

Markov Chains and Hidden Markov Models

= stochastic, generative models

(Drawing heavily from Durbin *et al.*, *Biological Sequence Analysis*)

**Systems Biology / Bioinformatics
Edward Marcotte, Univ of Texas at Austin**

**Markov Chains and Hidden Markov Models are important
probabilistic models in computational biology**

Some of their applications include:

- Finding genes in genomes
- Mapping introns, exons, and splice sites
- Identifying protein domains & families
- Detecting distant sequence homology
- Identifying secondary structures in proteins
- Identifying transmembrane segments in proteins
- Aligning sequences

& outside biology, they have many uses, including:

- Speech, handwriting, and gesture recognition
 - Tagging parts-of-speech
 - Language translation
 - Cryptanalysis
- and so on....

The key idea of both of these types of models is that:

Biological sequences can be modeled as series of stochastic (i.e., random) events.

It's easy to see how a random process might model stretches of DNA between genes and other important regions.

BUT, the idea of modeling something as structured and meaningful as a gene or protein sequence by a similar process might seem odd.

It's important to realize exactly what we're modeling.

The idea behind hidden Markov models is not that the sequence is random, but that the sequence we observe is one of many possible instances of some underlying process or object.

E.g., actin differs slightly from organism to organism.

Imagine an "ideal", but unobservable, actin, defined by specific underlying physico-chemical properties important for its function. What we see in nature is not this ideal gene, but many variants, all just a bit different.

In the hidden Markov model, the underlying process or structure is represented as hidden, unobservable **states** and the observed sequences represent possible sequences compatible with these states.

We say that the observed sequence is **emitted** from the hidden states.

Let's start with a easier case: **Markov chains**

We'll explore a simple non-biological model: a coin-toss

Flip a coin a bunch of times and observe the results, e.g.

TTTTTTHTTTHTTTHTHTHHHTTTTTTTTTHTTHTTTHTHHHTHH

We could model this process as two states:

H for heads,

T for tails,

and the probability of switching between them:



Let's start with a easier case: **Markov chains**

We'll explore a simple non-biological model: a coin-toss

Flip a coin a bunch of times and observe the results, e.g.

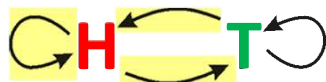
TTTTTTHTTTHTTTHTHTHHHTTTTTTTTTHTTHTTTHTHHHTHH

We could model this process as two states:

H for heads,

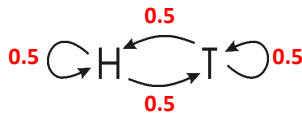
T for tails,

and the probability of switching between them:



A sequence is a walk along this graph:

H H T H ...



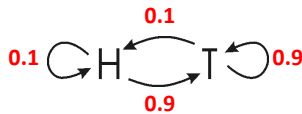
Important: All probabilities leading out of a state add up to 1!

With a **fair coin**:

The chance of seeing heads or tails is equal, and the chance of seeing heads following tails and vice versa is equal.

Therefore, the **transition probabilities** (corresponding to the arrows above) are:

Position i:	Position i+1:	
	Head	Tail
Head	0.5	0.5
Tail	0.5	0.5



Important: All probabilities leading out of a state add up to 1!

With a **biased coin** (e.g. tails comes up 90% of the time):

The chance of seeing heads or tails is not equal, nor is the chance of seeing heads following tails and vice versa.

We might have the same model, but with skewed **transition probabilities** :

Position i:	Position i+1:	
	Head	Tail
Head	0.1	0.9
Tail	0.1	0.9

How about a biological application? A classic example is **CpG islands**

In animal genomes, the dinucleotide CG is strongly underrepresented (note: NOT the base pair C:G, but rather 5'-CG-3')

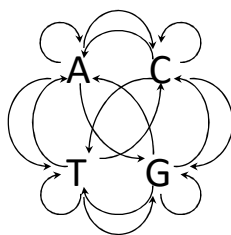
Why? C's are often methylated, and methylated C's mutate at higher rates into T's. So, over time, CG's convert to TG's EXCEPT around promoters, which tend not to be methylated.

