

Algorithms in Bioinformatics

Hidden Markov
Models



Recap

- HW4
- HW2 is graded
- HW5 + essay are up

Chapter 11: Hidden Markov Models

Least frequent dinucleotide: CG

- C is easily methylated, +
then a tendency to mutate
to a T.

However, this is suppressed around
genes in CG Islands.

An important problem is to
identify these islands.

Model this probabilistically:

Consider sliding windows
of some length, +
calculate some probability
of it being an island

First, we consider a related
"toy" problem

"Fair Bet Casino":

A crooked dealer is flipping a coin.
- either fair, or biased one that outputs H with probability $\frac{3}{4}$

The dealer needs to be sneaky, so only changes coins w/ probability $\frac{1}{10}$.

Goal: Try to figure out if he's using a fair or biased coin:

Fair Bet Casino Problem:

Given a sequence of coin tosses, determine when the dealer used a fair coin and when he used a biased coin.

Input: A sequence $x = x_1 x_2 x_3 \dots x_n$ of coin tosses (either H or T) made by two possible coins (F or B).

Output: A sequence $\pi = \pi_1 \pi_2 \pi_3 \dots \pi_n$, with each π_i being either F or B indicating that x_i is the result of tossing the fair or biased coin, respectively.

Unfortunately, doesn't quite make sense...

Consider if the dealer never changes coins first.

Fair coin: $P(H) = \frac{1}{2}$, $P(T) = \frac{1}{2}$
Biased coin: $P(H) = \frac{3}{4}$, $P(T) = \frac{1}{4}$

Let $x = x_1 x_2 x_3 \dots x_n$

be the sequence of flips

$$P(x \mid \text{fair coin}) =$$

$$\left(\frac{1}{2}\right)\left(\frac{1}{2}\right) \dots \left(\frac{1}{2}\right) = \left(\frac{1}{2}\right)^n \\ = \frac{1}{2^n}$$

If k heads are in x ,

$$P(x \mid \text{biased coin}) =$$

$$\left(\frac{1}{4}\right)^{n-k} \left(\frac{3}{4}\right)^k = \frac{3^k}{4^k} \cdot \frac{1}{4^{n-k}} = \frac{3^k}{4^n}$$

Idea: if $P(x \mid \text{fair}) > P(x \mid \text{biased})$,
then it probably
was a fair coin

If we solve for k :

$$4^n \cdot \frac{1}{2^n} \geq \frac{3^k}{4^n} \cdot 4^n$$

$$\frac{2^{2n}}{2^n} \geq 3^k \rightarrow 2^n \geq 3^k$$

$$\log_2 2^n \geq \log_2 3^k$$

$$n \cancel{\log_2 2} \geq k \log_2 3$$

So: if $k < \frac{n}{\log_2 3}$, then most likely a fair coin.

However, the dealer changes coins!

Log-odds ratio:

$$\begin{aligned} & \log_2 \frac{P(x \mid \text{fair coin})}{P(x \mid \text{biased coin})} \\ &= \sum_{i=1}^k \log_2 \frac{p^+(x_i)}{p^-(x_i)} \\ &= n - k \log_2 3 \end{aligned}$$

(where $p^+(x_i)$ = prob of unbiased
 $p^-(x_i)$ = prob of biased)

- So:
- slide a window of some length along H/T sequence
 - if it falls below 0, probably biased
 - if above 0, probably fair

Hidden Markov Models:

An abstract machine w/ ability to produce output using coin tossing.

Our machine is in one of K "hidden" states.

Each step, it decides

- What state will I move to next?
- What symbol from Σ (an alphabet) will I output?

Essentially, will decide randomly, but using a biased distribution

(Observer can see output, not states.)

