# Hidden Markov Models

Many pages are from Ben Langmead

# Problem two: given a genome, search for CpG islands in the genome.

```
ATCATGCTGGAAATTGAAGAGGAACGCAGCCAGGTCTTGGCAGAGATCTCTGAGGATAATCATGCTCAGGTTGCAGAAGACAGTGGGACAACATCACGT
GCTGGCCTGAGGCGCAGGTGGGAGCAGTTGTGGTGAAGCCATGCCCAAAGTACTTCAGGTTCCTCACCACTTTTCTCGGAAATGTGAGCAGGAATTGCAC
AAGTCAGGGCTGGACAGATGTGTACCCTGCACCGTATGCTGTTGCTTGTGGATACGATTCCAACAGCACACCAGGGGAAGAGCAAACGGCGTTCTATGGC
ACGGTCAAGACCGGCTACACCATTGGACATACCTTGTCACTCATTGCTCTCACAGCTGTATGATCATTTTCTGTCTCTTCAGAAAACTACACTGTACTC
GAAATTATATTCATATGCACCTCTTCATGTCCTTTATAATGAGAGCCATTGCAGTCTTCCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGC
GCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGC
GCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGaCTAAGCTTCAAGGAGGGTGACCTCATTACCCTGCTGGTGC
CTGAGGCCCGTGATGGCTGGCACTACGGAGAGAGTGAGAAGACCAAGATGCGGGGCTGGTTTCCCTTCTCCTACACCCGGGTCCTGGACAGCGACGGCAG
TGACAGGCTGCACATGAGCTTGCAGCAAGGGAAGAGCAGCAGCACGGGCAACCTCCTAGACAAGGACGACTTGGCCATCCCGCCCCCTGACTACGGCGCC
ACCTCCCGGGCCTTCCCTGCCCAGACAGCCAGTGGCTTCAAGCAGAGGCCCTACAGCGTGGCCGTGCCTGCCTTCTCCCAGGGTCTGGAAGACTATGGAG
                              CACGGTCCATGAGCAGTGGCAGCGGCACG
```

# Sequence models

Can we use Markov chains to pick out CpG islands from the rest of the genome?

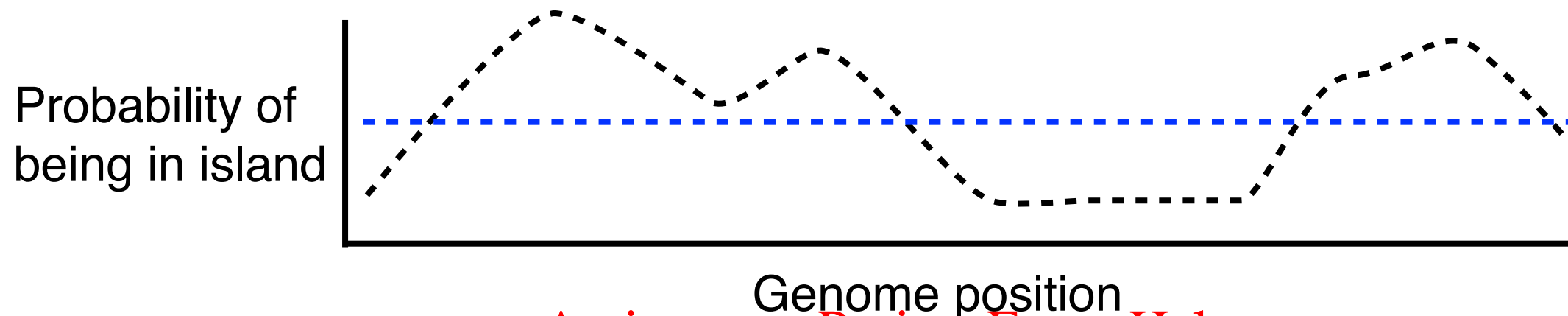Markov chain assigns a score to a string; doesn't naturally give a "running" score across a long sequence



Probability of being in island

Genome position

We could use a *sliding window*

(a) Pick window size *w*, (b) score every *w*-mer using Markov chains, (c) use a cutoff to find islands

Smoothing before (c) might also be a good idea

# Sequence models



Probability of being in island

Genome position

Choosing $w$ involves an assumption about how long the islands are

If $w$ is too large, we'll miss small islands

If $w$ is too small, we'll get many small islands where perhaps we should see fewer larger ones

In a sense, we want to switch *between Markov chains* when entering or exiting a CpG island

