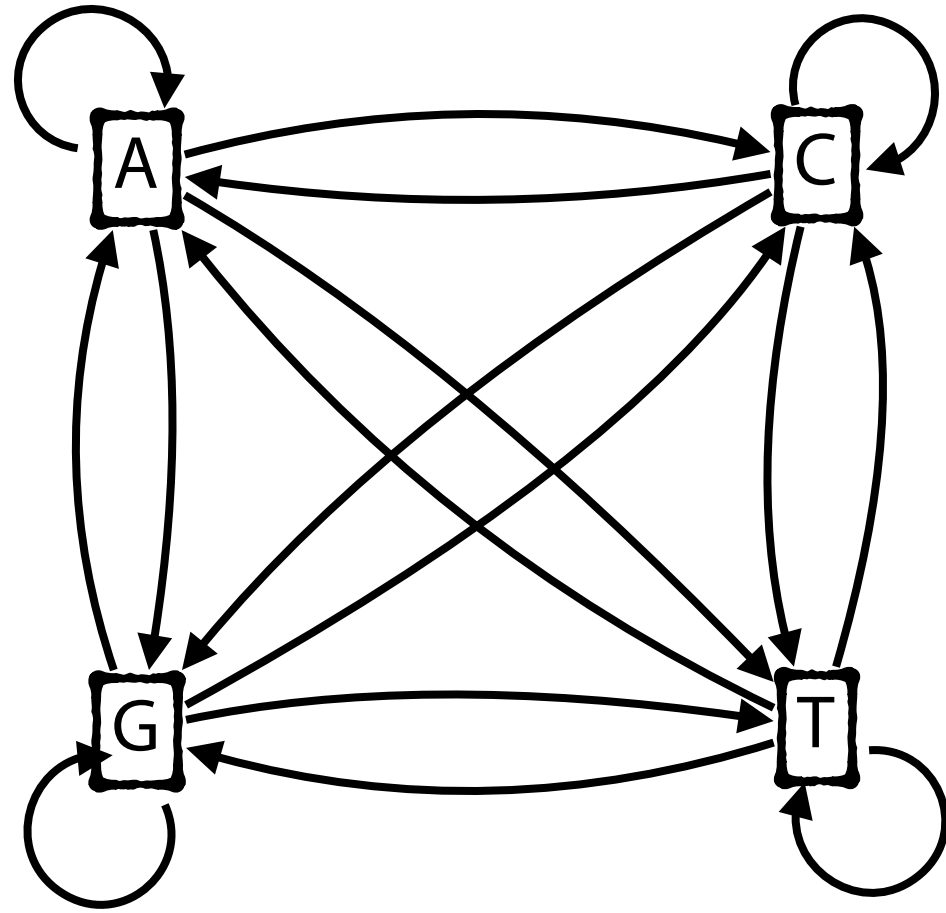
Modeling events with State Diagrams

A **state diagram** is a (usually weighted) directed graph where nodes are states and edges are transitions between them



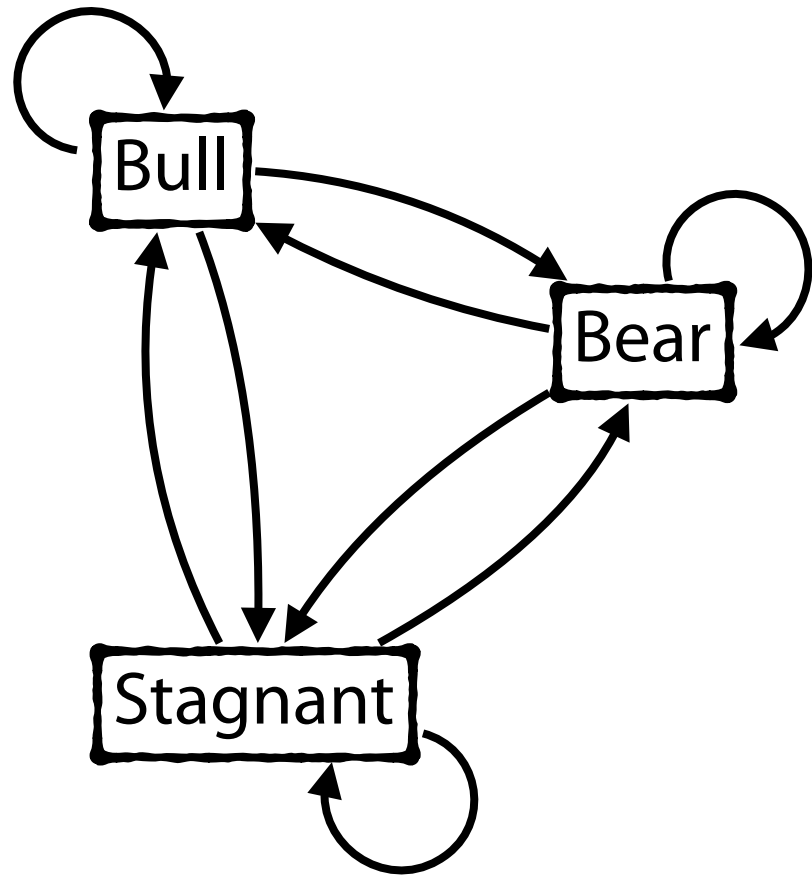
These diagrams are very useful in modeling many real world scenarios!

