

**Illumina claims \$1,000 genome win**

Illumina, of San Diego, California, has laid claim to a coveted prize in DNA sequencing technology: the \$1,000 genome. The economics of the HiSeq X Ten system, which the company unveiled on 14 January and expected to begin shipping in March, only work for population-scale, whole-genome sequencing initiatives, however. Access to the technology requires an initial outlay of \$10 million, as Illumina is only accepting orders for a minimum of ten systems. Its \$1,000-per-genome calculation is based on the use of ten systems over four years, which represents the equivalent of over 72,000 genomes—or about \$82 million in total capital and running costs. The new system offers a tenfold performance improvement over the company's current high-end HiSeq 2500 systems, which can sequence 600 gigabases (Gb) in ten days. A single run, which costs \$12,750, yields 16 Gb of data. An array of ten HiSeq X sequencers can sequence 1.8 terabases (1,800 Gb) in three days.

The performance senior manager in product development due to an increase in throughput, he says—and to a faster run time, which he attributes to a redesigned optical system and to faster oligonucleotide incorporation into Illumina's core sequencing-by-synthesis process.

*Cormac Sheridan*



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

NATURE BIOTECHNOLOGY VOLUME 32 NUMBER 2 FEBRUARY 2014

# Markov Chains and Hidden Markov Models

## = stochastic, generative models

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

**BIO337 Systems Biology / Bioinformatics – Spring 2014**  
**Edward Marcotte, Univ of Texas at Austin**

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

**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 domain 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....

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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.***

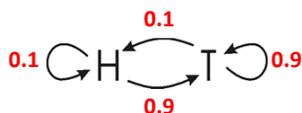
It's easy to see how a random process might model stretches of DNA between genes and other important regions.

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014







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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Now, imagine a scenario where the observed sequence of coin flips was actually generated by 2 coins, one fair and one biased.

To decide whether we are looking at a sequence of coin flips from the biased or fair coin, we could evaluate the ratio of the probabilities of observing the sequence by each model:

$$\frac{P(X \mid \text{fair coin})}{P(X \mid \text{biased coin})}$$

Does this remind you of something we've seen before?  
How might we test where the fair & biased coins were swapped along a long stretch of coin flips?

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

How might we test where the fair & biased coins were swapped along a long stretch of coin flips?

One way using our current Markov chain model is to calculate the ratio of probabilities (e.g. log odds ratio) in a **sliding window** along the sequence:



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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?

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

**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  $\frac{P(X | \text{CpG island})}{P(X | \text{not CpG island})}$  (or log ratio) along a sliding window, just like the fair/biased coin test

Edward Marcotte/Univ. of Texas/BI0337/Spring 2014

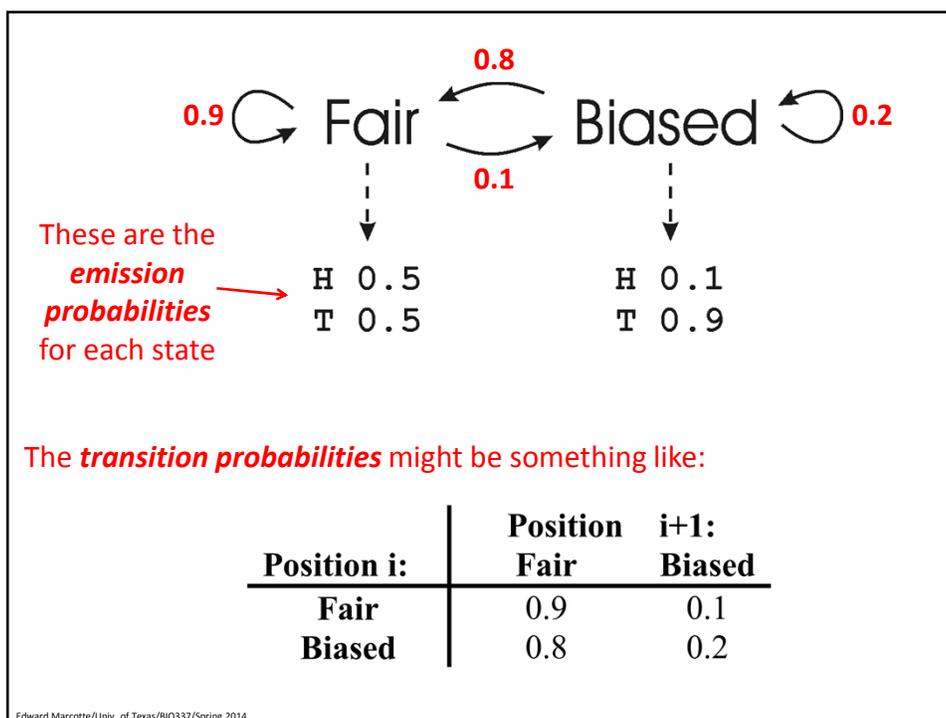
In these simple models, called **Markov chains**, we don't have hidden states.

