

Markov chains, part 1

Ben Langmead



For original Keynote files, email me (ben.langmead@gmail.com)

Markov chain

Assigning a probability to a sequence using Markov assumption:

$$P(x) \approx \text{Markov assumption}$$

Say x is a length- k string of nucleotides (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}

Markov chain

Assigning a probability to a sequence using Markov assumption:

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

Say x is a length- k string of nucleotides (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 we are given sequences of several CpG islands. How do we estimate, say, $P(G | C)$?

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

where X is any base

Markov chain

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

$$P(\mathbf{A} | \mathbf{A}) =$$

Markov chain

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

$$\begin{aligned} P(\mathbf{A} | \mathbf{A}) &= \# \text{ times } \mathbf{AA} \text{ occurs} / \# \text{ times } \mathbf{AX} \text{ occurs} \\ P(\mathbf{C} | \mathbf{A}) &= \# \text{ times } \mathbf{AC} \text{ occurs} / \# \text{ times } \mathbf{AX} \text{ occurs} \\ P(\mathbf{G} | \mathbf{A}) &= \# \text{ times } \mathbf{AG} \text{ occurs} / \# \text{ times } \mathbf{AX} \text{ occurs} \\ P(\mathbf{T} | \mathbf{A}) &= \# \text{ times } \mathbf{AT} \text{ occurs} / \# \text{ times } \mathbf{AX} \text{ occurs} \\ P(\mathbf{A} | \mathbf{C}) &= \# \text{ times } \mathbf{CA} \text{ occurs} / \# \text{ times } \mathbf{CX} \text{ occurs} \\ & \text{(etc)} \end{aligned}$$

where X is any base

Markov chain

Given example CpG island substrings we can estimate all $P(\text{base} \mid \text{previous base})$, $P(X_i \mid X_{i-1})$:

X_{i-1} (previous)	A	$P(\text{A} \mid \text{A})$	$P(\text{C} \mid \text{A})$	$P(\text{G} \mid \text{A})$	$P(\text{T} \mid \text{A})$
	C	$P(\text{A} \mid \text{C})$	$P(\text{C} \mid \text{C})$	$P(\text{G} \mid \text{C})$	$P(\text{T} \mid \text{C})$
	G	$P(\text{A} \mid \text{G})$	$P(\text{C} \mid \text{G})$	$P(\text{G} \mid \text{G})$	$P(\text{T} \mid \text{G})$
	T	$P(\text{A} \mid \text{T})$	$P(\text{C} \mid \text{T})$	$P(\text{G} \mid \text{T})$	$P(\text{T} \mid \text{T})$
	A	C	G	T	
	X_i				

Markov chain

Given example CpG island substrings we can estimate all $P(\text{base} \mid \text{previous base})$, $P(X_i \mid X_{i-1})$:

(real data from human chr1)

```
>>> iTab, nTab = islandTransitionTables(fn, ifn)
>>> print(iTab)
```

X_{i-1} (previous)	A	C	G	T
A	0.18427153	0.27129525	0.4055757	0.13885752
C	0.19081672	0.36113346	0.24897947	0.19907035
G	0.17440554	0.32764433	0.35676759	0.14118254
T	0.09348595	0.3474561	0.36885	0.19020795

A **C** **G** **T**

X_i

P(T | G)

