

Markov chains, part 2

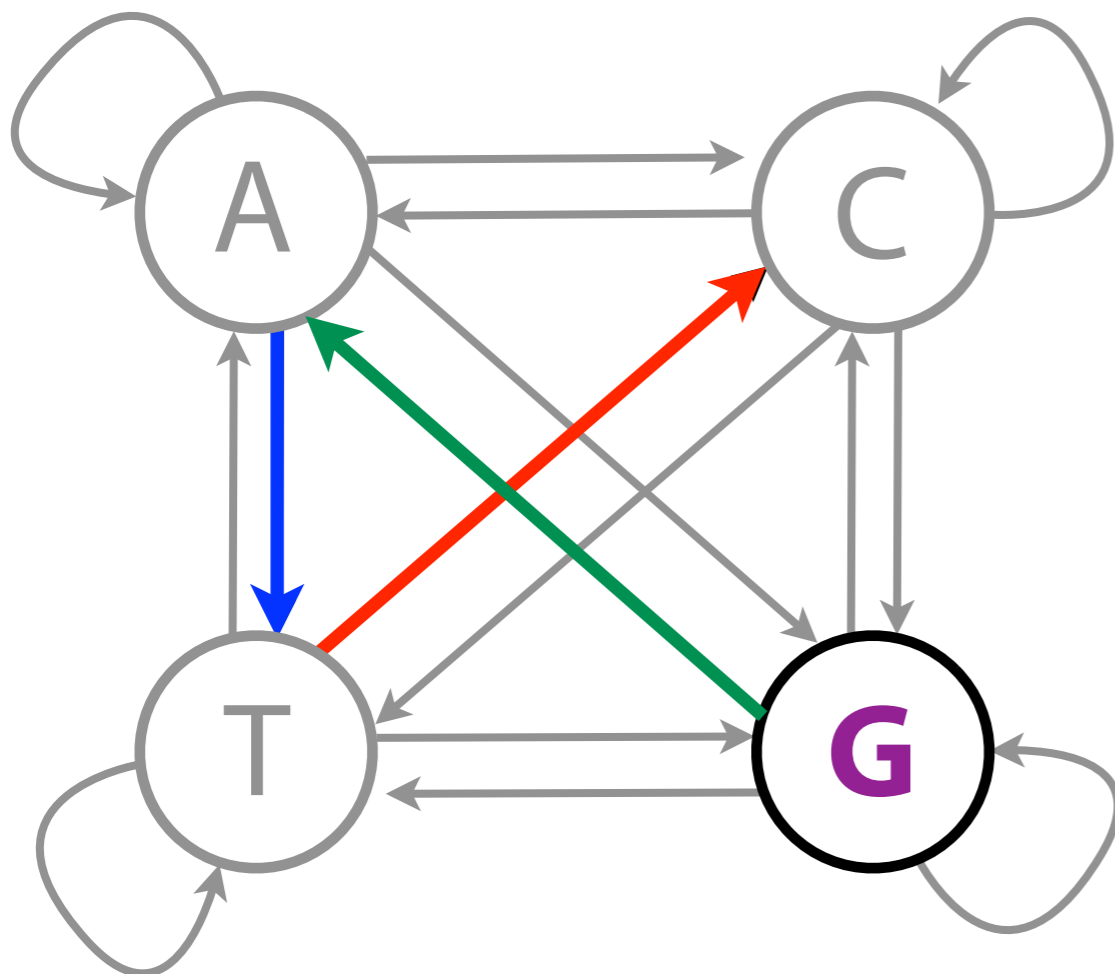
Ben Langmead



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

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
```



$x = \text{GATC}$

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

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

$$= 0.347 *$$

$$0.139 *$$

$$0.174 *$$

$$0.25$$

$$= 0.0021$$

Markov chain

** is exponentiation

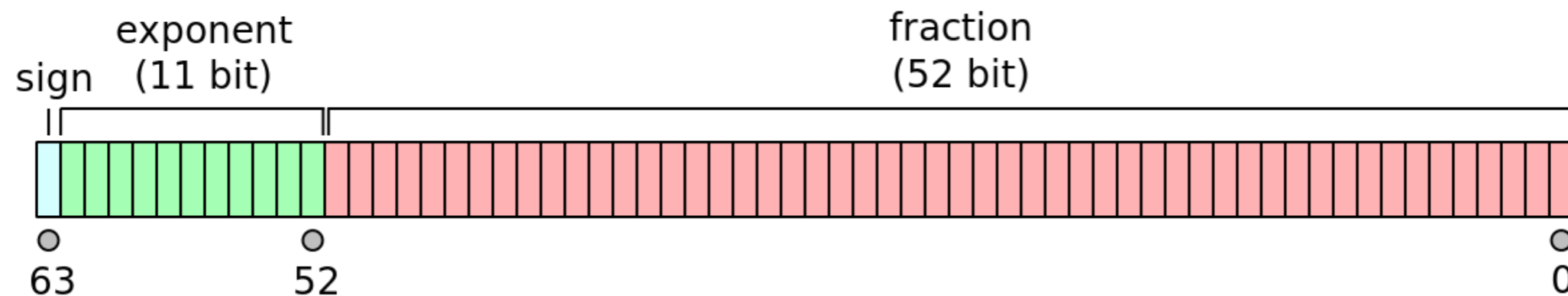
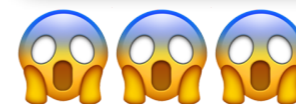
```
>>> 0.5 ** 1
0.5
>>> 0.5 ** 10
0.0009765625