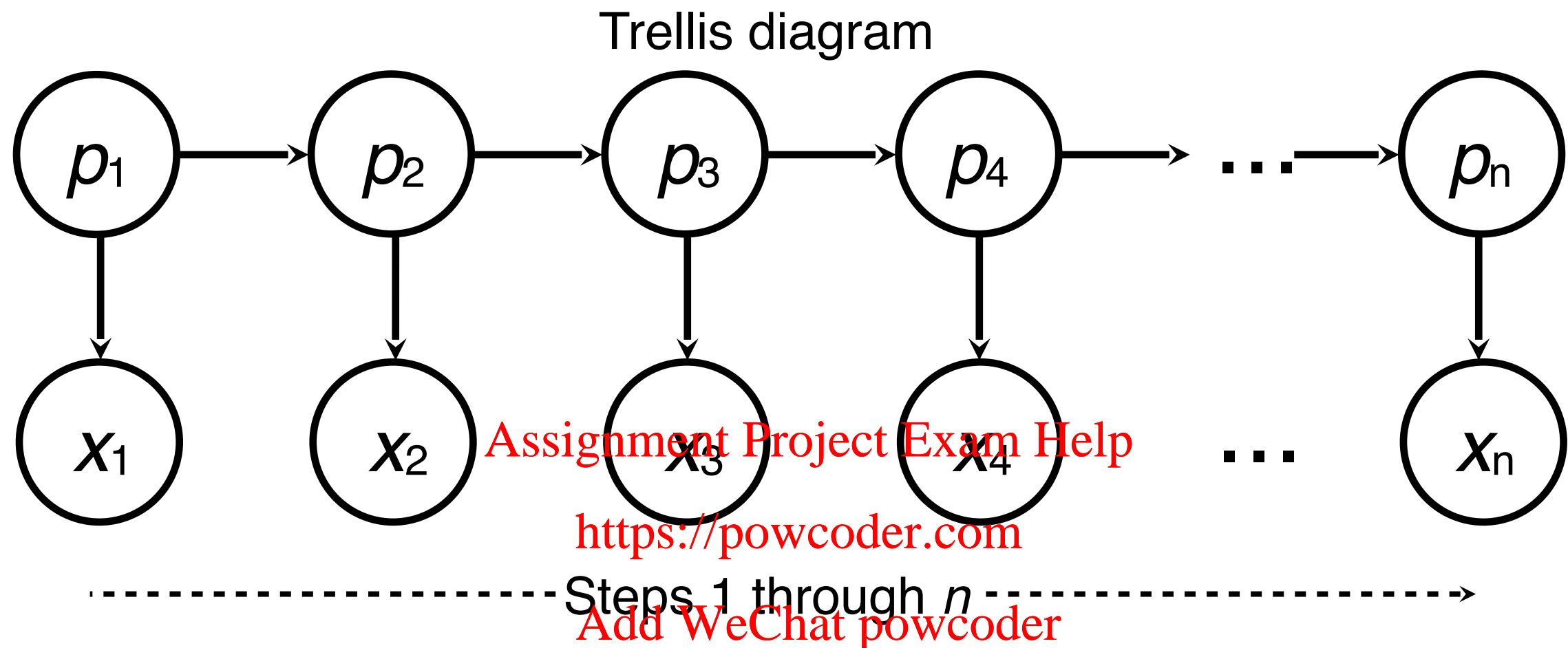# Ideal solution: every base (rather than a window) is assigned with a label, inside or outside

GAGAGCCATTGCAGTCTTCGCGCGCGCGCGCGCGCGGCGCGCGCGCGCGCTTACCCTGCGTG
CCT



out  in in                    in  out

# Hidden Markov Model
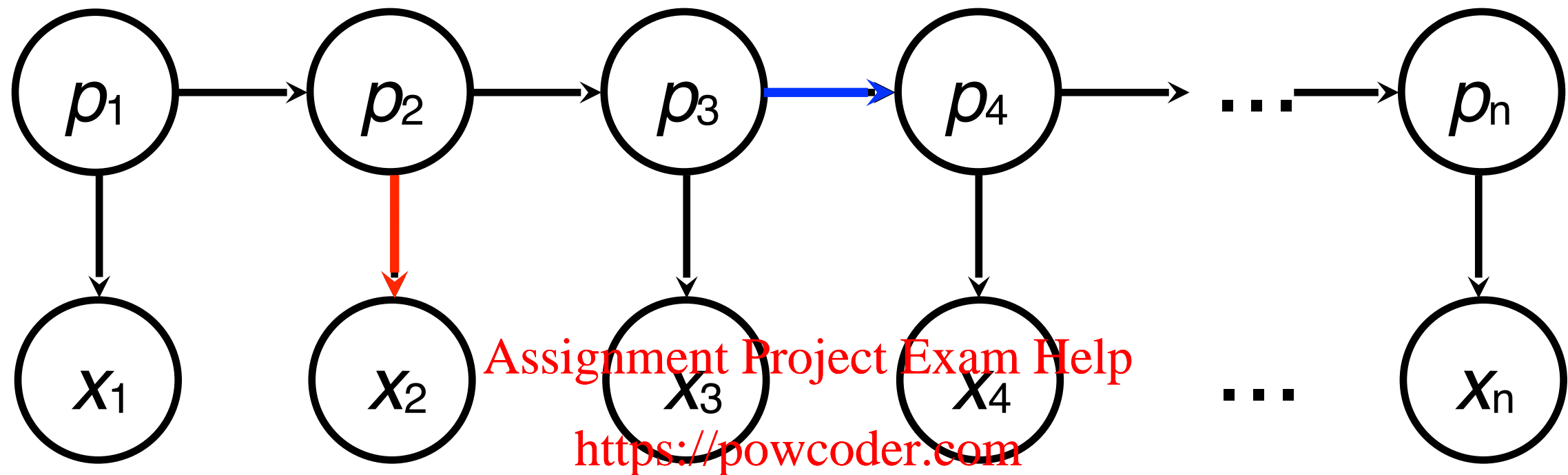
Trellis diagram

Steps 1 through $n$

$p = \{ p_1, p_2, ..., p_n \}$ is a sequence of *states* (AKA a *path*). Each $p_i$ takes a value from set $Q$. We **do not** observe $p$.

$x = \{ x_1, x_2, ..., x_n \}$ is a sequence of *emissions*. Each $x_i$ takes a value from set $\Sigma$. We **do** observe $x$.

Our goal: given the observation $x_1, x_2, \ldots x_n$, derive the hidden states $p_1, p_2, p_3, \ldots, p_n$

# Hidden Markov Model



Like for Markov chains, edges capture conditional independence:

$x_2$ is conditionally independent of everything else given $p_2$

$p_4$ is conditionally independent of everything else given $p_3$
Probability of being in a particular state at step $i$ is known once we know what state we were in at step $i$-1. Probability of seeing a particular emission at step $i$ is known once we know what state we were in at step $i$.