Sequence Modeling in Biology

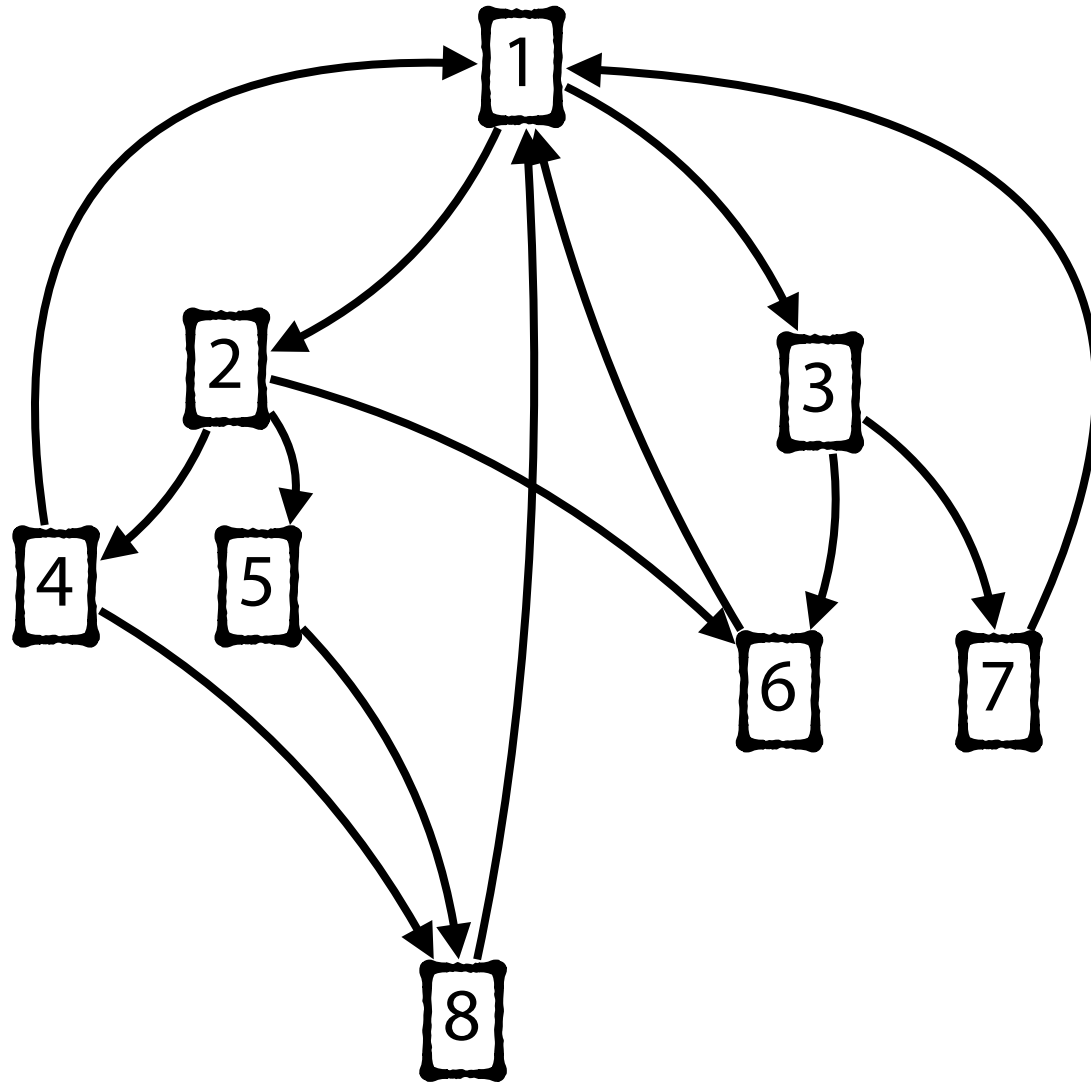


CATGACGTCGCGGACAACCCAGAATTGTCTTGAGCGATGGTAAGATCTAACCTCACTG
CTGGGGCTTTACTGATGTCATACCGTCTTGCACGGGGATAGAATGACGGTGCCCGTGT
ATTTTCTGAAAGTTACAGACTTCGATTA AAAAGATCGGACTGCGCGTGGGCCCCGGAG
TTTTTCGACGTGTCAAGGACTCAAGGGAATAGTTTGGCGGGAGCGTTACAGCTTCAATT
CGATAAAATTCAACTACTGGTTTCGGCCTAATAGGTCACGTTTTATGTGAAATAGAGGG
CCCTGGGTGTTCTATGATAAGTCCTGCTTTATAACACGGGGCGGTTAGGTTAAATGACT
ATCCAAGCGCCCCTAATTCTGTTCTGTTAATGTTTCATACCAATACTCACATCACATTA
AGCCCAGTCGCAAGGGTCTGCTGCTGTTGTCGACGCCTCATGTTACTCCTGGAATCTAC
GGTTAAGGCGTGTGATCGACGATGCAGGTATACATCGGCTCGGACCTACAGTGGTCGAT
TCGCGGTTTCGGCGCGTAGTTGAGTGCGATAACCCAACCGGTGGCAAGTAGCAAGAAGAC
AGACAACCTAACTAATAGTCTCTAACGGGGAATTACCTTTACCAGTCTCATGCCTCCAA
CAATGATATCGCCACAGAAAGTAGGGTCTCAGGTATCGCATACGCCGCGCCCCGGGTCC
GACAGTAGAGAGCTATTGTGTAATTCAGGCTCAGCATTTCATCGACCTTTCTGTTGTGA
TCTCGTCCGTAACGATCTGGGGGGCAAACCGAATATCCGTAATTCTCGTCCTACGGGTC
TGCGCGTGATCGTCAGTTAAGTTAAATTAATTCAGGCTACGGTAAACTTGTAGTGAGCT
ACGGGTTTCGCTACAGATGAACTGAATTTATACACGGGACAACCTCATCGCCCATTTGGGCG
AAAGTGGCAGATTAGGAGTGCTTGATCAGGTTAGCAGGTGGACTGTATCCAACAGCGCA
CCAAAGCGTTGTAGTGGTCTAAGCACCCCTGAACAGTGGCGCCCATCGTTAGCGTAGTA
AGGTGCGACATGGGGCCAGTTAGCCTGCCCTATATCCCTTGCACACGTTCAATAAGAGG
TTTTTAAATTAGGATGCCGACCCCATCATTGGTAACTGTATGTTTCATAGATATTTCTTC
AGCTGACACGCAAGGGTCAACAATAATTTCTACTATCACCCCGCTGAACGACTGTCTTT
CTTAGATTCGCGTCCTAACGTAGTGAGGGCCGAGTCATATCATAGATCAGGCATGAGAA
CACACGAGTTGTAAACA ACTTGATTGCTATACTGTAGCTACCGCAAGGATCTCCTACAT
ATCTGGATCCGAGTCAGAAATACGAGTTAATGCAAATTTACGTAGACCGGTGAAAACAC
AGACCGTAGTCAGAAGTGTGGCGCGCTATTCGTACCGAACCGGTGGAGTATACAGAATT
AGGAGCTCGGTCCCCAATGCACGCCAAAAAAGGAATAAAGTATTCAA ACTGCGCATGGT
CTATTATCCATCCGAACGTTGAACCTACTTCTCGGCTTATGCTGTCCTAACAGTATC
CGGCTGTGGATCTTAACGGCCACATTCTTAATTCCGACCGATCACCGATCGCCTTTCTT
ACTAAGTTATCCAGATCAAGGTTTGAACGGACTCGTATGACATGTGTGACTGAACCCGG
CTGTTTCAAGGCCTCTGCTTTGGTATCACTCAATATATTAGACCAGACAAGTGGCAAAA
CTAGGTATTACGCAACCGTTCGTAACATGCACTAAGGATAACTAGCGCCAGGGGGGCAT
AAAGACTACCCTATGGATTCCTTGGAGCGGGGACAATGCAGACCGGTTACGACACAATT
GGTATTATTAGCAAGACAATAAAGGACATTGCACAGAGACTTATTAGAATTCAACAAAC
GTGTTGGGTCGGGCAAGTCCCCGAAGCTCGGCCAAAAGATTCCGCCATGGAACCGTCTGG

