

Markov Chains

Ben Langmead



JOHNS HOPKINS

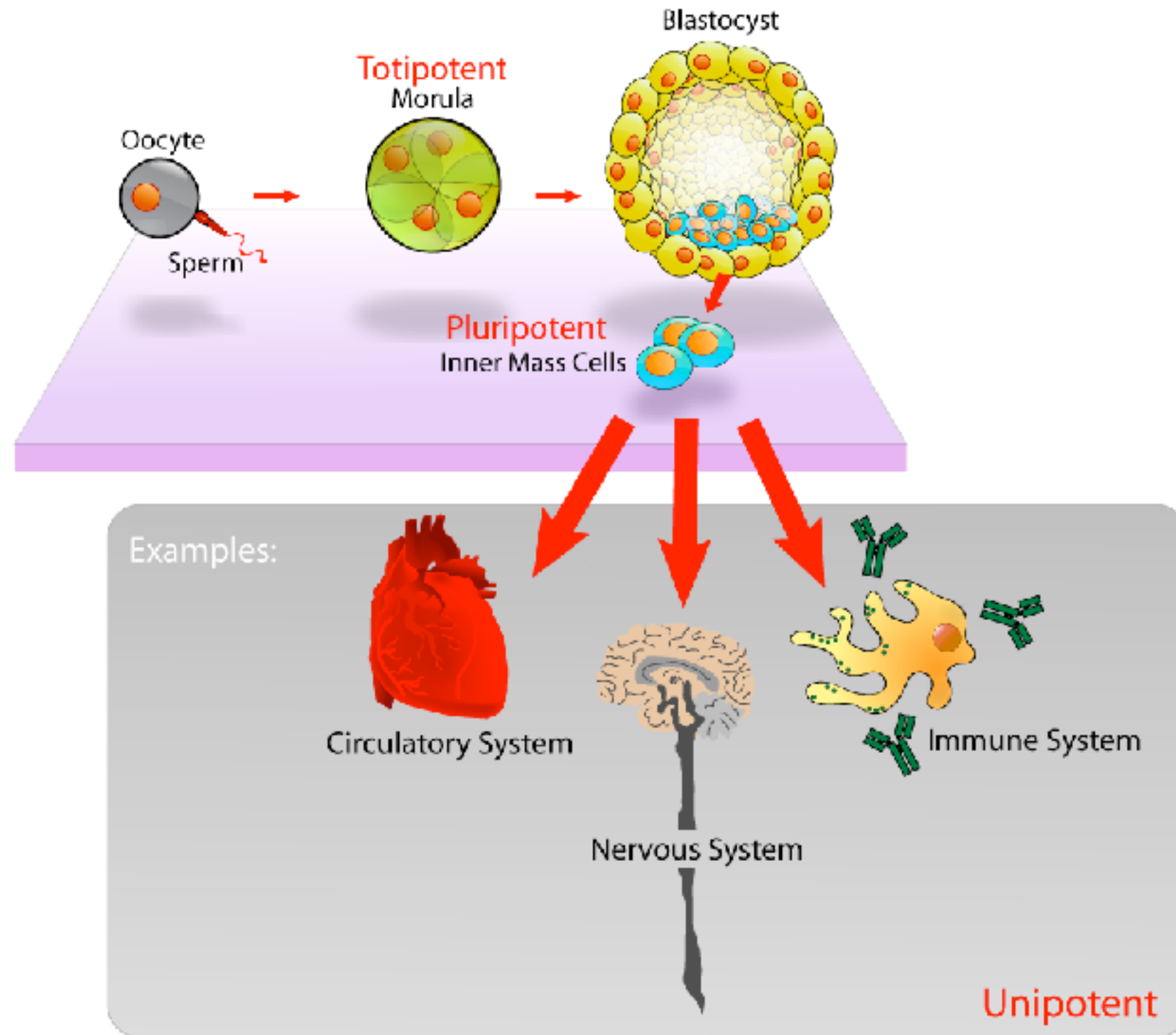
WHITING SCHOOL
of ENGINEERING

Department of Computer Science



Please sign the guestbook on my teaching materials page, or email me (ben.langmead@gmail.com) to tell me briefly how you are using the slides. For original Keynote files, email me.

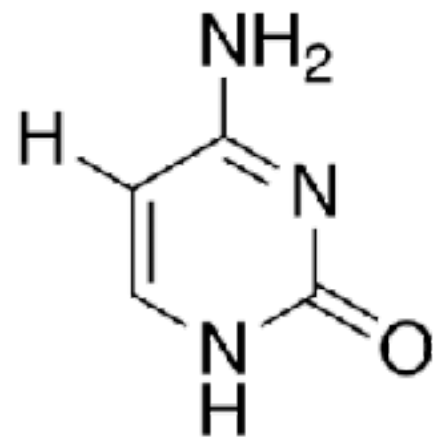
Epigenetics



http://en.wikipedia.org/wiki/File:Stem_cells_diagram.png

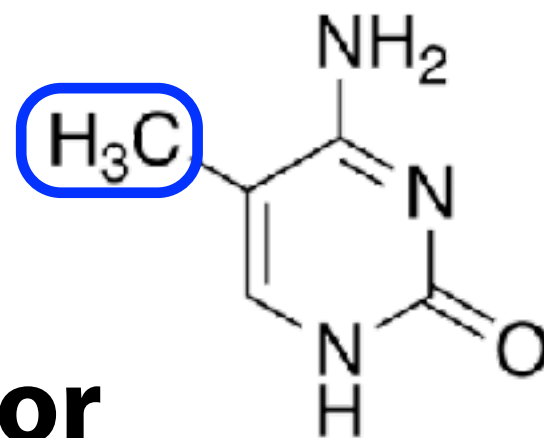
CpG Islands

Dinucleotide "CG" (AKA "CpG") is special because C can have a *methyl group* attached