Rows sum to ... **1**

Markov chain

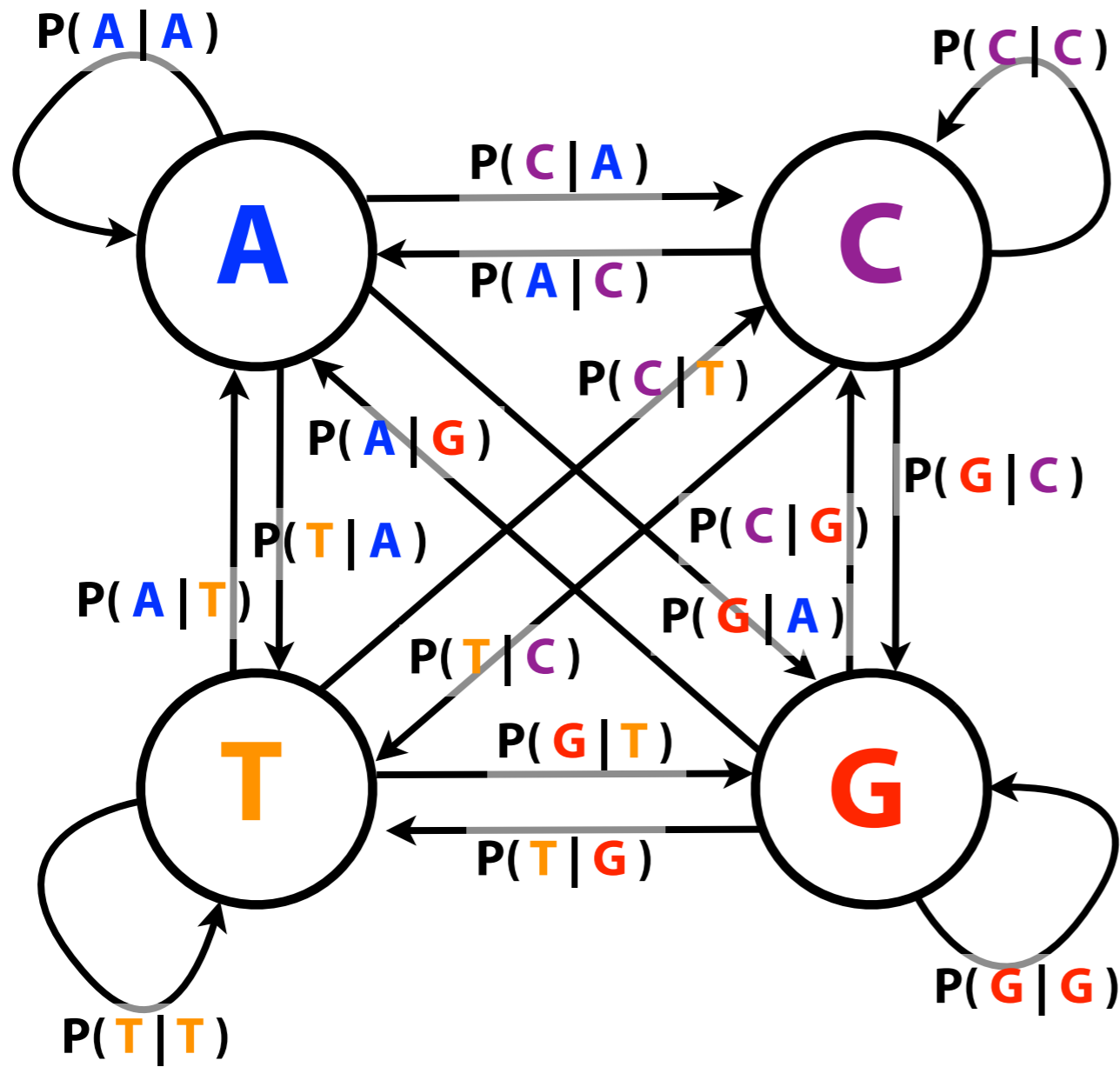
(real data from human chr1)

```
>>> iTab, oTab = islandTransitionTables(fn, ifn)
>>> print(iTab)
┌ Inside
├ A [[ 0.18427153, 0.27129525, 0.4055757 , 0.13885752],
├ C [ 0.19081672, 0.36113346, 0.24897947, 0.19907035],
├ G [ 0.17440554, 0.32764433, 0.35676759, 0.14118254],
└ T [ 0.09348595, 0.3474561 , 0.36885 , 0.19020795]]
>>> print(oTab)
┌ 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 about the **outside** conditional probabilities?

$P(\mathbf{G} \mid \mathbf{C})$ is low ; makes sense: outside CpG islands, C is rarely followed by G

Markov chain



Markov chain =
probabilistic automaton

Markov chain

```
>>> iTab, oTab = islandTransitionTables(fn, ifn)
>>> print(iTab)
```

X_{i-1}	A	[[0.18427153, 0.27129525, 0.4055757 , 0.13885752],		
C	[0.19081672, 0.36113346, 0.24897947, 0.19907035],			
G	[0.17440554, 0.32764433, 0.35676759, 0.14118254],			
T	[0.09348595, 0.3474561 , 0.36885 , 0.19020795]]			
	A	C	G	T
		X_i		

$x = \text{GATC}$

Markov chain