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

↻	Fair	↔	Biased	↻
	⋮		⋮	
	↓		↓	
	H 0.5		H 0.1	
	T 0.5		T 0.9	

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

Edward Marcotte/Univ. of Texas/BI0337/Spring 2014



### 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

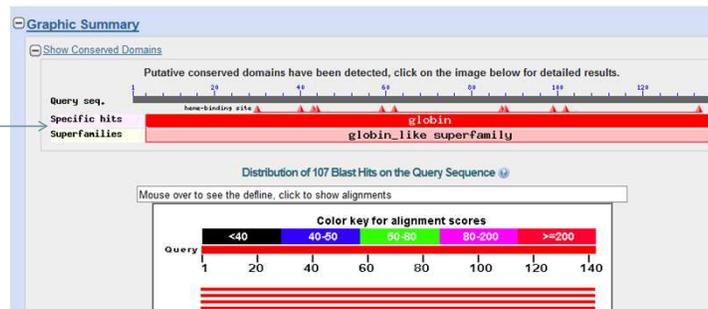
**Answer: Use the Viterbi algorithm.  
We'll see this shortly.**

Important questions we might like to ask:

- 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)



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Important questions we might like to ask:

- 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.**

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

**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.**

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

**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*).

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

**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.**

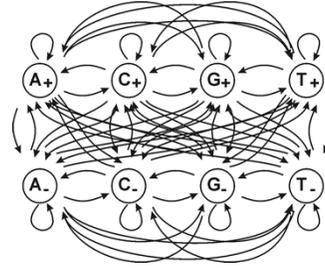
Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

**(just like aligning sequences)**

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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.

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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.

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

**The rules (stated formally):**

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

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

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

**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.

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

Edward Marcotte/Univ. of Texas/BI0337/Spring 2014

Step 1: Initialize the path matrix.

**Observed DNA sequence**  
 C G C G

**Possible states**

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

Edward Marcotte/Univ. of Texas/BI0337/Spring 2014

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

Edward Marcotte/Univ. of Texas/BI0337/Spring 2014

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

		C	G	C	G
$\mathcal{B}$	1	0	0	0	0
A+	0	0	0	0	0
C+	0	0.13	0	0.012	0
G+	0	0	0.034	0	0.0032
T+	0	0	0	0	0
A-	0	0	0	0	0
C-	0	0.13	0	0.0026	0
G-	0	0	0.01	0	0.00021
T-	0	0	0	0	0

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

Edward Marcotte/Univ. of Texas/BI0337/Spring 2014

**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+.**

$B$	1	0	0	0	0
A+	0	0	0	0	0
C+	0	0.13	0	0.012	0
G+	0	0	0.034	0	0.0032
T+	0	0	0	0	0
A-	0	0	0	0	0
C-	0	0.13	0	0.0026	0
G-	0	0	0.01	0	0.00021
T-	0	0	0	0	0

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

An HMM model of fair and loaded dice:

from Durbin et al.

Reconstructing which was used when, using the Viterbi algorithm:

```

Rolls 315116246446644245311321631164152133625144543631656626566666
Die   FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Viterbi FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF

Rolls 65116645313265124563666463163666316232645523626666625151631
Die   LLLLLLFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Viterbi LLLLLLFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF

Rolls 222555441666566563564324364131513465146353411126414626253356
Die   FFFFFFFFFLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL
Viterbi FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF

Rolls 366163666466232534413661661163252562462255265252266435353336
Die   LLLLLLLLLFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Viterbi LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLFFFF

Rolls 233121625364414432335163243633665562466662632666612355245242
Die   FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Viterbi FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
    
```

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

## 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.**

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

### The rules (stated formally):

Initialization ( $i=0$ ):

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

*f* is an entry in the forward algorithm path matrix

Recursion ( $i=1$  to  $L$ ):

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

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

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.

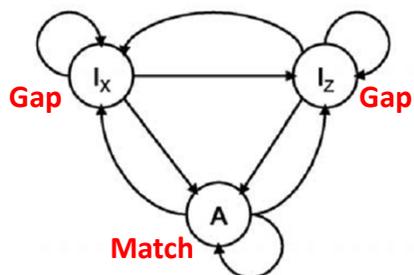
**For each Viterbi matrix entry:**

**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.**

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

### 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?

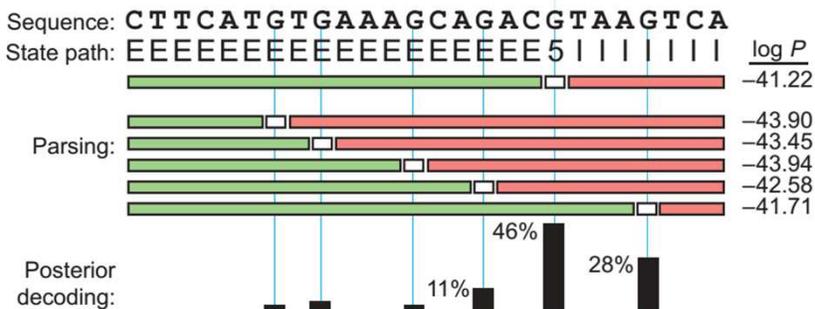
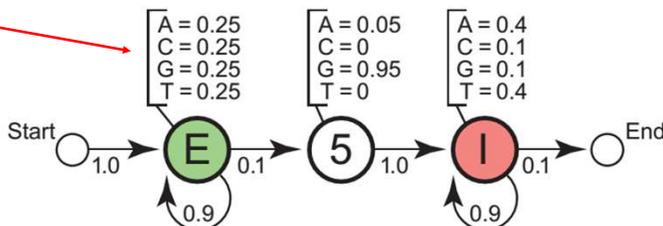
How could you change the model to perform the other kind of alignment?

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014cc

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?



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014cc

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014cc

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014cc

## What would this HMM produce?