Unmethylated

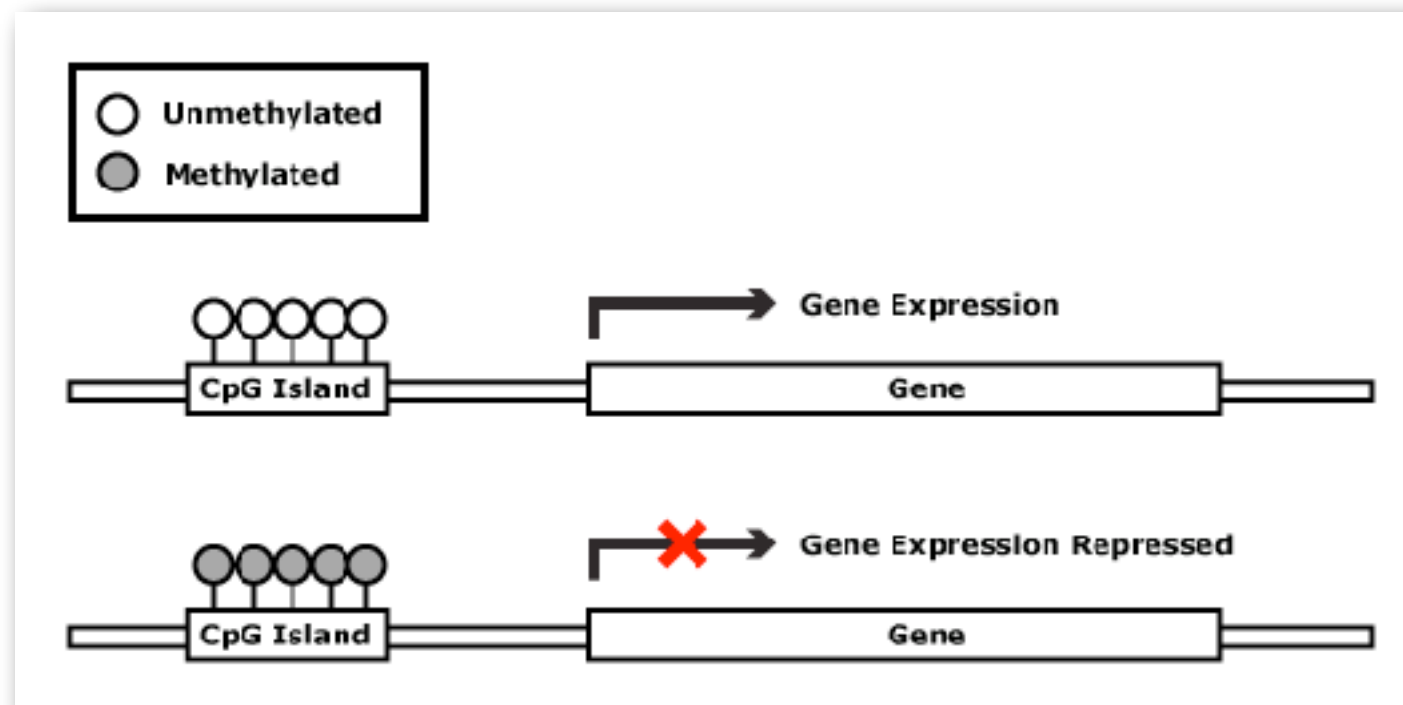
or



Methylated

CpG Islands

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



http://missinglink.ucsf.edu/lm/genes_and_genomes/methylation.html

Methylated CpG islands can suppress gene expression

CpG Islands

Goal: strategy for scoring a k -mer according to how confident we are it belongs to a CpG island

Scores should be *probabilities*

Probability review

Sample space (Ω) is set of all possible outcomes

E.g. $\Omega = \{ \text{all possible rolls of 2 dice} \}$

An *event* (A, B, C, \dots) is a subset of Ω

$A = \{ \text{rolls where first die is odd} \}$, $B = \{ \text{rolls where second die is even} \}$

We're often concerned with assigning a probability to an event

$P(A)$: fraction of all possible outcomes that are in A

$$P(A) = |A| / |\Omega| = 18 / 36 = 0.5$$

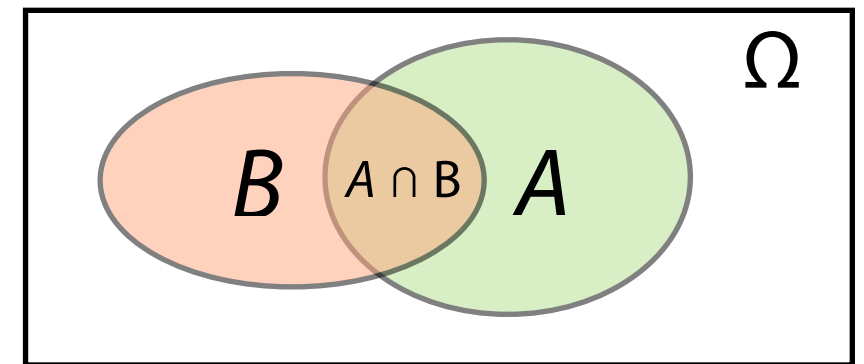
Probability review

$P(A, B)$: fraction of all possible outcomes that are in both A and B

$$P(A, B) = |A \cap B| / |\Omega| = 9 / 36 = 0.25$$

Sometimes written $P(A \cap B)$ or $P(AB)$

Joint probability of A and B



$P(A | B)$: fraction of outcomes in B that are also in A

conditional probability of A given B

$$P(A | B) = |A \cap B| / |B| = 9 / 18 = 0.5$$

$$P(A | B) = P(A, B) / P(B) \quad \text{Bayes rule}$$

$$P(A, B) = P(A | B) \cdot P(B) \quad \text{multiplication rule}$$

Probability review

Multiplication rule for joint prob with many variables:

$$P(A, B, C, D) = P(A \mid B, C, D) \cdot P(B, C, D)$$

$$P(A, B, C, D) = P(B \mid A, C, D) \cdot P(A, C, D)$$

$$P(A, B, C, D) = P(A \mid B, C, D) \cdot P(B \mid C, D) \cdot P(C, D)$$

$$P(A, B, C, D) = P(A \mid B, C, D) \cdot P(B \mid C, D) \cdot P(C \mid D) \cdot P(D)$$

*conditional
probabilities*

*marginal
probability*

Probability review

	A	Ω
B		

Events A and B are independent if $P(A \mid B) = P(A)$

So $P(A, B) = P(B) P(A \mid B) = P(A) P(B)$

Sequence models

Sequence model is a *probabilistic model* that associates probabilities with sequences

What *k*-mers do I see inside versus outside of a CpG island?

What's the probability of next character being A if previous characters were GATTAC?

Given a genome, where are the genes?

Right: model for eukaryotic gene finding

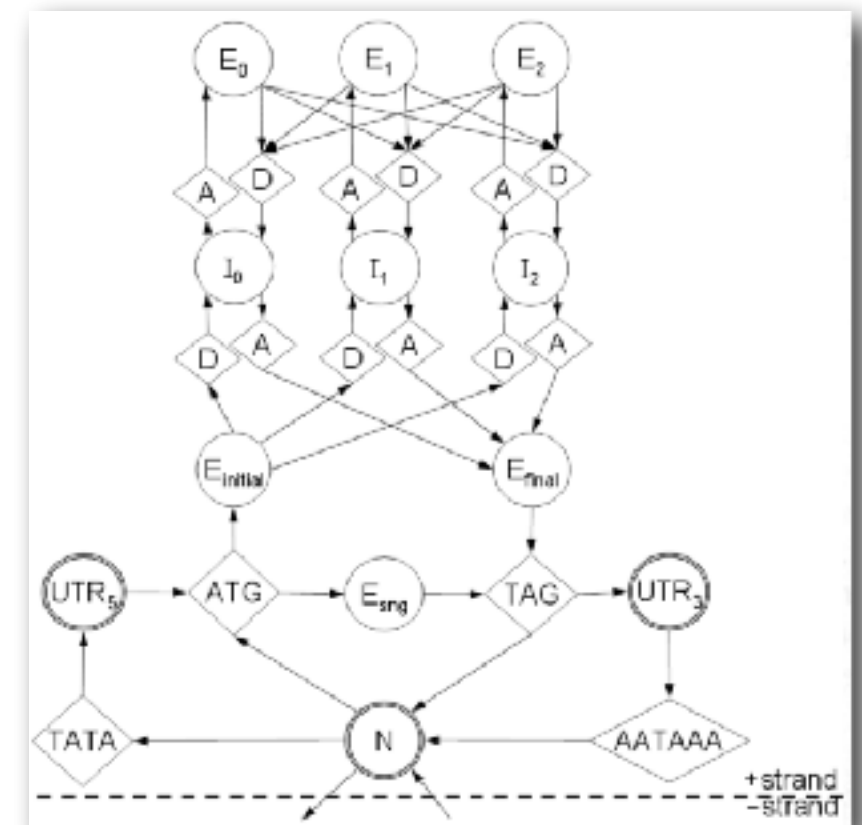


Image by Bill Majoros:
<http://www.genezilla.org/design.html>

Sequence models

Sequence models learn from examples

Say we have sampled 500K 5-mers from inside CpG islands and 500K 5-mers from outside

Can we guess whether CGCGC came from a CpG island?