```
>>> iTab, oTab = islandTransitionTables(fn, ifn)
>>> print(iTab)
```

X_{i-1}	A	[[0.18427153, 0.27129525, 0.4055757 , 0.13885752],		
C	[0.19081672, 0.36113346, 0.24897947, 0.19907035],			
G	[0.17440554, 0.32764433, 0.35676759, 0.14118254],			
T	[0.09348595, 0.3474561 , 0.36885 , 0.19020795]]			
	A	C	G	T
	X_i			

$x = \text{GATC}$

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

$$P(x) = P(\text{C} | \text{T}) P(\text{T} | \text{A}) P(\text{A} | \text{G}) P(\text{G})$$

$$= 0.347 \times$$

$$0.139 \times$$

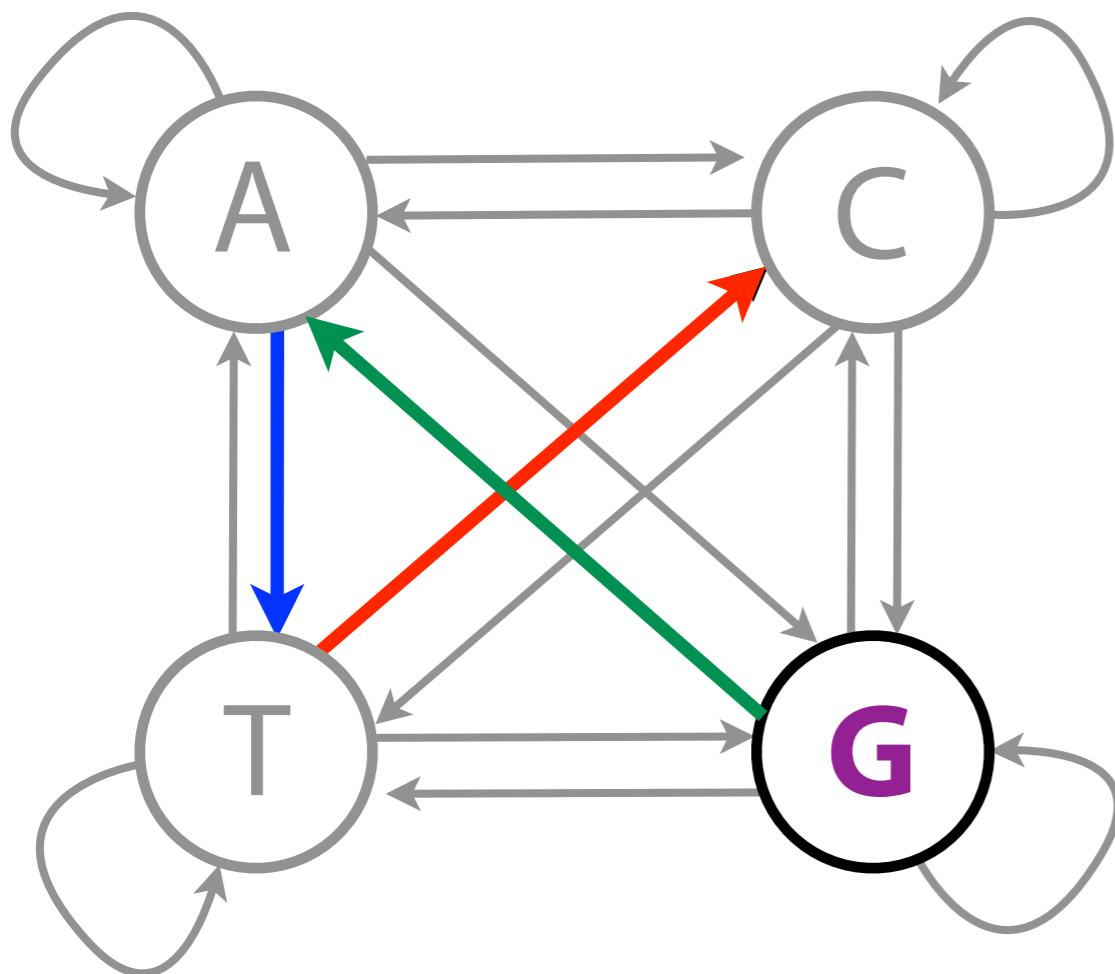
$$0.174 \times$$

$$0.25 \leftarrow (\text{Assuming A, C, G, T}$$

$$= 0.0021 \quad \text{equally likely})$$

Markov chain

```
>>> iTab, oTab = islandTransitionTables(fn, ifn)
>>> print(iTab)
A [[ 0.18427153, 0.27129525, 0.4055757 , 0.13885752],
C [[ 0.19081672, 0.36113346, 0.24897947, 0.19907035],
G [[ 0.17440554, 0.32764433, 0.35676759, 0.14118254],
T [[ 0.09348595, 0.3474561 , 0.36885 , 0.19020795]]
Xi-1      A      C      G      T
Xi
```



$x = \text{GATC}$

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

$$P(x) = P(\text{C} | \text{T}) P(\text{T} | \text{A}) P(\text{A} | \text{G}) P(\text{G})$$

$$= 0.347 *$$

$$0.139 *$$

$$0.174 *$$

$$0.25$$

$$= 0.0021$$