>>> 0.5 ** 100
7.888609052210118e-31
>>> 0.5 ** 1000
9.332636185032189e-302
>>> 0.5 ** 1100
```

Markov chain

Underflow: when result is a number so small, it cannot be represented in a limited-precision floating point number and is **rounded to 0**

** is exponentiation

```
>>> 0.5 ** 1
0.5
>>> 0.5 ** 10
0.0009765625
>>> 0.5 ** 100
7.888609052210118e-31
>>> 0.5 ** 1000
9.332636185032189e-302
>>> 0.5 ** 1100
0.0
```



https://en.wikipedia.org/wiki/Double-precision_floating-point_format

Markov chain

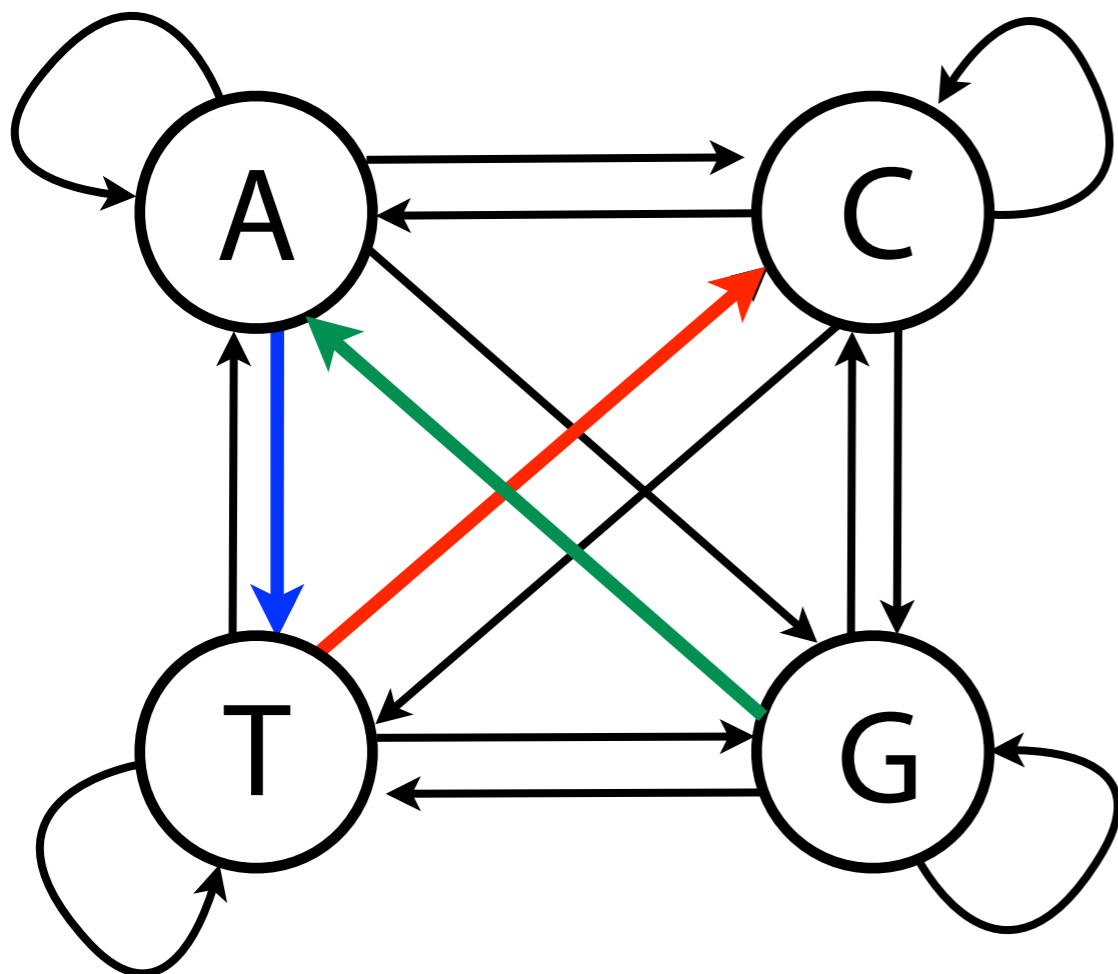
