

Markov Chains and Hidden Markov Models

Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder

Some slides are from Ben Langmead of the Johns Hopkins University.

Problem statement

Types of data from lecture
2

- Different types of data (not record or transaction types);
 - A different type of “classification” problem
- Record
 - Data Matrix
 - Document Data
 - Transaction Data
 - Graph
 - World Wide Web
 - Molecular Structures
 - Ordered
 - Spatial Data
 - Temporal Data
 - Sequential Data
 - **Genetic Sequence Data**

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

CpG islands

- Definition: CpG is the pair of nucleotides C and G appearing successively, in this order, along one DNA strand. (e.g. CGCGCGGCGGCGCG) But it can have “CC”, “GG” or even A and T in the islands.

Assignment Project Exam Help

- CpG relatively rare in most DNA sequences. However, in particular sub-sequences, which are a few hundred to a few thousand bases long, CpG is more frequent. Those sequences are called CpG islands.
<https://powcoder.com>
Add WeChat powcoder
- **Problem one:** given a short DNA sequence x, decide whether x is from CpG island.
- **Problem two:** given a genome, search for CpG islands in the genome.

Which one is harder ?

Classification problem

Given a short DNA sequence x , **decide whether x is from CpG island.**

S:AAGCCGGGAAGTTGTATG

Questions:

Is S from inside/outside CpG islands?

Training data: a bunch of sequences from inside and outside of CpG islands

<https://powcoder.com>

CpG island

CGCGCGGCCGCGGCCGA
CCGGCGCGCGCGTGCC
CCGGCGACCCGGGCGCG

.....

Not CpG island

CTACCGATCTCGCAAA
CGTAACATGACGATTGC
CCGTAATCCTTACTAG

.....

Decision tree? Naïve Bayes? kNN?
Bayes Belief Network?

Classification problem

Given a short DNA sequence x , **decide whether x is from CpG island.**

S: AAGCCGGGAAGTTGTATG

Questions:

Is S from inside/outside CpG islands?

Compare **$P(\text{inside}|S)$ and $P(\text{outside}|S)$**

<https://powcoder.com>

Training data: a bunch of sequences from inside and outside of CpG islands

Add WeChat powcoder

inside

CGCGCGGCCGCGGCCGA
CCGGCGCGCGCGTGCC
CCGGCGACCCGGGGCGCG

.....

outside

CTACCGATCTCGCAAA
CGTAACATGACGATTGC
CCGTAATCCTTACTAG

.....

How to compute the posterior probability

S2: CCGCGCGC

$$P(\text{inside}|S2) = P(\text{in} | \text{CCGCGCGC})$$
$$=$$

<https://powcoder.com>
P(in) is the prior probability, reflecting the fraction of CpG islands. Usually $P(\text{in}) < P(\text{out})$. P(in): the probability that the input sequence is inside a CpG island. P(out): the probability that the input sequence is outside a CpG island.

[Add WeChat powcoder](https://powcoder.com)

To focus on the interpretation of Markov chain model, we will assume that $P(\text{in})$ is equal to $P(\text{out})$. Thus

$P(\text{inside}|S2)$ and $P(\text{outside}|S2)$ ' ranking is determined by
 $P(S2|\text{in})$ and $P(S2|\text{out})$

Q: how to compute from the training data?

Sequence models

Let $P(x)$ be the probability of sequence x as assigned by the model

$$P(x) = P(\underbrace{x_k, x_{k-1}, \dots, x_1}_{\text{order matters}})$$

Joint probability of all sequence items appearing as they do
(**order matters**)

Assignment Project Exam Help

To estimate $P(x)$, count $\#$ times x appears in the training set
labeled *inside* divided by $\#$ times x appears in training set

But for sufficiently long k , we might not see *any* occurrences of x , or very few. Joint probabilities for rare events are hard to estimate well.

Sequence models

Bayes Belief network for 1st-order
Markov model

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

Re-write with conditional probability:

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

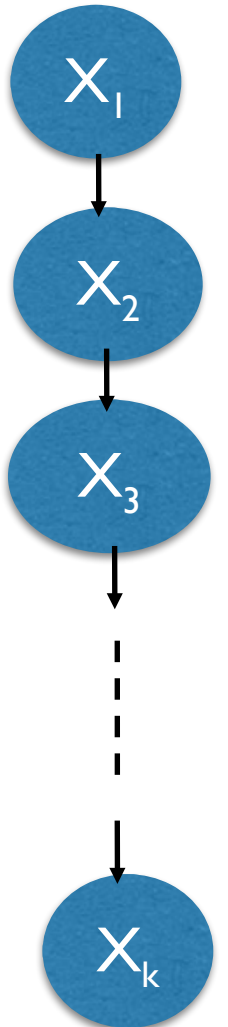
$$= P(x_k | x_{k-1}, \dots, x_1) P(x_{k-1} | x_{k-2}, \dots, x_1) P(x_{k-2}, \dots, x_1)$$