Thus, **CpG 'islands' often indicate nearby genes**. Finding them was a classic method for annotating genes.

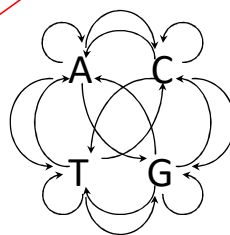
How could we make a CpG island finding model analogous to the fair/biased coin model?

A CpG island model might look like:

(of course, need the parameters, but maybe these are the most important....)



CpG island model



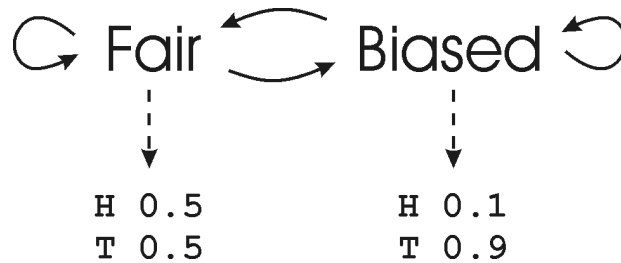
Not CpG island model

Could calculate $\frac{P(X \mid \text{CpG island})}{P(X \mid \text{not CpG island})}$

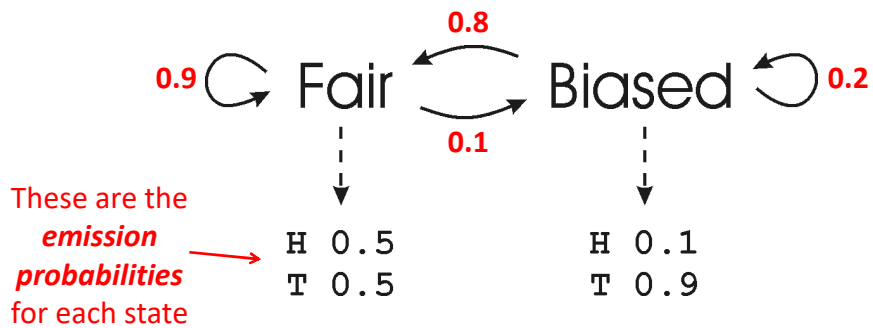
(or log ratio) along a sliding window, just like the fair/biased coin test

In these simple models, called *Markov chains*, we don't have hidden states.

BUT, we could have used a *hidden Markov model*:



Now, the underlying *state* (the choice of coin) is hidden. Each state *emits* H or T with different probabilities.



The *transition probabilities* might be something like:

Position i:	Position i+1:	
	Fair	Biased
Fair	0.9	0.1
Biased	0.8	0.2

Important questions we might like to ask:

- 1. Given an observed sequence and a model, what is the most likely sequence of hidden states?**

i.e., what is the path through the HMM that maximizes $P(p, X | I)$, where p is the sequence of states)?

In our coin example, we might be given an observed sequence:

HTHTHTHTTTTTTTTTTTTTTTTTTTTTHTHTHTHTHT

and want to identify when the biased coin was used:

FFFFFFFFFFFFBBBBBBBBBBBBBBBBBBBBFFFFFFFF

**Answer: Use the Viterbi algorithm.
We'll see this shortly.**

Important questions we might like to ask:

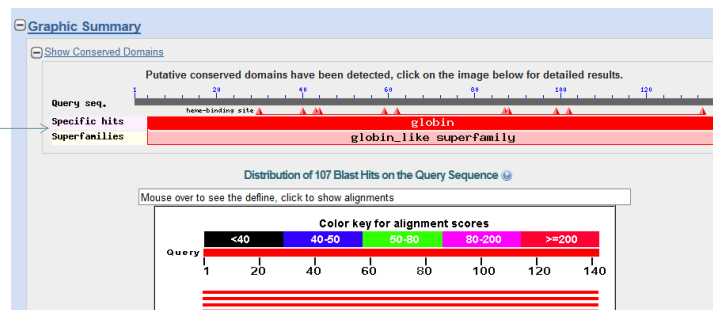
- 2. Given a *sequence of observations*, can we calculate the probability that the sequence was derived from our model ?**

i.e., can we calculate $P(X | I)$,

where X is our observed sequence, and I represents our HMM ?