# Hidden Markov Model

Example: occasionally dishonest casino

Dealer repeatedly flips a coin.  Sometimes the coin is *fair*, with P(heads) = 0.5, sometimes it's *loaded*, with P(heads) = 0.8.  Dealer occasionally switches coins, invisibly to you.

How does this map to an HMM?

$$p_1 \rightarrow p_2 \rightarrow p_3 \rightarrow p_4 \rightarrow \cdots \rightarrow p_n$$

$$p_1 \downarrow x_1 \quad p_2 \downarrow x_2 \quad p_3 \downarrow x_3 \quad p_4 \downarrow x_4 \quad \cdots \quad p_n \downarrow x_n$$

# Hidden Markov Model
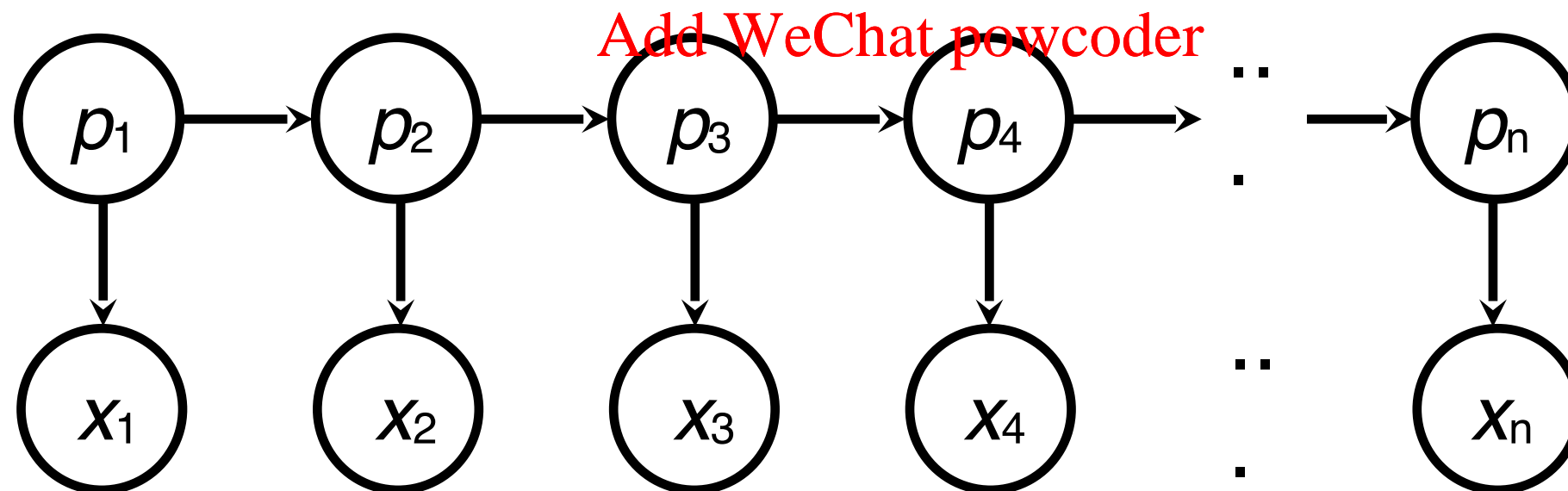
Example: occasionally dishonest casino

Dealer repeatedly flips a coin. Sometimes the coin is *fair*, with P(heads) = 0.5, sometimes it's *loaded*, with P(heads) = 0.8. Dealer occasionally switches coins, invisibly to you.
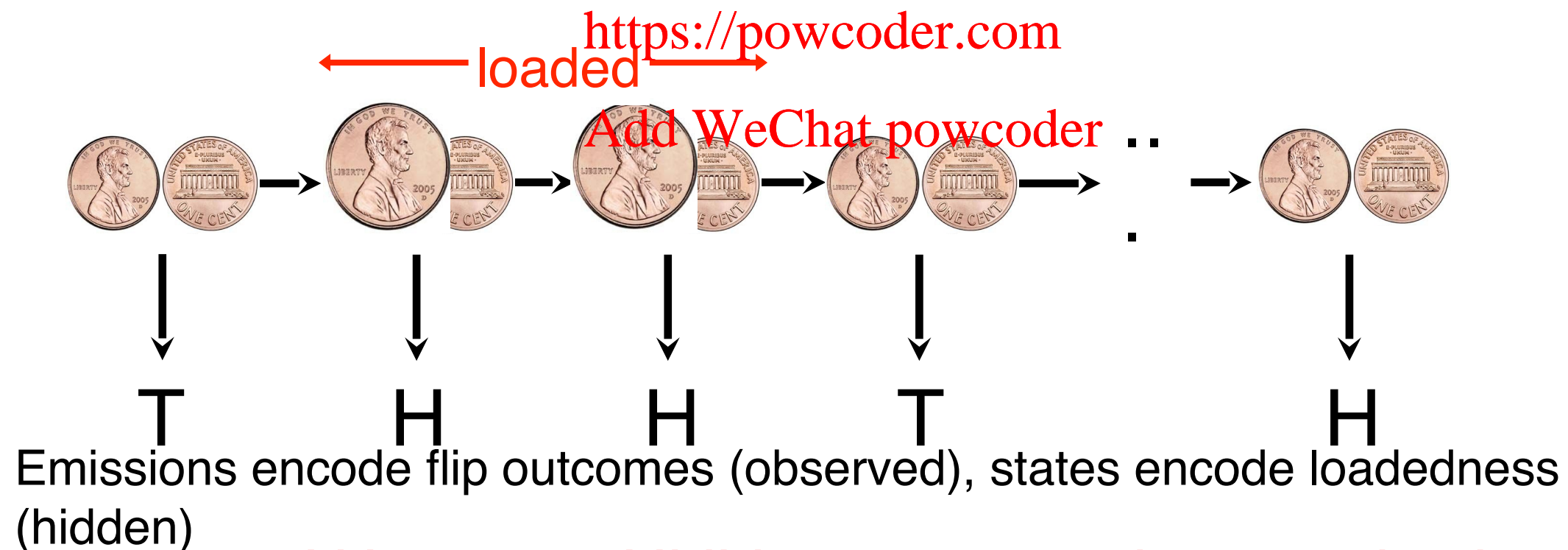
How does this map to an HMM?