(etc)

Add a **simplifying assumption**: to know the probability of having a particular item x_k , *we only have to know the previous item*: x_{k-1}

Formally: random variable 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

A **simplifying assumption**: to know the probability of having a particular item x_k , *we only have to know the previous item*: x_{k-1}

$$\begin{aligned} P(x) &= P(x_k, x_{k-1}, \dots x_1) \\ &= P(x_k \mid x_{k-1}, \dots x_1) P(x_{k-1}, \dots x_1) \\ &= P(x_k \mid x_{k-1}, \dots x_1) P(x_{k-1} \mid x_{k-2}, \dots x_1) P(x_{k-2}, \dots x_1) \\ &\quad \text{(etc)} \end{aligned}$$

Assignment Project Exam Help
<https://powcoder.com>
drop *drop* (bunch more drops once this is expanded)
Add WeChat powcoder

$$\approx P(x_k \mid x_{k-1}) P(x_{k-1} \mid x_{k-2}) \dots P(x_2 \mid x_1) P(x_1)$$

Markov property / Markov assumption

It's a big assumption, but it's often reasonable and it makes the model much easier to work with

Markov chain

Assigning a probability to a sequence using Markov property:

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

Assignment Project Exam Help

Say x is a nucleotide k -mer <https://powcoder.com>

Add WeChat powcoder

$P(x_i | x_{i-1})$ 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, $P(G | C)$?

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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Markov chain