For example, we might want to know if a given protein sequence is a member of a certain protein family.

e.g. as we saw before (although calculated a bit differently)



Important questions we might like to ask:

2. Given a *sequence of observations*, can we calculate the probability that the sequence was derived from our model ?

i.e., can we calculate $P(X|I)$,

where X is our observed sequence, and I represents our HMM ?

For example, we might want to know if a given protein sequence is a member of a certain protein family.

**Answer: Yes. Use the *forward algorithm*.
We'll see this shortly.**

Important questions we might like to ask:

3. Given a model, what is the most likely sequence of observations?

For example, after having trained an HMM to recognize a type of protein domain, what amino acid sequence best embodies that domain?

Answer: Follow the maximum transition and emission probability at each state in the model. This will give the most likely state sequence and observed sequence.

Important questions we might like to ask:

4. How do we train our HMM?

i.e., given some training observations, how do we set the emission and transition probabilities to maximize $P(X|I)$?

Answer: If the state sequence is known for your training set, just directly calculate the transition and emission frequencies. With sufficient data, these can be used as the probabilities.

This is what you will do in Problem Set #2.

With insufficient data, probabilities can be estimated from these (e.g., by adding pseudo-counts).

If the state path is unknown, use the *forward-backward algorithm* (also known as the *Baum-Welch algorithm*).

Important questions we might like to ask:

5. How do we choose the best HMM topology from the many possible choices?

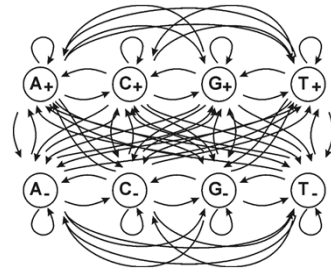
Answer: Good question. No great answer.

Often trial-and-error, and understanding the essential features of the system that you are modeling.

Each of these algorithms (the Viterbi, forward, and forward-backward) uses dynamic programming to find an optimal solution.

(just like aligning sequences)

Let's revisit the CpG islands using an HMM:



- 8 states: one per nucleotide inside CpG islands (+) and one per nucleotide outside CpG islands (-)
- All possible transition probabilities are represented as arrows
- This is a particularly simple model: each state emits the nucleotide indicated with probability of 1 and has zero probability of emitting a different nucleotide.

Just to be totally clear:

In this (greatly simplified) model, here is our table of emission probabilities:

		<u>Emission probabilities</u>				
		<u>States</u>	A	C	T	G
+ means a base inside a CpG island	←	A+	1	0	0	0
		A-	1	0	0	0
		C+	0	1	0	0
		C-	0	1	0	0
- means a base outside a CpG island	←	T+	0	0	1	0
		T-	0	0	1	0
		G+	0	0	0	1
		G-	0	0	0	1

i.e. an "A" can only be emitted from an A+ or A- state, not any other state

Emission probabilities from the same state must add up to 1

This is a particularly simple model: each state emits the nucleotide indicated with probability of 1 and has zero probability of emitting a different nucleotide.

**Given a DNA sequence X (e.g., CGATCGCG),
how do we find the most probable sequence of states
(e.g., ----++++)?**

→ ***The Viterbi algorithm***

We want to find the state path that maximizes the probability of observing that sequence from that HMM model.

Viterbi does this recursively using dynamic programming.

As with sequence alignment, we'll construct a path matrix that captures the best score (i.e., highest probability) along a single path through the HMM up to each position. We'll "grow" this matrix using a few simple recursion rule.