← loaded →



| T | H | H | T | H |

Emissions encode flip outcomes (observed), states encode loadedness (hidden)

We can use HMM to estimate what coin the dealer is using for each flip

# Hidden Markov Model

States encode which coin is used

**F** = fair
**L** = loaded

Emissions encode flip outcomes

**H** = heads
**T** = tails

$p_1$ F L

$p_2$ F L

$p_3$ F L

$p_4$ F L

.. .

$p_n$ F L

$x_1$ H T

$x_2$ H T

$x_3$ H T

$x_4$ H T

.. .

$x_n$ H T

# Hidden Markov Model

Example with six coin flips:



$p$ = FFLLLL

$x$ = HTHHHH

# Hidden Markov Model

$$P( p_1, p_2, ..., p_n, x_1, x_2, ..., x_n ) = ?$$

$$P( p_1, p_2, ..., p_n, x_1, x_2, ..., x_n ) = \prod_{k=1}^{n} P( x_k \mid p_k ) \prod_{k=2}^{n} P( p_k \mid p_{k-1} ) P( p_1 )$$

$| Q | \times | \textstyle\sum |$ emission matrix $E$ encodes $P( x_i \mid p_i )$  $E[p_i, x_i] = P( x_i \mid p_i )$

$| Q | \times | Q |$ transition matrix $A$ encodes $P( p_i \mid p_{i-1} )$  $A[p_{i-1}, p_i] = P( p_i \mid p_{i-1} )$

$| Q |$ array $I$ encodes initial probabilities of each state  $I[p_i] = P( p_1 )$

Q: the state set. $\sum$ : the symbol set

# Hidden Markov Model

Dealer repeatedly flips a coin.  Coin is sometimes *fair,* with P(heads) = 0.5, sometimes *loaded*, with P(heads) = 0.8.  Dealer occasionally switches coins, invisibly to you.

After each flip, dealer switches coins with probability 0.4

**A:**

transition

|  | F | L |
|---|---|---|
| F | 0.6 | 0.4 |
| L | 0.4 | 0.6 |

**E:**

emission

|  | H | T |
|---|---|---|
| F | 0.5 | 0.5 |
| L | 0.8 | 0.2 |

$|Q| \times |\Sigma|$ emission matrix $E$ encodes P( $x_i$ | $p_i$ )  $E[p_i, x_i] = $ P( $x_i$ | $p_i$ )

$|Q| \times |Q|$ transition matrix $A$ encodes P( $p_i$ | $p_{i-1}$ )  $A[p_{i-1}, p_i] = $ P( $p_i$ | $p_{i-1}$ )

# Hidden Markov Model

Given *A* & *E* (right), what is the joint probability of **p** & **x**?

| **A** | F | L |
|---|---|---|
| F | 0.6 | 0.4 |
| L | 0.4 | 0.6 |

| **E** | H | T |
|---|---|---|
| F | 0.5 | 0.5 |
| L | 0.8 | 0.2 |

| **p** | F | F | F | L | L | L | F | F | F | F | F |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **x** | T | H | T | H | H | H | T | H | T | T | H |
| **P( $x_i$ I $p_i$ )** | 0.5 | 0.5 | 0.5 | 0.8 | 0.8 | 0.8 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 |
| **P( $p_i$ I $p_{i-1}$ )** | - | 0.6 | 0.6 | 0.4 | 0.6 | 0.6 | 0.4 | 0.6 | 0.6 | 0.6 | 0.6 |

If P( $p_1$ = F ) = 0.5, then joint probability = $0.5^9 \ 0.8^3 \ 0.6^8 \ 0.4^2$

= 0.0000026874

# Hidden Markov Model

Given flips, can we say when the dealer was using the loaded coin? How many different paths can generate X?

| T | H | T | H | H | H | T | H | T | T | H |

X =

We want to find $p^*$, the most likely path given the emissions.

$$p^* = \operatorname*{argmax}_{p} P(\, p \mid x \,) = \operatorname*{argmax}_{p} P(\, p, x \,)$$

How to find p*? How many different state path p can produce x?

This is *decoding*.  *Viterbi* is a common decoding algorithm.

# Viterbi (from Wikipedia)

## Andrew Viterbi

From Wikipedia, the free encyclopedia

**Andrew James Viterbi** (born Andrea Giacomo Viterbi; March 9, 1935) is an American electrical engineer and businessman who co-founded Qualcomm Inc. and invented the Viterbi algorithm. He is currently Presidential Chair Professor of Electrical Engineering at the University of Southern California's Viterbi School of Engineering, which was named in his honor in 2004 in recognition of his $52 million gift.

**Contents** [hide]

1 Early life
2 Education
3 Further career
4 Personal life
5 References
6 Further reading
7 See also
8 External links

Viterbi School of Engineering, west wall

## Early life [ edit ]

Viterbi was born to Italian Jewish family[2] in Bergamo, Italy and emigrated with them to the United States two years before World War II. His original name was Andrea, but when he was naturalized in the US, his parents anglicized it to Andrew.