# CGCGC inside	1553
# CGCGC outside	45

$$p(\text{inside}) = 1553 / (1553 + 45) = 0.972$$

Python example: https://bit.ly/CG_MarkovChain

Sequence models

$P(x)$ = probability of sequence x

$$P(x) = P(\underbrace{x_k, x_{k-1}, \dots, x_1}_{\text{Joint probability of each base}})$$

Estimating $P(x)$: # occurrences *inside* \div # occurrences total

For large k , might see few or no occurrences of x . Joint probabilities for very rare events are hard to estimate well!

Sequence models

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

multiplication rule

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

multiplication rule

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

(etc)

Assumption: probability of item at position k depends only on item at previous position: x_{k-1}

Technically: x_k is *conditionally independent* of $x_1 \dots x_{k-2}$ given x_{k-1}

Informally: "the future is independent of the past given the present"

Sequence models

Assumption: probability of item at position k depends only on item at previous position: x_{k-1}

$$\begin{aligned} P(x) &= P(x_k, x_{k-1}, \dots x_1) \\ &= P(x_k \mid x_{k-1}, \dots x_1) P(x_{k-1}, \dots x_1) \\ &= P(x_k \mid x_{k-1}, \dots x_1) P(x_{k-1} \mid x_{k-2}, \dots x_1) P(x_{k-2}, \dots x_1) \\ &\quad \text{(etc)} \quad \underbrace{\hspace{1cm}}_{\text{drop}} \quad \underbrace{\hspace{1cm}}_{\text{drop}} \quad \text{(etc)} \\ &\approx P(x_k \mid x_{k-1}) P(x_{k-1} \mid x_{k-2}) \dots P(x_2 \mid x_1) P(x_1) \end{aligned}$$

Markov property / assumption

Big assumption, but often reasonable and opens the door to tractable, powerful algorithms

Markov assumption

“To predict next state of the Parcheesi game, just tell me the current state. I don’t care about any other previous states.”

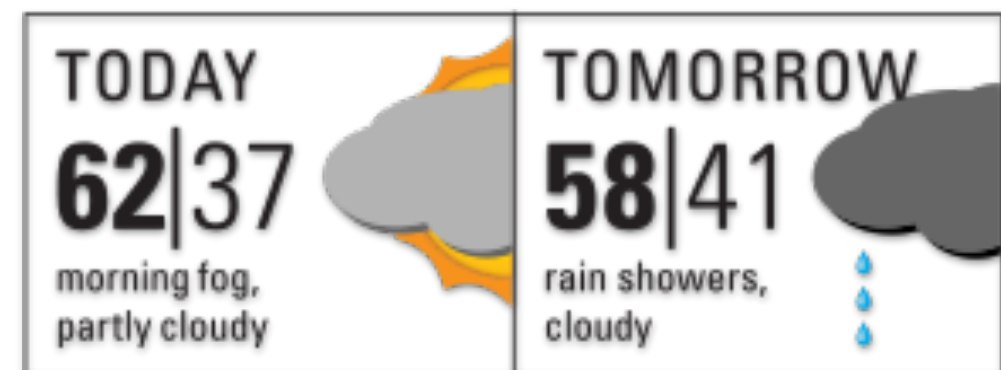
Reasonable assumption;
basically true

“To predict today’s weather, just tell me yesterday’s weather. I don’t care about any other previous days’ weather.”

It helps more to know more than just previous day’s weather. Still, fairly reasonable assumption.



<https://commons.wikimedia.org/wiki/File:Parcheesi-board.jpg>



https://en.wikipedia.org/wiki/Weather_forecasting#/media/File:Newspaper_weather_forecast_-_today_and_tomorrow.svg

Markov chain

Assigning a probability to a sequence using Markov property:

$$P(x) \underset{\substack{\text{Markov} \\ \text{property}}}{\approx} P(x_k | x_{k-1}) P(x_{k-1} | x_{k-2}) \dots P(x_2 | x_1) P(x_1)$$

Say x is a nucleotide k -mer

$P(x_i | x_{i-1})$ probability of seeing nucleotide x_i in i^{th} position given that previous nucleotide is x_{i-1}

Shorthand: $P(G | C) =$ probability of G given previous is C

Markov chain

Say someone gives us the sequences of several CpG islands. How do we estimate, say, $P(G | C)$?

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

where X is any base

Markov chain

Given CpG island sequences from human chromosome 1, count dinucleotide occurrences and estimate all 16 possible $P(x_i | x_{i-1})$:

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

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

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

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

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

(etc)


where X is any base

Markov chain

Given example CpG island substrings we can estimate all $P(\text{base} \mid \text{previous base})$:

```
>>> 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				
					$P(T \mid G)$



Rows sum to 1

http://bit.ly/CG_MarkovChain

Markov chain

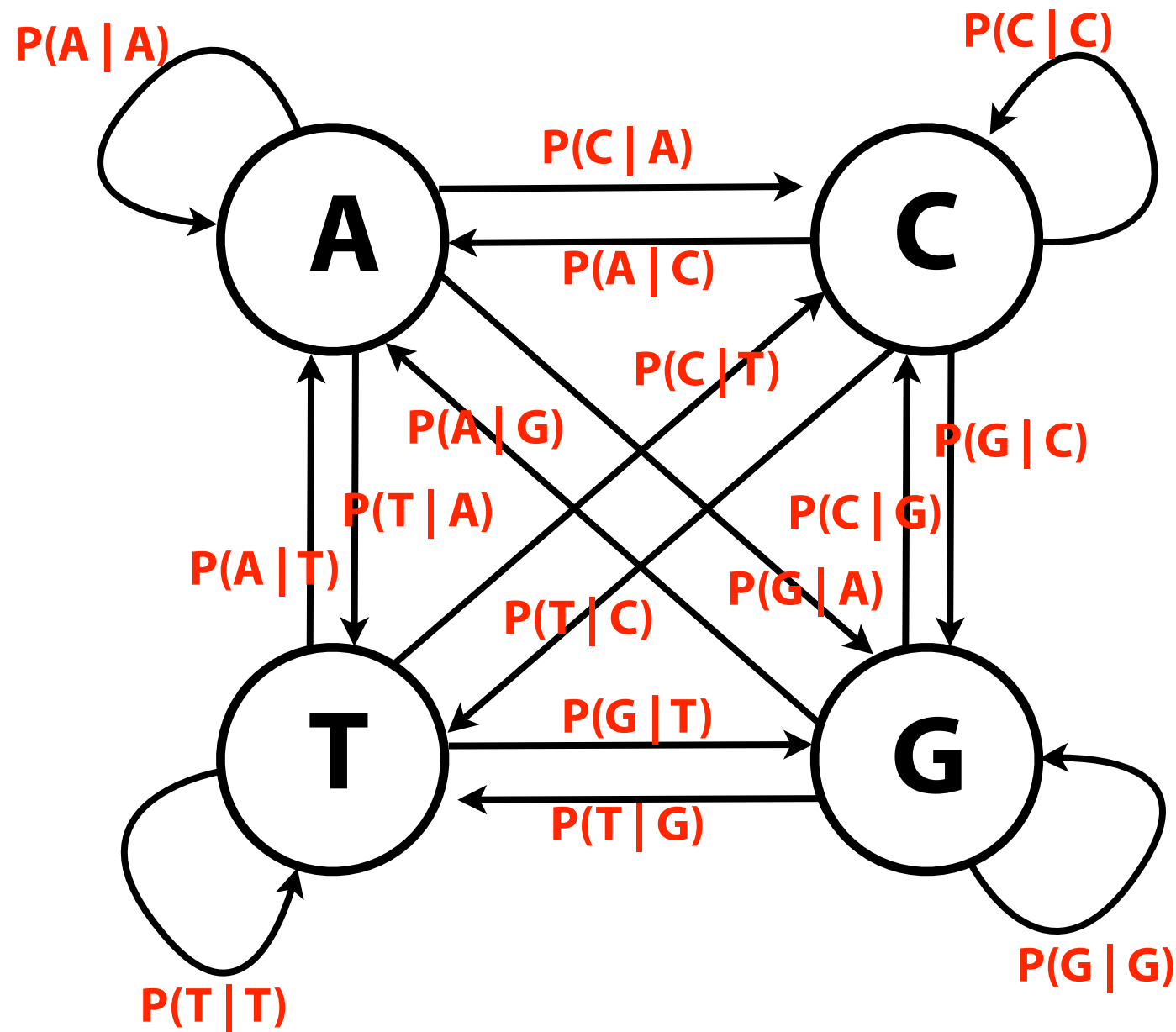
We can do the same for dinucleotides *outside* of CpG islands