Goal: Infer States

Formal defn: M (a HMM)
is described by

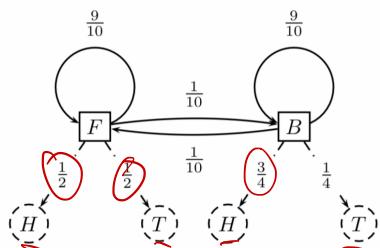
- Σ , an alphabet earlier H or T
- Q , a set of states, each of which outputs symbols (w/ some probability)
- a $|Q| \times |Q|$ matrix A
 $A[k][l] = \text{probability of going to state } l \text{ from state } k$
- E: a $|Q| \times |\Sigma|$ matrix
 $E[k][b] = \text{probability of outputting } b \text{ from state } k$

Ex: Fair bet casino HMM

$$\Sigma = \{H, T\}$$

$$Q = \{F, B\}$$

$$\left\{ \begin{array}{l} A[F][F] \\ = A[B][B] \\ A[F][B] \\ \{ = A[B][F] \end{array} \right. = \frac{9}{10}$$



$$E[F][H] = \frac{1}{2}$$

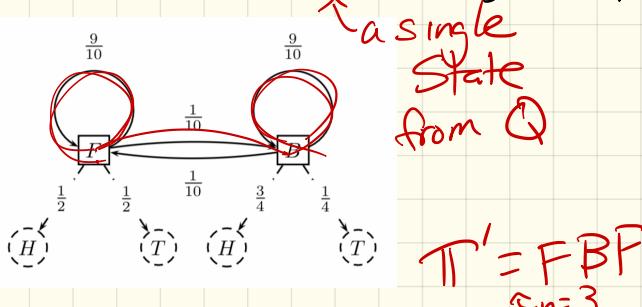
$$E[B][H] = \frac{3}{4}$$

$$E[F][T] = \frac{1}{2}$$

$$E[B][T] = \frac{1}{4}$$

A path in a HMM is just a sequence of states,

$$\pi = \pi_1 \pi_2 \pi_3 \dots \pi_n$$



i.e. $\pi = \underline{FFF} \underline{BBB} \underline{BBB} FFF$

$\hat{\cup}_{n=10}$
here

Suppose output was:
 $x = THTHHHHHHTHTT$

Can line up probabilities:

$$P(x_i | \pi_i) = \begin{pmatrix} T & H & T & H & H & H & T & H & T & T & H \\ 0 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 0 & 0 & 1 \\ F & F & F & B & B & B & B & B & F & F & F \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{3}{4} & \frac{3}{4} & \frac{3}{4} & \frac{1}{4} & \frac{3}{4} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \end{pmatrix}$$

\uparrow probability x_i was the output of state π_i
 (These were in E.)

Can also look at $P(\pi_i \rightarrow \pi_{i+1})$:
 the probability of going
 from state π_i to π_{i+1}
 (These were in matrix A).

Get larger view:

$$P(x_i | \pi_i) = \begin{pmatrix} 0 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 0 & 0 & 1 \\ F & F & F & B & B & B & B & B & F & F & F \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{3}{4} & \frac{3}{4} & \frac{3}{4} & \frac{1}{4} & \frac{3}{4} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{9}{10} & \frac{9}{10} \end{pmatrix}$$

↑ Pick 1st state randomly

Then the probability of generating
 x through path π is:

$$\left(\frac{1}{2} \cdot \frac{1}{2}\right) \left(\frac{1}{2}\right) \left(\frac{9}{10}\right) \dots$$

$$\approx 2.66 \times 10^{-6}$$

(Note: this assumed we knew π
and observed x .)

If we only know:

$$x = \text{THTHHHHTHTTHT} \quad (\text{not } \pi)$$

- how good is a particular π ?
- can we construct the most likely π ?