**Andrew J. Viterbi**[1]

| | |
|---|---|
| **Born** | Andrea Giacomo Viterbi<br>March 9, 1935 (age 85)<br>Bergamo, Italy |
| **Citizenship** | American |

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# Hidden Markov Model: Viterbi algorithm

Bottom-up dynamic programming

$S_{k, i}$ = score/probability of the most likely path up to step $i$ with $p_i = k$

Start at step 1, calculate successively longer $S_{k, i}$'s

Question for you: use $S_{k,l}$ to represent the maximum probability of seeing x

$S_{Fair, 3}$

# Hidden Markov Model: Viterbi algorithm

Given transition matrix *A* and emission matrix *E* (right), what is the most probable path *p* for the following *x*?

Initial probabilities of F/L are 0.5 (no preference for fair or loaded coin)

| *A* | F | L |
|-----|-----|-----|
| F | 0.6 | 0.4 |
| L | 0.4 | 0.6 |

| *E* | H | T |
|-----|-----|-----|
| F | 0.5 | 0.5 |
| L | 0.8 | 0.2 |

| *p* | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? |
|-----|---|---|---|---|---|---|---|---|---|---|---|
| *x* | T | H | T | H | H | H | T | H | T | T | H |
| $s_{Fair, i}$ | 0.25 | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| $s_{Loaded, i}$ | 0.1 | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? |

$S_{Fair, I} = 0.5*0.5$      Viterbi fills in all the question marks

$S_{Loaded, I} = 0.5*0.2$

# Hidden Markov Model: Viterbi algorithm

$$s_{\text{Fair},i} = P(\text{Heads} \mid \text{Fair}) \; X \quad \max_{k \in \{ \text{Fair, Loaded} \}} \{ \, s_{k,i-1} \; x \; P(\text{Fair} \mid k) \, \}$$

Emission prob

Transition prob



Assignment Project Exam Help

$s_{\text{Fair, }i-1}$ x P(Fair | Fair)

https://powcoder.com

max

Add WeChat powcoder

$s_{\text{Loaded, }i-1}$ x P(Fair | Loaded)

P(Heads | Fair)

$s_{k,\,i}$ = score/probability of the most likely path up to step $i$ with $p_i = k$

# Hidden Markov Model: Viterbi algorithm

| | Step 1 | Step 2 | Step 3 | Step 4 ... |
|---|---|---|---|---|
| $x$ | **T** | **H** | **T** | **H** |

**$S_{F, i}$**

p(T|F) = 0.5

p(F) = 0.5

$S_{F, 1} = 0.5^2$

p(H|F) = 0.5
max:
p(F|F) = 0.6
$0.5^3 \cdot 0.6$
p(F|L) = 0.4
$0.2 \cdot 0.5 \cdot 0.5^2$

F wins max
$S_{F, 2} = 0.5^3 \cdot 0.6$

p(T|F) = 0.5
max:
p(F|F) = 0.6
$0.5^4 \cdot 0.6^2$
p(F|L) = 0.4
$0.4 \cdot 0.4 \cdot 0.8 \cdot 0.5$

F wins max
$S_{F,3} = 0.5^4 \cdot 0.6^2$

etc

**$S_{L, i}$**

p(T|L) = 0.2
p(L) = 0.5

$S_{L, 1} = 0.2 \cdot 0.5$

p(H|L) = 0.8
max:
p(L|F) = 0.4
$0.4 \cdot 0.5^2 \cdot 0.8$
p(L|L) = 0.6
$0.2 \cdot 0.5 \cdot 0.6 \cdot 0.8$

F wins max
$S_{L, 2} = 0.4 \cdot 0.5^2 \cdot 0.8$

etc

etc

# Hidden Markov Model: Viterbi algorithm

Pick state in step *n* with highest score; *backtrace* for most likely path

Backtrace according to which state *k* "won" the max in:

| Step n-3 | Step n-2 | Step n-1 | Step n |
|---|---|---|---|

$s_{F, n-3}$   $s_{F, n-2}$   $s_{F, n-1}$   $s_{F, n}$

$s_{L, n-3}$   $s_{L, n-2}$   $s_{L, n-1}$   $s_{L, n}$

$$s_{L,n} > s_{F,n}$$

# Hidden Markov Model: Viterbi algorithm

How much work did we do, given *Q* is the set of states and *n* is the length of the sequence?

$$\xleftarrow{\hspace{4cm}} n \xrightarrow{\hspace{4cm}}$$

*l Q l*

\# $S_{k,\,i}$ values to calculate = *n* · l *Q* l  , each involves max over l *Q* l products

   $O(n \cdot l Q l^2)$

Matrix *A* has *l Q l²* elements, *E* has *l Q ll ∑ l* elements, *I* has *l Q l* elements

# Hidden Markov Model: Viterbi algorithm

```
>>> hmm = HMM({"FF":0.6, "FL":0.4, "LF":0.4,
"LL":0.6},
...            {"FH":0.5, "FT":0.5, "LH":0.8,
"LT":0.2},
...            {"F":0.5, "L":0.5})
>>> prob, _ = hmm.viterbi("THTHHHTHTTH")
>>> print prob
2.86654464e-06
>>> prob, _ = hmm.viterbi("THTHHHTHTTH" * 100)
>>> print prob
0.0
```

Occasionally
dishonest
casino setup

Repeat string
100 times

What happened?  Underflow!

# Hidden Markov Model: Viterbi algorithm

```
>>> hmm = HMM({"FF":0.6, "FL":0.4, "LF":0.4, "LL":0.6},
...           {"FH":0.5, "FT":0.5, "LH":0.8, "LT":0.2},
...           {"F":0.5, "L":0.5})
>>> prob, _ = hmm.viterbi("THTHHHTHTTH")
>>> print prob
2.86654464e-06
>>> prob, _ = hmm.viterbi("THTHHHTHTTH" * 100)
>>> print prob
0.0
>>> logprob, _ = hmm.viterbiL("THTHHHTHTTH" * 100)
>>> print logprob
-1824.4030071946879
```

When multiplying many numbers in (0, 1], we quickly approach the smallest number representable in a machine word.  Past that we have *underflow* and processor rounds down to 0.