```
>>> ins_conds, _ = markov_chain_from_dinucs(samp_in)
>>> print(ins_conds)
Inside
├ 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]]
>>> out_conds, _ = markov_chain_from_dinucs(samp_out)
>>> print(out_conds)
Outside
├ 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 ]])
      A           C           G           T
```

Notice anything interesting about the outside conditional probabilities?

$P(G | C)$ is low: outside CpG islands, G is rarely preceded by C

Markov chain



Markov chain is a probabilistic automaton

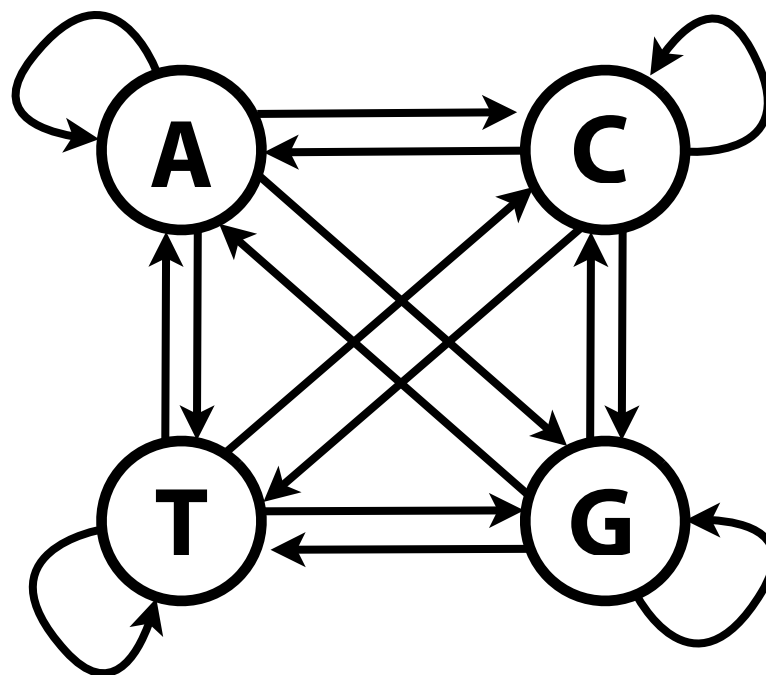
Edge has *transition probability*: probability that destination comes next after source

Markov chain

Recall how we assign a probability to a single string

$$P(x) \underset{\substack{\text{Markov} \\ \text{property}}}{\approx} P(x_k | x_{k-1}) P(x_{k-1} | x_{k-2}) \dots P(x_2 | x_1) P(x_1)$$

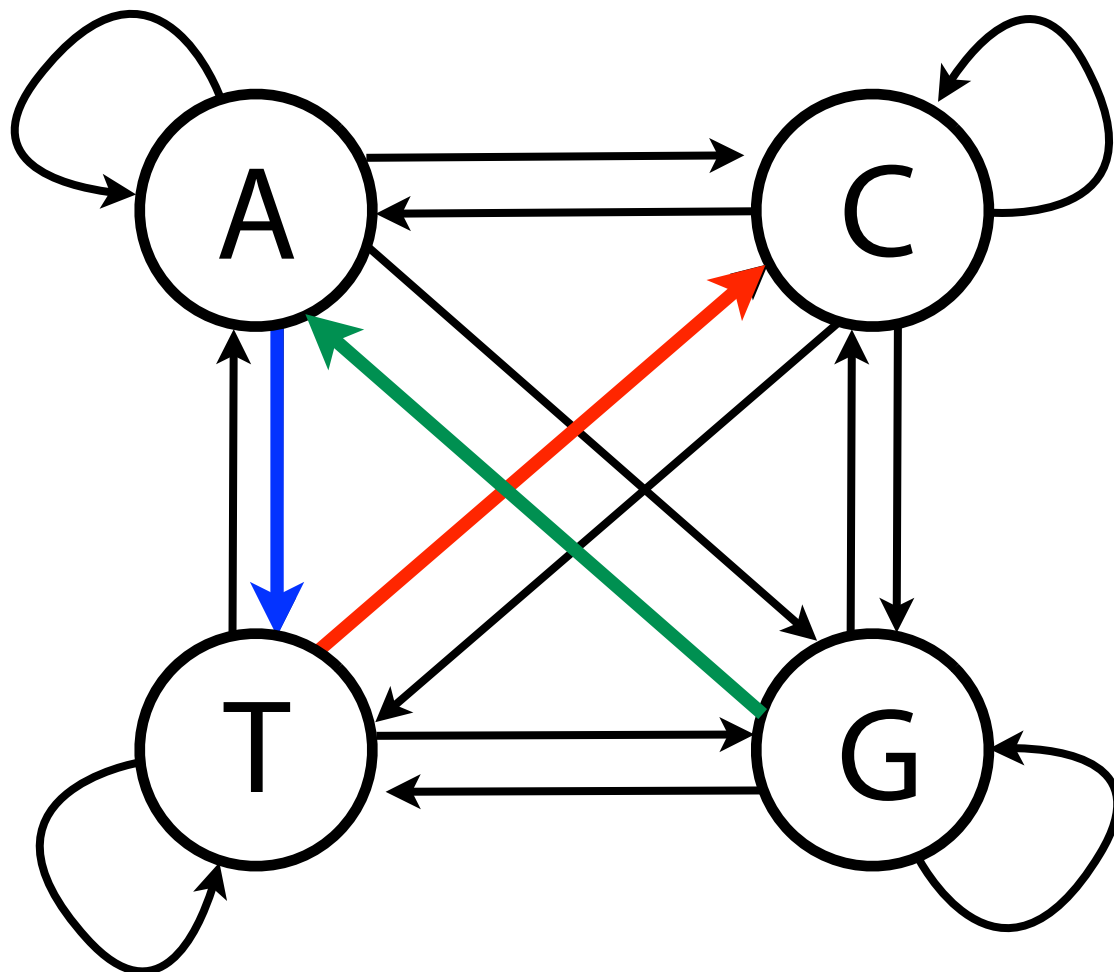
$P(x)$ equals product of Markov chain edge weights on our string-driven walk through the chain (...times $P(x_1)$)



Markov chain

```
>>> ins_conds, _ = markov_chain_from_dinucs(samp_in)
>>> 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	