Market Trends in Economics



PageRank in Graphs



Equilibrium State

1: $4/13$

2: $2/13$

3: $2/13$

4: $1/13$

5: $1/13$

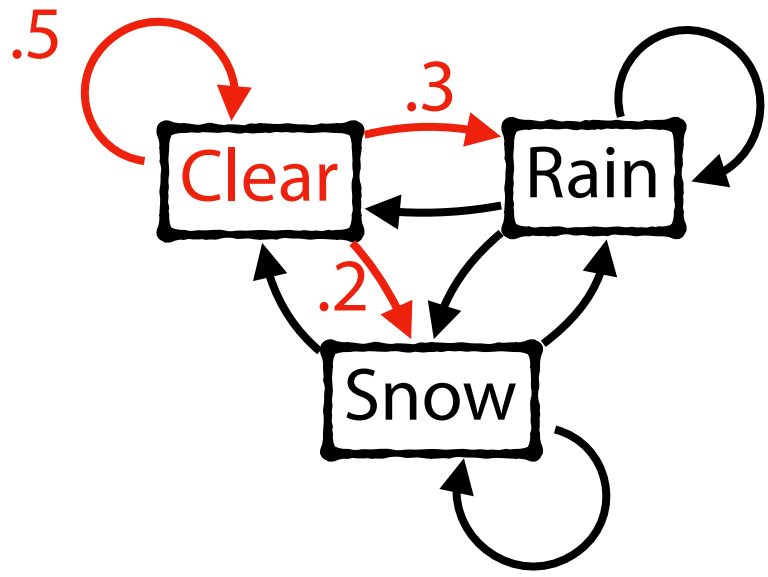
6: $1/13$

7: $1/13$

8: $1/13$

Markov Chain

A **finite Markov Chain** has a set of states S and a finite matrix M

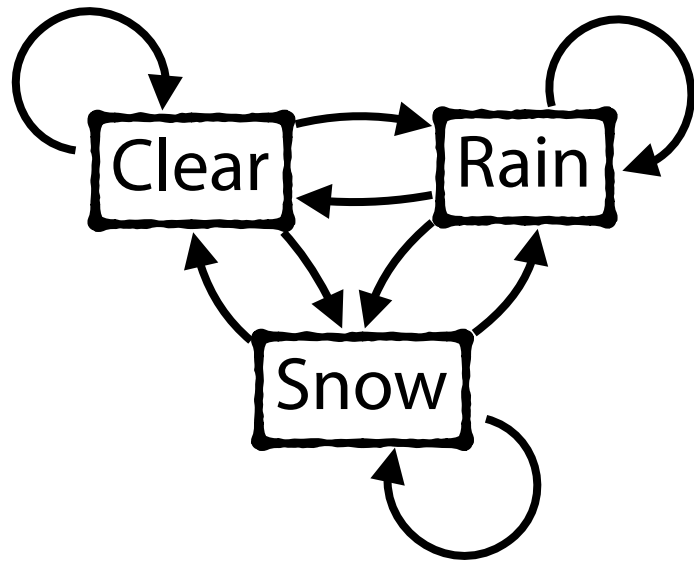


$$S = \{ \textit{Clear}, \textit{Rain}, \textit{Snow} \}$$

$$M = \begin{pmatrix} .5 & .3 & .2 \\ .4 & .1 & .1 \\ .2 & .4 & .7 \end{pmatrix}$$

Markov Chain

Given a Markov Chain and an initial state, all subsequent states can be represented either as **a series of random states** or a transition probability.



$$M = \begin{pmatrix} .5 & .3 & .2 \\ .5 & .4 & .1 \\ .2 & .1 & .7 \end{pmatrix}$$

$$X_0 = \text{Clear}$$

$$X_1 = \text{Clear}$$

$$X_2 = \text{Snow}$$

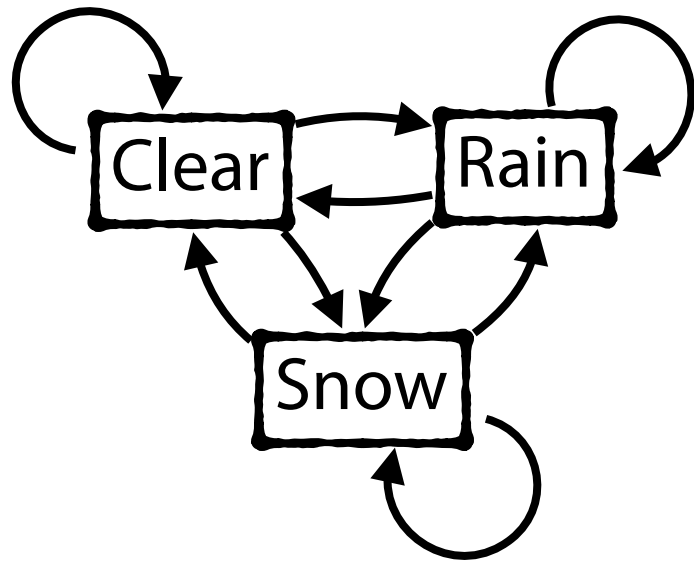
$$X_3 = \text{Snow}$$

$$X_4 = \text{Snow}$$

$$X_5 = \text{Rain}$$

Markov Chain

Given a Markov Chain and an initial state, all subsequent states can be represented either as a series of random states or a **transition probability**.



$$M = \begin{pmatrix} .5 & .3 & .2 \\ .5 & .4 & .1 \\ .2 & .1 & .7 \end{pmatrix}$$

$$M_0 = (.4 \quad .3 \quad .3)$$

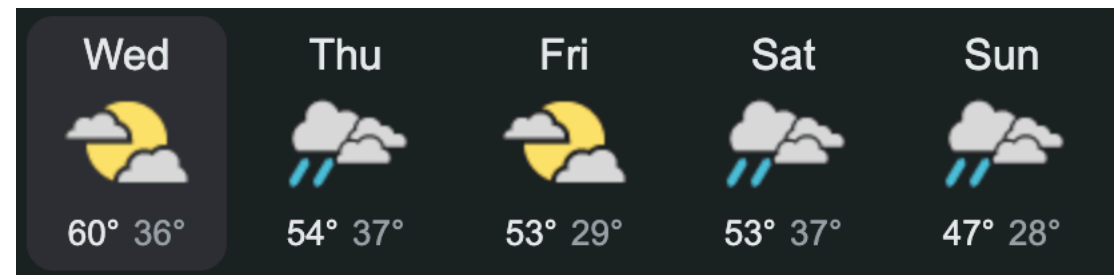
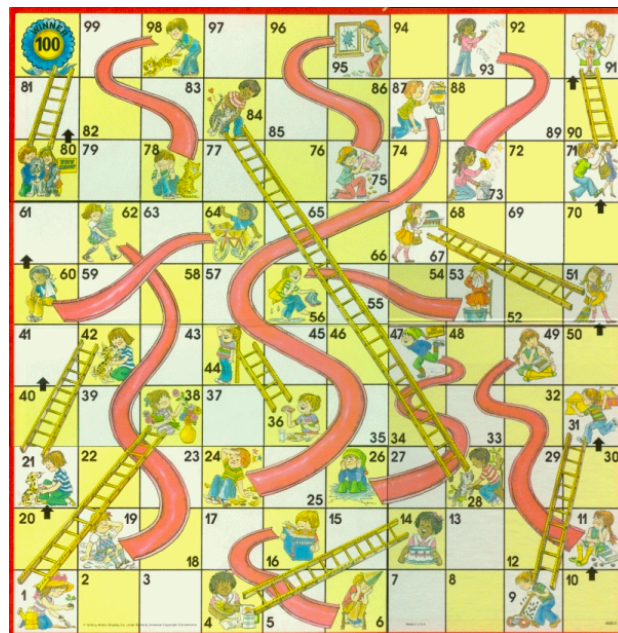
$$M_1 = (.41 \quad .27 \quad .32)$$