Switch to log space.  Multiplies become adds.

# Hidden Markov Model

Task: design an HMM for finding CpG islands?

# Hidden Markov Model

Idea 1: Q = { inside, outside }, ∑ = { A, C, G, T }

# Hidden Markov Model

Idea 1: Q = { inside, outside }, ∑ = { A, C, G, T }



inside                          outside

| A | I | O |
|---|---|---|
| I |   |   |
| O |   |   |

| E | A | C | G | T |
|---|---|---|---|---|
| I |   |   |   |   |
| O |   |   |   |   |

Estimate as
fraction of
positions where
we transition from
inside to outside

Transition matrix

Emission matrix

Estimate as
fraction of
nucleotides inside
islands that are C

# Hidden Markov Model

Example 1 using HMM idea 1:

| A | I | O |
|---|---|---|
| I | 0.8 | 0.2 |
| O | 0.2 | 0.8 |

| E | A | C | G | T |
|---|---|---|---|---|
| I | 0.1 | 0.4 | 0.4 | 0.1 |
| O | 0.25 | 0.25 | 0.25 | 0.25 |

x:

ATATATA**CGCGCGCGCGCG**ATATATATATATA
OOOOOOO**IIIIIIIIIIII**OOOOOOOOOOOOO

(from Viterbi)

# Hidden Markov Model

Example 2 using HMM idea 1:

| A | I | O |
|---|---|---|
| I | 0.8 | 0.2 |
| O | 0.2 | 0.8 |

| E | A | C | G | T |
|---|---|---|---|---|
| I | 0.1 | 0.4 | 0.4 | 0.1 |
| O | 0.25 | 0.25 | 0.25 | 0.25 |

x:
ATATCGCGCGCGATATATCGCGCGCGATATATAT
OOOOIIIIIIIIIOOOOOOIIIIIIIIIOOOOOOOO

(from Viterbi)

# Hidden Markov Model

Example 3 using HMM idea 1:

| A | I | O |
|---|---|---|
| I | 0.8 | 0.2 |
| O | 0.2 | 0.8 |

| E | A | C | G | T |
|---|---|---|---|---|
| I | 0.1 | 0.4 | 0.4 | 0.1 |
| O | 0.25 | 0.25 | 0.25 | 0.25 |

A: ATATATACCCCCCCCCCCCCCCCATATATATATATATA

O: OOOOOOOIIIIIIIIIIIIIIIIOOOOOOOOOOOOOOO

(from Viterbi)

Oops - not a CpG island!

# Hidden Markov Model

Idea 2: $Q = \{ A_i, C_i, G_i, T_i, A_o, C_o, G_o, T_o \}$, $\sum = \{ A, C, G, T \}$

# Hidden Markov Model

Idea 2: $Q = \{ A_i, C_i, G_i, T_i, A_o, C_o, G_o, T_o \}$, $\sum = \{ A, C, G, T \}$

All inside-outside edges

# Hidden Markov Model

Idea 2: $Q = \{ A_i, C_i, G_i, T_i, A_o, C_o, G_o, T_o \}$, $\sum = \{ A, C, G, T \}$

| **A** | $A_i$ | $C_i$ | $G_i$ | $T_i$ | $A_o$ | $C_o$ | $G_o$ | $T_o$ |
|---|---|---|---|---|---|---|---|---|
| $A_i$ | | | | | | | | |
| $C_i$ | | | | | | | | |
| $G_i$ | | | | | | | | |
| $T_i$ | | | | | | | | |
| $A_o$ | | | | | | | | |
| $C_o$ | | | | | | | | |
| $G_o$ | | | | | | | | |
| $T_o$ | | | | | | | | |

Estimate $P(C_i \mid T_i)$ as fraction of all dinucleotides where first is an inside T, second is an inside C

| **E** | A | C | G | T |
|---|---|---|---|---|
| $A_i$ | 1 | 0 | 0 | 0 |
| $C_i$ | 0 | 1 | 0 | 0 |
| $G_i$ | 0 | 0 | 1 | 0 |
| $T_i$ | 0 | 0 | 0 | 1 |
| $A_o$ | 1 | 0 | 0 | 0 |
| $C_o$ | 0 | 1 | 0 | 0 |
| $G_o$ | 0 | 0 | 1 | 0 |
| $T_o$ | 0 | 0 | 0 | 1 |

```
[[  1.85152516e-01   2.75974026e-01   4.00289017e-01
   1.37026750e-01
    3.19045117e-04   3.19045117e-04   6.38090233e-04
   2.81510397e-04]
 [  1.89303979e-01   3.58523577e-01   2.52868527e-01
   1.97836007e-01
    4.28792308e-04   5.72766368e-04   3.75584503e-05
   4.28792308e-04]
 [  1.72369088e-01   3.29501650e-01   3.55446538e-01
   1.40829292e-01
    3.39848138e-04   4.94038497e-04   7.64658311e-04
   2.54886104e-04]
 [  9.38783432e-02   3.46823149e-01   3.75970400e-01
   1.86949063e-01
    2.56686367e-04   5.57197233e-04   1.05804868e-03
   5.07112091e-04]
 [  0.00000000e+00   3.78291020e-05   0.00000000e+00
   0.00000000e+00
    2.94813496e-01   1.94641138e-01   2.86962055e-01
   2.23545482e-01]
 [  0.00000000e+00   7.57154865e-05   0.00000000e+00
   0.00000000e+00
    3.26811872e-01   2.94079570e-01   6.17258712e-02
   3.17306971e-01]
 [  0.00000000e+00   5.73810399e-05   0.00000000e+00
   0.00000000e+00
    2.57133507e-01   2.33483327e-01   2.94234944e-01
   2.15090841e-01]
 [  0.00000000e+00   3.11417347e-05   0.00000000e+00
```