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

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

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

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

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

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

(etc)

Markov chain (1st order)

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

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

X_{i-1}	A	0.18544138	0.27640458	0.40091352	0.13724053
	C	0.18958227	0.35905063	0.25324026	0.19812684
	G	0.17268916	0.33011349	0.35610656	0.14109079
	T	0.09410222	0.34163592	0.37686698	0.18739488
		A	C	G	T
		X_i			
		$P(T G)$			

Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder

Rows sum to 1

Markov chain

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

↑ *Inside* ↓

↑ *Outside* ↓

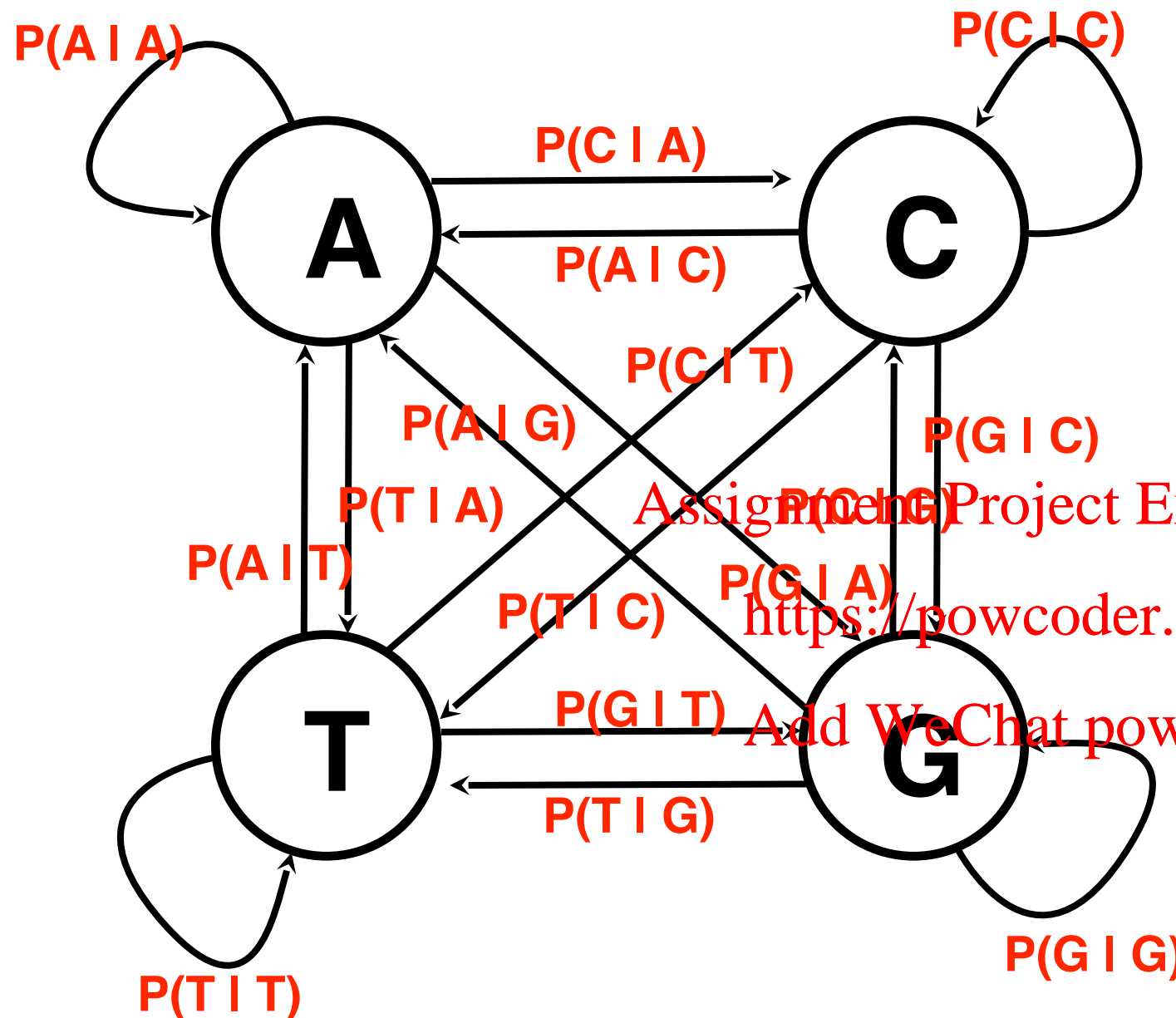
```
>>> iTab, nTab = islandTransitionTables(fn, ifn)
>>> print iTab
[[ 0.18544138  0.27640458  0.40091352  0.13724053]
 [ 0.18958227  0.35905063  0.25324026  0.19812684]
 [ 0.17268916  0.33011349  0.35610656  0.14109079]
 [ 0.09410222  0.34163592  0.37686698  0.18739488]]
>>> print nTab
[[ 0.2948135  0.19467897  0.28696205  0.22354548]
 [ 0.32681187  0.29415529  0.06172587  0.31730697]
 [ 0.25713351  0.23354071  0.29423494  0.21509084]
 [ 0.17956538  0.23250026  0.29462341  0.29331096]]
```