$$M_2 = (.404 \quad .263 \quad .333)$$

$$M_3 = (.401 \quad .259 \quad .340)$$

Markov Assumption

The probability of the next state depends only on our current state





Markov Assumption

Probability of state x_k depends only on previous state x_{k-1}

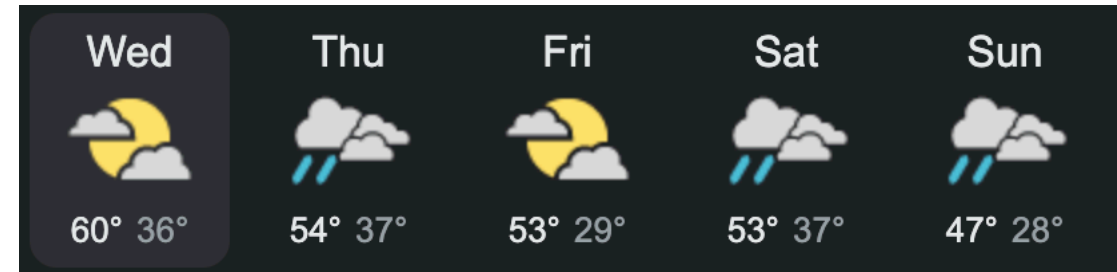
Ex: Let $x = \{C, R, C, R, R\}$

$$P(x) = P(x_k, x_{k-1}, \dots, x_1)$$

$$= P(x_k | x_{k-1}, \dots, x_1) P(x_{k-1}, \dots, x_1)$$

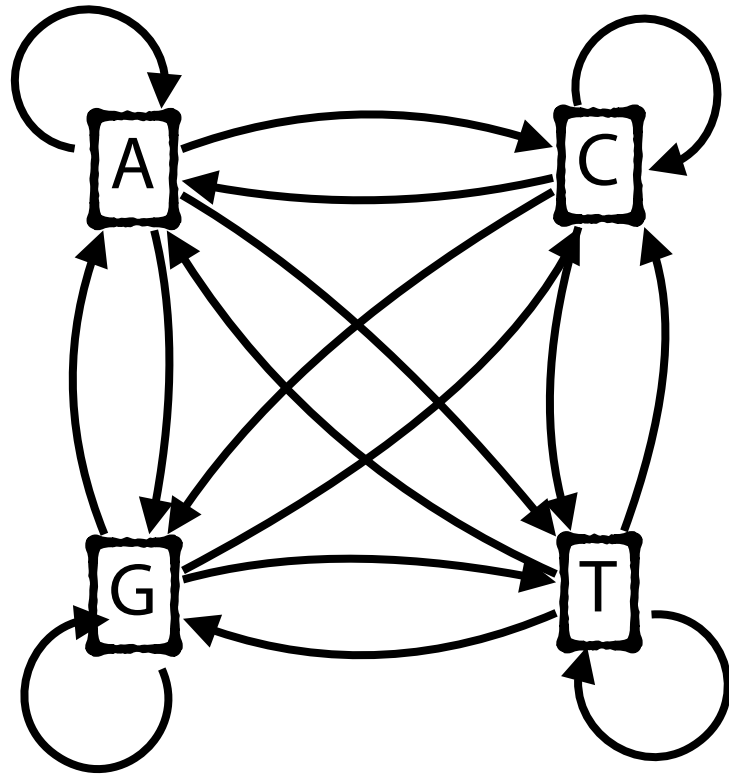
$$= P(x_k | x_{k-1}, \dots, x_1) P(x_{k-1} | x_{k-2}, \dots, x_1) \dots P(x_2 | x_1) P(x_1)$$

$$P(x) \approx$$



Markov Chain in Sequencing

Given a set of sequences, we can construct a model of transitions



$$P(A | A) = \# \text{ times } AA \text{ occurs} / \# \text{ times } AX \text{ occurs}$$

$$P(C | A) = \# \text{ times } AC \text{ occurs} / \# \text{ times } AX \text{ occurs}$$

$$P(G | A) = \# \text{ times } AG \text{ occurs} / \# \text{ times } AX \text{ occurs}$$

$$P(T | A) = \# \text{ times } AT \text{ occurs} / \# \text{ times } AX \text{ occurs}$$

$$P(A | C) = \# \text{ times } CA \text{ occurs} / \# \text{ times } CX \text{ occurs}$$

(etc)

where X is any base


Example by Ben Langmead

Markov Chain in Sequencing

Given a set of sequences, we can construct a model of transitions

```
>>> ins_conds, _ = markov_chain_from_dinucs(samp)
>>> print(ins_conds)
```

X_{i-1}	A	[[0.19152248, 0.27252589, 0.39998803, 0.1359636],		
	C	[0.18921984, 0.35832388, 0.25467081, 0.19778547],		
	G	[0.17322219, 0.33142737, 0.35571338, 0.13963706],		
	T	[0.09509721, 0.33836493, 0.37567927, 0.19085859]]		
	A	C	G	T
	X_i			
				$P(T G)$



Markov Chain in Sequencing

```

>>> ins_conds, _ = markov_chain_from_dinucs(samp)
>>> print(ins_conds)
A [[ 0.19152248, 0.27252589, 0.39998803, 0.1359636 ],
C [ 0.18921984, 0.35832388, 0.25467081, 0.19778547 ],
G [ 0.17322219, 0.33142737, 0.35571338, 0.13963706 ],
T [ 0.09509721, 0.33836493, 0.37567927, 0.19085859 ]]

```

X_{i-1}

A C G T

X_i

$x = GATC$

$$P(x) = P(x_4 | x_3) P(x_3 | x_2) P(x_2 | x_1) P(x_1)$$

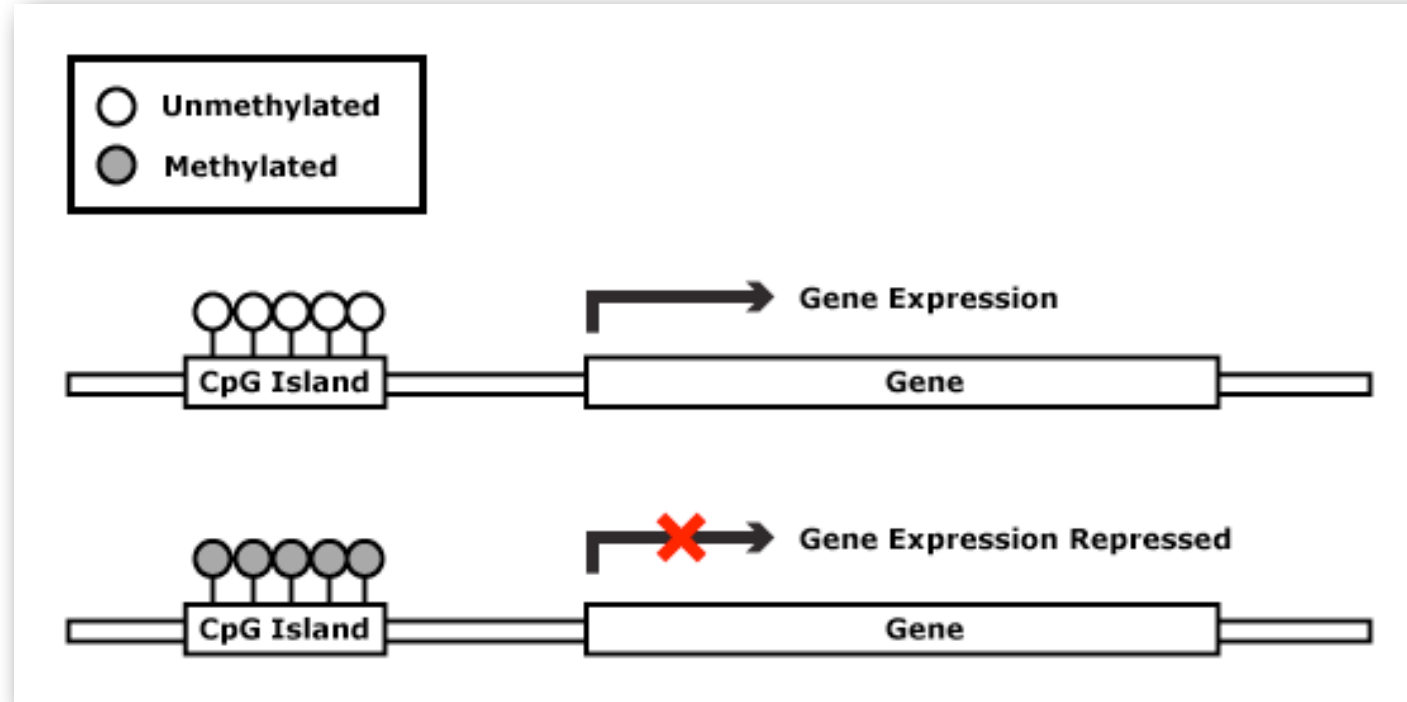
$$\begin{aligned}
 P(x) = & P(C | T) P(T | A) P(A | G) P(G) = 0.33836493 * 0.1359636 * 0.17322219 * 0.25 \\
 & = 0.001992
 \end{aligned}$$