To avoid underflow, switch to log domain

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

Assume logs are base 2

Markov chain

```
>>> iTab, nTab = islandTransitionTables(fn, ifn)
>>> print(numpy.log2(iTab))
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]]
Xi-1
A
C
G
T
Xi
```



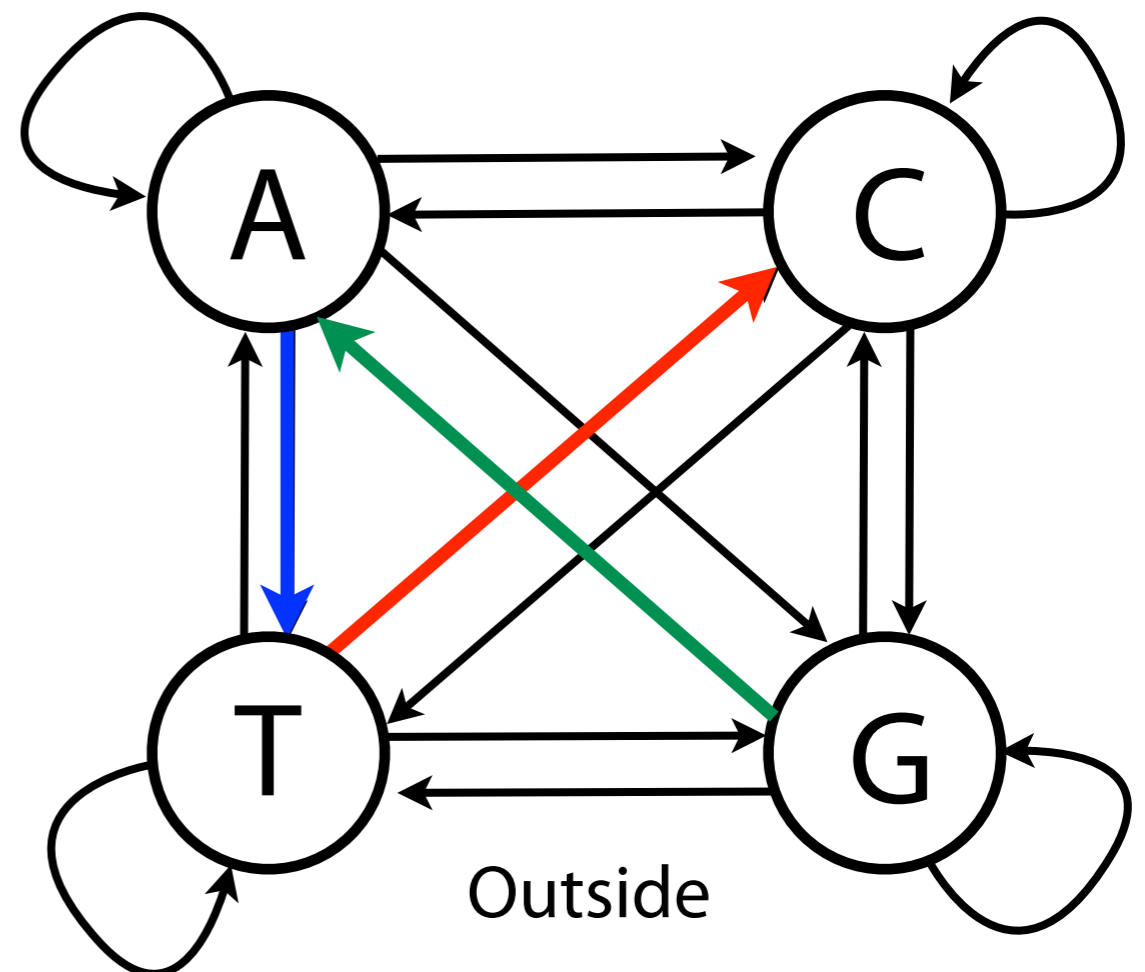
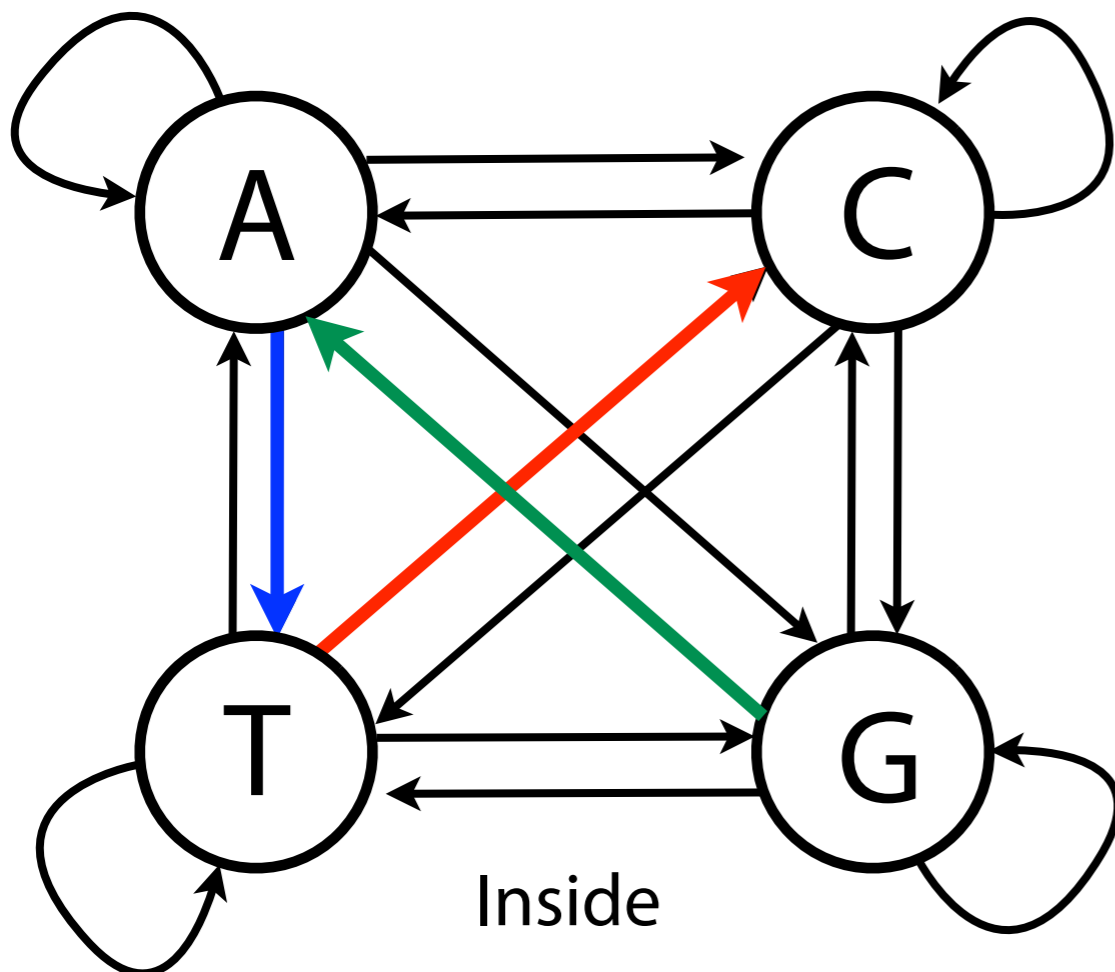
$x = \text{GATC}$

$$\begin{aligned} \log P(x) &= \sum_{i=2}^4 \log P(x_i | x_{i-1}) + \log P(x_1) \\ &= -1.525 + \\ &\quad -2.848 + \\ &\quad -2.519 + \\ &\quad -2.0 \\ &= -8.892 \end{aligned}$$