For above x , FFFBBBF~~FFFFF~~ is actually better:

$$\begin{array}{ll} x & \pi \\ \pi & P(x_i | \pi_i) \\ P(\pi_{i-1} \rightarrow \pi_i) & = \left(\begin{array}{cccccccccc} 0 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 0 & 0 & 1 \\ F & F & F & B & B & B & F & F & F & F & F \\ \frac{1}{2} & \frac{9}{10} & \frac{9}{10} & \frac{1}{2} & \frac{3}{4} & \frac{3}{4} & \frac{3}{4} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{9}{10} & \frac{9}{10} & \frac{1}{2} & \frac{1}{10} & \frac{9}{10} & \frac{9}{10} & \frac{1}{10} & \frac{9}{10} & \frac{9}{10} & \frac{9}{10} \end{array} \right) \end{array}$$

prob = $\left(\frac{1}{2} \times \frac{1}{2}\right) \left(\frac{1}{2} \cdot \frac{9}{10}\right) \dots$

$$\approx 3.54 \times 10^{-6}$$

(was 2, ...)

π is called "hidden".

Goal:

Decoding Problem:

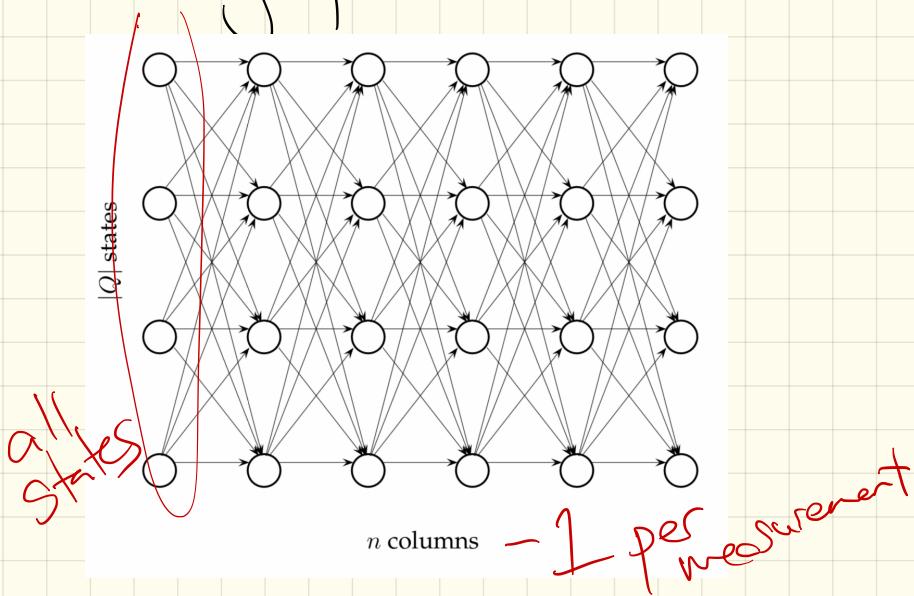
Find an optimal hidden path of states given observations.

Input: Sequence of observations $x = x_1 \dots x_n$ generated by an HMM $M(\sum, Q, A, E)$.

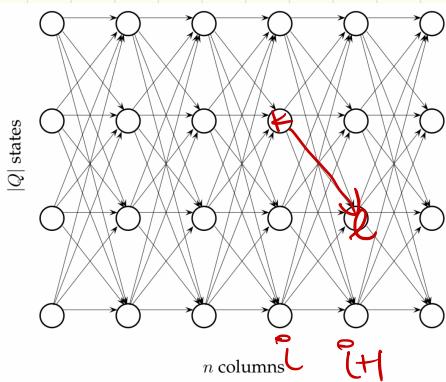
Output: A path that maximizes $P(x|\pi)$ over all possible paths π .

[Viterbi 1967] gave an approach
to solve

Build a graph:



Idea (in more detail):



Goal: Have weight on edges, so product of the weights for $\pi_1, \pi_2, \dots, \pi_n$
 $= P(x | \pi)$

How?