The rules (stated formally):

Initialization ($i=0$): $v_0(0) = 1, v_k(0)=0$ for $k>0$

Recursion ($i=1$ to L):
 $v_i(i) = e_i(x_i) \max_k (v_k(i-1) a_{ki})$
 $pointer_i(i) = \operatorname{argmax}_k (v_k(i-1) a_{ki})$

Termination:
 $P(X, p^*) = \max_k (v_k(L) a_{k0})$
 $p_L^* = \operatorname{argmax}_k (v_k(L) a_{k0})$

x indicates an observed character
e indicates an emission probability
i indicates our position in the sequence
v is an entry in the Viterbi path matrix
Find the best score among the alternatives at this position
a gives the transition probability between previous state k and current state l
i.e., draw the pointer back to the entry that gave rise to the current best score

For each Viterbi matrix entry:
We try to maximize the product of prior score and transition from that state to this one.
We then multiply that score times the emission probability for the current character.

Step 1: Initialize the path matrix.

Possible states

\mathcal{B}	1	0	0	0	0
A+	0				
C+	0				
G+	0				
T+	0				
A-	0				
C-	0				
G-	0				
T-	0				

Observed DNA sequence: C G C G

For simplicity, let's assume the transition probability from \mathcal{B} to each nucleotide is $1/8$. We'll also ignore all transition probabilities except these for now:

Position i:	Position i+1:			
	C+	G+	C-	G-
C+	0.37	0.27	small	small
G+	0.34	0.38	small	small
C-	small	small	0.3	0.08
G-	small	small	0.25	0.3

Step 2: Calculate the elements of the v_k matrix for $i = 1$.
Then keep going for $i = 2$, etc..

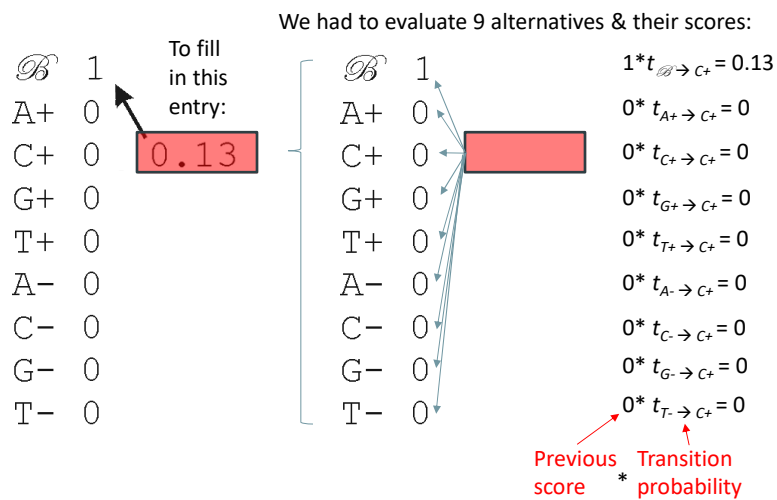
		C	G	C	G
\mathcal{B}	1	0	0	0	0
A+	0	0			
C+	0	0.13			
G+	0	0			
T+	0	0			
A-	0	0			
C-	0	0.13			
G-	0	0			
T-	0	0			

For simplicity, let's assume the transition probability from \mathcal{B} to each nucleotide is $1/8$. We'll also ignore all transition probabilities except these for now:

Position i:	Position i+1:		C-	G-
C+	0.37	0.27	small	small
G+	0.34	0.38	small	small
C-	small	small	0.3	0.08
G-	small	small	0.25	0.3

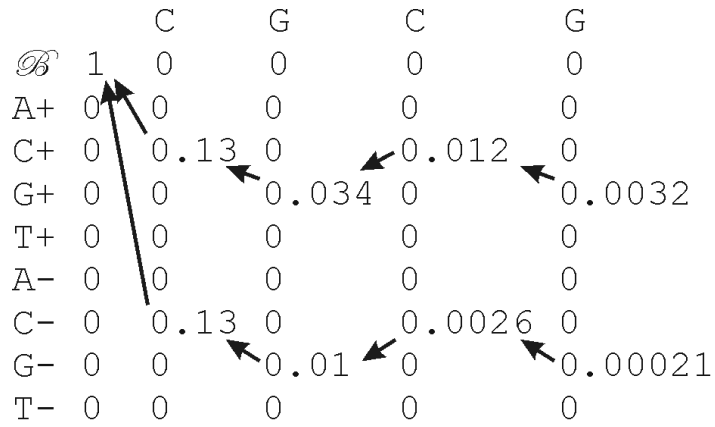
For example, the score $v_{C^+}(i=1) = 1 * \max_k \{1 * 1/8, 0 * a_{A^+,C^+}, 0 * a_{C^+,C^+}, \dots, 0 * a_{T^-,C^+}\} = 1/8$