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# Hidden Markov Model

Actual trained transition matrix A:

Red & orange: low probability

Yellow: high probability

White: probability = 0

|     | $A_i$ | $C_i$ | $G_i$ | $T_i$ | $A_o$ | $C_o$ | $G_o$ | $T_o$ |
|-----|-------|-------|-------|-------|-------|-------|-------|-------|
| $A_i$ |  |  |  |  |  |  |  |  |
| $C_i$ |  |  |  |  |  |  |  |  |
| $G_i$ |  |  |  |  |  |  |  |  |
| $T_i$ |  |  |  |  |  |  |  |  |
| $A_o$ |  |  |  |  |  |  |  |  |
| $C_o$ |  |  |  |  |  |  |  |  |
| $G_o$ |  |  |  |  |  |  |  |  |
| $T_o$ |  |  |  |  |  |  |  |  |

Once inside, we likely stay inside for a while

CpG islands end with G (in fact, they end with CG)

Same for outside

Ignore the small yellow blocks

CpG islands start with C (in fact, they start with CG)

# Hidden Markov Model

Viterbi result: lowercase = *outside*, uppercase = *inside*:

```
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
CGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
CGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
```

# Hidden Markov Model

Viterbi result: lowercase = *outside*; uppercase = *inside*.

atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatata
tatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatata
tatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatata
tatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatata
tatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatata
tatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat

atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatata
tatatatatatatatatatatatatatatatatatatatatatatatgggggggggggggggggggggggggggggggggggggggg
gggggggggggggggggggggggggggggggggggggggggggggggggggggggggggggggggggggggggggggggggggggggg
ggggggggggggggggggggggggggggggggggggggggggggggggggggggggggggggatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatata
tatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatata
tatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatata
tatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatata
tatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatata
tatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatata
tatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatata
tatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatat
atatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatatata
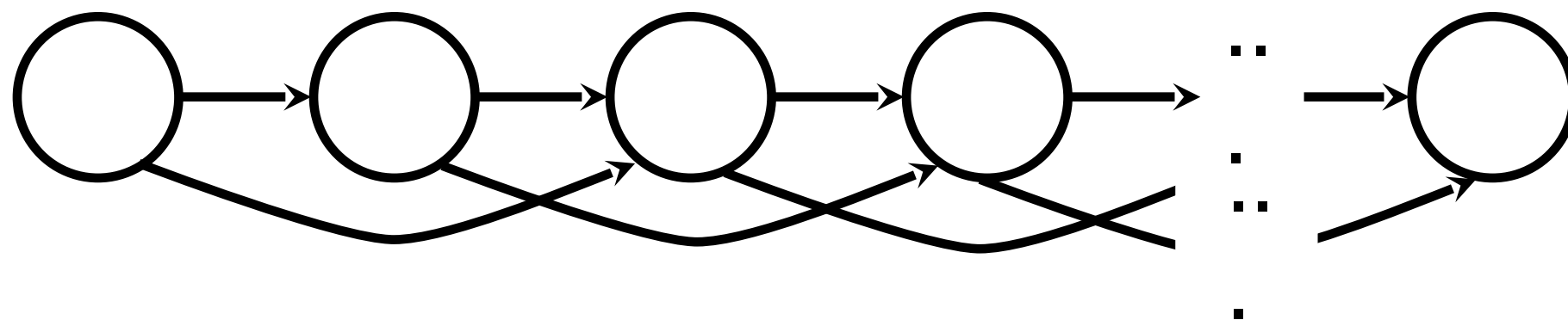
# Hidden Markov Model

Many of the Markov chains and HMMs we've discussed are *first order,* but we can also design models of higher orders
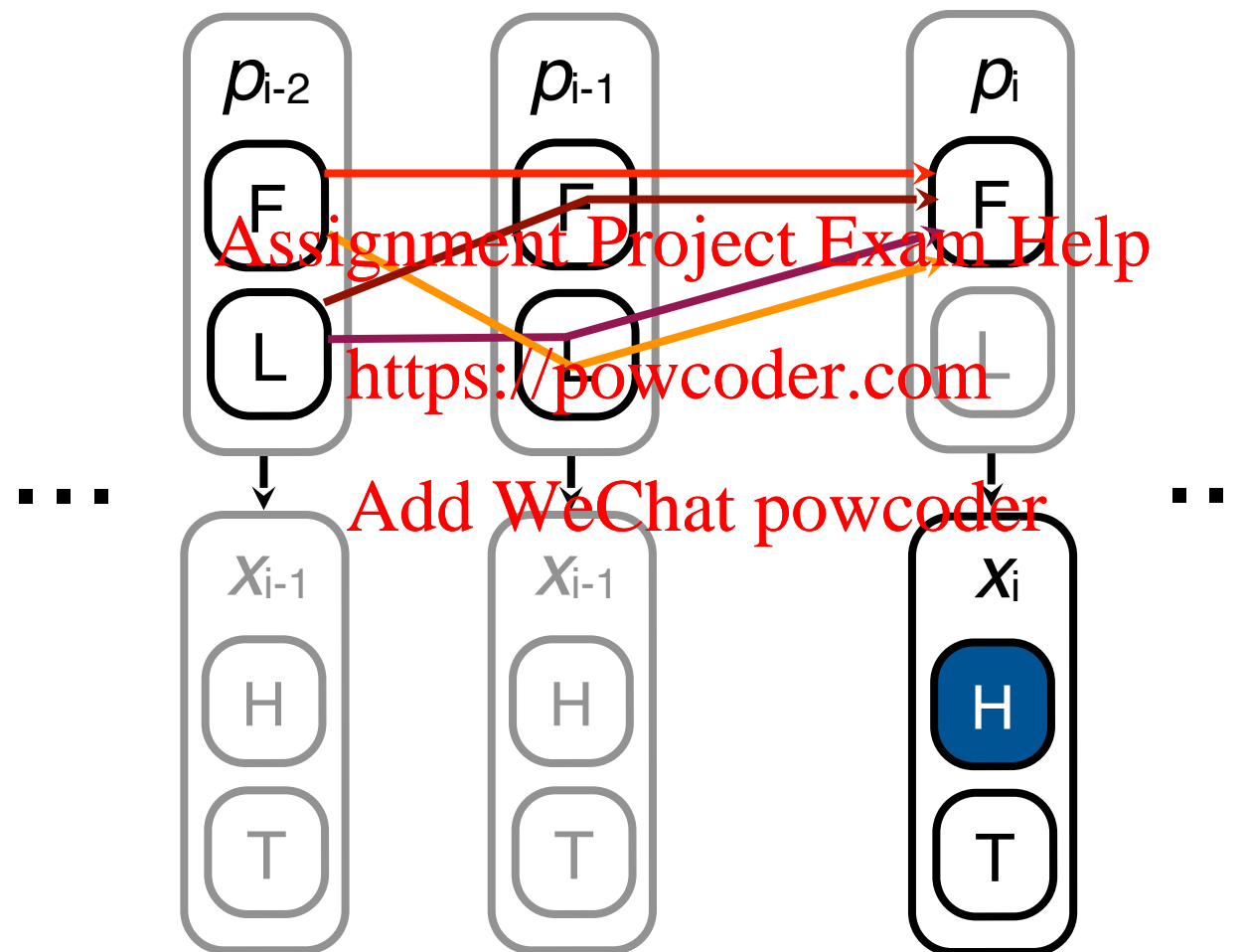
First-order Markov chain:

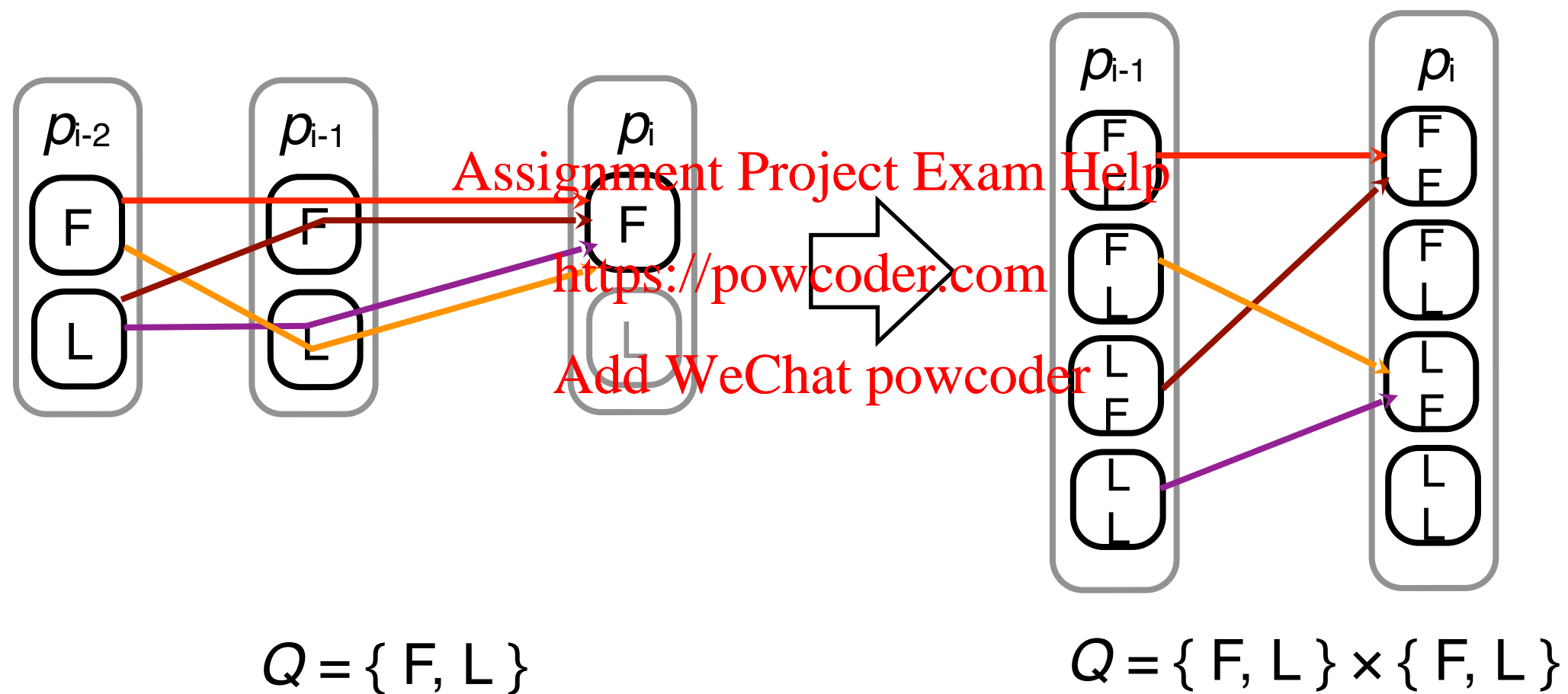Second-order Markov chain:

# Hidden Markov Model

For higher-order HMMs, Viterbi $S_{k,\,i}$ no longer depends on just the previous state assignment

Can sidestep the issue by expanding the state space...

# Hidden Markov Model

Now *one* state encodes the last *two* "loadedness"es of the coin

$Q = \{ F, L \}$

$Q = \{ F, L \} \times \{ F, L \}$

After expanding, usual Viterbi works fine.