A C G T

Notice anything interesting about the outside conditional probabilities?

$P(G | C)$ is low, matching our expectation that there are few CpGs outside islands

Markov chain



Markov chain is a probabilistic automaton

Each edge has a *transition probability*: probability that edge's destination is the next node visited after edge's source

Here, nodes labels are symbols and transition labels are conditional probabilities

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

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

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

Markov chain

Recall how we assign a probability to a single string

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

Assignment Project Exam Help

For simplicity, drop $P(x_1)$ <https://powcoder.com>

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

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

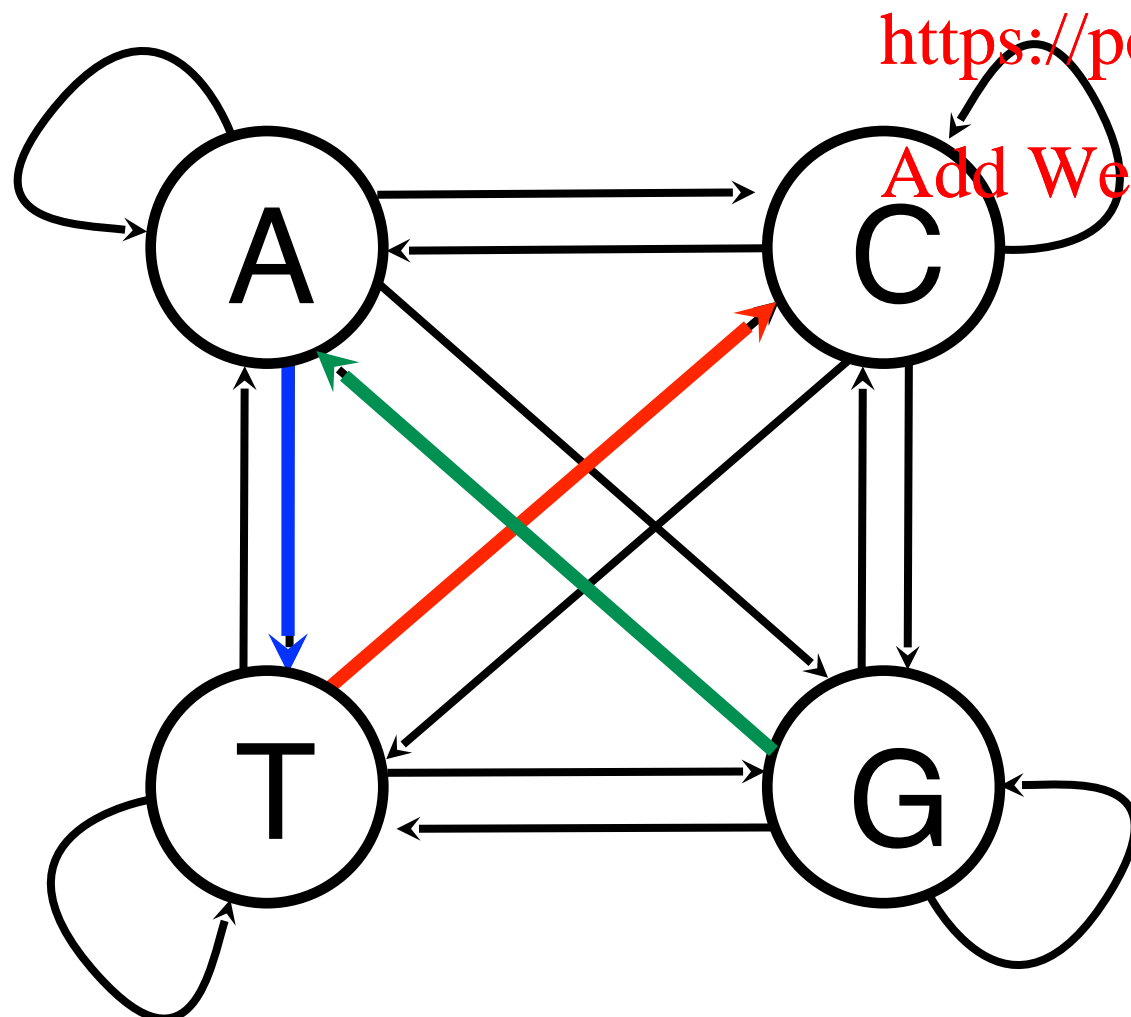
$P(x)$ now equals product of all the Markov chain edge weights on our string-driven walk through the chain