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*: $\frac{P(x) \text{ using model on the left}}{P(x) \text{ using model on the right}}$



Markov chain

```
>>> 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 ]
>>> combinedTab = numpy.log2(iTab) - numpy.log2(oTab)
>>> print(combinedTab)
    Log ratio
    | 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
```

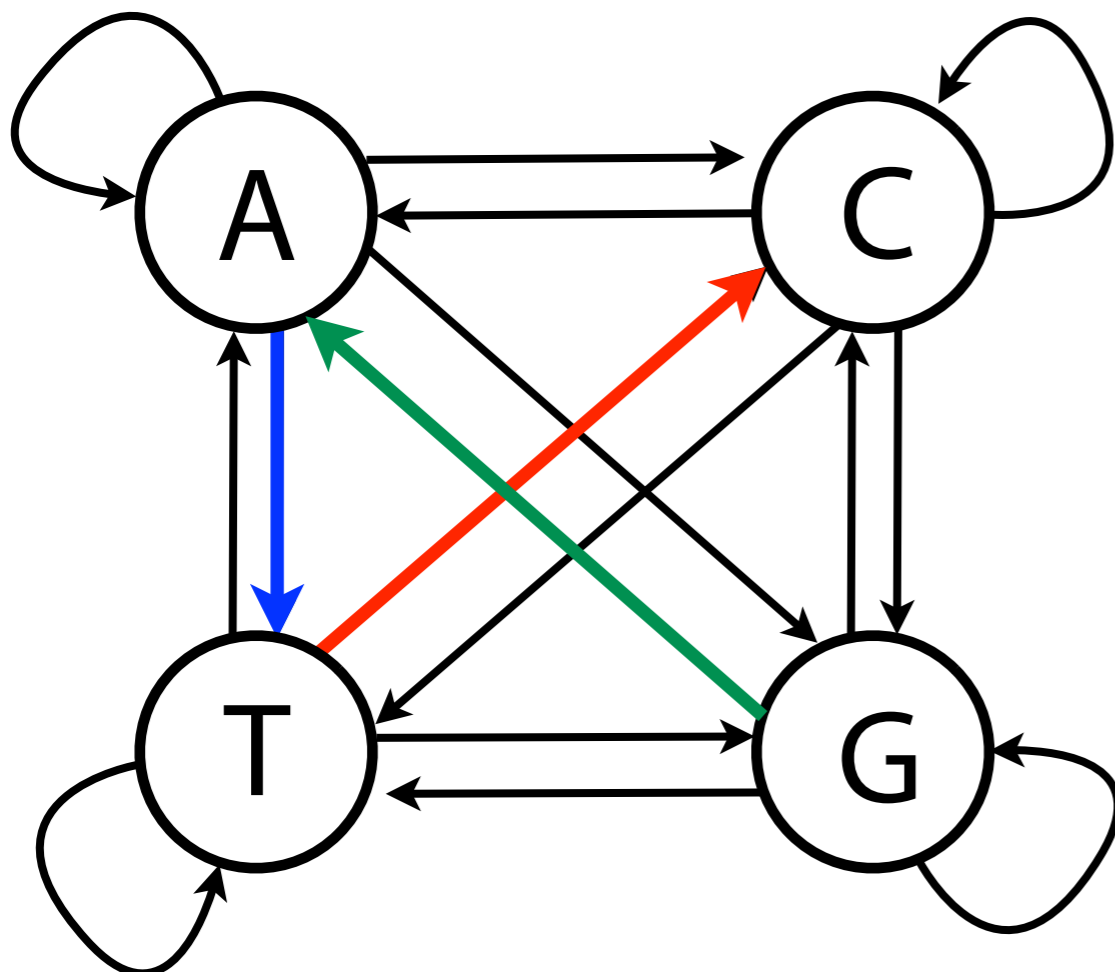
Markov chain

```

>>> iTab, oTab = islandTransitionTables(fn, ifn)
>>> combinedTab = numpy.log2(iTab) - numpy.log2(oTab)