- Set edge from (k, i) to $(l, i+1)$ to be $\rightarrow E[l][x_{i+1}] \times A[k][l]$
 Why? compare $p_{k,i}$ to $p_{l,i+1}$

↑ output prob. ↑ transition prob.

$$\begin{aligned}
 p_{l,i+1} &= \prod_{j=1}^{i+1} e_{\pi_j}(x_j) \cdot a_{\pi_{j-1}, \pi_j} \\
 &= \left(\prod_{j=1}^i e_{\pi_j}(x_j) \cdot a_{\pi_{j-1}, \pi_j} \right) \cdot (e_{\pi_{i+1}}(x_{i+1}) \cdot a_{\pi_i, \pi_{i+1}}) \\
 &= p_{k,i} \cdot e_l(x_{i+1}) \cdot a_{kl} \\
 &= p_{k,i} \cdot \text{weight of edge from } (k, i) \text{ to } (l, i+1)
 \end{aligned}$$

This changes the decoding problem into a longest path problem (in a DAG).
(ie Manhattan distance)

Then - use dynamic programming longest path solution from chapter 6.

Note: Now we need to multiply weights (& not add) to get length.

Trick: take log!

$$\log(ab) = \log a + \log b$$

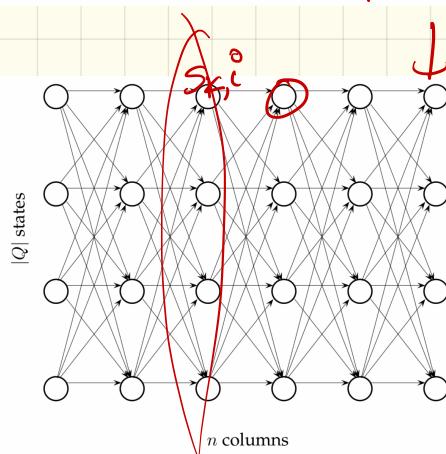
Dynamic program:

For any state l :

$$\begin{aligned}s_{l,i+1} &= \max_{k \in Q} \{ s_{k,i} \cdot \text{weight of edge between } (k,i) \text{ and } (l,i+1) \} \\&= \max_{k \in Q} \{ s_{k,i} \cdot a_{kl} \cdot e_l(x_{i+1}) \} \\&= e_l(x_{i+1}) \cdot \max_{k \in Q} \{ s_{k,i} \cdot a_{kl} \}\end{aligned}$$

Pick best weight here

Data Structure:



Runtime: $O(|Q| \times n)$ table
spots

Each takes $O(|Q|)$
time to fill in

$\Rightarrow O(|Q|^2 n)$

HMM Parameter estimation

Another problem:

probabilities are really often not known!

Need to be estimated from data.

So: we know Q (the states),
but not E or A .

Let Θ be a vector combining unknown transition & estimation ($E + A$) probabilities.

Given an observed x ,
let $P(x|\Theta) = \max$ probability
of x given Θ

Goal: $\max_{\Theta} P(x|\Theta)$

Instead of a single x , we often get a set of sequences in training x^1, \dots, x^m

Goal is then:

$$\max_{\Theta} \prod_{i=1}^m P(x^i | \Theta)$$

Hard! (optimization in a multi-dimensional space)

Heuristics & estimates:

If we knew states, +
 $a_{kl} = \#$ of transitions from k to l

$e_{kb} = \#$ times b outputted from state k

then set

$$A[k][l] = \frac{a_{kl}}{\sum_{q \in Q} a_{kq}}$$

$$E[k][b] = \frac{e_{kb}}{\sum_{\sigma \in \Sigma} e_{k\sigma}}$$

Note: We don't know States!

- Usually start with a guess
- Compute empirical estimates for transmission & estimation probabilities
- Solve problem to get a less wild guess

Ex: Baum-Welch algorithm

• GLIMMER

• GENSCAN

many others based on similar iterative, local improvement strategy

Next time — some applications of HMM

& on to Randomization