Markov chain

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

X_{i-1}	A	0.18544138	0.27640458	0.40091352	0.13724053
C		0.18958227	0.35905063	0.25324026	0.19812684
G		0.17268916	0.33011349	0.35610656	0.14109079
T		0.09410222	0.34163592	0.37686698	0.18739488
	A	C	G	T	

X_i



<https://powcoder.com>

Add WeChat powcoder

$x = \text{GATC}$

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

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

$$= 0.34163592 *$$

$$0.13724053 *$$

$$0.17268916$$

$$= 0.00809675$$

Markov chain

To avoid repeated multiplies yielding small numbers, we switch to log domain

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

Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder

Switching to log domain,
multiplies become adds

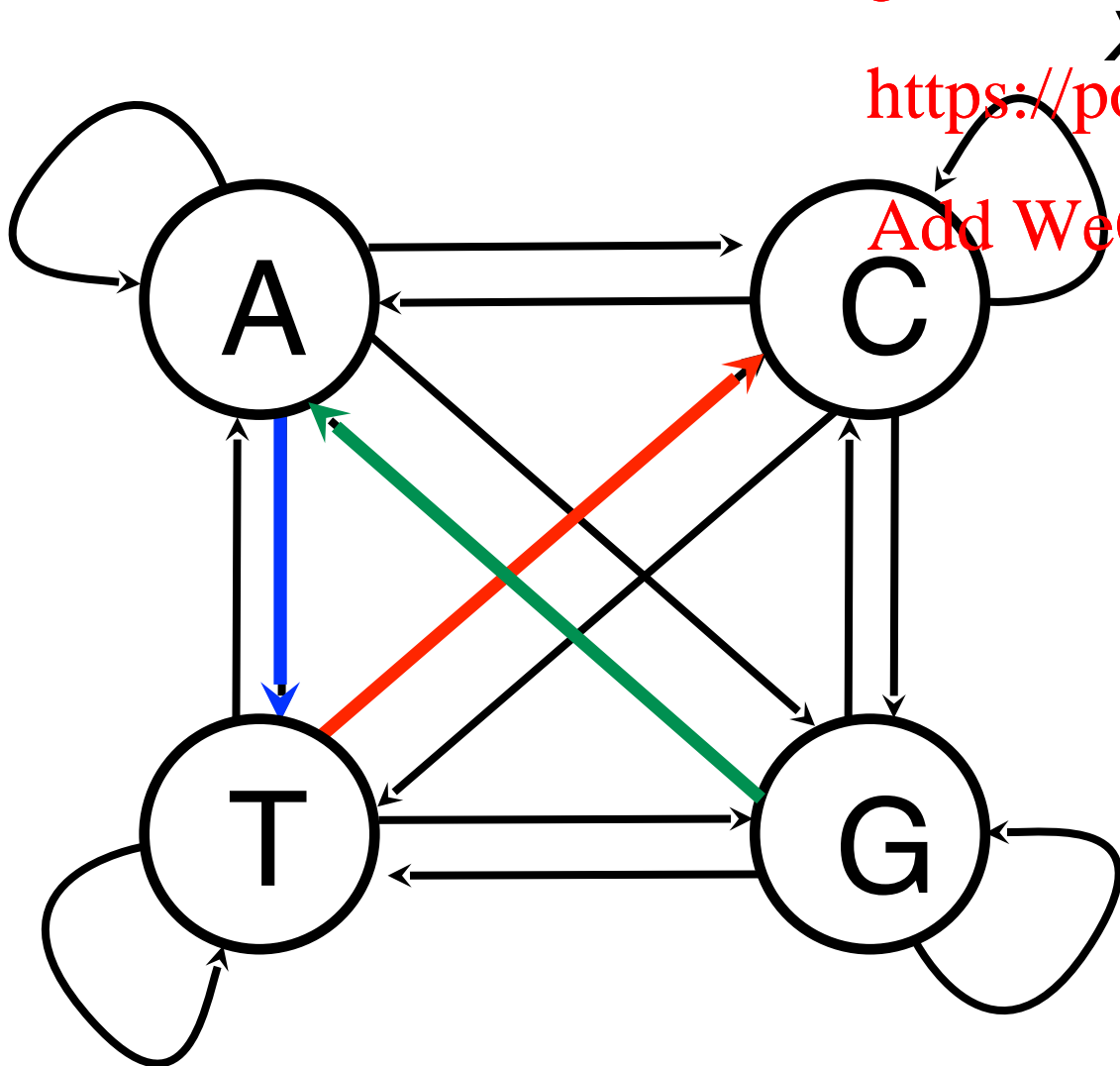
I'll use base-2 logs

Markov chain

```
>>> iTab, nTab = islandTransitionTables(fn, ifn)
>>> print numpy.log2(iTab)
```

	A	C	G	T
A	-2.43096492	-1.85514658	-1.31863704	-2.86522151
C	-2.39910406	-1.4777408	-1.98142131	-2.33550376
G	-2.53375061	-1.59896599	-1.48961909	-2.82530423
T	-3.40962748	-1.54946844	-1.40787269	-2.41584653