$x = \text{GATC}$

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

$$P(x) = P(C | T) P(T | A) P(A | G) P(G)$$

$$= 0.33836493 *$$

$$0.1359636 *$$

$$0.17322219 *$$

$$0.25$$

$$= 0.001992$$

Markov chain

To avoid underflow, switch to log domain

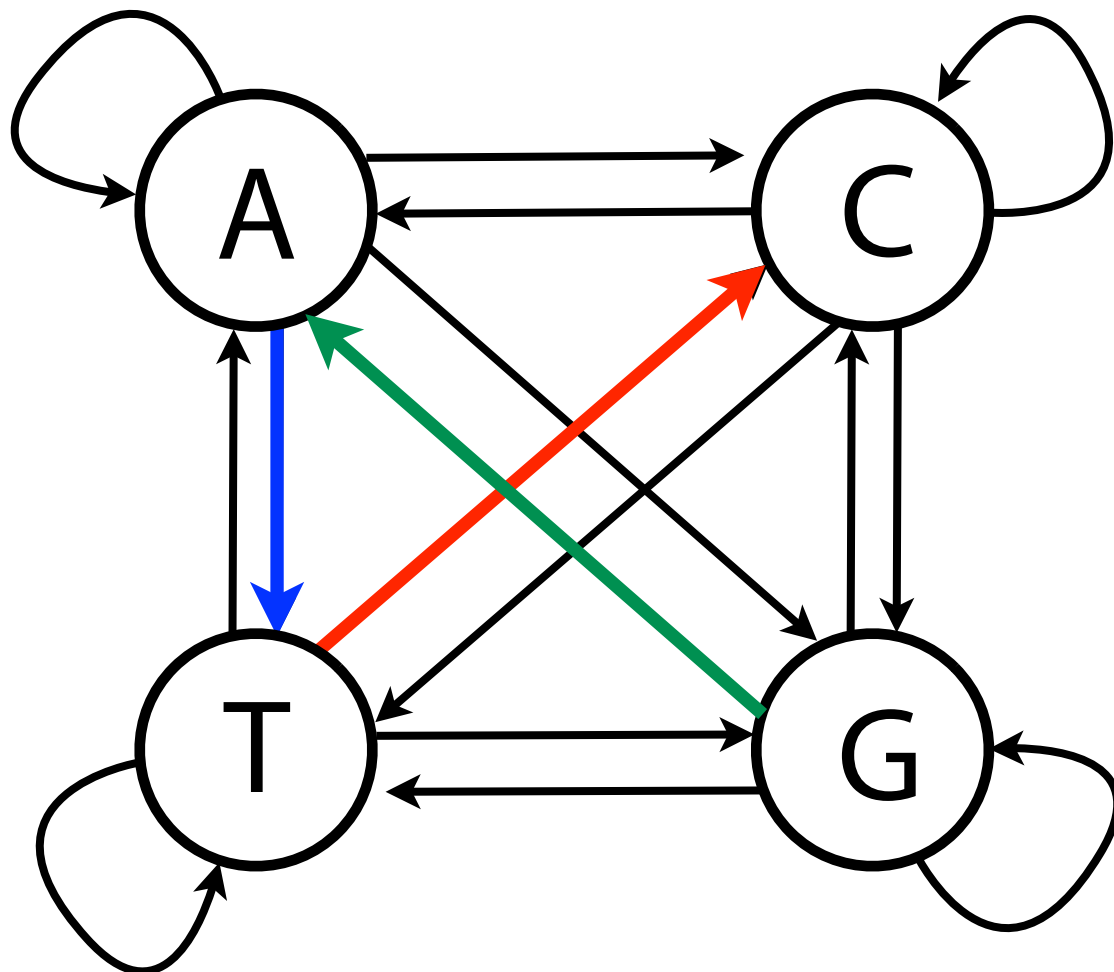
$$\begin{aligned}\log P(x) &\approx \log [P(x_k | x_{k-1}) P(x_{k-1} | x_{k-2}) \dots P(x_2 | x_1) P(x_1)] \\ &= \log P(x_k | x_{k-1}) + \log P(x_{k-1} | x_{k-2}) + \dots \\ &\quad \dots \text{product becomes sum!} \\ &= \sum_{i=2}^k \log P(x_i | x_{i-1}) + \log P(x_1)\end{aligned}$$

Assume logs are base 2

Markov chain

```
>>> ins_conds, _ = markov_chain_from_dinucs(samp_in)
>>> print(numpy.log2(ins_conds))
```

X_{i-1}	A	C	G	T
A	-2.44009488	-1.8820643	-1.30195688	-2.84832282
C	-2.38974049	-1.469396	-2.00590131	-2.32864974
G	-2.51948223	-1.60979755	-1.48694353	-2.82436637
T	-3.41910668	-1.52509737	-1.43889385	-2.39435058
	A	C	G	T
	X_i			