Example by Ben Langmead

Markov Chain in Sequencing

We can use this same approach to predict a *label* in our sequences as well

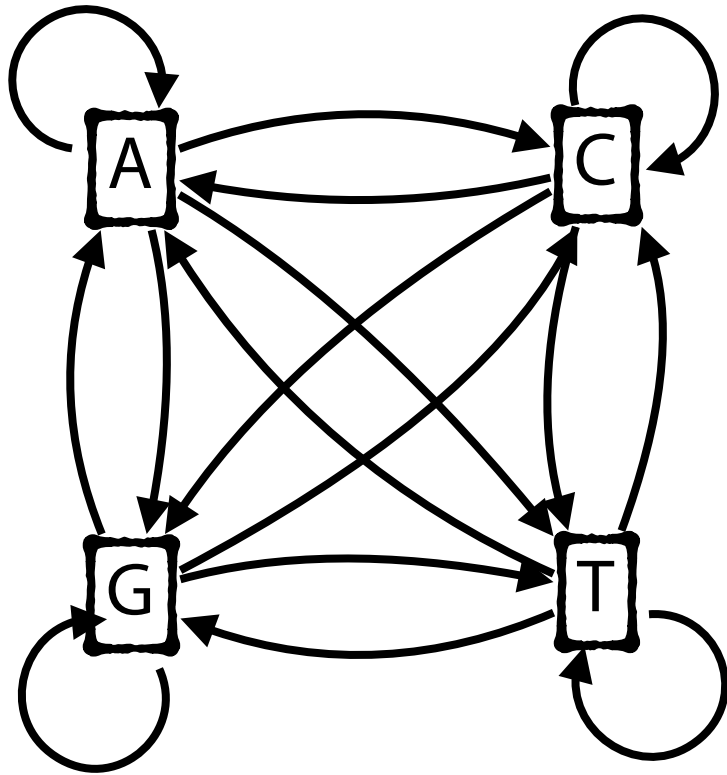
CpG island: part of the genome where CG occurs particularly frequently



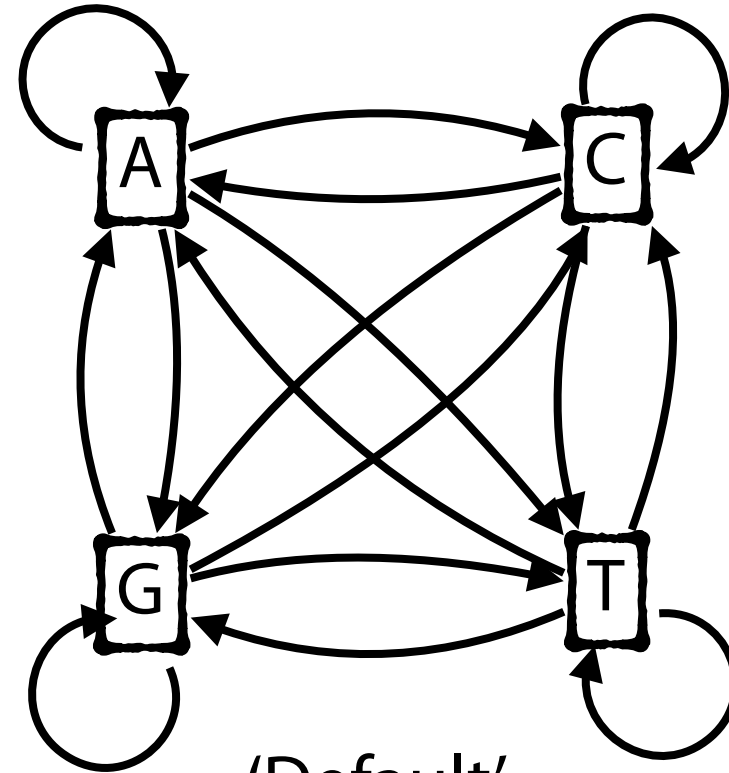
Example by Ben Langmead

Markov Chain in Sequencing

To predict a *label* of a sequencing region, make a Markov chain for both!



CpG Island



'Default'

Example by Ben Langmead


```

>>> cpg_conds, _ = markov_chain_from_dinucs(samp_cpg)
>>> print(cpg_conds)
    A [[ 0.19152248,  0.27252589,  0.39998803,  0.1359636 ],
    C [ 0.18921984,  0.35832388,  0.25467081,  0.19778547],
    G [ 0.17322219,  0.33142737,  0.35571338,  0.13963706],
    T [ 0.09509721,  0.33836493,  0.37567927,  0.19085859]]
>>> default_conds, _ = markov_chain_from_dinucs(samp_def)
>>> print(default_conds)
    A [[ 0.33804066,  0.17971034,  0.23104207,  0.25120694],
    C [ 0.37777025,  0.25612117,  0.03987225,  0.32623633],
    G [ 0.30257815,  0.20326794,  0.24910719,  0.24504672],
    T [ 0.21790184,  0.20942905,  0.2642385 ,  0.3084306 ]]

```

CpG

┆ A
┆ C
┆ G
┆ T

Default

┆ A
┆ C
┆ G
┆ T

Log ratio

┆ A
┆ C
┆ G
┆ T

```

>>> print(np.log2(cpg_conds) - np.log2(def_conds))
    A [[ -0.87536356,  0.59419041,  0.81181564, -0.85527103],
    C [ -0.98532149,  0.49570561,  2.64256972, -0.7126391 ],
    G [ -0.79486196,  0.68874785,  0.51821792, -0.79549511],
    T [ -1.22085697,  0.73036913,  0.48119354, -0.69736839]]

```

A

C

G

T

Markov Chain in Sequencing

```
>>> print(np.log2(cpg_conds) - np.log2(def_conds))
Xi-1 A [[ -0.87536356,  0.59419041,  0.81181564, -0.85527103],
        C [[ -0.98532149,  0.49570561,  2.64256972, -0.7126391 ],
        G [[ -0.79486196,  0.68874785,  0.51821792, -0.79549511],
        T [[ -1.22085697,  0.73036913,  0.48119354, -0.69736839]]
        A           C           G           T
        Xi
```