Let's blow up just one square of the matrix to understand a bit better:



Take the *max* (in this case, 0.13) and only then multiply that by the emission probability (here, 1) to give a final score for that cell of 0.13. Draw an arrow back to the cell it came from (here, \mathcal{B})

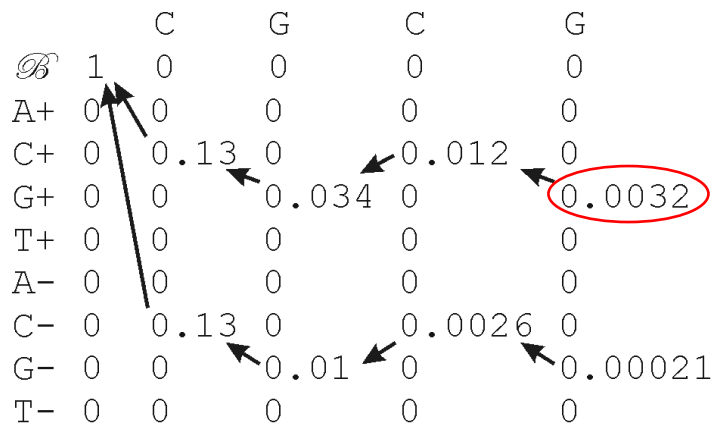
Step 3: Keep going for $i = 2$, etc..



Position i:	Position i+1:		C-	G-
	C+	G+		
C+	0.37	0.27	small	small
G+	0.34	0.38	small	small
C-	small	small	0.3	0.08
G-	small	small	0.25	0.3

The maximum scoring path scores 0.0032.

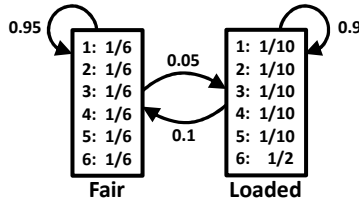
The most likely state path is found by traceback from the 0.0032 to give C+G+C+G+.



In a longer sequence, the model would switch back & forth between CpG and non-CpG states appropriately.

Can this really work? Here's a real example.

An HMM model of fair and loaded dice:



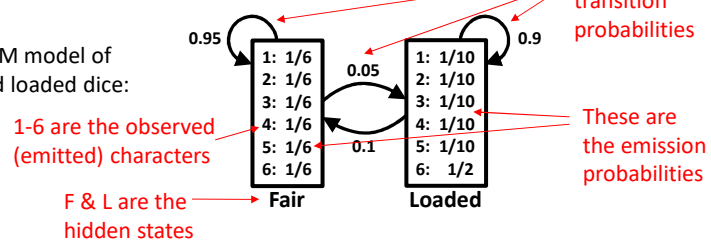
Reconstructing which was used when, using the Viterbi algorithm:

Rolls	315116246446644245311321631164152133625144543631656626566666
Die	FFL
Viterbi	FFL
Rolls	651166453132651245636664631636663162326455236266666625151631
Die	LLLLLFF
Viterbi	LLLLLFF
Rolls	222555441666566563564324364131513465146353411126414626253356
Die	FFFFFFFFLLL
Viterbi	FFL
Rolls	36616366646623253441366166116325256246225526525226643535336
Die	LLLLLLLLFF
Viterbi	LLLLLLLLLLLLLFF
Rolls	233121625364414432335163243633665562466662632666612355245242
Die	FF
Viterbi	FF

from Durbin et al.

Can this really work? Here's a real example.