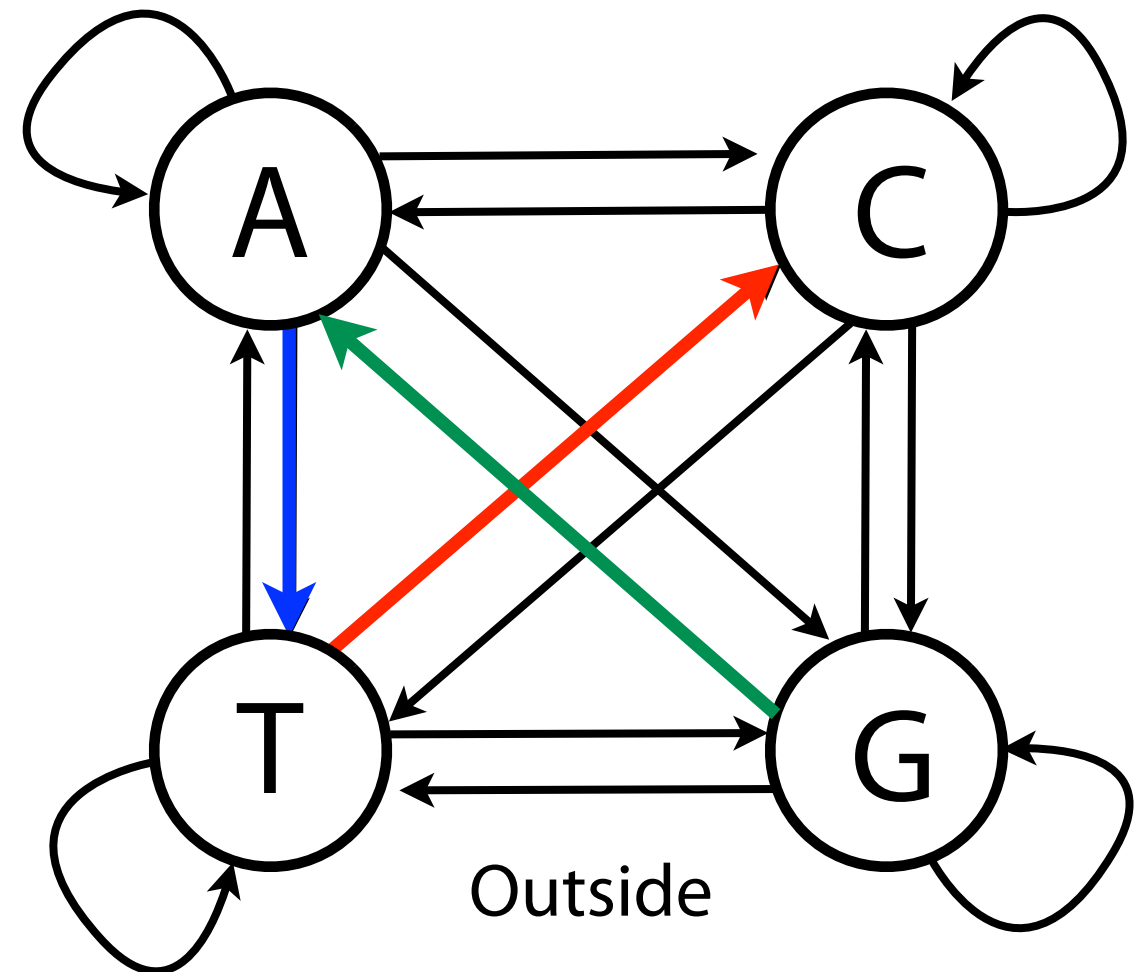
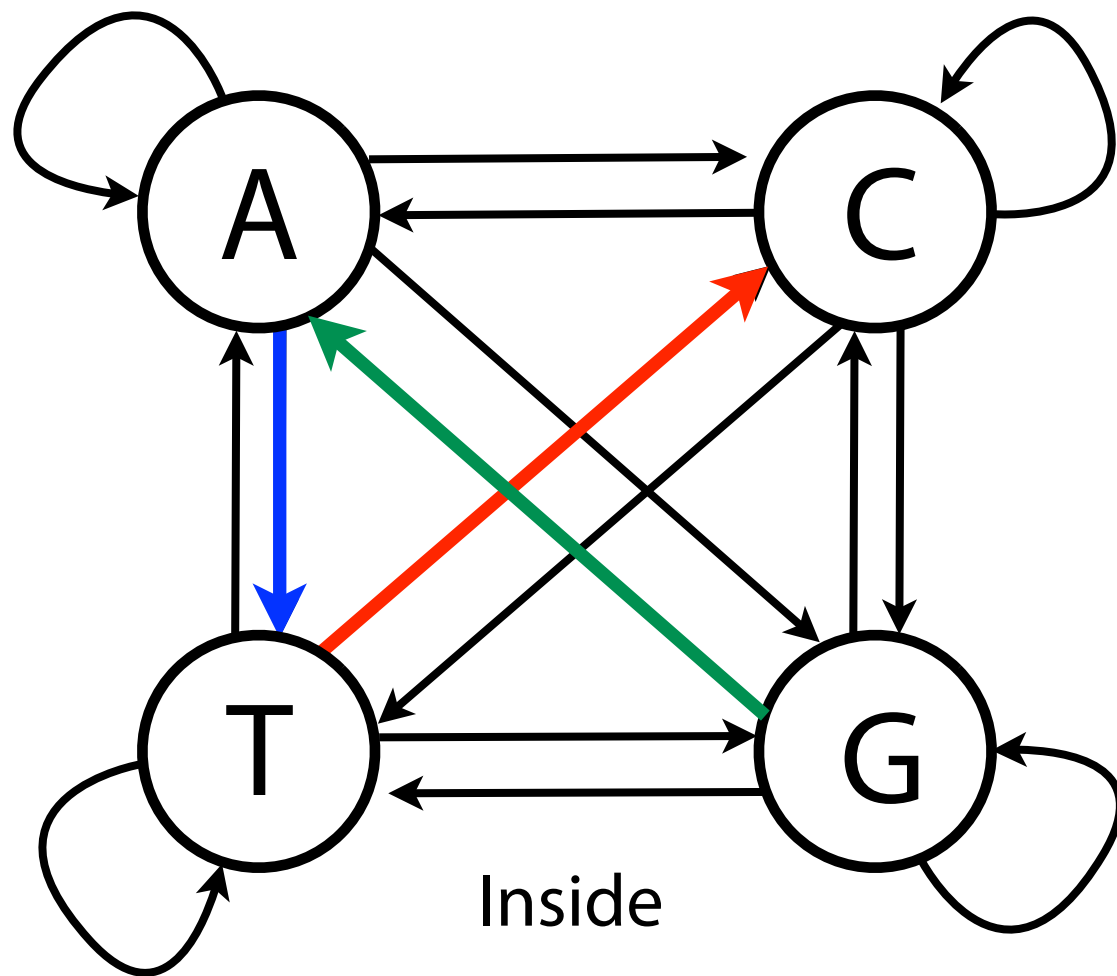
$x = \text{GATC}$

$$\begin{aligned}\log P(x) &= \sum_{i=2}^4 \log P(x_i | x_{i-1}) + \log P(x_1) \\ &= -1.52509737 + \\ &\quad -2.84832282 + \\ &\quad -2.51948223 + \\ &\quad -2.0 \\ &= -8.89290\end{aligned}$$

Markov chain

$P(x)$ given the inside-CpG model is helpful, but we really want to know which model is better, inside CpG or outside CpG?

Use *ratio*: $\frac{P(x) \text{ from inside model}}{P(x) \text{ from outside model}}$



Markov chain

Take log, get a *log ratio*:

$$S(x) = \log \frac{P(x) \text{ inside CpG}}{P(x) \text{ outside CpG}}$$

If inside more probable than outside, fraction is > 1 , log ratio is > 0 .
Otherwise, fraction is ≤ 1 and log ratio is ≤ 0 .

$$S(x) = \log \frac{P(x) \text{ inside CpG}}{P(x) \text{ outside CpG}}$$

$$\begin{aligned}
&= \log [P(x) \text{ inside CpG}] - \log [P(x) \text{ outside CpG}] && \text{(Marginal probabilities ignored here)} \\
&= \sum_{i=2}^k (\log [P(x_i | x_{i-1}) \text{ inside CpG}]) - \sum_{i=2}^k \log ([P(x_i | x_{i-1}) \text{ outside CpG}]) \\
&= \sum_{i=2}^k (\underbrace{ \log [P(x_i | x_{i-1}) \text{ inside CpG}] - \log [P(x_i | x_{i-1}) \text{ outside CpG}] }_{\text{log-odds}})
\end{aligned}$$

New table: elementwise log ratios between inside/outside

Markov chain

<i>Inside</i>	T	A	<pre>>>> ins_conds, _ = markov_chain_from_dinucs(samp_in)</pre>			
		C	<pre>>>> print(ins_conds)</pre>			
		G	[[0.19152248, 0.27252589, 0.39998803, 0.1359636],			
	I	T	[0.18921984, 0.35832388, 0.25467081, 0.19778547],			
			[0.17322219, 0.33142737, 0.35571338, 0.13963706],			
			[0.09509721, 0.33836493, 0.37567927, 0.19085859]]			
<i>Outside</i>	T	A	<pre>>>> out_conds, _ = markov_chain_from_dinucs(samp_out)</pre>			
		C	<pre>>>> print(out_conds)</pre>			
		G	[[0.33804066, 0.17971034, 0.23104207, 0.25120694],			
	I	T	[0.37777025, 0.25612117, 0.03987225, 0.32623633],			
			[0.30257815, 0.20326794, 0.24910719, 0.24504672],			
			[0.21790184, 0.20942905, 0.2642385 , 0.3084306]]			
<i>Log ratio</i>	T	A	<pre>>>> print(np.log2(ins_conds) - np.log2(out_conds))</pre>			
		C	[[-0.87536356, 0.59419041, 0.81181564, -0.85527103],			
		G	[-0.98532149, 0.49570561, 2.64256972, -0.7126391],			
	I	T	[-0.79486196, 0.68874785, 0.51821792, -0.79549511],			
			[-1.22085697, 0.73036913, 0.48119354, -0.69736839]]			
			A	C	G	T

Markov chain

Now, given a string x , we can easily assign it a log ratio “score” $S(x)$:

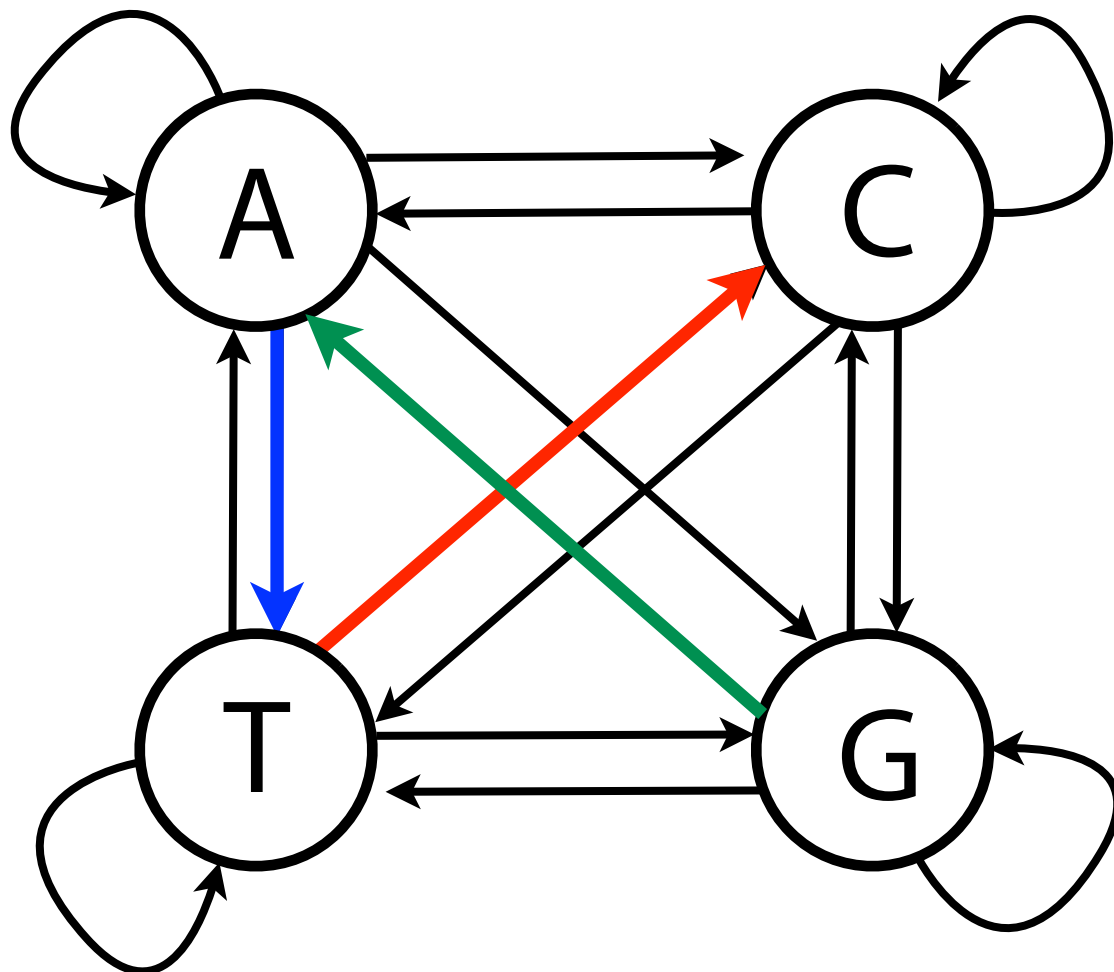
$$S(x) = \log \frac{P(x) \text{ inside CpG}}{P(x) \text{ outside CpG}}$$

$$\approx \sum_{i=2}^k \left(\log [P(x_i | x_{i-1}) \text{ inside CpG}] - \log [P(x_i | x_{i-1}) \text{ outside CpG}] \right)$$

Markov chain

```
>>> ins_conds, _ = markov_chain_from_dinucs(samp_in)
>>> out_conds, _ = markov_chain_from_dinucs(samp_out)
>>> print(np.log2(ins_conds) - np.log2(out_conds))
```

X_{i-1}	A	C	G	T
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
	X_i			



$x = \text{GATC}$

$$\begin{aligned} S(x) &= 0.73036913 + \\ &\quad -0.85527103 + \\ &\quad -0.79486196 \\ &= -0.919763 \end{aligned}$$

Negative, so probability with
outside model is greater

Markov chain

$$S(x) = \log \frac{P(x) \text{ inside CpG}}{P(x) \text{ outside CpG}}$$

$$S(\text{CGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG}) = 42.618$$

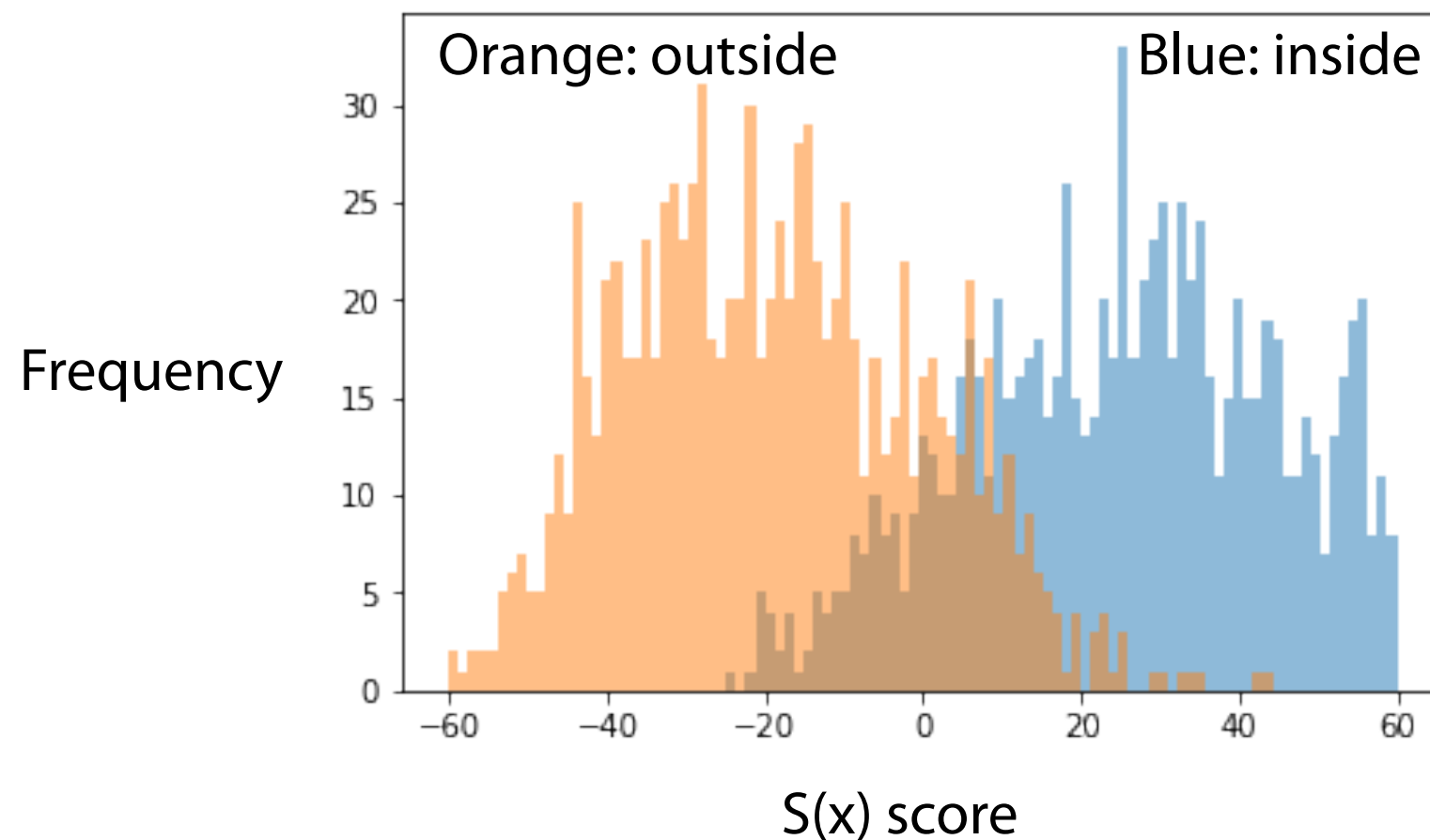
$$S(\text{ATTCTACTATCATCTATCTATCTTCT}) = -10.839$$

http://bit.ly/CG_MarkovChain

Markov chain

Drew 1,000 100-mers from inside CpG islands on chromosome 18, and another 1,000 from outside, and calculated $S(x)$ for all