>>> print(combinedTab)

```

X_{i-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	X_i



$x = \text{GATC}$ (Marginal probabilities ignored here)

$$\begin{aligned}
 S(x) &= 0.730 + \\
 &\quad -0.855 + \\
 &\quad -0.795 \\
 &= -0.92
 \end{aligned}$$

Negative, therefore probability with *outside* model is greater

Markov chain

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

$$S(\text{CGCGCGCGCGCGCGCGCGCGCGCGCGCGCG}) = 32.246609048$$

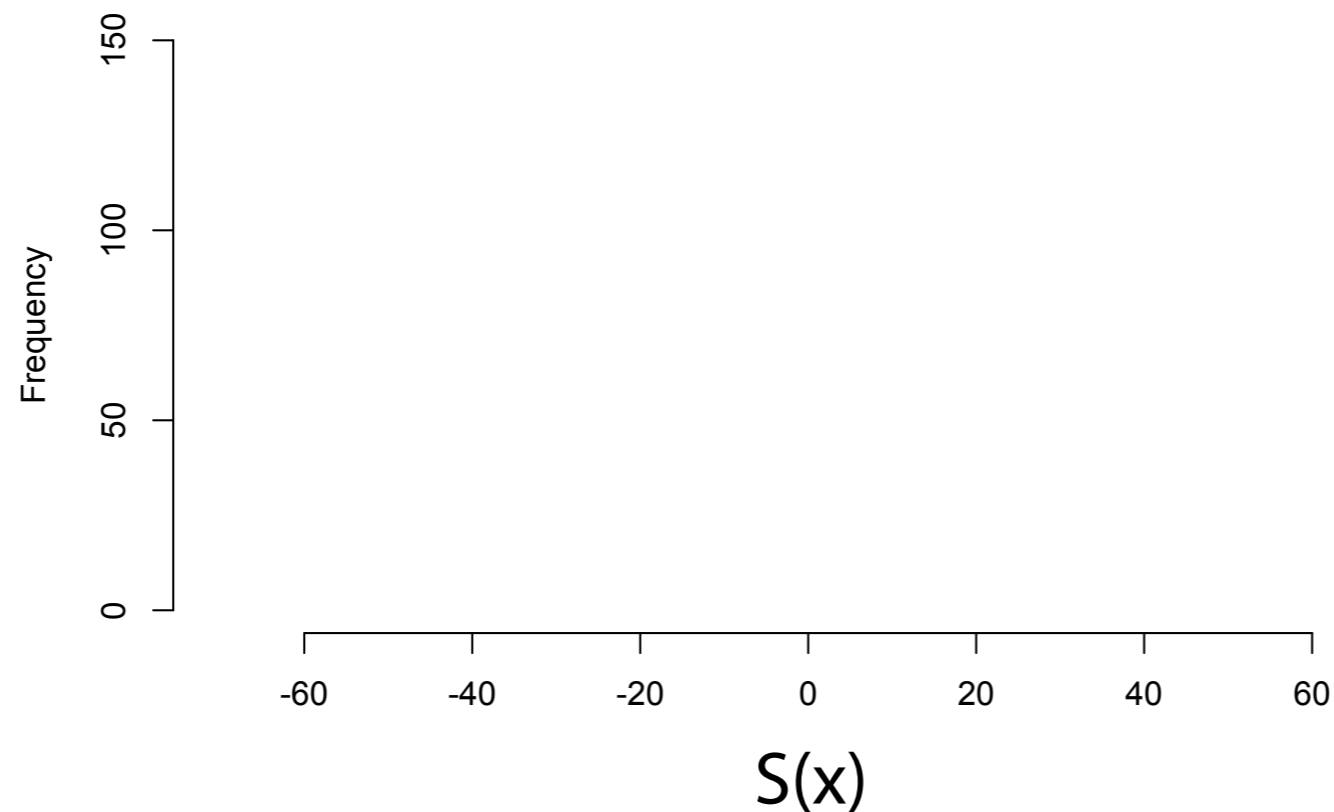
$$S(\text{ATTCTACTATCATCTATCTATCTTCT}) = -9.501209765$$

Markov chain

How well does $S(x)$ distinguish strings drawn from inside versus outside CpG islands?

Draw 1,000 100-mers from inside CpG islands on chromosome 1, and another 1,000 from outside, calculate log ratios for all

Markov chain trained on chromosome 22 islands



Markov chain

How well does $S(x)$ distinguish strings drawn from inside versus outside CpG islands?

Draw 1,000 100-mers from inside CpG islands on chromosome 1, and another 1,000 from outside, calculate log ratios for all

Markov chain trained on chromosome 22 islands