$x = \text{GATC}$

$$P(x) = P(x_4 | x_3) P(x_3 | x_2) P(x_2 | x_1) P(x_1)$$

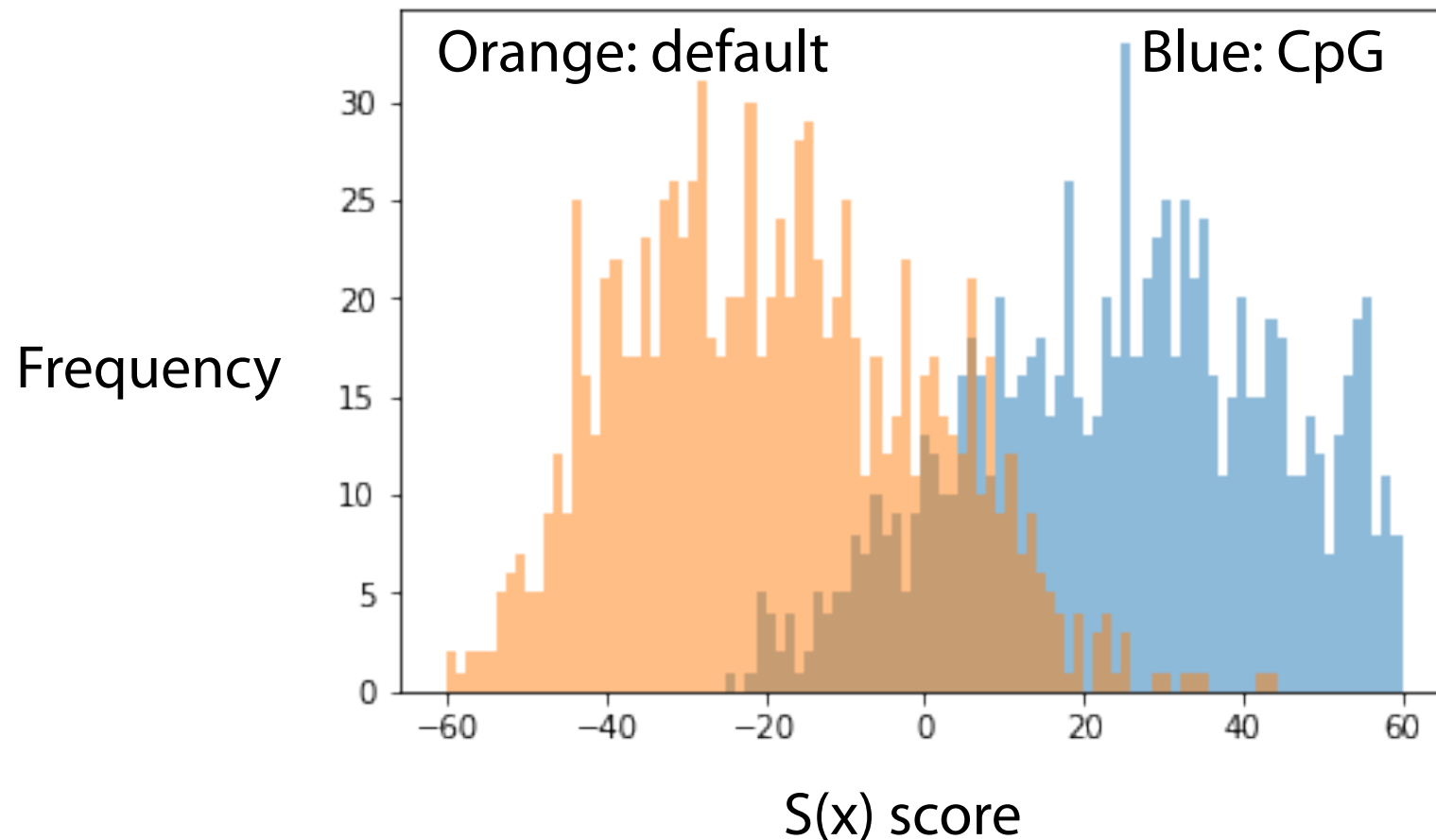
$$P(x) = P(\text{C} | \text{T}) P(\text{T} | \text{A}) P(\text{A} | \text{G}) P(\text{G}) = 0.73036913 + = -0.919763$$
$$-0.85527103 +$$
$$-0.79486196$$

Example by Ben Langmead

Markov Chain in Sequencing

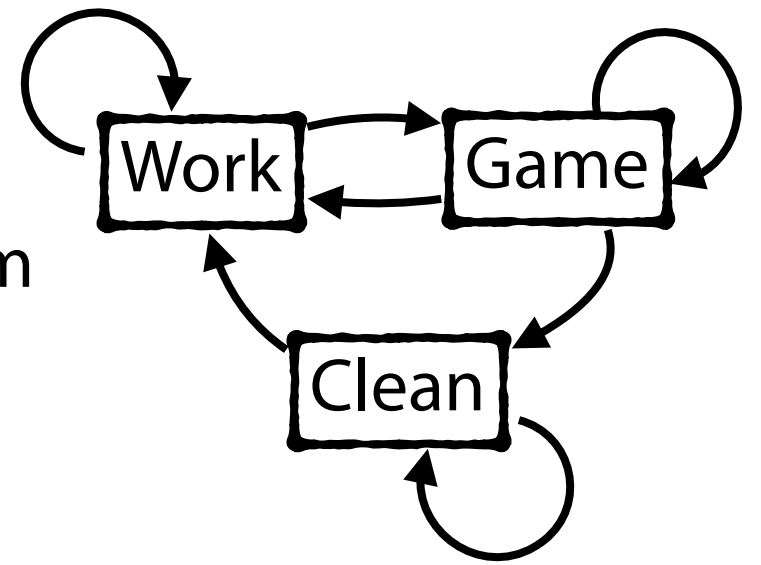


Drew 1,000 100-mers from inside CpG islands and another 1,000 from outside, and calculated $S(x)$ for all



Markov Chain Matrix

If I'm working at time 0, what is probability that I'm working at time t ?



Claim: $Pr(X_t = v | X_0 = u) = M^t[u, v]$

$$M = \begin{pmatrix} .4 & .6 & 0 \\ .1 & .6 & .3 \\ .5 & 0 & .5 \end{pmatrix}$$

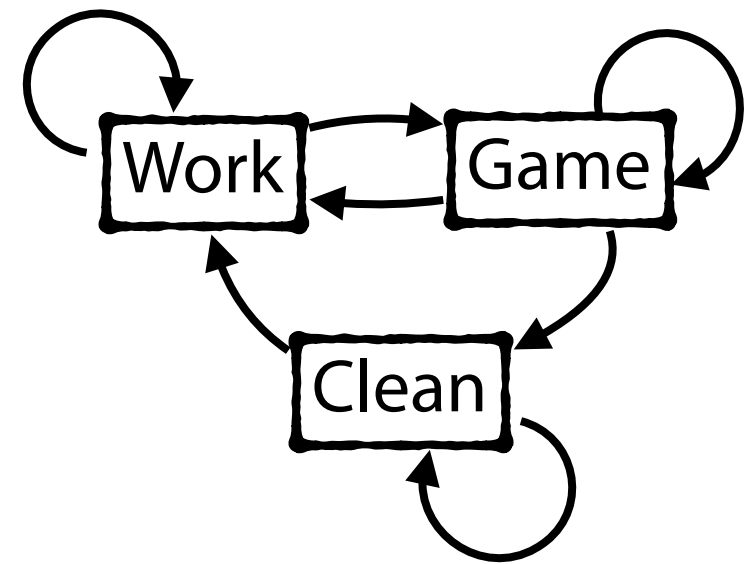
Markov Chain Matrix

Claim: $Pr(X_t = v | X_0 = u) = M^t[u, v]$

Base Case:

T=1:

T=2:



$$M = \begin{pmatrix} .4 & .6 & 0 \\ .1 & .6 & .3 \\ .5 & 0 & .5 \end{pmatrix}$$

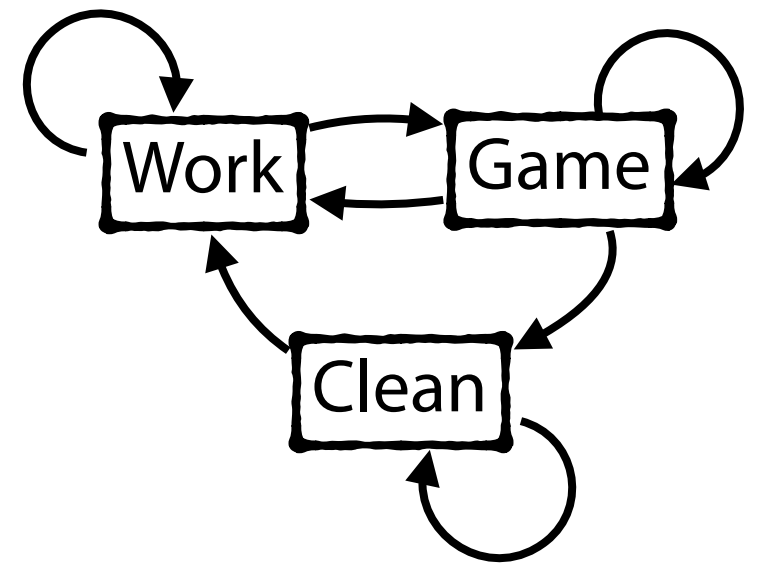
Markov Chain Matrix

Claim: $Pr(X_t = v | X_0 = u) = M^t[u, v]$

Induction:

Assume $Pr(X_{t-1} = v | X_0 = u) = M^{t-1}[u, v]$.

Show holds for $Pr(X_t = w | X_0 = u) = M^t[u, w]$



$$M = \begin{pmatrix} .4 & .6 & 0 \\ .1 & .6 & .3 \\ .5 & 0 & .5 \end{pmatrix}$$

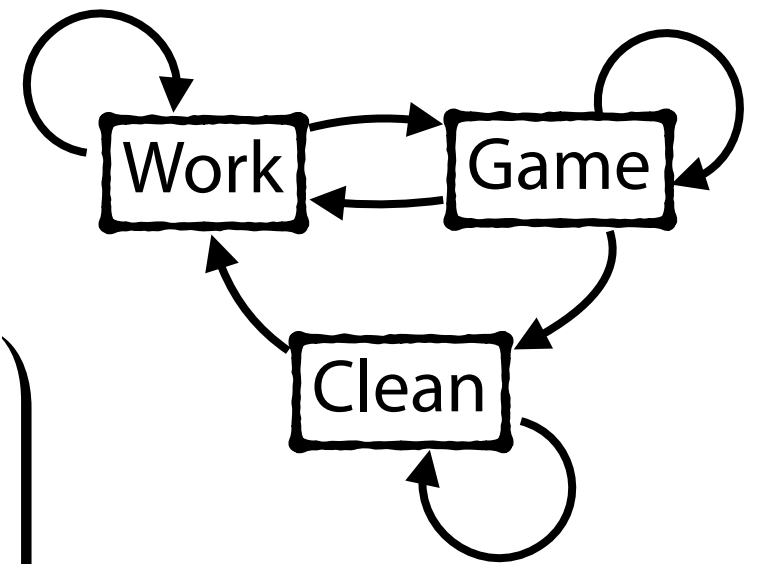
Markov Chain Matrix

What happens as $t \rightarrow \infty$?

$$M = \begin{pmatrix} .4 & .6 & 0 \\ .1 & .6 & .3 \\ .5 & 0 & .5 \end{pmatrix} \quad M^3 = \begin{pmatrix} .238 & .492 & .270 \\ .307 & .402 & .291 \\ .335 & .450 & .215 \end{pmatrix}$$

$$M^{10} = \begin{pmatrix} .2940 & .4413 & .2648 \\ .2942 & .4411 & .2648 \\ .2942 & .4413 & .2648 \end{pmatrix}$$

$$M^{60} = \begin{pmatrix} .2941 & .4412 & .2647 \\ .2941 & .4412 & .2647 \\ .2941 & .4412 & .2647 \end{pmatrix}$$



Markov Chain Stationary Distribution

A probability vector π is called a **stationary distribution** for a Markov Chain if it satisfies the stationary equation: $\pi = \pi M$

$$M = \begin{pmatrix} .4 & .6 & 0 \\ .1 & .6 & .3 \\ .5 & 0 & .5 \end{pmatrix}$$

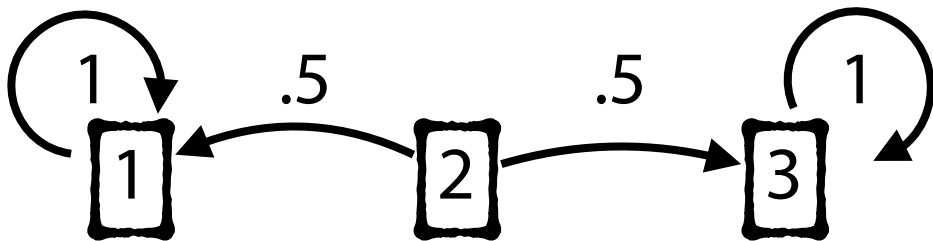
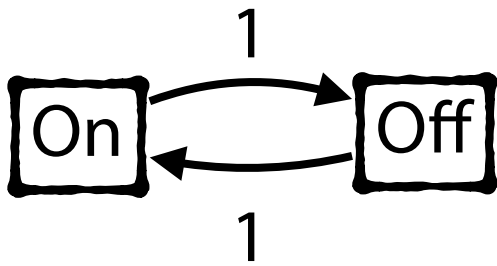
$$\pi[W] = .4\pi[W] + .1\pi[G] + .5\pi[C]$$

$$\pi[S] = .6\pi[W] + .6\pi[G] + 0\pi[C]$$

$$\pi[E] = 0\pi[W] + .3\pi[G] + .5\pi[C]$$

Markov Chain Stationary Distribution

Stationary distributions can be calculated using the system of equation (and that all probabilities sum to 1). **But not every Markov Chain has a steady state (and some have infinitely many)!**



Markov Chain Monte Carlo



There are ways to prove whether a Markov Chain has a stationary distribution, but several algorithms exist that approximate!

Gibbs Sampling:

Randomly assign values to a probability vector $\pi = (\theta_1, \theta_2, \dots, \theta_d)$.

Repeatedly:

Pick a random $1 \leq i \leq d$

Randomly update value $\theta_i \mid \theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_d$

Markov Chain Monte Carlo



There are ways to prove whether a Markov Chain has a stationary distribution, but several algorithms exist that approximate!