Trained markov chain on dinucleotides from chromosome 22



http://bit.ly/CG_MarkovChain

Markov chain

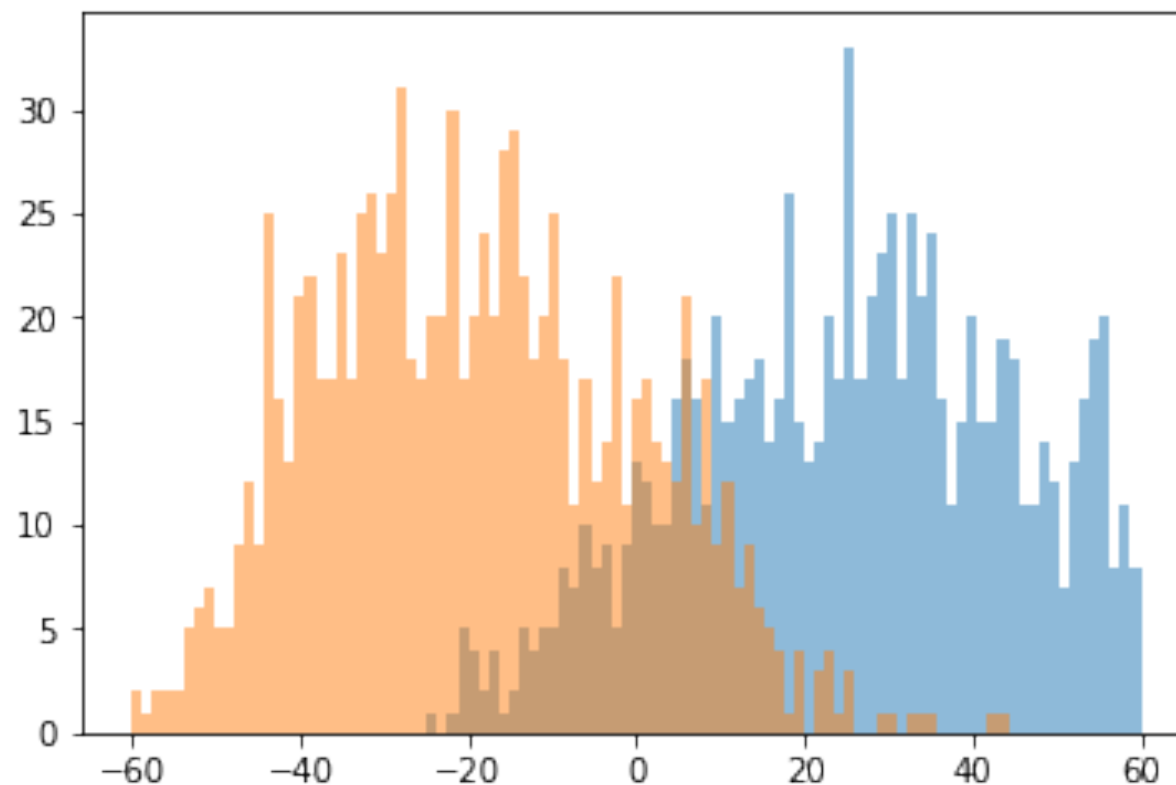
Markov property made our problem very tractable

$P(x_i | x_{i-1})$ s estimated in single, simple pass through training data

Transition probability tables have $|\Sigma|^2$ cells; fine for DNA & protein

Calculating $S(x)$ is $O(|x|)$; just lookups and additions

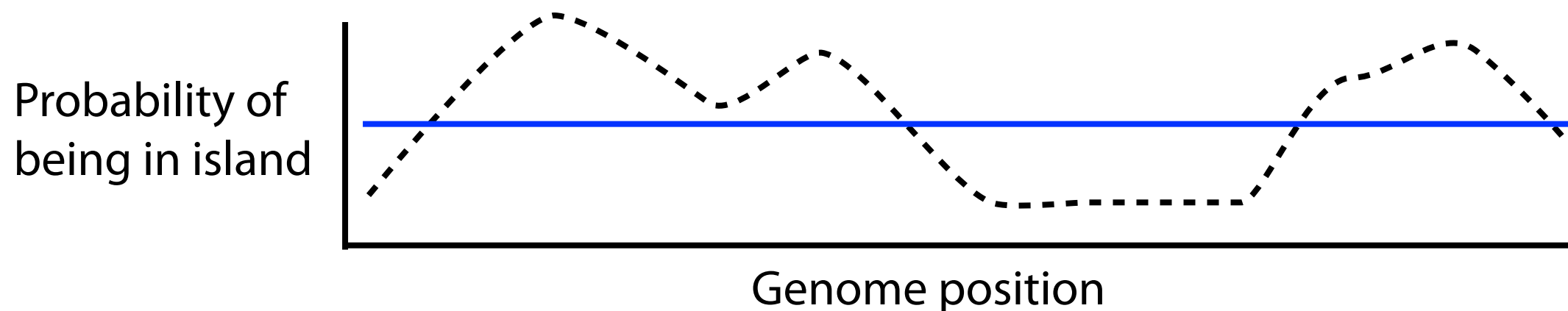
Discriminates fairly well between inside & outside:



Sequence models

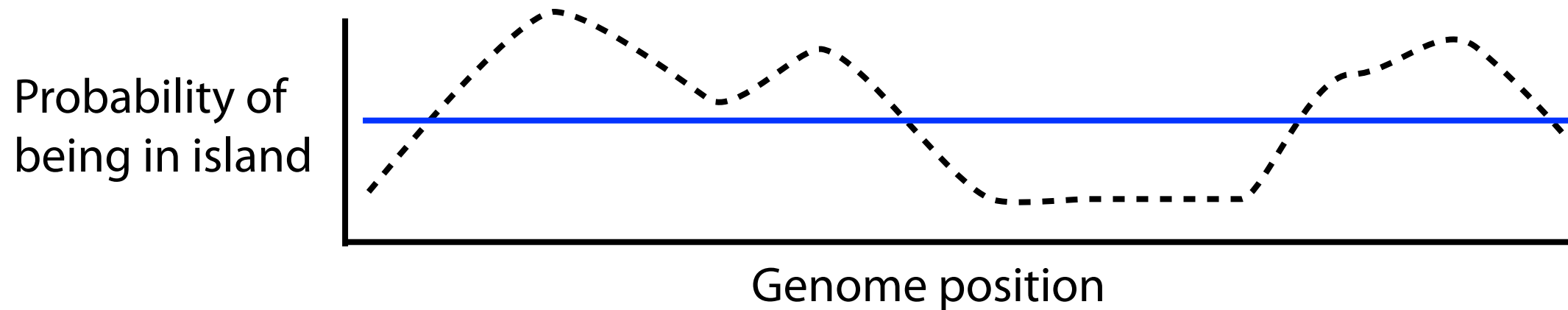
Can Markov chains find CpG islands in a “sea” of genome?

MC assigns a score to a string; doesn't naturally give a “running” score across a long sequence



But we can adapt it using a *sliding window*

Sequence models



Choice of k requires assumption about island lengths

If k is too large, we miss small islands

If k is too small, we see many small islands

We'd like a method that *switches between Markov chains* when entering or exiting a CpG island