Assignment Project Exam Help



$x = \text{GATC}$

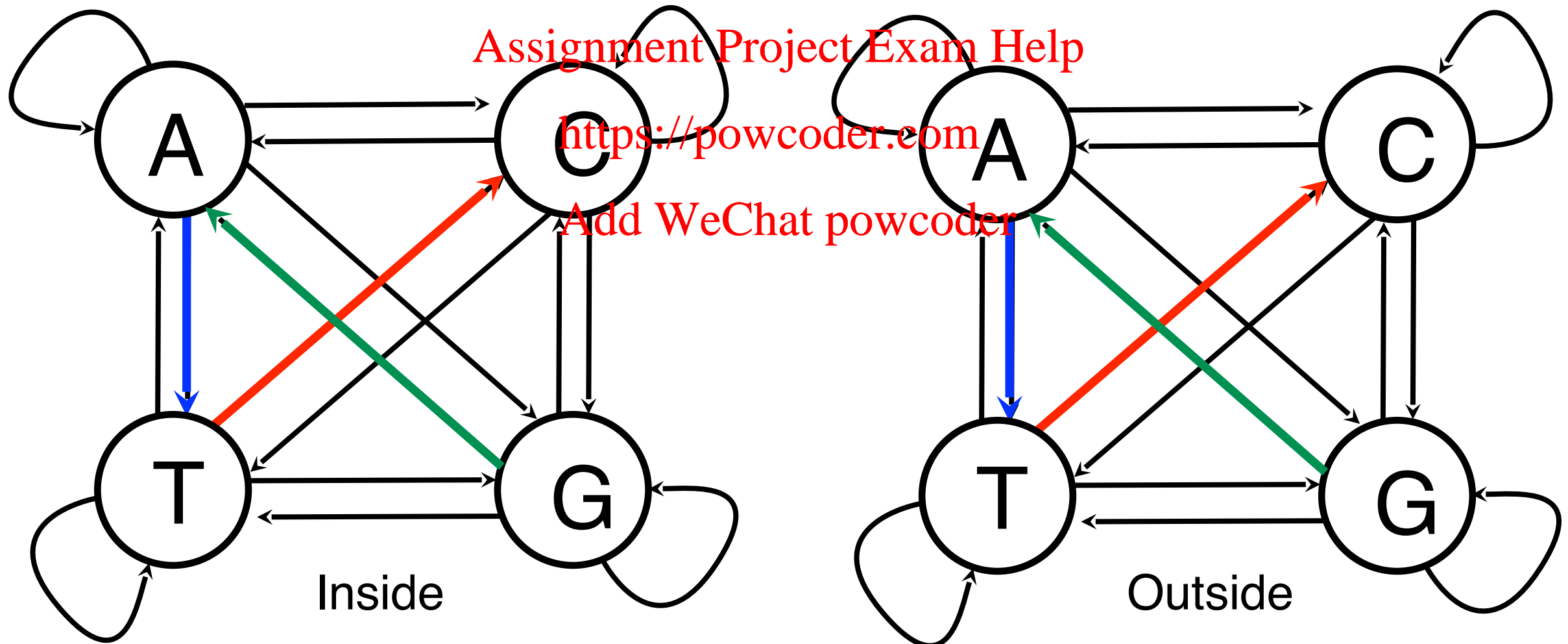
powcoder

$$\log P(x) = \sum_{i=2}^4 \log P(x_i / x_{i-1})$$
$$= -1.54946844 +$$
$$-2.86522151 +$$
$$-2.53375061$$
$$= -7.30174249$$

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

Taking log, we get *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 and log ratio is > 0 . Otherwise, fraction is ≤ 1 and log ratio is ≤ 0 .

$$\begin{aligned}
 S(x) &= \log \frac{P(x) \text{ inside CpG}}{P(x) \text{ outside CpG}} \\
 &= \log [P(x) \text{ inside CpG}] - \log [P(x) \text{ outside CpG}] \\
 &= \sum_{i=2}^k \left(\log [P(x_i | x_{i-1}) \text{ inside CpG}] \right) - \sum_{i=2}^k \log ([P(x_i | x_{i-1}) \text{ outside CpG}]) \\
 &= \sum_{i=2}^k \left(\log [P(x_i | x_{i-1}) \text{ inside CpG}] - \log [P(x_i | x_{i-1}) \text{ outside CpG}] \right)
 \end{aligned}$$

New table: take elementwise log of the inside/outside tables, subtract outside from inside

Markov chain

```
>>> iTab, nTab = islandTransitionTables(fn, ifn)
>>> print iTab
↑ A [[ 0.20328697  0.26144423  0.40629367  0.12897512]
Inside ↓ C [ 0.18175425  0.35880255  0.24915835  0.21028485]
      G [ 0.17900663  0.32594344  0.35910409  0.13594584]
      T [ 0.09718687  0.34541934  0.35518406  0.20220973]]
>>> print nTab
↑ A [[ 0.32756059  0.17183665  0.24355314  0.25704963]
Outside ↓ C [ 0.35218354  0.25880566  0.04404104  0.34496977]
      G [ 0.28883529  0.20906356  0.25862313  0.24347803]
      T [ 0.21890134  0.20417181  0.24905103  0.32789582]]
>>> lrTab = numpy.log2(iTab) - numpy.log2(nTab)
>>> print lrTab
↑ A [[ -0.68824404  0.6054655  0.73828635 -0.99495413]
Log ↓ C [ -0.95433841  0.471321  2.5001426 -0.71412499]
ratio ↓ G [ -0.69023394  0.64068002  0.47355078 -0.84075959]
      T [ -1.17144749  0.75856518  0.51224132 -0.69738508]]
      A          C          G          T
```

Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder

Markov chain

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

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

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

Assignment Project Exam Help

<https://powcoder.com>

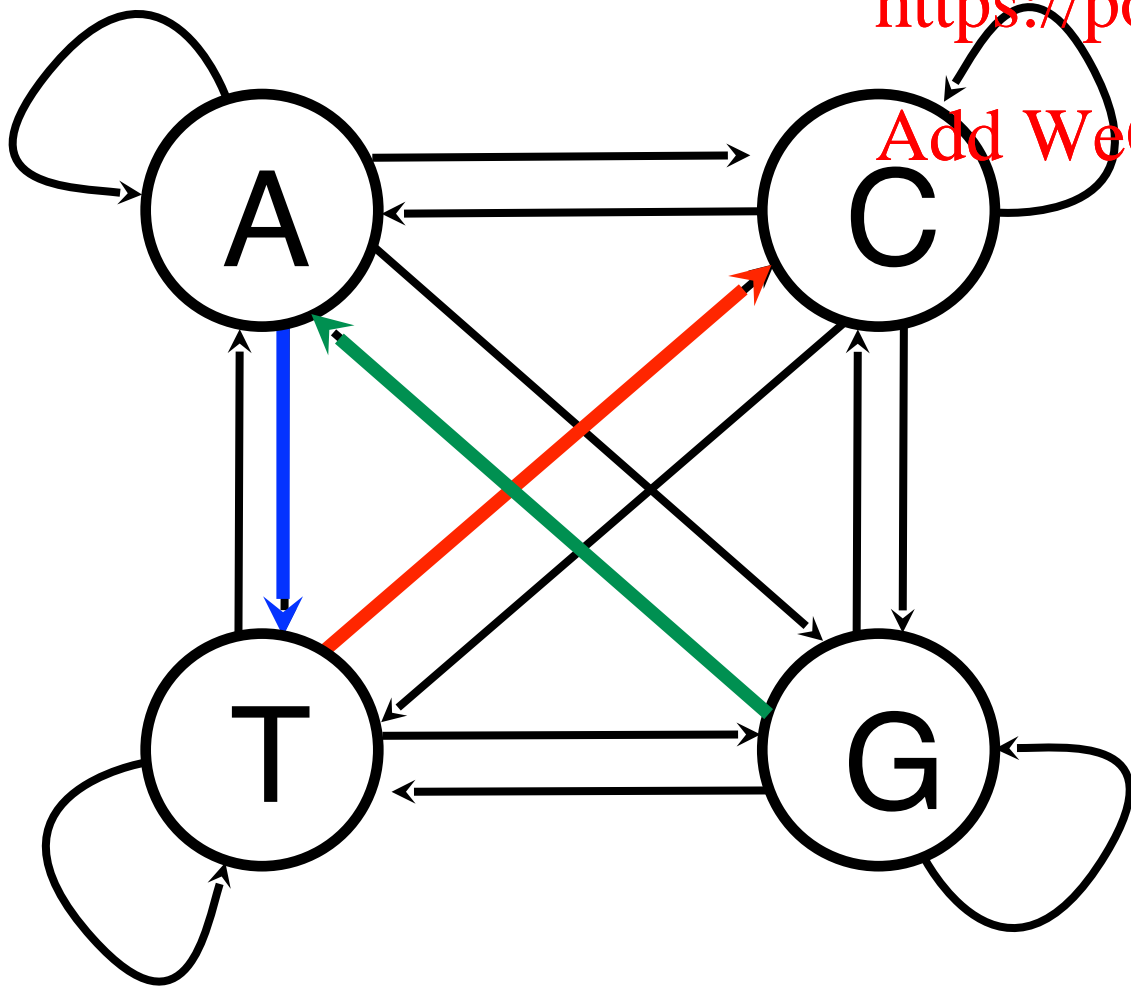
Add WeChat powcoder

Markov chain

```
>>> iTab, nTab = islandTransitionTables(fn, ifn)
>>> lrTab = numpy.log2(iTab) - numpy.log2(nTab)
>>> print lrTab
```

	A	C	G	T
A	-0.66883939	0.50568449	0.48243108	-0.70385181
C	-0.78563635	0.28760934	2.03655959	-0.67945489
G	-0.57434013	0.49928806	0.27534041	-0.60832223
T	-0.9322086	0.55522735	0.35518335	-0.6463494

Assignment Project Exam Help



<https://powcoder.com>

Add WeChat powcoder

$$S(x) = 0.55522735 + 0.70386181x + 0.57434013x^2 = -0.72297459$$

Negative, so probability with *outside* model is greater

Markov chain

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

$S(\text{CGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG}) \approx 2.246609048$