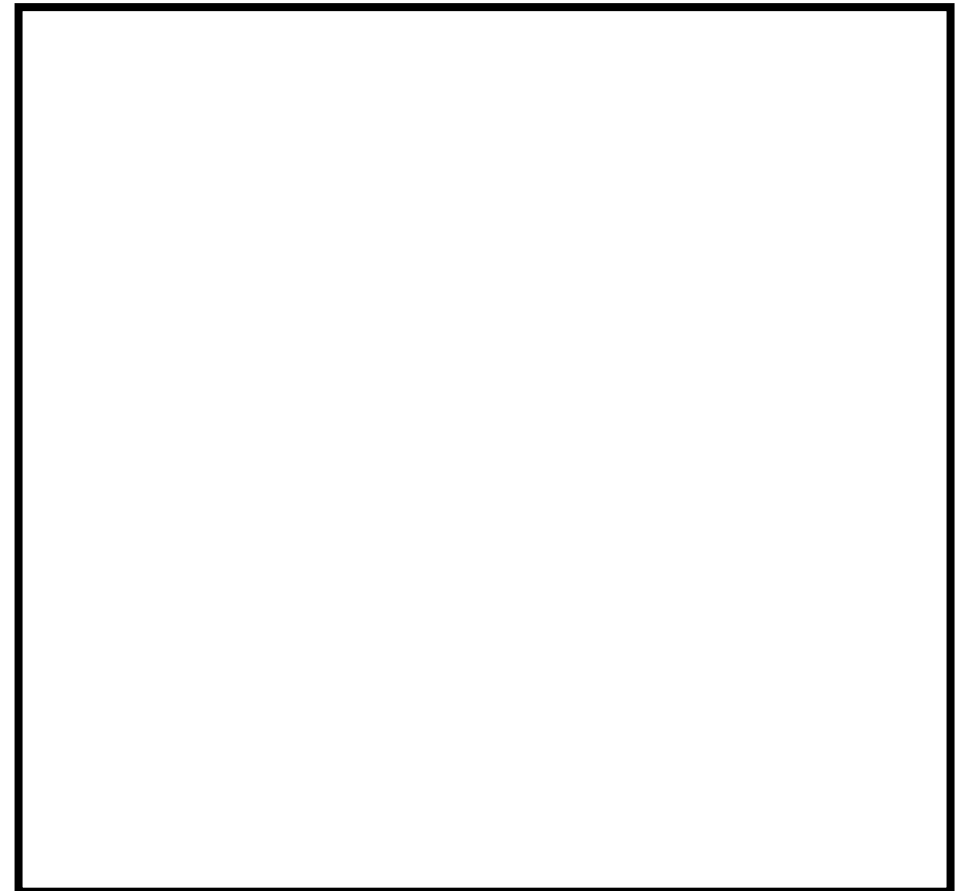
Gibbs Sampling:

Randomly assign values to a probability vector $\pi = (\theta_1, \theta_2, \dots, \theta_d)$.

Repeatedly:

Pick a random $1 \leq i \leq d$

Randomly update value θ_i based on $\theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_d$

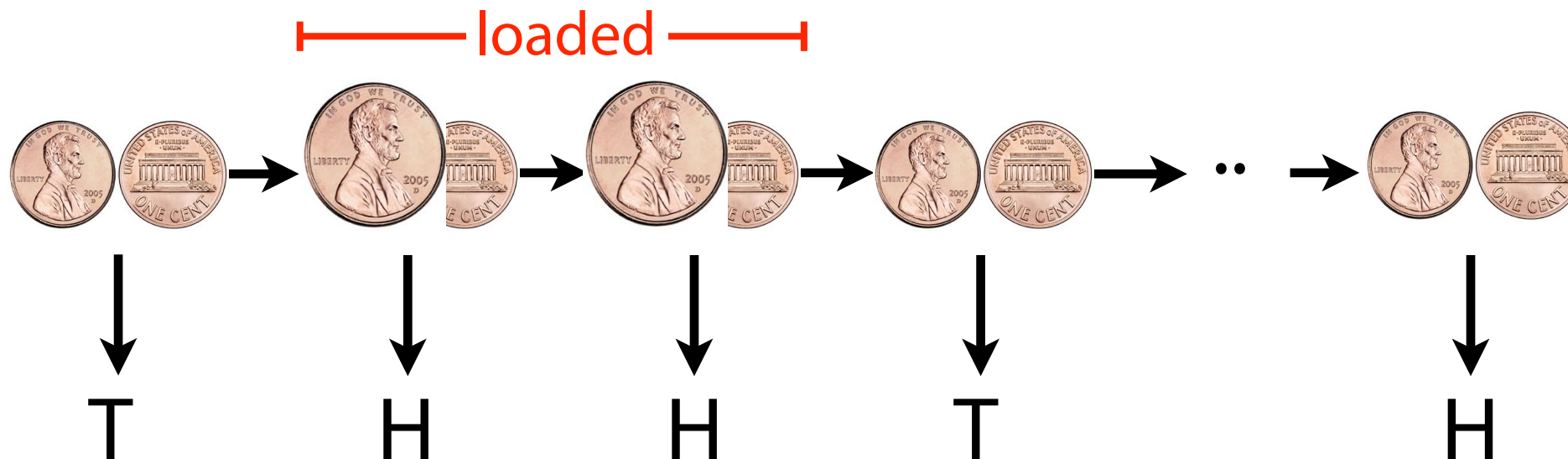


Hidden Markov Models

In the real world, we often don't know the underlying markov chain!

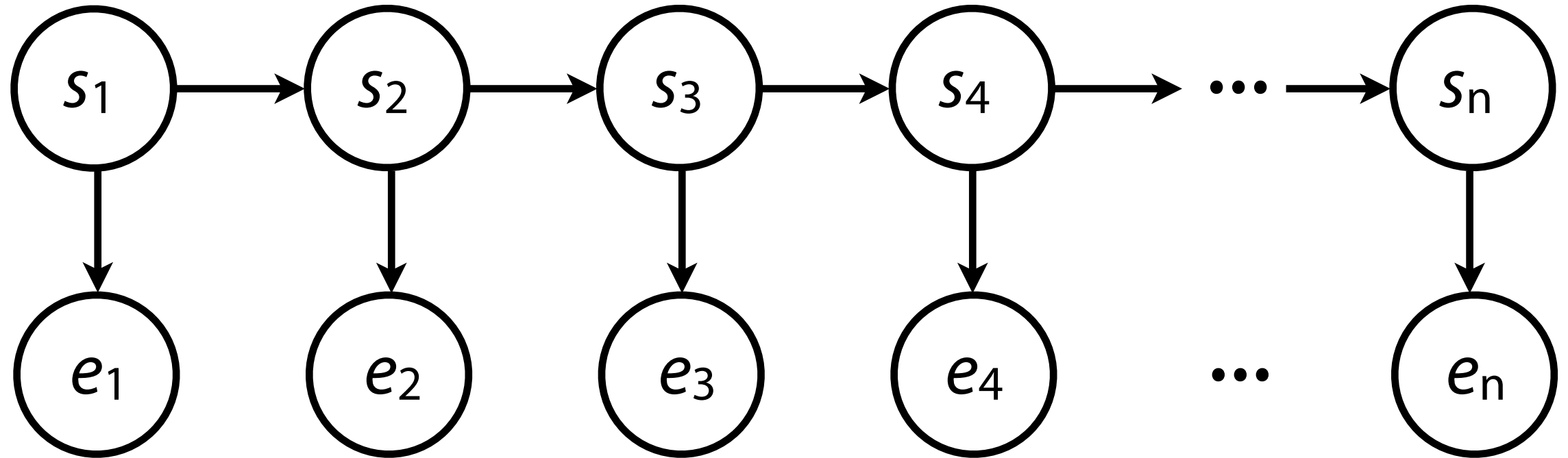
Instead, we have observations that can be used to predict our current state.

Ex: Repeated coin flips but *sometimes* I cheat and use a fixed coin.



Hidden Markov Model

Unobserved States



Observed Emissions