An HMM model of fair and loaded dice:



Reconstructing which was used when, using the Viterbi algorithm:

Rolls	315116246446644245311321631164152133625144543631656626566666
Die	FFL
Viterbi	FFL
Rolls	651166453132651245636664631636663162326455236266666625151631
Die	LLLLLFF
Viterbi	LLLLLFF
Rolls	222555441666566563564324364131513465146353411126414626253356
Die	FFFFFFFFLLL
Viterbi	FFL
Rolls	36616366646623253441366166116325256246225526525226643535336
Die	LLLLLLLLFF
Viterbi	LLLLLLLLLLLLLFF
Rolls	233121625364414432335163243633665562466662632666612355245242
Die	FF
Viterbi	FF

from Durbin et al.

How do we calculate the probability of a sequence given our HMM model?

→ *The forward algorithm*

Subtle difference from Viterbi:

Viterbi gives the probability of the sequence being derived from the model *given the optimal state path*.

The forward algorithm takes into account all possible state paths.

Again, it does this recursively using dynamic programming.

The rules (stated formally):

Initialization ($i=0$):

f is an entry in the forward algorithm path matrix

$$f_0(0) = 1, f_k(0) = 0 \text{ for } k > 0$$

Recursion ($i=1$ to L):

Same idea as Viterbi, but ADD the scores leading to the current position (not MAX)

$$f_l(i) = e_l(x_i) \sum_k f_k(i-1) a_{kl}$$

Termination:

$$P(x) = \sum_k f_k(L) a_{k0}$$

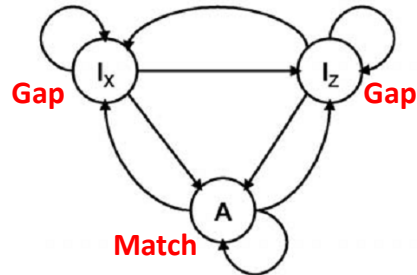
Note: No pointer! Just to calculate the probability of seeing this sequence from this model.

For each Viterbi matrix entry:

We try to maximize the product of prior score and transition from that state to this one.

We then multiply that score times the emission probability for the current character.

A toy HMM for sequence alignment



I_x : insertion in x (seq 1)
 I_z : insertion in z (seq 2)
 A: aligned symbols in x and z

x (seq 1) : T T C C G - -
 z (seq 2) : - - C C G T T

 y (states) : $I_x I_x A A A I_z I_z$

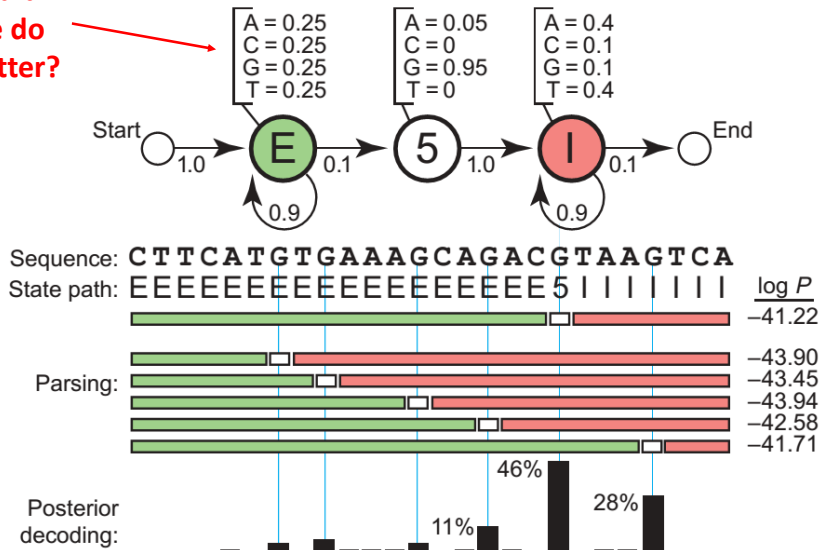
Is this global or local alignment?
How could you change the model to perform the other kind of alignment?

Curr Genomics (2009) 10(6): 402–415

A toy HMM for 5' splice site recognition (from Sean Eddy's NBT primer

linked on the course web page)

Could we do better?



How might you design an HMM to recognize a given type of protein domain?

How might we design HMMs to recognize sequences of a given length?

What would this HMM produce?