Assignment Project Exam Help

$S(\text{ATTCTACTATCATCTATCTATCTTCT}) \approx 9.501209765$

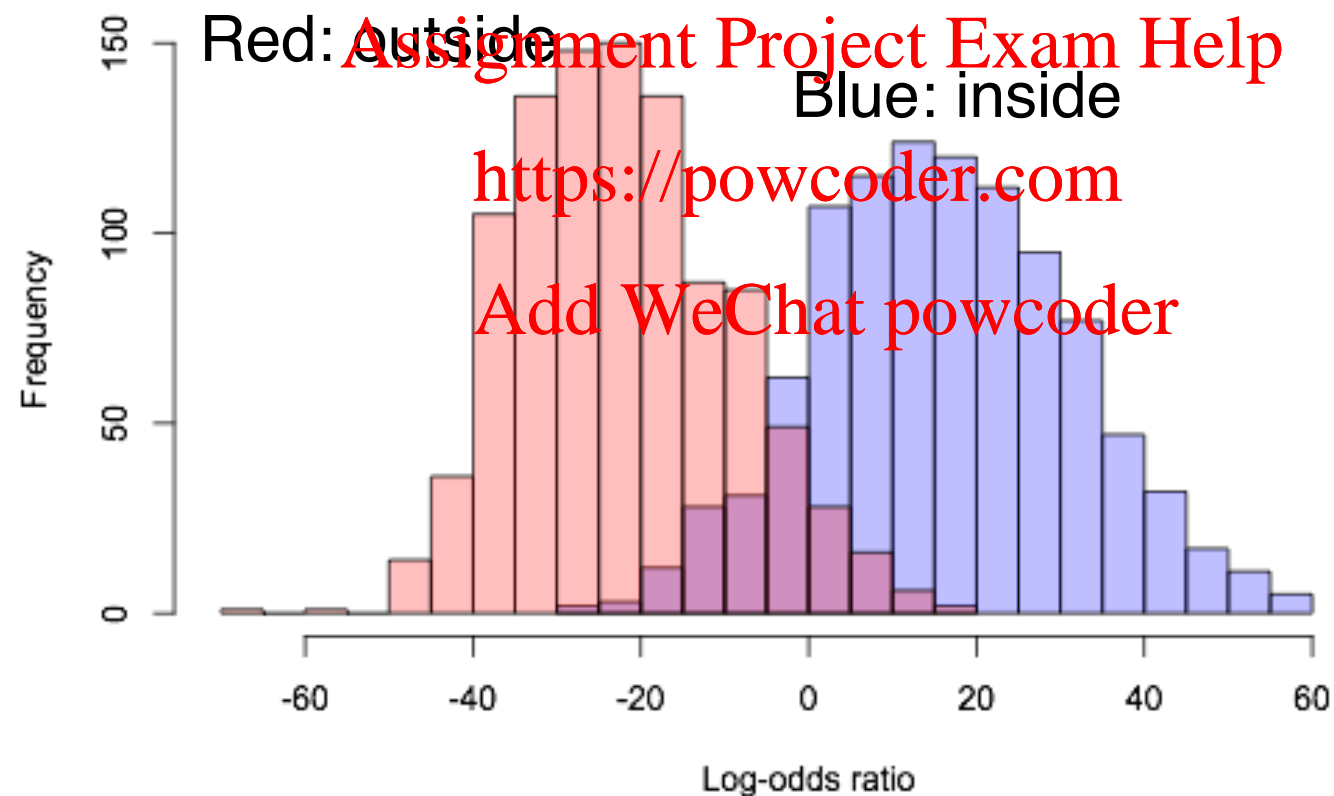
<https://powcoder.com>

Add WeChat powcoder

Markov chain: experiment

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

Trained markov chain on dinucleotides from chromosome 22



Markov chain

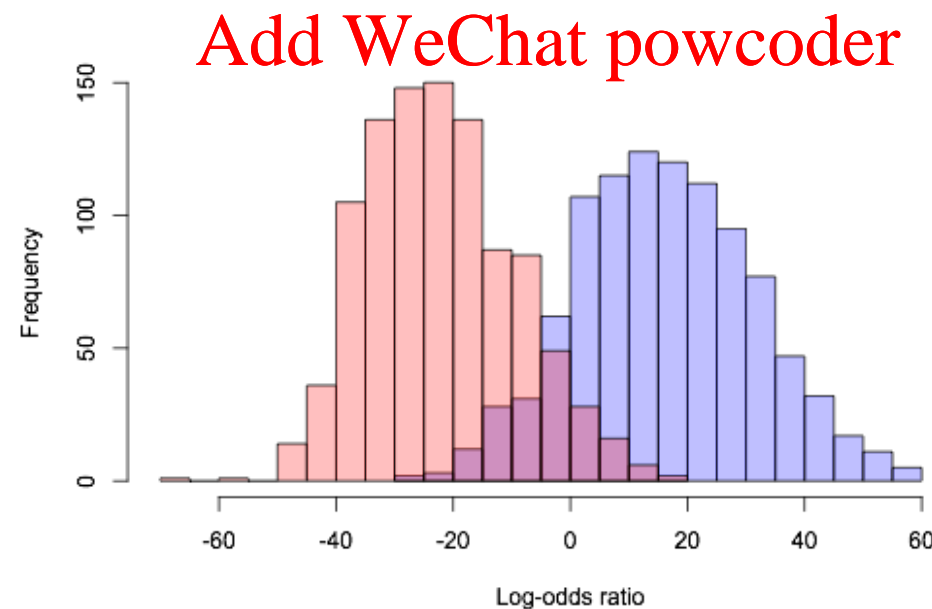
Markov property made our problem very tractable

$P(x_i | x_{i-1})$ 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

... and discriminates well between inside & outside examples in CpG island example



Higher order Markov model

$$P(x) \approx P(x_k | x_{k-1} x_{k-2}) P(x_{k-1} | x_{k-2} x_{k-3}) \dots P(x_3 | x_2 x_1) P(x_2 x_1)$$

2nd order

Assignment Project Exam Help

$$P(x) \approx P(x_k | x_{k-1} x_{k-2} x_{k-3}) P(x_{k-1} | x_{k-2} x_{k-3} x_{k-4}) \dots P(x_3 x_2 x_1)$$

3rd order

<https://powcoder.com>

Add WeChat powcoder